

Solutions to “Pattern Classification” by Duda et al.

tommyod @ github

December 11, 2018

Abstract

This document contains solutions to selected exercises from the book “Pattern Recognition” by Richard O. Duda, Peter E. Hart and David G. Stork. Although it was written in 2001, the second edition has truly stood the test of time—it’s a much-cited, well-written introductory text to the exciting field of *pattern recognition* (or simply *machine learning*). At the time of writing, the book has close to 40 000 citations according to Google.

While short chapter summaries are included in this document, they are not intended to substitute the book in any way. The summaries will largely be meaningless without the book, which I recommend buying if you’re interested in the subject.

The solutions and notes were typeset in \LaTeX to facilitate my own learning process. Machine learning has rightfully garnered considerable attention in recent years, and while many online resources are worthwhile it seems reasonable to favor books when attempting to learn the material thoroughly.

I hope you find my solutions helpful if you are stuck. Remember to make an attempt at solving the problems yourself before peeking. More likely than not, the solutions can be improved by a reader such as yourself. If you would like to contribute, please submit a pull request at <https://github.com/tommyod/lml/>.

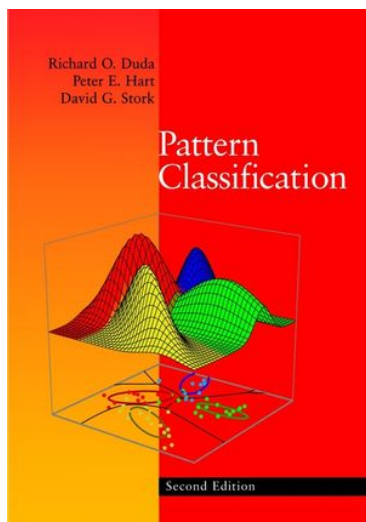


Figure 1: The front cover of [Duda et al., 2000].

Contents

- 1 Notes from “Pattern Classification” 3**
- 1.2 Bayesian Decision Theory 3
- 1.3 Maximum-likelihood and Bayesian parameter estimation 3
- 1.4 Nonparametric techniques 5
- 1.5 Linear discriminant functions 6
- 1.6 Multilayer Neural Networks 7
- 1.7 Stochastic methods 8
- 1.8 Nonmetric methods 9
- 1.9 Algorithm-independent machine learning 11
- 1.10 Unsupervised learning and clustering 12

- 2 Solutions to “Pattern Classification” 14**
- 2.2 Bayesian Decision Theory 14
- 2.3 Maximum-likelihood and Bayesian parameter estimation 23
- 2.4 Nonparametric techniques 34
- 2.5 Linear discriminant functions 40
- 2.6 Multilayer Neural Networks 48
- 2.7 Stochastic methods 55
- 2.8 Nonmetric methods 58
- 2.9 Algorithm-independent machine learning 64
- 2.10 Unsupervised learning and clustering 69

1 Notes from “Pattern Classification”

1.2 Bayesian Decision Theory

- Bayes theorem is

$$P(\omega_j | \mathbf{x}) = \frac{p(\mathbf{x} | \omega_j) \times P(\omega_j)}{p(\mathbf{x})} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}.$$

The Bayes decision rule is to choose the state of nature ω_m such that

$$\omega_m = \arg \max_j P(\omega_j | \mathbf{x}).$$

- Loss functions (or risk functions) with losses other than zero-one are possible. In general, we choose the action λ to minimize the risk $R(\lambda | \mathbf{x})$.
- The multivariate normal density (the Gaussian) is given by

$$p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right].$$

It is often analytically tractable, and closed form discriminant functions exist.

- If features \mathbf{y} are missing, we integrate them out (marginalize) using the sum rule

$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{y}) d\mathbf{y} = \int p(\mathbf{x} | \mathbf{y}) p(\mathbf{y}) d\mathbf{y}.$$

- In *Bayesian belief networks*, influences are represented by a directed network. If B is dependent on A , we add a directed edge $A \rightarrow B$ to the network.

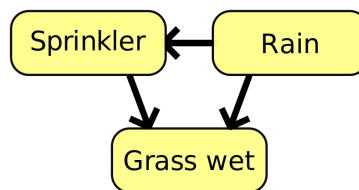


Figure 2: A Bayesian belief network. The source is Wikipedia.

1.3 Maximum-likelihood and Bayesian parameter estimation

- The *maximum likelihood* of a distribution $p(\mathbf{x} | \boldsymbol{\theta})$ is given by $\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} p(\mathcal{D} | \boldsymbol{\theta})$, assuming i.i.d. data points and maximizing the log-likelihood, we have

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \ln p(\mathcal{D} | \boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \ln \prod_{i=1}^n p(\mathbf{x}_i | \boldsymbol{\theta}).$$

Analytical solutions exist for the Gaussian. In general a maximum likelihood estimate may be biased, in the sense that $\mathbb{E}_{\mathbf{x}}[\hat{\boldsymbol{\theta}}] = \int \hat{\boldsymbol{\theta}} p(\mathbf{x}) d\mathbf{x} \neq \boldsymbol{\theta}$.

- In the Bayesian framework, the parameter $\boldsymbol{\theta}$ is expressed by a probability density function $p(\boldsymbol{\theta})$. This is called the *prior* distribution of $\boldsymbol{\theta}$, which is updated when new data is observed. The result is called the *posterior* distribution, given by

$$p(\boldsymbol{\theta} | \mathcal{D}) = \frac{p(\mathcal{D} | \boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathcal{D})} = \frac{p(\mathcal{D} | \boldsymbol{\theta})p(\boldsymbol{\theta})}{\int p(\mathcal{D} | \boldsymbol{\theta})p(\boldsymbol{\theta}) d\boldsymbol{\theta}}.$$

The estimate of \boldsymbol{x} then becomes

$$p(\boldsymbol{x} | \mathcal{D}) = \int p(\boldsymbol{x}, \boldsymbol{\theta} | \mathcal{D}) d\boldsymbol{\theta} = \int p(\boldsymbol{x} | \boldsymbol{\theta}, \mathcal{D})p(\boldsymbol{\theta} | \mathcal{D}) d\boldsymbol{\theta} = \int p(\boldsymbol{x} | \boldsymbol{\theta})p(\boldsymbol{\theta} | \mathcal{D}) d\boldsymbol{\theta},$$

which may be interpreted as a weighted average of models $p(\boldsymbol{x} | \boldsymbol{\theta})$, where $p(\boldsymbol{\theta} | \mathcal{D})$ is the weight associated with the model.

- The Bayesian framework is analytically tractable when using Gaussians. For instance, we can compute $p(\boldsymbol{\mu} | \mathcal{D})$ if we assume $p(\boldsymbol{\mu}) \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$. The distribution $p(\boldsymbol{\mu})$ is called a *conjugate prior* and $p(\boldsymbol{\mu} | \mathcal{D})$ is a *reproducing density*, since a normal prior transforms to a normal posterior (with different parameters) when new data is observed.
- In summary the Bayesian framework allows us to incorporate prior information, but the maximum-likelihood approach is simpler. Maximum likelihood gives us an estimate $\hat{\boldsymbol{\theta}}$, but the Bayesian framework gives us $p(\boldsymbol{\theta} | \mathcal{D})$ —the full distribution.
- *Principal Component Analysis* (PCA) yields components useful for *representation*. The covariance matrix is diagonalized, and low-variance directions in the hyperellipsoid are eliminated. The computation is often performed using the *Singular Value Decomposition* (SVD).
- *Discriminant Analysis* (DA) projects to a lower dimensional subspace with optimal *discrimination* (and not representation).
- *Expectation Maximization* (EM) is an iterative algorithm for finding the maximum-likelihood when data is missing (or latent).

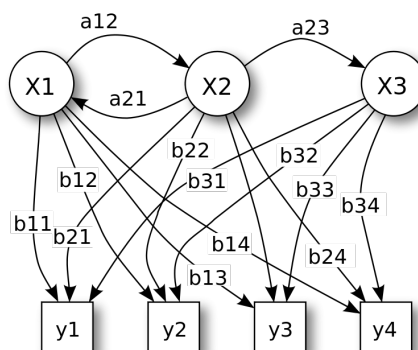


Figure 3: A hidden Markov model. The source is Wikipedia.

- A discrete, first order, hidden Markov model consists of a transition matrix \mathbf{A} and an emission matrix \mathbf{B} . The probability of transition from state i to state j is given by

a_{ij} , and the probability that state i emits signal j is given by b_{ij} . Three fundamental problems related to Markov models are:

- The evaluation problem - probability that \mathbf{V}^T was emitted, given \mathbf{A} and \mathbf{B} .
- The decoding problem - determine most likely sequence of hidden states $\boldsymbol{\omega}^T$, given emitted \mathbf{V}^T , \mathbf{A} and \mathbf{B} .
- The learning problem - determine \mathbf{A} and \mathbf{B} given training observations of \mathbf{V}^T and a coarse model.

1.4 Nonparametric techniques

- Two conceptually different approaches to nonparametric pattern recognition are:
 - Estimation of densities $p(\mathbf{x} | w_j)$, called the *generative* approach.
 - Estimation of $P(w_j | \mathbf{x})$, called the *discriminative* approach.
- Parzen-windows (kernel density estimation) is a generative method. It places a *kernel function* $\phi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ on every data point \mathbf{x}_i to create a density estimate

$$p_n(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \frac{1}{V_n} \phi \left(\frac{\|\mathbf{x} - \mathbf{x}_i\|_p}{h_n} \right),$$

where $\|\cdot\|_p : \mathbb{R}^d \rightarrow \mathbb{R}_+$ is the p -norm (which induces the so-called *Minkowski metric*) and $h_n > 0$ is the bandwidth.

- k -nearest neighbors is a discriminative method. It uses information about the k nearest neighbors of a point \mathbf{x} to compute $P(w_j | \mathbf{x})$. This automatically uses more of the surrounding space when data is sparse, and less of the surrounding space when data is dense. The k -nearest neighbor estimate is given by

$$P(w_j | \mathbf{x}) = \frac{\# \text{ samples labeled } w_j}{k}.$$

- The *nearest neighbor method* uses $k = 1$. It can be shown that the error rate P of the nearest neighbor method is never more than twice the Bayes error rate P^* in the limit of infinite data. More precisely, we have $P^* \leq P \leq P^*(2 - \frac{c}{c-1}P^*)$.
- In some applications, careful thought must be put into metrics. Examples include periodic data on \mathbb{R}/\mathbb{Z} and image data where the metric should be invariant to small shifts and rotations. One method to alleviate the problems of using the 2-norm as a metric on images is to introduce the *tangent distance*. For an image \mathbf{x}' , the tangent vector of a transformation \mathcal{F} (such as rotation by an angle α_i) is given by

$$\mathbf{TV}_i = \mathcal{F}(\mathbf{x}'; \alpha_i) - \mathbf{x}'.$$

If several transformations are available, their linear combination may be computed. For each test point \mathbf{x} , we search the tangent space for the linear combination minimizing the metric. This gives a metric $D(\mathbf{x}, \mathbf{x}')$ which is invariant to transformations such as small rotations and translations, compared to the 2-norm.

- *Reduced Coloumb energy networks* use ideas from both Parzen windows and k -nearest neighbors. It adjusts the size of the window so that it is less than some maximal radius, while not touching any observation of a different class. This creates “basins of attraction” for classification.

1.5 Linear discriminant functions

- A *linear discriminant function* splits the feature space in two using a hyperplane. The equation for a hyperplane is given by

$$g(\mathbf{x}) = \boldsymbol{\omega}^T \mathbf{x} + \omega_0 = \mathbf{a}^T \mathbf{y} = (\omega_0 \quad \boldsymbol{\omega}) \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix},$$

where ω_0 is the bias. The expression $\mathbf{a}^T \mathbf{y}$ is called the *augmented form*.

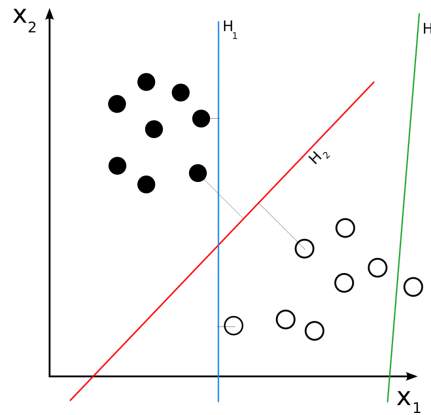


Figure 4: Linear discriminant functions. The source is Wikipedia.

- A *linear machine* assigns a point \mathbf{x} to state of nature ω_i if

$$g_i(\mathbf{x}) \geq g_j(\mathbf{x})$$

for every other class j . This leaves no ambiguous regions in the feature space.

- By introducing mappings $\mathbf{y} = h(\mathbf{x})$ to higher- or lower-dimensional spaces, nonlinearities in the original \mathbf{x} -space may be captured by linear classifiers working in \mathbf{y} -space. An example is $\mathbf{y} = h(\mathbf{x}) = \exp(-\mathbf{x}^T \mathbf{x})$ if data from one class is centered around the origin. Another example is transforming periodic data with period P from $0 \leq x < P$ to \mathbf{y} by use of the functions

$$y_1 = \cos(2\pi x/P) \quad y_2 = \sin(2\pi x/P).$$

- Several algorithms may be used to minimize an error function $J(\mathbf{a})$. Two popular choices are *gradient descent* and *Newton descent*.
 - Gradient descent moves in the direction of the negative gradient. It is often controlled by a step length parameter $\eta(k)$, which may decrease as the iteration counter k increases. The update rule is given by

$$\mathbf{a} \leftarrow \mathbf{a} - \eta(k) \nabla J(\mathbf{a}).$$

- Newton descent also moves in the direction of the negative gradient, but the optimal step length is computed by linearizing the function $\nabla J(\mathbf{a})$ (or, equivalently, a second order approximation of $J(\mathbf{a})$). The update rule is given by

$$\mathbf{a} \leftarrow \mathbf{a} - \mathbf{H}^{-1} \nabla J(\mathbf{a}).$$

- Criterion functions for linearly separable data sets include:
 - The Perceptron function $\sum_{y \in \mathcal{Y}} (-\mathbf{a}^T \mathbf{y})$, which is not smooth.
 - The squared error with margin, given by $\sum_{y \in \mathcal{Y}} (\mathbf{a}^T \mathbf{y} - b)^2 / \|\mathbf{y}\|^2$.
- The *Mean Squared Error* (MSE) approach may be used, but it is not guaranteed to yield a separating hyperplane—even if one exists.
 - The MSE solution is found analytically by the pseudoinverse $\mathbf{A}^\dagger = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$. The pseudoinverse should never be used explicitly because it's numerically wasteful and unstable. It represents the analytical solution to the problem

$$\min_x \mathbf{e}^T \mathbf{e} = \min_x (\mathbf{b} - \mathbf{A}\mathbf{x})^T (\mathbf{b} - \mathbf{A}\mathbf{x}),$$

which is solved by $\mathbf{x} = \mathbf{A}^\dagger \mathbf{b}$.

- The MSE approach is related to Fisher's linear discriminant for an appropriate choice of margin vector \mathbf{b} .
- LMS may be computed using matrix procedures (never use the pseudoinverse directly) or by the gradient descent algorithm.
- Ho-Kashyap procedures will return a separating hyperplane if one exists.
- *Linear programming* (LP) may also be used to find a separating hyperplane. Several reductions are possible by introducing *artificial variables*.
 - Minimizing the Perceptron criterion function may be formulated as an LP, and the result is typically decent even if a separating hyperplane does not exist.
- *Support Vector Machines* (SVM) find the minimum margin hyperplane. This is a *quadratic programming* (QP) problem, and the dual problem is easier to solve than the primal problem.

1.6 Multilayer Neural Networks

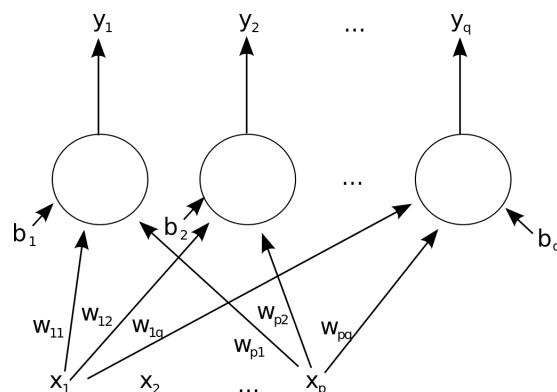


Figure 5: A three-layer neural network with bias. The source is Wikipedia.

- The feedforward operation on a $d - n_H - c$ three-layer neural network is defined by the following equation for the output

$$z_k = f \left(\sum_{j=1}^{n_H} w_{kj} f \left(\sum_{i=1}^d w_{ji} x_i + w_{j0} \right) + w_{k0} \right).$$

- The *Kolmogorov-Arnold representation theorem* implies that *any* continuous function from input to output may be expressed by a three layer $d - n_H - c$ neural network with sufficiently many hidden units.
- Backpropagation learns the weights \mathbf{w} by the gradient descent equation $\mathbf{w}_{n+1} = \mathbf{w}_n - \eta \nabla J(\mathbf{w}_n)$. The gradient, or derivative, is found by repeated application of the chain rule of calculus. Several protocols are available for backpropagation: stochastic, batch and on-line.
 - Backpropagation may be thought of as feature mapping. While the inputs x_i are not necessarily linearly separable, the outputs y_j of the hidden units become linearly separable as the weights are learned. The final linear discriminant works on this data instead of the x_i .
- Practical tips for improving learning in neural networks include: standardizing features, adding noise and data augmentations, initializing weights to random values in the range $-1/\sqrt{d} < w_{ji} < 1/\sqrt{d}$, using momentum in the gradient descent algorithm, adding weight decay (equivalent to regularization) while learning, training with hints (output units which are subsequently removed) and experimenting with various error functions.
- Second order methods for learning the weights include
 - **Newtons method** – uses \mathbf{H} in addition to $\nabla J(\mathbf{w})$.
 - **Quickprop** – two evaluations of $\nabla J(\mathbf{w})$ to approximate a quadratic.
 - **Conjugate gradient descent** – uses conjugate directions, which consists of a series of line searches. A given search direction does not spoil the result of the previous line searches. This is equivalent to a “smart momentum.”
- Other networks include:
 - Convolutional Neural Networks (CNNs) – translation invariant, has achieved great success on image data.
 - Recurrent Neural Networks (RNNs) – the output of the previous prediction is fed into the subsequent prediction. This simulates memory, and RNNs have been successful on time series data.
 - Cascade correlation – a technique where the topology is altered by adding more units until the performance is sufficiently good.

1.7 Stochastic methods

- Stochastic methods are used to search for optimal solutions when techniques such as gradient descent are not viable. For instance if the model is very complex, has a discrete nature where gradients are not available, or if there are time constraints.
- *Simulated annealing* is an optimization technique. As an example: to minimize the error $E(\mathbf{s})$, where $\mathbf{s} \in [-1, 1]^n$, we change a random entry of \mathbf{s} .
 - If the change produces a better result, then keep the new \mathbf{s} .
 - If the change does not produce a better result, we still might keep the change.

The probability of keeping a change which increases the error $E(\mathbf{s})$ is a function of the temperature T , which typically decreases exponentially as the algorithm pro-

gresses. Initially simulated annealing is a *random search*, and as the temperature progresses it becomes a *greedy search*.

- *Deterministic simulated annealing* replaces discrete s_i with analog (continuous) s_i . This forces the other magnets s_k ($k \neq i$) to determine s_i

$$s_i = f(T, \ell_i) = \tanh\left(\frac{\ell_i}{T}\right).$$

As $T \rightarrow 0$, the $\tanh(\cdot)$ sigmoid function converges to a step function.

- *Boltzmann networks* (or Boltzmann machines) employ simulated annealing in a network to make predictions. First, weights w_{ij} are learned so that inputs $s_j \in \alpha^i$ lead to correct outputs $s_k \in \alpha^o$ during classification. In the classification phase, the inputs α^i are *clamped* (fixed), and simulated annealing produces outputs α^o . If the weights w_{ij} are learned correctly, then the algorithm will produce good classifications.
 - Boltzmann networks are able to perform *pattern completion*.

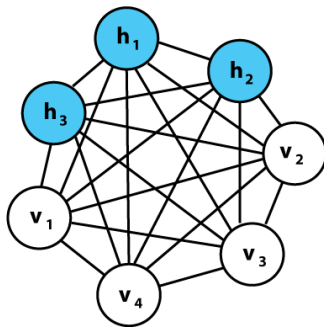


Figure 6: A non-restricted Boltzmann machine. The source is Wikipedia.

- *Evolutionary methods* take a population of classifiers through many generations. In each generation, new classifiers (offspring) are produced from the previous generation. The best classifiers are subject to (1) replication, (2) crossover and (3) mutation to produce offspring. The classifiers may be encoded as 2-bit chromosomes of length L . The bits represent some property of the classifier.
- Genetic programming is the process of modifying formulas such as

$$[(-x_1) + x_2] / [(\ln x_3) - x_2]$$

by evolutionary methods, mutating variables and operators and performing crossovers.

1.8 Nonmetric methods

- A *decision tree* typically splits the feature space along the axes if the data is numeric, and works well for non-metric (categorical, or nominal) data as well. To implement a decision tree, one must consider
 - The number of splits made per node (typically 2, since a higher branching factor B may be reduced to $B = 2$ anyway).

- How to choose an attribute to split on – often solved using *information gain*.
- When a node should be declared a leaf node.
- How to handle missing data.
- To decide which attribute to split on, an *impurity function* is defined for a node N consisting of training samples from various categories $\omega_1, \dots, \omega_C$. The impurity function should be 0 when all samples in node N are from the same category ω_j , and peak when the samples are uniformly drawn from the C categories.

Two examples of impurity functions are *entropy impurity* and *Gini impurity*, which are respectively defined as

$$i(N) = - \sum_{j=1}^C P(\omega_j) \ln P(\omega_j) \quad i(N) = \sum_{j=1}^C P(\omega_j) \sum_{k \neq j} P(\omega_k) = \sum_{j=1}^C P(\omega_j) [1 - P(\omega_j)].$$

A split is chosen so that it maximizes the decrease in impurity, i.e.

$$\Delta i(N) = i(N) - [P_L i(N_L) + (1 - P_L) i(N_R)].$$

The above equation says that the change in impurity equals the original impurity at node N minus the weighted average of the impurity of the left and right child node.

- Other considerations in decision trees include:
 - Pruning – simplifying the tree after training (bypassing the horizon effect).
 - Penalizing complexity – regularization of the tree structure.
 - Missing attributes – for instance using (1) surrogate splits or (2) sending a training sample down every path and then performing a weighted average.
- Four *string problems* in pattern classification are:
 - **Matching**: naive matching is slow, the *Boyer-Moore string matching algorithm* is much more efficient. It operates by increasing the shift s of \mathbf{x} using two heuristics in parallel: the *bad-character heuristic* and the *good-suffix heuristic*.
 - **Edit distance**: a way to compare the “distance” between strings by counting the number of insertions, deletions and substitutions required to transform \mathbf{x} to \mathbf{y} . If all costs are equal, then $D(\mathbf{x}, \mathbf{y})$ is a metric. A dynamic programming algorithm is used to compute edit distance.
 - **Matching with errors** is the same as matching, but using for instance the edit distance to find approximate matches. The problem is to find a shift s that minimizes the edit distance.
 - **Matching with the “don’t care”-symbol \emptyset** : same as matching, but the \emptyset -symbol matches any character in the alphabet \mathcal{A} .
- A grammar $G = (\mathcal{A}, \mathcal{I}, \mathcal{S}, \mathcal{P})$ consists of symbols \mathcal{A} , variables \mathcal{I} , a root symbol \mathcal{S} and productions \mathcal{P} . Concrete examples include English sentences and pronunciation of numbers. There are several types of grammars, and they constitute a hierarchy

$$\text{Type 3} \subset \text{Type 2} \subset \text{Type 1} \subset \text{Type 0}.$$

The types are respectively called regular, context free, context sensitive and free.

- A central question is whether a string \mathbf{x} is in the language \mathcal{L} generated by the grammar G , i.e. whether $\mathbf{x} \in \mathcal{L}(G)$. This can be answered using bottom-up parsing, which employs the product rules \mathcal{P} backwards.

1.9 Algorithm-independent machine learning

- The *no free lunch theorem* states that using the off-training set error, and assuming equally likely target functions $F(\mathbf{x})$, no algorithm is universally superior. Any statement about algorithms is a statement about the target function $F(\mathbf{x})$.
- The *ugly duckling theorem* states that no feature representation is universally superior. If pattern similarity is based on the number of possible shared predicates (f_1 OR f_2 , etc) then any two patterns are equally similar. The best feature representation consequently depends on the target function $F(\mathbf{x})$.
- The *minimum description length* is a version of Occam's razor. In this framework, the best hypothesis h (model) is the one compressing the data the most, i.e. the minimizer h^* of

$$K(h, \mathcal{D}) = K(h) + K(\mathcal{D} \text{ using } h),$$

where $K(\cdot)$ is the Kolmogorov complexity (length of smallest computer program).

- The *bias-variance trade-off* is exemplified by the equation

$$\mathbb{E}_{\mathcal{D}} [(g(\mathbf{x}; \mathcal{D}) - F(\mathbf{x}))^2] = \underbrace{\mathbb{E}_{\mathcal{D}} [g(\mathbf{x}; \mathcal{D}) - F(\mathbf{x})]^2}_{\text{bias}^2} + \underbrace{\mathbb{E}_{\mathcal{D}} [(g(\mathbf{x}; \mathcal{D}) - \mathbb{E}[g(\mathbf{x}; \mathcal{D})])^2]}_{\text{variance}},$$

and informally states that there is always a trade-off between bias and variance. In the regression setting, the bias is the average error over many data sets, while the variance is the variability of the estimate over many data sets. High variance typically implies many parameters (over-fitting), while high bias implies few parameters (under-fitting).

- The *Jackknife* re-sampling method involves removing the i th data point from $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ and computing a statistic. This is done for every data point, and the final result is averaged. From this the variance of the statistic may be assessed.
- The *bootstrap* method involves re-sampling n data points from \mathcal{D}^n with replacement B times. The statistic is computed B times, and is then averaged.

$$\hat{\theta}^{*(\cdot)} = \frac{1}{B} \sum_{b=1}^B \hat{\theta}^{*(b)} \quad \text{var}_{\text{boot}}[\hat{\theta}] = \frac{1}{B} \sum_{b=1}^B (\hat{\theta}^{*(b)} - \hat{\theta}^{*(\cdot)})^2$$

- *Bagging* consists of taking an average over model predictions. *Boosting* uses the result of model h^n to train model h^{n+1} . Informally, the next model prioritizes data in \mathcal{D} which the sequence of models so far has not performed well on.
 - Adaptive Boosting (*AdaBoost*) is a well known boosting algorithm. It samples data \mathbf{x}_i using probability weights w_i . If a model h^n does a poor job, the associated weights w_i^{n+1} are increased exponentially. This causes exponential error decay on training data.
- *Learning with queries* involves the model choosing the next data point to learn from. It assumes the existence of an oracle which can give the correct answer to any input, but using this oracle might be costly. Efficient learning involves choosing points where the classifier is uncertain, i.e. where $P(\omega_1 | \mathbf{x}) \approx P(\omega_2 | \mathbf{x})$.

- *Maximum-likelihood model comparison* picks the maximum of

$$P(\mathcal{D} | h_i) \simeq \underbrace{P(\mathcal{D} | \hat{\boldsymbol{\theta}}, h_i)}_{\text{best-fit likelihood}} \underbrace{P(\hat{\boldsymbol{\theta}} | h_i) \Delta \boldsymbol{\theta}}_{\text{Occam factor}} = P(\mathcal{D} | \hat{\boldsymbol{\theta}}, h_i) \frac{\Delta \boldsymbol{\theta}}{\Delta^0 \boldsymbol{\theta}},$$

where $\Delta \boldsymbol{\theta}$ is the volume of the possible parameter space given the data \mathcal{D} , and $\Delta^0 \boldsymbol{\theta}$ is the prior volume of the possible parameter before seeing data. Applying a few simplifications, this becomes a two step process: (1) compute maximum likelihood estimate $\hat{\boldsymbol{\theta}}$ from the data, then (2) compute the likelihood value given $\hat{\boldsymbol{\theta}}$.

- Mixture models can be used to combine classifiers using discriminant functions. The final discriminant function is $g = \sum_{i=1}^k w_i g_i(\mathbf{x}, \boldsymbol{\theta}_i)$, and an underlying mixture model of the following form is assumed:

$$p(\mathbf{y} | \mathbf{x}, \Theta^0) = \sum_{r=1}^k \underbrace{P(r | \mathbf{x}, \boldsymbol{\theta}_r^0)}_{\text{prob. of model } r} p(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}_r^0)$$

1.10 Unsupervised learning and clustering

- Unsupervised learning is a difficult problem. A mixture model may not be *identifiable*, which happens when $\boldsymbol{\theta}$ cannot be determined uniquely even in the limit of infinite data. There exist necessary conditions for ML solutions to $\hat{P}(w_i)$ and $\hat{\boldsymbol{\theta}}_i$, but they are not sufficient to guarantee that a maximum is found. Singular solutions may occur, since the likelihood can be made arbitrarily large when $\boldsymbol{\mu}_i$ is placed on a data point and $\boldsymbol{\sigma}_i \rightarrow 0$. A Bayesian approach is possible, but rarely feasible.
- The k -means algorithm is a simple procedure for finding $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k$. A fuzzy k -means algorithm also exists, in which membership categorization is not binary.
- It is better to find structure in data than to impose it. A metric $d(\cdot, \cdot)$ and a distance threshold d_0 may be defined, which may then be used to group points. This approach is sensitive to the effect of individual data points. Standardizing data using $\mathbf{y}_i = (\mathbf{x}_i - \boldsymbol{\mu}_k) / \boldsymbol{\sigma}_k$ may or may not be beneficial, and the same goes for PCA.
- Many criterion functions are available for clustering, e.g. sum of square error

$$J_e = \sum_{i=1}^c \sum_{\mathbf{x} \in \mathcal{D}_i} \|\mathbf{x} - \mathbf{m}_i\|^2.$$

Criterion functions based on the size of scatter matrices (LDA) are

$$J_e = \text{tr}[\mathbf{S}_W], \quad J_d = \det[\mathbf{S}_W] \quad \text{and} \quad J_f = \text{tr}[\mathbf{S}_T^{-1} \mathbf{S}_W].$$

- *Hierarchical clustering* represents hierarchical (clusters within clusters) structure in a *dendrogram*, see Figure 7 on page 13. Two approaches are possible: agglomerative (merging) or divisive (splitting). Many distance measures are available.
 - Using $d_{\min}(\mathcal{D}_i, \mathcal{D}_j)$ and agglomerative clustering, a *minimal spanning tree* (MST) is built from the data.

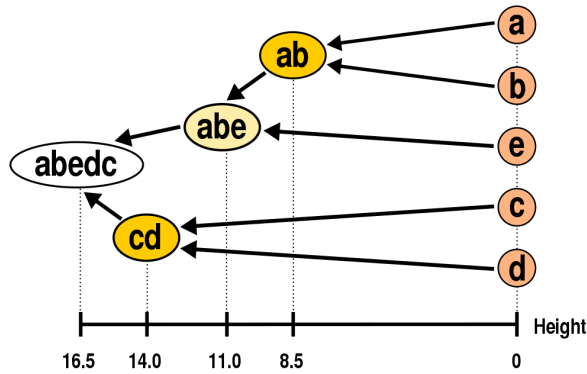


Figure 7: A dendrogram representation of hierarchical clustering. The source is Wikipedia.

- *Competitive learning* is a clustering algorithm implemented by a neural network. The data $\mathbf{x} \in \mathbb{R}^d$ is augmented to $\mathbf{x} := (1, \mathbf{x})$ and normalized so that it lies on a $(d + 1)$ -dimensional sphere. The weight $\mathbf{w}^* = \arg \max_{\mathbf{w}} \mathbf{w}^T \mathbf{x}$ is updated to be more similar to the sample \mathbf{x} . The algorithm can also create new clusters \mathbf{w}_j if certain conditions are met. It's suitable for on-line learning, but will not necessarily converge with constant learning rate.
- PCA can be implemented as a 3-layer neural network autoencoder. The first k eigenvalues/eigenvectors of \mathbf{X} solves the minimization of $\|\mathbf{X} - \mathbf{X}'_k\|$, where \mathbf{X}'_k is a rank k approximation of \mathbf{X} . *Non-linear PCA* is an extension of the idea, using a 5-layer non-linear neural network autoencoder.

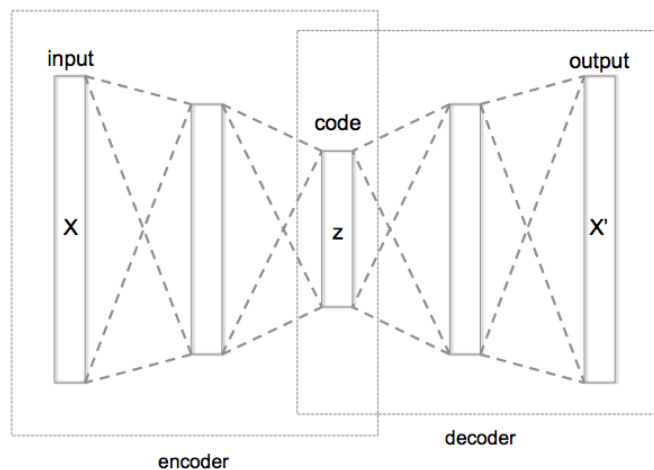


Figure 8: A neural network autoencoder with 5 layers. The source is Wikipedia.

- *Self-organizing feature maps* transform from a high dimensional space to low dimensional one, while preserving local topological information. They are implemented as neural networks, and a key idea is to update neighboring areas when learning the map. An additional benefit is that they automatically place more points in regions of high probability.

2 Solutions to “Pattern Classification”

This section contains solutions to problems in “Pattern Classification” by Duda et al. There are approximately 8 solved problems from each chapter.

2.2 Bayesian Decision Theory

Problem 2.6

- a) We want the probability of choosing action α_2 to be smaller than, or equal to, E_1 , given that the true state of nature is ω_1 . Let's assume that $\mu_1 < \mu_2$ and that the decision threshold is x^* , so we decide α_2 if $x > x^*$. We then have

$$\begin{aligned} P(\alpha_2 | \omega_1) &\leq E_1 \\ p(x > x^* | \omega_1) &\leq E_1 \\ \left[1 - \int_0^{x^*} p(x | \omega_1) dx \right] &\leq E_1 \end{aligned}$$

We let $\Phi : \mathbb{R} \rightarrow [0, 1]$ denote the cumulative Gaussian distribution, and $\Phi^{-1} : [0, 1] \rightarrow \mathbb{R}$ its inverse function. Making use of Φ we write

$$\begin{aligned} 1 - \Phi\left(\frac{x^* - \mu_1}{\sigma_1}\right) &\leq E_1 \\ x^* &\geq \mu_1 + \sigma_1 \Phi^{-1}(1 - E_1). \end{aligned}$$

If the desired error is close to zero, then x^* goes to positive infinity. If the desired error is close to one, then x^* goes to negative infinity.

- b) The error rate for classifying ω_2 as ω_1 is

$$P(\alpha_1 | \omega_2) = p(x \leq x^* | \omega_2) = \int_0^{x^*} p(x | \omega_2) dx = \Phi\left(\frac{x^* - \mu_2}{\sigma_2}\right).$$

Making use of x^* from the previous problem, we obtain

$$\Phi\left(\frac{\mu_1 + \sigma_1 \Phi^{-1}(1 - E_1) - \mu_2}{\sigma_2}\right) = \Phi\left(\frac{\mu_1 - \mu_2}{\sigma_2} + \frac{\sigma_1}{\sigma_2} \Phi^{-1}(1 - E_1)\right).$$

- c) The overall error rate becomes

$$\begin{aligned} P(\text{error}) &= P(\alpha_1, \omega_2) + P(\alpha_2, \omega_1) \\ &= P(\alpha_1 | \omega_2)P(\omega_2) + P(\alpha_2 | \omega_1)P(\omega_1) \\ &= \frac{1}{2} [P(\alpha_1 | \omega_2) + P(\alpha_2 | \omega_1)] \\ &= \frac{1}{2} \left[E_1 + \Phi\left(\frac{\mu_1 - \mu_2}{\sigma_2} + \frac{\sigma_1}{\sigma_2} \Phi^{-1}(1 - E_1)\right) \right]. \end{aligned}$$

In the last equality we used the results from the previous subproblems.

d) We substitute the given values into the equations, and obtain $x^* \approx 0.6449$. The total error rate is $P(\text{error}) \approx 0.2056$.

e) The Bayes error rate, as a function of x^* , is given by

$$\begin{aligned} P(\text{error}) &= P(\alpha_2 | \omega_1)P(\omega_1) + P(\alpha_1 | \omega_2)P(\omega_2) \\ &= \frac{1}{2} [p(x > x^* | \omega_1) + p(x < x^* | \omega_2)] \\ &= \frac{1}{2} \left[\left(1 - \Phi \left(\frac{x^* - \mu_1}{\sigma_1} \right) \right) + \Phi \left(\frac{x^* - \mu_2}{\sigma_2} \right) \right] \end{aligned}$$

The Bayes error rate for this problem is depicted in Figure 9.

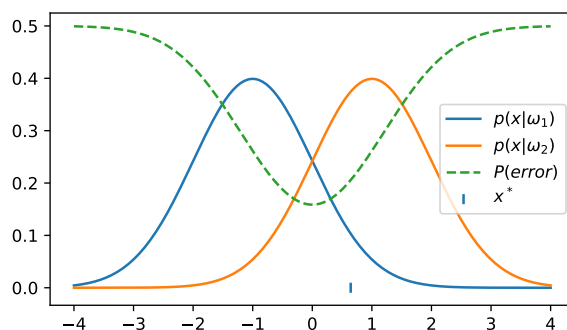


Figure 9: Graf accompanying problem 2.6.

Problem 2.12

a) A useful observation is that the maximal value $P(\omega_{\max} | \mathbf{x})$ is greater than, or equal to, the average. Therefore we obtain

$$P(\omega_{\max} | \mathbf{x}) \geq \frac{1}{c} \sum_{i=1}^c P(\omega_i | \mathbf{x}) = \frac{1}{c},$$

where the last equality is due to probabilities summing to unity.

b) The minimum error rate is achieved by choosing ω_{\max} , the most likely state of nature. The average probability of error over the data space is therefore the probability that ω_{\max} is *not* the true state of nature for a given \mathbf{x} , i.e.

$$P(\text{error}) = \mathbb{E}_{\mathbf{x}} [1 - P(\omega_{\max} | \mathbf{x})] = 1 - \int P(\omega_{\max} | \mathbf{x}) p(\mathbf{x}) d\mathbf{x}.$$

c) We see that

$$P(\text{error}) = 1 - \int P(\omega_{\max} | \mathbf{x}) p(\mathbf{x}) d\mathbf{x} \leq 1 - \int \frac{1}{c} p(\mathbf{x}) d\mathbf{x} = 1 - \frac{1}{c} = \frac{c-1}{c},$$

where we have used the fact that $\int p(\mathbf{x}) d\mathbf{x} = 1$.

d) A situation where $P(\text{error}) = (c-1)/c$ arises when $P(\omega_i) = 1/c$ for every i . Then the maximum value is equal to the average value, and the inequality in problem a) becomes an equality.

Problem 2.19

- a) The entropy is given by $H[p(x)] = - \int p(x) \ln p(x) dx$. The optimization problem gives the synthetic function (or Lagrange function)

$$H_s = - \int p(x) \ln p(x) dx + \sum_{k=1}^q \lambda_k \left(\int b_k(x) p(x) dx - a_k \right),$$

and since a probability density function has $\int p(x) dx = 1$ we add an additional constraint for $k = 0$ with $b_0(x) = 1$ and $a_k = 1$. Collecting terms we obtain

$$\begin{aligned} H_s &= - \int p(x) \ln p(x) dx + \sum_{k=0}^q \lambda_k \int b_k(x) p(x) dx - \sum_{k=0}^q \lambda_k a_k \\ &= - \int p(x) \left[\ln p(x) - \sum_{k=0}^q \lambda_k b_k(x) \right] dx - \sum_{k=0}^q \lambda_k a_k, \end{aligned}$$

which is what we were asked to show.

- b) Differentiating the equation above with respect to $p(x)$ and equating it to zero we obtain

$$- \int \left(1 \left[\ln p(x) - \sum_{k=0}^q \lambda_k b_k(x) \right] + p(x) \left[\frac{1}{p(x)} \right] \right) dx = 0.$$

This integral is zero if the integrand is zero for every x , so we require that

$$\ln p(x) - \sum_{k=0}^q \lambda_k b_k(x) + 1 = 0,$$

and solving this equation for $p(x)$ gives the desired answer.

Problem 2.21

We are asked to compute the entropy of the (1) Gaussian distribution, (2) triangle distribution and (3) uniform distribution. Every probability density function (pdf) has $\mu = 0$ and standard deviation σ , and we must write every pdf parameterized using σ .

Gaussian We use the definition $H[p(x)] = - \int p(x) \ln p(x) dx$ to compute

$$H[p(x)] = - \int \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2} \frac{x^2}{\sigma^2}\right) \left[\ln\left(\frac{1}{\sqrt{2\pi}\sigma}\right) - \frac{1}{2} \frac{x^2}{\sigma^2} \right] dx.$$

Let us denote the constant by $K = \frac{1}{\sqrt{2\pi}\sigma}$ to simplify notation. We obtain

$$\begin{aligned} & - \int K \exp\left(-\frac{1}{2} \frac{x^2}{\sigma^2}\right) \left[\ln K - \frac{1}{2} \frac{x^2}{\sigma^2} \right] dx = \\ & -K \ln K \int \exp\left(-\frac{1}{2} \frac{x^2}{\sigma^2}\right) dx + K \int \frac{1}{2} \frac{x^2}{\sigma^2} \exp\left(-\frac{1}{2} \frac{x^2}{\sigma^2}\right) dx \end{aligned}$$

The first term evaluates to $-\ln K$, since it's the normal distribution with an additional factor $-\ln K$. The second term is not as straightforward. We change variables to $y = x/(\sqrt{2}\sigma)$, and write it as

$$K \int y^2 \exp(-y^2) \sqrt{2}\sigma dy,$$

and this integral is solved by using the following observation (from integration by parts):

$$\int 1e^{-x^2} dx = \underbrace{xe^{-x^2}}_{0 \text{ at } \pm\infty} - \int x(-2x)e^{-x^2} dx.$$

Using the above equation in reverse, we integrate as follows:

$$K\sqrt{2}\sigma \int y^2 \exp(-y^2) dy = K\sqrt{2}\sigma \frac{1}{2} \int \exp(-y^2) dy = K\sqrt{2}\sigma \frac{1}{2} \sqrt{\pi} = \frac{1}{2}$$

To recap, the first integral evaluated to $-\ln K$, and the second evaluated to $\frac{1}{2}$. The entropy of the Gaussian is therefore $1/2 + \ln \sqrt{2\pi}\sigma$.

Triangle The triangle distribution may be written in the form

$$f(x) = \begin{cases} h - \frac{hx}{b} & \text{if } |x| < b \\ 0 & \text{if } |x| \geq b, \end{cases}$$

where h is the height and b is the width to the left of, and to the right of, $x = 0$.

Since the integral must evaluate to unity, we impose the constraint $hb = 1$ and obtain $f(x; b) = \frac{1}{b} \left(1 - \frac{x}{b}\right)$. We wish to parameterize the triangle distribution using the standard deviation σ instead of width b . We can use $\text{var}(X) = \mathbb{E}(X^2) - \mathbb{E}(X)^2$ to find the variance, since in this case $\mathbb{E}(X)^2 = \mu^2 = 0$ as a result of the function being centered on $x = 0$. Computing $\mathbb{E}(X^2)$ yields $b^2/6$, so $b^2 = 6\sigma^2$. The revised triangle distribution then becomes

$$f(x; \sigma) = \begin{cases} \frac{1}{\sqrt{6}\sigma} \left(1 - \frac{x}{\sqrt{6}\sigma}\right) & \text{if } |x| < \sqrt{6}\sigma \\ 0 & \text{if } |x| \geq \sqrt{6}\sigma. \end{cases}$$

We set $k = \frac{1}{\sqrt{6}\sigma}$ to ease notation. Due to symmetry, we compute the entropy as

$$\mathbb{H}[f(x; \sigma)] = -2 \int_0^{\sqrt{6}\sigma} k(1 - kx) \ln(k(1 - kx)) dx.$$

Changing variables to $y = 1 - kx$ we obtain

$$\begin{aligned} -2 \int_{x=0}^{x=\sqrt{6}\sigma} ky (\ln k + \ln y) dx &= -2 \int_{y=1}^{y=0} ky (\ln k + \ln y) \left(\frac{1}{-k}\right) dy \\ -2 \int_0^1 y (\ln k + \ln y) dy &= -2 \int_0^1 y \ln k dy - 2 \int_0^1 y \ln y dy = -2 \left(\ln k - \frac{1}{4}\right), \end{aligned}$$

where the last integral can be evaluated using integration by parts. The entropy of the triangle distribution turns out to be $1/2 + \ln \sqrt{6}\sigma$.

Uniform Using the same logic as with the triangle distribution, we first normalize a uniform distribution, and then parameterize by σ , to obtain

$$u(x; \sigma) = \begin{cases} \frac{1}{2b} & \text{if } |x| < b \\ 0 & \text{if } |x| \geq b \end{cases} = \begin{cases} \frac{1}{2\sqrt{3}\sigma} & \text{if } |x| < \sqrt{3}\sigma \\ 0 & \text{if } |x| \geq \sqrt{3}\sigma. \end{cases}$$

Computing the entropy is more straightforward than in the case of the Gaussian and the triangle distribution. We simply evaluate the integral as

$$H[p(x)] = 2 \int_0^{\sqrt{3}\sigma} \frac{1}{2\sqrt{3}\sigma} \ln \frac{1}{2\sqrt{3}\sigma} dx = \ln 2\sqrt{3}\sigma.$$

Let's briefly compare the results of our computations as follows:

$$H_{\text{Gaussian}}(\sigma) = 1/2 + \ln \sqrt{2\pi}\sigma = \frac{1}{2} + \ln \sqrt{2\pi} + \ln \sigma \approx 1.4189 + \ln \sigma$$

$$H_{\text{Triangle}}(\sigma) = 1/2 + \ln \sqrt{6}\sigma = \frac{1}{2} + \ln \sqrt{6} + \ln \sigma \approx 1.3959 + \ln \sigma$$

$$H_{\text{Uniform}}(\sigma) = \ln 2\sqrt{3}\sigma = 0 + \ln 2\sqrt{3} + \ln \sigma \approx 1.2425 + \ln \sigma$$

This verifies that out of the three distributions, the Gaussian has the maximal entropy. This was expected, since the Gaussian maximizes the entropy over *any* continuous probability density function having a prescribed mean and variance.

Problem 2.23

- a) To solve this problem, we need to find the inverse matrix, the determinant, and $\mathbf{w} = \mathbf{x} - \boldsymbol{\mu}$.

$$\boldsymbol{\Sigma}^{-1} = \frac{1}{21} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 5 & -2 \\ 0 & -2 & 5 \end{pmatrix} \quad \det \boldsymbol{\Sigma} = 21 \quad \mathbf{w} = \mathbf{x} - \boldsymbol{\mu} = \begin{pmatrix} -0.5 \\ -2 \\ -1 \end{pmatrix}$$

The number of dimension d is 3. The solution is

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{\frac{3}{2}} 21^{\frac{1}{2}}} \exp\left(-\frac{1}{2} \mathbf{w}^T \boldsymbol{\Sigma}^{-1} \mathbf{w}\right) = \frac{1}{(2\pi)^{\frac{3}{2}} 21^{\frac{1}{2}}} \exp\left(-\frac{1}{2} \frac{1}{21} \frac{69}{4}\right).$$

- b) The eigenvalues of $\boldsymbol{\Sigma}$ are $\lambda_1 = 3$, $\lambda_2 = 7$ and $\lambda_3 = 21$. The corresponding eigenvectors are $\mathbf{v}_1 = (0, 1, -1)^T/\sqrt{2}$, $\mathbf{v}_2 = (0, 1, 1)^T/\sqrt{2}$ and $\mathbf{v}_3 = (1, 0, 0)^T$. The whitening transformation is therefore given by

$$\mathbf{A}_w = \boldsymbol{\Phi} \boldsymbol{\Lambda}^{-1/2} = \frac{1}{\sqrt{2}} \begin{pmatrix} 0 & 0 & \sqrt{2} \\ 1 & 1 & 0 \\ -1 & 1 & 0 \end{pmatrix} \begin{pmatrix} -\sqrt{3} & 0 & 0 \\ 0 & -\sqrt{7} & 0 \\ 0 & 0 & -\sqrt{21} \end{pmatrix}.$$

The rest of the numerical computations are skipped.

- c) Skipped.
d) Skipped.
e) We are going to examine if the p.d.f is unchanged when vectors are transformed with $\mathbf{T}^T \mathbf{x}$ and matrices with $\mathbf{T}^T \boldsymbol{\Sigma} \mathbf{T}$. Let's consider the term $(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$ in the exponent first. Substituting $\mathbf{x} \mapsto \mathbf{T}^T \mathbf{x}$, $\boldsymbol{\mu} \mapsto \mathbf{T}^T \boldsymbol{\mu}$ and $\boldsymbol{\Sigma} \mapsto \mathbf{T}^T \boldsymbol{\Sigma} \mathbf{T}$, we observe that

$$\begin{aligned}
& (\mathbf{T}^T \mathbf{x} - \mathbf{T}^T \boldsymbol{\mu})^T (\mathbf{T}^T \boldsymbol{\Sigma} \mathbf{T})^{-1} (\mathbf{T}^T \mathbf{x} - \mathbf{T}^T \boldsymbol{\mu}) \\
& (\mathbf{T}^T (\mathbf{x} - \boldsymbol{\mu}))^T (\mathbf{T}^T \boldsymbol{\Sigma} \mathbf{T})^{-1} \mathbf{T}^T (\mathbf{x} - \boldsymbol{\mu}) \\
& (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{T} (\mathbf{T}^T \boldsymbol{\Sigma} \mathbf{T})^{-1} \mathbf{T}^T (\mathbf{x} - \boldsymbol{\mu}) \\
& (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{T} \mathbf{T}^{-1} \boldsymbol{\Sigma}^{-1} \mathbf{T}^{-T} \mathbf{T}^T (\mathbf{x} - \boldsymbol{\mu}) \\
& (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}),
\end{aligned}$$

where we have used $(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$ and $(\mathbf{AB})^{-1} = \mathbf{B}^{-1} \mathbf{A}^{-1}$, which are basic facts from linear algebra. The density remains proportional when applying a linear transformation, but not unscaled, since the proportionality term $|\boldsymbol{\Sigma}|^{1/2}$ becomes $|\mathbf{T}^T \boldsymbol{\Sigma} \mathbf{T}|^{1/2} = |\mathbf{T}^T|^{1/2} |\boldsymbol{\Sigma}|^{1/2} |\mathbf{T}|^{1/2} = |\mathbf{T}| |\boldsymbol{\Sigma}|^{1/2}$.

- f) Here we use the eigendecomposition of a symmetric matrix. We assume that $\boldsymbol{\Sigma}$ is positive definite such that every eigenvalue is positive. We write $\boldsymbol{\Sigma} = \boldsymbol{\Phi} \boldsymbol{\Lambda} \boldsymbol{\Phi}^T$ and apply the whitening transformation.

$$\mathbf{A}_w^T \boldsymbol{\Sigma} \mathbf{A}_w = \mathbf{A}_w^T \boldsymbol{\Phi} \boldsymbol{\Lambda} \boldsymbol{\Phi}^T \mathbf{A}_w = (\boldsymbol{\Phi} \boldsymbol{\Lambda}^{-1/2})^T \boldsymbol{\Phi} \boldsymbol{\Lambda} \boldsymbol{\Phi}^T (\boldsymbol{\Phi} \boldsymbol{\Lambda}^{-1/2})$$

The matrix $\boldsymbol{\Phi}$ is orthogonal, so its transpose is the inverse. Using this fact and proceeding, we obtain

$$(\boldsymbol{\Phi} \boldsymbol{\Lambda}^{-1/2})^T \boldsymbol{\Phi} \boldsymbol{\Lambda} \boldsymbol{\Phi}^T (\boldsymbol{\Phi} \boldsymbol{\Lambda}^{-1/2}) = (\boldsymbol{\Lambda}^{-1/2})^T \boldsymbol{\Lambda} \boldsymbol{\Lambda}^{-1/2} = \boldsymbol{\Lambda}^{-1/2} \boldsymbol{\Lambda} \boldsymbol{\Lambda}^{-1/2} = \mathbf{I},$$

so the covariance is proportional to the identity matrix, as we were tasked to show. The normalization constant becomes 1, since the proportionality term becomes $|\mathbf{T}| |\boldsymbol{\Sigma}|^{1/2}$ under the transformation, and

$$|\mathbf{T}| |\boldsymbol{\Sigma}|^{1/2} = |\boldsymbol{\Phi} \boldsymbol{\Lambda}^{-1/2}| |\boldsymbol{\Sigma}|^{1/2} = |\boldsymbol{\Phi} \boldsymbol{\Lambda}^{-1/2}| |\boldsymbol{\Phi} \boldsymbol{\Lambda} \boldsymbol{\Phi}^T|^{1/2} = |\mathbf{I}| = 1.$$

Problem 2.28

- a) We prove that if

$$p(x_i - \mu_i, x_j - \mu_j) = p(x_i - \mu_i) p(x_j - \mu_j),$$

then

$$\sigma_{ij} = \mathbb{E}[(x_i - \mu_i)(x_j - \mu_j)] = 0.$$

With words: we prove that statistical independence implies zero covariance.

$$\begin{aligned}\mathbb{E}[(x_i - \mu_i)(x_j - \mu_j)] &= \\ &= \iint p(x_i - \mu_i, x_j - \mu_j)(x_i - \mu_i)(x_j - \mu_j) dx_j dx_i = \\ &= \iint p(x_i - \mu_i)p(x_j - \mu_j)(x_i - \mu_i)(x_j - \mu_j) dx_j dx_i \\ &= \int p(x_i - \mu_i)(x_i - \mu_i) \left(\int p(x_j - \mu_j)(x_j - \mu_j) dx_j \right) dx_i\end{aligned}$$

If the term in the parenthesis is identically zero, then $\sigma_{ij} = 0$. This is indeed true, since the integral is

$$\int p(x_j - \mu_j)(x_j - \mu_j) dx_j = \mathbb{E}[(x_j - \mu_j)] = \mathbb{E}[x_j] - \mathbb{E}[\mu_j] = \mu_j - \mu_j = 0.$$

- b) We wish to prove the converse of a) in the Gaussian case. To achieve this, we must show that $\sigma_{ij} = 0$ when

$$p(x_i - \mu_i, x_j - \mu_j) = p(x_i - \mu_i)p(x_j - \mu_j).$$

Let's simplify the notation to x and y instead of x_i and x_j . If $\sigma_{xy} = 0$, then the covariance matrix is a diagonal matrix $\mathbf{D} = \text{diag}(\sigma_x^2, \sigma_y^2)$. We write the probability $p(x_i - \mu_i, x_j - \mu_j)$ as $p(x, y)$, where the means μ_x and μ_y are both zero. We write

$$\begin{aligned}p(x, y) &= \frac{1}{(2\pi)^{2/2}\sigma_x\sigma_y} \exp\left(-\frac{1}{2}\mathbf{x}^T\mathbf{D}^{-1}\mathbf{x}\right) = \frac{1}{(2\pi)^{2/2}\sigma_x\sigma_y} \exp\left(-\frac{1}{2}(x^2/\sigma_x^2 + y^2/\sigma_y^2)\right) \\ &= \frac{1}{(2\pi)^{1/2}\sigma_x} \exp\left(-\frac{1}{2}(x^2/\sigma_x^2)\right) \frac{1}{(2\pi)^{1/2}\sigma_y} \exp\left(-\frac{1}{2}(y^2/\sigma_y^2)\right) = p(x)p(y).\end{aligned}$$

This demonstrates that when $\sigma_{xy} = 0$, the covariance matrix is diagonal, and the Gaussian factors into products and we have statistical independence.

- c) This problem asks us to find a counterexample of the above, i.e. an example showing that $\sigma_{xy} \neq 0 \Rightarrow p(x, y) = p(x)p(y)$. The probability density function

$$p(x, y) = K \frac{1}{1 + x^2 + y^2}, \quad K^{-1} = \iint_{\mathbb{R}} \frac{1}{1 + x^2 + y^2} dx dy$$

achieves this. The covariance is zero, since

$$\sigma_{xy} = \mathbb{E}[(x - 0)(y - 0)] = \iint_{\mathbb{R}} \frac{xy}{1 + x^2 + y^2} dx dy = \iint_{\mathbb{R}} I(x, y) dx dy$$

is zero because the integrand $I(x, y)$ is an odd function.

On the other hand, $p(x, y)$ does *not* factor into $p(x)p(y)$. We have proved that $\sigma_{xy} \neq 0 \Rightarrow p(x, y) = p(x)p(y)$ by finding a counterexample.

Problem 2.31

- a) We'll assume that $\mu_1 < \mu_2$. Since $\sigma_1 = \sigma_2 = \sigma$, the minimum probability of error is achieved by setting the decision threshold to $x^* = (\mu_1 + \mu_2)/2$. When following the derivation below, it is illuminating to draw the real line and two Gaussians. The probability of error is

$$\begin{aligned} P_e &= P(x \in R_2, \omega_1) + P(x \in R_1, \omega_2) \\ &= P(x \in R_2 | \omega_1)P(\omega_1) + P(x \in R_1 | \omega_2)P(\omega_2) \\ &= \int_{R_2} p(x | \omega_1)P(\omega_1) dx + \int_{R_1} p(x | \omega_2)P(\omega_2) dx \\ &= \frac{1}{2} \left(\int_{x^*}^{\infty} p(x | \omega_1) dx + \int_0^{x^*} p(x | \omega_2) dx \right) = \int_{x=(\mu_1+\mu_2)/2}^{\infty} p(x | \omega_1) dx \\ &= \int_{x=(\mu_1+\mu_2)/2}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2} \frac{(x - \mu_1)^2}{\sigma^2}\right) dx. \end{aligned}$$

Changing variables to $u = (x - \mu_1)/\sigma$ and using $dx = \sigma du$ yields

$$P_e = \int_{u=a}^{\infty} \frac{1}{\sqrt{2\pi}} \exp(-u^2/\sigma^2) du,$$

where $a = (x - \mu_1)/\sigma = ((\mu_1 + \mu_2)/2 - \mu_1)/\sigma = (\mu_2 - \mu_1)/2\sigma$, as required.

- b) Using the inequality stated in the problem, it remains to show that

$$\lim_{a \rightarrow \infty} f(a) = \lim_{a \rightarrow \infty} \frac{1}{\sqrt{2\pi}a} \exp(-a^2/\sigma^2) = 0.$$

This holds if the derivative is negative as $a \rightarrow \infty$, since then the function decreases as $a \rightarrow \infty$. The derivative of $f(a)$ is

$$f'(x) = -\exp(-a^2/2) \left(1 - \frac{1}{a^2}\right),$$

which is negative as long as $|a| \geq 1$. Alternatively, we can simply note that both factors in $f(a)$ go to zero as $a \rightarrow \infty$.

Problem 2.43

- a) p_{ij} is the probability that the i 'th entry in the vector \mathbf{x} equals 1, given a state of nature ω_j .
 b) We decide ω_j if $P(\omega_j | \mathbf{x})$ is greater than $P(\omega_k | \mathbf{x})$ for every $k \neq j$.

$$P(\omega_j | \mathbf{x}) \propto p(\mathbf{x} | \omega_j)P(\omega_j)$$

We use the equation $p(\mathbf{x} | \omega_j) = \prod_{i=1}^d p(x_i | \omega_j)$, which follows from the fact that the entries are statistically independent. Furthermore, we see that

$$p(x_i | \omega_j) = \begin{cases} p_{ij} & \text{if } x_i = 1 \\ 1 - p_{ij} & \text{if } x_i = 0 \end{cases} = p_{ij}^{x_i} (1 - p_{ij})^{1-x_i}.$$

Now we take logarithms and obtain

$$\begin{aligned}\ln \left(\prod_{i=1}^d p(x_i | \omega_j) P(\omega_j) \right) &= \sum_{i=1}^d \ln p(x_i | \omega_j) + \ln P(\omega_j) \\ &= \sum_{i=1}^d \ln p_{ij}^{x_i} (1 - p_{ij})^{1-x_i} + \ln P(\omega_j) \\ &= \sum_{i=1}^d x_i \ln p_{ij} + (1 - x_i) \ln(1 - p_{ij}) + \ln P(\omega_j),\end{aligned}$$

which is easily arranged to correspond with the expression in the problem statement. In summary we choose the class ω_j if the probability of that class given the data point exceeds the probability of every other class.

2.3 Maximum-likelihood and Bayesian parameter estimation

Problem 3.2

- a) The maximum likelihood estimate for θ is $\max_{\theta} p(\mathcal{D} | \theta) = \max_{\theta} \prod_{i=1}^n p(x_i | \theta)$. The probability of a single sample $p(x_i | \theta)$ is given by the expression

$$p(x_i | \theta) = \begin{cases} 1/\theta & \text{if } 0 \leq x_i \leq \theta \\ 0 & \text{if } x_i > \theta. \end{cases}$$

Clearly the product $\prod_{i=1}^n p(x_i | \theta)$ is zero if any x_i is larger than θ . Therefore θ must be larger than, or equal to, $\max_k x_k$ for the likelihood to be non-zero.

On the other hand, the product equals $1/\theta^n$, and taking logarithms we obtain $-n \ln \theta$. This function is maximized when θ is as small as possible.

The conclusion is that θ must be greater than or equal to $\max_k x_k$ to avoid the likelihood being zero, and also as small as possible to maximize the likelihood. Therefore the maximum likelihood is given by $\hat{\theta} = \max_k x_k = \max \mathcal{D}$.

- b) Skipping this plot. The explanation of why the other data points are not needed is given in part a) of the problem.

Problem 3.4

The maximum likelihood estimate is

$$p(\mathcal{D} | \boldsymbol{\theta}) = \prod_{k=1}^n p(\mathbf{x} | \boldsymbol{\theta}) = \prod_{k=1}^n \prod_{i=1}^d \theta_i^{x_{ik}} (1 - \theta_i)^{(1-x_{ik})}.$$

The log likelihood $\ell(\boldsymbol{\theta})$ is $\ln p(\mathcal{D} | \boldsymbol{\theta})$, which is given by the expression

$$\ell(\boldsymbol{\theta}) = \sum_{k=1}^n \sum_{i=1}^d x_{ik} \ln \theta_i + (1 - x_{ik}) \ln (1 - \theta_i).$$

Differentiating $\ell(\boldsymbol{\theta})$ with respect to θ_i , every term in the sum $\sum_{i=1}^d$ vanishes except the i 'th. We perform the differentiation and equate the result to zero, which yields

$$\frac{d\ell(\boldsymbol{\theta})}{d\theta_i} = \sum_{k=1}^n \left[\frac{x_{ik}}{\theta_i} + \frac{x_{ik} - 1}{1 - \theta_i} \right] = \sum_{k=1}^n [x_{ik} - \theta_i] = 0.$$

Solving this for θ_i gives us $\theta_i = n^{-1} \sum_{k=1}^n x_{ik}$, or in vector notation, $\boldsymbol{\theta} = n^{-1} \sum_{k=1}^n \mathbf{x}_k$. This is what the problem asked us to show.

Problem 3.13

- a) Familiarity with summation notation helps when solving this problem. The matrix-vector product $\mathbf{A}\mathbf{a}$ may be written as $\sum_j A_{ij}a_j$. The sum is typically taken over repeated indices, but we will explicitly typeset the summation index.

Let's write the outer product as $\mathbf{a}\mathbf{b}^T = a_i \oplus b_j$, the trace as $\text{tr}(\mathbf{A}) = \sum_i A_{ii}$ and

$$\text{tr}(\mathbf{a}\mathbf{b}^T) = \sum_{i=j} a_i \oplus b_j = \sum_i a_i b_i.$$

We note that the effect of $\sum_{i=j}$ on a summand is to replace i by j , or vice versa.

In summation notation, $\mathbf{a}^T \mathbf{A} \mathbf{a} = \sum_i \sum_j A_{ij} a_j a_i$. Using the definitions above, and recalling that $\mathbf{A}\mathbf{a}$ is just a vector with value $\sum_j A_{ij} a_j$ in the i 'th index, we see that

$$\text{tr}(\mathbf{A}\mathbf{a}\mathbf{a}^T) = \sum_{i=k} \left(\sum_j A_{ij} a_j \right) \oplus a_k = \sum_i \sum_j A_{ij} a_j a_i.$$

- b) The likelihood is given by the expression

$$\begin{aligned} p(\mathcal{D} | \theta) &= \prod_{k=1}^n \frac{|\Sigma^{-1}|^{1/2}}{(2\pi)^{d/2}} \exp\left(-\frac{1}{2} (\mathbf{x}_k - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}_k - \boldsymbol{\mu})\right) \\ &= \frac{|\Sigma^{-1}|^{n/2}}{(2\pi)^{nd/2}} \exp\left(-\frac{1}{2} \sum_{k=1}^n (\mathbf{x}_k - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}_k - \boldsymbol{\mu})\right). \end{aligned}$$

Applying $\mathbf{a}^T \mathbf{A} \mathbf{a} = \text{tr}(\mathbf{A}\mathbf{a}\mathbf{a}^T)$ from problem a) yields

$$p(\mathcal{D} | \theta) = \frac{|\Sigma^{-1}|^{n/2}}{(2\pi)^{nd/2}} \exp\left(-\frac{1}{2} \sum_{k=1}^n \text{tr}\left(\Sigma^{-1} (\mathbf{x}_k - \boldsymbol{\mu}) (\mathbf{x}_k - \boldsymbol{\mu})^T\right)\right),$$

and by using $\text{tr}(\mathbf{A} + \mathbf{B}) = \text{tr}(\mathbf{A}) + \text{tr}(\mathbf{B})$ on the exponent we complete the problem.

- c) To solve this problem, we make use of two proofs from linear algebra.
- (i) The determinant is the product of the eigenvalues, i.e. $\det \mathbf{A} = \prod_{i=1}^d \lambda_i$. This can be proved by taking the determinant of the eigenvalue decomposition of \mathbf{A} , i.e. $\mathbf{A} = \mathbf{Q}\boldsymbol{\Lambda}\mathbf{Q}^T$, since the determinant of an orthogonal matrix is unity and $\boldsymbol{\Lambda}$ has eigenvalues on the diagonal, the result is that $\det \mathbf{A} = |\mathbf{A}| = \prod_{i=1}^d \lambda_i$.
 - (ii) The trace is the sum of the eigenvalues, i.e. $\text{tr} \mathbf{A} = \sum_{i=1}^d \lambda_i$. The proof for this involves characteristic polynomials, but is not sketched here.

The term involving Σ^{-1} is transformed in the following way: if $\mathbf{A} = \Sigma^{-1} \hat{\Sigma}$, then $|\mathbf{A}| = |\Sigma^{-1}| |\hat{\Sigma}|$. We then write the determinant of the inverse matrix as

$$|\Sigma^{-1}| = \frac{|\mathbf{A}|}{|\hat{\Sigma}|} = \frac{\prod_{i=1}^d \lambda_i}{|\hat{\Sigma}|}.$$

To transform the exponent, we write

$$\text{tr} \left(\sum_{k=1}^n \boldsymbol{\Sigma}^{-1} (\mathbf{x}_k - \boldsymbol{\mu}) (\mathbf{x}_k - \boldsymbol{\mu})^T \right) = \text{tr} \left(\boldsymbol{\Sigma}^{-1} n \hat{\boldsymbol{\Sigma}} \right) = n \text{tr} (\mathbf{A}) = n \sum_{i=1}^d \lambda_i.$$

Substituting these transformation into the previous equation yields the required expression, which is

$$p(\mathcal{D} | \theta) = \frac{\left(\prod_{i=1}^d \lambda_i \right)^{n/2}}{(2\pi)^{nd/2} |\hat{\boldsymbol{\Sigma}}|^{n/2}} \exp \left(-\frac{n}{2} \sum_{i=1}^d \lambda_i \right). \quad (1)$$

d) Taking logarithms of Equation (1) above gives us the log-likelihood function

$$\ln p(\mathcal{D} | \theta) = \frac{n}{2} \left[\sum_{i=1}^d \ln \lambda_i - d \ln 2\pi - \ln |\hat{\boldsymbol{\Sigma}}| \right] - \frac{n}{2} \sum_{i=1}^d \lambda_i,$$

and differentiating it with respect to the eigenvalue λ_j yields

$$\frac{\partial \ln p(\mathcal{D} | \theta)}{\partial \lambda_j} = \frac{n}{2} (1 - \lambda_j^2) = 0.$$

The Hessian matrix (second derivative) is a diagonal matrix with $-2\lambda_j$ on the (j, j) 'th entry, so it is negative definite (i.e. $\mathbf{x}^T \mathbf{H} \mathbf{x} < 0$ for every $\mathbf{x} \neq 0$) when all λ_j 's are positive. When the Hessian is negative definite, then the solution is a maximum. Therefore we take $\lambda_j = 1$ as the solution for every j (and not $\lambda_j = -1$, which is not a maximum point).

If every eigenvalues of \mathbf{A} is 1, then $\mathbf{A} = \mathbf{I}$. Recall that $\mathbf{A} = \boldsymbol{\Sigma}^{-1} \hat{\boldsymbol{\Sigma}}$, so when we substitute the identity for \mathbf{A} we obtain $\boldsymbol{\Sigma}^{-1} \hat{\boldsymbol{\Sigma}} = \mathbf{I}$. The likelihood is maximized when we take $\hat{\boldsymbol{\Sigma}}$ to be $n^{-1} \sum_{k=1}^n (\mathbf{x}_k - \boldsymbol{\mu}) (\mathbf{x}_k - \boldsymbol{\mu})^T$, and the proof is complete.

Problem 3.15

a) Starting with Equation (31) from the book, we get rid of σ_n^2 by substituting the expression given in Equation (32) to obtain

$$\mu_n = \underbrace{\left[\frac{1}{\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}} \right]}_{\sigma_n^2} \frac{n}{\sigma^2} \hat{\mu}_n + \underbrace{\left[\frac{1}{\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}} \right]}_{\sigma_n^2} \frac{\mu_0}{\sigma_0^2}.$$

We cancel terms, use the dogmatism $n_0 = \sigma^2/\sigma_0^2$, and realize that $\mu_0 = n_0^{-1} \sum_{k=-n_0+1}^0 x_k$.

$$\begin{aligned}\mu_n &= \left[\frac{1}{1 + \frac{\sigma^2}{n\sigma_0^2}} \right] \hat{\mu}_n + \left[\frac{1}{1 + \frac{n\sigma_0^2}{\sigma^2}} \right] \mu_0 \\ &= \left[\frac{1}{1 + \frac{n_0}{n}} \right] \frac{1}{n} \sum_{k=1}^n x_k + \left[\frac{1}{1 + \frac{n}{n_0}} \right] \frac{1}{n_0} \sum_{k=-n_0+1}^0 x_k \\ &= \frac{1}{n + n_0} \sum_{k=1}^n x_k + \frac{1}{n + n_0} \sum_{k=-n_0+1}^0 x_k = \frac{1}{n + n_0} \sum_{k=-n_0+1}^n x_k.\end{aligned}$$

- (a) One interpretation of μ_n is that it's a weighted average of the real samples and the fictitious samples, since

$$\mu_n = \frac{n}{n + n_0} \left[\frac{1}{n} \sum_{k=0}^n n \right] + \frac{n_0}{n + n_0} \left[\frac{1}{n_0} \sum_{k=-n_0+1}^0 x_k \right] = \frac{1}{n + n_0} \sum_{k=-n_0+1}^n x_k.$$

An interpretation of σ_n^2 is more straightforward if we consider it's reciprocal, the *precision* σ_n^{-2} . The precision can be written as $\sigma_n^{-2} = n\sigma^{-2} + n_0\sigma_0^{-2}$, so the number of fictitious samples provides an initial estimate for the precision, and as more real data points are added the precision is increased linearly.

Problem 3.16

- a) The trick is to multiply both terms by \mathbf{I} at an opportune moment, expanding the identity and using $(\mathbf{AB})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}$. We prove the theorem by writing

$$\begin{aligned}(\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1} &= (\mathbf{A}^{-1}\mathbf{I} + \mathbf{IB}^{-1})^{-1} = (\mathbf{A}^{-1}\mathbf{BB}^{-1} + \mathbf{A}^{-1}\mathbf{AB}^{-1})^{-1} \\ &= (\mathbf{A}^{-1}(\mathbf{B} + \mathbf{A})\mathbf{B}^{-1})^{-1} = \mathbf{B}(\mathbf{A} + \mathbf{B})^{-1}\mathbf{A}.\end{aligned}$$

Replacing the second equality by $\mathbf{IA}^{-1} + \mathbf{B}^{-1}\mathbf{I} = \mathbf{B}^{-1}\mathbf{BA}^{-1} + \mathbf{B}^{-1}\mathbf{AA}^{-1}$ would yield the second quality in the theorem as stated in the problem.

- b) The matrices must be square, since we use $\mathbf{I} = \mathbf{A}^{-1}\mathbf{A}$ to prove the first equality, and $\mathbf{I} = \mathbf{A}^{-1}\mathbf{A}$ to prove the second. The same logic applies to \mathbf{B} . In other words, we require that \mathbf{A} has a left-inverse and a right-inverse, and it must therefore be square.

An alternative approach is to consider the second equality in the theorem directly. Some algebra reveals that this equality requires \mathbf{A} to have a left and right inverse. It must therefore be square, against since no non-square matrix has a left-inverse and a right-inverse.

- c) The starting point are the equations

$$\Sigma_n^{-1} = n\Sigma^{-1} + \Sigma_0^{-1} \quad \text{and} \quad \Sigma_n^{-1}\boldsymbol{\mu}_n = n\Sigma^{-1}\hat{\boldsymbol{\mu}}_n + \Sigma_0^{-1}\boldsymbol{\mu}_0,$$

and we wish to solve these equations with respect to Σ_n and $\boldsymbol{\mu}_n$. In other words, we want the functional dependence to be $\Sigma_n = f(\hat{\boldsymbol{\mu}}_n, \boldsymbol{\mu}_0, \Sigma, \Sigma_0)$ and $\boldsymbol{\mu}_n = (\hat{\boldsymbol{\mu}}_n, \boldsymbol{\mu}_0, \Sigma, \Sigma_0)$.

We start with the covariance matrix. To solve for Σ_n , we write

$$\Sigma_n = (\Sigma_n^{-1})^{-1} = (n\Sigma^{-1} + \Sigma_0^{-1})^{-1} = \Sigma_0 \left(\Sigma_0 + \frac{1}{n}\Sigma \right)^{-1} \frac{1}{n}\Sigma,$$

where the last equality comes from using

$$(\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1} = \mathbf{B}^{-1} (\mathbf{B} + \mathbf{A})^{-1} \mathbf{A}^{-1}.$$

To solve for the mean μ_n , we write

$$\begin{aligned} \mu_n &= \Sigma_n n \Sigma^{-1} \hat{\mu}_n + \Sigma_n \Sigma_0^{-1} \mu_0 \\ &= \left[\Sigma_0 \left(\Sigma_0 + \frac{1}{n}\Sigma \right)^{-1} \frac{1}{n}\Sigma \right] n \Sigma^{-1} \hat{\mu}_n + \left[\frac{1}{n}\Sigma \left(\frac{1}{n}\Sigma + \Sigma_0 \right)^{-1} \Sigma_0 \right] \Sigma_0^{-1} \mu_0 \\ &= \Sigma_0 \left(\Sigma_0 + \frac{1}{n}\Sigma \right)^{-1} \hat{\mu}_n + \frac{1}{n}\Sigma \left(\Sigma_0 + \frac{1}{n}\Sigma \right)^{-1} \mu_0 \end{aligned}$$

Where we made use of both equalities given by the theorem in subproblem a).

Problem 3.24

Our task is to find the maximum likelihood estimate of θ in the Rayleigh distribution, which is given by the expression

$$p(x | \theta) = \begin{cases} 2\theta x \exp(-\theta x^2) & \text{if } x \geq 0 \\ 0 & \text{otherwise.} \end{cases}$$

The log-likelihood function is given by

$$\ell(\theta) = \ln(p(\mathcal{D} | \theta)) = \sum_{k=1}^n \ln(2\theta x_k \exp(-\theta x_k^2)) = \sum_{k=1}^n -\theta x_k^2 + \ln 2\theta x_k.$$

Differentiating the log-likelihood with respect to θ and equating it to zero yields

$$\ell'(\theta) = \sum_{k=1}^n \left(-x_k^2 + \frac{2x_k}{2\theta x_k} \right) = \sum_{k=1}^n \left(-x_k^2 + \frac{1}{\theta} \right) = \frac{n}{\theta} - \sum_{k=1}^n x_k^2 = 0,$$

and solving this equation for θ reveals the desired answer.

Problem 3.32

a) A simple $\Theta(n^2)$ algorithm is given by

```
f = 0
for i=0 to n-1 do:
    f = f + a[i] * x**i
end
```

Assuming that the computation of x^i requires $i - 1$ flops, the complexity for iteration i is $\Theta(i)$, and the full complexity becomes $1 + 2 + 3 + \dots + (n - 1) = \Theta(n^2)$.

- b) The waste in the algorithm from subproblem a) above stems from having to compute powers of x many times. If we know x^{k-1} , then $x^k = x^{k-1}x$, and there is no need to compute x^k as $\underbrace{x \cdot x \cdot \dots \cdot x}_{k \text{ times}}$. Define $S_k := \sum_{i=k}^{n-1} a_i x^{i-k}$, then

$$f(x) = S_0 = a_0 + xS_1 = a_0 + x(a_1 + xS_2) = a_0 + x(a_1 + x(a_2 + xS_3)).$$

Clearly $S_k = a_k + xS_{k+1}$, looping backwards we can use the following algorithm

```
f = a[n-1]
for i=n-2 to 0 do:
    f = a[i] + x * f
end
```

which requires a constant number flops in each iteration. The computational complexity is therefore $\Theta(n)$. A specific example when $n = 4$ is

$$a_0x^0 + a_1x^1 + a_2x^2 + a_3x^3 = a_0 + x(a_1 + x(a_2 + x(a_3 + x))),$$

and the algorithm evaluates this using the right-hand side, from the inside and out.

Problem 3.35

- a) The complexity of computing $\hat{\boldsymbol{\mu}}_n$ is $O(nd)$, where n is the number of data points and d is the dimensionality of the data. Adding n vectors of length d requires $d(n - 1)$ floating point operations (flops), dividing through by n requires another d flops. The resulting complexity is therefore $d(n - 1) + d = O(nd)$.

The complexity of computing \mathbf{C}_n is $O(nd^2)$. For each term in the sum, we compute the subtraction using d flops, followed by the outer product using d^2 flops. This is done of each of the n terms, so we need $n(d + d^2)$ flops to compute the terms. We then add the terms, requiring $d^2(n - 1)$ flops. Dividing requires d^2 flops. In total, the cost is $n(d + d^2) + d^2(n - 1) + d^2 \simeq 2nd^2 = O(nd^2)$ flops.

- b) We will show that these equations are correct by two different methods.

The recursive definition of $\hat{\boldsymbol{\mu}}_n$ may be proved by induction. It's clearly valid for $n = 0$, since then it reduces to $\hat{\boldsymbol{\mu}}_0 = \mathbf{x}_1$. The inductive step is

$$\begin{aligned} \hat{\boldsymbol{\mu}}_{n+1} &= \hat{\boldsymbol{\mu}}_n + \frac{1}{n+1}(\mathbf{x}_{n+1} - \hat{\boldsymbol{\mu}}_n) = \left(1 - \frac{1}{n+1}\right)\hat{\boldsymbol{\mu}}_n + \frac{1}{n+1}\mathbf{x}_{n+1} \\ &= \left(\frac{n}{n+1}\right)\hat{\boldsymbol{\mu}}_n + \frac{1}{n+1}\mathbf{x}_{n+1} = \frac{1}{n+1}\sum_{k=1}^n \mathbf{x}_k + \frac{1}{n+1}\mathbf{x}_{n+1} = \frac{1}{n+1}\sum_{k=1}^{n+1} \mathbf{x}_k, \end{aligned}$$

and this proves the recursive relation for $\hat{\boldsymbol{\mu}}_{n+1}$.

We will prove that the recursive equation for \mathbf{C}_{n+1} is correct by deriving it from the definition. This is a somewhat tedious computation, and the strategy involves

writing \mathbf{C}_{n+1} as a function of known quantities, i.e. \mathbf{C}_n , $\hat{\boldsymbol{\mu}}_n$, $\hat{\boldsymbol{\mu}}_{n+1}$ and \mathbf{x}_{n+1} . We begin by writing the definition of \mathbf{C}_{n+1} , which is

$$\mathbf{C}_{n+1} = \frac{1}{n} \sum_{k=1}^{n+1} (\mathbf{x}_k - \hat{\boldsymbol{\mu}}_{n+1}) (\mathbf{x}_k - \hat{\boldsymbol{\mu}}_{n+1})^T.$$

We then subtract and add $\hat{\boldsymbol{\mu}}_n$ to write the outer product as

$$((\mathbf{x}_k - \hat{\boldsymbol{\mu}}_n) - (\hat{\boldsymbol{\mu}}_{n+1} - \hat{\boldsymbol{\mu}}_n)) ((\mathbf{x}_k - \hat{\boldsymbol{\mu}}_n) - (\hat{\boldsymbol{\mu}}_{n+1} - \hat{\boldsymbol{\mu}}_n))^T. \quad (2)$$

The last term in each factor above may be written as

$$\hat{\boldsymbol{\mu}}_{n+1} - \hat{\boldsymbol{\mu}}_n = \frac{1}{n+1} (\mathbf{x}_{n+1} - \hat{\boldsymbol{\mu}}_n).$$

Equation (2) has the same functional form as

$$(\mathbf{a} - \mathbf{b})(\mathbf{a} - \mathbf{b})^T = \mathbf{a}\mathbf{a}^T - \mathbf{a}\mathbf{b}^T - \mathbf{b}\mathbf{a}^T + \mathbf{b}\mathbf{b}^T,$$

and we expand it as such to obtain the sum

$$\begin{aligned} \mathbf{C}_{n+1} &= \frac{1}{n} \sum_{k=1}^{n+1} (\mathbf{x}_k - \hat{\boldsymbol{\mu}}_n) (\mathbf{x}_k - \hat{\boldsymbol{\mu}}_n)^T + (\mathbf{x}_k - \hat{\boldsymbol{\mu}}_n) \frac{1}{n+1} (\mathbf{x}_{n+1} - \hat{\boldsymbol{\mu}}_n)^T \\ &\quad + \frac{1}{n+1} (\mathbf{x}_{n+1} - \hat{\boldsymbol{\mu}}_n) (\mathbf{x}_k - \hat{\boldsymbol{\mu}}_n) + \frac{1}{n+1} (\mathbf{x}_{n+1} - \hat{\boldsymbol{\mu}}_n) \frac{1}{n+1} (\mathbf{x}_{n+1} - \hat{\boldsymbol{\mu}}_n)^T. \end{aligned} \quad (3)$$

We will study each of the four preceding terms in order.

The **first term** in Equation (3) may be written as

$$\frac{1}{n} \sum_{k=1}^{n+1} (\mathbf{x}_k - \hat{\boldsymbol{\mu}}_n) (\mathbf{x}_k - \hat{\boldsymbol{\mu}}_n)^T = \frac{n-1}{n} \mathbf{C}_n + \frac{1}{n} (\mathbf{x}_{n+1} - \hat{\boldsymbol{\mu}}_n) (\mathbf{x}_{n+1} - \hat{\boldsymbol{\mu}}_n)^T,$$

where we stripped out the last term in the sum and used the definition of \mathbf{C}_n .

The **second term** in Equation (3) may be written as

$$\frac{1}{n(n+1)} \left[\sum_{k=1}^{n+1} (\mathbf{x}_k - \hat{\boldsymbol{\mu}}_n) \right] (\mathbf{x}_{n+1} - \hat{\boldsymbol{\mu}}_n)^T = \frac{1}{n(n+1)} (\mathbf{x}_{n+1} - \hat{\boldsymbol{\mu}}_n) (\mathbf{x}_{n+1} - \hat{\boldsymbol{\mu}}_n)^T,$$

since $\sum_{k=1}^{n+1} (\mathbf{x}_k - \hat{\boldsymbol{\mu}}_n) = (\sum_{k=1}^n \mathbf{x}_k) + \mathbf{x}_{n+1} - (n+1)\hat{\boldsymbol{\mu}}_n = n\hat{\boldsymbol{\mu}}_n + \mathbf{x}_{n+1} - n\hat{\boldsymbol{\mu}}_n - \hat{\boldsymbol{\mu}}_n$.

The **third term** in Equation (3) may also be written as

$$\frac{1}{n(n+1)} (\mathbf{x}_{n+1} - \hat{\boldsymbol{\mu}}_n) (\mathbf{x}_{n+1} - \hat{\boldsymbol{\mu}}_n)^T,$$

where we used the same logic as in the computations related to the second term.

The **fourth term** in Equation (3) may be written as

$$\frac{1}{n(n+1)} (\mathbf{x}_{n+1} - \hat{\boldsymbol{\mu}}_n) (\mathbf{x}_{n+1} - \hat{\boldsymbol{\mu}}_n)^T,$$

since the entire term is constant with respect to the index k . We multiplied the last term with $n+1$ to get it out of the sum, canceled part of the fraction, and applied the n^{-1} fraction from outside the sum.

Let's return Equation (3) again, and write $\mathbf{w}_k = (\mathbf{x}_{n+1} - \hat{\boldsymbol{\mu}}_n)$ to ease notation. Using our findings from above, the sum becomes

$$\begin{aligned} \mathbf{C}_{n+1} &= \frac{n-1}{n} \mathbf{C}_n + \frac{1}{n} \mathbf{w} \mathbf{w}^T - \frac{1}{n(n+1)} \mathbf{w} \mathbf{w}^T \\ &\quad - \frac{1}{n(n+1)} \mathbf{w} \mathbf{w}^T + \frac{1}{n(n+1)} \mathbf{w} \mathbf{w}^T, \end{aligned}$$

and since $1/n - 1/(n(n+1)) = 1/(n+1)$ we obtain the desired result

$$\mathbf{C}_{n+1} = \frac{n-1}{n} \mathbf{C}_n + \frac{1}{n+1} \mathbf{w} \mathbf{w}^T.$$

This concludes our derivation of the recursive formulas.

- c) The calculation of $\hat{\boldsymbol{\mu}}_{n+1}$ is dominated by vector addition and subtraction, and the complexity is $O(d)$. The calculation of \mathbf{C}_{n+1} is dominated by an outer product, which has $O(d^2)$ complexity, and a matrix addition, which is $O(d^2)$ too. The overall complexity is for the iterative sample covariance matrix \mathbf{C}_{n+1} is therefore $O(d^2)$.
- d) When data arrives in a stream (on-line learning), the recursive methods are clearly superior to application of the naive formulas. Consider data iteratively entering a learning system. If we compute $\hat{\boldsymbol{\mu}}_1, \hat{\boldsymbol{\mu}}_2, \dots, \hat{\boldsymbol{\mu}}_n$ naively, the overall complexity is

$$1d + 2d + \dots + nd = (1 + 2 + \dots + n) d = O(n^2 d).$$

If we compute the sequence using the recursive equation, the complexity resolves to

$$1d + 1d + \dots + 1d = (1 + 1 + \dots + 1) d = O(nd).$$

Using the same logic, naive computation of $\mathbf{C}_1, \mathbf{C}_2, \dots, \mathbf{C}_n$ will have complexity $O(n^2 d^2)$, while using the recursive formula result in $O(nd^2)$ complexity.

If data arrives sequentially and predictions are to be made before all of the data has arrived, then the recursive formulas are superior. They also require less intermediate storage. It would be interesting to examine the numerical stability of both approaches.

Problem 3.36

- a) We'll prove the Sherman-Morrison-Woodbury matrix identity by demonstrating that it reduces to $I = I$. Recall that $1 + \mathbf{y}^T \mathbf{A}^{-1} \mathbf{x}$ is a scalar. The algebra is

$$\begin{aligned} (A + \mathbf{x}\mathbf{y}^T)^{-1} &= \mathbf{A}^{-1} - \frac{\mathbf{A}^{-1}\mathbf{x}\mathbf{y}^T\mathbf{A}^{-1}}{1 + \mathbf{y}^T\mathbf{A}^{-1}\mathbf{x}} \\ (A + \mathbf{x}\mathbf{y}^T)(A + \mathbf{x}\mathbf{y}^T)^{-1} &= (A + \mathbf{x}\mathbf{y}^T)\mathbf{A}^{-1} - (A + \mathbf{x}\mathbf{y}^T)\frac{\mathbf{A}^{-1}\mathbf{x}\mathbf{y}^T\mathbf{A}^{-1}}{1 + \mathbf{y}^T\mathbf{A}^{-1}\mathbf{x}} \\ I &= I + \mathbf{x}\mathbf{y}^T\mathbf{A}^{-1} - \frac{\mathbf{x}\mathbf{y}^T\mathbf{A}^{-1} + \mathbf{x}\mathbf{y}^T\mathbf{A}^{-1}\mathbf{x}\mathbf{y}^T\mathbf{A}^{-1}}{1 + \mathbf{y}^T\mathbf{A}^{-1}\mathbf{x}} \\ I &= I + \mathbf{x}\mathbf{y}^T\mathbf{A}^{-1} - \frac{\mathbf{x}(1 + \mathbf{y}^T\mathbf{A}^{-1}\mathbf{x})\mathbf{y}^T\mathbf{A}^{-1}}{1 + \mathbf{y}^T\mathbf{A}^{-1}\mathbf{x}} \\ I &= I + \mathbf{x}\mathbf{y}^T\mathbf{A}^{-1} - \mathbf{x}\mathbf{y}^T\mathbf{A}^{-1}. \end{aligned}$$

- b) Recall the recursive equation for the sample covariance matrix, where we define a , b and \mathbf{w} to ease notation in the following derivation.

$$\mathbf{C}_{n+1} = \underbrace{\frac{n-1}{n}}_a \mathbf{C}_n + \underbrace{\frac{1}{n+1}}_b \underbrace{(\mathbf{x}_{n+1} - \hat{\boldsymbol{\mu}}_n)}_{\mathbf{w}} \underbrace{(\mathbf{x}_{n+1} - \hat{\boldsymbol{\mu}}_n)^T}_{\mathbf{w}^T} = a\mathbf{C}_n + b\mathbf{w}\mathbf{w}^T$$

The inverse is given by

$$\mathbf{C}_{n+1}^{-1} = (a\mathbf{C}_n + b\mathbf{w}\mathbf{w}^T)^{-1},$$

and we now apply the Sherman-Morrison-Woodbury matrix identity to obtain

$$\begin{aligned} \mathbf{C}_{n+1}^{-1} &= (a\mathbf{C}_n + b\mathbf{w}\mathbf{w}^T)^{-1} = (a\mathbf{C}_n)^{-1} - \frac{(a\mathbf{C}_n)^{-1}b\mathbf{w}\mathbf{w}^T(a\mathbf{C}_n)^{-1}}{1 + \mathbf{w}^T(a\mathbf{C}_n)^{-1}b\mathbf{w}} \\ &= \frac{1}{a} \left(\mathbf{C}_n^{-1} - \frac{b}{a} \left(\frac{\mathbf{C}_n^{-1}\mathbf{w}\mathbf{w}^T\mathbf{C}_n^{-1}}{1 + \frac{b}{a}\mathbf{w}^T\mathbf{C}_n^{-1}\mathbf{w}} \right) \right) \\ &= \frac{1}{a} \left(\mathbf{C}_n^{-1} - \frac{\mathbf{C}_n^{-1}\mathbf{w}\mathbf{w}^T\mathbf{C}_n^{-1}}{\frac{a}{b} + \mathbf{w}^T\mathbf{C}_n^{-1}\mathbf{w}} \right). \end{aligned}$$

Substituting back the relations $a^{-1} = n/(n-1)$ and $a/b = (n^2-1)/n$, we have shown that the inverse of \mathbf{C}_{n+1} is indeed given by

$$\mathbf{C}_{n+1}^{-1} = \frac{n}{n-1} \left(\mathbf{C}_n^{-1} - \frac{\mathbf{C}_n^{-1}\mathbf{w}\mathbf{w}^T\mathbf{C}_n^{-1}}{\frac{n^2-1}{n} + \mathbf{w}^T\mathbf{C}_n^{-1}\mathbf{w}} \right).$$

- c) The computational complexity is $O(d^2)$. First \mathbf{w} is computed using $O(d)$ flops. In the numerator, the matrix-vector product $\mathbf{x} = \mathbf{C}_n^{-1}\mathbf{w}$ is computed in $O(d^2)$ time, and so is $\mathbf{y}^T = \mathbf{w}^T\mathbf{C}_n^{-1}$. Then the outer product $\mathbf{x}\mathbf{y}^T = (\mathbf{C}_n^{-1}\mathbf{w})(\mathbf{w}^T\mathbf{C}_n^{-1})$ may be computed in $O(d^2)$ time too.

The denominator is also computed in $O(d^2)$ time, and so is the matrix subtraction. Therefore the overall computational complexity is $O(d^2)$. Notice that if $\mathbf{w}\mathbf{w}^T$ is

computed first in the denominator, the computational complexity would be $O(d^3)$, since matrix-matrix products are $O(d^3)$, and must be avoided if possible.

d) See the answer to Problem 35 d) given above.

Problem 3.38

a) The criterion function is given by

$$J_1(\mathbf{w}) = \frac{(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^2}{\sigma_1^2 + \sigma_2^2} = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T (\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2) \mathbf{w}},$$

since the numerator is given by

$$\begin{aligned} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^2 &= (\mathbf{w}^T \boldsymbol{\mu}_1 - \mathbf{w}^T \boldsymbol{\mu}_2) (\mathbf{w}^T \boldsymbol{\mu}_1 - \mathbf{w}^T \boldsymbol{\mu}_2)^T \\ &= \mathbf{w}^T (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \mathbf{w} = \mathbf{w}^T \mathbf{S}_B \mathbf{w} \end{aligned}$$

and the denominator is given by

$$\sigma_1^2 + \sigma_2^2 = \mathbf{w}^T (\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2) \mathbf{w}.$$

Now we compare this with Equation (96) in Chapter 3 in the book. By the same logic used there, the solution is given by Equation (106) from the book, which is

$$\mathbf{w} = (\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2)^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) = (\boldsymbol{\Sigma}_1^{-1} + \boldsymbol{\Sigma}_2^{-1}) (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2).$$

b) Simply substitute $\sigma_i^2 \rightarrow P(\omega_i)\sigma_i^2$ into problem a).

c) Recall that $\tilde{s}_i^2 = \sum_{y \in \mathcal{Y}_i} (y_i - \tilde{m}_i)^2$. This expression is n_i times the population variance, so $\tilde{s}_i^2 \approx n_i \sigma_i^2$. Equation (96) in the book therefore has $n_1 \sigma_1^2 + n_2 \sigma_2^2$ in the denominator, $J_1(\mathbf{w})$ from subproblem a) had $\sigma_1^2 + \sigma_2^2$ in the denominator and the denominator of $J_2(\mathbf{w})$ is

$$P(\omega_1)\sigma_1^2 + P(\omega_2)\sigma_2^2 = \frac{n_1}{n}\sigma_1^2 + \frac{n_2}{n}\sigma_2^2 = \frac{1}{n} (n_1\sigma_1^2 + n_2\sigma_2^2).$$

The denominator in Equation (96) in the book is proportional to $J_2(\mathbf{w})$, so they are the most similar. Their optimization results in the same value of \mathbf{w} , since constants make no difference when optimizing $J(\mathbf{w})$.

Problem 3.39

a) We start by expanding terms

$$\begin{aligned} J_1 &= \frac{1}{n_1 n_2} \sum_i \sum_j (y_i^2 - 2y_i y_j + y_j^2) \\ &= \frac{1}{n_1 n_2} \left(n_2 \sum_i y_i^2 - 2 \left(\sum_i y_i \right) \left(\sum_j y_j \right) + n_1 \sum_j y_j^2 \right). \end{aligned}$$

From $\text{var}(Y) = \mathbb{E}(Y^2) - \mathbb{E}(Y)^2$ we note that

$$\sum_i y_i^2 = \sum_i (y_i - m_i)^2 + \frac{1}{n_i} \left(\sum_i y_i \right)^2 = s_1^2 + \frac{1}{n_i} \left(\sum_i y_i \right)^2.$$

Now we denote $k_i = \sum_i y_i$ and $k_j = \sum_j y_j$, then we substitute the equation above into the equation for J_1 to obtain

$$\begin{aligned} J_1 &= \frac{1}{n_1 n_2} \left[n_2 \left(s_1^2 + \frac{1}{n_1} k_i^2 \right) - 2k_i k_j + n_1 \left(s_2^2 + \frac{1}{n_2} k_j^2 \right) \right] \\ &= \frac{s_1^2}{n_1} + \frac{s_2^2}{n_2} + \frac{k_i^2}{n_1^2} - 2 \frac{k_i k_j}{n_1 n_2} + \frac{k_j^2}{n_2^2} \\ &= \frac{s_1^2}{n_1} + \frac{s_2^2}{n_2} + m_1^2 - 2m_1 m_2 + m_2^2 = \frac{s_1^2}{n_1} + \frac{s_2^2}{n_2} + (m_1 - m_2)^2. \end{aligned}$$

- b) Not sure about this one. Skipped.
c) Not sure about this one. Skipped.

Problem 3.40

- a) We have $\tilde{\mathbf{S}}_W = \mathbf{W}^T \mathbf{S}_W \mathbf{W}$, using the fact that \mathbf{W} contains eigenvector columns with $\mathbf{e}_i^T \mathbf{S}_W \mathbf{e}_j^T = \delta_{ij}$, we observe that

$$\begin{aligned} \mathbf{W}^T \mathbf{S}_W \mathbf{W} &= \begin{pmatrix} \mathbf{e}_1^T \\ \vdots \\ \mathbf{e}_n^T \end{pmatrix} (\mathbf{S}_W \mathbf{e}_1 \quad \dots \quad \mathbf{S}_W \mathbf{e}_n) = \begin{pmatrix} \mathbf{e}_1^T \\ \vdots \\ \mathbf{e}_n^T \end{pmatrix} \mathbf{S}_W (\mathbf{e}_1 \quad \dots \quad \mathbf{e}_n) \\ &= \begin{pmatrix} \mathbf{e}_1^T \mathbf{S}_W \mathbf{e}_1 & \mathbf{e}_1^T \mathbf{S}_W \mathbf{e}_2 & \dots \\ \mathbf{e}_2^T \mathbf{S}_W \mathbf{e}_1 & \mathbf{e}_2^T \mathbf{S}_W \mathbf{e}_2 & \vdots \\ \vdots & \dots & \mathbf{e}_n^T \mathbf{S}_W \mathbf{e}_n \end{pmatrix} = \delta_{ij} = \mathbf{I}. \end{aligned}$$

The same exact procedure may be applied to $\tilde{\mathbf{S}}_B$ to show that it's a diagonal matrix with eigenvalues, Λ . We will not devote space to this computation.

- b) Applying the result from subproblem a), we see that the value of J is

$$J = \frac{|\tilde{\mathbf{S}}_B|}{|\tilde{\mathbf{S}}_W|} = \frac{|\mathbf{W}^T \mathbf{S}_B \mathbf{W}|}{|\mathbf{W}^T \mathbf{S}_W \mathbf{W}|} = \frac{|\Lambda|}{|\mathbf{I}|} = \prod_{i=1}^n \lambda_i.$$

- c) The relation is $\mathbf{y} = \mathbf{W}^T \mathbf{x}$, and the relation after scaling and rotating becomes $\mathbf{y}' = \mathbf{Q} \mathbf{D} \mathbf{W}^T \mathbf{x}$. We replace \mathbf{W}^T by $\mathbf{Q} \mathbf{D} \mathbf{W}^T$ and obtain

$$J' = \frac{|(\mathbf{Q} \mathbf{D} \mathbf{W}^T) \mathbf{S}_B (\mathbf{Q} \mathbf{D} \mathbf{W}^T)^T|}{|(\mathbf{Q} \mathbf{D} \mathbf{W}^T) \mathbf{S}_W (\mathbf{Q} \mathbf{D} \mathbf{W}^T)^T|} = \frac{|\mathbf{Q}| |\mathbf{D}| |\Lambda| |\mathbf{D}| |\mathbf{Q}^{-1}|}{|\mathbf{Q}| |\mathbf{D}| |\mathbf{I}| |\mathbf{D}| |\mathbf{Q}^{-1}|} = \frac{|\Lambda|}{|\mathbf{I}|} = J,$$

where we used $\mathbf{W}^T \mathbf{S}_W \mathbf{W} = \mathbf{I}$ and $\mathbf{W}^T \mathbf{S}_B \mathbf{W} = \Lambda$, as well as $\mathbf{Q}^T = \mathbf{Q}^{-1}$. Clearly $\mathbf{W}^T \rightarrow \mathbf{Q} \mathbf{D} \mathbf{W}^T$ leaves J unchanged, so it's invariant to this transformation.

2.4 Nonparametric techniques

Problem 4.2

- a) We'll prove $\bar{p}_n(x) \sim \mathcal{N}(\mu, \sigma^2 + h_n^2)$ by direct computation. An alternative approach would be to use the Fourier convolution theorem, which states that convolution becomes pointwise multiplication in the Fourier basis.

We start with

$$\bar{p}_n(x) = \int \frac{1}{h_n} \phi\left(\frac{x-v}{h_n}\right) p(v) dv,$$

and if $\phi(x) \sim \mathcal{N}(0, 1)$ then it is easy to see that $\frac{1}{h_n} \phi\left(\frac{x-v}{h_n}\right) \sim \mathcal{N}(v, h_n^2)$. We'll expand the exponents, write it as a quadratic function of v , integrate out the v -variable and using algebra transform the result into a more well-known form. The integral is

$$\frac{1}{\sqrt{2\pi}\sigma} \frac{1}{\sqrt{2\pi}h_n} \int \exp\left[-\frac{1}{2} \left(\left(\frac{x-v}{h_n}\right)^2 + \left(\frac{v-\mu}{\sigma}\right)^2 \right)\right] dv, \quad (4)$$

and we wish to integrate out the v . To do so, we write the exponent as a function of v . Defining $\beta = x^2\sigma^2 + h_n^2\mu^2$, the exponent may be written as

$$\frac{\sigma^2(x^2 - 2xv + v^2) + h_n^2(v^2 - 2v\mu + \mu^2)}{(h_n\sigma)^2} = \frac{v^2(\sigma^2 + h_n^2) - 2v(x\sigma^2 + \mu h_n^2) + \beta}{(h_n\sigma)^2}.$$

Now we introduce $y = v\sqrt{\sigma^2 + h_n^2}$, since we are aiming to write the exponent as $(y-a)^2/b^2 + c$ for some y . Defining α and completing the square, the right hand side of the equation above may be written as

$$\frac{y^2 - 2v \overbrace{\left(\frac{x\sigma^2 + \mu h_n^2}{\sqrt{\sigma^2 + h_n^2}}\right)}^{\alpha} + \beta}{(h_n\sigma)^2} = \frac{y^2 - 2v\alpha + \beta}{(h_n\sigma)^2} = \frac{(y - \alpha)^2 - \alpha^2 + \beta}{(h_n\sigma)^2}.$$

We return to the integral in Equation (4), use $dv = (\sigma^2 + h_n^2)^{-1/2} dy$ and write

$$\frac{1}{\sqrt{2\pi}\sigma} \frac{1}{\sqrt{2\pi}h_n} \frac{1}{\sqrt{\sigma^2 + h_n^2}} \exp\left[-\frac{1}{2} \left(\frac{-\alpha^2 + \beta}{(h_n\sigma)^2} \right)\right] \underbrace{\int \exp\left[-\frac{1}{2} \left(\frac{y - \alpha}{h_n\sigma} \right)^2\right] dy}_{=\sqrt{2\pi}h_n\sigma}$$

where the integral is evaluated easily since it's merely $\mathcal{N}(\alpha, h_n^2\sigma^2)$. We have

$$\bar{p}_n(x) = \frac{1}{\sqrt{2\pi}\sqrt{\sigma^2 + h_n^2}} \exp\left[-\frac{1}{2} \left(\frac{\beta - \alpha^2}{(h_n\sigma)^2} \right)\right]$$

and all that remains is to now is clean up the exponent. We write

$$\begin{aligned}\frac{-\alpha^2 + \beta}{(h_n\sigma)^2} &= \frac{1}{(h_n\sigma)^2} \left[\frac{(\sigma^2 + h_n^2)(x^2\sigma^2 + h_n^2\mu^2)}{\sigma^2 + h_n^2} - \frac{(x\sigma^2 + \mu h_n^2)^2}{\sigma^2 + h_n^2} \right] \\ &= \frac{1}{\sigma^2 + h_n^2} (x^2 + 2\mu x + \mu^2) = \frac{(x - \mu)^2}{\sigma^2 + h_n^2},\end{aligned}$$

where some tedious algebra was omitted. Finally the integral becomes

$$\bar{p}_n(x) = \frac{1}{\sqrt{2\pi}\sqrt{\sigma^2 + h_n^2}} \exp \left[-\frac{1}{2} \frac{(x - \mu)^2}{\sigma^2 + h_n^2} \right],$$

which is clearly $\mathcal{N}(\mu, \sigma^2 + h_n^2)$, and this is what we were asked to show.

- b) Did not figure this one out. Skipped.
c) Did not figure this one out. Skipped.

Problem 4.3

- a) The mean of the Parzen-window estimate is given by

$$\bar{p}(x) = \frac{1}{V_n} \int p(v) \phi \left(\frac{x - v}{h_n} \right) dv,$$

where $p(v) = U(0, a)$. Two observations about when this integral is zero are needed.

- (i) $p(v) = U(0, a)$ is zero outside of $0 < v < a$.
(ii) $\phi \left(\frac{x-v}{h_n} \right)$ is zero when $x - v > 0 \Leftrightarrow v < x$.

Let's consider the integral in every case.

When $x < 0$, v must be 0 too since $v < x$, and the integral is zero.

When $0 < x < a$, v can range from 0 to x . We obtain

$$\begin{aligned}\bar{p}(x) &= \frac{1}{V_n} \int p(v) \phi \left(\frac{x - v}{h_n} \right) dv = \frac{1}{h_n} \int_0^x \frac{1}{a} \exp \left(\frac{v - x}{h_n} \right) dv \\ &= \frac{1}{ah_n} h_n \exp \left(\frac{v - x}{h_n} \right) \Big|_{v=0}^{v=x} = \frac{1}{a} \left(1 - \exp \left(\frac{-x}{h_n} \right) \right).\end{aligned}$$

When $x > a$, v is not affected by x and ranges from 0 to a . We obtain

$$\begin{aligned}\bar{p}(x) &= \frac{1}{V_n} \int p(v) \phi \left(\frac{x - v}{h_n} \right) dv = \frac{1}{h_n} \int_a^x \frac{1}{a} \exp \left(\frac{v - x}{h_n} \right) dv \\ &= \frac{1}{ah_n} h_n \exp \left(\frac{v - x}{h_n} \right) \Big|_{v=0}^{v=a} = \frac{1}{a} \left(\exp \left(\frac{a - x}{h_n} \right) - \exp \left(\frac{-x}{h_n} \right) \right) \\ &= \frac{1}{a} \left(\exp \left(\frac{a}{h_n} \right) - 1 \right) \exp \left(\frac{-x}{h_n} \right).\end{aligned}$$

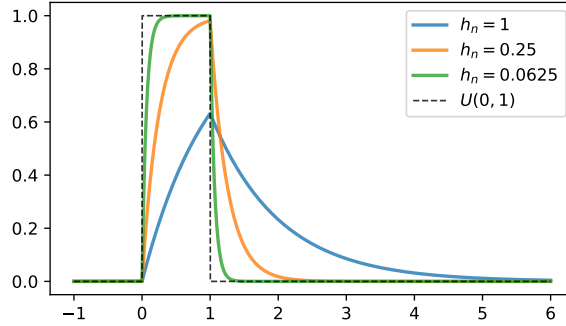


Figure 10: Plot of $\bar{p}(x)$ versus x for $a = 1$ and $h = 1, 1/4$ and $1/16$.

- b) The plot is found in Figure 10.
c) The bias is $\mathbb{E}(p(x) - \hat{p}(x)) = p(x) - \bar{p}(x)$, and we obtain the relative bias (in percentage) by dividing with $p(x)$ so that

$$\text{bias}(x) = \frac{|p(x) - \bar{p}(x)|}{p(x)} = \frac{\frac{1}{a} - \bar{p}(x)}{\frac{1}{a}} = 1 - ap(x) = e^{-x/h_n}.$$

The bias is decreasing on $0 < x < a$, so if we want the bias to be less than 0.01 on 99% of the interval, it amounts to requiring that

$$\text{bias}\left(\frac{a}{100}\right) = 0.01 \quad \Leftrightarrow \quad \exp\left(-\frac{a}{100h_n}\right) = 0.01.$$

Solving this equation for h_n reveals that $h_n = a/(100 \ln 100)$.

- d) When $a = 0$, h_n becomes $h_n = 1/(100 \ln 100) \approx 0.0022$. See Figure 11 for a plot. Although it is hard to tell exactly from the plot, observe that when $x = 0.01$, the relative bias is indeed close to 0.01 so that $\bar{p}(0.01) = 0.99$.

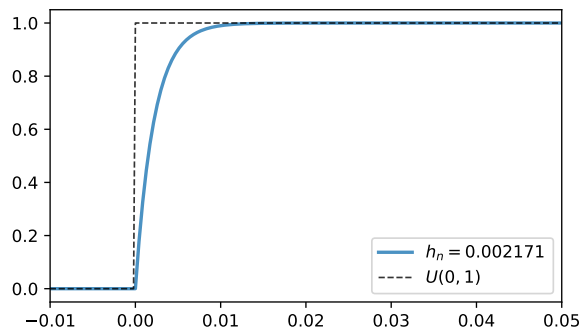


Figure 11: Plot accompanying problem 4.3d).

Problem 4.8

a) We are given that $P(\omega_i) = 1/c$ and

$$p(x | \omega_i) = \begin{cases} 1 & 0 \leq x \leq \frac{cr}{c-1} \\ 1 & i \leq x \leq i + 1 - \frac{cr}{c-1} \\ 0 & \text{elsewhere} \end{cases}.$$

It helps to visualize this function. The first line in the function definition above is for a common, overlapping region of the probability space. The second line defines a region of the space where each category ω_i has a non-overlapping positive area of the probability density function.

In other words: the value of a random x is in an overlapping part of the probability space with probability $\frac{cr}{c-1}$, and in a non-overlapping part of the space with probability $1 - \frac{cr}{c-1}$. The value of r is between 0 and $\frac{c-1}{c}$.

Every point x has a conditional probability $p(x | \omega_i)$ equal to 0 or 1. To find the probability of error, we consider regions of the space with overlapping distributions, denoted o , and non-overlapping regions, in turn. The Bayes error is

$$\begin{aligned} P^* &= P(e) = \int P(e | x)p(x) dx \\ &= P(e | o)P(o) + P(e | \text{not } o)P(\text{not } o) \\ &= \frac{c-1}{c} \frac{cr}{c-1} + 0 \left(1 - \frac{cr}{c-1}\right) = r \end{aligned}$$

b) The nearest neighbor error rate is given by Equation (45) in Chapter 4, which is

$$P = \int \left[1 - \sum_{i=1}^c P^2(\omega_i | x)\right] p(x) dx.$$

We'll integrate over the non-zero regions of the probability density. In the non-overlapping regions, the integral becomes

$$P_{\text{non-overlapping}} = \int \left[1 - \sum_{i=1}^c P^2(\omega_i | x)\right] p(x) dx = \int [1 - 1] p(x) dx = 0.$$

This makes intuitive sense, there is no error in non-overlapping regions in the limit of many data points. In the overlapping regions, we have

$$P_{\text{overlapping}} = \int \left[1 - \sum_{i=1}^c \frac{1}{c^2}\right] p(x) dx = \int_0^{\frac{cr}{c-1}} \left[1 - \frac{1}{c}\right] dx = r = P^*.$$

In other words, the Bayes error P^* equals the nearest neighbor error P , which both equal r .

Problem 4.17

We assume that $p(\omega_i) = 1/c$ and $p(\mathbf{x} | \omega_i) = p(\mathbf{x})$. In this case, for any point \mathbf{x} , any guess is as good as any other, and the Bayes error rate is clearly

$$P^* = \frac{c-1}{c}.$$

To prove that the bound

$$P \leq P^* \left(2 - \frac{c}{c-1} P^* \right) \quad (5)$$

is achieved, we calculate the error rate P . In the following calculation, we use the fact that $p(\mathbf{x} | \omega_i) = p(\mathbf{x})$ and $\int p(\mathbf{x}) d\mathbf{x} = 1$. We observe that

$$\begin{aligned} P &= \int \left[1 - \sum_{i=1}^c P^2(\omega_i | \mathbf{x}) \right] p(\mathbf{x}) d\mathbf{x} \\ &= \int \left[1 - \sum_{i=1}^c \left(\frac{p(\mathbf{x} | \omega_i) P(\omega_i)}{p(\mathbf{x})} \right)^2 \right] p(\mathbf{x}) d\mathbf{x} \\ &= \int \left[1 - \sum_{i=1}^c P^2(\omega_i) \right] p(\mathbf{x}) d\mathbf{x} \\ &= \left[1 - \sum_{i=1}^c \frac{1}{c^2} \right] \int p(\mathbf{x}) d\mathbf{x} \\ &= 1 - c \frac{1}{c^2} = \frac{c-1}{c}. \end{aligned}$$

In other words, P^* and P are equal. When we substitute P^* and P into Equation (5), we see that the bound is achieved since

$$P \leq P^* \left(2 - \frac{c}{c-1} P^* \right) \Rightarrow \frac{c-1}{c} \leq \frac{c-1}{c} \left(2 - \frac{c}{c-1} \frac{c-1}{c} \right) = \frac{c-1}{c},$$

which completes the proof.

Problem 4.27

- Every property of a metric is easy to prove, except the triangle inequality. I was unable to prove this.
- There is a typo in the book, since $\binom{6}{2} = 15$. The pairings are listed in Table 1 on page 39.
- Skipped this one.

Table 1: Ranked order of combinations of words for problem 4.27.

Word 1	Word 2	$D_{\text{Tantimoto}}$
pots	stop	0.0
pattern	elementary	0.444
pattern	pat	0.5
taxonomy	elementary	0.5
pat	pots	0.6
pat	stop	0.6
pattern	taxonomy	0.7
pattern	pots	0.75
pattern	stop	0.75
pat	taxonomy	0.75
pat	elementary	0.778
pots	taxonomy	0.778
stop	taxonomy	0.778
pots	elementary	0.909
stop	elementary	0.909

2.5 Linear discriminant functions

Problem 5.4

a) We wish to solve the problem

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{x}_a\|^2 = \mathbf{x}_a^T \mathbf{x}_a - 2\mathbf{x}^T \mathbf{x}_a + \mathbf{x}^T \mathbf{x} \\ & \text{subject to} \quad g(\mathbf{x}) = \boldsymbol{\omega}^T \mathbf{x} + \omega_0 = 0. \end{aligned}$$

To accomplish this, we start by constructing the Lagrange function

$$L(\mathbf{x}, \lambda) = \mathbf{x}_a^T \mathbf{x}_a - 2\mathbf{x}^T \mathbf{x}_a + \mathbf{x}^T \mathbf{x} - \lambda(\boldsymbol{\omega}^T \mathbf{x} + \omega_0 - 0),$$

which we differentiate with respect to \mathbf{x} and λ to obtain:

$$\begin{aligned} L_{\mathbf{x}} &= -2\mathbf{x}_a + 2\mathbf{x} - \lambda\boldsymbol{\omega} = 0 \\ L_{\lambda} &= \boldsymbol{\omega}^T \mathbf{x} + \omega_0 = 0 \end{aligned}$$

We wish to solve these equations for \mathbf{x} .

Solving the first equation for \mathbf{x} yields $\mathbf{x} = \lambda\boldsymbol{\omega}/2 + \mathbf{x}_a$. It remains to solve for λ . If we left-multiply by $\boldsymbol{\omega}^T$ and compare with the second equation, we observe that

$$\boldsymbol{\omega}^T \mathbf{x} = -\omega_0 \quad \text{and} \quad \boldsymbol{\omega}^T \mathbf{x} = \frac{\|\boldsymbol{\omega}\|^2 \lambda}{2} + \boldsymbol{\omega}^T \mathbf{x}_a.$$

This implies that

$$-\omega_0 = \frac{\|\boldsymbol{\omega}\|^2 \lambda}{2} + \boldsymbol{\omega}^T \mathbf{x}_a \quad \Leftrightarrow \quad \lambda = -\frac{2}{\|\boldsymbol{\omega}\|^2} (\omega_0 + \boldsymbol{\omega}^T \mathbf{x}_a),$$

and substituting this into $\mathbf{x} = \lambda\boldsymbol{\omega}/2 + \mathbf{x}_a$ yields the optimal answer

$$\mathbf{x}^* = -\frac{\boldsymbol{\omega}}{\|\boldsymbol{\omega}\|^2} (\boldsymbol{\omega}^T \mathbf{x}_a + \omega_0) + \mathbf{x}_a = -\frac{\boldsymbol{\omega}}{\|\boldsymbol{\omega}\|^2} g(\mathbf{x}_a) + \mathbf{x}_a. \quad (6)$$

Inserting this into $\|\mathbf{x}^* - \mathbf{x}_a\|$ yields

$$\|\mathbf{x}^* - \mathbf{x}_a\| = \left\| \frac{\boldsymbol{\omega}}{\|\boldsymbol{\omega}\|^2} g(\mathbf{x}_a) \right\| = \frac{|g(\mathbf{x}_a)|}{\|\boldsymbol{\omega}\|}.$$

b) The projection onto the plane is the minimizer \mathbf{x}^* from Equation (6) in the previous sub-problem, so we immediately see that

$$\mathbf{x}^* = -\frac{\boldsymbol{\omega}}{\|\boldsymbol{\omega}\|^2} (\boldsymbol{\omega}^T \mathbf{x}_a + \omega_0) + \mathbf{x}_a = -\frac{\boldsymbol{\omega}}{\|\boldsymbol{\omega}\|^2} g(\mathbf{x}_a) + \mathbf{x}_a = \mathbf{x}_a - \frac{g(\mathbf{x}_a)}{\|\boldsymbol{\omega}\|^2} \boldsymbol{\omega},$$

which is what we were required to show.

Problem 5.13

We wish to choose $\eta(k)$ to minimize the quadratic function

$$J(\mathbf{a}(k+1)) \simeq J(\mathbf{a}(k)) - \eta(k) \|\nabla \mathbf{J}\|^2 + \frac{1}{2} \eta^2(k) \nabla \mathbf{J}^T \mathbf{H} \nabla \mathbf{J}.$$

Differentiating and setting the result equal to zero yields

$$\frac{\partial J(\mathbf{a}(k+1))}{\partial \eta(k)} = -\|\nabla \mathbf{J}\|^2 + \eta(k) \nabla \mathbf{J}^T \mathbf{H} \nabla \mathbf{J} = 0,$$

and solving this equation for $\eta(k)$ yields the desired answer, which is

$$\eta(k) = \frac{\|\nabla \mathbf{J}\|^2}{\nabla \mathbf{J}^T \mathbf{H} \nabla \mathbf{J}}.$$

Problem 5.15

a) If $\alpha > \beta^2/2\gamma$, then $-2\alpha\gamma + \beta^2 < 0$ and

$$\|\mathbf{a}(k+1) - \alpha \hat{\mathbf{a}}\|^2 \leq \|\mathbf{a}(k) - \alpha \hat{\mathbf{a}}\|^2 - 2\alpha\gamma + \beta^2$$

represents error which decreases at each step. After k corrections we obtain

$$\|\mathbf{a}(k+1) - \alpha \hat{\mathbf{a}}\|^2 \leq \|\mathbf{a}(1) - \alpha \hat{\mathbf{a}}\|^2 + k(-2\alpha\gamma + \beta^2),$$

and the error is zero when $\|\mathbf{a}(1) - \alpha \hat{\mathbf{a}}\|^2 + k_0(-2\alpha\gamma + \beta^2) = 0$, which implies that

$$k_0 = \frac{\|\mathbf{a}(1) - \alpha \hat{\mathbf{a}}\|^2}{2\alpha\gamma - \beta^2}.$$

This is what we were asked to show.

b) Skipped. Somehow I did not see this problem when originally solving.

Problem 5.21

To ease the notation, let us write $\mathbf{m} := \mathbf{m}_1 - \mathbf{m}_2$ ¹. Starting with Equation (53) from Chapter 5 in [Duda et al., 2000], we left-multiply by the bracketed term to isolate \mathbf{w} as

$$\mathbf{w} = \left[\frac{1}{n} \mathbf{S}_W + \frac{n_1 n_2}{n^2} \mathbf{m} \mathbf{m}^T \right]^{-1} \mathbf{m}. \quad (7)$$

Recall from Problem 3.36 that the Sherman-Morrison-Woodbury matrix identity is

$$(\mathbf{A} + \mathbf{x} \mathbf{y}^T)^{-1} = \mathbf{A}^{-1} - \frac{\mathbf{A}^{-1} \mathbf{x} \mathbf{y}^T \mathbf{A}^{-1}}{1 + \mathbf{y}^T \mathbf{A}^{-1} \mathbf{x}}.$$

¹This definition of \mathbf{m} is not the same as the one used in the book, where \mathbf{m} is the grand mean.

We now apply the identity to the bracketed term in Equation (7). In doing so, we identify

$$\mathbf{A} \cong \frac{1}{n} \mathbf{S}_W \quad \text{and} \quad \mathbf{xy}^T \cong \frac{n_1 n_2}{n^2} \mathbf{mm}^T.$$

Applying the matrix identity, we obtain

$$\begin{aligned} \mathbf{w} &= \left[n \mathbf{S}_W^{-1} - \frac{n^2 \left(\frac{n_1 n_2}{n^2} \right) \mathbf{S}_W^{-1} \mathbf{mm}^T \mathbf{S}_W^{-1}}{1 + n \left(\frac{n_1 n_2}{n^2} \right) \mathbf{m}^T \mathbf{S}_W^{-1} \mathbf{m}} \right] \mathbf{m} \\ &= n \mathbf{S}_W^{-1} \mathbf{m} - \frac{n \left(\frac{n_1 n_2}{n} \right) \mathbf{S}_W^{-1} \mathbf{mm}^T \mathbf{S}_W^{-1} \mathbf{m}}{1 + \left(\frac{n_1 n_2}{n} \right) \mathbf{m}^T \mathbf{S}_W^{-1} \mathbf{m}}. \end{aligned}$$

To simplify the notation and remind ourselves that some of these quantities are simply scalars, let us denote $a := n_1 n_2 / n$ and $b := \mathbf{m}^T \mathbf{S}_W^{-1} \mathbf{m}$. We simplify and factor our $n \mathbf{S}_W^{-1} \mathbf{m}$ to obtain

$$\mathbf{w} = n \mathbf{S}_W^{-1} \mathbf{m} - \frac{na \mathbf{S}_W^{-1} \mathbf{m} b}{1 + ab} = n \mathbf{S}_W^{-1} \mathbf{m} \left[1 - \frac{ab}{1 + ab} \right] = n \mathbf{S}_W^{-1} \mathbf{m} [1 + ab]^{-1}.$$

Recalling now that $a := n_1 n_2 / n$ and $b := \mathbf{m}^T \mathbf{S}_W^{-1} \mathbf{m}$, we have accomplished the goal. The result is that

$$\mathbf{w} = n \mathbf{S}_W^{-1} \mathbf{m} \alpha = n \mathbf{S}_W^{-1} \mathbf{m} \left[1 + \underbrace{\left(\frac{n_1 n_2}{n} \right)}_a \underbrace{\left(\mathbf{m}^T \mathbf{S}_W^{-1} \mathbf{m} \right)}_b \right]^{-1},$$

which shows that α is indeed given by the quantity in the problem statement.

Problem 5.25

- a) We supply a proof by induction: we first show the base case, and then the inductive step.

The base case is verified by checking that the relation holds for $\eta^{-1}(2) = \eta^{-1}(1) + y_1^2$. This is true, since it implies

$$\eta(2) = \frac{1}{\eta^{-1}(1) + y_1^2} = \frac{\eta(1)}{1 + \eta(1)y_1^2} = \frac{\eta(1)}{1 + \eta(1) \sum_{i=1}^1 y_i^2},$$

which is clearly the given formula for $k = 2$.

In the inductive step we assume that the relation holds for $\eta(k)$, and show that this implies that it holds for $\eta(k+1)$ too. The required algebra is

$$\begin{aligned} \eta^{-1}(k+1) &= \eta^{-1}(k) + y_k^2 \\ &= \left(\frac{\eta(1)}{1 + \eta(1) \sum_{i=1}^{k-1} y_i^2} \right)^{-1} + y_k^2 \\ &= \frac{1 + \eta(1) \sum_{i=1}^{k-1} y_i^2}{\eta(1)} + \frac{\eta(1)}{\eta(1)} y_k^2 \\ &= \frac{1 + \eta(1) \sum_{i=1}^k y_i^2}{\eta(1)}. \end{aligned}$$

Inverting this shows that $\eta(k+1) = \eta(1) / \left(1 + \eta(1) \sum_{i=1}^k y_i^2\right)$, as required.

- b) To show why the sequence of coefficients will satisfy the sums, we will first bound the terms, and then convert the problems to integrals.

If $0 < a \leq y_i^2 \leq b < \infty$ for every i , then the sum is bounded by

$$a(k-1) \leq \sum_{i=1}^{k-1} y_i^2 \leq b(k-1),$$

and this in turn implies the expression $\eta(k)$ that may be bounded by

$$\frac{\eta(1)}{1 + \eta(1)b(k-1)} \leq \eta(k) \leq \frac{\eta(1)}{1 + \eta(1)a(k-1)}.$$

To show that $\sum \eta(k) \rightarrow \infty$, we note that $\sum \eta(k) \simeq \lim_{\alpha \rightarrow \infty} \int_{x=1}^{x=\alpha} \eta(x) dx$. We observe that the integral

$$\lim_{\alpha \rightarrow \infty} \int_{x=1}^{x=\alpha} \frac{\eta(1)}{1 + \eta(1)(x-1)b} dx = \lim_{\alpha \rightarrow \infty} \frac{1}{b} \ln |u| \Big|_{u=1}^{u=1+\eta(1)(\alpha-1)b}$$

diverges for any value of b , where we used the substitution $u = 1 + \eta(1)(x-1)b$. Since b represents the maximal value of the terms $\eta(k)$, any other value of y_i^2 will diverge too, and the sum $\sum \eta(k)$ diverges to infinity.

To show that $\sum \eta^2(k) \rightarrow L < \infty$, we again use $\sum \eta^2(k) \simeq \lim_{\alpha \rightarrow \infty} \int_{x=1}^{x=\alpha} \eta^2(x) dx$. Now the integral converges, since for any value of a the integral

$$\lim_{\alpha \rightarrow \infty} \int_{x=1}^{x=\alpha} \frac{\eta^2(1)}{(1 + \eta(1)(x-1)b)^2} dx = \lim_{\alpha \rightarrow \infty} \frac{\eta(1)}{bu} \Big|_{u=1+\eta(1)(\alpha-1)a}^{u=1} \leq \frac{\eta(1)}{b}$$

converges, and a represents the maximal bound on $\eta(k)$. The sum $\sum \eta(k)$ diverges to infinity, for any $0 < a \leq y_i^2 \leq b < \infty$.

Problem 5.27

- a) The data points are plotted in Figure 12 on page 44, and as seen in the plot they are not linearly separable.
b) From Equation (95) we observe that the optimal choice of η is given by

$$\eta(k) = \frac{\|\mathbf{Y}^T \mathbf{e}(k)\|^2}{\|\mathbf{Y}\mathbf{Y}^T \mathbf{e}(k)\|^2} = \frac{\mathbf{e}^T \mathbf{Y}\mathbf{Y}^T \mathbf{e}}{\mathbf{e}^T \mathbf{Y}\mathbf{Y}^T \mathbf{Y}\mathbf{Y}^T \mathbf{e}}.$$

This value varies from loop to loop, depending on \mathbf{e} . For this specific data set, the value of \mathbf{Y} and $\mathbf{Y}\mathbf{Y}^T$ are given in equation (b)).

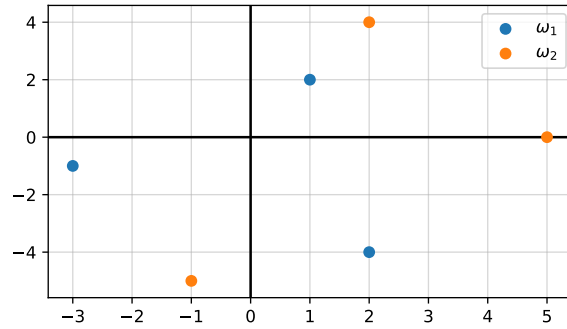


Figure 12: Graf accompanying problem 5.27.

$$\mathbf{Y} = \begin{pmatrix} 1.0 & 1.0 & 2.0 \\ 1.0 & 2.0 & -4.0 \\ 1.0 & -3.0 & -1.0 \\ -1.0 & -2.0 & -4.0 \\ -1.0 & 1.0 & 5.0 \\ -1.0 & -5.0 & 0.0 \end{pmatrix}$$

$$\mathbf{Y}\mathbf{Y}^T = \begin{pmatrix} 6.0 & -5.0 & -4.0 & -11.0 & 10.0 & -6.0 \\ -5.0 & 21.0 & -1.0 & 11.0 & -19.0 & -11.0 \\ -4.0 & -1.0 & 11.0 & 9.0 & -9.0 & 14.0 \\ -11.0 & 11.0 & 9.0 & 21.0 & -21.0 & 11.0 \\ 10.0 & -19.0 & -9.0 & -21.0 & 27.0 & -4.0 \\ -6.0 & -11.0 & 14.0 & 11.0 & -4.0 & 26.0 \end{pmatrix}$$

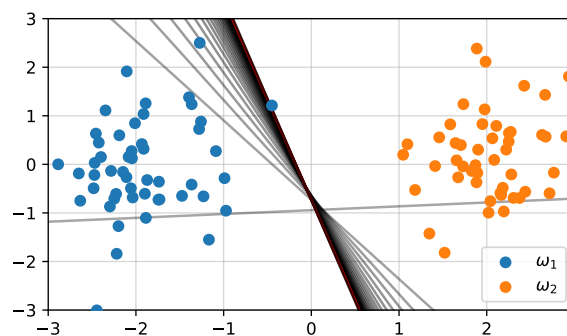


Figure 13: Convergence of the Ho Kashyap algorithm.

Problem 5.29

To show that there always exists a mapping to a higher dimensional space which leaves points from two classes linearly separable, we will explicitly provide such a mapping. The mapping would be very inefficient in practice, but provides a constructive proof.

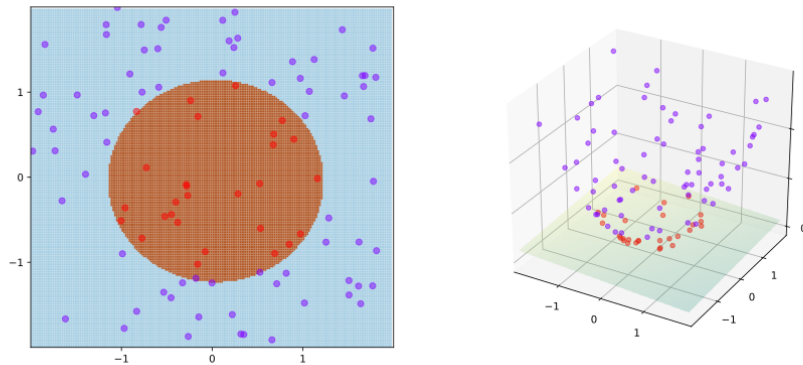


Figure 14: Using a kernel to raise points in a new dimension. Source: https://commons.wikimedia.org/wiki/File:Kernel_trick_idea.svg

Observe first that to apply a linear classifier to points where data from one class is close to the origin, a function such as $y = \phi(\mathbf{x}) = \exp(-\mathbf{x}^T \mathbf{x})$ may be used. This leaves the data linearly separable in the new space, as illustrated in Figure 14. Below we will extend this idea by introducing many distinct mappings $\phi(\mathbf{x})$.

Consider now points $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ in \mathbb{R}^d . Assume that some points $S \subseteq \mathcal{D}$ belong to ω_1 , and that we know exactly which points. If we know which points belong to ω_1 , then a kernel density estimate (Parzen window) such as

$$y = \text{Parzen}(S) = \frac{1}{|S|} \sum_{\mathbf{x}_i \in S} \frac{1}{h} \phi\left(\frac{\mathbf{x} - \mathbf{x}_i}{h}\right)$$

with a sufficiently small value of the bandwidth h will raise these points in the new y feature space. When $h \ll 1$, neighboring points not belonging to ω_1 will be unaffected. A plane such as $y = 0.1$ will then perfectly separate the points.

Since we do not know the true subset $S \subseteq \mathcal{D}$, we can apply this transformation on *every* possible subset of \mathcal{D} (the *power set* of \mathcal{D}). There are $2^{|S|}$ such subsets. We use

$$\mathbf{y} = (\mathbf{x} \quad \text{Parzen}(S_1) \quad \text{Parzen}(S_1) \quad \dots \quad \text{Parzen}(S_{2^{|S|}}))$$

to map $\mathbf{x} \in \mathbb{R}^d$ to a $d + 2^{|S|}$ space. In other words: if $\{\mathbf{x}_1\}$ is “raised” in one new feature dimension, $\{\mathbf{x}_1, \mathbf{x}_2\}$ in another, $\{\mathbf{x}_1, \mathbf{x}_3\}$ in yet another, etc. for *every* combination of points, then in some dimension in the feature space the points are linearly separable.

Problem 5.32

- a) The plot is show in Figure 15. By inspection the weight vector is

$$\mathbf{a} = (a_0, a_1, a_2) = (-1.5, 1, 1).$$

This corresponds to the line $y = 1.5 - x$. The optimal margin is the distance from the line $y = 1.5 - x$ to a point, say $(1, 1)$, and this distance is $\sqrt{2}/4 \approx 0.3536$.

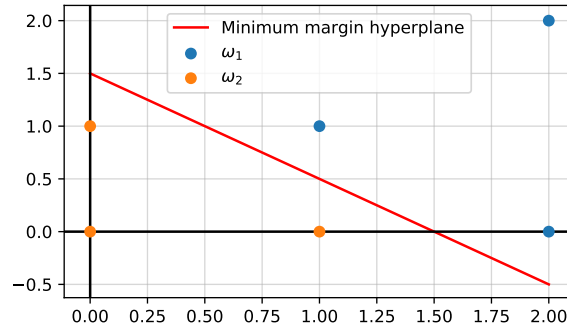


Figure 15: Plot related to problem 5.32a).

b) From inspection there are four support vectors, and they are

$$\begin{array}{cc} (0, 1)^T & (1, 0)^T \\ (1, 1)^T & (2, 0)^T. \end{array}$$

c) This laborious computation is omitted.

Problem 5.33

a) The optimization problem

$$L(\mathbf{a}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{a}\|^2 - \sum_{k=1}^n \alpha_k [z_k \mathbf{a}^T \mathbf{y}_k - 1] \quad (8)$$

has a solution which is a saddle point, since we wish to maximize with respect to $\boldsymbol{\alpha}$ and minimize with respect to \mathbf{a} .

b) Here I believe there to be a slight error in the text, which might lead to confusion. In the initial pages of the chapter, the distance from a point to a hyperplane is given by $r = (\boldsymbol{\omega}^T \mathbf{x} + \omega_0) / \|\boldsymbol{\omega}\|$, and this is indeed correct.

In the context of SVMs, however, the distance is said to be

$$\frac{g(\mathbf{y})}{\|\mathbf{a}\|} = \frac{\mathbf{a}^T \mathbf{y}}{\|\mathbf{a}\|} = \frac{(\omega_0, \mathbf{a}') (1, \mathbf{y}')^T}{\|\mathbf{a}\|} = \frac{\mathbf{a}'^T \mathbf{y} + \omega_0}{\|\mathbf{a}\|},$$

which is different because \mathbf{a} is in the denominator, not \mathbf{a}' . Amending Equation (8) with this information, the $\frac{1}{2} \|\mathbf{a}\|^2$ should be replaced by $\frac{1}{2} \|\mathbf{a}'\|^2$, i.e. dropping $\mathbf{a}_0 = \omega_0$. If we do this and differentiate with respect to the first component of \mathbf{a} , we obtain

$$\frac{\partial L(\mathbf{a}, \boldsymbol{\alpha})}{\partial \mathbf{a}_0} = \sum_k \alpha_k^* z_k \mathbf{y}_0 = 0,$$

which gives the desired result since $\mathbf{y}_0 = 1$ due to the augmentation of the vector.

c) To prove this, we differentiate with respect to \mathbf{a} and obtain

$$\frac{\partial L(\mathbf{a}, \boldsymbol{\alpha})}{\partial \mathbf{a}} = \mathbf{a}^* - \sum_k \alpha_k^* z_k \mathbf{y}_k = 0.$$

d) If the Lagrange multiplier (or *undetermined multiplier*) α_k^* is zero, then it's said to be *inactive*. At the optimum, the constraint is not used. The optimum of $L(\mathbf{a}, \boldsymbol{\alpha})$ is then the same with or without this constraint.

If the Lagrange multiplier α_k^* is non-zero, then the constraint is said to be *active*. The constrained solution is different from the unconstrained solution, and the optimum lies on the boundary of the constraint. Since $\alpha_k^* z_k \mathbf{y}_k \geq 1$ in the feasible region, the optimal solution is on the boundary if $\alpha_k^* z_k \mathbf{y}_k = 1$, but then α_k^* is non-zero since the constraint is active.

In conclusion, either $\alpha_k^* z_k \mathbf{y}_k = 1$ if the constraint is active, or $\alpha_k^* = 0$ if the constraint is inactive. This is one of the Karush–Kuhn–Tucker (KKT) conditions, and it may be expressed as

$$\alpha_k^* [\alpha_k^* z_k \mathbf{y}_k - 1] = 0 \quad k = 1, \dots, n.$$

e) We simply multiply the brackets in Equation (8) from subproblem a).

$$L(\mathbf{a}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{a}\|^2 - \sum_{k=1}^n \alpha_k [z_k \mathbf{a}^T \mathbf{y}_k - 1] = \frac{1}{2} \|\mathbf{a}\|^2 - \sum_{k=1}^n \alpha_k z_k \mathbf{a}^T \mathbf{y}_k + \sum_{k=1}^n \alpha_k$$

f) Using $\mathbf{a}^* = \sum_j \alpha_j^* z_j \mathbf{y}_j$ we observe that

$$\begin{aligned} L(\boldsymbol{\alpha}) &= \frac{1}{2} \|\mathbf{a}^*\|^2 - \sum_{k=1}^n \alpha_k z_k \mathbf{a}^{*T} \mathbf{y}_k + \sum_{k=1}^n \alpha_k \\ &= \frac{1}{2} \left(\sum_j \alpha_j^* z_j \mathbf{y}_j^T \right) \left(\sum_k \alpha_k^* z_k \mathbf{y}_k \right) - \sum_{k=1}^n \alpha_k z_k \left(\sum_j \alpha_j^* z_j \mathbf{y}_j^T \right) \mathbf{y}_k + \sum_{k=1}^n \alpha_k \end{aligned}$$

and since the first and second terms are equal, we obtain

$$-\frac{1}{2} \left(\sum_j \alpha_j^* z_j \mathbf{y}_j^T \right) \left(\sum_k \alpha_k^* z_k \mathbf{y}_k \right) + \sum_{k=1}^n \alpha_k = -\frac{1}{2} \sum_j \sum_k \alpha_j^* \alpha_k^* z_j z_k \mathbf{y}_j^T \mathbf{y}_k + \sum_{k=1}^n \alpha_k$$

as desired. We have formulated the problem as a maximization over $L(\boldsymbol{\alpha})$.

2.6 Multilayer Neural Networks

Problem 6.5

The backpropagation rule for input-to-hidden weights is given by Equation (21), which is

$$\Delta w_{ji} = \eta \delta_j x_i = \eta \left[\sum_{k=1}^c w_{kj} \delta_k \right] f'(\text{net}_j) x_i.$$

Notice that if $z_k = t_k$ for every value of k , there is nothing to learn since the error is zero. In this case, $\delta_k = 0$ and as a result $\Delta w_{ji} = 0$. The learning rate is proportional to the error ($t_k - z_k$), but it's also weighted by w_{kj} . The higher this weight is, the more an error contributes and the further it is adjusted in the backpropagation rule.

Problem 6.8

- a) In this problem it helps to make a drawing, or study Figures 6.4 and 6.2 in the book.
 - (i) The bias is connected to $n_H + c$ weights.
 - (ii) Every one of the d inputs is connected to every one of the n_H hidden units, for a total of dn_H weights.
 - (iii) Every hidden unit n_H is connected to every one of the c outputs, for $n_H c$ weights.

The total number of weights is therefore given by

$$n_H + c + dn_H + n_H c = n_H(1 + d + c) + c.$$

- b) Consider the equation

$$z_k = f \left[\sum_{j=1}^{n_H} w_{kj} f \left(\sum_{i=1}^d w_{ji} x_i + w_{j0} \right) + w_{k0} \right]$$

for a neural network output. If the sign is flipped on every weight going into a hidden unit, and the weights leading out from that same unit are also flipped, then the net result is no change if the activation function obeys $f(-x) = -f(x)$. In other words, if the $w_{ji} \mapsto -w_{ji}$ and $w_{j0} \mapsto -w_{j0}$ in the equation above, then

$$f \left(\sum_{i=1}^d -w_{ji} x_i - w_{j0} \right) = -f \left(\sum_{i=1}^d w_{ji} x_i + w_{j0} \right)$$

and a mapping $w_{kj} \mapsto -w_{kj}$ in the outer sum will cancel the sign flip. Note that this result only applies to odd functions where $f(-x) = -f(x)$.

- c) If there are n_H hidden units, there are $n_H!$ ways to permute them. For each of these permutations, every associated weight might have its sign flipped or not, for a total of 2^{n_H} possibilities with respect to sign flipping. The total count is therefore $n_H! 2^{n_H}$ exchange symmetries (not $n_H 2^{n_H}$ as [Duda et al., 2000] claims). This result is verified in [Bishop, 2011] on page 232, so we assume that [Duda et al., 2000] contains a typo here and the “!” sign was forgotten.

Problem 6.10

- a) The rules of differentiation in calculus shows that when $f(x) = 1/(1 + e^{ax})$, then the derivative may be expressed in terms of the function as

$$f'(x) = -ae^{ax}f^2(x).$$

- b) We'll study $f(x) = a(1 - e^{-2bx})/(1 + e^{-2bx})$, and it pays off to ease notation somewhat. Letting $g(x) := e^{-2bx}$, we have $g'(x) = -2bg(x)$ and we expedite notation by writing

$$f(x) = a \frac{1 - e^{-2bx}}{1 + e^{-2bx}} := a \frac{1 - g(x)}{1 + g(x)} = a(1 - g(x))(1 + g(x))^{-1}.$$

Differentiating this and using the chain rule and product rule of calculus, we obtain

$$\begin{aligned} f'(x) &= \frac{-ag'(x)(1 + g(x)) - ag'(x)(1 - g(x))}{(1 + g(x))^2} \\ &= \frac{4abg(x)}{(1 + g(x))^2} = a \underbrace{\frac{1 - g(x)}{1 + g(x)}}_{f'(x)} \frac{4bg(x)}{(1 - g(x))(1 + g(x))}. \end{aligned}$$

Substituting back $g(x) := e^{-2bx}$, we see that the derivative may be expressed in terms of the function as

$$f'(x) = f(x) \frac{4be^{-2bx}}{1 - e^{-4bx}}.$$

Problem 6.16

- a) Without loss of generality we assume that $J(\mathbf{w}) = \frac{1}{2}\mathbf{w}^T \mathbf{H} \mathbf{w}$, since by a transformation of the type $\mathbf{w} = \mathbf{w}' - \mathbf{c}$ we could transform a general, non-centered quadratic equation to this form. This is analogous to scaling $f(x) = ax^2 + bx + c$ in order to formulate it as $f(x') = a'x'^2$.

To study convergence, we examine the learning rule, which may be written as

$$\begin{aligned} \mathbf{w}_{n+1} &= \mathbf{w}_n - \eta \nabla J(\mathbf{w}_n) \\ &= \mathbf{w}_n - \eta \mathbf{H} \mathbf{w}_n \\ &= (\mathbf{I} - \eta \mathbf{H}) \mathbf{w}_n, \end{aligned}$$

and observe that $\mathbf{w}_n = (\mathbf{I} - \eta \mathbf{H})^n \mathbf{w}_0$. Convergence is ensured if $(\mathbf{I} - \eta \mathbf{H})^n$ goes to zero as $n \rightarrow \infty$. Informally, repeated application of the matrix should bring the vector to zero. The convergence of repeated application of a matrix to a vector is controlled by eigenvectors and eigenvalues.

We write $(\mathbf{I} - \eta \mathbf{H}) = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^T$ to express the symmetric, positive definite matrix in it's eigenvalue decomposition. Then we define $\mathbf{w}'_k = \mathbf{Q} \mathbf{w}_k = \sum_{i=1}^d \alpha_i \mathbf{v}_i$ to express \mathbf{w}_k in the *eigenvalue basis* of \mathbf{H} , where the \mathbf{v}_i s are the orthonormal eigenvectors. This transformation is a rotation of the space, since \mathbf{Q} is orthonormal.

The expression on the left below is for the vector, and the expression on the right captures the behavior of a single eigenvector component \mathbf{v}_i .

$$\left(\sum_{i=1}^d \alpha_i \mathbf{v}_i \right)_{n+1} = \mathbf{Q} \Lambda \left(\sum_{i=1}^d \alpha_i \mathbf{v}_i \right)_n \Leftrightarrow (\alpha_i \mathbf{v}_i)_{n+1} = \mathbf{Q} \lambda_i (\alpha_i \mathbf{v}_i)_{n+1}$$

The matrix \mathbf{Q} is orthogonal, so it merely rotates or reflects a vector, but $\|\mathbf{v}_i\| = \|\mathbf{Q}\mathbf{v}_i\|$ so the magnitude is unchanged. Therefore we ignore it when considering convergence, and see that learning is ensured if $|\lambda_i| < 1$ for every $i = 1, 2, \dots, d$. Since every eigenvalue is positive, the maximal eigenvalue controls convergence.

The only thing left to do is to relate the eigenvalues of $(\mathbf{I} - \eta\mathbf{H})$ to the eigenvalues of \mathbf{H} . If a matrix is multiplied by a scalar, then so are its eigenvalues. The effect of subtraction by the identity matrix is not as obvious, but if we write out the equation for the characteristic polynomials we observe that

$$\begin{aligned} \text{eig. vals of } (\mathbf{I} - \eta\mathbf{H}) &\Leftrightarrow \text{zeros of } \det(\mathbf{I} - \eta\mathbf{H} - \lambda\mathbf{I}) = \det(\mathbf{I}(1 - \lambda) - \eta\mathbf{H}) \\ \text{eig. vals of } \eta\mathbf{H} &\Leftrightarrow \text{zeros of } \det(\mathbf{I}\lambda' - \eta\mathbf{H}), \end{aligned}$$

and comparing we observe that $\lambda' = \lambda - 1$. We're really after the eigenvalues of \mathbf{H} , which we'll now denote by $\lambda_{\mathbf{H}} = \lambda'/\eta$. From this, we have convergence if

$$\begin{aligned} -1 &< \lambda_{\max} < 1 \\ -1 &< 1 - \eta\lambda_{\mathbf{H}\max} < 1 \\ \eta\lambda_{\mathbf{H}\max} &< 2, \end{aligned}$$

which is the same as requiring that $\eta < \lambda_{\mathbf{H}\max}/2$.

- b) The best learning rate is achieved if we manage to choose η such that the maximal eigenvalue (in absolute value) of $(\mathbf{I} - \eta\mathbf{H})$ is close to zero. The optimal rate is given by²

$$\text{rate}(\eta) = \max_i |1 - \eta\lambda_i| = \max(|1 - \eta\lambda_1|, |1 - \eta\lambda_d|)$$

where λ_1 is the smallest eigenvalue of \mathbf{H} and λ_d is the largest. Convergence is fastest when these arguments are equal in absolute value, so we require that

$$(1 - \eta\lambda_1) = -(1 - \eta\lambda_d)$$

and solve for η to obtain $\eta^* = 2/(\lambda_1 + \lambda_d)$. The optimal learning rate is obtained by substituting this η^* into the function $\text{rate}(\cdot)$. We consider the positive value $(1 - \eta^*\lambda_1)$, which becomes

$$\text{rate}(\eta^*) = (1 - \eta^*\lambda_1) = \frac{\lambda_d - \lambda_1}{\lambda_d + \lambda_1} = \frac{\lambda_d/\lambda_1 - 1}{\lambda_d/\lambda_1 + 1},$$

and recognize that the learning rate is indeed dependent on the ratio of the largest to the smallest eigenvalue of \mathbf{H} .

²See <https://distill.pub/2017/momentum/> for details.

- c) An informal argument is as follows: Consider unstandardized data which is not highly correlated (i.e. the covariance matrix is nearly diagonal), but where data is from different dimensions and is on vastly different scales (i.e. the diagonal of the covariance matrix has entries on different scales). Assume that this data is used to train a neural network, and that the weights are initialized to sensible, small, random values. Changing one weight value w_1 might then potentially change the error drastically compared to changing another weight w_2 , and $J(\mathbf{w}) \approx \frac{1}{2} \mathbf{w}^T \mathbf{H} \mathbf{w}$ will have diagonals on vastly different scales. The eigenvalues are $\lambda_1 \approx \min \text{diag } \mathbf{H}$ and $\lambda_d \approx \max \text{diag } \mathbf{H}$, so learning will be slow since $\text{rate}(\eta^*)$ will be close to 1.

Standardizing alleviates these problems. The data is on the same scale, so the diagonals entries of \mathbf{H} will likely be on the same scale and learning will be faster.

- d) Standardizing the data by subtracting means and diving by standard deviations in every dimension individually transforms the data to the same scale. The data might still be highly correlated. The whitening transform consists of scaling *and* rotating the data. The rotation diagonalizes the covariance matrix, and the scaling makes the spectrum of eigenvalues uniform, so that the variances are all unity.

Problem 6.21

- a) The softmax function is given by Equation (30) in Chapter 6 in the book, which is

$$z_k = f(\text{net}_k) = \frac{e^{\text{net}_k}}{\sum_{m=1}^c e^{\text{net}_m}} \propto e^{\text{net}_k}.$$

The learning rule when the sum squared error function is used becomes dependent on $\frac{\partial J}{\partial \text{net}_k}$ and $\frac{\partial \text{net}_k}{\partial w_{kj}}$, see Equation (13) in the book. The only change is that now $f(\text{net}_k)$ is not a sigmoid, but the softmax function. To ease notation, we define $k := \text{net}_k$, and differentiate.

$$\frac{\partial f(k)}{\partial k} = e^k \left[\sum_{m=1}^c e^m \right]^{-1} - e^k \left[\sum_{m=1}^c e^m \right]^{-2} e^k = \frac{e^k \left[\sum_{m \neq k}^c e^m \right]}{\left[\sum_{m=1}^c e^m \right]^2}$$

The update rules becomes a simple modification of Equation (17) from the book:

$$\Delta w_{kj} = \eta \delta_k y_j = \eta (t_k - z_k) \frac{e^{\text{net}_k} \left[\sum_{m \neq k}^c e^{\text{net}_m} \right]}{\left[\sum_{m=1}^c e^{\text{net}_m} \right]^2} y_j.$$

Extending the result to Δw_{ji} is straightforward, since Δw_{ji} a function of δ_k above.

- b) When employing cross entropy, the learning rule is identical to the sum squared error above, save for $\frac{\partial J}{\partial z_k}$, which now becomes

$$\frac{\partial J(\mathbf{w})}{\partial z_k} = \frac{\partial}{\partial z_k} \left[\sum_{k=1}^c t_k \ln \frac{t_k}{z_k} \right] = t_k \frac{z_k}{t_k} \frac{\partial}{\partial z_k} (t_k z_k^{-1}) = -\frac{t_k}{z_k}.$$

Everything else is exactly equal to the solution to sub-problem a) above.

Problem 6.24

We compare Equation (7) for the 3-layer neural network with Equation (32) for general additive models (GAM). The equations are, respectively, given by

$$z_k = f \left[\sum_{j=1}^{n_H} w_{kj} f \left(\sum_{i=1}^d w_{ji} x_i + w_{j0} \right) + w_{k0} \right] \quad \text{and}$$

$$z = f \left[\sum_{j=1}^d w_{kj} f_i(\mathbf{x}_i) + w_{k0} \right].$$

The functions $f_i(\mathbf{x}_i)$ in GAM may in general be multivariate. According to Wikipedia “The GAM model class is quite broad, given that *smooth function* is a rather broad category. For example, a covariate \mathbf{x}_i may be multivariate and the corresponding f_i a smooth function of several variables, or f_i might be the function mapping the level of a factor to the value of a random effect.”

Comparing the equations, we observe that

$$\sum_{j=1}^d w_{kj} f_i(\mathbf{x}_i) \cong \sum_{j=1}^{n_H} w_{kj} f \left(\sum_{i=1}^d w_{ji} x_i + w_{j0} \right).$$

In the typical case, the f on the right hand side a sigmoid function. The right hand side is a weighted sum of sigmoids, where the input to the sigmoids are again a linear function of the neural network inputs. The number of terms in the sum has no significance in the left-hand sum, since a function $f_i(\mathbf{x})$ could be defined as $f_1(\mathbf{x}) + f_2(\mathbf{x})$. The *Kolmogorov-Arnold representation theorem* guarantees that both of these constructions can, in theory, approximate any continuous function.

Problem 6.39

a) If we write out the sums, we obtain

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{K} \mathbf{x} = \sum_i \sum_j x_i K_{ij} x_j.$$

Of course, the function $f(\cdot)$ above is a mapping from a vector $\mathbf{x} \in \mathbb{R}^d$ to a real number \mathbb{R} . In this setting, the derivative is the gradient, i.e. $\nabla f(\mathbf{x}) = f'(\mathbf{x})$. To find the gradient using explicit components, we write out the sums as

$$\begin{aligned} \frac{d}{dx_k} (\mathbf{x}^T \mathbf{K} \mathbf{x}) &= \\ \frac{d}{dx_k} \left(\sum_i x_i \sum_j K_{ij} x_j \right) &= \\ \frac{d}{dx_k} \left(x_1 \left(\sum_j K_{1j} x_j \right) + x_2 \left(\sum_j K_{2j} x_j \right) + \dots + x_k \left(\sum_j K_{kj} x_j \right) + \dots \right). \end{aligned} \quad (9)$$

For every term $i \neq k$, the derivative is $x_i K_{ik}$ since every other term in the sums vanish. On the k 'th term, we apply the product rule of differentiation to obtain

$$\frac{d}{dx_k} \left(x_k \left(\sum_j K_{kj} x_j \right) \right) = \sum_j K_{kj} x_j + x_k (K_{kk}).$$

Applying these results to the non- k 'th terms and the k 'th term respectively, Equation (9) may be written in a more readable form as

$$\begin{aligned} \frac{d}{dx_k} (\mathbf{x}^T \mathbf{K} \mathbf{x}) &= x_1 K_{1k} + x_2 K_{2k} + x_3 K_{3k} + \dots + \left(\sum_j K_{kj} x_j + x_k K_{kk} \right) + \dots \\ &= \sum_i x_i K_{ik} + \sum_j K_{kj} x_j = \mathbf{K}^T \mathbf{x} + \mathbf{K} \mathbf{x} = (\mathbf{K}^T + \mathbf{K}) \mathbf{x}, \end{aligned}$$

where the last equality follows from

$$\begin{aligned} \mathbf{b} = \mathbf{A} \mathbf{x} &\Leftrightarrow b_i = \sum_j A_{ij} x_j \\ \mathbf{b} = \mathbf{A}^T \mathbf{x} &\Leftrightarrow b_i = \sum_j A_{ji} x_j. \end{aligned}$$

- b) Here's an approach which does not require the use of indices, and is therefore more expedient. Let $f(\mathbf{x}) = \mathbf{x}^T \mathbf{H} \mathbf{x}$, where \mathbf{H} is symmetric. We have

$$f(\mathbf{x} + \delta \mathbf{x}) - f(\mathbf{x}) = (\delta \mathbf{x})^T \mathbf{H} \mathbf{x} + \mathbf{x}^T \mathbf{H} (\delta \mathbf{x}) + O(\|\delta \mathbf{x}\|^2),$$

and when \mathbf{H} is symmetric this becomes

$$f(\mathbf{x} + \delta \mathbf{x}) - f(\mathbf{x}) = 2 (\delta \mathbf{x})^T \mathbf{H} \mathbf{x} + O(\|\delta \mathbf{x}\|^2).$$

Observe that $2\mathbf{H}\mathbf{x}$ is the first-order approximation to $f(\mathbf{x} + \delta \mathbf{x})$, i.e. the derivative.

Problem 6.42

The weight decay rule of Equation (38) in the book does not exactly lead to Equation (39). The rule multiplies the weight found by gradient descent by the factor $1 - \epsilon$, such that

$$\mathbf{w}_{n+1} = (\mathbf{w}_n - \eta \nabla J(\mathbf{w}_n)) (1 - \epsilon).$$

After some algebra, we observe that this is equivalent to

$$\mathbf{w}_{n+1} = \mathbf{w}_n - \underbrace{\eta \left((1 - \epsilon) \nabla J(\mathbf{w}_n) + \frac{\epsilon}{\eta} \mathbf{w}_n \right)}_{\nabla J_{ef}(\mathbf{w}_n)}.$$

Since $\nabla \mathbf{w}^T \mathbf{w} = 2\mathbf{w}$, we see that

$$J_{ef}(\mathbf{w}_n) = (1 - \epsilon) J(\mathbf{w}_n) + \epsilon \frac{1}{2\eta} \mathbf{w}_n^T \mathbf{w}_n.$$

This is a weighted average of the ordinary error function $J(\mathbf{w}_n)$ and the regularization term $\frac{1}{2\eta}\mathbf{w}_n^T\mathbf{w}_n$, weighted by $(1 - \epsilon)$ and ϵ respectively. We can easily verify that

- i) When $\epsilon = 0$, everything reduces to gradient descent with no weight decay.
- ii) When $\epsilon = 1$, the weight vector \mathbf{w}_{n+1} is forever stuck at $\mathbf{0}$.

This also shows that Equation (39) in the book is wrong with respect to Equation (38), since if $\epsilon = 1$ in the book, Equation (38) would always set $\mathbf{w}_{n+1} = \mathbf{0}$, while optimizing Equation (39) would not—it would amount to optimizing $J(\mathbf{w})$ with a regularization term.

2.7 Stochastic methods

Problem 7.4

Every magnet can be turned up or down, so there are 2^N possible configurations. However, flipping the sign of a state vector \mathbf{s} does not alter E . For instance $\mathbf{s} = (1, 1, -1)$ results in the same energy as $\mathbf{s} = (-1, -1, 1)$. Therefore, only half of the 2^N states correspond to (possibly) unique energy levels, i.e. 2^{N-1} unique energy levels.

The total time required for exhaustive search will therefore be $T(N) = 2^{N-1} \cdot 10^{-8}$ seconds. Searching the space for $N = 100$ units would take $6.338 \cdot 10^{21}$ seconds. Searching the space for $N = 1000$ units would take $5.357 \cdot 10^{292}$ seconds. Both of these numbers are huge. Searching the space would take millions and millions of years.

Problem 7.5

- a) The formula is $E \propto \sum_{i=1}^n \sum_{j=1}^n w_{ij} s_i s_j$, which is naively computed using n^2 flops. Since we have symmetry and $w_{ij} = w_{ji}$, we can write the sum as $\sum_{i=1}^n \sum_{j=i+1}^n w_{ij} s_i s_j$. The number of flops is then

$$(n-1) + (n-2) + \dots + 2 + 1 = \frac{n(n-1)}{2}.$$

Furthermore, there are 2^{n-1} possibly unique energy levels. So the total time is given by

$$T(n) = 2^{n-2} n(n-1) \cdot 10^{-10}.$$

- b) See figure 16 for a plot.

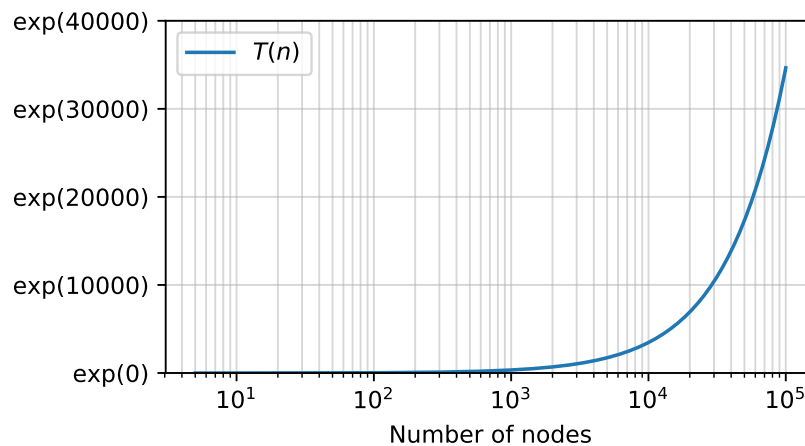


Figure 16: Solution to problem 7.5b).

- c) To find the value of n which may be computed in a day, we find n such that

$$T(n) = 3600 \cdot 24 \quad \text{or, alternatively} \quad \ln T(n) = \ln(3600 \cdot 24).$$

Working in log-space avoids numerical overflow, since $T(n)$ grows very large very quickly. The answers are:

$$\begin{aligned} \text{A day:} \quad \ln T(n) &= \ln(3600 \cdot 24) && \Rightarrow n = 40 \\ \text{A year:} \quad \ln T(n) &= \ln(3600 \cdot 24 \cdot 365) && \Rightarrow n = 48 \\ \text{A century:} \quad \ln T(n) &= \ln(3600 \cdot 24 \cdot 365 \cdot 100) && \Rightarrow n = 55 \end{aligned}$$

Problem 7.6

We wish to show that at high temperature, every configuration is equally likely. The probability that a system is in a configuration of energy E_γ is given by Equation (3) in chapter 7 the book, i.e. $P(\gamma) \propto \exp(-E_\gamma/T)$. As T goes to infinity, we have

$$\lim_{T \rightarrow \infty} P(\gamma) = \lim_{T \rightarrow \infty} (e^{-1/T})^{E_\gamma} = 1^{E_\gamma} = 1,$$

and so the probability of every state γ is equally likely. The normalization constant $Z(T)$ makes it a true probability function, i.e. ensures that it sums to unity.

Computer exercise 7.2

A simple implementation of Algorithm 1 in from Chapter 7 is found in the Python file `simulated_annealing.py`. The answers to sub-problem a) and b) are shown in Figures 17 and 18 respectively. A plot showing the average reduction in energy over 500 runs in show in Figure 19, and it's interesting to see that the average reduction appears to be exponential.

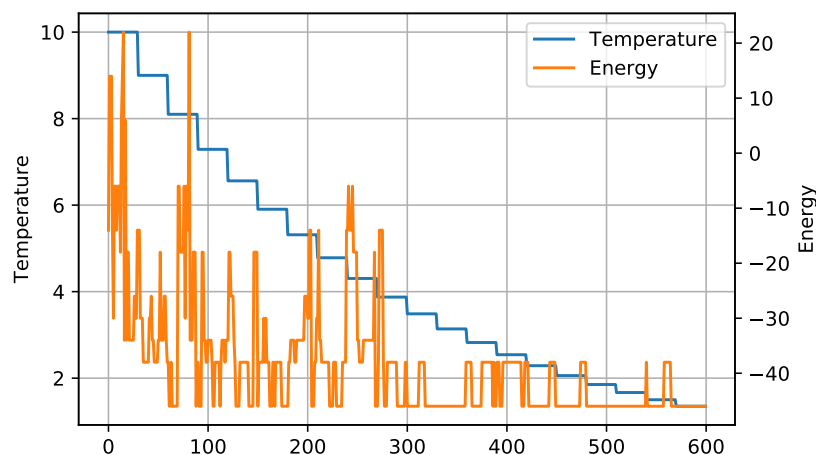


Figure 17: Simulated annealing with $T(1) = 10$ and $c = 0.9$. A total of 20 main iterations (with fixed temperature), and $5n$ sub iterations (with the same temperature) gives $100n = 600$ iterations.

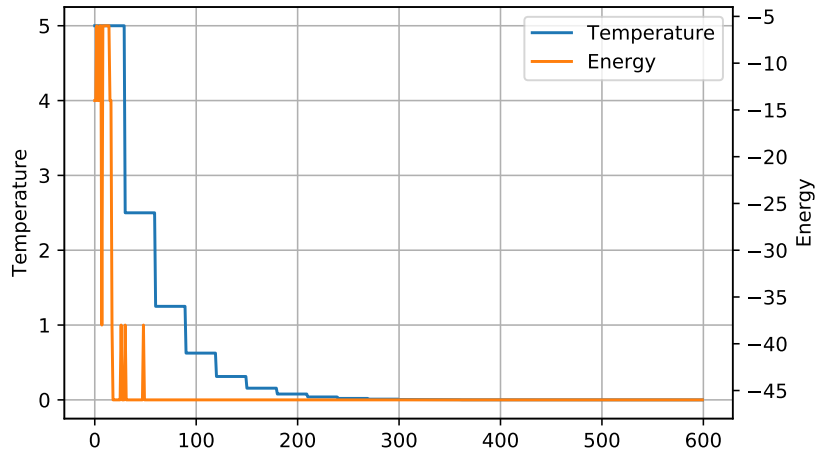


Figure 18: Simulated annealing with $T(1) = 5$ and $c = 0.5$. A total of 20 main iterations (with fixed temperature), and $5n$ sub iterations (with the same temperature) gives $100n = 600$ iterations.

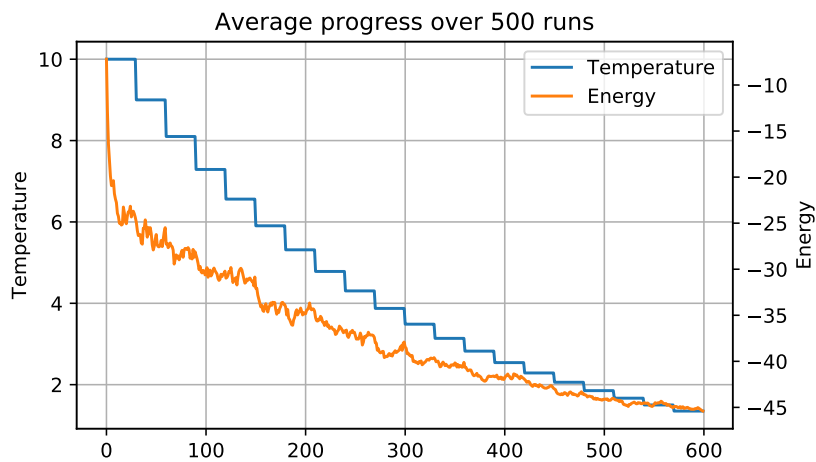


Figure 19: Simulated annealing over 500 runs with $T(1) = 10$ and $c = 0.9$. A total of 20 main iterations (with fixed temperature), and $5n$ sub iterations (with the same temperature) gives $100n = 600$ iterations.

2.8 Nonmetric methods

Problem 8.10

- a) One of the defining properties of an impurity function is that it should be zero when only ω_1 is present, or when only ω_2 is present. For a polynomial to incorporate this information, it must have two zero points, and therefore it must at least be quadratic.
- b) The simplest quadratic form of $P(\omega_1)$ obeying the boundary conditions is

$$i(P(\omega_1)) \propto P(\omega_1) [1 - P(\omega_1)] = P(\omega_1)P(\omega_2).$$

- c) Let X be a Bernoulli variable, with $X \in \{0, 1\} = \{\omega_2, \omega_1\} = C$ and $P(X = 1) = P(\omega_1)$. In the given problem, we do not know the true probabilities—but they may be estimated from the fractions. The variance of X is given by

$$\begin{aligned} \text{var}[X] &= \sum_{\omega_j \in C} P(\omega_j) [\omega_j - \mu]^2 \\ &= P(\omega_1) [\omega_1 - \mu]^2 + P(\omega_2) [\omega_2 - \mu]^2 \\ &= P(\omega_1) [1 - \mu]^2 + (1 - P(\omega_1)) [0 - \mu]^2 \\ &= P(\omega_1) [1 - P(\omega_1)]^2 + (1 - P(\omega_1)) P(\omega_1)^2 \\ &= [1 - P(\omega_1)] [P(\omega_1) (1 - P(\omega_1)) + P(\omega_1)^2] \\ &= (1 - P(\omega_1)) P(\omega_1). \end{aligned}$$

In other words, sample variance is proportional to the impurity estimate defined in subproblem b). If the variance is high, the data is impure.

Problem 8.14

We generalize the problem of a single missing attribute in a single training point to several missing attributes, and to several deficient training points.

In this problem, a pattern $\mathbf{x} = (x_1, x_2, \dots, x_d)$ with d attributes can have between 0 and $d - 1$ missing attributes. Furthermore, the data set $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ may have many training points with missing attributes. The data might look something like what is presented in Table 2 on page 59.

High level Python-like pseudocode is given below. The principal difference between the code below and the code for the same problem without missing attributes, is that there might now be $n - m_i - 1$ possible splits, where m_i are the number of data points missing from attribute i . This is in contrast to the original algorithm, where there are in general $n - 1$ possible splits for every attribute i .

```
for attribute=1 to d do:
    attribute_data = data[attribute, :]
    num_missing = count_missing(attribute_data)
```

```

best_split = None
maximal_impurity_gain = -inf

for possible_split=1 to (n - 1 - num_missing) do:
    impurity_gain = compute_gain(attribute_data, possible_split)

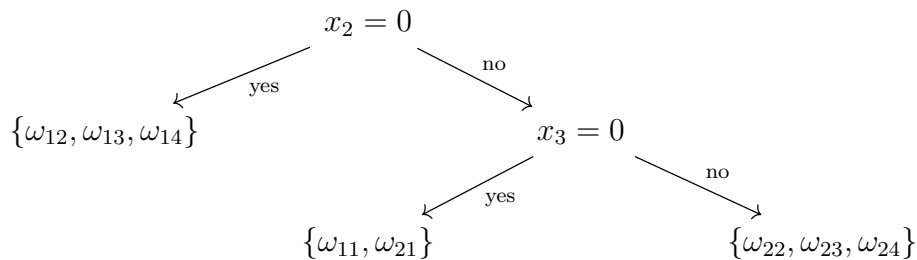
    if impurity_gain > maximal_impurity_gain:
        best_split = (attribute, possible_split)
        maximal_impurity_gain = impurity_gain

split_along_maximal_impurity_gain(best_split)

```

Problem 8.16

- a) The following tree was constructed by hand. Recall that YES is to the left, and NO is to the right. The notation ω_{ij} signifies a pattern j from category i .



The leftmost leaf node assigns to ω_1 with probability 1, the middle leaf node assigns to either ω_1 or ω_2 with probability 0.5, and the rightmost node assigns to ω_2 with probability 1.

- b) The resulting tree would be identical to the one given in subproblem a), but the assignment from the middle node should be ω_1 with probability $2/3$.

Table 2: Example of what data for problem 8.14 might look like.

x_1	x_2	...	x_d
0	\emptyset	...	1
\emptyset	\emptyset	...	1
1	\emptyset	...	\emptyset
\vdots	\vdots	\vdots	\vdots
0	\emptyset	...	1

Problem 8.18

We will consider every alignment of \mathbf{x} and $text$. For each of these alignments, there will be between 1 and m comparisons before termination, depending on whether matches are found or not. There is a probability d^{-1} of a match for each comparison. We will see that expected number of comparisons may be expressed as a *arithmetico-geometric* sum, i.e. a sum of terms which are the product of an arithmetic and geometric sum.

We start out by observing that there are $n - m + 1$ alignments of \mathbf{x} and $text$. For every alignment, at most m characters must be compared. Let $1 \leq \ell \leq m$ be the number of comparisons made before the loop terminates. The loop terminates when a comparison yields no match. Fixing \mathbf{x} and considering every character of $text$ to be randomly drawn from an alphabet \mathcal{A} of d characters, there is a probability of d^{-1} of a match in each comparison.

Let us now study $P(\ell)$, the probability of ℓ comparisons before termination of the loop. For instance $P(\ell) = 1 - 1/d$, since it represents the probability of one comparison before exiting the loop, and this happens if there is no match on the first comparison. More generally, we observe that

$$\begin{aligned}
 P(\ell = 1) &= 1 - \frac{1}{d} && \text{(no match on the first)} \\
 P(\ell = 2) &= \left(1 - \frac{1}{d}\right) \frac{1}{d} && \text{(match, then no match)} \\
 P(\ell = 2) &= \left(1 - \frac{1}{d}\right) \frac{1}{d^2} && \text{(match, match, then no match)} \\
 &\vdots = \vdots && \\
 P(\ell = m) &= \left(1 - \frac{1}{d}\right) \frac{1}{d^{m-1}} + \frac{1}{d^m} && \text{(no match on final, or match on final)}
 \end{aligned}$$

We now take the expected value of ℓ to obtain the expected number of comparisons, and introduce a variable $r := 1/d$ for notational convenience.

$$\begin{aligned}
 \mathbb{E}[\ell] &= \sum_{k=1}^m P(\ell = k) k = \left(1 - \frac{1}{d}\right) \left(1 + \frac{2}{d} + \frac{3}{d^2} + \cdots + \frac{m}{d^{m-1}}\right) + \frac{m}{d^m} \\
 &= (1 - r) (1 + 2r + 3r^2 + \cdots + mr^{m-1}) + mr^m
 \end{aligned}$$

The sum in the second parenthesis in the first term is a arithmetico-geometric series.

Using the summation formula³ we obtain

$$\begin{aligned}\mathbb{E}[\ell] &= \sum_{k=1}^m P(\ell = k) k = (1-r)(1+2r+3r^2+\dots+mr^{m-1}) + mr^m \\ &= (1-r) \left(\frac{1-(1-m)r^m}{1-r} + \frac{r(1-r^m)}{(1-r)^2} \right) + mr^m \\ &= \frac{(1-r) - (1-r)(1+m)r^m + r(1-r^m) + (1-r)mr^m}{1-r}.\end{aligned}$$

After some algebraic operations on this final fraction, we obtain the desired result

$$\mathbb{E}[\ell] = \sum_{k=1}^m P(\ell = k) k = \frac{1-r^m}{1-r} = \frac{1-d^{-m}}{1-d^{-1}}.$$

The problem is solved, since the expected number of comparisons is

$$\text{alignments} \times \text{comparisons} = (n-m+1) \times \frac{1-d^{-m}}{1-d^{-1}}.$$

The inequality $(1-d^{-m})/(1-d^{-1}) \leq 2$ presented in the problem description stems from the fact that the fraction is the sum of a geometric series

$$\frac{1-d^{-m}}{1-d^{-1}} = 1 + \frac{1}{d} + \frac{1}{d^2} + \dots + \frac{1}{d^{m-1}}.$$

When $d = 2$ above, the sum equals 2. Clearly, when $d > 2$ the terms become smaller, and the sum must therefore be < 2 . Therefore, as long as $d \geq 2$, the sum is ≤ 2 . This proves the inequality when $d \geq 2$, which is a reasonable assumption considering the problem.

Problem 8.22

A naive algorithm which loops through the alphabet (of size d) for every letter in \mathbf{x} (of size m) would give an $O(dm)$ algorithm. This is not very efficient, and faster algorithms may easily be constructed.

We can construct a **better algorithm** with runtime $O(d+m)$ as follows.

1. Construct an empty lookup table for every character in the alphabet
2. For every character and position in \mathbf{x} , overwrite the lookup table

Below is actual Python code implementing the pseudocode above.

```
def last_occurrence(alphabet, word):
    """
    Returns a mapping F[char] -> last occurrence in word.
    """
    mapping = {i:-1 for i in alphabet} # This is O(d)
```

³See Wikipedia for the summation formula of an arithmetico-geometric series.

```

for index, char in enumerate(word): # This is O(m)
    mapping[char] = index

return mapping

```

It's possible to make the algorithm closer to $O(\min(m, d))$, if we know the values of m and d before the algorithm starts. If $d \gg m$ or $m \gg d$, such an algorithm can avoid looping over the longest string, but we will not pursue this any further in this problem.

We are now in a position to answer the questions posed in the problem.

- a) The time complexity of the algorithm is $O(m + d)$.
- b) The space complexity is $O(d)$ in the algorithm above. We can avoid storing the empty matches (for instance by using the `collections.defaultdict` data structure in Python) and get the storage down to $O(\text{unique_chars}(m))$.
- c) Running the algorithm above requires $d+m$ operations, so that $\mathbf{x} = \text{"bonbon"}$ yields $26 + 6$ operations. Since $m < d$ here, we can use a `collections.defaultdict` data structure in Python and bring it down to $m = 6$ operations.

Below is a print showing that `collections.defaultdict` indeed leads to faster runtime in practice. The `collections.defaultdict` data structure is a dictionary (hash-table) which calls a *default factory function* if a value of a given key is not found.

```

%timeit last_occurrence(alphabet, word)
2.32 microseconds +/- 42 ns per loop
%timeit last_occurrence_defaultdict(alphabet, word)
1.17 microseconds +/- 3.17 ns per loop

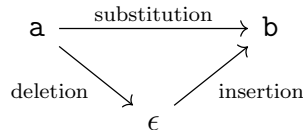
```

Problem 8.26

- a) The properties of a metric $D(\mathbf{a}, \mathbf{b})$ are (i) Non-negativity; (ii) Reflexivity; (iii) Symmetry; and (iv) Triangle inequality. These properties are listed on page 187 in [Duda et al., 2000], which corresponds to Section 4.6.1. In order, we see that
 - (i) Non-negativity is always obeyed, since every cost is positive.
 - (ii) Reflexivity is always obeyed, since if $\mathbf{a} = \mathbf{b}$ no operation is performed and the total cost of transformation is zero.
 - (iii) Symmetry is not always obeyed. A counterexample follows.
 - (iv) The triangle inequality is not always obeyed. A counterexample follows.
- b) We provide counterexamples for symmetry and the triangle inequality.

Counterexample for symmetry Let $\mathbf{a} = \text{dog}$ and $\mathbf{b} = \text{dogs}$. Define the distance $D(\mathbf{a}, \mathbf{b})$ as the cost of transforming \mathbf{a} to \mathbf{b} . Then $D(\mathbf{a}, \mathbf{b})$ is the distance of a single insertion to transform $\text{dog} \mapsto \text{dogs}$, while $D(\mathbf{b}, \mathbf{a})$ is the distance of a single deletion $\text{dogs} \mapsto \text{dog}$. If the cost of insertion and deletion differ, then $D(\mathbf{a}, \mathbf{b}) \neq D(\mathbf{b}, \mathbf{a})$. This provides a counterexample of the symmetry property of a metric when the costs of operations may be different.

Counterexample for the triangle inequality Let ϵ be the empty string. Consider the transformation shown in the diagram below. If the cost of substitution is much greater than the cost of deletion and insertion, then $D(\mathbf{a}, \epsilon) + D(\epsilon, \mathbf{b}) \ll D(\mathbf{a}, \mathbf{b})$. This provides a counterexample of the triangle inequality of a metric.

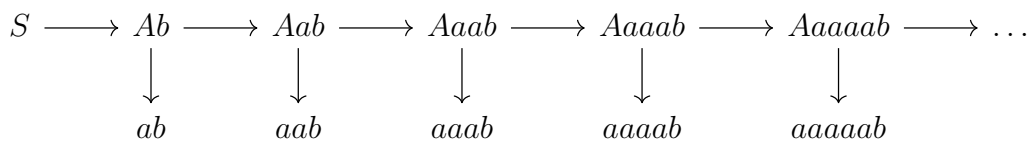


Problem 8.31

a) A grammar generating $\mathcal{L}(G) = \{\mathbf{a}^n\mathbf{b} \mid n \geq 1\}$ is, for instance, given by

$$\begin{aligned} \mathcal{A} &= \{a, b\} & \mathcal{S} &= S & \mathcal{I} &= \{A\} \\ \mathcal{P} &= \{S \rightarrow Ab, A \rightarrow a, A \rightarrow Aa\}. \end{aligned}$$

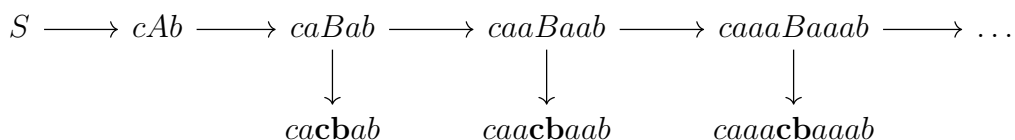
b) Below is a derivation tree showing how ab and $aaaaab$ are generated.



Problem 8.32

a) Since **Type 3** \subset **Type 2** \subset **Type 1**, we check if it's a **Type 3** (regular) grammar first. It's not a regular grammar, since the rule $B \rightarrow aBa$ is not of the form $\alpha \rightarrow z\beta$ or of the form $\alpha \rightarrow z$. It is, however, a **Type 2** (context-free) grammar since every rule is of the form $I \rightarrow \gamma$. In other words, every rewrite rule is from an intermediate symbol I to a string γ made up of intermediate or terminal symbols.

b) To see that the grammar generates the language $\mathcal{L}(G) = \{\mathbf{ca}^n\mathbf{cba}^n\mathbf{b} \mid n \geq 1\}$, we draw a derivation tree with intermediate states in the first row, and final states in the second row. The rule $B \rightarrow \mathbf{cb}$ is highlighted in boldface.



Clearly moving down from $caBab$ produces $\{\mathbf{ca}^n\mathbf{cba}^n\mathbf{b} \mid n = 1\}$, and moving ℓ to the right from $caBab$, and then down, produces $\{\mathbf{ca}^n\mathbf{cba}^n\mathbf{b} \mid n = \ell + 1\}$.

c) See the solution to subproblem b) above for derivation trees.

2.9 Algorithm-independent machine learning

Problem 9.9

a) The binomial theorem is

$$(x + y)^n = \sum_{r=0}^n \binom{n}{r} x^r y^{n-r} = x^n + nx^{n-1}y + \cdots + nxy^{n-1} + y^n,$$

and when we substitute $x = y = 1$ the result immediately follows, since

$$2^n = (1 + 1)^n = \sum_{r=0}^n \binom{n}{r} 1^r 1^{n-r} = \sum_{r=0}^n \binom{n}{r}.$$

b) The *base case* of the inductive argument is straightforward, we have

$$K(1) = \sum_{r=0}^1 \binom{1}{r} = \binom{1}{0} + \binom{1}{1} = 1 + 1 = 2.$$

To solve the problem of the *inductive step*, we will need to use following three observations:

$$\binom{n}{r} = \binom{n-1}{r} + \binom{n-1}{r-1} \quad \binom{n}{n} = 1 \quad \binom{n}{0} = 1$$

The first one might be familiar from the construction of Pascals triangle, and the final two are easily derived when writing out $\binom{n}{r} := n!/(n-r)!r!$.

The strategy is to write out $K(n+1)$, strip off the first and last term in the sum, apply the first identity above, and then collect terms to obtain $K(n)$ twice. We proceed:

$$\begin{aligned} K(n+1) &= \sum_{r=0}^{n+1} \binom{n+1}{r} \\ &= \binom{n+1}{0} + \sum_{r=1}^n \left[\binom{n+1}{r} \right] + \binom{n+1}{n+1} \\ &= 1 + \sum_{r=1}^n \left[\binom{n}{r} + \binom{n}{r-1} \right] + 1 \\ &= \left[1 + \sum_{r=1}^n \binom{n}{r} \right] + \left[\sum_{r=1}^n \binom{n}{r-1} + 1 \right] \\ &= \sum_{r=0}^n \binom{n}{r} + \sum_{r=0}^n \binom{n}{r} = K(n) + K(n) = 2K(n) \end{aligned}$$

If $K(n+1) = 2K(n)$ and $K(1) = 2$, then $K(n) = 2^n$ for all n , which is what we wanted to prove.

Problem 9.13

- a) The n length string $010110111011110\dots$ has the pattern $01-011-0111-01111$, which may be generated using a for-loop. To generate n bits, we need a counter c keeping track of the number of 1s. For instance, to generate the $2 + 3 + 4 + 5 = 14$ bit string $01-011-0111-01111$ the counter must go to 4. Similarly, to generate the $2 + 3 + 4 + 5 + 6 = 20$ bit string $01-011-0111-01111-011111$ the counter must go to 5.

More generally, to generate n the counter c must go to approximately

$$1 + 2 + 3 + \dots + c = n = \frac{c(c-1)}{2}.$$

Therefore, we need a counter going to $\sim \sqrt{n}$ to generate a n bit string. Storing the number c requires $\log_2 c$ bits, so the complexity is $O(\log_2 \sqrt{n})$.

- b) We must iterate over $\sim n/2$ zeros, then output a 1, then iterate over $\sim n/2$ zeros again. To do this, we need a counter going to $n/2$, storing this requires $\log_2(n/2)$ bits, so the Kolmogorov complexity is $O(\log_2 n)$.
- c) The complexity is $O(1)$, since there exists programs of finite size which will yield any arbitrarily large number of consecutive digits of e . This is the same reason as why π is $O(1)$, as explained in Section 9.2.3 in [Duda et al., 2000].
- d) Also constant, for the same reason as the previous sub-problem.
- e) Assuming a finite size program which generates the digits of π one-by-one exists, we use a counter to substitute every 100th digit of π in the base-2 number system to a 1. The counter resets after hitting 100, so its storage is constant. The total complexity is $O(1)$, i.e. constant.
- f) Same as in the previous question, but storing the counter requires $\log_2 n$ bits. Therefore the complexity is $O(\log_2 n)$.

Problem 9.18

We use the relation $\text{var}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$, solve for $\mathbb{E}[X^2]$ and substitute $X \mapsto g(\mathbf{x}; \mathcal{D}) - F(\mathbf{x})$. Doing this, we observe that

$$\begin{aligned} \mathbb{E}[X^2] &= \mathbb{E}[X]^2 + \text{var}[X] \\ \mathbb{E}[(g(\mathbf{x}; \mathcal{D}) - F(\mathbf{x}))^2] &= \underbrace{\mathbb{E}[g(\mathbf{x}; \mathcal{D}) - F(\mathbf{x})]^2}_{\text{bias}^2} + \underbrace{\text{var}[g(\mathbf{x}; \mathcal{D}) - F(\mathbf{x})]}_{\text{variance}}. \end{aligned}$$

The last term for the variance is not immediately recognizable as the term given in the book, but we note that

$$\begin{aligned} \text{var}[g(\mathbf{x}; \mathcal{D}) - F(\mathbf{x})] &= \mathbb{E}[(g(\mathbf{x}; \mathcal{D}) - F(\mathbf{x}) - \mathbb{E}[g(\mathbf{x}; \mathcal{D}) - F(\mathbf{x})])^2] \\ &= \mathbb{E}[(g(\mathbf{x}; \mathcal{D}) - \mathbb{E}[g(\mathbf{x}; \mathcal{D})])^2] \end{aligned}$$

since $\mathbb{E}[F(\mathbf{x})] = F(\mathbf{x})$. In other words, $F(\mathbf{x})$ is the true function and it's expected value is itself. Combining the results above, we have

$$\mathbb{E}[(g(\mathbf{x}; \mathcal{D}) - F(\mathbf{x}))^2] = \underbrace{\mathbb{E}[g(\mathbf{x}; \mathcal{D}) - F(\mathbf{x})]^2}_{\text{bias}^2} + \underbrace{\mathbb{E}[(g(\mathbf{x}; \mathcal{D}) - \mathbb{E}[g(\mathbf{x}; \mathcal{D})])^2]}_{\text{variance}}$$

as stated.

The *bias* can be negative, for instance if the estimate is much lower in value than the true function. However, the *squared bias* is always non-negative. The *variance* is non-negative, since

$$\text{var}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2]$$

by definition, and the squared term is always non-negative.

Problem 9.23

We prove that the average of the leave-one-out means $\mu_{(\cdot)}$ is equivalent to the sample mean $\hat{\mu} = \bar{x}$ by substitution. Starting with the definition of $\mu_{(\cdot)}$, we have

$$\begin{aligned} \mu_{(\cdot)} &= \frac{1}{n} \sum_{i=1}^n \mu_{(i)} = \frac{1}{n} \sum_{i=1}^n \left[\frac{n\bar{x} - x_i}{n-1} \right] = \frac{1}{n(n-1)} \sum_{i=1}^n [n\bar{x} - x_i] \\ &= \frac{1}{n(n-1)} [n^2\bar{x} - n\bar{x}] = \frac{n\bar{x} - \bar{x}}{n-1} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i. \end{aligned}$$

Problem 9.26

There is an error in Equation (23) in chapter 9 in the book. The equation cannot possibly measure the variance (the uncertainty) of a mean value. To see this, consider two data sets and the results of applying Equation (23):

$$\mathcal{D}_1 = \{0, 0, 1, 1\} \mapsto_{(23)} 0.75$$

$$\mathcal{D}_2 = \{0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1\} \mapsto_{(23)} 4.75$$

Clearly we are more certain about the mean of \mathcal{D}_2 than \mathcal{D}_1 , but this is not reflected in the estimate of the variance of the mean. Equation (23) should be

$$\hat{\sigma}^2 = \frac{1}{n(n-1)} \sum_{i=1}^n (x_i - \hat{\mu})^2,$$

since

$$\text{var}[\bar{X}] = \text{var}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n^2} \sum_{i=1}^n \text{var}[X_i] = \frac{1}{n} \text{var}[X] = \frac{1}{n(n-1)} \sum_{i=1}^n (x_i - \hat{\mu})^2.$$

Then $\mathcal{D}_1 \mapsto 0.08333$ and $\mathcal{D}_2 \mapsto 0.0131578$, which is more in line with what we'd expect—more observations reduce the uncertainty.

Onto the problem at hand: to show that Equation (26) reduces to the above, we write

$$\mu_{(i)} - \mu_{(\cdot)} = \left[\frac{n\bar{x} - x_i}{n-1} \right] - \bar{x} = \frac{\bar{x} - x_i}{n-1}.$$

Squaring this and pulling $(n-1)^2$ out of the sum yields the relationship

$$\text{var}[\hat{\mu}] = \frac{n-1}{n} \sum_{i=1}^n (\mu_{(i)} - \mu_{(\cdot)})^2 = \frac{n-1}{n} \sum_{i=1}^n \left(\frac{\bar{x} - x_i}{n-1} \right)^2 = \frac{1}{n(n-1)} \sum_{i=1}^n (\bar{x} - x_i)^2$$

as desired.

Problem 9.34

See Figure 20 on page 68 for a plot accompanying this solution.

- a) We find the maximum likelihood estimate for x_ℓ and x_u to be

$$\hat{x}_\ell = \min \mathcal{D} = 0.2 \quad \hat{x}_u = \max \mathcal{D} = 0.9.$$

Note that these maximum likelihood estimates are biased. The value of \hat{x}_ℓ will be greater than the true value of x_ℓ , and the value of \hat{x}_u will be smaller than the true value of x_u . The unbiased estimates are

$$\begin{aligned} \hat{x}_u &= \min \mathcal{D} + \frac{n}{n-1} (\max \mathcal{D} - \min \mathcal{D}) \\ \hat{x}_\ell &= \min \mathcal{D} - \frac{1}{n-1} (\max \mathcal{D} - \min \mathcal{D}) \end{aligned}$$

- b) The estimates are $\hat{\mu} = 0.514$ and $\hat{\sigma} = 0.223$, and the maximum likelihood estimate of σ is biased.
- c) The problem states that there is no prior reason to assume one model over the other, so $P(h_1) = P(h_2) = 0.5$. We will not make any assumptions on the range of parameter values in θ , so the term $p(\hat{\theta} | h_i) \Delta \theta$ in Equation (42) in the book is set to 1. What's left is to evaluate $\ln P(\mathcal{D} | \hat{\theta}, h_i)$, which becomes

$$\begin{aligned} \ln P(\mathcal{D} | \hat{\theta}, h_1) &= \sum_{i=1}^n \ln p(x_i | \hat{x}_\ell = 0.2, \hat{x}_u = 0.9) = \left(\frac{1}{0.7} \right)^7 \approx 2.497 \\ \ln P(\mathcal{D} | \hat{\theta}, h_2) &= \sum_{i=1}^n \ln p(x_i | \hat{\mu} = 0.514, \hat{\sigma} = 0.223) \approx 0.567 \end{aligned}$$

We decide on model h_1 , since it's associated with much greater likelihood.

Problem 9.37

The plot for this problem found in Figure 21 on page 68. To create the plot, and outer loop for values of n was used. For each value of n , a true value of p was set. Then given (n, p) , a range of values for a 95% confidence interval was computed. These values were stored and plotted. Notice that the graph is somewhat jagged when n is small because of the discrete structure of the problem.

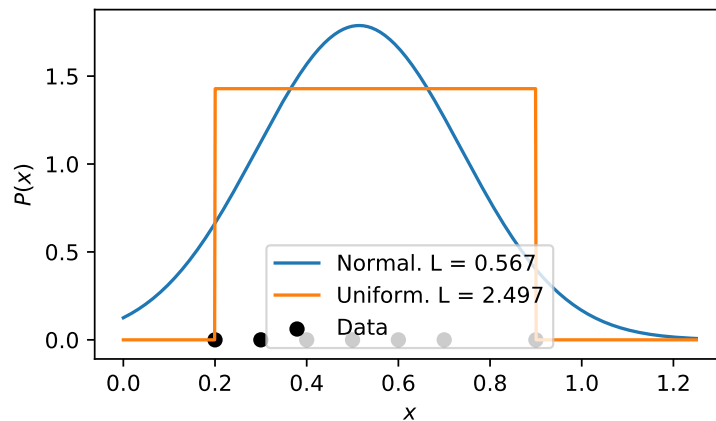


Figure 20: Plot of data and models for Problem 9.34.

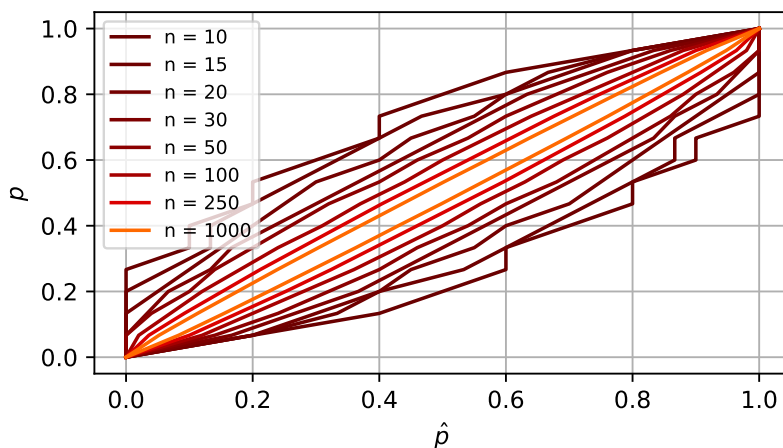


Figure 21: Plot associated with Problem 9.37.

Problem 9.40

We wish to maximize the likelihood $P(\mathcal{D} | \theta) = P(k | p)$, where the model is given by the binomial distribution $P(k | p) = \binom{n}{k} p^k (1 - p)^{n-k}$. The log-likelihood becomes

$$\ln P(k | p) = \ln \binom{n}{k} + k \ln p + (n - k) \ln(1 - p),$$

and differentiating with respect to p and setting this to zero yields

$$\frac{\partial \ln P(k | p)}{\partial p} = \frac{k}{p} + \frac{n - k}{p - 1} = 0.$$

The maximum likelihood is the \hat{p} such that the above expression is zero, and solving for p yields $\hat{p} = k/n$ as stated in the problem text.

2.10 Unsupervised learning and clustering

Problem 10.7

The log-likelihood of a mixture is given by the expression

$$\ell(\alpha) = \sum_{k=1}^n \ln p(\mathbf{x}_k | \alpha) = \sum_{k=1}^n \ln \left(\sum_{j=1}^c p(\mathbf{x}_k | \omega_j, \alpha) P(\omega_j) \right).$$

We'll use the first expression above, differentiate with respect to the parameter α , and then expand the sum and differentiate through it with respect to α .

Executing the plan outlined above, we observe that

$$\begin{aligned} \frac{\partial \ell(\alpha)}{\partial \alpha} &= \sum_{k=1}^n \frac{1}{p(\mathbf{x}_k | \alpha)} \frac{\partial}{\partial \alpha} p(\mathbf{x}_k | \alpha) \\ &= \sum_{k=1}^n \frac{1}{p(\mathbf{x}_k | \alpha)} \frac{\partial}{\partial \alpha} \left(\sum_{j=1}^c p(\mathbf{x}_k | \omega_j, \alpha) P(\omega_j) \right) && \text{(by definition)} \\ &= \sum_{k=1}^n \frac{1}{p(\mathbf{x}_k | \alpha)} \sum_{j=1}^c P(\omega_j) \frac{\partial p(\mathbf{x}_k | \omega_j, \alpha)}{\partial \alpha} \\ &= \sum_{k=1}^n \sum_{j=1}^c \frac{P(\omega_j)}{p(\mathbf{x}_k | \alpha)} \frac{\partial p(\mathbf{x}_k | \omega_j, \alpha)}{\partial \alpha}. && \text{(constant into sum)} \end{aligned}$$

Since $(\ln x)' = x'/x$ by the chain rule of calculus, the equation $x' = x(\ln x)'$ holds too. Using this on the last factor, we observe that

$$\sum_{k=1}^n \sum_{j=1}^c \frac{P(\omega_j)}{p(\mathbf{x}_k | \alpha)} \frac{\partial p(\mathbf{x}_k | \omega_j, \alpha)}{\partial \alpha} = \sum_{k=1}^n \sum_{j=1}^c \frac{P(\omega_j)}{p(\mathbf{x}_k | \alpha)} p(\mathbf{x}_k | \omega_j, \alpha) \frac{\partial \ln p(\mathbf{x}_k | \omega_j, \alpha)}{\partial \alpha}.$$

As stated in the problem text, the first two factors may be simplified using Bayes theorem, since

$$P(\omega_j | \mathbf{x}_k, \alpha) = \frac{p(\mathbf{x}_k | \omega_j, \alpha) P(\omega_j | \alpha)}{p(\mathbf{x}_k, | \alpha)} = \frac{p(\mathbf{x}_k | \omega_j, \alpha) P(\omega_j)}{p(\mathbf{x}_k, | \alpha)}$$

when ω_j is independent of α , i.e. $P(\omega_j | \alpha) = P(\omega_j)$. We have completed the problem, since applying the simplification due to Bayes theorem yields the desired result.

Problem 10.11

- a) Solving this problem requires knowledge of matrix calculus. We will differentiate the second term, and then the first term, and finally combine the results. The reason we start with the second term is that it is easier.

Second term Let's start by differentiating $\mathbf{x}^T \mathbf{A} \mathbf{x}$ with respect to an element A_{ij} . This is functionally equivalent to the second term in the problem statement, and solving this solves the original problem. If \mathbf{A} is symmetric, then

$$\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial A_{ij}} = \frac{\partial \sum_i \sum_j x_i A_{ij} x_j}{\partial A_{ij}} = \begin{cases} x_i x_j & \text{if } i = j \\ 2x_i x_j & \text{if } i \neq j \end{cases} = (2 - \delta_{ij}) x_i x_j.$$

First term Let's now consider differentiating $f(\mathbf{A}) = \ln(\det \mathbf{A})$ with respect to an element A_{ij} . This is functionally equivalent to the first term in the problem statement. We'll make use of Jacobi's formula⁴, which is given by

$$\frac{d}{dt} \det \mathbf{A}(t) = \text{tr} \left(\text{adj}(\mathbf{A}(t)) \frac{d\mathbf{A}(t)}{dt} \right),$$

where $\text{tr}(\mathbf{A})$ is the trace of \mathbf{A} , and $\text{adj}(\mathbf{A})$ is the adjugate of \mathbf{A} .

We differentiate using the chain rule and apply Jacobi's formula to obtain

$$\frac{\partial \ln(\det \mathbf{A})}{\partial A_{ij}} = \frac{1}{\det \mathbf{A}} \frac{\partial \det \mathbf{A}}{\partial A_{ij}} = \frac{1}{\det \mathbf{A}} \text{tr} \left(\text{adj}(\mathbf{A}) \frac{\partial \mathbf{A}}{\partial A_{ij}} \right).$$

Now comes the part where we consider the diagonal elements $i = j$ and off-diagonal elements $i \neq j$ separately. We carry out the differentiation $d\mathbf{A}/\partial A_{ij}$ and obtain

$$\frac{1}{\det \mathbf{A}} \text{tr} \left(\text{adj}(\mathbf{A}) \frac{\partial \mathbf{A}}{\partial A_{ij}} \right) = \begin{cases} \frac{1}{\det \mathbf{A}} \text{tr}(\text{adj}(\mathbf{A}) [\mathbf{E}_{ii}]) & \text{if } i = j \\ \frac{1}{\det \mathbf{A}} \text{tr}(\text{adj}(\mathbf{A}) [\mathbf{E}_{ij} + \mathbf{E}_{ji}]) & \text{if } i \neq j \end{cases},$$

where by \mathbf{E}_{ij} we mean a matrix of zeros everywhere except in the ij th position. We do not show it here, but by the definition of \mathbf{E}_{ij} and the trace, it should not be too hard to see that $\text{tr}(\mathbf{A} \mathbf{E}_{ij}) = A_{ji}$. Using the above in combination with the fact that $\text{adj}(\mathbf{A})_{ji} / \det(\mathbf{A}) = \mathbf{A}_{ji}^{-1}$ and the fact that the inverse of a symmetric matrix is symmetric, we observe that the above expression may be written using the Kronecker delta symbol δ_{ij} as

$$\begin{cases} \frac{\text{adj}(\mathbf{A})_{ii}}{\det \mathbf{A}} & \text{if } i = j \\ \frac{\text{adj}(\mathbf{A})_{ji}}{\det \mathbf{A}} + \frac{\text{adj}(\mathbf{A})_{ij}}{\det \mathbf{A}} & \text{if } i \neq j \end{cases} = \begin{cases} \mathbf{A}_{ii}^{-1} & \text{if } i = j \\ \mathbf{A}_{ji}^{-1} + \mathbf{A}_{ij}^{-1} & \text{if } i \neq j \end{cases} = (2 - \delta_{ij}) \mathbf{A}_{ij}^{-1}.$$

Putting it all together We have established that if \mathbf{A} is a symmetric matrix, then

$$\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial A_{ij}} = (2 - \delta_{ij}) x_i x_j \quad \text{and} \quad \frac{\partial \ln(\det \mathbf{A})}{\partial A_{ij}} = (2 - \delta_{ij}) \mathbf{A}_{ij}^{-1}.$$

⁴See Wikipedia for details: https://en.wikipedia.org/wiki/Jacobi's_formula

Applying the two results above to the problem at hand yields

$$\begin{aligned}
\frac{\partial}{\partial \sigma^{pq}(i)} (\ln p(\mathbf{x}_k | \omega_i, \boldsymbol{\theta}_i)) &= \frac{\partial}{\partial \sigma^{pq}(i)} \left(\ln \frac{|\boldsymbol{\Sigma}_i^{-1}|^{1/2}}{(2\pi)^{d/2}} - \frac{1}{2} (\mathbf{x}_k - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{x}_k - \boldsymbol{\mu}_i) \right) \\
&= \frac{\partial}{\partial \sigma^{pq}(i)} \left(\frac{1}{2} \ln |\boldsymbol{\Sigma}_i^{-1}| - \frac{1}{2} (\mathbf{x}_k - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{x}_k - \boldsymbol{\mu}_i) \right) \\
&= \left(1 - \frac{\delta_{pq}}{2} \right) \sigma_{pq}(i) - \left(1 - \frac{\delta_{pq}}{2} \right) (x_p(k) - \mu_p(i)) (x_q(k) - \mu_q(i)) \\
&= \left(1 - \frac{\delta_{pq}}{2} \right) [\sigma_{pq}(i) - (x_p(k) - \mu_p(i)) (x_q(k) - \mu_q(i))],
\end{aligned}$$

which corresponds with Equation (23) in chapter 10 in the book as expected.

- b) We will extend the result from the previous sub-problem to every entry in the matrix, then solve for the maximum likelihood estimate for $\boldsymbol{\Sigma}$.

If we extend the differentiation formulas from the previous sub-problem to every entry in the matrix, we recognize that

$$\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{A}} = \mathbf{x} \mathbf{x}^T \quad \text{and} \quad \frac{\partial \ln(\det \mathbf{A})}{\partial \mathbf{A}} = \mathbf{A}^{-1}.$$

Applying these equations to the second factor in the solution to problem 7, and ignoring constants, we obtain a result proportional to

$$\frac{\partial}{\partial \boldsymbol{\Sigma}^{-1}} (\ln p(\mathbf{x}_k | \omega_i, \boldsymbol{\theta}_i)) = \boldsymbol{\Sigma} - (\mathbf{x}_k - \boldsymbol{\mu}_i) (\mathbf{x}_k - \boldsymbol{\mu}_i)^T.$$

Our problem becomes that of solving the equation

$$\begin{aligned}
\frac{\partial \ell}{\partial \boldsymbol{\Sigma}^{-1}} &= \sum_{k=1}^n \sum_{i=1}^c P(\omega_i | \mathbf{x}_k, \boldsymbol{\theta}) \left[\boldsymbol{\Sigma} - (\mathbf{x}_k - \boldsymbol{\mu}_i) (\mathbf{x}_k - \boldsymbol{\mu}_i)^T \right] = \mathbf{0} \\
\Rightarrow \boldsymbol{\Sigma} \sum_{k=1}^n \sum_{i=1}^c P(\omega_i | \mathbf{x}_k, \boldsymbol{\theta}) &= \sum_{k=1}^n \sum_{i=1}^c P(\omega_i | \mathbf{x}_k, \boldsymbol{\theta}) \left[\mathbf{x}_k \mathbf{x}_k^T - 2 \mathbf{x}_k \boldsymbol{\mu}_i^T + \boldsymbol{\mu}_i \boldsymbol{\mu}_i^T \right] \quad (10)
\end{aligned}$$

for $\boldsymbol{\Sigma}$. We will now examine each term of Equation (10) separately.

Notice first that the left hand side of equation (10) becomes

$$\boldsymbol{\Sigma} \sum_{k=1}^n \sum_{i=1}^c P(\omega_i | \mathbf{x}_k, \boldsymbol{\theta}) = \boldsymbol{\Sigma} n \sum_{k=1}^n \sum_{i=1}^c p(\mathbf{x}_k, \omega_i | \boldsymbol{\theta}) = \boldsymbol{\Sigma} n,$$

where we have made use of Bayes theorem.

Moving to the right hand side of Equation (10), we see that

$$\sum_{k=1}^n \sum_{i=1}^c P(\omega_i | \mathbf{x}_k, \boldsymbol{\theta}) \mathbf{x}_k \mathbf{x}_k^T = \sum_{k=1}^n \mathbf{x}_k \mathbf{x}_k^T,$$

since the sum of the probabilities of each class c must be 1 for each data point k .

The cross term of the right-hand side of Equation (10) becomes

$$\begin{aligned} -2 \sum_{k=1}^n \sum_{i=1}^c P(\omega_i | \mathbf{x}_k, \boldsymbol{\theta}) \mathbf{x}_k \boldsymbol{\mu}_i^T &= -2 \sum_{k=1}^n \left(\sum_{i=1}^c P(\omega_i | \mathbf{x}_k, \boldsymbol{\theta}) \mathbf{x}_k \right) \boldsymbol{\mu}_i^T \\ &= -2 \sum_{k=1}^n \left(\sum_{i=1}^c n P(\mathbf{x}_k | \omega_i, \boldsymbol{\theta}) P(\omega_i | \boldsymbol{\theta}) \mathbf{x}_k \right) \boldsymbol{\mu}_i^T \\ &= -2 \sum_{i=1}^c n \boldsymbol{\mu}_i P(\omega_i | \boldsymbol{\theta}) \boldsymbol{\mu}_i^T = -2n \sum_{i=1}^c P(\omega_i | \boldsymbol{\theta}) \boldsymbol{\mu}_i \boldsymbol{\mu}_i^T, \end{aligned}$$

where the second equality comes from Bayes theorem. In the the second to last equality we used $\sum_{k=1}^n p(\mathbf{x}_k | \omega_i) \mathbf{x}_k = \boldsymbol{\mu}_i$.

Finally, the rightmost term of the right-hand side of Equation (10) becomes

$$\sum_{k=1}^n \sum_{i=1}^c P(\omega_i | \mathbf{x}_k, \boldsymbol{\theta}) \boldsymbol{\mu}_i \boldsymbol{\mu}_i^T = n \sum_{i=1}^c P(\omega_i | \boldsymbol{\theta}) \boldsymbol{\mu}_i \boldsymbol{\mu}_i^T,$$

where again Bayes theorem was used.

Putting all these results back into Equation (10), we finally see that

$$\boldsymbol{\Sigma} n = \sum_{k=1}^n \mathbf{x}_k \mathbf{x}_k^T - n \sum_{i=1}^c P(\omega_i | \boldsymbol{\theta}) \boldsymbol{\mu}_i \boldsymbol{\mu}_i^T,$$

and dividing through by n gives the desired result.

Problem 10.12

- a) Creating a plot of $p(x | \omega_1) = \mathcal{N}(0, 1)$ and $p(x | \omega_2) = \mathcal{N}(0, 1/2)$ is illuminating. See figure 22 for such a plot, which reveals the significance of $\ln 2$ immediately.

The likelihood of the mixture is given by

$$\ell(x_1 | P(\omega_1)) = p(x_1) = \frac{P(\omega_1)}{\sqrt{2\pi}} e^{-x_1^2/2} + \frac{(1 - P(\omega_1))}{\sqrt{\pi}} e^{-x_1^2},$$

and using the logarithm will not aid us in our attempt to maximize this function.

Notice that for a given value of x_1 , the likelihood is linear in $P(\omega_1)$, since

$$f(P(\omega_1)) = p(x_1) = P(\omega_1) p(x | \omega_1) + (1 - P(\omega_1)) p(x | \omega_2),$$

and it's derivative with respect to $P(\omega_1)$ is $f'(P(\omega_1)) = p(x | \omega_1) - p(x | \omega_2)$. Since $0 \leq P(\omega_1) \leq 1$, there are two cases to consider:

- (i) If $f'(P(\omega_1)) > 0$, the likelihood is maximized by the choice $P(\omega_1) = 1$.

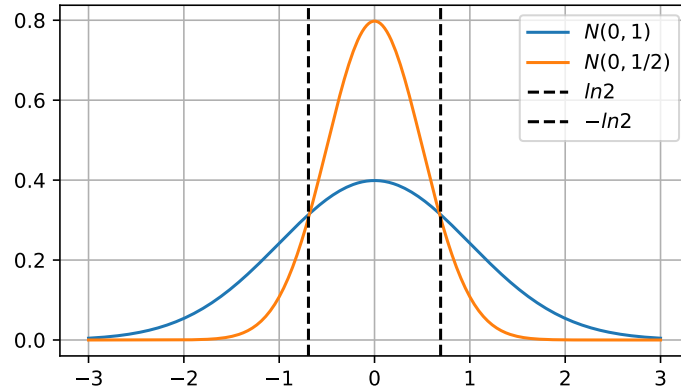


Figure 22: Plot for problem 12, chapter 10.

(ii) If $f'(P(\omega_1)) < 0$, the likelihood is maximized by the choice $P(\omega_1) = 0$.

It remains to investigate for which values of x_1 the derivative is positive or negative.

We take logarithms to study when $p(x | \omega_1) > p(x | \omega_2)$ and vice versa, since this determines the sign of the derivative. We compare $p(x | \omega_1) = \mathcal{N}(0, 1)$ and $p(x | \omega_2) = \mathcal{N}(0, 1/2)$ by equating them, taking logarithms and using some algebra:

$$\frac{1}{\sqrt{2}}e^{-x_1^2/2} = e^{-x_1^2} \quad \Leftrightarrow \quad \ln 2 = x_1^2.$$

When $x_1^2 < \ln 2$, we find that $e^{-x_1^2} > \frac{1}{\sqrt{2}}e^{-x_1^2/2}$. This implies that $p(x | \omega_2) > p(x | \omega_1)$ and the derivative of the linear function is negative. In that case the likelihood is maximized by $P(\omega_1) = 0$, as claimed in the problem.

- b) This is simply the reverse case of what we saw in the previous problem. When $x_1^2 > \ln 2$ we find that $p(x | \omega_1) > p(x | \omega_2)$ and the derivative is positive. In this case the likelihood is maximized by $P(\omega_1) = 1$.
- c) With two Gaussians having the same mean, but different standard deviations, the maximum likelihood estimate of $P(\omega_1)$ is either 0 or 1 when one data point x_1 is seen. The maximum likelihood assigns all weight to the prior corresponding to the class with the largest class-conditional density $p(x | \omega_i)$.

Problem 10.14

From the solution to Problem 6.39, we know that the derivative of $\mathbf{x}^T \mathbf{A} \mathbf{x}$ is $(\mathbf{A} + \mathbf{A}^T)\mathbf{x}$. This result is also in Appendix A.2 in [Duda et al., 2000]. To find the minimizer, we differentiate with respect to \mathbf{x} and find that

$$\frac{\partial J(\mathbf{x})}{\partial \mathbf{x}} = \sum_{k=1}^m (\boldsymbol{\Sigma}^{-1} + \boldsymbol{\Sigma}^{-T}) (\mathbf{x}_k - \mathbf{x}^*) = \mathbf{0}.$$

Since Σ may be moved outside of the sum and is non-singular, we require that

$$\sum_{k=1}^m (\mathbf{x}_k - \mathbf{x}^*) = \mathbf{0} \quad \Leftrightarrow \quad \sum_{k=1}^m \mathbf{x}_k = m\mathbf{x}^*.$$

Dividing by m , we see that $\mathbf{x}^* = \bar{\mathbf{x}}$, which is what we wanted to demonstrate.

Problem 10.16

A distance computation between two d -dimensional vectors is $O(d)$. Pseudocode makes the runtime evident, since we essentially only need to count the number of loops.

```
for each iteration (T):
  for each data point (n):
    for each cluster center (c):
      compute distance dist(x_j, mu_i) (d)
```

The algorithm has runtime $O(T) \cdot O(n) \cdot O(c) \cdot O(d) = O(Tncd)$ as claimed.

Problem 10.21

Let us assume that $n \geq c$ and that there exists an empty subset. Then there must exist at least one non-empty subset with more than one sample. An example would be $\{\{\mathbf{x}_1\}, \{\}, \{\mathbf{x}_2, \mathbf{x}_3\}\}$, where the second subset is empty, and the third subset has more than one sample. We show that moving a sample from a non-empty subset $\mathcal{D}_{\text{non-empty}}$ to an empty subset $\mathcal{D}_{\text{empty}}$ always decreases J_e .

The sum-of-squared error over the empty subset $\mathcal{D}_{\text{empty}}$ does not contribute to J_e , since it's only defined over non-empty subsets. Moving a single sample to $\mathcal{D}_{\text{empty}}$ contributes

$$\sum_{i=1}^1 \|\mathbf{x}_i - \mathbf{m}\|^2 = \|\mathbf{x}_1 - \mathbf{x}\|^2 = \|\mathbf{x}_1 - \mathbf{x}_1\|^2 = 0$$

to J_e , i.e. no contribution at all. In summary, moving a sample to $\mathcal{D}_{\text{empty}}$ does not alter the contribution to J_e from the subset $\mathcal{D}_{\text{empty}}$ —it's zero in both cases.

The sum-of-squared error over the non-empty subset $\mathcal{D}_{\text{non-empty}}$ does contribute to J_e . If $\mathcal{D}_{\text{non-empty}}$ contains n samples whose mean is \mathbf{m} , the contribution is

$$\sum_{i=1}^n \|\mathbf{x}_i - \mathbf{m}\|^2.$$

Now we remove a sample \mathbf{x}_n from the subset, to put in the empty subset instead. As long as the sample \mathbf{x}_n that we remove is not identical to the mean, we have

$$\sum_{i=1}^n \|\mathbf{x}_i - \mathbf{m}\|^2 = \sum_{i=1}^{n-1} \|\mathbf{x}_i - \mathbf{m}\|^2 + \|\mathbf{x}_n - \mathbf{m}\|^2 > \sum_{i=1}^{n-1} \|\mathbf{x}_i - \mathbf{m}\|^2.$$

In summary, moving a sample from $\mathcal{D}_{\text{non-empty}}$ *does* alter the contribution to J_e from the subset $\mathcal{D}_{\text{non-empty}}$ —it always decreases.

Combining the statements above reveals that if any empty subset exists, we can always decrease J_e . As a result there are no empty subsets in a partition that minimizes J_e . As an example, the partition $\{\{\mathbf{x}_1\}, \{\mathbf{x}_2\}, \{\mathbf{x}_3\}\}$ will be superior to $\{\{\mathbf{x}_1\}, \{\}, \{\mathbf{x}_2, \mathbf{x}_3\}\}$.

Problem 10.32

- a) Let's study the similarity measure under the assumption that $x_i, y_i \in \{-1, 1\}$. From the definition, we have that

$$s(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} = \frac{\sum_{i=1}^d x_i y_i}{\sqrt{d} \sqrt{d}} = \frac{\sum_{i=1}^d [x_i = y_i] - \sum_{i=1}^d [x_i \neq y_i]}{d},$$

where $[\cdot]$ is the Iverson bracket, returning 1 if the condition in the bracket is true, and 0 if it's false. An interpretation is therefore

$$s(\mathbf{x}, \mathbf{y}) = \frac{\sum_{i=1}^d [x_i = y_i] - \sum_{i=1}^d [x_i \neq y_i]}{d} = \frac{\text{num equal} - \text{num different}}{\text{total number}},$$

and we also note that $-1 \leq s(\mathbf{x}, \mathbf{y}) \leq 1$ are tight bounds, realized when every element is different, or element entry is equal, respectively.

- b) This is shown using algebra. We solve the problem by realizing that

$$\begin{aligned} \|\mathbf{x} - \mathbf{y}\|^2 &= \sum_{i=1}^d (x_i - y_i)^2 = \sum_{i=1}^d x_i^2 - \sum_{i=1}^d 2x_i y_i + \sum_{i=1}^d y_i^2 \\ &= d - 2\mathbf{x}^T \mathbf{y} - d = 2d - 2d \frac{\mathbf{x}^T \mathbf{y}}{d} = 2s(1 - s(\mathbf{x}, \mathbf{y})). \end{aligned}$$

Problem 10.35

The smallest increase happens when we merge clusters so that their contribution to J_e increases by as little as possible. We consider the error from clusters i and j before and after merging, and show that the difference in contribution to J_e is given by the expression in the problem statement.

The contribution before merging is

$$\sum_{\mathbf{x}_i \in \mathcal{D}_i} \|\mathbf{x}_i - \mathbf{m}_i\|^2 + \sum_{\mathbf{x}_j \in \mathcal{D}_j} \|\mathbf{x}_j - \mathbf{m}_j\|^2.$$

The contribution after merging is

$$\sum_{\mathbf{x}_k \in \mathcal{D}_i \cup \mathcal{D}_j} \|\mathbf{x}_k - \mathbf{m}\|^2 = \sum_{\mathbf{x}_i \in \mathcal{D}_i} \|\mathbf{x}_i - \mathbf{m}\|^2 + \sum_{\mathbf{x}_j \in \mathcal{D}_j} \|\mathbf{x}_j - \mathbf{m}\|^2,$$

where $\mathbf{m} = \mathbf{m}_i n_i / (n_i + n_j) + \mathbf{m}_j n_j / (n_i + n_j)$ is the new mean after merging the clusters.

We now compute the difference, which becomes

diff = after merging – before merging

$$\begin{aligned} &= \left(\sum_{\mathbf{x}_i \in \mathcal{D}_i} \|\mathbf{x}_i - \mathbf{m}\|^2 + \sum_{\mathbf{x}_j \in \mathcal{D}_j} \|\mathbf{x}_j - \mathbf{m}\|^2 \right) - \left(\sum_{\mathbf{x}_i \in \mathcal{D}_i} \|\mathbf{x}_i - \mathbf{m}_i\|^2 + \sum_{\mathbf{x}_j \in \mathcal{D}_j} \|\mathbf{x}_j - \mathbf{m}_j\|^2 \right) \\ &= \sum_{\mathbf{x}_i \in \mathcal{D}_i} (\|\mathbf{x}_i - \mathbf{m}\|^2 - \|\mathbf{x}_i - \mathbf{m}_i\|^2) + \sum_{\mathbf{x}_j \in \mathcal{D}_j} (\|\mathbf{x}_j - \mathbf{m}\|^2 - \|\mathbf{x}_j - \mathbf{m}_j\|^2) \end{aligned} \quad (11)$$

To make progress with Equation (11), we note that

$$\begin{aligned} \sum_{\mathbf{x}_i \in \mathcal{D}_i} \|\mathbf{x}_i - \mathbf{m}\|^2 - \|\mathbf{x}_i - \mathbf{m}_i\|^2 &= \sum_{\mathbf{x}_i \in \mathcal{D}_i} \|(\mathbf{x}_i - \mathbf{m}_i) - (\mathbf{m} - \mathbf{m}_i)\|^2 - \|\mathbf{x}_i - \mathbf{m}_i\|^2 \\ &= \sum_{\mathbf{x}_i \in \mathcal{D}_i} \|\mathbf{x}_i - \mathbf{m}_i\|^2 + \|\mathbf{m} - \mathbf{m}_i\|^2 - \|\mathbf{x}_i - \mathbf{m}_i\|^2 \\ &= \sum_{\mathbf{x}_i \in \mathcal{D}_i} \|\mathbf{m} - \mathbf{m}_i\|^2, \end{aligned}$$

where we have used the fact that the cross term $\sum_{\mathbf{x}_j \in \mathcal{D}_j} -2(\mathbf{x}_i - \mathbf{m}_i)^T (\mathbf{m} - \mathbf{m}_i)$ is zero.

Substituting this result back into Equation (11) we obtain

$$\begin{aligned} \text{diff} &= \text{after merging} - \text{before merging} \\ &= \sum_{\mathbf{x}_i \in \mathcal{D}_i} \|\mathbf{m} - \mathbf{m}_i\|^2 + \sum_{\mathbf{x}_j \in \mathcal{D}_j} \|\mathbf{m} - \mathbf{m}_j\|^2 \\ &= n_i \|\mathbf{m} - \mathbf{m}_i\|^2 + n_j \|\mathbf{m} - \mathbf{m}_j\|^2. \end{aligned} \quad (12)$$

Now we compute the difference

$$\mathbf{m} - \mathbf{m}_i = \left(\mathbf{m}_i \frac{n_i}{n_i + n_j} + \mathbf{m}_j \frac{n_j}{n_i + n_j} \right) - \mathbf{m}_i = \frac{n_j}{n_i + n_j} (\mathbf{m}_j - \mathbf{m}_i),$$

and use this result in Equation (12) to see that

$$\begin{aligned} \text{diff} &= \text{after merging} - \text{before merging} \\ &= n_i \frac{n_j^2}{(n_i + n_j)^2} \|\mathbf{m}_j - \mathbf{m}_i\|^2 + n_j \frac{n_i^2}{(n_i + n_j)^2} \|\mathbf{m}_i - \mathbf{m}_j\|^2 \\ &= \frac{n_i n_j}{(n_i + n_j)^2} (n_j + n_i) \|\mathbf{m}_i - \mathbf{m}_j\|^2 \\ &= \frac{n_i n_j}{n_i + n_j} \|\mathbf{m}_i - \mathbf{m}_j\|^2. \end{aligned}$$

This is the difference in the sum-of-squared error before and after merging. Of course it is never negative, since merging can only increase J_e . The smallest increase results from merging clusters i and j so that this difference is minimized.

References

- [Bishop, 2011] Bishop, C. M. (2011). *Pattern Recognition and Machine Learning*. Springer, New York.
- [Duda et al., 2000] Duda, R. O., Hart, P. E., and Stork, D. G. (2000). *Pattern Classification*. Wiley-Interscience, New York, 2 edition edition.