# SC4/SM8 Advanced Topics in Statistical Machine Learning

**Department of Statistics, University of Oxford**
https://github.com/ywteh/advml2020

Lecturer: Yee Whye Teh

Hilary Term 2020

# Contents

# Contents

# 1 Review of Fundamentals

## 1.1 Unsupervised Learning Basics

### Notational remarks

- We will typically assume that we have collected $p$ variables (features/attributes/dimensions) on $n$ examples (items/observations) which can be represented as an $n \times p$ *data matrix* $\mathbf{X} = (x_{ij})$, where $x_{ij}$ is the observed value of the $j$-th variable for the $i$-th example:

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1j} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2j} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{i1} & x_{i2} & \dots & x_{ij} & \dots & x_{ip} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nj} & \dots & x_{np} \end{bmatrix}. \tag{1.1}$$

- We will denote the *rows* of $\mathbf{X}$ as $x_i \in \mathbb{R}^p$ and treat them as *column vectors*: i.e., the $i$-th item/example/observation $x_i$ is the transpose of the $i$-th row of the data matrix $\mathbf{X}$:

$$x_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{bmatrix} = [x_{i1}, x_{i2}, \dots, x_{ip}]^\top, \quad i = 1, \dots, n. \tag{1.2}$$

- We often assume that $x_1, \dots, x_n$ are *independent and identically distributed (iid)* samples of a *random vector* $X$ over $\mathbb{R}^p$. When referring to the $j$-th dimension of random vector $X$, we will write $X^{(j)}$.

Unsupervised learning is a broad and arguably more challenging part of machine learning. The goal of unsupervised learning is to extract key features of the "unlabelled" dataset. While in supervised learning our data items $\{x_i\}_{i=1}^n$ come with an extra piece of information which we are trying to predict, in unsupervised learning we are trying to understand the process which generated data $\{x_i\}_{i=1}^n$ itself.

We will here review two basic unsupervised learning tasks: *dimensionality reduction* and *clustering*. Broadly speaking, **dimensionality reduction** aims to, for each data item $x_i \in \mathbb{R}^p$, find a lower dimensional representation $z_i \in \mathbb{R}^k$ with $k \ll p$ such that the map $x \mapsto z$ preserves certain *interesting statistical properties* in data. **Clustering** on the other hand, partitions the set of $n$ data items into $K$ disjoint groups. We will also review a simple algorithm for each: *Principal Components Analysis (PCA)* for dimensionality reduction and the $K$-*means* algorithm for clustering.

### 1.1.1 Dimensionality Reduction with PCA

**Principal Components Analysis (PCA)** is a dimensionality reduction technique which aims to preserve *variance* in the data. PCA is a *linear* dimensionality reduction technique: it essentially looks for a *new basis* to represent a noisy dataset.

For simplicity, we will assume for PCA that our dataset is **centred**, i.e., that its average is $\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i = 0$. If not, we can always subtract it from each $x_i$ (this is called **data centering**). Thus, we can write the **sample covariance matrix** $S$ as

$$S = \widehat{\mathrm{Cov}}(X) = \frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{x})(x_i - \bar{x})^\top = \frac{1}{n-1} \sum_{i=1}^{n} x_i x_i^\top = \frac{1}{n-1} \mathbf{X}^\top \mathbf{X}. \quad (1.3)$$

Note that the matrix $S$ is symmetric and positive semi-definite.

PCA recovers an orthonormal basis $v_1, v_2, \ldots, v_p$ in $\mathbb{R}^p$ – vectors $v_i$ are called **principal components** (PC) or **loading vectors** – such that:

- The first principal component (PC) $v_1 \in \mathbb{R}^p$ is (a unit length vector in) the *direction of greatest variance* of data.

- The $j$-th PC $v_j$ is the *direction orthogonal to $v_1, v_2, \ldots, v_{j-1}$ of greatest variance*, for $j = 2, \ldots, p$.

Given this basis, the $k$-dimensional representation of data item $x_i$ is the vector of projections of $x_i$ onto the first $k$ PCs:

$$z_i = V_{1:k}^\top x_i = \left[ v_1^\top x_i, \ldots, v_k^\top x_i \right]^\top \in \mathbb{R}^k,$$

where $V_{1:k} = [v_1, \ldots, v_k]$ is a $p \times k$ matrix. This gives us the *transformed data matrix*, also called the **score matrix**

$$\mathbf{Z} = \mathbf{X} V_{1:k} \in \mathbb{R}^{n \times k}. \quad (1.4)$$

#### Deriving the first principal component

Recall that we model our dataset as an iid sample $\{x_i\}_{i=1}^{n}$ of a random vector $X = \left[ X^{(1)} \ldots X^{(p)} \right]^\top$. Projections to PCs define a linear transformation of $X$ given by a random vector $Z = V_{1:k}^\top X$ which is a $k$-dimensional random vector. Dimensions of $Z$ are called **derived variables**. Consider the first dimension of $Z$:

$$Z^{(1)} = v_1^\top X = v_{11} X^{(1)} + v_{12} X^{(2)} + \cdots + v_{1p} X^{(p)}. \quad (1.5)$$

The first PC $v_1 = [v_{11}, \ldots, v_{1p}]^\top \in \mathbb{R}^p$ is chosen to maximise the sample variance $\widehat{\mathrm{Var}}(Z^{(1)}) = v_1^\top \widehat{\mathrm{Cov}}(X) v_1$, i.e. it is defined as the solution to

$$\max_{v_1} v_1^\top S v_1$$

$$\text{subject to: } v_1^\top v_1 = 1.$$

By considering the Lagrangian:

$$\mathcal{L}(v_1, \lambda_1) = v_1^\top S v_1 - \lambda_1 \left( v_1^\top v_1 - 1 \right) \tag{1.6}$$

and the corresponding vector of partial derivatives

$$\frac{\partial \mathcal{L}(v_1, \lambda_1)}{\partial v_1} = 2 S v_1 - 2 \lambda_1 v_1 \tag{1.7}$$

we obtain the eigenvector equation $S v_1 = \lambda_1 v_1$, i.e. $v_1$ must be an eigenvector of $S$ and the dual variable $\lambda_1$ is the corresponding eigenvalue. Since $v_1^\top S v_1 = \lambda_1 v_1^\top v_1 = \lambda_1$, the first PC must be the eigenvector associated with the *largest eigenvalue* of $S$.

Similarly, we can show that each subsequent PC is given by the eigenvector corresponding to the next largest eigenvalue (problem sheet shows this for the second PC). As a result, we obtain the **eigenvalue decomposition** of $S$ given by

$$S = V \Lambda V^\top \tag{1.8}$$

where $\Lambda$ is a diagonal matrix with eigenvalues

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_p \geq 0 \tag{1.9}$$

on the diagonal and $V$ is a $p \times p$ orthogonal matrix (i.e. $VV^\top = V^\top V = I$) whose *columns* are the $p$ eigenvectors of $S$, i.e. the principal components $v_1, \ldots, v_p$.

In summary,

- Derived scalar variable (projection to the $j$-th principal component) $Z^{(j)} = v_j^\top X$ has sample variance $\lambda_j$, for $j = 1, \ldots, p$.

- Derived variables are *uncorrelated*: $\mathrm{Cov}(Z^{(i)}, Z^{(j)}) \approx v_i^\top S v_j = \lambda_j v_i^\top v_j = 0$, for $i \neq j$.

- The **total sample variance** is given by $\mathrm{Tr}(S) = \sum_{i=1}^p S_{ii} = \lambda_1 + \ldots + \lambda_p$, so the **proportion of total variance** explained by the $j^{th}$ PC is $\frac{\lambda_j}{\lambda_1 + \lambda_2 + \ldots + \lambda_p}$

### Reconstruction view of PCA

We can map back to the original $p$-dimensional space using

$$\hat{x}_i = V_{1:k} V_{1:k}^\top x_i. \tag{1.10}$$

This is a **reconstruction** of data item $x_i$. It can be shown (exercise) that PCA gives the *optimal linear reconstruction* based on a $k$-dimensional compression.

## PCA via the Singular Value Decomposition

PCA can also be understood using the **singular value decomposition** (SVD) of data matrix $\mathbf{X}$. Recall that any real-valued $n \times p$ matrix $\mathbf{X}$ can be written as $\mathbf{X} = UDV^\top$ where

- $U$ is an $n \times n$ orthogonal matrix: $UU^\top = U^\top U = I_n$.

- $D$ is a $n \times p$ matrix with decreasing *non-negative* elements on the diagonal (the singular values of $\mathbf{X}$) and zero off-diagonal elements.

- $V$ is a $p \times p$ orthogonal matrix: $VV^\top = V^\top V = I_p$.

Note that

$$(n-1)S = \mathbf{X}^\top\mathbf{X} = (UDV^\top)^\top(UDV^\top) = VD^\top U^\top UDV^\top = VD^\top DV^\top,$$

using orthogonality of $U$. The eigenvalues of $S$ are thus the diagonal entries of $\Lambda = \frac{1}{n-1}D^\top D$.

We also have

$$\mathbf{X}\mathbf{X}^\top = (UDV^\top)(UDV^\top)^\top = UDV^\top VD^\top U^\top = UDD^\top U^\top,$$

using orthogonality of $V$.

The $n \times n$ matrix $\mathbf{K} = \mathbf{X}\mathbf{X}^\top$ with entries $\mathbf{K}_{ij} = x_i^\top x_j$ is called the **Gram matrix** of dataset $\mathbf{X}$. Note that $\mathbf{K}$ and $(n-1)S = \mathbf{X}^\top\mathbf{X}$ have the same nonzero eigenvalues, equal to the non-zero squared singular values of $\mathbf{X}$ (non-zero entries on the diagonals of $D^\top D$ and $DD^\top$).

If we consider projections to *all principal components*, the transformed data matrix is

$$\mathbf{Z} = \mathbf{X}V = UDV^\top V = UD, \tag{1.11}$$

If $p \le n$ this means

$$z_i = [U_{i1}D_{11}, \ldots, U_{ip}D_{pp}]^\top, \tag{1.12}$$

and if $p > n$ only the first $n$ projections are defined (sample covariance will be at most rank $n$):

$$z_i = [U_{i1}D_{11}, \ldots, U_{in}D_{nn}, 0, \ldots, 0]^\top. \tag{1.13}$$

Thus, $\mathbf{Z}$ can be obtained from the eigendecomposition of Gram matrix $\mathbf{K}$. When $p \gg n$, eigendecomposition of $\mathbf{K}$ requires much less computation, $O(n^3)$, than the eigendecomposition of the covariance matrix, $O(p^3)$, so is the preferred method for PCA in that case.

## 1.1.2 Clustering

*Clustering* is one of the fundamental and ubiquitous tasks in exploratory data analysis – a first intuition about the data is often based on identifying meaningful disjoint groups among the data items. In *partition-based* clustering, which we consider in this note, one divides $n$ data items into $K$ clusters $C_1, \ldots, C_K$ where for all $k, k' \in \{1, \ldots, K\}$,

$$C_k \subset \{1, \ldots, n\}, \qquad C_k \cap C_{k'} = \emptyset \ \ \forall k \neq k', \qquad \bigcup_{k=1}^{K} C_k = \{1, \ldots, n\}.$$

Central to the goals of clustering is the notion of similarity/dissimilarity between data items. There will be many ways to define the notion of similarity, and the choice will depend on the dataset being analyzed and dictated by domain specific knowledge.

Intuitively, clustering aims to group similar items together and to place separate dissimilar items into different groups. However, note that these two objectives in many cases contradict each other (similarity is not a transitive relation, while being in the same cluster is an equivalence relation). One could imagine a long sequence of items such that each next item is very similar to the previous one so that they should all belong to the same cluster – but that would also mean that the endpoints are potentially highly dissimilar. Hence, there are also different clustering techniques which emphasize different aspects of these goals, i.e. whether to keep similar points together or dissimilar points apart.

### Lack of Axiomatic Definition

There have been several attempts to construct an axiomatic definition of clustering, but it is surprisingly difficult to put on rigorous footing. Consider the following three basic properties required of a clustering method $\mathcal{F} : (\mathcal{D} = \{x_i\}_{i=1}^{n}, \rho) \mapsto \{C_1, \ldots, C_K\}$ which takes as an input dataset $\mathcal{D}$ and a dissimilarity function $\rho$ and returns a partition of $\mathcal{D}$:

- **Scale invariance.** For any $\alpha > 0$, $\mathcal{F}(\mathcal{D}, \alpha\rho) = \mathcal{F}(\mathcal{D}, \rho)$, i.e. partition should not depend on units in which dissimilarity is measured.

- **Richness.** For any partition $C = \{C_1, \ldots, C_K\}$ of $\mathcal{D}$, there exists dissimilarity $\rho$, such that $\mathcal{F}(\mathcal{D}, \rho) = C$, i.e. the outcome is fully controlled by the dissimilarity function.

- **Consistency.** If $\rho$ and $\rho'$ are two dissimilarities such that for all $x_i, x_j \in \mathcal{D}$ the following holds:

$$x_i, x_j \text{ belong to the same cluster in } \mathcal{F}(\mathcal{D}, \rho) \implies \rho'(x_i, x_j) \leq \rho(x_i, x_j)$$
$$x_i, x_j \text{ belong to different clusters in } \mathcal{F}(\mathcal{D}, \rho) \implies \rho'(x_i, x_j) \geq \rho(x_i, x_j),$$

  then $\mathcal{F}(\mathcal{D}, \rho') = \mathcal{F}(\mathcal{D}, \rho)$. In other words, if the items in the same cluster become more similar and the items already separated become less similar, then the clustering should not change.

While all three properties appear natural, and there are algorithms satisfying any two of them, Kleinberg's impossibility theorem [13] states that there exists no clustering method that satisfies all three properties, implying that every clustering method will have some undesirable properties. For further discussion, see Section 22.5 in [21].

We will consider here the simplest widely used clustering method: $K$-means algorithm and two extensions.

### $K$-means algorithm

$K$-*means* is the simplest partition-based clustering algorithm. It uses a preassigned number of clusters and represents each cluster using a **prototype** or **cluster centroid** $\mu_k$.

The idea of $K$-means is to measure the quality of each cluster $C_k$ using its **within-cluster deviance** from the cluster centroids

$$W(C_k, \mu_k) = \sum_{i \in C_k} \|x_i - \mu_k\|_2^2.$$

The overall quality of the clustering is then given by the total within-cluster deviance:

$$W(\{C_k\}, \{\mu_k\}) = \sum_{k=1}^{K} W(C_k, \mu_k) = \sum_{k=1}^{K} \sum_{i \in C_k} \|x_i - \mu_k\|_2^2 = \sum_{i=1}^{n} \|x_i - \mu_{c_i}\|_2^2,$$

where $c_i = k$ if and only if $i \in C_k$. This is now the overall objective function used to select both the cluster centroids and the assignment of points to clusters. The joint optimization over both the partition $\{C_k\}$ and centroids $\{\mu_k\}$ is a combinatorial optimization problem and is computationally hard. However, note that

- Given partition $\{C_k\}$, we can easily find the optimal centroids by differentiating $W$ with respect to $\mu_k$:

$$\frac{\partial W}{\partial \mu_k} = 2 \sum_{i \in C_k} (x_i - \mu_k) = 0 \qquad \Rightarrow \mu_k = \frac{1}{|C_k|} \sum_{i \in C_k} x_i$$

- Given prototypes, we can easily find the optimal partition by assigning each data point to the closest cluster prototype:

$$c_i = \operatorname{argmin}_k \|x_i - \mu_k\|_2^2.$$

Thus one can employ an iterative alternating optimization, which is exactly the $K$-means algorithm given in Algorithm 1.1.

$K$-means is a heuristic search algorithm so it can (and often will) get stuck at local optima. The result depends on the starting configurations. Typically one performs a number of runs from different random initial values of centroids, and then chooses the end result with minimum $W$. Since each step does not increase the objective function and the number of possible partitions is finite, the algorithm will converge to a local optimum. However, note that there could be ties in the cluster assignment, which need to be broken in a systematic fashion.

---

**Algorithm 1.1** K-means algorithm

---

**Input**: dataset $\mathcal{D} = \{x_i\}_{i=1}^n$, desired number of clusters $K$
**Output**: partition $\{C_1, \ldots, C_K\}$

Randomly initialize $K$ cluster centroids $\mu_1, \ldots, \mu_K$.
**while** the partition has not converged **do**

- *Cluster assignment:* For each $i = 1, \ldots, n$, assign each $x_i$ to the cluster with the nearest centroid,

$$c_i := \mathrm{argmin}_{k=1,\ldots,K} \|x_i - \mu_k\|_2^2$$

  Set $C_k := \{i : c_i = k\}$ for each $k$.

- *Move centroids:* Set $\mu_1, \ldots, \mu_K$ to the averages of the new clusters:

$$\mu_k := \frac{1}{|C_k|} \sum_{i \in C_k} x_i$$

**return** partition $\{C_1, \ldots, C_K\}$

---

### K-means++

A simple yet provably effective solution to the problem of initialization of centroids in the K-means algorithm was proposed by [1]. The method starts with sampling a data item from $\mathcal{D} = \{x_i\}_{i=1}^n$ uniformly at random and making it centroid $\mu_1$. We then compute the squared distances $\rho_i^2 = \|x_i - \mu_1\|_2^2$. Centroid $\mu_2$ is then initialized to another data item sampled using the probability mass function $p(i) = \rho_i^2 / \sum_{j=1}^n \rho_j^2$ and the process continues with the probability mass function being updated at each step, i.e. to initialize $k$-th centroid $\mu_k$, we compute

$$\rho_i^2 = \min \left\{ \|x_i - \mu_1\|_2^2, \ldots, \|x_i - \mu_{k-1}\|_2^2 \right\}. \tag{1.14}$$

Remarkably, this method comes with a precise theoretical guarantee. In particular, [1] show that if partition $\{C_k^{++}\}$ is obtained using K-means++ then

$$\mathbb{E} \left[ W \left( \{C_k^{++}\} \right) \right] \leq 8 \left( \log K + 2 \right) W^*, \tag{1.15}$$

where $W^*$ is the within-cluster deviance of the globally optimal clustering and the expectation is taken over the random sampling used in the initialisation.

### DP-means

$K$-means is intuitive and straightforward to implement, but how do we select the number of clusters $K$ in the first place? Clearly, the objective function is minimized (and equals zero) if we let $K = n$, but this is not a meaningful clustering.

One elegant approach is the *DP-means algorithm* [14] that comes from the interpretation of $K$-means using small variance asymptotics of the Expectation Maximization (EM) algorithm for mixture modelling (see Chapter 5). We will discuss mixture modelling and EM algorithm later in the course. DP-means starts from a single cluster, i.e. $K = 1$ and modifies the cluster assignment step as follows:

1. Initialize $K = 1$ and $\mu_1 = \frac{1}{n} \sum_{i=1}^{n} x_i$ (the global mean).

2. *DP-means cluster assignment:* For each $i = 1, \ldots, n$,
    - if $\min_{k=1,\ldots,K} \|x_i - \mu_k\|_2^2 > \lambda$, set $K \leftarrow K + 1$, $c_i \leftarrow K$, $\mu_K \leftarrow x_i$
    - otherwise, set $c_i = \operatorname{argmin}_{k=1,\ldots,K} \|x_i - \mu_k\|_2^2$, and update the previous and current cluster centroids that item $i$ belongs to.

3. Repeat Step 2 until it stops changing the cluster assignments for all items.

The tuning parameter $\lambda$ controls the tradeoff between the traditional $K$-means objective and the number of clusters, and can be interpreted as a cost associated with each cluster used. Like $K$-means, DP-means can be shown to terminate after a finite number of iterations (problem sheet).

## 1.2 Supervised Learning Basics

### 1.2.1 Empirical Risk Minimisation (ERM)

**Loss and Risk**

In **supervised learning**, we are trying to learn a function $f : \mathcal{X} \rightarrow \mathcal{Y}$ from an input space $\mathcal{X}$ into an output space $\mathcal{Y}$ based on a set of paired examples $(x_1, y_1), \ldots (x_n, y_n)$ and a given **loss function** $L$. It is typically assumed that examples $(x_1, y_1), \ldots (x_n, y_n)$ are i.i.d. samples from an *unknown joint probability distribution* $P_{X,Y}$ on $\mathcal{X} \times \mathcal{Y}$. The goal of supervised learning is to find the function $f$ which *minimizes the expectation of the loss* over $P_{X,Y}$, which is called *risk*. In machine learning, if the output space is discrete, this is called **classification**, while **regression** refers to the case the output space is continuous.

**Empirical Risk Minimisation (ERM)**

A **loss function** is any function

$$L : \mathcal{Y} \times \mathcal{Y} \times \mathcal{X} \rightarrow \mathbb{R}^+. \tag{1.16}$$

The loss $L(y_i, f(x_i), x_i)$ measure the discrepancy between predicted output values $f(x_i)$ and the true output values $y_i$ at the inputs $x_i$.

The **risk** of a function $f : \mathcal{X} \rightarrow \mathcal{Y}$ is the expected loss:

$$R(f) = \mathbb{E}_{X,Y} L(Y, f(X), X). \tag{1.17}$$

For a given dataset $(x_1, y_1), \ldots (x_n, y_n)$, the **empirical risk** of $f$ is given by

$$\hat{R}(f) = \frac{1}{n} \sum_{i=1}^{n} L(y_i, f(x_i), x_i). \tag{1.18}$$

The **empirical risk minimisation** (ERM) problem aims to find the function with minimum risk:

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{H}} \hat{R}(f), \tag{1.19}$$

where $\mathcal{H}$ is a given class of functions (called the **hypothesis class**).

*Remark* 1. The ultimate goal of learning is to minimise the true risk - *not* the empirical risk, which is only an estimate of the true risk. But the true risk of any given function is unknown because the distribution $P_{X,Y}$ is unknown.

*Remark* 2. Loss functions typically depend on the input $x$ only through $f(x)$, so that with some abuse of notation we often write $L(y, f(x))$ instead of $L(y, f(x), x)$. $L(y, f(x))$ is usually some notion of distance between the true output $y$ and the predicted output $f(x)$.

### Examples of hypothesis classes.

Hypothesis classes can be very simple, e.g. for $\mathcal{X} = \mathbb{R}^p$, we can consider all linear functions $f(x) = w^\top x + b$, parametrized by $w \in \mathbb{R}^p$ and $b \in \mathbb{R}$, or we could consider a specific **nonlinear feature expansion** $\varphi : \mathcal{X} \to \mathbb{R}^D$, and a model linear in those features: $f(x) = w^\top \varphi(x) + b$, but nonlinear in the original inputs $\mathcal{X}$, parametrized by $w \in \mathbb{R}^D$ and $b \in \mathbb{R}$. For example, starting with $\mathcal{X} = \mathbb{R}^2$, we can consider $\varphi\left(\begin{bmatrix} x_{i1} \\ x_{i2} \end{bmatrix}\right) = [x_{i1}, x_{i2}, x_{i1}^2, \sqrt{2}x_{i1}x_{i2}, x_{i2}^2]^\top$, such that the resulting function can depend on quadratic and interaction terms as well. An important type of hypothesis class we will consider in this course are *Reproducing Kernel Hilbert Spaces (RKHS)*, which are also linear in certain feature expansions but those feature expansions could potentially be infinite-dimensional; see Chapter 3.

### Examples of loss functions.

Loss functions come in many different forms. One of the main considerations for selecting loss functions is the type of outputs we are trying to predict, i.e., whether it is real-valued or discrete/categorical. Note that even if outputs are discrete, the function $f(x)$ we are trying to learn is typically real-valued. For example, in binary classification, the common convention is that the two classes are denoted by $-1$ and $+1$. One associates predictions of these classes with $\operatorname{sign}(f(x))$, whereas the magnitude of $f(x)$ can be thought of as the confidence in those predictions (not necessarily in a probabilistic sense). The loss

can penalize misclassification (wrong sign) as well as the overconfident misclassification (wrong sign and large magnitude) and even underconfident correct classification (correct sign but small magnitude). Thus, the loss functions can often be expressed as a function of $yf(x)$.



Figure 1.1: Loss functions for binary classification

Below are some loss functions commonly used in binary classification and regression.

- Binary classification:
  - *0/1 loss*: $L(y, f(x)) = \mathbf{1}\{yf(x) \leq 0\}$,
    (also called **misclassification loss**. The optimal solution is called the **Bayes classifier** and is given by $f(x) = \text{argmax}_{k \in \{0,1\}} \mathbb{P}(Y = k | X = x)$),
  - **hinge loss**: $L(y, f(x)) = (1 - yf(x))_+$
    (used in *support vector machines* - leads to sparse solutions),
  - **exponential loss**: $L(y, f(x)) = e^{-yf(x)}$
    (used in *boosting* algorithms, e.g. Adaboost),
  - **logistic loss**: $L(y, f(x)) = \log\left(1 + e^{-yf(x)}\right)$
    (used in *logistic regression*, and associated with a linear log-odds probabilistic model).

- Regression:
  - **squared loss**: $L(y, f(x)) = (y - f(x))^2$
    (**least squares regression**: optimal $f$ is the conditional mean $\mathbb{E}[Y | X = x]$),
  - **absolute loss**: $L(y, f(x)) = |y - f(x)|$
    (**least absolute deviations regression**, which is less sensitive to outliers: optimal $f$ is the conditional median $\text{med}[Y | X = x]$),

Figure 1.2: Loss functions for regression

- $\tau$-**pinball loss**: $L(y, f(x)) = 2\max\{\tau(y - f(x)), (\tau - 1)(y - f(x))\}$ for $\tau \in (0, 1)$
  (**quantile regression**: optimal $f$ is the $\tau$-quantile of $p(y|X = x)$),

- $\epsilon$-**insensitive loss** or **Vapnik loss**: $L(y, f(x)) = \begin{cases} 0, \text{ if } |y - f(x)| \leq \epsilon, \\ |y - f(x)| - \epsilon, \text{ otherwise.} \end{cases}$
  (**support vector regression**, which leads to sparse solutions).

In binary classification, $0/1$ is an idealised version of loss which penalizes misclassification regardless of the magnitude of $f(x)$. However, ERM under $0/1$ loss is NP hard[1]. Therefore, we typically use *convex upper bound surrogate losses* (hinge, exponential, logistic[2]). What is the importance of the convexity of loss as a function of $yf(x)$ as shown in Fig. 1.1? Consider the hypothesis class $f(x) = w^\top \varphi(x)$, with $w \in \mathbb{R}^D$ (we ignore the intercept to simplify notation) and assume that $L(y, f(x)) = \rho(yf(x))$ for a convex differentiable function $\rho$. Then the empirical risk and its gradient are given by

$$\hat{R}(w) = \frac{1}{n}\sum_{i=1}^{n} \rho\left(y_i w^\top \varphi(x_i)\right), \quad \frac{\partial \hat{R}}{\partial w} = \frac{1}{n}\sum_{i=1}^{n} \rho'\left(y_i w^\top \varphi(x_i)\right) y_i \varphi(x_i).$$

Furthermore, the Hessian matrix of the empirical risk is given by

$$\frac{\partial^2 \hat{R}}{\partial w \partial w^\top} = \frac{1}{n}\sum_{i=1}^{n} \rho''\left(y_i w^\top \varphi(x_i)\right) \varphi(x_i)\varphi(x_i)^\top, \quad (1.20)$$

---

[1]It is NP-hard to even approximately minimize the ERM under $0/1$ loss - i.e. there is no known polynomial-time algorithm to obtain a solution which is a small constant worse than the optimum.

[2]to make it into an upper bound on $0/1$, divide the logistic loss by $\log(2)$ - rescaling of the loss does not change the ERM problem

using $y_i^2 = 1$. This Hessian is now a positive semidefinite matrix which can be seen from $\rho''(t) \geq 0 \; \forall t$ and

$$\alpha^\top \frac{\partial^2 \hat{R}}{\partial w \partial w^\top} \alpha \;=\; \frac{1}{n} \sum_{i=1}^{n} \rho'' \left( y_i w^\top \varphi(x_i) \right) \left( \alpha^\top \varphi(x_i) \right)^2 \geq 0.$$

for any $\alpha \in \mathbb{R}^D$. Thus, empirical risk is a convex function of $w$ and thus has a *unique minimum*. Typically, there is no closed form solution for $w$ and iterative optimisation techniques like *gradient ascent* or *Newton-Raphson algorithm* are used.

### 1.2.2 Regularisation

Recall that we are not ultimately interested in the exact minimizer of the *empirical risk* but in that of the *true risk*. ERM thus risks **overfitting**: when the hypothesis class is complex, one can easily find a function that matches the observed examples exactly but does not *generalise* to other examples drawn from $P_{X,Y}$.

The idea behind **regularisation** is to limit the flexibility of the hypothesis class in order to prevent overfitting. For example, for the hypothesis space $\mathcal{H} = \{ f_\theta : \theta \in \mathbb{R}^p \}$ consisting of functions parameterised by a vector $\theta$, this can be achieved by adding a term which *penalises large values for the parameters $\theta$* to the ERM criterion:

$$\min_\theta \hat{R}(f_\theta) + \lambda \|\theta\|_\rho^\rho = \min_\theta \frac{1}{n} \sum_{i=1}^{n} L(y_i, f_\theta(x_i)) + \lambda \|\theta\|_\rho^\rho$$

where $\rho \geq 1$, and $\|\theta\|_\rho = (\sum_{j=1}^{p} |\theta_j|^\rho)^{1/\rho}$ is the $L_\rho$ norm of $\theta$ (also of interest when $\rho \in [0, 1)$, but this is no longer a norm). These methods are also known as **shrinkage** methods since their effect is to shrink parameter estimates towards 0. Note that we have an additional **tuning parameter** (or **hyperparameter**) $\lambda$ which controls the amount of regularisation, and, as a result, also controls the complexity of the model.

The most common forms of regularisation include **ridge regression / Tikhonov regularization**: $\rho = 2$, **LASSO** penalty: $\rho = 1$, and **elastic net** regularization with a mixed $L_1/L_2$ penalty:

$$\min_\theta \frac{1}{n} \sum_{i=1}^{n} L(y_i, f_\theta(x_i)) + \lambda \left[ (1 - \alpha) \|\theta\|_2^2 + \alpha \|\theta\|_1 \right].$$

In some hypothesis classes, it is possible to directly penalise some notion of *smoothness* of the function $f$ we are trying to learn, e.g. for $\mathcal{X} = \mathbb{R}$, the regularisation term can consist of the **Sobolev norm**

$$\|f\|_{W^1}^2 = \int_{-\infty}^{+\infty} f(x)^2 dx + \int_{-\infty}^{+\infty} f'(x)^2 dx, \tag{1.21}$$

which penalises functions with large derivative values.

### 1.2.3 Examples of ERM

**Regularised Least Squares / Ridge Regression**

This corresponds to the squared loss $L(y, f(x)) = (y - f(x))^2$. For linear functions $f(x) = w^\top x + b$, we have

$$\min_{w,b} \frac{1}{n} \sum_{i=1}^{n} (y_i - w^\top x_i - b)^2 + \frac{\lambda}{n} \|w\|_2^2. \tag{1.22}$$

Note the rescaling of the regularisation term and that the bias term $b$ is not included in the regularisation. This is important as otherwise the predictions would depend on the origin for the response variables $y$ (i.e. adding a constant $c$ to each target would result in different predictions from simply shifting the original predictions by $c$). Fortunately, when using centred inputs, i.e., $\sum_{i=1}^{n} x_i = 0$, $b$ can be estimated by $\bar{y} = \frac{1}{n} \sum_{i=1}^{n} y_i$, so we can also assume that the responses are centred and remove the intercept from the model. We obtain the problem

$$\min_w \|\mathbf{y} - \mathbf{X}w\|_2^2 + \lambda \|w\|_2^2. \tag{1.23}$$

Differentiating and setting to zero gives the closed form solution

$$w = \left( \mathbf{X}^\top \mathbf{X} + \lambda I \right)^{-1} \mathbf{X}^\top \mathbf{y}. \tag{1.24}$$

**Logistic Regression**

Despite the name, *logistic regression is a method for classification*. It uses the logistic loss $L(y, f(x)) = \log \left(1 + e^{-yf(x)}\right)$. Hence, again for a linear classifier $f(x) = w^\top x + b$,

$$\min_{w,b} \frac{1}{n} \sum_{i=1}^{n} \log \left( 1 + e^{-y_i(w^\top x_i + b)} \right) + \frac{\lambda}{n} \|w\|_2^2. \tag{1.25}$$

Logistic regression can also be associated to a probabilistic model. Namely, assume that the function of interest $f(x) = w^\top x + b$ models the log-odds ratio:

$$\log \frac{p(y_i = +1 | w, b, x_i)}{p(y_i = -1 | w, b, x_i)} = w^\top x_i + b. \tag{1.26}$$

Then the conditional distribution of $Y|X$ is given by

$$p(y_i = +1 | w, b, x_i) = \frac{1}{1 + e^{-(w^\top x_i + b)}} = \sigma(w^\top x_i + b), \tag{1.27}$$

$$p(y_i = -1 | w, b, x_i) = \frac{1}{1 + e^{w^\top x_i + b}} = \sigma(-w^\top x_i - b), \tag{1.28}$$

where we denoted by $\sigma(t) = 1/(1 + e^{-t})$ the **logistic function** (aka **sigmoid function**) which maps the real line to the $(0, 1)$ interval. Note that the logistic function satisfies

$\sigma(-t) = 1 - \sigma(t)$. Thus, we can write (1.27) and (1.28) as $p(y_i|w, b, x_i) = \sigma(y_i(w^\top x_i + b))$ and the conditional log-likelihood of the outputs given the inputs is

$$\log p(\mathbf{y}|w, b, \mathbf{X}) = \log \prod_{i=1}^{n} \sigma(y_i(w^\top x_i + b)) = -\sum_{i=1}^{n} \log\left(1 + e^{-y_i(w^\top x_i + b)}\right).$$

Thus finding the parameters $w$ and $b$ that maximise the conditional log-likelihood is equivalent to minimising the empirical risk corresponding to the logistic loss, which is the negative log-likelihood of the linear log-odds model. Moreover, the regularisation term can be interpreted as a normal prior on $w$ in *Bayesian logistic regression*. Again, there is no closed form solution for logistic regression, but the objective is convex and differentiable and the numerical optimisation via gradient ascent or Newton-Raphson algorithm can be used.

The connection between maximisation of the log-likelihood and minimisation of the empirical risk extends beyond logistic regression. Indeed, in the context of classification, whenever $p(y_i|x_i, \theta)$ is a log-concave function of $y_i f_\theta(x_i)$, we can define a convex loss $\rho(y f_\theta(x)) = -\log p(y_i|x_i, \theta)$. But the converse is not true, e.g. hinge loss used in the SVMs below does not correspond to a negative log-likelihood in any probabilistic model (unless additional artificial classes are introduced).

## Support Vector Machines

Support Vector Machines (SVMs) for classification use hinge loss, $L(y, f(x)) = \max\{0, 1 - yf(x)\}$. Thus, for a linear classifier $f(x) = w^\top x + b$, we obtain

$$\min_{w, b} \frac{1}{n} \sum_{i=1}^{n} \max\{0, 1 - y_i(w^\top x_i + b)\} + \frac{\lambda}{n}\|w\|_2^2. \tag{1.29}$$

This does not have a closed form solution and requires numerical optimisation.

Eq. (1.29) is not how you would typically see an SVM written in the literature, though. In the next chapter, we will make a deep dive into SVMs, introducing it from the perspective of *maximum margin classification*.

# 2 Support Vector Machines

These notes are a revised version of lecture notes from the UCL course "Reproducing Kernel Hilbert Spaces in Machine Learning" [10], reproduced here by courtesy of Arthur Gretton.

## 2.1 Duality in Convex Optimization

We will need some basic results from **duality** in **convex optimization** in order to study support vector machines, one of the fundamental techniques for classification. This review covers the material from [5, Sections 5.1-5.5].

### 2.1.1 The Lagrangian

Consider a **constrained optimization** problem of an objective function $f_0 : \mathbb{R}^n \to \mathbb{R}$, with $m$ inequality and $r$ equality constraints:

$$
\begin{aligned}
\text{minimize} \quad & f_0(x) & & \\
\text{subject to} \quad & f_i(x) \le 0 & & i = 1, \ldots, m \\
& h_j(x) = 0 & & j = 1, \ldots r.
\end{aligned} \tag{2.1}
$$

We denote $\mathcal{D} := \bigcap_{i=0}^m \mathrm{dom} f_i \cap \bigcap_{j=1}^r \mathrm{dom} h_j$, and require the domain $\mathcal{D} \subseteq \mathbb{R}^n$ where the objective function $f_0$ and the constraint functions $f_1, \ldots, f_m, h_1, \ldots, h_r$ are all defined to be nonempty. We will refer to (2.1) as the **primal problem** and denote by $p^* = f_0(x^*)$ its optimal value. Any point $\tilde{x} \in \mathcal{D}$ for which constraints are satisfied , i.e. $f_i(\tilde{x}) \le 0$, $h_j(\tilde{x}) = 0$, is called a **primal feasible** point.

The **Lagrangian** $L : \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^r \to \mathbb{R}$ associated with problem (2.1) is given by

$$
L(x, \lambda, \nu) := f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{j=1}^r \nu_j h_j(x).
$$

The vectors $\lambda \in \mathbb{R}^m$ and $\nu \in \mathbb{R}^r$ are called **Lagrange multipliers** or **dual variables**. The **Lagrange dual function** (or just "dual function") is written

$$
g(\lambda, \nu) = \inf_{x \in \mathcal{D}} L(x, \lambda, \nu).
$$

The domain of $g$, $\mathrm{dom} g$, is the set of values $(\lambda, \mu)$ for which the Lagrangian is bounded from below, i.e. $g > -\infty$. The dual function is a pointwise infimum of **affine**[1] functions of $(\lambda, \nu)$, hence it is concave in $(\lambda, \nu)$ [5, p. 83]. A **dual feasible** pair $(\lambda, \nu)$ is a pair for which $\lambda \succeq 0$ and $(\lambda, \nu) \in \mathrm{dom} g$.

---

[1] A function $f : \mathbb{R}^n \to \mathbb{R}^m$ is affine if it takes the form $f(x) = Ax + b$.

Figure 2.1: Example: Lagrangian with one inequality constraint, $L(x, \lambda) = f_0(x) + \lambda f_1(x)$, where $x$ here can take one of four values for ease of illustration. The infimum of the resulting set of four affine functions is concave in $\lambda$.

**Proposition 3.** *When $\lambda \succeq 0$, then for all $\nu$ we have*

$$g(\lambda, \nu) \le p^*. \tag{2.2}$$

*Proof.* Assume $\tilde{x} \in \mathcal{D}$ is feasible, i.e. $f_i(\tilde{x}) \le 0$, $h_j(\tilde{x}) = 0$, and assume $\lambda \succeq 0$. Then

$$\sum_{i=1}^{m} \lambda_i f_i(\tilde{x}) + \sum_{j=1}^{r} \nu_j h_j(\tilde{x}) \le 0$$

and so

$$
\begin{aligned}
g(\lambda, \nu) \quad &:= \quad \inf_{x \in \mathcal{D}} \left( f_0(x) + \sum_{i=1}^{m} \lambda_i f_i(x) + \sum_{j=1}^{r} \nu_j h_j(x) \right) \\
&\le \quad f_0(\tilde{x}) + \sum_{i=1}^{m} \lambda_i f_i(\tilde{x}) + \sum_{j=1}^{r} \nu_j h_j(\tilde{x}) \\
&\le \quad f_0(\tilde{x}).
\end{aligned}
$$

This holds for every feasible $\tilde{x}$, and thus also for $x^*$, hence (2.2) holds. $\qquad\square$

The Lagrangian can be interpreted as a lower bound on the original optimization problem. Ideally we would write the problem (2.1) as the unconstrained problem

$$\text{minimize } f_0(x) + \sum_{i=1}^{m} I_- \left( f_i(x) \right) + \sum_{j=1}^{r} I_0 \left( h_j(x) \right),$$

Figure 2.2: Linear lower bounds on indicator functions. Blue functions represent linear lower bounds for different slopes $\lambda$ and $\nu$, for the inequality and equality constraints, respectively.

where

$$I_-(u) = \begin{cases} 0 & u \leq 0 \\ \infty & u > 0, \end{cases} \qquad I_0(u) = \begin{cases} 0 & u = 0 \\ \infty & u \neq 0, \end{cases}$$

i.e. giving an infinite penalty when any constraint is violated. Instead of these infinite penalty constraints (which are hard to optimize), we replace the constraints with a set of soft linear constraints, as shown in Fig. 2.2. It is now clear why $\lambda$ must be positive for the inequality constraint: a negative $\lambda$ would not yield a lower bound. Note also that as well as being penalized for $f_i > 0$, the linear lower bounds reward us for achieving $f_i < 0$. This is illustrated in Fig. 2.3.

### 2.1.2 The dual problem

The dual problem attempts to find the best lower bound $g(\lambda, \nu)$ on the optimal solution $p^*$ of (2.1). This results in the **dual problem**

$$\begin{aligned} \text{maximize} \quad & g(\lambda, \nu) \\ \text{subject to} \quad & \lambda \succeq 0. \end{aligned} \tag{2.3}$$

Denote by $(\lambda^*, \nu^*)$ the arguments optimizing (2.3) and by $d^*$ the optimal value of the dual problem. Note that (2.3) is always a convex optimization problem, since the function being maximized is concave and the constraint set is convex. The property of **weak duality** is immediate:

$$d^* \leq p^*.$$

The difference $p^* - d^*$ is called the **optimal duality gap**. If the duality gap is zero, then **strong duality** holds:

$$d^* = p^*.$$

Figure 2.3: Illustration of the Lagrangian on a simple problem with one inequality constraint (from [5, Fig. 5.1]).

Conditions under which strong duality holds are called **constraint qualifications**. As an important case: strong duality holds if the primal problem is convex,[2] i.e. of the form

$$
\begin{aligned}
\text{minimize} \quad & f_0(x) \\
\text{subject to} \quad & f_i(x) \leq 0 \qquad\qquad\qquad\qquad i = 1, \ldots, n \qquad\qquad (2.4) \\
& Ax = b
\end{aligned}
$$

for convex $f_0, \ldots, f_m$, *and* if **Slater's condition** holds: there exists some *strictly* feasible point[3] $\tilde{x} \in \mathrm{relint}(\mathcal{D})$ such that

$$
f_i(\tilde{x}) < 0 \quad i = 1, \ldots, m \quad A\tilde{x} = b.
$$

A weaker version of Slater's condition is sufficient for strong convexity when some of the constraint functions $f_1, \ldots, f_k$ are affine (note the inequality constraints are no longer strict):

$$
f_i(\tilde{x}) \leq 0 \quad i = 1, \ldots, k \quad f_i(\tilde{x}) < 0 \quad i = k+1, \ldots, m \quad A\tilde{x} = b.
$$

A proof of this result is given in [5, Section 5.3.2].

---

[2] Strong duality can also hold for non-convex problems: see e.g. [5, p. 229].

[3] We denote by $\mathrm{relint}(\mathcal{D})$ the relative interior of the set $\mathcal{D}$. This looks like the interior of the set, but is non-empty even when the set is a subspace of a larger space. See [5, Section 2.1.3] for the formal defintion.

## 2.1.3 A saddlepoint/game characterization of weak and strong duality

In this section, we ignore equality constraints for ease of discussion. We write the solution to the primal problem as an optimization

$$
\begin{aligned}
\sup_{\lambda \succeq 0} L(x, \lambda) &= \sup_{\lambda \succeq 0} \left( f_0(x) + \sum_{i=1}^{m} \lambda_i f_i(x) \right) \\
&= \begin{cases} f_0(x) & \text{if } f_i(x) \leq 0, \ i = 1, \dots, m \\ \infty & \text{otherwise.} \end{cases}
\end{aligned}
$$

In other words, we recover the primal problem when the inequality constraint holds, and get infinity otherwise. We can therefore write

$$
p^* = \inf_x \sup_{\lambda \succeq 0} L(x, \lambda).
$$

We already know

$$
d^* = \sup_{\lambda \succeq 0} \inf_x L(x, \lambda).
$$

Weak duality therefore corresponds to the **max-min inequality**:

$$
\sup_{\lambda \succeq 0} \inf_x L(x, \lambda) \leq \inf_x \sup_{\lambda \succeq 0} L(x, \lambda). \tag{2.5}
$$

which holds for general functions, and not just $L(x, \lambda)$. Strong duality occurs at a saddlepoint, and the inequality becomes an equality.

There is also a game interpretation: $L(x, \lambda)$ is a sum that must be paid by the person adjusting $x$ to the person adjusting $\lambda$. On the right hand side of (2.5), player $x$ plays first. Knowing that player 2 ($\lambda$) will maximize their return, player 1 ($x$) chooses their setting to give player 2 the worst possible options over all $\lambda$. The max-min inequality says that whoever plays second has the advantage.

## 2.1.4 Optimality conditions

If the primal is equal to the dual, we can make some interesting observations about the duality constraints. Denote by $x^*$ the optimum solution of the original problem (the minimum of $f_0$ under its constraints), and by $(\lambda^*, \nu^*)$ the solutions to the dual. Then

$$
\begin{aligned}
f_0(x^*) &= g(\lambda^*, \nu^*) \\
&\underset{(a)}{=} \inf_{x \in \mathcal{D}} \left( f_0(x) + \sum_{i=1}^{m} \lambda_i^* f_i(x) + \sum_{i=1}^{r} \nu_i^* h_i(x) \right) \\
&\underset{(b)}{\leq} f_0(x^*) + \sum_{i=1}^{m} \lambda_i^* f_i(x^*) + \sum_{i=1}^{r} \nu_i^* h_i(x^*) \\
&\leq f_0(x^*),
\end{aligned}
$$

where in (a) we use the definition of $g$, in (b) we use that $\inf_{x \in \mathcal{D}}$ of the expression in the parentheses is necessarily no greater than its value at $x^*$, and the last line we use that at $(x^*, \lambda^*, \nu^*)$ we have $\lambda^* \succeq 0$, $f_i(x^*) \leq 0$, and $h_i(x^*) = 0$. From this chain of reasoning, it follows that

$$\sum_{i=1}^{m} \lambda_i^* f_i(x^*) = 0, \tag{2.6}$$

which is the condition of **complementary slackness**. This means

$$\lambda_i^* > 0 \quad \Longrightarrow \quad f_i(x^*) = 0,$$
$$f_i(x^*) < 0 \quad \Longrightarrow \quad \lambda_i^* = 0.$$

Consider now the case where the functions $f_i, h_i$ are differentiable, and the duality gap is zero. Since $x^*$ minimizes $L(x, \lambda^*, \nu^*)$, the derivative at $x^*$ should be zero,

$$\nabla f_0(x^*) + \sum_{i=1}^{m} \lambda_i^* \nabla f_i(x^*) + \sum_{i=1}^{r} \nu_i^* \nabla h_i(x^*) = 0.$$

We now gather the various conditions for optimality we have discussed. The **KKT conditions** for the primal and dual variables $(x, \lambda, \nu)$ are

$$
\begin{aligned}
f_i(x) &\leq 0, \ i = 1, \ldots, m, \\
h_i(x) &= 0, \ i = 1, \ldots, r, \\
\lambda_i &\geq 0, \ i = 1, \ldots, m, \\
\lambda_i f_i(x) &= 0, \ i = 1, \ldots, m, \\
\nabla f_0(x) + \sum_{i=1}^{m} \lambda_i \nabla f_i(x) + \sum_{i=1}^{r} \nu_i \nabla h_i(x) &= 0.
\end{aligned}
$$

If a convex optimization problem with differentiable objective and constraint functions satisfies Slater's conditions, then the KKT conditions are necessary and sufficient for global optimality.

## 2.2 Support vector classification

### 2.2.1 The linearly separable case

We first consider the problem of classifying two clouds of points, where there exists a hyperplane which linearly separates one cloud from the other without error. This is illustrated in Figure 2.4 for a 2-dimensional classification problem. As can be seen, there are infinitely many possible hyperplanes that solve this problem: the question is then: which one to choose? The principle behind support vector machines is that we choose the one which has the largest **margin**, which is defined as twice the *smallest* distance from each class to the separating hyperplane (see Figure 2.4).

Figure 2.4: The linearly separable case. There are many linear separating hyperplanes, but only one maximum margin separating hyperplane.

This problem can be expressed as follows:[4]

$$\max_{w,b} (\text{margin}) = \max_{w,b} \left( \frac{2}{\|w\|} \right) \tag{2.8}$$

subject to

$$\begin{cases} \min \left( w^\top x_i + b \right) = 1 & i \ : \ y_i = +1, \\ \max \left( w^\top x_i + b \right) = -1 & i \ : \ y_i = -1. \end{cases} \tag{2.9}$$

The resulting classifier is

$$y = \text{sign}(w^\top x + b),$$

where sign takes value $+1$ for a positive argument, and $-1$ for a negative argument (its value at zero is not important, since for non-pathological cases we will not need to evaluate it there). We can rewrite to obtain

$$\max_{w,b} \frac{1}{\|w\|} \quad \text{or} \quad \min_{w,b} \|w\|^2$$

subject to

$$y_i(w^\top x_i + b) \geq 1. \tag{2.10}$$

---

[4]It's easy to see why the equation below is the margin (the distance between the positive and negative classes): consider two points exactly opposite each other and located on the margins, such that $(x_i - x_j) = \beta w$ for some scalar $\beta$ (where we recall $w$ is orthogonal to the decision boundary, hence aligned with $x_i - x_j$). Then the distance between them (which is the width of the margin) is

$$\|x_i - x_j\| = (x_i - x_j)^\top \frac{(x_i - x_j)}{\|x_i - x_j\|} = (x_i - x_j)^\top \frac{w}{\|w\|} \tag{2.7}$$

Subtracting the two equations in the constraints (2.9) from each other, we get

$$w^\top (x_i - x_j) = 2.$$

Substituting this into (2.7) proves the result.

$(1 - \alpha)_+$

$\mathbb{I}(\alpha < 0)$

$\alpha$

Figure 2.5: The hinge loss is an upper bound on the 0/1 loss.

## 2.2.2 When no linear separator exists (or we want a larger margin)

If the classes are not linearly separable, we may wish to allow a certain number of errors in the classifier (points within the margin, or even on the wrong side of the decision boudary). We therefore want to trade off such "margin errors" vs maximising the margin. Ideally, we would optimise

$$\min_{w,b} \left( \frac{1}{2} \|w\|^2 + C \sum_{i=1}^{n} \mathbf{1} \left\{ y_i \left( w^\top x_i + b \right) < 1 \right\} \right),$$

where $C$ controls the tradeoff between maximum margin and loss (the factor of $1/2$ is to simplify the algebra later, and is not important: we can adjust $C$ accordingly). This is a combinatorial optimization problem, which would be very expensive to solve. Instead, we replace the indicator function with a convex upper bound,

$$\min_{w,b} \left( \frac{1}{2} \|w\|^2 + C \sum_{i=1}^{n} h \left( y_i \left( w^\top x_i + b \right) \right) \right).$$

We use the **hinge loss**,

$$h(\alpha) = (1 - \alpha)_+ = \begin{cases} 1 - \alpha & 1 - \alpha > 0 \\ 0 & \text{otherwise.} \end{cases}$$

although obviously other choices are possible (e.g. a quadratic upper bound). See Figure 2.5.

Figure 2.6: The nonseparable case. Note the red point which is a distance $\xi/\|w\|$ from the margin.

Substituting in the hinge loss, we get

$$\min_{w,b} \left( \frac{1}{2}\|w\|^2 + C \sum_{i=1}^{n} h\left( y_i \left( w^\top x_i + b \right) \right) \right).$$

or equivalently the constrained problem

$$\min_{w,b,\xi} \left( \frac{1}{2}\|w\|^2 + C \sum_{i=1}^{n} \xi_i \right) \tag{2.11}$$

subject to[5]

$$\xi_i \geq 0 \qquad y_i \left( w^\top x_i + b \right) \geq 1 - \xi_i$$

(compare with (2.10)). See Figure 2.6. This formulation is known as the $C$-SVM.

Now let's write the Lagrangian for this problem, and solve it.

$$L(w,b,\xi,\alpha,\lambda) = \frac{1}{2}\|w\|^2 + C\sum_{i=1}^{n}\xi_i + \sum_{i=1}^{n}\alpha_i \left( 1 - y_i \left( w^\top x_i + b \right) - \xi_i \right) + \sum_{i=1}^{n}\lambda_i(-\xi_i) \tag{2.12}$$

with dual variable constraints

$$\alpha_i \geq 0, \qquad \lambda_i \geq 0.$$

We minimize wrt the primal variables $w$, $b$, and $\xi$.

---

[5]To see this, we can write it as $\xi_i \geq 1 - y_i \left( w^\top x_i + b \right)$. Thus either $\xi_i = 0$, and $y_i \left( w^\top x_i + b \right) \geq 1$ as before, or $\xi_i > 0$, in which case to minimize (2.11), we'd use the smallest possible $\xi_i$ satisfying the inequality, and we'd have $\xi_i = 1 - y_i \left( w^\top x_i + b \right)$.

Derivative wrt $w$:

$$\frac{\partial L}{\partial w} = w - \sum_{i=1}^{n} \alpha_i y_i x_i = 0 \qquad w = \sum_{i=1}^{n} \alpha_i y_i x_i. \tag{2.13}$$

Derivative wrt $b$:

$$\frac{\partial L}{\partial b} = \sum_i y_i \alpha_i = 0. \tag{2.14}$$

Derivative wrt $\xi_i$:

$$\frac{\partial L}{\partial \xi_i} = C - \alpha_i - \lambda_i = 0 \qquad \alpha_i = C - \lambda_i. \tag{2.15}$$

We can replace the final constraint by noting $\lambda_i \geq 0$, hence

$$\alpha_i \leq C.$$

Before writing the dual, we look at what these conditions imply about the scalars $\alpha_i$ that define the solution (2.13) due to complementary slackness.

**Non-margin SVs:** $\alpha_i = C > 0$:

1. We immediately have $1 - \xi_i = y_i \left( w^\top x_i + b \right)$.

2. Also, from condition $\alpha_i = C - \lambda_i$, we have $\lambda_i = 0$, hence $\xi_i \geq 0$.

**Margin SVs:** $0 < \alpha_i < C$:

1. We again have $1 - \xi_i = y_i \left( w^\top x_i + b \right)$

2. This time, from $\alpha_i = C - \lambda_i$, we have $\lambda_i > 0$, hence $\xi_i = 0$.

**Non-SVs:** $\alpha_i = 0$

1. This time we have: $y_i \left( w^\top x_i + b \right) \geq 1 - \xi_i$

2. From $\alpha_i = C - \lambda_i$, we have $\lambda_i > 0$, hence $\xi_i = 0$.

This means that the solution is *sparse*: all the points which are not either on the margin, or "margin errors", contribute nothing to the solution. In other words, only those points on the decision boundary, or which are margin errors, contribute. Furthermore, the influence of the non-margin SVs is bounded, since their weight cannot exceed $C$: thus, severe outliers will not overwhelm the solution.

We now write the dual function, by substituting equations (2.13), (2.14), and (2.15) into (2.12), to get

$$
\begin{aligned}
g(\alpha, \lambda) &= \frac{1}{2}\|w\|^2 + C\sum_{i=1}^{n}\xi_i + \sum_{i=1}^{n}\alpha_i\left(1 - y_i\left(w^\top x_i + b\right) - \xi_i\right) + \sum_{i=1}^{n}\lambda_i(-\xi_i) \\
&= \frac{1}{2}\sum_{i=1}^{m}\sum_{j=1}^{m}\alpha_i\alpha_j y_i y_j x_i^\top x_j + C\sum_{i=1}^{m}\xi_i - \sum_{i=1}^{m}\sum_{j=1}^{m}\alpha_i\alpha_j y_i y_j x_i^\top x_j - b\underbrace{\sum_{i=1}^{m}\alpha_i y_i}_{0} \\
&\quad + \sum_{i=1}^{m}\alpha_i - \sum_{i=1}^{m}\alpha_i\xi_i - \sum_{i=1}^{m}(C - \alpha_i)\xi_i \\
&= \sum_{i=1}^{m}\alpha_i - \frac{1}{2}\sum_{i=1}^{m}\sum_{j=1}^{m}\alpha_i\alpha_j y_i y_j x_i^\top x_j.
\end{aligned}
$$

Thus, our goal is to maximize the dual,

$$
g(\alpha) = \sum_{i=1}^{m}\alpha_i - \frac{1}{2}\sum_{i=1}^{m}\sum_{j=1}^{m}\alpha_i\alpha_j y_i y_j x_i^\top x_j,
$$

subject to the constraints

$$
0 \le \alpha_i \le C, \quad \sum_{i=1}^{n} y_i\alpha_i = 0.
$$

So far we have defined the solution for $w$, but not for the offset $b$. This is simple to compute: for the margin SVs, i.e., those $x_i$ for which $0 < \alpha_i < C$, we have $1 = y_i\left(w^\top x_i + b\right)$. Thus, we can obtain $b$ from any of these, or take an average for greater numerical stability.

### 2.2.3 The $\nu$-SVM

It can be hard to interpret $C$. Therefore we modify the formulation to get a more intuitive parameter. Again, we drop $b$ for simplicity. Solve

$$
\min_{w,\rho,\xi}\left(\frac{1}{2}\|w\|^2 - \nu\rho + \frac{1}{n}\sum_{i=1}^{n}\xi_i\right)
$$

subject to

$$
\begin{aligned}
\rho &\ge 0 \\
\xi_i &\ge 0 \\
y_i w^\top x_i &\ge \rho - \xi_i,
\end{aligned}
$$

where we see that we now optimize the margin width $\rho$. Thus, rather than choosing $C$, we now choose $\nu$ as a hyperparameter; the meaning of the latter will become clear shortly.

The Lagrangian is

$$\frac{1}{2}\|w\|_{\mathcal{H}}^2 + \frac{1}{n}\sum_{i=1}^{n}\xi_i - \nu\rho + \sum_{i=1}^{n}\alpha_i\left(\rho - y_i w^\top x_i - \xi_i\right) + \sum_{i=1}^{n}\beta_i(-\xi_i) + \gamma(-\rho)$$

for $\alpha_i \geq 0$, $\beta_i \geq 0$, and $\gamma \geq 0$. Differentiating wrt each of the primal variables $w$, $\xi$, $\rho$, and setting to zero, we get

$$w = \sum_{i=1}^{n}\alpha_i y_i x_i$$

$$\alpha_i + \beta_i = \frac{1}{n} \tag{2.16}$$

$$\nu = \sum_{i=1}^{n}\alpha_i - \gamma \tag{2.17}$$

From $\beta_i \geq 0$, equation (2.16) implies

$$0 \leq \alpha_i \leq n^{-1}.$$

From $\gamma \geq 0$ and (2.17), we get

$$\nu \leq \sum_{i=1}^{n}\alpha_i.$$

We typically have $\rho > 0$ at the global solution (i.e. non-zero margin) and hence $\gamma = 0$, and (2.17) becomes

$$\sum_{i=1}^{n}\alpha_i = \nu. \tag{2.18}$$

Complementary slackness conditions now lead to a very convenient interpretation of parameter $\nu$. In particular, if we denote by $N(\alpha)$ the set of non-margin support vectors, i.e. margin errors, and by $M(\alpha)$ the set of margin support vectors, then (problem sheet):

$$\frac{|N(\alpha)|}{n} \leq \nu \leq \frac{|N(\alpha)| + |M(\alpha)|}{n}.$$

Thus $\nu$ corresponds to an upper bound on the proportion of margin errors and a lower bound on the proportion of the overall number of support vectors - tuning $\nu$ is hence much more interpretable than tuning $C$.

Substituting into the Lagrangian, we can also obtain the dual formulation of $\nu$-SVM, i.e.

$$\frac{1}{2}\sum_{i=1}^{m}\sum_{j=1}^{m}\alpha_i\alpha_j y_i y_j x_i^\top x_j + \frac{1}{n}\sum_{i=1}^{n}\xi_i - \rho\nu - \sum_{i=1}^{m}\sum_{j=1}^{m}\alpha_i\alpha_j y_i y_j x_i^\top x_j + \sum_{i=1}^{n}\alpha_i\rho - \sum_{i=1}^{n}\alpha_i\xi_i$$

$$- \sum_{i=1}^{n}\left(\frac{1}{n} - \alpha_i\right)\xi_i - \rho\left(\sum_{i=1}^{n}\alpha_i - \nu\right)$$

$$= -\frac{1}{2}\sum_{i=1}^{m}\sum_{j=1}^{m}\alpha_i\alpha_j y_i y_j x_i^\top x_j$$

Thus, we must maximize

$$g(\alpha) = -\frac{1}{2} \sum_{i=1}^{m} \sum_{j=1}^{m} \alpha_i \alpha_j y_i y_j x_i^\top x_j,$$

subject to

$$\sum_{i=1}^{n} \alpha_i \geq \nu \qquad 0 \leq \alpha_i \leq \frac{1}{n}.$$

# 3 Kernel Methods

## 3.1 Feature Maps and Feature Spaces

Kernel methods are a versatile algorithmic framework which allows construction of non-linear machine learning algorithms (for a variety of both supervised and unsupervised learning tasks: clustering, dimensionality reduction, classification, regression) by employing linear tools in a nonlinearly transformed feature space. Let us first recall the definition of an abstract inner product, which is central to kernel methods.

**Definition 4.** Let $\mathcal{H}$ be a vector space over $\mathbb{R}$. A function $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H} \times \mathcal{H} \to \mathbb{R}$ is said to be an **inner product** on $\mathcal{H}$ if

1. $\langle \alpha_1 f_1 + \alpha_2 f_2, g \rangle_{\mathcal{H}} = \alpha_1 \langle f_1, g \rangle_{\mathcal{H}} + \alpha_2 \langle f_2, g \rangle_{\mathcal{H}}$

2. $\langle f, g \rangle_{\mathcal{H}} = \langle g, f \rangle_{\mathcal{H}}$

3. $\langle f, f \rangle_{\mathcal{H}} \geq 0$ and $\langle f, f \rangle_{\mathcal{H}} = 0$ if and only if $f = 0$.

We can define a **norm** using the inner product as $\|f\|_{\mathcal{H}} := \sqrt{\langle f, f \rangle_{\mathcal{H}}}$. A **Hilbert space** is a vector space on which an inner product is defined, along with an additional technical condition.[1] We are now ready to define the notion of a *kernel*.

**Definition 5.** Let $\mathcal{X}$ be a non-empty set. A function $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is called a **kernel** if there exists a Hilbert space $\mathcal{H}$ and a map $\varphi : \mathcal{X} \to \mathcal{H}$ such that $\forall x, x' \in \mathcal{X}$,

$$k(x, x') := \left\langle \varphi(x), \varphi(x') \right\rangle_{\mathcal{H}}.$$

We will call the Hilbert space associated with kernel $k$ a **feature space** and the map $\varphi$ will be called a **feature map**. Note that we imposed almost no conditions on $\mathcal{X}$: in particular, we do not require there to be an inner product defined on the elements of $\mathcal{X}$. The case of text documents is an instructive example: one cannot take an inner product between two books, but can take an inner product between vector-valued features of the text in those books.

Clearly, a single kernel can correspond to multiple pairs of underlying feature maps and feature spaces. For a simple example, consider $\mathcal{X} := \mathbb{R}^p$:

$$\phi_1(x) = x \qquad \text{and} \qquad \phi_2(x) = \left[ \frac{x_1}{\sqrt{2}}, \cdots, \frac{x_p}{\sqrt{2}}, \frac{x_1}{\sqrt{2}}, \cdots, \frac{x_p}{\sqrt{2}} \right]^{\top}.$$

---

[1]Specifically, a Hilbert space must be *complete*, i.e. it must contain the limits of all Cauchy sequences with respect to the norm defined by its inner product.

Both $\phi_1$ and $\phi_2$ are valid feature maps (with feature spaces $\mathcal{H}_1 = \mathbb{R}^p$ and $\mathcal{H}_2 = \mathbb{R}^{2p}$) of kernel $k(x, x') = x^\top x'$.

It turns out that all kernel functions (defined as inner products between some features) are *positive definite*.

**Definition 6.** A symmetric function $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is **positive definite** if $\forall n \geq 1$, $\forall (a_1, \ldots a_n) \in \mathbb{R}^n$, $\forall (x_1, \ldots, x_n) \in \mathcal{X}^n$,

$$\sum_{i=1}^{n} \sum_{j=1}^{n} a_i a_j k(x_i, x_j) \geq 0.$$

The function $k(\cdot, \cdot)$ is **strictly positive definite** if for mutually distinct $x_i$, the equality holds only when all the $a_i$ are zero.[2]

Every inner product is a positive definite function, and so is every inner product between feature maps.

**Lemma 7.** *Let $\mathcal{H}$ be any Hilbert space, $\mathcal{X}$ a non-empty set and $\phi : \mathcal{X} \to \mathcal{H}$. Then $k(x, y) := \langle \phi(x), \phi(y) \rangle_{\mathcal{H}}$ is a positive definite function.*

*Proof.*

$$\begin{aligned}
\sum_{i=1}^{n} \sum_{j=1}^{n} a_i a_j k(x_i, x_j) &= \sum_{i=1}^{n} \sum_{j=1}^{n} \langle a_i \phi(x_i), a_j \phi(x_j) \rangle_{\mathcal{H}} \\
&= \left\| \sum_{i=1}^{n} a_i \phi(x_i) \right\|_{\mathcal{H}}^2 \geq 0.
\end{aligned}$$

$\square$

## 3.2 Reproducing Kernel Hilbert Spaces

We have introduced the notation of feature spaces, and kernels on these feature spaces. What's more, we've determined that these kernels are positive definite. In this section, we use these kernels to define *functions* on $\mathcal{X}$. The space of such functions is known as a **reproducing kernel Hilbert space** (RKHS).

**Definition 8.** Let $\mathcal{H}$ be a *Hilbert space of functions $f : \mathcal{X} \to \mathbb{R}$* defined on a non-empty set $\mathcal{X}$. A function $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is called a **reproducing kernel** of $\mathcal{H}$ if it satisfies

- $\forall x \in \mathcal{X}$, $k_x = k(\cdot, x) \in \mathcal{H}$,

- $\forall x \in \mathcal{X}, \forall f \in \mathcal{H}$, $\langle f, k(\cdot, x) \rangle_{\mathcal{H}} = f(x)$ (*the reproducing property*).

If $\mathcal{H}$ has a reproducing kernel, it is called a reproducing kernel Hilbert space (RKHS).

---

[2]The corresponding terminology used for matrices is "positive semi-definite" vs "positive definite".

In particular, note that for any $x, y \in \mathcal{X}$, reproducing kernel satisfies $k(x, y) = \langle k(\cdot, y), k(\cdot, x) \rangle_{\mathcal{H}} = \langle k(\cdot, x), k(\cdot, y) \rangle_{\mathcal{H}}$. Thus, reproducing kernel is clearly a kernel, i.e. an inner product between features with a feature space $\mathcal{H}$ and a feature map $\phi \colon x \mapsto k(\cdot, x)$. This way of writing feature mapping is called the **canonical feature map**. Note that these features are not specified explicitly in a vector form, but rather as functions on $\mathcal{X}$.

We have seen that any reproducing kernel is a kernel and that every kernel is a positive definite function. Remarkably, the **Moore-Aronszajn theorem** [3] shows that *for every positive definite function $k$, there exists a unique RKHS with kernel $k$.* The theorem is outside of the scope of this course, but it provides an insight into the structure of the RKHS corresponding to $k$. It turns out RKHS can be written as $\overline{\text{span} \{k(\cdot, x) : x \in \mathcal{X}\}}$, i.e. the space of all linear combinations of canonical features, completed with respect to an inner product on these linear combinations defined as

$$\left\langle \sum_{i=1}^{r} \alpha_i k(\cdot, x_i), \sum_{j=1}^{s} \beta_j k(\cdot, y_j) \right\rangle := \sum_{i=1}^{r} \sum_{j=1}^{s} \alpha_i \beta_j k(x_i, y_j).$$

Thus, all three notions: (1) reproducing kernel, (2) kernel as inner product between features and (3) positive definite function, are equivalent. Recall that the feature space of a kernel is not unique - but its RKHS (feature space as a space of functions) is - we will henceforth denote the RKHS of kernel $k$ by $\mathcal{H}_k$. For example, for the **linear kernel** $k(x, y) = x^\top y$ considered earlier, many possible feature representations exist but the canonical feature representation that associates to each $x$ the function $k(\cdot, x) \colon y \mapsto x^\top y$ is what determines the structure of its RKHS. In particular, linear kernel $k(x, y) = x^\top y$ corresponds to the RKHS $\mathcal{H}_k$ which is the space of all linear functions $f(x) = w^\top x$ (*why?*).

## 3.3 Representer Theorem

Now that we have defined an RKHS, we can consider it as a hypothesis class for empirical risk minimisation (ERM). In particular, we are looking for the function $f^*$ in the RKHS $\mathcal{H}_k$ which solves the regularised ERM problem

$$\min_{f \in \mathcal{H}_k} \hat{R}(f) + \Omega\left(\|f\|_{\mathcal{H}_k}^2\right),$$

for empirical risk $\hat{R}(f) = \frac{1}{n} \sum_{i=1}^{n} L(y_i, f(x_i), x_i)$, a loss function $L \colon \mathcal{Y} \times \mathcal{Y} \times \mathcal{X} \to \mathbb{R}_+$ and any non-decreasing function $\Omega$.

**Theorem 9.** *There is a solution to*

$$\min_{f \in \mathcal{H}_k} \hat{R}(f) + \Omega\left(\|f\|_{\mathcal{H}_k}^2\right) \tag{3.1}$$

*that takes the form $f^* = \sum_{i=1}^{n} \alpha_i k(\cdot, x_i)$. If $\Omega$ is strictly increasing, all solutions have this form.*

*Proof.* Let $f$ be any minimiser of (3.1). Denote by $f_s$ the projection of $f$ onto the subspace

$$\text{span}\,\{k(\cdot, x_i):\ i = 1, \ldots, n\}$$

such that

$$f = f_s + f_\perp,$$

where $f_s = \sum_{i=1}^n \alpha_i k(\cdot, x_i)$ and $f_\perp$ is orthogonal to the subspace $\text{span}\,\{k(\cdot, x_i):\ i = 1, \ldots, n\}$. Since

$$\|f\|_{\mathcal{H}_k}^2 = \|f_s\|_{\mathcal{H}_k}^2 + \|f_\perp\|_{\mathcal{H}_k}^2 \geq \|f_s\|_{\mathcal{H}_k}^2,$$

we have

$$\Omega\left(\|f\|_{\mathcal{H}_k}^2\right) \geq \Omega\left(\|f_s\|_{\mathcal{H}_k}^2\right).$$

On the other hand, the individual terms $f(x_i)$ in the loss are given by

$$f(x_i) = \langle f, k(\cdot, x_i)\rangle_{\mathcal{H}_k} = \langle f_s + f_\perp, k(\cdot, x_i)\rangle_{\mathcal{H}_k} = \langle f_s, k(\cdot, x_i)\rangle_{\mathcal{H}_k} = f_s(x_i),$$

so

$$L(y_i, f(x_i), x_i) = L(y_i, f_s(x_i), x_i) \quad \forall i = 1, \ldots, n.$$

and thus empirical risks must be the same: $\hat{R}(f) = \hat{R}(f_s)$. Thus $f_s$ is also a minimiser of (3.1) and if $\Omega$ is strictly increasing, it must be that $f_\perp = 0$. $\qquad\square$

We see that the key parts of the theorem are the fact that the empirical risk only depends on the components of $f$ lying in the subspace spanned by the canonical features and that the regulariser $\Omega(\cdot)$ is minimised when $f = f_s$ (adding additional orthogonal components to the function makes it more complex but does not change the empirical risk). Moreover, if $\Omega$ is strictly increasing, then $\|f_\perp\|_{\mathcal{H}_k} = 0$ is required at the minimum.

## 3.4 Operations with Kernels

Kernels can be combined and modified to get new kernels. For example,

**Lemma 10.** *[Sums of kernels are kernels] Given $\alpha > 0$ and $k$, $k_1$ and $k_2$ all kernels on $\mathcal{X}$, then $\alpha k$ and $k_1 + k_2$ are kernels on $\mathcal{X}$.*

To prove the above, just check *positive definiteness*. Note that a difference between two kernels need not be a kernel: if $k_1(x, x) - k_2(x, x) < 0$, then condition 3 of inner product definition 4 may be violated.

**Lemma 11.** *[Mappings between spaces] Let $\mathcal{X}$ and $\widetilde{\mathcal{X}}$ be non-empty sets, and define a map $A:\ \mathcal{X} \to \widetilde{\mathcal{X}}$. Define the kernel $k$ on $\widetilde{\mathcal{X}}$. Then $k(A(x), A(x'))$ is a kernel on $\mathcal{X}$.*

**Lemma 12.** *[Products of kernels are kernels] Given $k$ on $\mathcal{X}$ and $l$ on $\mathcal{Y}$, then*

$$\kappa\left((x, y), (x', y')\right) = k\left(x, x'\right) l\left(y, y'\right)$$

*is a kernel on $\mathcal{X} \times \mathcal{Y}$. Moreover, if $\mathcal{X} = \mathcal{Y}$, then*

$$\kappa\left(x, x'\right) = k\left(x, x'\right) l\left(x, x'\right)$$

is a kernel on $\mathcal{X}$.

The general proof would require some technical details about Hilbert space tensor products, but the main idea can be understood with some simple linear algebra. We consider the case where $\mathcal{H}$ corresponding to $k$ is $\mathbb{R}^M$, and $\mathcal{G}$ corresponding to $l$ is $\mathbb{R}^N$. Write $k\left(x, x'\right) = \varphi(x)^\top \varphi(x')$ and $l\left(y, y'\right) = \psi(y)^\top \psi(y')$. We will use that a notion of inner product between matrices $A \in \mathbb{R}^{M \times N}$ and $B \in \mathbb{R}^{M \times N}$ is given by

$$\langle A, B \rangle = \text{trace}(A^\top B). \tag{3.2}$$

Then

$$
\begin{aligned}
k\left(x, x'\right) l\left(y, y'\right) &= \varphi(x)^\top \varphi(x') \psi(y')^\top \psi(y) \\
&= \text{tr}(\psi(y)\varphi(x)^\top \varphi(x')\psi(y')^\top) \\
&= \left\langle \varphi(x)\psi(y)^\top, \varphi(x')\psi(y')^\top \right\rangle,
\end{aligned}
$$

thus we can define features $A(x, y) = \varphi(x)\psi(y)^\top$ of the product kernel.

The sum and product rules allow us to define a huge variety of kernels.

**Lemma 13.** *[**polynomial kernel**] Let $x, x' \in \mathbb{R}^p$ for $p \geq 1$, and let $m \geq 1$ be an integer and $c \geq 0$. Then*

$$k(x, x') := \left(\left\langle x, x' \right\rangle + c\right)^m$$

*is a valid kernel.*

To prove: expand out this expression into a sum (with non-negative scalars) of kernels $\langle x, x' \rangle$ raised to integer powers. These individual terms are valid kernels by the product rule.

Can we extend this combination of sum and product rule to sums with infinitely many terms? Consider for example the exponential function applied to an inner product $k(x, x') = \exp\left(\langle x, x' \rangle\right)$. Since addition and multiplication preserve positive definiteness and since all the coefficients in the Taylor series expansion of the exponential function are nonnegative, $k_m(x, x') = \sum_{r=1}^m \frac{\langle x, x' \rangle^r}{r!}$ is a valid kernel $\forall m \in \mathbb{N}$. Fix some $\{\alpha_i\}$ and $\{x_i\}$. Then $A_m = \sum_{i,j} \alpha_i \alpha_j k_m(x_i, x_j) \geq 0 \; \forall m$ since $k_m$ is positive definite. But $A_m \to \sum_{i,j} \alpha_i \alpha_j \exp\left(\langle x_i, x_j \rangle\right)$ as $m \to \infty$, so $\sum_{i,j} \alpha_i \alpha_j \exp\left(\langle x_i, x_j \rangle\right) \geq 0$ as well. Thus, $\exp\left(\langle x, x' \rangle\right)$ is also a valid kernel (it is called **exponential kernel**). We may combine all the results above (*exercise*) to show that the following kernel, in practice widely used and known under various names: **Gaussian kernel**, **RBF kernel**, **squared exponential kernel** or **exponentiated quadratic kernel**, is valid on $\mathbb{R}^p$:

$$k(x, x') := \exp\left(-\frac{1}{2\gamma^2} \left\| x - x' \right\|^2\right).$$

The RKHS of this kernel is infinite-dimensional. Moreover, if the domain $\mathcal{X}$ is a compact subset of $\mathbb{R}^p$, its RKHS is dense in the space of all bounded continuous functions with

respect to the uniform norm. Despite that, since all functions in its RKHS are infinitely differentiable, Gaussian kernel is often considered to be excessively smooth. A less smooth alternative is the **Matérn kernel**, given by

$$k(x, x') = \frac{2^{1-\nu}}{\Gamma(\nu)} \left( \frac{\sqrt{2\nu}}{\gamma} \|x - x'\| \right)^\nu K_\nu \left( \frac{\sqrt{2\nu}}{\gamma} \|x - x'\| \right), \quad \nu > 0, \gamma > 0,$$

where $K_\nu$ is the modified Bessel function of the second kind of order $\nu$. The Matérn kernels corresponding to the values $\nu = s + \frac{1}{2}$ for non-negative integers $s$ take a simpler form, in particular:

- $\nu = 1/2$: $k(x, x') = \exp\left( -\frac{1}{\gamma} \|x - x'\| \right)$,

- $\nu = 3/2$: $k(x, x') = \left( 1 + \frac{\sqrt{3}}{\gamma} \|x - x'\| \right) \exp\left( -\frac{\sqrt{3}}{\gamma} \|x - x'\| \right)$,

- $\nu = 5/2$: $k(x, x') = \left( 1 + \frac{\sqrt{5}}{\gamma} \|x - x'\| + \frac{5}{3\gamma^2} \|x - x'\|^2 \right) \exp\left( -\frac{\sqrt{5}}{\gamma} \|x - x'\| \right)$.

For $\nu = s + \frac{1}{2}$, its RKHS consists of $s + 1$ times differentiable functions with square integrable derivatives of order up to $s + 1$. Moreover, the RKHS norms directly penalize the derivatives of $f$, e.g. for $\nu = 3/2$ and in one dimension, it can be shown that

$$\|f\|_{\mathcal{H}_k}^2 \propto \int f''(x)^2 dx + \frac{6}{\gamma^2} \int f'(x)^2 dx + \frac{9}{\gamma^4} \int f(x)^2 dx.$$

As $\nu \to \infty$, Matérn kernel converges to the Gaussian RBF, i.e. $k(x, x') = \exp\left( -\frac{1}{2\gamma^2} \|x - x'\|^2 \right)$.

Another popular choice is the **rational quadratic kernel** which arises as a scale mixture of Gaussian kernels. In particular, consider Gaussian RBF parametrisation $k_\theta(x, x') = \exp\left( -\theta \|x - x'\|^2 \right)$ and a Gamma density placed on $\theta$, i.e. $p(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{\alpha-1} \exp(-\beta\theta)$, with shape $\alpha$ and rate $\beta$. Then, we define

$$\begin{aligned}
\kappa(x, x') &= \int_0^\infty k_\theta(x, x') p(\theta) d\theta \\
&= \frac{\beta^\alpha}{\Gamma(\alpha)} \int_0^\infty \exp\left( -\theta \left( \|x - x'\|^2 + \beta \right) \right) \theta^{\alpha-1} d\theta \\
&= \frac{\beta^\alpha}{\Gamma(\alpha)} \frac{\Gamma(\alpha)}{\left( \|x - x'\|^2 + \beta \right)^\alpha} \\
&= \left( 1 + \frac{\|x - x'\|^2}{\beta} \right)^{-\alpha}.
\end{aligned}$$

Rational quadratic RKHSs contain functions which vary smoothly across multiple lengthscales. If we write $\beta = 2\alpha\gamma^2$ and let $\alpha \to \infty$ we again recover Gaussian RBF, i.e. $\exp\left( -\frac{1}{2\gamma^2} \|x - x'\|^2 \right)$.

## 3.5 Kernel SVM

We can straightforwardly define a maximum margin classifier, i.e. a Support Vector Machine (SVM) in the RKHS. We write the original hinge loss formulation of the regularized empirical risk minimization (ignoring the offset $b$ here for simplicity[3]):

$$\min_{w \in \mathcal{H}} \left( \frac{1}{2} \|w\|_{\mathcal{H}}^2 + C \sum_{i=1}^{n} \left( 1 - y_i \langle w, k(x_i, \cdot) \rangle_{\mathcal{H}} \right)_+ \right)$$

for the RKHS $\mathcal{H}$ with kernel $k(x, x')$. This **kernel SVM** satisfies the assumption of the representer theorem, so we are looking for the solutions of the form

$$w = \sum_{i=1}^{n} \beta_i k(x_i, \cdot). \tag{3.3}$$

In this case, maximizing the margin in the RKHS is equivalent to minimizing $\|w\|_{\mathcal{H}}^2$: as we have seen, for many RKHSs (e.g. the RKHS corresponding to a Gaussian kernel), this corresponds to enforcing smoothness of the learned functions.

Substituting (3.3) and introducing the $\xi_i$ variables as before, we get

$$\min_{\beta, \xi} \left( \frac{1}{2} \beta^\top K \beta + C \sum_{i=1}^{n} \xi_i \right) \tag{3.4}$$

$$\text{subject to } \xi_i \geq 0 \qquad y_i \sum_{j=1}^{n} \beta_j k(x_i, x_j) \geq 1 - \xi_i$$

where the Gram matrix $K$ has $i, j$th entry $K_{ij} = k(x_i, x_j)$. Thus, the primal variables $w$ are replaced with coefficients $\beta$. Note that the problem remains convex since matrix $K$ is positive definite. With an easy calculation (left for exercise), we can verify that the dual takes the form

$$g(\alpha) = \sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i=1}^{m} \sum_{j=1}^{m} \alpha_i \alpha_j y_i y_j k(x_i, x_j),$$

subject to the constraints

$$0 \leq \alpha_i \leq C,$$

and the decision function takes the form

$$w = \sum_{i=1}^{n} y_i \alpha_i k(x_i, \cdot).$$

This is analogous to the original dual SVM, with inner products replaced with the kernel $k$.

---

[3]Note that it suffices to add a constant feature or equivalently use the kernel $k(x, x') + 1$ to include the offset.

## 3.6 Kernel PCA

Kernel PCA is a popular nonlinear dimensionality reduction technique [20]. Assume we have a dataset $\{x_i\}_{i=1}^n$, where $x_i \in \mathbb{R}^p$. Consider an explicit feature transformation $x \mapsto \varphi(x) \in \mathcal{H}$, and assume that we are interested in performing PCA in the feature space $\mathcal{H}$. Assume that the features $\{\varphi(x_i)\}_{i=1}^n$ are centred. Assume for the moment that the feature space is finite-dimensional, i.e. $\mathcal{H} = \mathbb{R}^M$. Then the $M \times M$ sample covariance matrix in the feature space is given by

$$\mathbf{S} = \frac{1}{n-1} \sum_{i=1}^n \varphi(x_i)\varphi(x_i)^\top = \frac{1}{n-1}\Phi^\top\Phi,$$

where $\Phi \in \mathbb{R}^{n \times M}$ is the feature representation of the data. To perform PCA, recall that we are interested in solving the eigenvalue problem $\mathbf{S}v_m = \lambda_m v_m$, $m = 1, \ldots, M$, and we need the top $k \ll \min\{n, M\}$ eigenvectors $v_m$, $m = 1, \ldots, k$, to construct the PC projections $z_i^{(m)} = v_m^\top \varphi(x_i)$. A property analogous to the representer theorem holds here: whenever $\lambda_m > 0$, the eigenvectors lie in the linear span of feature vectors span $\{\varphi(x_i) : i = 1, \ldots, n\}$, i.e.

$$v_m = \sum_{i=1}^n a_{mi}\varphi(x_i) \tag{3.5}$$

for some scalars $a_{mi}$. To see this, note that

$$\lambda_m v_m = \mathbf{S}v_m = \frac{1}{n-1} \sum_{i=1}^n \varphi(x_i)\left(\varphi(x_i)^\top v_m\right)$$

and since $\lambda_m > 0$, it suffices to take $a_{mi} = \frac{1}{\lambda_m(n-1)}\left(\varphi(x_i)^\top v_m\right)$ and clearly $v_m$ has form (3.5). Thus eigenvectors can also be recovered in the *dual space*. Consider now the $n \times n$ kernel matrix $\mathbf{K}$ with $\mathbf{K}_{ij} = k(x_i, x_j) = \varphi(x_i)^\top\varphi(x_j)$. By substituting $v_m = \sum_{i=1}^n a_{mi}\varphi(x_i)$ back into the eigenvalue problem, we have:

$$\mathbf{S}v_m = \frac{1}{n-1} \sum_{i=1}^n \varphi(x_i) \sum_{\ell=1}^n a_{m\ell}k(x_i, x_\ell) = \lambda_m \sum_{i=1}^n a_{mi}\varphi(x_i).$$

To express the above in terms of the kernel matrix, we project both sides onto $\varphi(x_j)$, for each $j = 1, \ldots, n$. This gives

$$\frac{1}{n-1} \sum_{i=1}^n k(x_j, x_i) \sum_{\ell=1}^n a_{m\ell}k(x_i, x_\ell) = \lambda_m \sum_{i=1}^n a_{mi}k(x_j, x_i), \quad j = 1, \ldots, n,$$

which in matrix notation can be written as

$$\mathbf{K}^2 a_m = \lambda_m(n-1)\mathbf{K}a_m.$$

Assuming that $\mathbf{K}$ is invertible, $a_m$ vectors can be found as the eigenvectors of the kernel matrix $\mathbf{K}$ with corresponding eigenvalues given by $\lambda_m(n-1)$.

But if we simply perform the eigendecomposition of $\mathbf{K}$, we will obtain $n$-dimensional eigenvectors of unit norm, and we are after the $M$-dimensional eigenvectors $v_m$ of $\mathbf{S}$ which have unit norm. We see that $1 = v_m^\top v_m = a_m^\top \mathbf{K} a_m = \lambda_m(n-1)a_m^\top a_m$. Thus, if $u_m$ denotes the $m$-th eigenvector of $\mathbf{K}$ with unit norm, to ensure that $v_m$ has unit norm, we need to rescale $a_m = u_m/\sqrt{\lambda_m(n-1)}$. Now, we have an implicit representation of eigenvectors in terms of their dual coefficients. The PC projections are

$$z_i^{(m)} = v_m^\top \varphi(x_i) = \left( \sum_{j=1}^n a_{mj}\varphi(x_j) \right)^\top \varphi(x_i) = \sum_{j=1}^n a_{mj}k(x_j, x_i),$$

or equivalently, the $m$-th dimension of the PC projections is given by

$$\mathbf{z}^{(m)} = \mathbf{K}a_m = \lambda_m(n-1)a_m = \sqrt{\lambda_m(n-1)}u_m. \tag{3.6}$$

We have seen this before! Note that PC projections can be discovered from the SVD $\Phi = UDV^\top$ as either $\mathbf{Z} = \Phi V$ or $\mathbf{Z} = UD$. The latter expression is exactly (3.6), since $u_m$ are the eigenvectors of kernel matrix $\mathbf{K}$ (i.e. the left singular vectors of the feature matrix $\Phi$) and $D_{mm} = \sqrt{\lambda_m(n-1)}$ (*why?*). But note that the eigendecomposition of $\mathbf{K}$ and these projections do not require explicit feature transformations - thus, all the computation is happening in the dual representation and $\varphi(x_i)$ need not be computed, only the kernel matrix $\mathbf{K}$ with $\mathbf{K}_{ij} = k(x_i, x_j)$. The kernel formalism also allows us to compute the projection $v_m^\top \varphi(\tilde{x})$ of a new (previously unseen) data vector $\tilde{x} \in \mathbb{R}^p$ to the $m$-th kernel principal component using

$$\left( \sum_{i=1}^n a_{mi}\varphi(x_i) \right)^\top \varphi(\tilde{x}) = \sum_{i=1}^n a_{mi}k(x_i, \tilde{x}) = a_m^\top \mathbf{k}_{\tilde{x}},$$

where $\mathbf{k}_{\tilde{x}} = [k(x_1, \tilde{x}), \ldots, k(x_n, \tilde{x})]^\top$, so again no explicit feature transformations are needed.

Recall that the above all assumes that the features are **centred**, i.e. that $\frac{1}{n}\sum_{i=1}^n \varphi(x_i) = 0$, but if we are just given a kernel function $k(x, x')$, there is no reason to believe that the features would be centred. Fortunately, it is straightforward to transform *any* kernel matrix into a centred form. Note that the squared distance matrix in the feature space, i.e. matrix $\mathbf{D}$ for which

$$\mathbf{D}_{ij} = \|\varphi(x_i) - \varphi(x_j)\|_{\mathcal{H}}^2 = k(x_i, x_i) + k(x_j, x_j) - 2k(x_i, x_j)$$

can easily be recovered from the Gram/kernel matrix. In matrix form,

$$\mathbf{D} = \operatorname{diag}(\mathbf{K})\mathbf{1}^\top + \mathbf{1}\operatorname{diag}(\mathbf{K})^\top - 2\mathbf{K}.$$

But distances are invariant to centering and the Gram matrix corresponding to centred features can then also be recovered from the distance matrix (exercise).

## 3.7 Representation of probabilities in RKHS

We have seen that kernel methods effectively work on implicit representations of individual data points, via the canonical feature map $\phi\colon x \mapsto k\left(\cdot, x\right)$, such that every data point is represented as a point in the RKHS $\mathcal{H}_k$. One can similarly represent probability distributions $P$ in the RKHSs by considering the **kernel mean embedding**

$$P \mapsto \mu_k\left(P\right) = \mathbb{E}_{X \sim P} k\left(\cdot, X\right) \in \mathcal{H}_k.$$

This is a potentially infinite-dimensional representation of $P$ akin to a characteristic function of a probability distribution. Kernel mean embedding represents expectations over RKHS:

$$\langle f, \mu_k\left(P\right) \rangle_{\mathcal{H}_k} = \mathbb{E}_{X \sim P} f(X), \quad \forall f \in \mathcal{H}_k$$

and exists whenever $f \mapsto \mathbb{E}_{X \sim P} f(X)$ is a bounded functional. Note that this is always true if the kernel function itself is bounded, i.e. $k(x, y) \leq M < \infty \; \forall x, y$. Namely, by Cauchy-Schwarz

$$\mathbb{E}_{X \sim P} f(X) = \mathbb{E}_{X \sim P} \langle f, k\left(\cdot, X\right) \rangle \leq \|f\|_{\mathcal{H}_k} \mathbb{E}_{X \sim P} \|k\left(\cdot, X\right)\|_{\mathcal{H}_k} \leq \sqrt{M} \|f\|_{\mathcal{H}_k}$$

Such representation imposes a simple Hilbert space structure on probability distributions. In particular, inner products between kernel mean embeddings can be computed as

$$\langle \mu_k\left(P\right), \mu_k\left(Q\right) \rangle_{\mathcal{H}_k} = \mathbb{E}_{X \sim P} \mathbb{E}_{Y \sim Q} k(X, Y).$$

**MMD.** We can easily estimate the (squared) distances between probability measures induced by this RKHS representation since they correspond to simple expectations. Such distances are called **maximum mean discrepancy** (MMD):

$$
\begin{aligned}
\mathrm{MMD}_k^2\left(P, Q\right) &= \|\mu_k\left(P\right) - \mu_k\left(Q\right)\|_{\mathcal{H}_k}^2 \qquad\qquad\qquad (3.7)\\
&= \mathbb{E}_{X, X' \overset{iid}{\sim} P} k(X, X') + \mathbb{E}_{Y, Y' \overset{iid}{\sim} Q} k(Y, Y') - 2\mathbb{E}_{X \sim P, Y \sim Q} k(X, Y),
\end{aligned}
$$

where $X$ and $X'$ denote independent copies of random variables with law $P$, and similarly for $Y$ and $Y'$.

The name MMD comes from the following interpretation: it can also be written as the largest discrepancy between expectations of the unit norm RKHS functions with respect to two distributions (problem sheet):

$$\mathrm{MMD}_k\left(P, Q\right) = \sup_{f \in \mathcal{H}_k \colon \|f\|_{\mathcal{H}_k} \leq 1} \left|\mathbb{E}_{X \sim P} f(X) - \mathbb{E}_{Y \sim Q} f(Y)\right|.$$

As a consequence, the function $f$ where the supremum is attained, which can be shown to be proportional to the difference between embeddings, i.e. $\mu_k\left(P\right) - \mu_k\left(Q\right)$, can be thought of as the **witness function** for the difference between distributions $P$ and
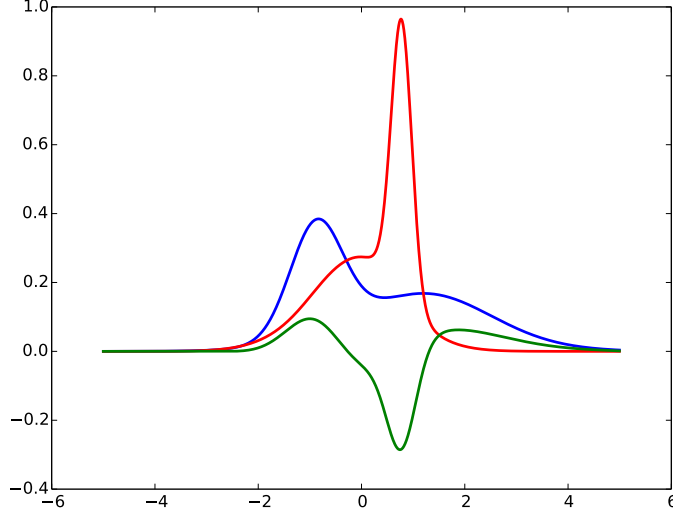
Figure 3.1: Witness function for a difference between two univariate densities

$Q$. An example of such a witness function is shown in green in Fig. 3.1, where $P$ and $Q$ correspond to distributions on the real line whose densities are drawn in blue and red, respectively. We can see that the witness function is large in amplitude where the difference between two densities is large and it can thus be used to discover regions in the space where two distributions disagree.

For a large class of kernels, including Gaussian, Matern family and rational quadratic, MMD is a proper metric on probability distributions, in the sense that $MMD_k(P, Q) = 0$ implies $P = Q$. Such kernels are called **characteristic**. MMD is a popular probability metric, used for nonparametric hypothesis testing [11] and in various machine learning applications, e.g. training deep generative models [9]. Given two samples $\{x_i\}_{i=1}^{n_x} \sim P$ and $\{y_i\}_{i=1}^{n_y} \sim Q$, a simple unbiased estimator of the squared MMD in (3.7) is given by

$$\widehat{\text{MMD}}_k^2(P, Q) = \frac{1}{n_x(n_x - 1)} \sum_{i \neq j} k(x_i, x_j) + \frac{1}{n_y(n_y - 1)} \sum_{i \neq j} k(y_i, y_j) - \frac{2}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} k(x_i, y_j),$$

which can be interpreted as the difference between within-sample average similarity (self-similarity excluded) and the between-sample average similarity.

**HSIC.** Another use of kernel embeddings is in measuring dependence between random variables taking values in some generic domains (e.g. random vectors, strings, or graphs). Recall that for any kernels $k_{\mathcal{X}}$ and $k_{\mathcal{Y}}$ on the respective domains $\mathcal{X}$ and $\mathcal{Y}$, we can define $k = k_{\mathcal{X}} \otimes k_{\mathcal{Y}}$, given by

$$k\left((x, y), (x', y')\right) = k_{\mathcal{X}}(x, x') k_{\mathcal{Y}}(y, y') \tag{3.8}$$

which is a valid kernel on the product domain $\mathcal{X} \times \mathcal{Y}$ by Lemma 12. The tensor notation signifies that the canonical feature map of $k$ is $(x, y) \mapsto k_{\mathcal{X}}(\cdot, x) \otimes k_{\mathcal{Y}}(\cdot, y)$. Here the

feature of pair $(x, y)$, $\varphi_{x,y} = k_{\mathcal{X}}(\cdot, x) \otimes k_{\mathcal{Y}}(\cdot, y)$ is understood as a function on $\mathcal{X} \times \mathcal{Y}$, i.e. $\varphi_{x,y}(x', y') = k_{\mathcal{X}}(x', x)k_{\mathcal{Y}}(y', y)$. The RKHS of the product kernel $k = k_{\mathcal{X}} \otimes k_{\mathcal{Y}}$ is in fact isometric to $\mathcal{H}_{k_{\mathcal{X}}} \otimes \mathcal{H}_{k_{\mathcal{Y}}}$, which can be viewed as the space of **Hilbert-Schmidt operators** between $\mathcal{H}_{k_{\mathcal{Y}}}$ and $\mathcal{H}_{k_{\mathcal{X}}}$. We are now ready to define an RKHS-based measure of dependence between random variables $X$ and $Y$.

**Definition 14.** Let $X$ and $Y$ be random variables on domains $\mathcal{X}$ and $\mathcal{Y}$ (non-empty topological spaces). Let $k_{\mathcal{X}}$ and $k_{\mathcal{Y}}$ be kernels on $\mathcal{X}$ and $\mathcal{Y}$ respectively. The **Hilbert-Schmidt independence criterion** (HSIC) $\Xi_{k_{\mathcal{X}}, k_{\mathcal{Y}}}(X, Y)$ of $X$ and $Y$ is the squared MMD between the joint measure $P_{XY}$ and the product of marginals $P_X P_Y$, computed with the product kernel $k = k_{\mathcal{X}} \otimes k_{\mathcal{Y}}$, i.e.,

$$
\begin{aligned}
\Xi_{k_{\mathcal{X}}, k_{\mathcal{Y}}}(X, Y) &= \|\mu_k(P_{XY}) - \mu_k(P_X P_Y)\|^2_{\mathcal{H}_k} \\
&= \|\mathbb{E}_{XY}[k_{\mathcal{X}}(., X) \otimes k_{\mathcal{Y}}(., Y)] - \mathbb{E}_X k_{\mathcal{X}}(., X) \otimes \mathbb{E}_Y k_{\mathcal{Y}}(., Y)\|^2_{\mathcal{H}_k}.
\end{aligned}
$$

A sufficient condition for HSIC to be well defined is that both kernels $k_{\mathcal{X}}$ and $k_{\mathcal{Y}}$ are bounded. The name of HSIC comes from the operator view of the RKHS $\mathcal{H}_{k_{\mathcal{X}} \otimes k_{\mathcal{Y}}}$. Namely, by repeated use of the reproducing property, it can be verified (**exercise**) that the difference between embeddings $\mu_k(P_{XY}) - \mu_k(P_X P_Y)$ can be identified with the **cross-covariance operator** $C_{XY} : \mathcal{H}_{k_{\mathcal{Y}}} \to \mathcal{H}_{k_{\mathcal{X}}}$ for which

$$
\langle f, C_{XY} g \rangle_{\mathcal{H}_{k_{\mathcal{X}}}} = \text{Cov}\left[f(X)g(Y)\right], \quad \forall f \in \mathcal{H}_{k_{\mathcal{X}}}, g \in \mathcal{H}_{k_{\mathcal{Y}}}.
$$

Note that this is analogous to the finite-dimensional property $f^\top C_{XY} g = \text{Cov}\left[f^\top X, g^\top Y\right]$, where $X$ and $Y$ are random vectors and $C_{XY}$ is their cross-covariance matrix, i.e. $[C_{XY}]_{ij} = \text{Cov}\left[X^{(i)}, Y^{(j)}\right]$. HSIC is then simply the squared Hilbert-Schmidt norm $\|C_{XY}\|^2_{HS}$ of this operator.

To obtain an estimator of the HSIC, we first express it in terms of the expectations of kernels. Starting from the definition, and expanding the Hilbert space norm into inner products:

$$
\begin{aligned}
\Xi_{k_{\mathcal{X}}, k_{\mathcal{Y}}}(X, Y) &= \|\mathbb{E}_{XY}(k_{\mathcal{X}}(\cdot, X) \otimes k_{\mathcal{Y}}(\cdot, Y)) \\
&\quad - \mathbb{E}_X(k(\cdot, X)) \otimes \mathbb{E}_Y(k(\cdot, Y))\|^2_{\mathcal{H}_k} \\
&= \langle \mathbb{E}_{XY}(k_{\mathcal{X}}(\cdot, X) \otimes k_{\mathcal{Y}}(\cdot, Y)), \mathbb{E}_{XY}(k_{\mathcal{X}}(\cdot, X) \otimes k_{\mathcal{Y}}(\cdot, Y)) \rangle_{\mathcal{H}_k} \\
&\quad + \langle \mathbb{E}_X(\mathbb{E}_Y(k_{\mathcal{X}}(\cdot, X) \otimes k_{\mathcal{Y}}(\cdot, Y))), \mathbb{E}_X(\mathbb{E}_Y(k_{\mathcal{X}}(\cdot, X) \otimes k_{\mathcal{Y}}(\cdot, Y))) \rangle_{\mathcal{H}_k} \\
&\quad - 2\langle \mathbb{E}_{XY}(k_{\mathcal{X}}(\cdot, X) \otimes k_{\mathcal{Y}}(\cdot, Y)), \mathbb{E}_X(\mathbb{E}_Y(k_{\mathcal{X}}(\cdot, X) \otimes k_{\mathcal{Y}}(\cdot, Y))) \rangle_{\mathcal{H}_k} \\
&= \mathbb{E}_{XY}\left(\mathbb{E}_{X'Y'}\left(\langle k_{\mathcal{X}}(\cdot, X) \otimes k_{\mathcal{Y}}(\cdot, Y), k_{\mathcal{X}}(\cdot, X') \otimes k_{\mathcal{Y}}(\cdot, Y') \rangle_{\mathcal{H}_k}\right)\right) \\
&\quad + \mathbb{E}_{XX'}\left(\mathbb{E}_{YY'}\left(\langle k_{\mathcal{X}}(\cdot, X) \otimes k_{\mathcal{Y}}(\cdot, Y), k_{\mathcal{X}}(\cdot, X') \otimes k_{\mathcal{Y}}(\cdot, Y') \rangle_{\mathcal{H}_k}\right)\right) \\
&\quad - 2\mathbb{E}_{XY}\left(\mathbb{E}_{X'}\left(\mathbb{E}_{Y'}\left(\langle k_{\mathcal{X}}(\cdot, X) \otimes k_{\mathcal{Y}}(\cdot, Y), k_{\mathcal{X}}(\cdot, X') \otimes k_{\mathcal{Y}}(\cdot, Y') \rangle_{\mathcal{H}_k}\right)\right)\right) \\
&= \mathbb{E}_{XY}\left(\mathbb{E}_{X'Y'}\left(k_{\mathcal{X}}(X, X') k_{\mathcal{Y}}(Y, Y')\right)\right) \\
&\quad + \mathbb{E}_{XX'}\left(k_{\mathcal{X}}(X, X')\right)\mathbb{E}_{YY'}\left(k_{\mathcal{Y}}(Y, Y')\right) \\
&\quad - 2\mathbb{E}_{XY}\left(\mathbb{E}_{X'}\left(k_{\mathcal{X}}(X, X')\right)\mathbb{E}_{Y''}\left(k_{\mathcal{Y}}(Y, Y'')\right)\right)
\end{aligned}
$$

$$(3.9)$$

Here the first expectation is taken over two independent copies $(X, Y)$, $(X', Y') \sim P_{XY}$, the second over two independent $X, X' \sim P_X$ and two independent $Y, Y' \sim P_Y$ and the third over a pair $(X, Y) \sim P_{XY}$ sampled from the joint and an independent pair $X' \sim P_X$, $Y'' \sim P_Y$. Now, given a sample $Z = \{z_i\}_{i=1}^m = \{(x_i, y_i)\}_{i=1}^m$, where each $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$, we can derive an estimator of the HSIC by estimating each of the three terms in the expansion.

Denote for convenience $k_{ij} = k_{\mathcal{X}}(x_i, x_j)$ and $l_{ij} = k_{\mathcal{Y}}(y_i, y_j)$ for $i, j \in \{1, 2, ..., m\}$ and define the kernel matrices $K = (k_{ij})_{i,j=1}^m$ and $L = (l_{ij})_{i,j=1}^m$ (recall that they are symmetric and positive-definite). Following, we estimate:

$$\widehat{\text{first term}} = \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m k_{ij} l_{ij} = \frac{1}{m^2} \text{tr}(KL)$$

$$\widehat{\text{second term}} = \frac{1}{m^4} \left( \sum_{i=1}^m \sum_{j=1}^m k_{ij} \right) \left( \sum_{i=1}^m \sum_{j=1}^m l_{ij} \right)$$

$$= \frac{1}{m^4} \left( 1_m^T K 1_m \right) \left( 1_m^T L 1_m \right)$$

$$\widehat{\text{third term}} = \frac{1}{m^3} \sum_{i=1}^m \sum_{j=1}^m \sum_{q=1}^m k_{ij} l_{iq} = \frac{1}{m^3} 1_m^T K L 1_m$$

$$= \frac{1}{m^3} 1_m^T L K 1_m$$

Here $1_m$ is the vector with m entries equal to 1. Therefore an estimator for the HSIC can be written as:

$$\widehat{\Xi_{k_{\mathcal{X}}, k_{\mathcal{Y}}}}(X, Y) = \frac{1}{m^2} \left( \text{tr}(KL) - \frac{2}{m} 1_m^T K L 1_m + \frac{1}{m^2} \left( 1_m^T K 1_m \right) \left( 1_m^T L 1_m \right) \right)$$

$$= \frac{1}{m^2} \left( \text{tr}(KL) - \frac{1}{m} \text{tr} \left( 1_m 1_m^T K L \right) - \frac{1}{m} \text{tr} \left( K 1_m 1_m^T L \right) \right.$$

$$\left. + \frac{1}{m^2} \text{tr} \left( 1_m 1_m^T K 1_m 1_m^T L \right) \right)$$

$$= \frac{1}{m^2} \text{tr} \left( \left( I - \frac{1}{m} 1_m 1_m^T \right) K \left( I - \frac{1}{m} 1_m 1_m^T \right) L \right)$$

$$= \frac{1}{m^2} \text{tr}(KHLH).$$

Here we used that $\text{tr}(AB) = \text{tr}(BA)$, $\text{tr}(A) = \text{tr}(A^T)$ and that any real number is equal to its own trace. We also defined

$$H := I - \frac{1}{m} 1_m 1_m^T$$

which is the **centering matrix**. Namely, if $A$ is any $m \times m$-matrix, $AH$ centers the rows of $A$ and $HA$ centers the columns of $A$. Note also that $H$ is symmetric and idempotent,

i.e. $H^2 = H$. Hence, $\text{tr}\left((HKH)(HLH)\right) = \text{tr}\left(H(KHLH)\right) = \text{tr}\left(KHLHH\right) = \text{tr}\left(KHLH\right)$.

Recall that the kernel is an inner product between features of the inputs and that inner products are bilinear. Therefore, the matrices $\widetilde{K} = HKH$ and $\widetilde{L} = HLH$ are the kernel matrices for the variables centered in feature space. We therefore arrive at the expression for the estimator:

$$\widehat{\Xi_{k_{\mathcal{X}},k_{\mathcal{Y}}}}(X,Y) = \frac{1}{m^2}\text{tr}\left(\widetilde{K}\widetilde{L}\right), \tag{3.10}$$

which has an intuitive explanation of how it measures the dependence between $X$ and $Y$. Namely, the function $(A,B) \to \text{tr}\left(A^T B\right)$ is an inner product on the vector space of real $m \times m$ matrices. Therefore, our estimate measures the similarity between the (centered) kernel matrices, which in turn measure the "similarity patterns" between the individual observations. If there is some dependence between the $X$ and $Y$, we also expect that the kernel matrices will have a similar structure and hence the inner product between them (and hence our HSIC estimator in (3.10)), will be larger.

# 4 Deep Learning

## 4.1 Why Deep Learning

Consider empirical risk minimisation framework for linear binary classification:

$$\min_{w \in \mathbb{R}^p} \left( \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n L(y_i, w^\top x_i) \right)$$

We have seen how we can generalise the linear method to a non-linear one using a kernel. The resulting ERM objective is:

$$\min_{w \in \mathcal{H}_k} \left( \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n L(y_i, \langle w, k(x_i, \cdot) \rangle_{\mathcal{H}_k}) \right)$$

This implicitly and non-linearly maps inputs $x_i$ into an RKHS $k(x_i, \cdot)$ element. A choice for the kernel corresponds effectively to a choice of features to use, and often domain knowledge is used to make an appropriate choice.

There are a number of problems with this approach:

1. For many problems it may not be obvious what kernel/features to use for best performance.

2. Even if a kernel family can be identified, we still need to perform cross-validation to choose hyperparameters (if a small number), which can be computationally expensive.

3. Kernel methods can be computationally slow. Using the Representer Theorem means that optimisation is based on the Gram matrix which is $n^2$ in size, and can require $O(n^3)$ computational cost to optimise. This is infeasible for large $n$.

From the perspective of ERM, deep learning addresses both the modelling (kernel/feature choice) and computational issues faced by kernel methods. Important ideas are:

1. Instead of "feature engineering", we should do **feature learning** using **end-to-end learning**. That is, as much of the system as possible is learnt by minimising the empirical risk (with regularisation). This ensures that all parts of the system is optimised for what we want it to perform.

2. Simple stochastic gradient descent (SGD) optimisation algorithms allows for scalability to significantly larger datasets and model sizes.

In addition, deep learning has a number of important advantage over past ML frameworks:

1. Domain knowledge can be reflected in the design choices for the neural architecture and other components of the system (including data pre-processing, regularisation and optimiser) for good generalisation properties.

2. Highly developed software and hardware stack enables

   a) Hardware acceleration using GPUs, TPUs and other ASICs, along with distributed computing, which means even faster and more scalable learning.

   b) Modularity: ease of building complex models from a large range of simple modules, as well as ease of building reusable modules.

   c) Differentiable programming: automated learning (using SGD) of models.

   All the above mean that it is much easier to critique models, diagnose problems, and improve models than before.

## 4.2 Basic Deep Learning

The basic idea of deep learning is to construct functions (**neural networks**) as compositions of simpler functions (layers or modules):

$$x^{(0)} = x$$
$$x^{(l)} = f^{(l)}(x^{(l-1)}) \quad \text{for } l = 1, \ldots, m$$
$$f(x) = x^{(m)} = f^{(m)} \circ f^{(m-1)} \circ \cdots f^{(2)} \circ f^{(1)}(x)$$

Each layer (say $l$) can be thought of as producing a **representation** $x^{(l)}$ of the input $x$, with the representation of each subsequent layer being more useful for making the final prediction of the output $y$. Typically each $x^{(l)}$ is high-dimensional, as we want to learn **rich representations** of the input $x$ that are useful to modelling complex input-output relationships.

Each $f^{(l)}$ is a simple function and can have a vector of learnable parameters $\theta^{(l)}$ of length $p_l$. The overall function $f$ has parameters $\theta = [\theta^{(1)}, \ldots, \theta^{(m)}]$ of length $p = \sum_{l=1}^{m} p_l$. We often write $f_\theta$ to make explicit dependence on the parameters $\theta$. The ERM problem is then,

$$\min_{\theta} \quad \mathcal{L}(\theta) \qquad\qquad \mathcal{L}(\theta) = \frac{\lambda}{2}\|\theta\|_2^2 + \frac{1}{n}\sum_{i=1}^{n} L(y_i, f_\theta(x_i)) \qquad (4.1)$$

A key assumption of deep learning is that all functions we work with should be differentiable (except on sets of measure 0), and gradient descent is almost exclusively used for minimising (4.1). Starting with an initial $\theta_0$, gradient descent steps are:

$$\theta_t = \theta_{t-1} - \alpha \frac{d\mathcal{L}}{d\theta}(\theta_{t-1}) \qquad (4.2)$$

where $\alpha$ is a step size.

The gradients can be computed by applying the chain rule multiple times. Letting $L_i = L(y_i, f_\theta(x_i))$, we see that:

$$\frac{d\mathcal{L}}{d\theta} = \lambda\theta + \frac{1}{n}\sum_{i=1}^{n}\frac{\partial\mathcal{L}}{\partial L_i}\frac{\partial L_i}{\partial\theta}$$

$$\frac{\partial L_i}{\partial\theta^{(l)}} = \frac{\partial L_i}{\partial x_i^{(m)}}\frac{\partial x_i^{(m)}}{\partial x_i^{(m-1)}}\cdots\frac{\partial x_i^{(l+1)}}{\partial x_i^{(l)}}\frac{\partial x_i^{(l)}}{\partial\theta^{(l)}}$$

$$= \left(\left(\left(\frac{\partial L_i}{\partial x_i^{(m)}}\frac{\partial x_i^{(m)}}{\partial x_i^{(m-1)}}\right)\cdots\frac{\partial x_i^{(l+1)}}{\partial x_i^{(l)}}\right)\frac{\partial x_i^{(l)}}{\partial\theta^{(l)}}\right) \tag{4.3}$$

$$= \left(\frac{\partial L_i}{\partial x_i^{(m)}}\left(\frac{\partial x_i^{(m)}}{\partial x_i^{(m-1)}}\cdots\left(\frac{\partial x_i^{(l+1)}}{\partial x_i^{(l)}}\frac{\partial x_i^{(l)}}{\partial\theta^{(l)}}\right)\right)\right) \tag{4.4}$$

While the two orderings of computing the product of matrices give the same answer, they differ significant in computational cost. For (4.3), it is $O(\sum_{j=l}^{m-1}d_{j+1}d_j + d_j p_j)$ while for (4.4), it is $O(d_m p_j + \sum_{j=l}^{m-1}d_{j+1}d_j p_l)$, so (4.4) is $p_l$ times more costly. (4.3) computes gradients of the loss with respect to the variables and parameters in a *backward* fashion, and is called **backpropagation**. (4.4) computes the Jacobians in a forward fashion. the term forward-propagation typically refers to the computation of the function $f(x_i)$ in a forward fashion. The gradients with respect to all parameters can be computed in a single backward propagation phase, by computing $\frac{\partial L_i}{\partial x_i^{(l)}}$ and $\frac{\partial L_i}{\partial\theta^{(l)}}$ iteratively for $l = m, m-1, \ldots, 1$, with

$$\frac{\partial L_i}{\partial x_i^{(l)}} = \frac{\partial L_i}{\partial x_i^{(l+1)}}\frac{\partial x_i^{(l+1)}}{\partial x_i^{(l)}}$$

for $l \neq m$.

Note however that back-propagation requires **activations** $x_i^{(l)}$ to be stored in the forward pass to be used in the backward pass, so is more memory intensive than computing the Jacobians in the forward pass.

## 4.3 Computation Graphs and Automatic Differentiation

In the basic example, $f(x)$ is obtained by composing a chain of simpler functions. In general, $f(x)$ can be computed using a **computation graph**. This is a graph, with nodes corresponding either to inputs or operators (ops in short; simpler functions) and directed edges connecting them. We assume that the graph is acyclic. Each input is a multi-dimensional array (often referred to as a **tensor**[1]), and each op is a function $(y_1, \ldots, y_n) = \mathsf{op}(x_1, \ldots, x_m)$ taking a list of input tensors $x_1, \ldots, x_m$ and outputs a list

---

[1]These are not real tensors used in physics, in that we do not distinguish between covariant and contravariant indices.

of tensors $y_1, \ldots, y_n$. An edge $\mathsf{op}_1 \to \mathsf{op}_2$ connecting two ops means that an output of $\mathsf{op}_1$ is used an input of $\mathsf{op}_2$ (technically we need to specify which of the outputs of $\mathsf{op}_1$ and inputs of $\mathsf{op}_2$, but this is typically suppressed for notational convenience).

The implementation of each op needs to construct two things: $\mathsf{op}$ itself, and $\mathsf{vjp}$ which computes the backpropagated gradients. Given tensors $v_1, \ldots, v_n$ of same shape as $y_1, \ldots, y_n$, interpreted as gradients of some objective function with respect to $y_j$, $\mathsf{vjp}$ computes the gradients with respect to the $x_i$'s:

$$\mathsf{vjp}((x_i)_{i=1}^m, (v_j)_{j=1}^n) = \left( \sum_{j=1}^n v_j \cdot \frac{\partial y_j}{\partial x_i} \right)_{i=1}^m$$

where the $i$th entry of the resulting output is the same shape as $x_i$. $\mathsf{vjp}$ stands for vector-Jacobian products (the $v_j$'s are vectors and we multiply them with the Jacobians $\frac{\partial y_j}{\partial x_i}$'s). In packages that do automatic differentiation, there is also a forward mode differentiation, which implements a Jacobian-vector product:

$$\mathsf{jvp}((x_i)_{i=1}^m, (u_i)_{i=1}^m) = \left( \sum_{i=1}^n \frac{\partial y_j}{\partial x_i} \cdot u_i \right)_{j=1}^n$$

Here $u_i$ is a tensor of the same shape as $x_i$ and interpreted as the gradient of $x_i$ with respect to some scalar.

The op-level computations can now be chained together for computations over the whole graph. The forward pass computes each op once its inputs are all given or computed (in a so-called topological order), and the backward pass computes the gradients using the corresponding VJPs, in reverse order on the graph. For forward mode differentiation, both the op and the JCPs are computed together. However the gradients of each $x_i$ with respect to each element of the input nodes that feed into $x_i$ needs to be computed and kept track of.

As before, the overall computational cost for forward mode differentiation is linear in the dimensionality of the inputs of the graph, while backward mode differentiation has computational cost linear in the dimensionality of the outputs. In machine learning the final output is almost always a single dimension (the objective) so backward mode differentiation is much more predominant.

Note that each op in the computation graph can itself be implemented using its own computation graph, by specifying the inputs and outputs and abstracting away the inner ops. Once a procedure is implemented to create a computation graph, this can then be used to create multiple ops that are used in a larger computation graph. This notion of **modularity** is important in allowing very complex computation graphs to be constructed and used in deep learning. In fact the whole endeavour can be thought of as a new way of engineering complex software, called **differentiable programming**, which has wider implications for computer science as a field and driving much research across many fronts.

## 4.4 Basic Modules

There are a huge variety of layers/modules that are in use across deep learning. Here we will describe a selection of the more popular ones.

To make the description clearer, we need to distinguish among 3 types of objects, which I will call: ops, module s and factories. An op is a node in a computation graph. A module is a function can be applied multiple times to different inputs, each time creating a new op. A module can have parameters, with these parameters being shared across all ops created by it. A factory is a procedure that creates modules, each module having its own set of parameters.

For example, consider a factory $\mathsf{Factory}(n, \sigma)$ that generates functions of the form $f : \mathbb{R}^n \to \mathbb{R}^n$, $f(x) = \sigma(Wx)$ where $W \in \mathbb{R}^{n \times n}$ and $\sigma$ is an elementwise non-linearity. The following generates two such functions, and applies each twice to a vector x:

$$\mathsf{module}_1 \sim \mathsf{Factory}(n, \sigma)$$
$$\mathsf{module}_2 \sim \mathsf{Factory}(n, \sigma)$$
$$\mathsf{op}_1 = \mathsf{module}_1(x)$$
$$\mathsf{op}_2 = \mathsf{module}_1(\mathsf{op}_1)$$
$$\mathsf{op}_3 = \mathsf{module}_2(\mathsf{op}_2)$$
$$\mathsf{op}_4 = \mathsf{module}_2(\mathsf{op}_3)$$

The two modules are functions of the same form, but with two different sets of parameters, say $W_1$ and $W_2$. We use $\sim$ in the above as $\mathsf{module}_1$ and $\mathsf{module}_2$ are different objects, with different parameters that are initially randomly drawn from the same distribution given in $\mathsf{Factory}(n, \sigma)$. The computation graph consists of 4 ops, two sets of parameters, and 1 input vector, and produces:

$$\sigma(W_2 \sigma(W_2 \sigma(W_1 \sigma(W_1 x))))$$

- **fully-connected**, **dense**, **feed-forward**, or **linear layer**:

$$f \sim \mathsf{Dense}(m, n)$$
$$f(x) = Wx + b$$

  where $x \in \mathbb{R}^n$ is the input, and $W \in \mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^n$ are parameters called **weights** and **biases** respectively. It is specified by the input and output dimensions $n$ and $m$. In deep learning software frameworks $n$ is typically not directly provided and inferred from the inputs into $f$.

- **nonlinearity**. These are simply take an array as input and apply a nonlinear function to each element of the array. They don't have learnable parameters. For

example:

$$\mathsf{sigmoid}(x) = 1/(1 + \exp(-x))$$
$$\mathsf{tanh}(x) = (\exp(x) - \exp(-x))/(\exp(x) + \exp(-x))$$
$$\mathsf{ReLU}(x) = \max(0, x)$$
$$\mathsf{softplus}(x) = \log(1 + \exp(x))$$
$$\mathsf{swish}(x) = x\,\mathsf{sigmoid}(x)$$
$$\mathsf{ELU}_\alpha(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(\exp(x) - 1) & \text{if } x \le 0 \end{cases}$$

A **softmax** is also a nonlinearity (but not elementwise):

$$\mathsf{softmax}([x_1, \ldots, x_d]) = \left[ \frac{\exp(x_1)}{\sum_{i=1}^d \exp(x_i)}, \ldots, \frac{\exp(x_d)}{\sum_{i=1}^d \exp(x_i)} \right]$$

- **2D convolutional layer** see 2D convolutional layer ($\mathsf{Conv2D}$). This layer takes as input a 3D array $x$ of size $(H, W, C)$ (height, width, and number of channels), and outputs a 3D array $x'$ of size $(H', W', C')$, by convolving $x$ with a convolutional kernel $k$. In the simplest case, $k$ is of size $(h, w, C, C')$, with $h = H - H' + 1$, $w = W - W' + 1$, and:

$$f \sim \mathsf{Conv2D}(h, w, C, C')$$
$$x' = f(x)$$
$$x'_{i'j'c'} = \sum_{i=1}^h \sum_{j=1}^w \sum_{c=1}^C k_{ijlc} x_{i'+i-1, j'+j-1, c, c'}$$

Technically, this is a cross-correlation, not a convolution! The indexing above assumes it starts at 1, following standard vector/matrix convention. In Python indexing starts at 0 so the indices are slightly different. Additional variations include:

- There can be a bias for each output channel.

- There can be padding around edges of input by 0s, so that output has same height and width as input.

- There can be striding, e.g. for 2x3 strides only every 2nd row and 3rd column of $x'$ is computed (and $x'$ is reduced 2x3 times in size).

A **depthwise spatial convolution** applies a separate convolution kernel to each channel of $x$ separately. A **1x1 convolution** or **pointwise convolution** is one whose kernel has height 1 and width 1. It is equivalently to taking the channels at each $(i, j)$ location as a vector, and multiplying by a parameter matrix of size $C' \times C$. There are also 1D and 3D convolutions.

- **max-pooling**: For each channel, this takes the maximum value over the elements of $x$ that would have participated in a 2D depthwise convolution:

$$x' = \mathsf{MaxPool}_{h,w}(x)$$

$$x'_{i'j'c} = \max_{i=1}^{h} \max_{j=1}^{w} x_{i'+i-1,j'+j-1,c}$$

Max-pooling technically can be considered as a nonlinearity. **Average-pooling** applies an average instead of max.

## 4.5 Larger Modules

New and more complex modules can be constructed from simpler ones. Popular examples are:

- **Multi-layer perceptron** (MLP) simply chains together multiple fully-connected layers and nonlinearities. It is specified by the number of fully-connected layers, the sizes of each layer, and nonlinearity.

$$f \sim \mathsf{MLP}(d_m, d_{m-1}, \ldots, d_0, \sigma) \quad \begin{cases} f^{(l)} \sim \mathsf{Dense}(d_l, d_{l-1}) & \text{for } l = 1, \ldots, m \\ f = f^{(m)} \circ \sigma \circ f^{(m-1)} \circ \cdots \circ \sigma \circ f^{(1)} \end{cases}$$

$$f(x) = W_m \sigma(W_{m-1}\sigma(\cdots W_2\sigma(W_1 x + b_1) + b_2 \cdots) + b_{m-1}) + b_m \tag{4.5}$$

The outputs of each $\sigma$ op is a vector, each entry of which is called an **activation**, corresponding to a **hidden unit** of the MLP. The input into the $\sigma$ op is a vector of the **pre-activations**.

- **Recurrent neural networks** (RNNs) are NNs that are applied to time series data. In the simplest setting, we have a sequence of inputs $x_1, x_2, \ldots$ and want to predict a corresponding sequence of outputs $y_1, y_2, \ldots$. We would like each prediction for $y_t$ to depend on the sequence of inputs so far $x_{1:t}$. RNNs do this by producing a sequence of representations $h_t$ of the sequence history $x_{1:t}$. $h_t$ is a function of $x_{1:t}$ that can be computed recursively from $h_{t-1}$ and $x_t$. RNNs use two modules, an encoder $\mathsf{encode}$ and a decoder $\mathsf{decode}$, and, starting with an initial state $h_0$, computes as follows:

$$f \sim \mathsf{RNN}(h_t, \hat{y}_t) \quad = f(h_{t-1}, x_t) \quad \begin{cases} \mathsf{encode} \sim \mathsf{EncoderFactory} \\ \mathsf{decode} \sim \mathsf{DecoderFactory} \\ \quad h_t = \mathsf{encode}(h_{t-1}, x_t) \\ \quad \hat{y}_t = \mathsf{decode}(h_t) \end{cases}$$

We can also think of $h_t$ as the system's memory, with the encoder being the RNN's procedure to update its memory based on new information in $x_t$. In the simplest case both $\mathsf{encode}$ and $\mathsf{decode}$ are dense layers followed by a nonlinearity (or MLPs).

- **Long short-term memory** (LSTM) is a particular type of RNN. It was designed specifically to address gradient explosion/vanishing problems in standard RNNs (see Section 4.6). In addition to $h_t$, it adds another hidden state called the cell state $C_t$. Both $h_t$ and $C_t$ can be thought of as memories of the LSTM, with $C_t$ playing the role of the long term memory, and $h_t$ the short term memory. The LSTM module is defined as follows:

$$
\begin{aligned}
\mathsf{forget\_gate}_t &= \mathsf{sigmoid}(W_f[h_{t-1}, x_t] + b_f) \\
\mathsf{insert\_gate}_t &= \mathsf{sigmoid}(W_i[h_{t-1}, x_t] + b_i) \\
C'_t &= \mathsf{tanh}(W_C[h_{t-1}, x_t] + b_C) \\
C_t &= \mathsf{forget\_gate}_t * C_{t-1} + \mathsf{insert\_gate}_t * C'_t \\
\mathsf{output\_gate}_t &= \mathsf{sigmoid}(W_o[h_{t-1}, x_t] + b_o) \\
h_t &= \mathsf{output\_gate}_t * \mathsf{tanh}(C_t)
\end{aligned}
$$

The basic operation is as follows: given the current hidden state $h_{t-1}$ and input $x_t$, the system decides what parts of the cell state $C_t$ to forget via $\mathsf{forget\_gate}_t$, what new information $C'_t$ to insert into it via $\mathsf{insert\_gate}_t$, and what to output to the new state $h_t$ via the $\mathsf{output\_gate}_t$.

## 4.6 Optimisation

Optimising neural networks is difficult, and a lot of research effort has been devoted to optimisation algorithms and techniques specially tailored to deep learning settings.

Recall the ERM problem (4.1), reproduced here:

$$
\min_\theta \quad \mathcal{L}(\theta) \qquad\qquad \mathcal{L}(\theta) = \frac{\lambda}{2}\|\theta\|_2^2 + \frac{1}{n}\sum_{i=1}^n L(y_i, f_\theta(x_i)) \tag{4.6}
$$

where $f_\theta$ is a neural network with parameters $\theta$. Automatic differentiation tools allows the gradient $\frac{d\mathcal{L}(\theta)}{d\theta}$ to be computed effectively, so gradient based optimisation algorithms can be employed. However there are still many difficulties, both computationally, and in terms of the complexity of navigating the "loss landscape".

The first and more tractable difficulty is computation. Computing the gradient of (4.1) scales linearly in the dataset size $n$. Many datasets of interest in deep learning are too large, such that computing $\frac{d\mathcal{L}}{d\theta}$ is impractical even with modern GPUs and TPUs. The solution is to replace the $\frac{d\mathcal{L}}{d\theta}$ with a stochastic estimate, say $\widehat{\frac{d\mathcal{L}}{d\theta}}$, so the **stochastic gradient descent** (SGD) algorithm is:

$$
\theta_t = \theta_{t-1} - \alpha_t \widehat{\frac{d\mathcal{L}}{d\theta}}(\theta_{t-1}) \tag{4.7}
$$

where $\alpha_t$ not a iteration dependent step size. In expectation, the SGD steps are moving in the "downhill' direction. An unbiased estimate of the gradient $\frac{d\mathcal{L}}{d\theta}(\theta)$ can be obtained

by using a small subsample of the dataset, called a **minibatch**. Let $B$ be a random subset of $\{1, \ldots, n\}$ with $|B| \ll n$. Then:

$$\frac{\widehat{d\mathcal{L}}}{d\theta}(\theta) = \lambda\theta + \frac{1}{|B|} \sum_{i \in B} \frac{dL(y_i, f_\theta(x_i))}{d\theta}$$

Note that the variance of the stochastic gradient scales as $1/|B|$. Typically the dataset is split into $n/|B|$ minibatches, and each is used in turn. Once the whole dataset has been processed once (called an epoch), the split into minibatches is randomised for the next epoch.

SGD is very well studied in the optimisation literature and there is a wealth of powerful results, particularly in well-behaved objectives (e.g. strongly convex and smooth). While some of the intuitions gained are relevant in the deep learning setting, many problems faced in optimising deep learning models are precisely because the objectives are not well-behaved. Unlike linear regression or kernel methods, $f_\theta$ for NNs can be a highly complex nonlinear function of $\theta$, and $\mathcal{L}(\theta)$ is not convex in $\theta$. It can have multi-modality, be ill-conditioned with ridges, valleys and flat plateaus.

As an example of the complexity, consider the gradients for a MLP (4.5) , which are given by (4.3). This is given by a product of a sequence of matrices. If the matrices have eigenvalues that are above 1, there could be directions where the product become very large, this is called **exploding gradients**. Conversely, if there are eigenvalues that are small, the gradients can **vanish**. Such exploding/vanishing gradient problems are particularly severe in deep models with many layers since the resulting gradients are products of many matrices. The chain of encode ops in the RNN can be thought of as an MLP with each time step being a layer, which is why its gradients are so ill-behaved. Part of the design of the LSTM cell state is as a way to mitigate exploding/vanishing gradients. Each cell state $C_t$ relates to the previous $C_{t-1}$ only via the forget gates, so the resulting gradient is just a product of the forget gates, which take values in $(0, 1)$ and will not explode. They can vanish over long time horizons though, so "long term credit assignment" is still difficult in LSTMs.

### 4.6.1 SGD Guidance

In practice, the step size or learning rate $\alpha_t$ is one of the most important tuning parameters of a deep learning system. Too small a step size and convergence can become very slow, and too large and the optimisation can become unstable or even diverge. A good rule of thumb is find the smallest step size such that SGD becomes unstable, then reduce it by a factor of 2-10. Alternatively, find the largest stable step size.

SGD tends to behave as follows: initially, when the parameters are far away from a good mode, it is easy to estimate the gradient from the minibatch and the stochastic gradient tends to point in the same direction as the gradient, so that SGD tends to get to the vicinity of a mode quickly. Once it is close to a minimum the magnitude of the true gradient decreases, the behaviour of the algorithm starts being dominated by the stochasticity, and SGD will "jump around" the local minimum. Small step sizes suppress

the variance of the stochastic gradients, and it is typical to decrease the step size over the course of the optimisation. For example decreasing it by a constant factor, say by half every 10 epochs, or whenever SGD "looks like it got stuck".

Early analysis of SGD used step sizes that decrease slowly, for example such that

$$\sum_{t=1}^{\infty} \alpha_t = \infty \qquad\qquad \sum_{t=1}^{\infty} \alpha_t^2 < \infty$$

to guarantee convergence. For example $\alpha_t = \alpha/t$. However such step size schedules are now not popular anymore. This is for two reasons. Firstly, while such schedules decrease slowly for large $t$, they decrease rapidly initially so learning can be very slow. Secondly, the slow decrease has the effect of simulated annealing with a slowly decreasing temperature, and the converged minimum tends to be "sharp and deep". There is much evidence that such minima tend to generalise poorly. In fact it has been argued that the stochasticity in SGD has a regularising effect. If the step size is large enough, SGD will be oblivious to such bad minima and converge to wider minima that generalise better. Note that the step size schedule in the previous paragraph decreases exponentially quickly and does not satisfy the requirements here. Other recent schedules like cosine and linear decay also do not satisfy the requirements.

### 4.6.2 Momentum, ADAM

While the stochasticity of SGD can have a regularising effect, it is still useful to reduce the stochasticity, so that larger learning rates can be used (with effectively the same variance, the convergence can be faster). A popular approach is using **momentum** to average out the stochastic gradients over previous minibatches:

$$\Delta_t = \beta \Delta_{t-1} + \widehat{\frac{d\mathcal{L}}{d\theta}}(\theta_t)$$
$$\theta_t = \theta_t - \alpha_t \Delta_t$$

The average is using an exponentially decaying weight over past minibatch gradients. Momentum can also help in narrow valleys, as the direction along the valley tends to get small gradients which can be accumulated by the momentum.

Another issue with standard SGD is that it uses a single step size for all parameters. Because of the nonlinearity of the function $f_\theta$, different parameters may have different scales and may need different step sizes. **ADAM** is a popular optimiser that does this:

$$\Delta_t = \beta_1 \Delta_{t-1} + (1 - \beta_1) \widehat{\frac{d\mathcal{L}}{d\theta}}(\theta_t) \qquad\qquad \tilde{\Delta}_t = \frac{\Delta_t}{1 - \beta_1^t}$$

$$V_t = \beta_2 V_{t-1} + (1 - \beta_2) \widehat{\frac{d\mathcal{L}}{d\theta}}(\theta_t)^2 \qquad\qquad \tilde{V}_t = \frac{V_t}{1 - \beta_2^t}$$

$$\theta_t = \theta_t - \frac{\alpha_t}{\sqrt{\tilde{V}_t} + \epsilon} \tilde{\Delta}_t$$

Both $\Delta_0$ and $V_0$ are initialised at 0, and the division by $1 - \beta^t$ ensures that the exponentially decaying weights for averaging sum to 1. $V_t$ is an estimate of the expected squared gradients, which is related to the curvature of the loss landscape.

## 4.7 Initialisation and Parameter Scales

The parameters of a neural network are initialised randomly, typically with zero-mean Gaussian distributions. Random initialisation is important to "break symmetry" so that each unit in a layer will learn differently. Variances of the initialisation are chosen so that at initialisation the network does something reasonable.

Consider again the simple MLP (4.5). Ignoring the nonlinearities, the pre-activation going into a unit $x_h^{(l)}$ has the form:

$$\tilde{x}_h^{(l)} = \sum_{j=1}^{d_{l-1}} W_{hj}^{(l)} x_j^{(l-1)} + b_h$$

It is generally a good idea to control the typical scale of $\tilde{x}_h^{(l)}$ to be $O(1)$. Many non-linearities (e.g. $\mathsf{sigmoid}, \mathsf{tanh}, \mathsf{softplus}$) are designed to have their nonlinear regimes around $[-2, 2]$, so if $\tilde{x}_h^{(l)}$ is much bigger or small than $O(1)$, the nonlinearities will look like linear functions, or worse, constant functions (constant functions will have zero gradient so the unit will not be able to learn due to vanishing gradient problem).

We can control the typical scales of pre-activations layer by layer. Suppose $x_j^{(l-1)} = O(1)$, and suppose they are close to being independent (this is called a mean-field assumption). Then $\mathrm{Var}(\tilde{x}_h^{(l)}) = O(d_{l-1}\sigma_{lw}^2 + \sigma_{lb}^2)$ where $\sigma_{lw}^2$ is the variance of the weights and $\sigma_{lb}^2$ the variance of the biases, for layer $l$. If we set $\sigma_{lw}^2 = O(1/d_{l-1})$, we can get that $\mathrm{Var}(\tilde{x}_h^{(l)}) = O(1)$. Using $\sigma_{lw}^2 = 1/d_{l-1}$ at initialisation is called **Xavier initialisation**.

## 4.8 Regularisation

The standard $L_2$ norm regularisation term in ERM is called **weight decay**, as its effect is to decay each parameter by $(1 - \alpha_t \lambda)$ at each iteration.

While it is frequently used, often other more algorithmic forms of regularisation are used in deep learning. For example, we have seen that the stochasticity of SGD can be thought of as a form of regularisation.

Another popular regularisation via randomisation method is called **dropout**, which works as follows: during training, for each data item, and for each unit in the network, randomly remove it from the network (with probability $p$, typically $p = 0.5$) by multiplying its activation by 0. The reduced networks are then trained using standard SGD. Dropout can be thought of as a cheap and exponentially large ensemble of neural networks, with each network being constructed by randomly dropping out units of a main network.
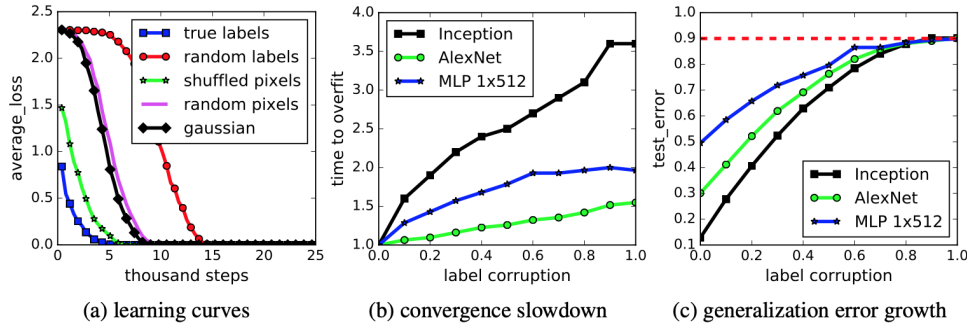
Figure 1: Fitting random labels and random pixels on CIFAR10. (a) shows the training loss of various experiment settings decaying with the training steps. (b) shows the relative convergence time with different label corruption ratio. (c) shows the test error (also the generalization error since training error is 0) under different label corruptions.

Figure 4.1: Figure from [23].

Another popular regularisation technique is **early stopping**. Monitor validation loss during training, and stop when it starts going up. Controlling amount of computation is a generic and actively studied way to control the complexity of models, in both deep learning and convex/kernel methods.

## 4.9 Other Observations

Tips for getting deep learning systems working:

- Start by making sure that there are no optimisation problems, model under-capacity problems or data pre-processing problems. Use a small subset of the data (as in a single minibatch), and make sure the network can overfit to it.

- Start with simple baselines before trying fancier methods.

Andrew Ng's advice from Nuts and Bolts of Applying Deep Learning:

1. Choose a large enough model and training regime so that it can overfit on the whole training set.

2. Use the validation set to tune the model regularisation such that it performs as well as possible. If this is still not good enough (on validation set), collect more data and go to Step 1.

3. If validation result is good enough, check on test data.

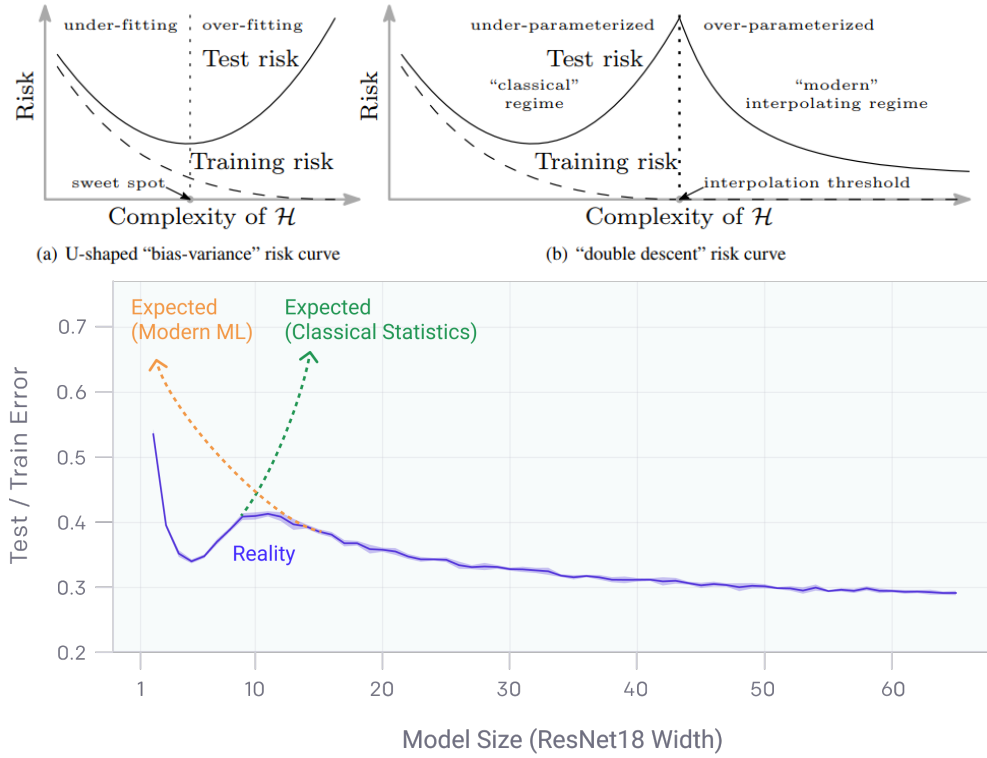4. If test result is good enough, done!

Figure 4.2: Figures from [2] and OpenAI blog on "Deep Double Descent".

A well known classical result is that simple MLPs with a single hidden layer that is wide enough can approximate any smooth function arbitrarily well. So there is common folk wisdom that the hypothesis classes of most modern neural networks are very rich. In fact, while technically the loss landscape is highly complex and multi-modal, recent empirical observations show that modern deep neural networks can always overfit any training data, even random datasets (see Figure 4.1).

Generally larger models never hurt. In fact recent empirical and theoretical results show a curious "double descent" style learning curve for large models [2] (see Figure 4.2), where the best generalisation performance is obtained for overparameterised models. All these results show that classical statistical and learning theory fail to explain the behaviour of modern deep learning models, and these empirical findings have driven much exciting theoretical research in recent years.

Important points about deep learning:

- More data is almost always better than fancier models.

- Learning system has to be scalable, to be able to take advantage of data.

- Learning system has to be flexible and expressive, to be able to learn complex

relationships and representations from data.

- Whole system learning, or end-to-end learning, is a good idea.

- There is no clear distinction between model and computation.

- Everything is a regulariser.

# 5 Latent Variable Models, Variational Inference and the EM algorithm

## 5.1 Latent Variable Models

$K$-means and PCA are two popular unsupervised learning methods. Their goals are for summarising and visualising data. A wider class of methods for unsupervised learning are based on **probabilistic generative models**, where the goal is to model the data generating process using latent variables, with the latent variables representing a summary of the data. For example, in mixture models, the latent variables correspond to assignments of data items to clusters, while in probabilistic PCA and variational autoencoders, the latent variables are lower dimensional representations of data.

The probabilistic approach provides a rich toolbox of techniques, and a coherent principle to reason about and model complex data and to derive new algorithms:

- Conditional distributions over latent variables reflects the model's uncertainties about the summary representations of data.

- The probabilistic method can be extended to Bayesian learning, where model parameters are themselves treated as random, and where uncertainties about model parameters can be expressed as well.

- Probabilistic models allow handling missing data, outliers, heterogenous data, structured data,

- Prior knowledge and inductive biases can be introduced into the model in a principle way.

- Algorithms can be derived from basic principles.

The disadvantages of the probabilistic and Bayesian approach include:

- Exact computation of posterior distributions can be intractable for complex models, and approximations are necessary.

- Automated tools only work well for specialised subclasses of probabilistic models.

### 5.1.1 Clustering and Mixture Modelling

$K$-means and hierarchical clustering are non-probabilistic algorithms – based on intuitive notions of clustering "similar" instances together and "dissimilar" instances apart. Their

goal is not to model the probability of the observed data items. In contrast, a **mixture model** is a probabilistic counterpart which described a full generative process of data, and can produce "soft" assignment of data items to clusters. Mixture models assume that our dataset **x** was created by sampling iid from $K$ distinct populations (called **mixture components**). In other words, data come from a mixture of several sources and the model for the data can be viewed as a convex combination of several distinct probability distributions, often modelled with a given parametric family.

Samples in population $k$ can be modelled using a distribution $F_{\mu_k}$ with density $f(x|\mu_k)$, where $\mu_k$ is the *model parameter* for the $k$-th component. For a concrete example, consider a $p$-dimensional multivariate normal density with unknown mean $\mu_k$ and *known diagonal* covariance $\sigma^2 I$,

$$f(x|\mu_k) = |2\pi\sigma^2|^{-\frac{p}{2}} \exp\left(-\frac{1}{2\sigma^2}\|x - \mu_k\|_2^2\right). \tag{5.1}$$

Such model corresponds to the following generative model, whereby for each data item $i = 1, 2, \ldots, n$, we

(i) first determine the **assignment variable** (independently for each data item $i$):

$$Z_i \overset{iid}{\sim} \text{Mult}(\pi_1, \ldots, \pi_K) \qquad\qquad \text{i.e., } \mathbb{P}(Z_i = k) = \pi_k$$

where for $k = 1, \ldots, K$, $\pi_k \geq 0$, such that $\sum_{k=1}^{K} \pi_k = 1$, are the **mixing proportions**, additional model parameters to be inferred;

(ii) then, given the assignment $Z_i = k$ of the mixture component, $X_i = (X_i^{(1)}, \ldots, X_i^{(p)})^\top$ is sampled (independently) from the corresponding $k$-th component:

$$X_i|(Z_i = k) \sim f(x|\mu_k).$$

We observe $X_i = x_i$ for each $i$ but do not observe its assignment $Z_i$ (**latent variables**), and would like to learn the parameters $\theta = (\mu_1, \ldots, \mu_K, \pi_1, \ldots, \pi_K)$ as well as infer the latent variables.

Note that the **complete log likelihood** of the model

$$\log p(\mathbf{z}, \mathbf{x}|\theta) = \log\left(\prod_{i=1}^{n} \pi_{z_i} f(x_i|\mu_{z_i})\right) = \sum_{i=1}^{n} \left(\log \pi_{z_i} + \log f(x_i|\mu_{z_i})\right) \tag{5.2}$$

is not available as $z_i$ is not observed. We can consider marginalising over the latent variables, giving the **marginal distribution**

$$p(\mathbf{x}|\theta) = \sum_{z_1=1}^{K} \cdots \sum_{z_n=1}^{K} \prod_{i=1}^{n} \pi_{z_i} f(x_i|\mu_{z_i}) = \prod_{i=1}^{n} \left(\sum_{k=1}^{K} \pi_k f(x_i|\mu_k)\right). \tag{5.3}$$

The **marginal log likelihood** of the parameters given the observations,

$$\ell(\theta) = \log p(\mathbf{x}|\theta) = \sum_{i=1}^{n} \log \sum_{k=1}^{K} \pi_k f(x_i|\mu_k).$$

Another important quantity is the **posterior distribution**,

$$p(\mathbf{z}|\mathbf{x}, \theta) = \frac{p(\mathbf{z}, \mathbf{x}|\theta)}{p(\mathbf{x}|\theta)} = \prod_{i=1}^{n} p(z_i|X_i, \theta) \tag{5.4}$$

which factorises into posterior distributions of individual latent variables $z_i$ given the corresponding observations $X_i$ in case of $K$-means, as our data is assumed iid.

In the context of latent variable models, there are two important processes: **inference**, which is the computation or approximation of the posterior distribution over latent variables, and learning, which is the estimation of the model parameters $\theta$. Typically the learning algorithm is done via maximising the marginal log likelihood, and requires inference as a subroutine. However, direct maximisation is often not feasible and the marginal log likelihood will often have many local optima. Fortunately, there is a general approach for both learning and inference based on the formulation of an objective function called the *variational free energy*, which is a lower bound on the marginal log likelihood. From the variational free energy we can derive large classes of algorithms, including the *Expectation Maximisation (EM) algorithm*, which we will describe in Section 5.3, and *amortised variational inference*, which is used for learning in variational autoencoders.

## 5.2 KL Divergence and Gibbs' Inequality

Before we describe the variational formulation, we will review the notion of **Kullback-Leibler (KL) divergence** or **relative entropy** between probability distributions $P$ and $Q$.

**KL divergence.**

- Let $P$ and $Q$ be two absolutely continuous probability distributions on $\mathcal{X} \subseteq \mathbb{R}^d$ with densities $p$ and $q$ respectively. Then the KL divergence *from $Q$ to $P$* is defined as

$$D_{KL}\left(P \parallel Q\right) = \int_{\mathcal{X}} p(x) \log \frac{p(x)}{q(x)} dx. \tag{5.5}$$

- Let $P$ and $Q$ be two discrete probability distributions with probability mass functions $p$ and $q$ respectively. Then the KL divergence *from $Q$ to $P$* is defined as

$$D_{KL}\left(P \parallel Q\right) = \sum_{i} p(x_i) \log \frac{p(x_i)}{q(x_i)}. \tag{5.6}$$

In both cases, we can write

$$D_{KL}\left(P \parallel Q\right) = \mathbb{E}_p \left[ \log \frac{p(X)}{q(X)} \right], \tag{5.7}$$

where $\mathbb{E}_p$ denotes that expectation is taken over $p$. By convexity of $f(x) = -\log(x)$ and Jensen's inequality (5.9), we have that

$$D_{KL}\left(P \parallel Q\right) = \mathbb{E}_p\left[-\log\frac{q(X)}{p(X)}\right] \geq -\log\mathbb{E}_p\frac{q(X)}{p(X)} = 0, \tag{5.8}$$

where in the last step we used that $\int_{\mathcal{X}} q(x)dx = 1$ in continuous case and $\sum_i q(x_i) = 1$ in discrete case.

**Jensen's inequality.** Let $f$ be a convex function and $X$ be a random variable. Then

$$\mathbb{E}\left[f(X)\right] \geq f\left(\mathbb{E}X\right). \tag{5.9}$$

If $f$ is strictly convex, then equality holds if and only if $X$ is almost surely a constant.

Thus, we conclude that KL-divergence is always non-negative. This consequence of Jensen's inequality is called **Gibbs' inequality**. Moreover, since $f(x) = -\log(x)$ is strictly convex on $x > 0$, the equality holds if and only if $p(x) = q(x)$ almost everywhere, i.e. $P = Q$. Note that in general KL-divergence is *not symmetric*: $D_{KL}\left(P \parallel Q\right) \neq D_{KL}\left(Q \parallel P\right)$.

## 5.3 Variational Free Energy and EM Algorithm

The **EM algorithm** is a general purpose iterative strategy for local maximisation of the likelihood under missing data/hidden variables. The method has been proposed many times for specific models– it was given its name and studied as a general framework by [8].

Let $(\mathbf{x}, \mathbf{z})$ be a pair of observed variables $\mathbf{x}$, and latent variables $\mathbf{z}$. Our probabilistic model is given by $p(\mathbf{x}, \mathbf{z}|\theta)$, but we have no access to $\mathbf{z}$. Therefore, we would like to maximise the observed data log-likelihood (marginal log-likelihood) $\ell(\theta) = \log p(\mathbf{x}|\theta) = \log \int p(\mathbf{x}, \mathbf{z}|\theta)d\mathbf{z}$ over $\theta$. However, marginalisation of latent variables typically results in an intractable optimization problem and we need to resort to approximations.

Now, assume for a moment that we have access to another objective function $\mathcal{F}(\theta, q)$, where $q(\mathbf{z})$ is a certain distribution on latent variables $\mathbf{z}$, which we are free to choose and will call **variational distribution**. Moreover, assume that $\mathcal{F}$ satisfies

$$\mathcal{F}(\theta, q) \leq \ell(\theta) \text{ for all } \theta, q, \tag{5.10}$$

$$\max_q \mathcal{F}(\theta, q) = \ell(\theta), \tag{5.11}$$

i.e. $\mathcal{F}(\theta, q)$ is a *lower bound on the log-likelihood* for any variational distribution $q$ (5.10), which also *matches the log-likelihood* at a particular choice of $q$ (5.11).

Given these two properities, we can construct an alternating maximisation: *coordinate ascent* algorithm as follows:

**Coordinate ascent on the lower bound.** For $t = 1, 2 \ldots$ until convergence:

$$q^{(t)} := \operatorname{argmax}_q \mathcal{F}(\theta^{(t-1)}, q)$$

$$\theta^{(t)} := \operatorname{argmax}_\theta \mathcal{F}(\theta, q^{(t)}).$$

**Theorem 15.** *Assuming* (5.10) *and* (5.11), *coordinate ascent on the lower bound* $\mathcal{F}(\theta, q)$ *does not decrease the log likelihood* $\ell(\theta)$.

*Proof.* $\ell(\theta^{(t-1)}) = \mathcal{F}(\theta^{(t-1)}, q^{(t)}) \leq \mathcal{F}(\theta^{(t)}, q^{(t)}) \leq \mathcal{F}(\theta^{(t)}, q^{(t+1)}) = \ell(\theta^{(t)})$. Additional assumption, that $\nabla^2_\theta \mathcal{F}(\theta^{(t)}, q^{(t)})$ are negative definite with eigenvalues $< -\epsilon < 0$, implies that $\theta^{(t)} \to \theta^*$ where $\theta^*$ is a local MLE. $\qquad\square$

But how to find such lower bound $\mathcal{F}$? It is given by the so called *variational free energy*, which we define next.

**Definition 16. Variational free energy** in a latent variable model $p(\mathbf{x}, \mathbf{z}|\theta)$ is defined as

$$\mathcal{F}(\theta, q) = \mathbb{E}_q[\log p(\mathbf{x}, \mathbf{z}|\theta) - \log q(\mathbf{z})], \tag{5.12}$$

where $q$ is any probability density/mass function over the latent variables $\mathbf{z}$.

The variational free energy is also referred to as **ELBO**, which stands for evidence lower bound, where **evidence** is an alternative term for the marginal log likelihood.

Consider the KL divergence between $q(\mathbf{z})$ and the true conditional based on our model $p(\mathbf{z}|\mathbf{x}, \theta) = p(\mathbf{x}, \mathbf{z}|\theta)/p(\mathbf{x}|\theta)$ for the observations $\mathbf{x}$ and a fixed parameter vector $\theta$. Since KL is non-negative,

$$
\begin{aligned}
0 \leq D_{KL}\left[q(\mathbf{z}) \parallel p(\mathbf{z}|\mathbf{x}, \theta)\right] &= \mathbb{E}_{\mathbf{z} \sim q} \log \frac{q(\mathbf{z})}{p(\mathbf{z}|\mathbf{x}, \theta)} \\
&= \log p(\mathbf{x}|\theta) + \mathbb{E}_{\mathbf{z} \sim q} \log \frac{q(\mathbf{z})}{p(\mathbf{x}, \mathbf{z}|\theta)}.
\end{aligned}
$$

Thus, we have obtained a lower bound on the marginal log-likelihood which holds true for *any parameter value* $\theta$ and *any choice of the variational distribution* $q$:

$$\ell(\theta) = \log p(\mathbf{x}|\theta) \geq \mathbb{E}_{\mathbf{z} \sim q} \log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z})} = \underbrace{\mathbb{E}_{\mathbf{z} \sim q} \log p(\mathbf{x}, \mathbf{z}|\theta)}_{\text{energy}} \overbrace{- \mathbb{E}_{\mathbf{z} \sim q} \log q(\mathbf{z})}^{\text{entropy}}. \tag{5.13}$$

The right hand side in (5.13) is precisely the variational free energy - we see it decomposes in two terms. The first term is usually referred to as *energy*, more precisely it is

the *expected complete data log-likelihood* (if we observed $\mathbf{z}$, we would just maximise the complete data log-likelihood $\log p(\mathbf{x}, \mathbf{z}|\theta)$, but since $\mathbf{z}$ is not observed we need to integrate it out - but recall that $q$ here is *any* distribution over latent variables). The second term is the Shannon entropy $H(q) = -\mathbb{E}_q \log q(\mathbf{z})$ of the variational distribution $q(\mathbf{z})$, and does not depend on $\theta$ (it can be thought of as the complexity penalty on $q$).

The terms "energy" and "entropy" comes from statistical physics. The energy term here should in fact be the *negative* of the average energy of a system in thermal equilibrium, and "free energy" means energy that is accessible or useful (free as opposed to being tied up in entropy so is not accessible). Again the free energy term here should in fact be the *negative* free energy. Physical systems prefer low energy states, while here we are trying to maximise the variational free energy. We keep this inverted usage as it is common in machine learning.

The inequality becomes an equality when KL divergence is zero, i.e. when $q(\mathbf{z}) = p(\mathbf{z}|\mathbf{x}, \theta)$ which means that the optimal choice of variational distribution $q$ for fixed parameter value $\theta$ is *the true conditional of the latent variables given the observations and that $\theta$.*

Thus, we have proved the following lemma:

**Lemma 17.** *Let $\mathcal{F}$ be the variational free energy in a latent variable model $p(\mathbf{x}, \mathbf{z}|\theta)$. Then (a) $\mathcal{F}(\theta, q) \leq \ell(\theta)$ for all $q$ and for all $\theta$, and (b) for any $\theta$, $\mathcal{F}(\theta, q) = \ell(\theta)$ iff $q(\mathbf{z}) = p(\mathbf{z}|\mathbf{x}, \theta)$.*

Thus, properties (5.10) and (5.11) are satisfied and we can recast the alternating maximisation of the variational free energy into iterative updates of $q$ (E-step, via the plug-in full conditional of $\mathbf{z}$ using the current estimate of $\theta$) and the updates of $\theta$ (M-step, by maximising the 'energy' for the current estimate of $q$). Provided that both E-step and M-step can be solved exactly, the EM algorithm converges to a local maximum likelihood solution.

**EM algorithm.** Initialize $\theta^{(0)}$. At time $t \geq 1$:

- E-step: Set $q^{(t)}(\mathbf{z}) = p(\mathbf{z}|\mathbf{x}, \theta^{(t-1)})$

- M-step: Set $\theta^{(t)} = \arg\max_\theta \mathbb{E}_{\mathbf{z} \sim q^{(t)}} \log p(\mathbf{x}, \mathbf{z}|\theta)$.

## 5.4 EM Algorithm for Mixtures

Consider again our mixture model from Section 5.1.1 with

$$p(\mathbf{z}, \mathbf{x}|\theta) = \prod_{i=1}^{n} \pi_{z_i} f(x_i|\mu_{z_i}).$$

Recall that our latent variables $\mathbf{z}$ are discrete (they correspond to cluster assignments) so $q$ is a probability mass function over $\mathbf{z} := (z_i)_{i=1}^n$. Using the expression (5.2), we can write the variational free energy as

$$
\begin{aligned}
\mathcal{F}(\theta, q) =& \mathbb{E}_q[\log p(\mathbf{x}, \mathbf{z}|\theta) - \log q(\mathbf{z})] \\
=& \mathbb{E}_q\left[\left(\sum_{i=1}^n \sum_{k=1}^K \mathbf{1}(z_i = k)\left(\log \pi_k + \log f(x_i|\mu_k)\right)\right) - \log q(\mathbf{z})\right] \\
=& \sum_{\mathbf{z}} q(\mathbf{z})\left[\left(\sum_{i=1}^n \sum_{k=1}^K \mathbf{1}(z_i = k)\left(\log \pi_k + \log f(x_i|\mu_k)\right)\right) - \log q(\mathbf{z})\right] \\
=& \sum_{i=1}^n \sum_{k=1}^K q(z_i = k)\left(\log \pi_k + \log f(x_i|\mu_k)\right) + H(q).
\end{aligned}
$$

We will denote $Q_{ik} = q(z_i = k)$, which is called the **responsibility** of cluster $k$ for data item $i$.

Now, the E-step simplifies because

$$
\begin{aligned}
p(\mathbf{z}|\mathbf{x}, \theta) = \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{p(\mathbf{x}|\theta)} = \frac{\prod_{i=1}^n \pi_{z_i} f(x_i|\mu_{z_i})}{\sum_{\mathbf{z}'} \prod_{i=1}^n \pi_{Z_i'} f(x_i|\mu_{Z_i'})} &= \prod_{i=1}^n \frac{\pi_{z_i} f(x_i|\mu_{z_i})}{\sum_k \pi_k f(x_i|\mu_k)} \\
&= \prod_{i=1}^n p(z_i|x_i, \theta).
\end{aligned}
$$

Thus, for a fixed $\theta^{(t-1)} = (\mu_1^{(t-1)}, \ldots, \mu_K^{(t-1)}, \pi_1^{(t-1)}, \ldots, \pi_K^{(t-1)})$ we can set

$$
Q_{ik}^{(t)} = p(z_i = k|x_i, \theta^{(t-1)}) = \frac{\pi_k^{(t-1)} f(x_i|\mu_k^{(t-1)})}{\sum_{j=1}^K \pi_j^{(t-1)} f(x_i|\mu_j^{(t-1)})}. \tag{5.14}
$$

Now, consider the M-step. For mixing proportions we have a constraint that $\sum_{j=1}^K \pi_j = 1$, so we introduce the Lagrange multiplier and obtain

$$
\begin{aligned}
\nabla_{\pi_k}&\left(\mathcal{F}(\theta, q) - \lambda(\textstyle\sum_{j=1}^K \pi_j - 1)\right) \\
&= \sum_{i=1}^n \frac{Q_{ik}}{\pi_k} - \lambda = 0 \quad \Rightarrow \quad \pi_k \propto \sum_{i=1}^n Q_{ik}.
\end{aligned}
$$

Since

$$
\sum_{k=1}^K \sum_{i=1}^n Q_{ik} = \sum_{i=1}^n \underbrace{\sum_{k=1}^K Q_{ik}}_{=1} = n,
$$

the M-step update for mixing proportions is

$$
\pi_k^{(t)} = \frac{\sum_{i=1}^n Q_{ik}^{(t)}}{n}, \tag{5.15}
$$

i.e., they are simply given by the total responsibility of each cluster. Note that this update holds regardless of the form of the parametric family $f(\cdot|\mu_k)$ used for mixture components.

Setting derivative with respect to $\mu_k$ to 0, we obtain

$$\nabla_{\mu_k}\mathcal{F}(\theta, q) = \sum_{i=1}^{n} Q_{ik}\nabla_{\mu_k}\log f(x_i|\mu_k) = 0. \tag{5.16}$$

This equation can be solved quite easily for mixture of normals in (5.1), giving the M-step update

$$\mu_k^{(t)} = \frac{\sum_{i=1}^{n} Q_{ik}^{(t)} x_i}{\sum_{i=1}^{n} Q_{ik}^{(t)}}, \tag{5.17}$$

which implies that the $k$-th cluster mean estimate is simply a weighted average of all the data items, where the weights correspond to the responsibilities of cluster $k$ for these points.

Put together, the EM for normal mixture model with known (fixed) covariance is very similar to K-means algorithm where cluster assignments are soft, i.e. rather than assigning each data item $x_i$ to a single cluster at each iteration, we carry forward a responsibility vector $(Q_{i1}, \ldots, Q_{iK})$ giving probabilities of $x_i$ belonging to each cluster. Indeed, K-means algorithm can be undestood as EM where $\sigma^2 \to 0$, such that E-step will assign exactly one entry in $(Q_{i1}, \ldots, Q_{iK})$ to one (corresponding to the nearest mean vector) and the rest to zero.

**EM for Normal Mixtures (known covariance) – "Soft K-means"**

1. Initialize $K$ cluster means $\mu_1, \ldots, \mu_K$ and mixing proportions $\pi_1, \ldots, \pi_K$.

2. *Update responsibilites (E-step):* For each $i = 1, \ldots, n$, $k = 1, \ldots, K$:

$$Q_{ik} = \frac{\pi_k \exp\left(-\frac{1}{2\sigma^2}\|x_i - \mu_k\|_2^2\right)}{\sum_{j=1}^{K} \pi_j \exp\left(-\frac{1}{2\sigma^2}\|x_i - \mu_j\|_2^2\right)} \tag{5.18}$$

3. *Update parameters (M-step):* Set $\mu_1, \ldots, \mu_K$ and $\pi_1, \ldots, \pi_K$ and based on the new cluster responsibilities:

$$\pi_k = \frac{\sum_{i=1}^{n} Q_{ik}}{n}, \qquad \mu_k = \frac{\sum_{i=1}^{n} Q_{ik}x_i}{\sum_{i=1}^{n} Q_{ik}}. \tag{5.19}$$

4. Repeat steps 2-3 until convergence.

5. Return the responsibilites $\{Q_{ik}\}$ and parameters $\mu_1, \ldots, \mu_K, \pi_1, \ldots, \pi_K$.

In some cases, depending on the form of the parametric family $f(\cdot|\mu_k)$ the M-step update for mixtures cannot be solved exactly. In these cases, we can use *gradient ascent* algorithm *inside the M-step*:

$$\mu_k^{(r+1)} = \mu_k^{(r)} + \alpha \sum_{i=1}^{n} Q_{ik} \nabla_{\mu_k} \log f(x_i|\mu_k^{(r)}).$$

This leads to a **generalized EM algorithm**.

Mixture models can be applied to supervised learning settings as well. In a **mixture of experts**, we assume that a dataset of input/output pairs $\{(x_i, y_i)\}_{i=1}^{n}$ can be modelled using $K$ experts. Expert $k$ predicts $y$ given $x$ using a neural network which outputs a distribution $f_{\theta_k}(y|x)$ with parameters $\theta_k$. Given $x$, an expert is stochastically chosen using a gating network $g_\phi(x)$ which produces a vector of probabilities over experts. The marginal probability of $y$ given $x$ is then:

$$p(y|x) = \sum_{k=1}^{K} g_\phi(x)_k f_{\theta_k}(y|x).$$

The parameters of the mixture of experts can be learnt to maximise the log marginal likelihood given the dataset using the generalised EM algorithm.

## 5.5 Probabilistic PCA

We describe **probabilistic PCA** [22], a latent variable model for probabilistic dimensionality reduction. Just like in PCA, we try to model a collection of $n$ $p$-dimensional vectors using a $k$-dimensional representation with $k < p$. Probabilistic PCA corresponds to the following generative model.

For each data item $i = 1, 2, \ldots, n$:

- Let $Z_i$ be a (latent) $k$-dimensional normally distributed random vector with mean 0 and identity covariance:
$$Z_i \sim \mathcal{N}(0, I_k),$$

- Given $Z_i$, the distribution of the $i$-th data item is a $p$-dimensional normal:
$$X_i|Z_i \sim \mathcal{N}(\mu + LZ_i, \sigma^2 I)$$
where the parameters $\theta = (\mu, L, \sigma^2)$ correspond to a vector $\mu \in \mathbb{R}^p$, a matrix $L \in \mathbb{R}^{p \times k}$ and $\sigma^2 > 0$.

Note that unlike in clustering, the latent variables $Z_1, \ldots, Z_n$ are now continuous.

From an equivalent representation $X_i = \mu + LZ_i + \epsilon$, where $\epsilon \sim \mathcal{N}(0, \sigma^2 I_p)$ and is independent of $Z_i$, we see that the marginal model on $X_i$'s is

$$f(x|\theta) \;=\; \mathcal{N}\left(x; \mu, LL^\top + \sigma^2 I\right),$$

where parameters are denoted $\theta = \left(\mu, L, \sigma^2\right)$. From here it is clear that the maximum marginal likelihood estimator of $\mu$ is available directly as $\hat{\mu} = \frac{1}{n}\sum_{i=1}^{n} X_i$ and thus, we do not require EM to estimate $\mu$. We will henceforth assume that the data is centred, to simplify notation and remove $\mu$ from the parameters.

On the other hand, maximum marginal likelihood solution for $L$ is unique only up to orthonormal transformations, which is why a certain form of $L$ is usually enforced (e.g. lower-triangular, orthogonal columns). [22] shows that the MLE for PPCA has the following form. Let $\lambda_1 \geq \cdots \geq \lambda_p$ be the eigenvalues of the sample covariance and $V_{1:k} \in \mathbb{R}^{p \times k}$ the top $k$ eigenvectors as before. Let $Q \in \mathbb{R}^{k \times k}$ be any orthogonal matrix. Then we have:

$$\mu^{\text{MLE}} = \bar{x} \qquad (\sigma^2)^{\text{MLE}} = \frac{1}{p-k}\sum_{j=k+1}^{p} \lambda_j$$
$$L^{\text{MLE}} = V_{1:k}\text{diag}((\lambda_1 - (\sigma^2)^{\text{MLE}})^{\frac{1}{2}}, \ldots, (\lambda_k - (\sigma^2)^{\text{MLE}})^{\frac{1}{2}})Q.$$

We note that the standard PCA is recovered when $\sigma^2 \to 0$. While the EM algorithm is not necessarily preferable to SVD for standard PCA, we will now proceed by deriving the EM algorithm.

**E-step.** By Gaussian conditioning (*exercise*),

$$q(z_i) = p(z_i|x_i, \theta) = \mathcal{N}\left(z_i|b_i, R\right),$$

where

$$b_i = \left(L^\top L + \sigma^2 I\right)^{-1} L^\top x_i, \tag{5.20}$$

$$R = \sigma^2 \left(L^\top L + \sigma^2 I\right)^{-1}. \tag{5.21}$$

**M-step. [not examinable]** Recall that the parameters of interest are $\theta = \left(L, \sigma^2\right)$ (since the marginal maximum likelihood estimate of the mean parameter $\mu$ is directly available). We would like to maximise the variational free energy given by:

$$\mathcal{F}(\theta, q) = \mathbb{E}_{\mathbf{Z} \sim q}\left[\sum_{i=1}^{n} \log p(x_i, z_i|\theta)\right] + \text{const.}$$

By ignoring terms that do not depend on $\theta$ and denoting $=_c$ to mean "equal up to a constant independent on $\theta$"

$$\begin{aligned}
\log p(x_i, z_i|\theta) &=_c -\frac{p}{2}\log \sigma^2 - \frac{1}{2\sigma^2}(x_i - Lz_i)^\top (x_i - Lz_i) \\
&=_c -\frac{p}{2}\log \sigma^2 - \frac{1}{2\sigma^2}\left\{x_i^\top x_i - 2x_i^\top Lz_i + z_i^\top L^\top Lz_i\right\} \\
&=_c -\frac{p}{2}\log \sigma^2 - \frac{1}{2\sigma^2}\left\{x_i^\top x_i - 2x_i^\top Lz_i + \text{Tr}\left[L^\top Lz_iz_i^\top\right]\right\}.
\end{aligned}$$

Taking expectation over $q(z_i) = \mathcal{N}(z_i|b_i, R)$ gives

$$\mathbb{E}_{z_i \sim q}(\log p(x_i, z_i|\theta)) =_c -\frac{p}{2}\log \sigma^2 - \frac{1}{2\sigma^2}\left\{x_i^\top x_i - 2x_i^\top L b_i + \mathrm{Tr}\left[L^\top L\left(b_i b_i^\top + R\right)\right]\right\}.$$

It remains to sum over all observations to get

$$\begin{aligned}\mathcal{F}(\theta, q) \quad =_c \quad & -\frac{np}{2}\log \sigma^2 \\ & -\frac{1}{2\sigma^2}\left\{\sum_{i=1}^n x_i^\top x_i - 2\sum_{i=1}^n x_i^\top L b_i + \mathrm{Tr}\left[L^\top L\left(\sum_{i=1}^n b_i b_i^\top + nR\right)\right]\right\}.\end{aligned}$$

Now, we have

$$\frac{\partial \mathcal{F}}{\partial L} \quad = \quad \frac{1}{\sigma^2}\left\{\sum_{i=1}^n x_i b_i^\top - L\left(\sum_{i=1}^n b_i b_i^\top + nR\right)\right\},$$

which by setting to 0 gives the update rule

$$L^{(\mathrm{new})} = \left(\sum_{i=1}^n x_i b_i^\top\right)\left(\sum_{i=1}^n b_i b_i^\top + nR\right)^{-1}. \tag{5.22}$$

Letting $\tau = \sigma^{-2}$, we have:

$$\frac{\partial \mathcal{F}}{\partial \tau} \quad = \quad \frac{np}{2}\frac{1}{\tau} - \frac{1}{2}\left\{\sum_{i=1}^n x_i^\top x_i - 2\sum_{i=1}^n x_i^\top L b_i + \mathrm{Tr}\left[L^\top L\left(\sum_{i=1}^n b_i b_i^\top + nR\right)\right]\right\},$$

and thus

$$\left(\sigma^2\right)^{(\mathrm{new})} = \frac{1}{np}\left\{\sum_{i=1}^n x_i^\top x_i - 2\sum_{i=1}^n x_i^\top L^{(\mathrm{new})} b_i + \mathrm{Tr}\left[L^{(\mathrm{new})\top}L^{(\mathrm{new})}\left(\sum_{i=1}^n b_i b_i^\top + nR\right)\right]\right\}. \tag{5.23}$$

Both Probabilistic PCA and normal mixtures are examples of linear Gaussian models, all of which have the corresponding learning algorithms based on EM. For a unifying review of these and a number of other models from the same family, including factor analysis and hidden Markov models, cf. [19].

## 5.6 Variational Autoencoders

Probabilistic PCA and PCA are linear dimensionality reduction techniques. We can use neural networks to define a probabilistic model for nonlinear dimensionality reduction. Such models are called **deep generative models**, and the **variational autoencoder** (VAE) [12, 18] is the most popular example of such models.

The basic idea is to keep the assumption in probabilistic PCA of a standard Gaussian prior over a latent variable

$$Z_i \sim N(0, I_k)$$

with $k < p$, but to make the **generative model**, or **decoder**, nonlinear and non-Gaussian

$$X_i | Z_i \sim p_\theta(\cdot | Z_i)$$

For example, $p_\theta(\cdot | Z_i)$ can be a Gaussian with mean $\mu_\theta(Z_i)$ and diagonal variance $\sigma_\theta^2(Z_i)$, where both $\mu_\theta$ and $\sigma_\theta^2$ are outputs of a neural network with parameters $\theta$. For another example, $p_\theta(X_i | Z_i)$ can be an independent Bernoulli for each element of $X_i$, with mean $\mu_\theta(Z_i) \in (0,1)$, with the range restricted to $(0,1)$ using a sigmoid output nonlinearity.

To learn a VAE, we again consider optimising the variational free energy/evidence lower bound (ELBO). Given the neural network in the generative model, the exact posterior distribution will typically not be tractable anymore. Instead we consider a variational distribution $q_\phi(z|x)$, called a **recognition model** or **encoder**, which is itself given by a neural network with parameters $\phi$. The encoder is typically a Gaussian with a diagonal covariance matrix,

$$q_\phi(z|x) = \mathcal{N}(z; \mu_\phi(x), \Sigma_\phi(x)) \tag{5.24}$$

We will also use the notation $\sigma_\phi^2(x)$ for the vector of diagonal variances in $\Sigma_\phi(x)$. Other options for encoders are possible, including augmenting VAE posterior approximations by transforming drawn samples through mappings (so called *flows*) with additional trainable parameters to achieve richer variational families.

The approach of learning a function approximator to approximate the posterior distribution is called **amortised inference**. The word amortised comes, via the probabilistic programming literature, from the compilers literature, where amortisation refers to pre-computing functions at compile time so that function evaluations can be performed efficiently at run time. Amortised inference uses a function approximator which is learnt at training time, so that inference for new data items at test time can be performed efficiently using the function approximator. This is as opposed to iterative inference procedures, e.g. Markov chain Monte Carlo or coordinate ascent variational inference, which can be much more expensive at test time.

**VAE ELBO.** Consider different ways in which we can write the ELBO for a single observation $x$:

$$
\begin{aligned}
\mathcal{L}(x, \theta, \phi) &= \mathbb{E}_{q_\phi(z|x)} \left[ \log p_\theta(x, z) \right] + H\left( q_\phi(\cdot | x) \right) \\
&= \mathbb{E}_{q_\phi(z|x)} \left[ \log \frac{p_\theta(x, z)}{q_\phi(z|x)} \right] \\
&= \mathbb{E}_{q_\phi(z|x)} \left[ \log p_\theta(x|z) \right] + \mathbb{E}_{q_\phi(z|x)} \left[ \log \frac{p(z)}{q_\phi(z|x)} \right] \\
&= \mathbb{E}_{q_\phi(z|x)} \left[ \log p_\theta(x|z) \right] - KL\left( q_\phi(z|x) \| p(z) \right).
\end{aligned}
\tag{5.25}
$$

The final expression is particularly convenient as for a normal variational posterior in (5.24) and a standard normal prior $p(z)$, the KL divergence is available in closed form:

$$KL\left(q_\phi(z|x)||p(z)\right) = \frac{1}{2}\left[\mu_\phi(x)^\top \mu_\phi(x) + \text{tr}\left(\Sigma_\phi(x)\right) - \log\det\left(\Sigma_\phi(x)\right) - k\right]$$

$$= \frac{1}{2}\sum_{j=1}^{k}\left[\mu_{\phi,j}^2(x) + \sigma_{\phi,j}^2(x) - \log\left(\sigma_{\phi,j}^2(x)\right) - 1\right]. \quad (5.26)$$

To obtain the ELBO objective on the whole set of observations $\{x_i\}_{i=1}^n$ ,we average the individual terms in (5.25), i.e.

$$\mathcal{L}(\theta,\phi) = \frac{1}{n}\sum_{i=1}^{n}\left\{\mathbb{E}_{q_\phi(z|x_i)}\left[\log p_\theta\left(x_i|z\right)\right] - KL\left(q_\phi\left(z|x_i\right)||p(z)\right)\right\}. \quad (5.27)$$

Assuming that $\{(x_i, z_i)\}_{i=1}^n$ are drawn i.i.d. from $p_\theta(x, z)$, $\mathcal{L}(\theta, \phi)$ is a lower bound on the (scaled) model evidence $\frac{1}{n}\log p_\theta\left(\{x_i\}_{i=1}^n\right) = \frac{1}{n}\sum_{i=1}^n\log p_\theta\left(x_i\right)$, since $\mathcal{L}(x_i, \theta, \phi) \leq \log p_\theta\left(x_i\right)$, for each $i$.

As opposed to the EM algorithm, which has two separate steps for inference and learning, we proceed with stochastic gradient descent step which jointly optimise (5.27) with respect to $\theta$ and $\phi$, using minibatches of observations $x_i$ in order to efficiently compute unbiased estimators of the gradients of ELBO. We note, however, that the terms $\mathbb{E}_{q_\phi(z|x_i)}\left[\log p_\theta\left(x_i|z\right)\right]$ are generally not tractable, which we address next.

**Reparametrization trick.** In order to optimize the ELBO objective over $\theta$ and $\phi$, it remains to estimate the terms $\mathbb{E}_{q_\phi(z|x_i)}\left[\log p_\theta\left(x_i|z\right)\right]$. Getting a good Monte Carlo estimator would require drawing many samples of codes for each obervation $x_i$, which is prohibitive. Thus, one is tempted to proceed with drawing a single $z_i \sim q_\phi\left(z|x_i\right)$ and estimating

$$\hat{\mathbb{E}}_{q_\phi(z|x_i)}\left[\log p_\theta\left(x_i|z\right)\right] = \log p_\theta\left(x_i|z_i\right).$$

However, this is clearly problematic, as explicit dependence of this estimator on the variational parameters $\phi$ has been lost and we cannot compute its gradients with respect to $\phi$. There is a simple solution, however, known as the **reparametrization trick**.

For example, in the case of a normal variational posterior, a draw $z_i \sim \mathcal{N}\left(z|\mu_\phi(x), \sigma_\phi^2(x)\right)$ can be written as $z_i = \mu_\phi(x) + \sigma_\phi(x)\epsilon_i$, with $\epsilon_i \sim \mathcal{N}(0, I)$, and the multiplication with $\sigma_\phi(x)$ being elementwise. We can rewrite

$$\mathbb{E}_{q_\phi(z|x_i)}\left[\log p_\theta\left(x_i|z\right)\right] = \mathbb{E}_\epsilon\left[\log p_\theta\left(x_i|\mu_\phi(x) + \sigma_\phi(x)\epsilon\right)\right],$$

and use estimator of the form $\log p_\theta(x_i|\mu_\phi(x) + \sigma_\phi(x)\epsilon_i)$, based on a single draw $\epsilon_i \sim \mathcal{N}(0, I)$. Now, it is possible to compute gradients (using automatic differentiation) $\nabla_\theta\log p_\theta(x_i|\mu_\phi(x) + \Sigma_\phi^{1/2}(x)\epsilon_i)$ and $\nabla_\phi\log p_\theta(x_i|\mu_\phi(x) + \sigma_\phi(x)\epsilon_i)$ with respect to both $\theta$ and $\phi$, and, importantly, they are unbiased estimators of $\nabla_\theta\mathbb{E}_{q_\phi(z|x_i)}\left[\log p_\theta(x_i|z)\right]$ and $\nabla_\phi\mathbb{E}_{q_\phi(z|x_i)}\left[\log p_\theta(x_i|z)\right]$, respectively.

In general, the reparametrization trick is possible whenever we can sample from $q_\phi(z|x_i)$ by computing some differentiable function $h(\epsilon, \phi, x_i)$, where $\epsilon$ is a draw from a standard distribution with no parameters.

In the general case, for example if either the prior or the variational distribution are not Gaussians, the KL term need not be tractable. However we can still form an unbiased estimator for it using $\log q_\phi(z_i|x_i) - \log p(z_i)$ using a single sample $z_i \sim q_\phi(\cdot|x_i)$, and again computed an unbiased estimator of the gradient with respect to both $\phi$ and $\theta$. Note that this works if the estimator is differentiable with respect to $z_i$. For example, it won't work if $z_i$ is discrete.

**VAEs and unbiased estimators of $p_\theta(x)$.** Assume that we have access to some strictly positive unbiased estimator $\hat{p}_\theta(x)$ of $p_\theta(x)$, with

$$\int \hat{p}_\theta(x)q_{\theta,\phi}(u|x)\, du = p_\theta(x),$$

where $u \sim q_{\theta,\phi}(\cdot|x)$ denotes all random variables used to compute the estimator $\hat{p}_\theta(x) = \hat{p}_\theta(x; u, \phi)$, with $\phi$ denoting any additional parameters of the sampling distribution. Jensen's inequality then tells us that

$$\int q_{\theta,\phi}(u|x) \log \hat{p}_\theta(x)du \quad \leq \quad \log \int \hat{p}_\theta(x)q_{\theta,\phi}(u|x)\, du = \log p_\theta(x).$$

Thus, the term

$$\mathcal{L}(x, \theta, \phi) := \int q_{\theta,\phi}(u|x) \log \hat{p}_\theta(x; u, \phi)\, du \tag{5.28}$$

is a lower bound on evidence for any choice of an unbiased estimator $\hat{p}_\theta(x; u, \phi)$. VAE framework we described above corresponds to simply setting $u = z$ and the estimator is of the form $\hat{p}_\theta(x) = p_\theta(x, z)/q_\phi(z|x)$, which is the importance weight if $q_\phi(z|x)$ is the proposal distribution and the target is the posterior $p_\theta(z|x)$.

One can consider other types of estimators. For example, **importance weighted autoencoder** (IWAE) [7] uses an estimator based on $s$ importance samples $u = \{z_j\}_{j=1}^s$, with $z_j \overset{iid}{\sim} q_\phi(\cdot|x)$

$$\hat{p}_\theta(x) = \frac{1}{s} \sum_{j=1}^s \frac{p_\theta(x, z_j)}{q_\phi(z_j|x)}.$$

While the lower bound in (5.28) may not be tractable in this case, it suffices to only compute unbiased estimators of its gradients with respect to parameters $\theta$ and $\phi$. Similarly to before, we can still use the reparameterisation trick. The IWAE ELBO is better than the VAE ELBO, since the corresponding unbiased estimator of $p_\theta(x)$ has lower variance, and it monotonically improved and approaches the exact log marginal likelihood $\log p_\theta(x)$ as $s \to \infty$.

## 5.7 Score gradient identity [not examinable]

There is an alternative approach to obtain unbiased estimates of gradients with respect to variational parameters, termed *score gradient* or *REINFORCE gradient*. First recall the score identity, i.e. that the expected value of the score is zero. Given a model $q_\phi(z)$, with parameters $\phi$, the gradient of the log-likelihood with respect to parameters, $\frac{\partial}{\partial \phi}[\log q_\phi(z)]$ is termed the score. Under mild regularity condition, the expectation of the score is zero:

$$
\begin{aligned}
\mathbb{E}_{q_\phi(z)} \frac{\partial}{\partial \phi}[\log q_\phi(z)] &= \int q_\phi(z) \frac{\partial}{\partial \phi}[\log q_\phi(z)] \, dz \\
&= \int \frac{\partial}{\partial \phi} q_\phi(z) \, dz \\
&= \frac{\partial}{\partial \phi} \int q_\phi(z) \, dz \\
&= \frac{\partial}{\partial \phi} 1 = 0.
\end{aligned}
$$

Score gradient is based on the following observation

$$
\begin{aligned}
\frac{\partial}{\partial \phi} \mathcal{L} &= \frac{\partial}{\partial \phi} \int q_\phi(z|x) \log \frac{p(x,z)}{q_\phi(z|x)} \, dz \\
&= \int \frac{\partial}{\partial \phi}[q_\phi(z|x) \log p(x,z) - q_\phi(z|x) \log q_\phi(z|x)] \, dz \\
&= \int \left\{ \frac{\partial}{\partial \phi}[q_\phi(z|x)] \log p(x,z) - \frac{\partial}{\partial \phi}[q_\phi(z|x)] \log q_\phi(z|x) \right\} dz - \underbrace{\int q_\phi(z|x) \frac{\partial}{\partial \phi}[\log q_\phi(z|x)] \, dz}_{=0} \\
&= \int \left\{ \frac{\partial}{\partial \phi}[\log q_\phi(z|x)] \log p(x,z) - \frac{\partial}{\partial \phi}[\log q_\phi(z|x)] \log q_\phi(z|x) \right\} q_\phi(z|x) \, dz \\
&= \mathbb{E}_{q_\phi(z|x)} \left\{ \frac{\partial}{\partial \phi}[\log q_\phi(z|x)] \log \frac{p(x,z)}{q_\phi(z|x)} \right\}.
\end{aligned}
$$

In the third line, we used the score identity. Now we can use unbiased estimator of the gradient of the ELBO by sampling $\{z_j\}_{j=1}^s \overset{iid}{\sim} q_\phi(\cdot|x)$ and computing

$$
\frac{1}{s} \sum_{j=1}^s \frac{\partial}{\partial \phi}[\log q_\phi(z_j|x)] \log \frac{p(x, z_j)}{q_\phi(z_j|x)}.
$$

This approach has an advantage of working for discrete latent variable models. The variance can be a problem for this vanilla version, however, and hence one can consider various approaches for controlling the variance of the gradients, using e.g. Rao-Blackwellization, control variates or importance sampling.

# 6 Bayesian Machine Learning

## 6.1 Bayesian Inference

So far, our treatment of probabilistic machine learning models has been *frequentist*, i.e. we used one set of tools to reason about latent variables $\mathbf{z}$ (e.g. cluster indicators in a mixture model) and another to reason about model parameters $\theta$ (e.g. parameters of mixture components defining those clusters). The generative processes we considered define the *likelihood function*: the joint distribution $p(\mathcal{D}|\theta)$ of all the observed data $\mathcal{D}$ given the model parameters $\theta$ and the learning consists in computing the maximum likelihood estimator

$$\hat{\theta} = \arg\max_{\theta \in \Theta} p(\mathcal{D}|\theta).$$

For example, in the EM algorithm (which is a frequentist method aimed at locally maximising the likelihood function), we were placing a variational distribution $q$ on latent variables but not on $\theta$, which was inferred using point estimates at each iteration.

In **Bayesian machine learning**, we also treat the model parameters $\theta$ as random variables and the process of learning is then simply the computation of the **posterior distribution** $p(\theta|\mathcal{D})$. In addition to the **likelihood** $p(\mathcal{D}|\theta)$ specified by the generative model, one needs to also specify a **prior distribution** $p(\theta)$. Posterior distribution is then given by **Bayes Theorem**:

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})},$$

where the denominator is the **marginal likelihood** or **evidence**:

$$p(\mathcal{D}) = \int_{\Theta} p(\mathcal{D}|\theta)p(\theta)d\theta.$$

All the questions about model parameters can be addressed based on the posterior. We can, for example, consider

- *Posterior mode*: $\widehat{\theta}^{\mathrm{MAP}} = \arg\max_{\theta \in \Theta} p(\theta|\mathcal{D})$ (maximum a posteriori).

- *Posterior mean*: $\widehat{\theta}^{\mathrm{mean}} = \mathbb{E}[\theta|\mathcal{D}]$.

- *Posterior variance*: $\mathrm{Var}[\theta|\mathcal{D}]$.

- *Posterior expectations of functions of parameters*: $\mathbb{E}[g(\theta)|\mathcal{D}]$ for some $g : \Theta \to \mathbb{R}^s$.

A particularly convenient choice of prior distributions are **conjugate priors** to a given likelihood function. A prior and likelihood are said to be conjugate if they result in a posterior that lies in the same parametric family as the prior.

**Example: Bayesian inference on a categorical distribution.**

Suppose we observe $\mathcal{D} = \{y_i\}_{i=1}^n$, with $y_i \in \{1, \ldots, K\}$, and model them as i.i.d. with the probability mass function $\pi = (\pi_1, \ldots, \pi_K)$:

$$p(\mathcal{D}|\pi) = \prod_{i=1}^n \pi_{y_i} = \prod_{k=1}^K \pi_k^{n_k}$$

with $n_k = \sum_{i=1}^n \mathbf{1}(y_i = k)$ and $\pi_k > 0$, $\sum_{k=1}^K \pi_k = 1$. The conjugate prior on $\pi$ is the **Dirichlet distribution** $\mathrm{Dir}(\alpha_1, \ldots, \alpha_K)$ with parameters $\alpha_k > 0$, and density

$$p(\pi) = \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K \pi_k^{\alpha_k - 1}$$

on the probability simplex $\{\pi : \pi_k > 0, \sum_{k=1}^K \pi_k = 1\}$. Since

$$p(\pi|\mathcal{D}) \quad \propto \quad \prod_{k=1}^K \pi_k^{n_k + \alpha_k - 1},$$

the posterior is also Dirichlet $\mathrm{Dir}(\alpha_1 + n_1, \ldots, \alpha_K + n_K)$. Posterior mean is given by

$$\widehat{\pi}_k^{\mathrm{mean}} = \frac{\alpha_k + n_k}{\sum_{j=1}^K \alpha_j + n_j}.$$

Notice how parameters of the prior (hyperparameters) are essentially playing the role of **pseudocounts** for each of the classes $1, \ldots, K$ (but they need not be integer-valued). They are reflecting prior beliefs about class proportions as well as the strength of the prior beliefs.

For the case of two classes, this is equivalent to a **Beta distribution** $\mathrm{Beta}(\alpha_1, \alpha_2)$ prior on $\pi_1$, i.e. $p(\pi_1) = \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_1)} \pi_1^{\alpha_1} (1 - \pi_1)^{\alpha_2}$.

## 6.2 Predictive distributions and Bayesian decision theory

How do we construct predictions based on the posterior distributions? Write the observations as $\mathcal{D} = \{x_i\}_{i=1}^n$ and assume the generative model specifies $p(x|\theta)$, e.g. it can be a mixture model $p(x|\theta) = \sum_{k=1}^K \pi_k f(x|\phi_k)$, with $\theta = (\pi_1, \ldots, \pi_K; \phi_1, \ldots, \phi_K)$. The **posterior predictive distribution** is the conditional distribution of $x_{n+1}$ given $\mathcal{D} = \{x_i\}_{i=1}^n$:

$$p(x_{n+1}|\mathcal{D}) = \int_\Theta p(x_{n+1}|\theta, \mathcal{D}) p(\theta|\mathcal{D}) d\theta$$

$$= \int_\Theta p(x_{n+1}|\theta) p(\theta|\mathcal{D}) d\theta.$$

Thus, we predict new data by *averaging the predictive distribution over the posterior*. This is fundamentally different than predicting using a point estimate of $\theta$, i.e. $p(x_{n+1}|\hat{\theta})$ as it takes into account the posterior uncertainty in parameters.

In case of supervised learning, we are interested typically only in the conditional distribution of outputs given inputs, and not in the marginal distribution over inputs. We can do this with a conditional model $p(y|x, \theta)$ and a prior $p(\theta)$. the posterior distribution given dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{n}$ is:

$$p(\theta|\mathcal{D}) = \frac{p(\theta) \prod_{i=1}^{n} p(y_i|x_i, \theta)}{p(\mathcal{D})}$$

and the posterior predictive distribution is given as

$$p(y_{n+1}|x_{n+1}, \mathcal{D}) = \int p(y|x, \theta) p(\theta|\mathcal{D}) d\theta.$$

## 6.2.1 Example: Bayesian treatment of naïve Bayes classifer.

Consider a $K$-class classification problem with binary input vectors, i.e. $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{n}$, $x_i \in \{0, 1\}^p$ and $y_i \in \{1, \ldots, K\}$. ***Naïve Bayes*** classifier uses the following model:

$$p(y_i = k|\theta) = \pi_k, \quad p(x_i|y_i = k, \theta) = \prod_{j=1}^{p} \phi_{kj}^{x_i^{(j)}} (1 - \phi_{kj})^{1 - x_i^{(j)}},$$

i.e. it assumes that given the class labels, individual dimensions in input vectors are *independent*. The parameters of the model are collated into $\theta = ((\pi_k), (\phi_{kj}))$. It is often used in text classification where data items correspond to documents and $x_i^{(j)}$ indicates whether word $j$ from a list of $p$ words has appeared in document $i$. Class labels correspond to e.g. topics of the documents. Despite the name, naïve Bayes is often treated in a frequentist way, i.e. using maximum likelihood estimation of parameters. If we set $n_k = \sum_{i=1}^{n} \mathbf{1}\{y_i = k\}$, $n_{kj} = \sum_{i=1}^{n} \mathbf{1}(y_i = k, x_i^{(j)} = 1)$, the MLE can be written as

$$\hat{\pi}_k = \frac{n_k}{n}, \qquad\qquad \hat{\phi}_{kj} = \frac{\sum_{i:y_i=k} x_i^{(j)}}{n_k} = \frac{n_{kj}}{n_k}.$$

But the MLEs can be problematic in some cases. For example, if the $\ell$-th word did not appear in any documents labelled as class $k$ ($n_{kl} = 0$), then $\hat{\phi}_{k\ell} = 0$. But if we then wish to compute the predictive probability once for a new document $\tilde{x}$ which contains $\ell$-th word, we have:

$$p(\tilde{y} = k|\tilde{x} \text{ with } \ell\text{-th entry equal to } 1, \hat{\theta})$$

$$\propto \hat{\pi}_k \prod_{j=1}^{p} \left(\hat{\phi}_{kj}\right)^{\tilde{x}^{(j)}} \left(1 - \hat{\phi}_{kj}\right)^{1 - \tilde{x}^{(j)}} = 0,$$

since $\hat{\phi}_{k\ell} = 0$. This means that we will never attribute a new document containing word $\ell$ to class $k$ (regardless of what other words in it may be!). Moreover, probability of a document under all classes can be 0 by the same reasoning.

Let us consider a Bayesian approach to the same model. We can write the likelihood as

$$p(\mathcal{D}|\theta) = \prod_{i=1}^{n} p(x_i, y_i|\theta) = \prod_{i=1}^{n} \prod_{k=1}^{K} \left( \pi_k \prod_{j=1}^{p} \phi_{kj}^{x_i^{(j)}} (1 - \phi_{kj})^{1-x_i^{(j)}} \right)^{\mathbf{1}(y_i=k)}$$

$$= \prod_{k=1}^{K} \pi_k^{n_k} \prod_{j=1}^{p} \phi_{kj}^{n_{kj}} (1 - \phi_{kj})^{n_k - n_{kj}}.$$

For a conjugate prior, we can use $\text{Dir}((\alpha_k)_{k=1}^{K})$ for $\pi$, and $\text{Beta}(a,b)$ for $\phi_{kj}$ independently. Now, because the likelihood factorises, the posterior distribution over $\pi$ and $(\phi_{kj})$ also factorises, and posterior for $\pi$ is $\text{Dir}((\alpha_k + n_k)_{k=1}^{K})$, and for $\phi_{kj}$ is $\text{Beta}(a + n_{kj}, b + n_k - n_{kj})$. If we want to predict a label $\tilde{y}$ for a new document $\tilde{x}$, we obtain

$$p(\tilde{x}, \tilde{y} = k|\mathcal{D}) = p(\tilde{y} = k|\mathcal{D}) p(\tilde{x}|\tilde{y} = k, \mathcal{D})$$

with

$$p(\tilde{y} = k|\mathcal{D}) = \frac{\alpha_k + n_k}{\sum_{l=1}^{K} \alpha_l + n}$$

$$p(\tilde{x}^{(j)} = 1|\tilde{y} = k, \mathcal{D}) = \frac{a + n_{kj}}{a + b + n_k}$$

and the predicted class is

$$p(\tilde{y} = k|\tilde{x}, \mathcal{D}) = \frac{p(\tilde{y} = k|\mathcal{D}) p(\tilde{x}|\tilde{y} = k, \mathcal{D})}{p(\tilde{x}|\mathcal{D})} \propto \frac{\alpha_k + n_k}{\sum_{l=1}^{K} \alpha_l + n} \prod_{j=1}^{p} \left( \frac{a + n_{kj}}{a + b + n_k} \right)^{\tilde{x}^{(j)}} \left( \frac{b + n_k - n_{kj}}{a + b + n_k} \right)^{1-\tilde{x}^{(j)}}.$$

Compared to the MLE plug-in predictions, pseudocounts help to "regularise" probabilities away from the extreme values.

## 6.2.2 Bayesian decision theory

Suppose we have a supervised learning setting with dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{n}$. Given a new input $x_{n+1}$, a system is required to make a prediction or take an action, say $f(x_{n+1})$, which incurs a loss $L(y_{n+1}, f(x_{n+1}), x_{n+1})$ if output is $y_{n+1}$. Empirical risk minimisation directly parameterises and optimises for $f(x)$ as a function of $x$ to reduce empirical risk.

Bayesian approach to decision theory differs from empirical risk minimisation. It starts with a model, parameterised by $\theta$, describing either the joint or the conditional distribution over inputs and outputs. From this model, we get a predictive distribution $p(y|x, \mathcal{D})$. Now a prediction/action $f(x)$ incurs expected loss:

$$\mathbb{E}[L(y, f(x), x)|x, \mathcal{D}] = \int L(y, f(x), x) p(y|x, \mathcal{D}) dy$$

so the optimal prediction/action is:

$$\arg \min_{f(x)} \mathbb{E}[L(y, f(x), x)|x, \mathcal{D}]$$

For example, if the loss is squared loss, $|y - f(x)|^2$, the optimal prediction is the predictive mean $\mathbb{E}[y|x, \mathcal{D}]$, and if the loss is $|y - f(x)|$, the optimal prediction is the predictive median.

Bayesian decision theory separates out the process of modelling and inference from the process of decision making. Advantages include: that one can ascertain and critique the model's fit to data prior to making decisions, that one can quantify the model's uncertainties, and that one can use the same prediction distribution to take different actions given different losses. Disadvantages include: Bayesian inference is typically significantly more computationally expensive and less scalable, and if approximations are required it is unclear how to calibrate the approximation to the loss at hand. In other words, different approximations emphasise differ aspects of the posterior, and the choice of approximations should reflect the aspects of the posterior that are important for the loss function. This then breaks the separation between modelling/inference and decision making.

## 6.3 Approximate Bayesian Inference

In all but the simplest models, the posterior distribution is intractable, and there is a large literature on approaches to approximate it. The important classes include:

- **Markov chain Monte carlo** (MCMC): construct a Markov chain $(X_1, X_2, \ldots)$ whose equilibrium distribution can be shown to be the posterior, and where each Markov kernel $X_{t+1}|X_t$ can be efficiently simulated. Simulate the Markov chain for a large enough number of time steps, and use ergodic averages to approximate posterior statistics of interest.

- **Sequential Monte Carlo** (SMC): construct a finite sequence of importance samplers targeting a sequence of distributions with the last being the posterior distribution. Use resampling and rejuvenation steps to improve the effective sample size of each importance sampler.

- **Approximate Bayesian computation** (ABC): A Monte Carlo approach that is applicable in cases where the likelihood is intractable or unknown, but can be simulated. These methods are sometimes also called **likelihood-free** inference as they do not assume knowledge of the likelihood function.

- **Laplace approximation**: Use a local Gaussian approximation constructed at a maximum a posteriori (MAP) parameter.

- **Variational approximation**: Frame approximation of the posterior as optimising a lower bound on the evidence with respect to a variational distribution.

It is beyond the scope of this course to survey the literature (note that there are two other Statistics Part C courses covering this literature).

## 6.4 Laplace Approximation

Bayesian approach to learning is conceptually very elegant, but the posterior distributions are intractable in almost all interesting cases, and we therefore need to resort to various approximations. One of the techniques for approximation of intractable posterior distributions is the **Laplace approximation** also known as **saddlepoint approximation**. The idea is to simply approximate the posterior distribution $p(\theta|\mathcal{D})$ with a (multivariate) Gaussian distribution. Given the ease of manipulating Gaussians, this is a convenient choice, since the various posterior expectations and predictive distributions can be easier to calculate when we have Gaussian approximate posteriors.

Consider for simplicity the case where parameter $\theta$ is a scalar and assume that posterior mode $\widehat{\theta}^{\mathrm{MAP}}$ is available. Often, the posterior mode can be found even if the normalising constant $p(\mathcal{D})$ is intractable since it suffices to maximise $p(\theta|\mathcal{D}) \propto p(\theta, \mathcal{D}) = p(\mathcal{D}|\theta)p(\theta)$ using a numerical method. Then, we can use a Taylor expansion of $\log p(\theta|\mathcal{D})$ around the posterior mode $\widehat{\theta}^{\mathrm{MAP}}$:

$$
\begin{aligned}
\log p(\theta|\mathcal{D}) \;=\; & \log p(\widehat{\theta}^{\mathrm{MAP}}|\mathcal{D}) + \left.\frac{\partial \log p(\theta|\mathcal{D})}{\partial \theta}\right|_{\theta=\widehat{\theta}^{\mathrm{MAP}}} \left(\theta - \widehat{\theta}^{\mathrm{MAP}}\right) \\
& + \left.\frac{\partial^2 \log p(\theta|\mathcal{D})}{\partial \theta^2}\right|_{\theta=\widehat{\theta}^{\mathrm{MAP}}} \frac{\left(\theta - \widehat{\theta}^{\mathrm{MAP}}\right)^2}{2} + \mathcal{O}\left(\left(\theta - \widehat{\theta}^{\mathrm{MAP}}\right)^3\right).
\end{aligned}
$$

By ignoring the third and higher order terms and noticing that the the first derivative at the mode must be zero, we have an approximation:

$$
\log p(\theta|\mathcal{D}) \approx \log p(\widehat{\theta}^{\mathrm{MAP}}|\mathcal{D}) - \frac{\tau}{2}\left(\theta - \widehat{\theta}^{\mathrm{MAP}}\right)^2, \tag{6.1}
$$

where we write $\tau = -\frac{\partial^2 \log p(\theta|\mathcal{D})}{\partial \theta^2} \geq 0$. But recall that $\log \mathcal{N}\left(\theta|\mu, \sigma^2\right) = \log\left(\left(2\pi\sigma^2\right)^{-1/2}\right) - \frac{1}{2\sigma^2}\left(\theta - \mu\right)^2$, so this second order Taylor approximation has exactly the form of a normal log-density with mean $\mu = \widehat{\theta}^{\mathrm{MAP}}$ and variance $\sigma^2 = \tau^{-1}$ so we can approximate the posterior with $\mathcal{N}\left(\widehat{\theta}^{\mathrm{MAP}}, \left(-\frac{\partial^2 \log p(\theta|\mathcal{D})}{\partial \theta^2}\right)^{-1}\right)$.

This idea easily extends to multivariate densities. In particular, the Laplace approximation of $p(\theta|\mathcal{D})$ is a multivariate Gaussian $\mathcal{N}\left(\widehat{\theta}^{\mathrm{MAP}}, \Sigma\right)$, where *the inverse covariance matrix is given by the negative Hessian of the log-posterior* evaluated at the posterior mode:

$$
\Sigma^{-1} \;=\; -\left.\frac{\partial^2 \log p(\theta|\mathcal{D})}{\partial \theta \partial \theta^\top}\right|_{\theta=\widehat{\theta}^{\mathrm{MAP}}}.
$$

Since $\log p(\theta|\mathcal{D})$ agrees with $\log p(\theta, \mathcal{D})$ up to a constant, they have the same derivatives, so often we work with the *energy function* $J(\theta) = -\log p(\theta, \mathcal{D})$, which is the negative logarithm of the unnormalised posterior. Then we can write

$$\Sigma^{-1} = \left.\frac{\partial^2 J(\theta)}{\partial\theta\partial\theta^\top}\right|_{\theta=\widehat{\theta}^{\mathrm{MAP}}}.$$

## 6.5 Variational Bayes

One of the workhorses of Bayesian machine learning are *variational approximations*, which turn posterior inference in intractable Bayesian models into optimization. We have seen that Bayesian model selection proceeds by optimizing (maximizing) the model evidence. While model evidence is almost always intractable, using the same principles as in the EM algorithm (Gibbs inequality), lower bounds may be available which can be optimized instead. The variational approach to Bayesian machine learning is often referred to as **variational Bayes**.

### 6.5.1 ELBO

Assume that we are taking a Bayesian approach to inference in a latent variable model $p(\mathbf{X}, \mathbf{z}|\theta)$ with observations $\mathbf{X}$, latent variables $\mathbf{z}$ and parameters $\theta$. Now, because we are using a Bayesian model, our treatment of latent variables and model parameters is exactly the same. We can now consider some joint distribution $q(\mathbf{z}, \theta)$ of latent variables and parameters, called variational distribution (similarly to EM, but note that EM was not allowed to place a distribution over $\theta$!). We claim that the quantity

$$\mathcal{F}(q) = \mathbb{E}_q\left[\log p(\mathbf{X}, \mathbf{z}, \theta)\right] + H(q) \tag{6.2}$$

is a lower bound on log-evidence $\log p(\mathbf{X})$. Namely, we can write

$$
\begin{aligned}
\mathcal{F}(q) &= \mathbb{E}_q\left[\log p(\mathbf{X}, \mathbf{z}, \theta)\right] - \mathbb{E}_q[\log q(\mathbf{z}, \theta)] \\
&= \mathbb{E}_q\left[\log p(\mathbf{z}, \theta|\mathbf{X})\right] + \log p(\mathbf{X}) - \mathbb{E}_q[\log q(\mathbf{z}, \theta)] \\
&= -\mathrm{KL}\left(q(\mathbf{z}, \theta)||p(\mathbf{z}, \theta|\mathbf{X})\right) + \log p(\mathbf{X}),
\end{aligned}
$$

which is by Gibbs inequality maximised (and equal to log-evidence) when KL is zero, i.e. when $q(\mathbf{z}, \theta) = p(\mathbf{z}, \theta|\mathbf{X})$. Thus, for any variational distribution $q$, $\mathcal{F}(q) \leq \log p(x)$. Expression 6.2 is called the **evidence lower bound** (ELBO).

To reason about all the unknowns in the model, we would simply need to compute the joint posterior $p(\mathbf{z}, \theta|\mathbf{X})$, but this is almost always intractable. Hence, variational Bayesian inference *approximates* the posterior by starting with a family $\mathcal{Q}$ of tractable variational distributions $q(\mathbf{z}, \theta)$ (e.g. $q(\mathbf{z}, \theta|\nu)$ where $\nu$ are the **variational parameters**), and aims to minimize the divergence $\mathrm{KL}\left(q(\mathbf{z}, \theta)||p(\mathbf{z}, \theta|\mathbf{X})\right)$ over $\mathcal{Q}$ or, equivalently, maximise the ELBO, i.e. find the tightest lower bound on the log-evidence.

In a nutshell, variational Bayes projects the (intractable) posterior $p(\mathbf{z}, \theta|\mathbf{X})$ onto a tractable family $\mathcal{Q}$ with respect to the KL divergence $\mathrm{KL}\left(q(\mathbf{z}, \theta)||p(\mathbf{z}, \theta|\mathbf{X})\right)$. Alternative divergences are possible in this context. In particular, since KL is not symmetric, minimization of the "reverse" divergence $\mathrm{KL}\left(p(\mathbf{z}, \theta|\mathbf{X})||q(\mathbf{z}, \theta)\right)$ results in a different family of approximate Bayesian methods, known as expectation propagation (EP).

## 6.5.2 Bayesian EM and Mean-Field Variational Family

Variational approximation requires specifying the variational family $\mathcal{Q}$. The complexity of $\mathcal{Q}$ determines the difficulty of the optimization; it is more difficult to optimize over a large family $\mathcal{Q}$ than over a simpler, smaller one. Consider family $\mathcal{Q}$ of variational distributions which factorize across the latents and the parameters: $q(\mathbf{z}, \theta) = q_{\mathbf{Z}}(\mathbf{z}) q_{\Theta}(\theta)$. For a fixed $q_{\Theta}$, we can solve for $q_{\mathbf{Z}}$ which maximises ELBO (*exercise):*

$$q_{\mathbf{Z}}(\mathbf{z}) \propto \exp\left(\int \log p(\mathbf{X}, \mathbf{z}, \theta) q_{\Theta}(\theta) \, d\theta\right),$$

and by symmetry, for a fixed $q_{\mathbf{Z}}$, we can solve for $q_{\Theta}$ which maximises ELBO:

$$q_{\Theta}(\theta) \propto \exp\left(\int \log p(\mathbf{X}, \mathbf{z}, \theta) q_{\mathbf{Z}}(\mathbf{z}) \, d\mathbf{z}\right).$$

Now, one can formulate an algorithm similar to EM, which alternates between optimising $q_{\mathbf{Z}}$ and $q_{\Theta}$, such that each iteration increases ELBO and thus decreases the KL divergence from the posterior. Hence such an algorithm is sometimes called a **variational Bayesian EM** algorithm or VBEM.

We notice the symmetry between $\mathbf{z}$ and $\theta$. Indeed, the distinction between parameters and latent variables disappears in Bayesian modelling, as all unobserved quantities in the model are treated in the same way and our goal is to approximate their posterior distribution. In the rest of this section, we will drop $\theta$ from the notation and treat them as a part of the set of all unobserved quantities $\mathbf{z}$.

The further simplification we often make in Variational Bayes is to focus on a **mean-field approximation** where the variational distribution fully factorizes

$$q(\mathbf{z}) = \prod_{j=1}^{m} q_j(z_j),$$

i.e. all latent variables are mutually independent and each latent $z_j$ is governed by its own variational factor $q_j$. Note that there could be a mix between categorical and continuous latents, each having the appropriate factor $q_j$. Also, $z_j$ itself need not be a univariate latent – see an example with LDA below. Using the mean-field family implies that we will not be able to capture any posterior correlations between the latent variables $z_j$ and $z_{j'}$ for $j \neq j'$ and that the best we can hope for is a rich representations of the posterior marginals.

The iterative procedure similar to Bayesian EM can now be applied to each individual factor, giving rise to the algorithm 6.1 called **coordinate ascent variational**

---

**Algorithm 6.1** Coordinate Ascent Variational Inference (CAVI)

---

**Input**: a model $p(\mathbf{z}, \mathbf{x})$, dataset $\mathbf{x}$
**Output**: a variational posterior $q(\mathbf{z})$

**while** the ELBO has not converged **do**

- **for** $j = 1, \dots, m$
    - $q_j(z_j) \propto \exp\left[\mathbb{E}_{\mathbf{z}_{-j} \sim q} \log p\left(z_j | \mathbf{z}_{-j}, \mathbf{x}\right)\right]$
- $\text{ELBO}(q) = \mathbb{E}_{\mathbf{z} \sim q}\left[\log p(\mathbf{x}, \mathbf{z})\right] + H(q)$

**return** $q\left(\mathbf{z}\right) = \prod_{j=1}^{m} q_j\left(z_j\right)$

---

inference (CAVI), closely related to Gibbs sampling (a popular class of MCMC algorithms), i.e. it also uses full conditionals $p\left(z_j | \mathbf{z}_{-j}, \mathbf{x}\right) \propto p(\mathbf{z}, \mathbf{x})$, where we denoted $\mathbf{z}_{-j} = [z_1, \dots, z_{j-1}, z_{j+1}, \dots, z_n]$.

### 6.5.3 Complete conditionals in the exponential family

When the complete conditionals $p\left(z_j | \mathbf{z}_{-j}, \mathbf{x}\right)$ belong to an exponential family of distributions, i.e. they are given by

$$p(z_j | \mathbf{z}_{-j}, \mathbf{x}) = h\left(z_j\right) \exp\left[\eta_j\left(\mathbf{z}_{-j}, \mathbf{x}\right)^{\top} z_j - A\left(\eta_j\left(\mathbf{z}_{-j}, \mathbf{x}\right)\right)\right],$$

a particular convenient form of CAVI is available, i.e. the updates in Algorithm 6.1 are available in closed form. Above, we assume $z_j$ is already transformed to its appropriate sufficient statistic, $h\left(\cdot\right)$ is a base measure, $A\left(\cdot\right)$ is the log-normalizer and $\eta_j$ are the natural parameters (which depend on the conditioning set). Now, the CAVI update reads

$$\begin{aligned} q_j(z_j) &\propto \exp\left[\mathbb{E}_{-j} \log p\left(z_j | \mathbf{z}_{-j}, \mathbf{x}\right)\right] \\ &= \exp\left[\log h\left(z_j\right) + \{\mathbb{E}_{-j}\eta_j\left(\mathbf{z}_{-j}, \mathbf{x}\right)\}^{\top} z_j - \mathbb{E}_{-j}A\left(\eta_j\left(\mathbf{z}_{-j}, \mathbf{x}\right)\right)\right] \\ &\propto h\left(z_j\right) \exp\left[\{\mathbb{E}_{-j}\eta_j\left(\mathbf{z}_{-j}, \mathbf{x}\right)\}^{\top} z_j\right] \end{aligned}$$

and thus, the variational factors are in the same exponential family as the complete conditionals with natural parameter being the expected natural parameter of the complete conditional

$$\nu_j = \mathbb{E}_{-j}\eta_j\left(\mathbf{z}_{-j}, \mathbf{x}\right).$$

This setup describes many models, including Bayesian Gaussian mixtures, Dirichlet process mixtures, matrix factorization, multilevel regression and latent Dirichlet allocation, giving thus one classical overarching CAVI algorithm with closed-form updates for many instances of Variational Bayes.

While the distinction between parameters and latent variables disappears in Bayesian modelling, there is still a relevant distinction in terms of where in the model hierarchy these unobserved quantities appear. In that respect, we can differentiate *global latent vector* $\theta$ which is associated to all observations and the *local latents* $\{z_i\}_{i=1}^n$ each of which is associated to an individual observation $x_i$, such that the observation $x_i$ is conditionally independent of $\{z_1, \ldots, z_{i-1}, z_{i+1}, \ldots, z_n\}$ given $\theta$ and $z_i$. The joint density is then

$$p(\theta, \mathbf{z}, \mathbf{x}) = p(\theta) \prod_{i=1}^n p(z_i, x_i | \theta).$$

The normal mixtures are an example of this, where the mixture parameters are the global latents, while cluster assignments are the local latents. The impact of such hierarchy is that not all updates in Algorithm 6.1 need to be performed sequentially. Multiple levels of hierarchy are also possible.

We next study a concrete example of this in topic modelling, Latent Dirichlet Allocation [4].

## 6.5.4 Example: Topic Modelling

Topic models are a class of probabilistic models of text that lead to parsimonious representations of hidden thematic structure of a collection of documents. A popular approach to topic modelling is Latent Dirichlet Allocation (LDA[1]) [4].

### Latent Dirichlet Allocation

**Latent Dirichlet allocation** (LDA) captures the intuition that a text document typically exhibits multiple topics and blends them in a particular way. In LDA, each topic is modelled as a probability distribution over words and each document as a mixture of corpus-wide topics (i.e. it can be identified with a distribution over topics). Each observed word in a document is then treated as a draw from the mixture, i.e. it belongs to one of the topics (mixture components). Mixture proportions are thus unique for each document, i.e. they are local latents, but mixture components are shared across the whole collection - they are global latents.

This setting is also called a **mixed membership model**. Another important example of mixed membership model is the STRUCTURE model in population genetics. There, the DNA of each individual in a population is modelled as a mixture over ancestral populations, with each individual having different proportions of their DNA coming from each ancestral population. This is most clear for individuals of mixed ethnic ancestry, but such approaches are also applicable to, e.g. "native" British, where the the inferred population structure has been used to understand the history and migration patterns of the peoples of the British isles before modern times [15].

LDA posits the following conditionally conjugate model, Let $K$ be the number of topics and $V$ the size of the vocabulary.

---

[1]Do not confuse it with Linear Discriminant Analysis which shares the acronym - these two models are not related.
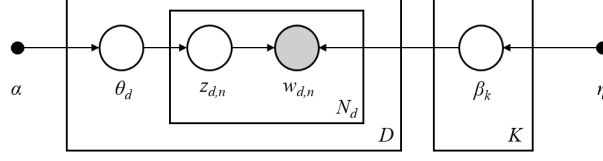
Figure 6.1: Graphical model representation of LDA. Plates represent replication, for example there are $D$ documents each having a topic proportion vector $\theta_d$

1. For each topic in $k = 1, \ldots, K$,

   a) Draw a distribution over $V$ words $\beta_k \sim \mathrm{Dir}_V(\eta)$

2. For each document in $d = 1, \ldots, D$,

   a) Draw a vector of topic proportions $\theta_d \sim \mathrm{Dir}_K(\alpha)$

   b) For each word in $n = 1, \ldots, N_d$,

       i. Draw a topic assignment $z_{dn} \sim \mathrm{Mult}(\theta_d)$, i.e. $p(z_{dn} = k | \theta_d) = \theta_{dk}$

       ii. Draw a word $w_{dn} \sim \mathrm{Mult}(\beta_{z_{dn}})$, i.e. $p(w_{dn} = v | \beta, z) = \beta_{z_{dn}v}$

The goal of LDA is to find the posterior

$$p(\text{topics,proportions,assignments}|\text{observed words})$$

Note that a corpus of text to be analyzed may consist of millions of documents, thus having possibly billions of latent variables. We can write the joint distribution as

$$
\begin{aligned}
p(\beta, \theta, z, w) &= \prod_{k=1}^{K} p(\beta_k; \eta) \prod_{d=1}^{D} \left\{ p(\theta_d; \alpha) \prod_{n=1}^{N_d} p(z_{dn}|\theta_d) p(w_{dn}|\beta, z) \right\} \\
&= \frac{1}{B(\eta)^K B(\alpha)^D} \prod_{k=1}^{K} \prod_{v=1}^{V} \beta_{kv}^{\eta_v - 1} \prod_{d=1}^{D} \left\{ \prod_{k=1}^{K} \theta_{dk}^{\alpha_k - 1} \prod_{n=1}^{N_d} \theta_{d,z_{dn}} \beta_{z_{dn}, w_{dn}} \right\} (6.3)
\end{aligned}
$$

The model has the following latents: $\beta$ (topics), $\theta$ (proportions), and $z$ (assignments). Note that there are also hyperparameter vectors $\eta \in \mathbb{R}_+^V$ and $\alpha \in \mathbb{R}_+^K$ in the Dirichlet priors - these are assumed fixed. Data are the observed words $\{w_{dn}\}$. There will be some abuse of notation here - we denote by $w_{dn}$ both the appropriate draw from vocabulary $\{1, \ldots, V\}$ - to be used for indexing, and its "one-hot" encoding i.e. a binary $V$-vector with $w_{dn}[v] = 1$ if $w_{dn} = v$ and zero otherwise. We will write $w_{dn}[\cdot]$ in the case of the latter. Similarly, we denote by $z_{dn}$ both the appropriate topic assignment from $\{1, \ldots, K\}$ and its "one-hot" encoding i.e. a binary K-vector with $z_{dn}[k] = 1$ if $z_{dn} = k$ and zero otherwise. We will write $z_{dn}[\cdot]$ in the case of the latter.

We will use a mean-field family of the form

$$
q(\beta, \theta, z) = \prod_{k=1}^{K} q(\beta_k; \lambda_k) \prod_{d=1}^{D} \left\{ q(\theta_d; \gamma_d) \prod_{n=1}^{N_d} q(z_{dn}; \phi_{dn}) \right\}.
$$

The complete conditionals are proportional to the joint distribution in (6.3):

1. Complete conditional on the topic assignment is a multinomial with

$$p\left(z_{dn} = k | \theta_d, \beta, w_d\right) \propto \theta_{dk}\beta_{k,w_{dn}} = \exp\left(\log\theta_{dk} + \log\beta_{k,w_{dn}}\right). \qquad (6.4)$$

Thus, for the variational approximation we also use a multinomial but with a "free parameter" $\phi_{dn}$, where we denote $\phi_{dn}[k] = q\left(z_{dn} = k\right)$, i.e. $\phi_{dn}$ is simply a probability mass function over $K$ topics.

2. Complete conditional on the topic proportions depends only on the assignments and is given by

$$p\left(\theta_d | z_d\right) = \text{Dir}_K\left(\theta_d; \alpha + \sum_{n=1}^{N_d} z_{dn}\left[\cdot\right]\right). \qquad (6.5)$$

For the variational approximation we also use Dirichlet, with parameter vector $\gamma_d \in \mathbb{R}_+^K$.

3. Complete conditional on the topics is

$$p\left(\beta_k | z, w\right) = \text{Dir}_V\left(\beta_k; \eta + \sum_{d=1}^{D}\sum_{n=1}^{N_d} z_{dn}\left[k\right] w_{dn}\left[\cdot\right]\right). \qquad (6.6)$$

For the variational approximation we also use Dirichlet, with parameter vector $\lambda_k \in \mathbb{R}_+^V$.

With these full conditionals we can derive the CAVI updates in the LDA model. We will need the following fact about the Dirichlet distribution given here without proof.

**Fact 18.** *If $\pi \sim Dir_L\left(\alpha\right)$, then*

$$\mathbb{E}\left[\log\pi_j\right] = \psi\left(\alpha_j\right) - \psi\left(\sum_{\ell=1}^{L}\alpha_\ell\right),$$

*where $\psi\left(u\right) = \frac{\Gamma'(u)}{\Gamma(u)} = \int_0^\infty\left(\frac{e^{-t}}{t} - \frac{e^{-ut}}{1-e^{-t}}\right)dt$ is the digamma function.*

Now we can obtain the closed-form updates for each set of the latents.

**Proposition 19.** *CAVI updates in the LDA model are given by*

*1. $\phi_{dn}[k] \propto \exp\left(\psi\left(\gamma_{dk}\right) + \psi\left(\lambda_{k,w_{dn}}\right) - \psi\left(\sum_{v=1}^{V}\lambda_{k,v}\right)\right),$*

*2. $\gamma_d = \alpha + \sum_{n=1}^{N_d}\phi_{dn},$*

*3. $\lambda_k = \eta + \sum_{d=1}^{D}\sum_{n=1}^{N_d}\phi_{dn}\left[k\right]w_{dn}\left[\cdot\right],$*

*where $\psi$ is the digamma function.*

*Proof.* Steps (2) and (3) directly follow from the exponential family properties of Dirichlet distribution and are left for exercise. For (1), we make use of Fact 18 and write

$$
\begin{aligned}
\phi_{dn}[k] \quad \propto \quad & \exp\left(\mathbb{E}_{\theta_d,\beta \sim q} \log p\left(z_{dn} = k | \theta_d, \beta, w_d\right)\right) \\
\propto \quad & \exp\left(\mathbb{E}_{\theta_d \sim q} \log \theta_{dk} + \mathbb{E}_{\beta_k \sim q} \log \beta_{k,w_{dn}}\right) \\
\propto \quad & \exp\left(\psi\left(\gamma_{dk}\right) - \psi\left(\sum_{\ell=1}^{K} \gamma_{d\ell}\right) + \psi\left(\lambda_{k,w_{dn}}\right) - \psi\left(\sum_{v=1}^{V} \lambda_{k,v}\right)\right) \\
\propto \quad & \exp\left(\psi\left(\gamma_{dk}\right) + \psi\left(\lambda_{k,w_{dn}}\right) - \psi\left(\sum_{v=1}^{V} \lambda_{k,v}\right)\right),
\end{aligned}
$$

as required. $\qquad\square$

# 7 Gaussian Processes and Bayesian Optimisation

## 7.1 Different views of regression

Regression with least squares loss $L(y, f(x)) = (y - f(x))^2$ implies that we are fitting the *conditional mean* function $f^*(x) = \mathbb{E}[Y|X = x]$. This loss also corresponds to the probabilistic model where $y_i$ is a noisy version of the underlying function $f$ evaluated at input $x_i$:

$$y_i | f(x_i) \sim \mathcal{N}(f(x_i), \sigma^2), \quad \text{independently for} \quad i = 1, \dots, n. \tag{7.1}$$

There are different ways to model the class of functions $f$.

- *Frequentist Parametric* approach: model $f$ as $f_\theta$ for some parameter vector $\theta$. Fit $\theta$ by ML / ERM with squared loss (**linear regression**).

- *Frequentist Nonparametric* approach: model $f$ as the unknown parameter taking values in an infinite-dimensional space of functions (RKHS). Fit $f$ by *regularized* ML / ERM with squared loss (**kernel ridge regression**)

- *Bayesian Parametric* approach: model $f$ as $f_\theta$ for some parameter vector $\theta$. Put a prior on $\theta$ and compute a posterior $p(\theta|\mathcal{D})$ (**Bayesian linear regression**).

- *Bayesian Nonparametric* approach: treat $f$ as the random variable taking values in an infinite-dimensional space of functions. Put a prior over functions $f \in \mathcal{F}$, and compute a posterior $p(f|\mathcal{D})$ (**Gaussian Process regression**).

## 7.2 Gaussian Process Regression

**Gaussian processes** (GPs) are a widely used class of models that allow us to place a prior distribution directly on the space of functions rather than on parameters in a particular family of functions. This prior can then be converted into a posterior distribution once we have seen some data.

One can think of a Gaussian process as an infinite-dimensional generalisation of a multivariate normal distribution. Namely, given an *index set* $\mathcal{X}$, a collection of random variables $\{A_x\}_{x \in \mathcal{X}}$ is said to be a Gaussian process if and only if for every finite set of indices $x_1, \dots, x_n$, vector $[A_{x_1}, \dots, A_{x_n}]^\top$ has a multivariate normal distribution on $\mathbb{R}^n$. Thus, to any Gaussian process, we can associate a random function $f \colon \mathcal{X} \to \mathbb{R}$ by setting

$f(x) = A_x$, for all $x \in \mathcal{X}$. Gaussian process is fully specified by its **mean function** and **covariance function**, i.e.

$$
\begin{aligned}
m(x) &= \mathbb{E}[f(x)], \\
k(x, x') &= \mathbb{E}[(f(x) - m(x))(f(x') - m(x'))],
\end{aligned}
$$

where expectations are taken over $f$ ($x$ and $x'$ are fixed elements in the index set $\mathcal{X}$). This means that for any finite set $x_1, \ldots, x_n$, $\mathbf{f} = [f(x_1), \ldots, f(x_n)]^\top \in \mathbb{R}^n$ has a distribution $\mathcal{N}(\mathbf{m}, \mathbf{K})$, where $\mathbf{m}_i = m(x_i)$ and $\mathbf{K}$ is the covariance matrix given by $\mathbf{K}_{ij} = k(x_i, x_j)$.

We will typically assume that the mean function $m(x)$ is zero under the Gaussian process prior. If we know before seeing any data that the distribution of the function evaluations should be centered around some other mean, we could easily include that into the model. Equivalently, we could also subtract that known mean from the data and just use the zero mean model. If we are looking at the data to estimate the mean function, then often the zero mean GP suffices – in fact, structural information about mean functions (constant, linear) can be included into the choice of the covariance function (*exercises*). Covariance functions $k \colon \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ obviously has to be positive definite so they are essentially equivalent to the kernel functions we have seen before. In fact, there is a rich connection between RKHS methods and Gaussian processes, an example of which we will discuss below.

## 7.2.1 Gaussian Conditioning and Regression Model

The convenience of manipulating multivariate normal distributions carries over to Gaussian processes. Let us review the rules for Gaussian conditioning, which are key to Gaussian process regression.

**Gaussian Conditioning**. Let $\mathbf{z} \sim \mathcal{N}(\mu, \Sigma)$ be a multivariate normal random vector and let us split its dimensions into two parts, i.e.

$$
\mathbf{z} = \begin{bmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \end{bmatrix}, \quad \mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}. \tag{7.2}
$$

Note that $\Sigma_{21} = \Sigma_{12}^\top$ due to symmetry of covariance matrices. Then the conditional density of $\mathbf{z}_2$ given $\mathbf{z}_1$ is also normal and given by

$$
p(\mathbf{z}_2 | \mathbf{z}_1) = \mathcal{N}(\mathbf{z}_2; \mu_2 + \Sigma_{21}\Sigma_{11}^{-1}(\mathbf{z}_1 - \mu_1), \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}). \tag{7.3}
$$

For a given set of inputs $\mathbf{x} = \{x_i\}_{i=1}^n$, we denote the vector of evaluations of $f$ by $\mathbf{f} = [f(x_1), \ldots, f(x_n)]^\top \in \mathbb{R}^n$ and the vector of observed outputs by $\mathbf{y} = [y_1, \ldots, y_n]^\top \in \mathbb{R}^n$. Note that since we treat $f$ as a random function, $\mathbf{f}$ is a random $n$-dimensional vector. The Gaussian process regression model, assuming likelihood function in (7.1), is then

given by

$$\mathbf{f} \sim \mathcal{N}(0, \mathbf{K})$$
$$\mathbf{y}|\mathbf{f} \sim \mathcal{N}(\mathbf{f}, \sigma^2 I),$$

where $\mathbf{K}$ is the covariance (kernel) matrix given by $\mathbf{K}_{ij} = k(x_i, x_j)$. But because both the prior and the likelihood are normal this simply means that $\mathbf{f}$ and $\mathbf{y}$ are *jointly normal* with

$$\begin{bmatrix} \mathbf{f} \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N}\left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{K} & \mathbf{K} \\ \mathbf{K} & \mathbf{K} + \sigma^2 I \end{bmatrix} \right). \tag{7.4}$$

For example, to find the cross-covariance between $\mathbf{f}$ and $\mathbf{y}$ note that

$$\mathbb{E}\left[\mathbf{f}\mathbf{y}^\top\right] = \mathbb{E}\left[\mathbf{f}(\mathbf{f} + \sigma\epsilon)^\top\right] = \mathbb{E}\left[\mathbf{f}\mathbf{f}^\top\right] + \sigma\mathbb{E}\left[\mathbf{f}\epsilon^\top\right] = \mathbf{K}, \tag{7.5}$$

where $\epsilon \sim \mathcal{N}(0, I)$ is independent of $\mathbf{f}$. Now, we can simply apply the Gaussian conditioning to find the posterior distribution

$$\mathbf{f}|\mathbf{y} \sim \mathcal{N}(\mathbf{K}(\mathbf{K} + \sigma^2 I)^{-1}\mathbf{y}, \mathbf{K} - \mathbf{K}(\mathbf{K} + \sigma^2 I)^{-1}\mathbf{K}).$$

This gives as the posterior distribution of the evaluations of the unknown function at the set of inputs where we have observed noisy evaluations $\mathbf{y}$.

## 7.2.2 Posterior Predictive Distribution

But we can continue with this formalism further and construct the *posterior predictive distribution*. Suppose $\mathbf{x}' = \{x'_j\}_{j=1}^m$ is a test set. We can extend our model to include the function values $\mathbf{f}' = [f(x'_1), \ldots, f(x'_m)]^\top \in \mathbb{R}^m$ at the test set. The prior can now be extended to include $\mathbf{f}'$ (recall that our prior was on the whole function – not on its values at specific locations!), so that the model reads:

$$\begin{bmatrix} \mathbf{f} \\ \mathbf{f}' \end{bmatrix} | \mathbf{x}, \mathbf{x}' \sim \mathcal{N}\left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{\mathbf{xx}} & \mathbf{K}_{\mathbf{xx}'} \\ \mathbf{K}_{\mathbf{x}'\mathbf{x}} & \mathbf{K}_{\mathbf{x}'\mathbf{x}'} \end{bmatrix} \right)$$
$$\mathbf{y}|\mathbf{f} \sim \mathcal{N}(\mathbf{f}, \sigma^2 I)$$

where $(\mathbf{K}_{\mathbf{xx}})_{ij} = k(x_i, x_j)$, $(\mathbf{K}_{\mathbf{x}'\mathbf{x}'})_{ij} = k(x'_i, x'_j)$, $\mathbf{K}_{\mathbf{xx}'}$ is an $n \times m$ matrix with $(i, j)$-th entry $k(x_i, x'_j)$ and $\mathbf{K}_{\mathbf{x}'\mathbf{x}} = \mathbf{K}_{\mathbf{xx}'}^\top$ We are now making use of the joint normality of $\mathbf{f}'$ and $\mathbf{y}$:

$$\begin{bmatrix} \mathbf{f}' \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N}\left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{\mathbf{x}'\mathbf{x}'} & \mathbf{K}_{\mathbf{x}'\mathbf{x}} \\ \mathbf{K}_{\mathbf{xx}'} & \mathbf{K}_{\mathbf{xx}} + \sigma^2 I \end{bmatrix} \right) \tag{7.6}$$

and from Gaussian conditioning rules again, we can read off the posterior predictive distribution as

$$\mathbf{f}'|\mathbf{y} \sim \mathcal{N}\left( \mathbf{K}_{\mathbf{x}'\mathbf{x}}(\mathbf{K}_{\mathbf{xx}} + \sigma^2 I)^{-1}\mathbf{y}, \mathbf{K}_{\mathbf{x}'\mathbf{x}'} - \mathbf{K}_{\mathbf{x}'\mathbf{x}}(\mathbf{K}_{\mathbf{xx}} + \sigma^2 I)^{-1}\mathbf{K}_{\mathbf{xx}'} \right). \tag{7.7}$$

Thus, we also have a closed form expression for the posterior distribution of the evaluations of the unknown function at any collection of inputs in $\mathcal{X}$. While this follows

directly from the joint normality and Gaussian conditioning rules, it is instructive to notice that we could have arrived at the posterior predictive by integrating $p(\mathbf{f}'|\mathbf{f})$ through the posterior $p(\mathbf{f}|\mathbf{y})$, i.e.

$$p(\mathbf{f}'|\mathbf{y}) = \int p(\mathbf{f}'|\mathbf{f}) p(\mathbf{f}|\mathbf{y}) d\mathbf{f}. \tag{7.8}$$

This follows from $\int \mathcal{N}(a|Bc, D)\mathcal{N}(c|e, F)dc = \mathcal{N}(a|Be, D + BFB^\top)$ (*exercises*). Namely, even if we no longer have the Gaussian observation model (7.1) and $\mathbf{y}$ and $\mathbf{f}$ are no longer jointly normal, we can still use (7.8) to reason about the posterior predictive distribution.

### 7.2.3 Kernel Ridge Regression vs Gaussian Process Regression

If kernel ridge regression (KRR) uses the same kernel as the covariance function in Gaussian process regression (GPR) and moreover, if the regularisation parameter $\lambda$ in KRR is the same as the noise variance $\sigma^2$ in GPR, KRR estimate of the function coincides with the GPR posterior mean. Indeed, recall that in KRR we are solving empirical risk minimisation

$$\min_{f \in \mathcal{H}_k} \sum_{i=1}^{n} (y_i - f(x_i))^2 + \sigma^2 \|f\|_{\mathcal{H}_k}^2,$$

and are fitting a function of the form $f(x) = \sum_{i=1}^{n} \alpha_i k(\cdot, x_i)$. Closed form solution is given by $\alpha = \left(\mathbf{K_{xx}} + \sigma^2 I\right)^{-1} \mathbf{y}$. But then if we wish to predict function values at a new set $\mathbf{x}' = \{x'_j\}_{j=1}^{m}$ of input vectors, we have

$$f(x'_j) = \sum_{i=1}^{n} \alpha_i k(x'_j, x_i) = \left[ k(x'_j, x_1), \dots, k(x'_j, x_n) \right] (\mathbf{K_{xx}} + \sigma^2 I)^{-1} \mathbf{y},$$

and $\left[ k(x'_j, x_1), \dots, k(x'_j, x_n) \right]$ is the $j$-th row of $\mathbf{K_{x'x}}$, so this is the same as the mean in (7.7). Note that GPR also gives predictive variance, a measure of uncertainty, which can be important when making predictions far away from the input data. There are other important differences between the two approaches: KRR is frequentist, while GPR is Bayesian, and thus the hyperparameters are fitted in different ways. KRR typically uses cross-validation and grid search, while GPR, as we discuss next, uses maximum marginal likelihood or a fully Bayesian treatment with hyperparameters integrated out.

## 7.3 Hyperparameter Selection

Probabilistic model given by Gaussian processes allows principled selection of hyperparameters in the model (parameters of the kernel function and the noise variance in the likelihood (7.1)) using *maximum marginal likelihood.*

Marginal likelihood of the hyperparameter vector $\theta = (\nu, \sigma^2)$ which would generally include kernel parameters $\nu$ as well as the standard deviation $\sigma^2$ of the noise in the observation model, is given by

$$p(\mathbf{y}|\theta) = \int p(\mathbf{y}|\mathbf{f}, \theta) p(\mathbf{f}|\theta) d\mathbf{f} = \mathcal{N}\left(\mathbf{y}; 0, \mathbf{K}_\nu + \sigma^2 I\right).$$

We will introduce the shorthand $\mathbf{K}_{\theta+} = \mathbf{K}_\nu + \sigma^2 I$. Thus, we can write the marginal log-likelihood as

$$\log p(\mathbf{y}|\theta) = -\frac{1}{2}\log |\mathbf{K}_{\theta+}| - \frac{1}{2}\mathbf{y}^\top \mathbf{K}_{\theta+}^{-1}\mathbf{y} - \frac{n}{2}\log(2\pi). \tag{7.9}$$

In general, marginal log-likelihood is a nonconvex function of the parameter vector $\theta$ and it can have multiple maxima - thus we typically resort to numerical optimisation methods, such as gradient ascent. The derivative with respect to $\theta_i$ (*exercise*) has the form

$$\frac{\partial}{\partial \theta_i}\log p(\mathbf{y}|\theta) = -\frac{1}{2}\mathrm{Tr}\left(\mathbf{K}_{\theta+}^{-1}\frac{\partial \mathbf{K}_{\theta+}}{\partial \theta_i}\right) + \frac{1}{2}\mathbf{y}^\top \mathbf{K}_{\theta+}^{-1}\frac{\partial \mathbf{K}_{\theta+}}{\partial \theta_i}\mathbf{K}_{\theta+}^{-1}\mathbf{y}. \tag{7.10}$$

Some common kernel choices in this context involve *Automatic Relevance Determination (ARD)* kernel

$$k(x, x') = \tau^2 \exp\left(-\sum_{j=1}^{p}\frac{(x^{(j)} - x'^{(j)})^2}{\eta_j^2}\right), \tag{7.11}$$

which has a global scale parameter $\tau$ as well as one *bandwidth* parameter $\eta_j$ per covariate dimension $j$. If in the hyperparameter selection, very large values of $\eta_j$ are selected, this essentially means that the dimension $j$ is switched off (does not contribute to the kernel function). This is very useful in applications where it is likely that not all dimensions will be relevant.

In addition to maximum marginal likelihood, we can also perform full Bayesian inference for hyperparameters. Namely, we could start with a prior $p(\theta)$ on $\theta$ and draw samples from the posterior

$$p(\theta|\mathbf{y}) \quad \propto \quad p(\theta)p(\mathbf{y}|\theta) = p(\theta)\int p(\mathbf{y}|\mathbf{f}, \theta)p(\mathbf{f}|\theta)d\mathbf{f}.$$

This means that we can integrate uncertainty over hyperparameters into predictions as well, and approximate (integral is typically not available in closed form)

$$p(\mathbf{f}'|\mathbf{y}) = \int p(\mathbf{f}'|\mathbf{y}, \theta)p(\theta|\mathbf{y})d\theta.$$

## 7.4 Gaussian Processes for Classification

In Bayesian classification problems, we are interested in modelling the posterior probabilities of the categorical response variable given a set of training examples and a new input vector. These probabilities must lie in the interval $(0, 1)$ while a Gaussian process models functions that have output on the entire real axis. Thus, it is necessary to adapt

Gaussian processes by transforming their outputs using an appropriate nonlinear activation/link function. Consider the binary classification model with classes $-1$ and $+1$, using the logistic sigmoid:

$$p(y_i = +1|f(x_i)) = \sigma(f(x_i)) = \frac{1}{1 + e^{-f(x_i)}}. \tag{7.12}$$

This non-Gaussian form of the likelihood function, however, renders exact posterior inference intractable and approximate methods are needed. There are a number of approximate schemes that can be used but we will focus here on Laplace approximation. We know that

$$
\begin{aligned}
\log p(\mathbf{f}|\mathbf{y}) &= \text{const} + \log p(\mathbf{f}) + \log p(\mathbf{y}|\mathbf{f}) \\
&= \text{const} - \frac{1}{2}\mathbf{f}^\top \mathbf{K}^{-1}\mathbf{f} + \sum_{i=1}^{n} \log \sigma(y_i f(x_i)).
\end{aligned}
$$

Thus, we can compute the gradient

$$\frac{\partial \log p(\mathbf{f}|\mathbf{y})}{\partial \mathbf{f}} = -\mathbf{K}^{-1}\mathbf{f} + \mathbf{g_f}, \tag{7.13}$$

where the gradient of the likelihood is $\mathbf{g_f} = \frac{\partial \log p(\mathbf{y}|\mathbf{f})}{\partial \mathbf{f}}$ with $[\mathbf{g_f}]_i = \frac{\partial \log p(\mathbf{y}|\mathbf{f})}{\partial f_i} = \sigma(-y_i f(x_i))y_i$. The Hessian is given by

$$\frac{\partial^2 \log p(\mathbf{f}|\mathbf{y})}{\partial \mathbf{f} \partial \mathbf{f}^\top} = -\mathbf{K}^{-1} - \mathbf{D}_f, \tag{7.14}$$

where $\mathbf{D}_f = -\frac{\partial^2 \log p(\mathbf{y}|\mathbf{f})}{\partial \mathbf{f} \partial \mathbf{f}^\top}$ is the negative Hessian of the log-likelihood, which is an $n \times n$ diagonal matrix[1] with $(\mathbf{D}_f)_{ii} = \sigma(f(x_i))\sigma(-f(x_i))$. The overall Hessian of the log-posterior is negative definite, since $\mathbf{K}^{-1}$ is positive definite and $(\mathbf{D}_f)_{ii} \geq 0$. Thus, there is a unique posterior mode. Note also that $\mathbf{D}_f$ depends on $\mathbf{f} = [f_1, \ldots, f_n]^\top$ but not on the labels $\mathbf{y}$. We can now employ numerical optimisation (gradient ascent or Newton-Raphson method) to find the posterior mode $\hat{\mathbf{f}}^{\text{MAP}}$ and approximate the posterior $p(\mathbf{f}|\mathbf{y})$ with a normal distribution:

$$\tilde{p}(\mathbf{f}|\mathbf{y}) = \mathcal{N}\left(\mathbf{f}\,\middle|\,\hat{\mathbf{f}}^{\text{MAP}}, \left(\mathbf{K}^{-1} + \mathbf{D}_{\hat{\mathbf{f}}^{\text{MAP}}}\right)^{-1}\right).$$

Note that this can be rewritten as

$$\tilde{p}(\mathbf{f}|\mathbf{y}) = \mathcal{N}\left(\mathbf{f}\,\middle|\,\hat{\mathbf{f}}^{\text{MAP}}, \mathbf{K} - \mathbf{K}\left(\mathbf{K} + \mathbf{D}_{\hat{\mathbf{f}}^{\text{MAP}}}^{-1}\right)^{-1}\mathbf{K}\right),$$

---

[1] Note that $\frac{\partial^2 \log p(\mathbf{y}|\mathbf{f})}{\partial \mathbf{f} \partial \mathbf{f}^\top}$ is a diagonal matrix in any GP model regardless of the form of the likelihood function as long as it factorizes across observations, i.e. $p(\mathbf{y}|\mathbf{f}) = \prod_{i=1}^{n} p(y_i|f_i)$ with $(\mathbf{D}_f)_{ii} = -\frac{\partial^2 \log p(y_i|f_i)}{\partial f_i^2}$. In addition, if $\log p(y|f)$ is concave in $f$, $(\mathbf{D}_f)_{ii} \geq 0$.

using the Woodbury identity[2] $\left(\mathbf{K}^{-1} + \mathbf{D}\right)^{-1} = \mathbf{K} - \mathbf{K}\left(\mathbf{K} + \mathbf{D}^{-1}\right)^{-1}\mathbf{K}$ for invertible matrices $\mathbf{K}$ and $\mathbf{D}$.

We can use the Laplace approximation further to construct an approximation of the predictive posterior at a test set $\mathbf{x}' = \{x'_j\}_{j=1}^m$, writing

$$\tilde{p}(\mathbf{f}'|\mathbf{y}) = \int p(\mathbf{f}'|\mathbf{f})\tilde{p}(\mathbf{f}|\mathbf{y})d\mathbf{f}, \tag{7.15}$$

which can now be solved in the closed form since $p(\mathbf{f}'|\mathbf{f})$ is also normal,

$$p\left(\mathbf{f}'|\mathbf{f}\right) = \mathcal{N}\left(\mathbf{f}' \mid \mathbf{K}_{\mathbf{x}'\mathbf{x}}\mathbf{K}_{\mathbf{xx}}^{-1}\mathbf{f}, \mathbf{K}_{\mathbf{x}'\mathbf{x}'} - \mathbf{K}_{\mathbf{x}'\mathbf{x}}\mathbf{K}_{\mathbf{xx}}^{-1}\mathbf{K}_{\mathbf{xx}'}\right),$$

giving

$$\tilde{p}\left(\mathbf{f}'|\mathbf{y}\right) = \mathcal{N}\left(\mathbf{f}' \mid \mathbf{K}_{\mathbf{x}'\mathbf{x}}\mathbf{K}_{\mathbf{xx}}^{-1}\hat{\mathbf{f}}^{\mathrm{MAP}}, \mathbf{K}_{\mathbf{x}'\mathbf{x}'} - \mathbf{K}_{\mathbf{x}'\mathbf{x}}\left(\mathbf{K}_{\mathbf{xx}} + \mathbf{D}_{\hat{\mathbf{f}}^{\mathrm{MAP}}}^{-1}\right)^{-1}\mathbf{K}_{\mathbf{xx}'}\right). \tag{7.16}$$

Let us compare this to the predictive distribution $p\left(\mathbf{f}'|\hat{\mathbf{f}}^{\mathrm{MAP}}\right)$ based on simply plugging in the point estimate $\hat{\mathbf{f}}^{\mathrm{MAP}}$ at the training points, which is

$$p\left(\mathbf{f}'|\hat{\mathbf{f}}^{\mathrm{MAP}}\right) = \mathcal{N}\left(\mathbf{f}' \mid \mathbf{K}_{\mathbf{x}'\mathbf{x}}\mathbf{K}_{\mathbf{xx}}^{-1}\hat{\mathbf{f}}^{\mathrm{MAP}}, \mathbf{K}_{\mathbf{x}'\mathbf{x}'} - \mathbf{K}_{\mathbf{x}'\mathbf{x}}\mathbf{K}_{\mathbf{xx}}^{-1}\mathbf{K}_{\mathbf{xx}'}\right).$$

Two distributions have the same mean but the plug-in predictive underestimates the variance. To see this, note that for any test point $x_\star$, the predictive variances are

$$\mathrm{var}\left[f\left(x_\star\right)|\mathbf{y}\right] = k_{\star\star} - \mathbf{k}_{\star\mathbf{x}}\left(\mathbf{K}_{\mathbf{xx}} + \mathbf{D}_{\hat{\mathbf{f}}^{\mathrm{MAP}}}^{-1}\right)^{-1}\mathbf{k}_{\mathbf{x}\star},$$

$$\mathrm{var}\left[f\left(x_\star\right)|\hat{\mathbf{f}}^{\mathrm{MAP}}\right] = k_{\star\star} - \mathbf{k}_{\star\mathbf{x}}\mathbf{K}_{\mathbf{xx}}^{-1}\mathbf{k}_{\mathbf{x}\star},$$

and positive-definiteness of $\mathbf{D}_{\hat{\mathbf{f}}^{\mathrm{MAP}}}^{-1}$ implies $\left(\mathbf{K}_{\mathbf{xx}} + \mathbf{D}_{\hat{\mathbf{f}}^{\mathrm{MAP}}}^{-1}\right)^{-1} \preceq \mathbf{K}_{\mathbf{xx}}^{-1}$, whereby $\mathrm{var}\left[f\left(x_\star\right)|\mathbf{y}\right] \geq \mathrm{var}\left[f\left(x_\star\right)|\hat{\mathbf{f}}^{\mathrm{MAP}}\right]$.

An alternative to the logistic link is the probit model, i.e.

$$p(y_i = +1|f(x_i)) = \Phi(f(x_i)), \tag{7.17}$$

where $\Phi(z) = \frac{1}{\sqrt{2\pi}}\int_{-\infty}^z e^{-t^2/2}dt$ is the standard normal cdf. Derivations proceed similarly by considering the gradient and Hessian of the log-posterior

$$\log p(\mathbf{f}|\mathbf{y}) = \mathrm{const} - \frac{1}{2}\mathbf{f}^\top\mathbf{K}^{-1}\mathbf{f} + \sum_{i=1}^n \log\Phi(y_i f(x_i)).$$

---

[2]Woodbury matrix identity or matrix inversion lemma in its general form is $(A + UCV)^{-1} = A^{-1} - A^{-1}U\left(C^{-1} + VA^{-1}U\right)^{-1}VA^{-1}$ for matrices $A,U,C,V$ of conformable sizes.

Thus, it suffices to replace

$$
\begin{aligned}
(\mathbf{g_f})_i &= \frac{y_i \phi(f_i)}{\Phi(y_i f_i)}, \\
(\mathbf{D}_f)_{ii} &= \frac{\phi(f_i)^2}{\Phi(y_i f_i)^2} + \frac{y_i f_i \phi(f_i)}{\Phi(y_i f_i)}
\end{aligned}
$$

in (7.13) and (7.14), where $\phi(z) = \Phi'(z)$ is the standard normal pdf. The overall Hessian is again negative definite as a consequence of the log-concavity of $\Phi$.

## 7.5 Numerically stable implementation

Kernel matrix $\mathbf{K}$ can in practice have eigenvalues close to zero and thus be numerically unstable to invert. Fortunately, the direct inversion of $\mathbf{K}$ can be avoided. Consider, for example, the Newton iteration for finding the MAP given by

$$
\begin{aligned}
\mathbf{f}^{new} &= \mathbf{f} - \left( \frac{\partial^2 \log p(\mathbf{f}|\mathbf{y})}{\partial \mathbf{f} \partial \mathbf{f}^\top} \right)^{-1} \frac{\partial \log p(\mathbf{f}|\mathbf{y})}{\partial \mathbf{f}} \\
&= \mathbf{f} + \left( \mathbf{K}^{-1} + \mathbf{D}_f \right)^{-1} \left( \mathbf{g_f} - \mathbf{K}^{-1} \mathbf{f} \right) \\
&= \left( \mathbf{K}^{-1} + \mathbf{D}_f \right)^{-1} \left( \mathbf{K}^{-1} \mathbf{f} + \mathbf{D}_f \mathbf{f} \right) + \left( \mathbf{K}^{-1} + \mathbf{D}_f \right)^{-1} \left( \mathbf{g_f} - \mathbf{K}^{-1} \mathbf{f} \right) \\
&= \left( \mathbf{K}^{-1} + \mathbf{D}_f \right)^{-1} \left( \mathbf{D}_f \mathbf{f} + \mathbf{g_f} \right) \\
&= \left[ \mathbf{K} - \mathbf{K} \left( \mathbf{K} + \mathbf{D_f}^{-1} \right)^{-1} \mathbf{K} \right] \left( \mathbf{D}_f \mathbf{f} + \mathbf{g_f} \right),
\end{aligned}
$$

with the last expression involving only the inverse of $\mathbf{K} + \mathbf{D_f}^{-1}$ and not of $\mathbf{K}$. A recommended implementation [17, Section 3.4.3] is to consider the matrix

$$
\mathbf{B} = \mathbf{D_f}^{1/2} \left( \mathbf{K} + \mathbf{D_f}^{-1} \right) \mathbf{D_f}^{1/2} = \mathbf{D_f}^{1/2} \mathbf{K} \mathbf{D_f}^{1/2} + I
$$

which is guaranteed to be well conditioned for most kernel functions since its eigenvalues are between 1 and $1 + n \max_{i,j} \mathbf{K}_{ij}/4$ and perform its Cholesky decomposition $\mathbf{B} = \mathbf{L} \mathbf{L}^\top$. The Newton update can then be implemented as

$$
\mathbf{f}^{new} = \mathbf{K} \left( I - \mathbf{D_f}^{1/2} \mathbf{L}^{-\top} \mathbf{L}^{-1} \mathbf{D_f}^{1/2} \mathbf{K} \right) \left( \mathbf{D}_f \mathbf{f} + \mathbf{g_f} \right).
$$

## 7.6 Large-Scale Kernel Approximations

Gaussian processes and kernel methods require computational cost that scales at least as $O(n^2)$ and often as $O(n^3)$ in the number of observations $n$ (due to the need to compute, store and invert the $n \times n$ kernel matrix $\mathbf{K}$). This is the price we pay for having a nonparametric model, i.e. for performing the computation in terms of the dual coefficients. For large datasets, e.g. where $n \sim 10^5$, this becomes a prohibitive computational cost and memory requirement, however. Many methods have been proposed to deal with this issue, here we will overview the basic approaches based on the reduced-rank approximation of $\mathbf{K_{xx}}$ - see Chapter 8 of [17] for an in-depth overview.

### 7.6.1 Nyström method

GP regression and kernel ridge regression both require inversion of the matrix $\mathbf{K_{xx}} + \sigma^2 I$. Let us assume for the moment that $\mathbf{K_{xx}}$ can be approximated by a rank $m$ matrix, with $m \ll n$, i.e. $\mathbf{K_{xx}} \approx QQ^\top$, where $Q$ is an $n \times m$ matrix. Then we can apply the matrix inversion lemma and write

$$\left( QQ^\top + \sigma^2 I \right)^{-1} = \sigma^{-2} I - \sigma^{-2} Q \left( \sigma^2 I + Q^\top Q \right)^{-1} Q^\top, \tag{7.18}$$

such that the inversion of an $n \times n$ matrix has been transformed into an inversion of an $m \times m$ matrix. However, in order to derive the optimal reduced-rank approximation to $\mathbf{K_{xx}}$, we need to perform the eigendecomposition of $\mathbf{K_{xx}}$ which is itself a costly operation, requiring $O(n^3)$ computation. Instead, an often used approach is the **Nyström approximation**:

$$\tilde{\mathbf{K}}_{\mathbf{xx}} = \mathbf{K_{xz}} \mathbf{K_{zz}^{-1}} \mathbf{K_{zx}},$$

where $\{z_j\}_{j=1}^m$ is a collection of a small number of inputs in $\mathcal{X}$ (which could be a subset of the training set, but could also be some auxiliary pseudo-inputs) called **inducing points** or **landmark points**, and we denoted as usual $(\mathbf{K_{zz}})_{ij} = k(z_i, z_j)$, $(\mathbf{K_{xz}})_{ij} = k(x_i, z_j)$ and $\mathbf{K_{zx}} = \mathbf{K_{xz}^\top}$. Now, we can set $Q = \mathbf{K_{xz}} \mathbf{K_{zz}^{-1/2}}$ and apply the formula (7.18). Note that this is equivalent to using a finite-dimensional feature map $\phi : x \mapsto \mathbf{K_{zz}^{-1}} [k(z_1, x), \ldots, k(z_m, x)]^\top$ and an *approximate kernel*:

$$\tilde{k}(x, x') = \phi(x)^\top \phi(x').$$

### 7.6.2 Random Fourier Features

Another popular method within the frequentist kernel methods are random Fourier features (RFF) of [16]. The idea behind RFF is to use **Bochner's Theorem**, which gives a representation of **translation-invariant** or **stationary** kernels on $\mathbb{R}^p$. If a real-valued kernel $k(x, x')$ depends only on the difference $x - x'$, then it can be written as

$$
\begin{aligned}
k(x, x') &= \int_{\mathbb{R}^p} \exp\left( i\omega^\top (x - x') \right) d\Lambda(\omega) \\
&= \int_{\mathbb{R}^p} \left\{ \cos\left( \omega^\top x \right) \cos\left( \omega^\top x' \right) + \sin\left( \omega^\top x \right) \sin\left( \omega^\top x' \right) \right\} d\Lambda(\omega) \tag{7.19}
\end{aligned}
$$

for some positive measure (w.l.o.g. a probability distribution) $\Lambda$ called **spectral measure** of $k$. For many widely used kernels, spectral measure takes a simple form, e.g. if $k(x, x') = \exp\left( -\frac{1}{2\gamma^2} \|x - x'\|_2^2 \right)$, then $\Lambda$ is a multivariate normal $\mathcal{N}(0, \gamma^{-2} I)$. Now, for a given $\Lambda$, we sample $m$ frequencies $\{\omega_j\} \sim \Lambda$ and use a Monte Carlo estimator of the kernel function given by the integral in (7.19):

$$
\begin{aligned}
\tilde{k}(x, y) &= \frac{1}{m} \sum_{j=1}^m \left\{ \cos\left( \omega_j^\top x \right) \cos\left( \omega_j^\top y \right) + \sin\left( \omega_j^\top x \right) \sin\left( \omega_j^\top y \right) \right\} \\
&= \langle \phi_\omega(x), \phi_\omega(y) \rangle_{\mathbb{R}^{2m}},
\end{aligned}
$$

which is an approximate kernel corresponding to an explicit set of features $\phi_\omega(x) \in \mathbb{R}^{2m}$ given by

$$x \mapsto \frac{1}{\sqrt{m}} \left[ \cos\left(\omega_1^\top x\right), \sin\left(\omega_1^\top x\right), \ldots, \cos\left(\omega_m^\top x\right), \sin\left(\omega_m^\top x\right) \right]$$

With this set of features, we can now run algorithms in the primal representation which is less costly than paying the computational cost in $n$.

## 7.7 Bayesian Optimization

### 7.7.1 Tuning hyperparameters as optimizing "black-box" functions

We have considered several families of machine learning models which have complex inference algorithms and often require tuning of a number of hyperparameters in order to make them work in practice. These could for example be kernel parameters, the number of layers and units per layer in a deep neural network, learning rates, regularization parameters, or batch sizes in stochastic optimizers. Many important implementation details of hyperparameter tuning are often missing and can be extremely time-consuming as the algorithm needs to be run repeatedly for a large number of hyperparameter configurations – one often uses some form of grid search or random search based on a particular objective function, such as cross-validated empirical risk. Can this be done in a more principled and automated way, i.e. without having "human in the loop"? Ideally, we would want a fully automated machine learning pipeline where (nearly) optimal model configuration is selected with a small number of algorithm runs.

More broadly, we are interested in optimizing a particular 'well behaved' (i.e. it exhibits some degree of smoothness but is nonconvex and possibly multimodal) function $f : \mathcal{X} \to \mathbb{R}$ over some bounded domain $\mathcal{X} \subset \mathbb{R}^d$, i.e. in solving[3]

$$x_\star = \operatorname{argmin}_{x \in \mathcal{X}} f(x).$$

However, $f$ is not known explicitly, i.e. it is a *black-box* function. We can only observe its potentially noisy pointwise evaluations $y_i = f(x_i) + \epsilon_i$ at selected locations $x_i$. Moreover, these pointwise evaluations may be extremely expensive (i.e. they correspond to training of a large machine learning model or even running a complex physical experiment - see [6] for further details).

### 7.7.2 Surrogate Gaussian Process models

The principle of Bayesian optimization is to use a surrogate probabilistic model of the black-box function $f$ in order to carry out the optimization. The default choice for a surrogate model is a Gaussian process (GP) – being a flexible prior over functions (but other models are possible, based on random forests or student-t processes). Further,

---

[3]We will phrase the optimization, w.l.o.g. as minimization here. For maximization of $f$, we can just use minimization of $-f$.

based on the GP model, we need to define a criterion, which we call an *acquisition function*, which determines the next location in $\mathcal{X}$ where $f$ will be evaluated. These locations will in certain sense be *most informative about the optimum of $f$*.

Good acquisition functions will need to balance exploration (learning more about $f$ based on new evaluations) vs exploitation (finding the maximum based on the current model of $f$). Exploration-exploitation tradeoff will be based on our estimates of the uncertainty in the values of $f$ - which is why GP regression models with closed-form uncertainty estimates are most commonly used in this context. Namely, it is customary to assume that the noise $\epsilon_i$ in the evaluations of the black-box function is i.i.d. $\mathcal{N}\left(0, \delta^2\right)$, to bring us to the vanilla GP regression context. We have seen in the previous chapter that the GP model

$$\mathbf{f} \sim \mathcal{N}(0, \mathbf{K})$$
$$\mathbf{y}|\mathbf{f} \sim \mathcal{N}(\mathbf{f}, \delta^2 I),$$

gives us a closed form expression for the posterior predictive mean $\mu\left(x\right)$ and the posterior predictive marginal standard deviation $\sigma\left(x\right) = \sqrt{\kappa\left(x, x\right)}$ at any new location $x$, i.e.

$$f\left(x\right)|\mathcal{D} \sim \mathcal{N}\left(\mu\left(x\right), \kappa\left(x, x\right)\right),$$

where

$$
\begin{aligned}
\mu\left(x\right) &= \mathbf{k}_{x\mathbf{x}}(\mathbf{K} + \delta^2 I)^{-1}\mathbf{y}, \\
\kappa\left(x, x\right) &= k(x, x) - \mathbf{k}_{x\mathbf{x}}(\mathbf{K} + \delta^2 I)^{-1}\mathbf{k}_{\mathbf{x}x}
\end{aligned}
$$

Exploration in this context means that we are seeking locations with high posterior variance $\kappa\left(x, x\right)$, exploitation that we are seeking location with low posterior mean $\mu\left(x\right)$.

### 7.7.3 Acquisition Functions

Most commonly used acquisition functions are: *GP-LCB, probability of improvement (PI)* and *expected improvement (EI)*.

### GP-LCB

GP-LCB follows the principle of the "optimism in the phase of uncertainty" and simply seeks to minimize the lower $(1 - \alpha)$-credible bound of the posterior of the unknown function values $f(x)$, i.e.

$$\alpha_{LCB}\left(x\right) = \mu\left(x\right) - z_{1-\alpha}\sigma\left(x\right),$$

where $z_{1-\alpha} = \Phi^{-1}\left(1 - \alpha\right)$ is the desired quantile of the standard normal distibution. Note that the same principle is used extensively in the theory of K-armed bandits which deals with a similar setup of the optimization of uncertain objectives over a discrete number of choices (not covered in this course), but is typically referred to as UCB (upper confidence/credible bound) since the convention is to consider maximization rather than minimization.

**Probability of improvement (PI)**

If we denote by $\tilde{x}$ the optimal location so far, i.e. the observed $\tilde{y}$ is the minimum among $(y_1, \ldots, y_n)$. Consider $u(x) = \mathbf{1}\{f(x) < \tilde{y}\}$, i.e. $u$ represents the indicator of the event that the function value at any given location is below the observed minimum. Then the probability of improvement is simply

$$\alpha_{PI}(x) = \mathbb{E}\left[u(x)|\mathcal{D}\right] = \int_{-\infty}^{\tilde{y}} \mathcal{N}\left(f; \mu(x), \sigma^2(x)\right) df = \Phi\left(\frac{\tilde{y} - \mu(x)}{\sigma(x)}\right).$$

**Expected Improvement (EI)**

Note that in the probability of improvement, we do not take into account the size of an improvement in the objective function. Thus, the PI method can be viewed as preferring exploitation - locations that have a high probability of being infinitesimally smaller than $\tilde{y}$ will be drawn over points that offer possibly larger improvement but less certainty. We wish to now use as the utility the expected improvement at the location $x$, i.e. we define

$$u(x) = \max(0, \tilde{y} - f(x)),$$

which measures the size of an improvement from the current (estimated) minimum $\tilde{y}$ and we seek to maximize $\mathbb{E}[u(x)|\mathcal{D}]$.

We will need the following simple result about a truncated normal random variable.

**Lemma 20.** *Let $S \sim \mathcal{N}\left(m, \tau^2\right)$. Denote by $\phi$ the density and by $\Phi$ the cdf of a standard normal random variable. Then*

$$\mathbb{E}\left[\max(0, S)\right] = \int_0^\infty s\mathcal{N}\left(s; m, \tau^2\right) ds = \Phi\left(\frac{m}{\tau}\right) m + \phi\left(\frac{m}{\tau}\right) \tau.$$

As a shorthand, we will write $\quad \gamma(x) = \frac{\tilde{y} - \mu(x)}{\sigma(x)}$.

Now,

$$\begin{aligned}
\alpha_{EI}(x) = \mathbb{E}\left[u(x)|\mathcal{D}\right] &= \int_{-\infty}^{\tilde{y}} (\tilde{y} - f)\mathcal{N}\left(f; \mu(x), \sigma^2(x)\right) df \\
&= \int_0^\infty s\mathcal{N}\left(s; \tilde{y} - \mu(x), \sigma^2(x)\right) ds
\end{aligned}$$

after the substitution $s = \tilde{y} - f$. Now, by Lemma 20 we get the most commonly used expression for the expected improvement acquisition function:

$$\begin{aligned}
\alpha_{EI}(x) &= \Phi\left(\gamma(x)\right)\left(\tilde{y} - \mu(x)\right) + \phi\left(\gamma(x)\right)\sigma(x) \\
&= \sigma(x)\left(\gamma(x)\Phi\left(\gamma(x)\right) + \phi\left(\gamma(x)\right)\right).
\end{aligned}$$

While it is the most commonly used, the above expression is ignoring the noise in $f(\tilde{x})$, i.e. it is treating $\tilde{y}$ as the actual value of the objective which can be problematic if the noise level in evaluations is not negligible. Naively, we could simply replace $\tilde{y}$ with the predictive mean $\mu(\tilde{x})$. However, the predictive variance is still ignored.

**Exercise 21.** Derive the alternative expression for both the probability of improvement and for the expected improvement acquisition functions by considering the distribution of $u(x) = \mathbf{1}\{f(x) < f(\tilde{x})\}$. In particular, show that

$$\mathbb{E}\left[u(x)|\mathcal{D}\right] \quad = \quad \Phi\left(\gamma(x)\right)\left(\mu(\tilde{x}) - \mu(x)\right) + \phi\left(\gamma(x)\right)\rho(x,\tilde{x}), \qquad (7.20)$$

where $\rho(x,\tilde{x}) = \sqrt{\kappa(x,x) + \kappa(\tilde{x},\tilde{x}) - 2\kappa(x,\tilde{x})}$.

Note that we can interpret the two terms in 7.20 as trading-off exploration vs. exploitation. We can increase the acquisition either by decreasing $\mu(x)$ (exploitation) or by increasing $\rho(x,\tilde{x})$, i.e. a notion of a distance from the current optimum $\tilde{x}$ (exploration).

# Bibliography

[1] David Arthur and Sergei Vassilvitskii. K-means++: The Advantages of Careful Seeding. In *Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '07, pages 1027–1035, 2007.

[2] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854, 2019.

[3] A. Berlinet and C. Thomas-Agnan. *Reproducing Kernel Hilbert Spaces in Probability and Statistics*. Kluwer, 2004.

[4] David M. Blei. Probabilistic topic models. *Commun. ACM*, 55(4):77–84, April 2012.

[5] Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, New York, NY, USA, 2004.

[6] Eric Brochu, Vlad M Cora, and Nando de Freitas. A tutorial on Bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning. eprint arXiv:1012.2599, arXiv.org, December 2010.

[7] Yuri Burda, Roger Grosse, and Ruslan Salakhutdinov. Importance weighted autoencoders. In *The 4th International Conference on Learning Representations (ICLR)*, 2016.

[8] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, 39(1):1–38, 1977.

[9] Gintare Karolina Dziugaite, Daniel M. Roy, and Zoubin Ghahramani. Training generative neural networks via Maximum Mean Discrepancy optimization. In *Proc. Uncertainty in Artificial Intelligence (UAI)*, 2015.

[10] Arthur Gretton. Lecture notes on Reproducing kernel Hilbert spaces in Machine Learning, University College London, 2017. `http://www.gatsby.ucl.ac.uk/~gretton/coursefiles/rkhscourse.html`.

[11] Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13:723–773, 2012.

[12] Diederik P Kingma and Max Welling. Auto-encoding variational Bayes. In *The 2nd International Conference on Learning Representations (ICLR)*, 2014.

[13] Jon M. Kleinberg. An impossibility theorem for clustering. In *Advances in Neural Information Processing Systems 15*, pages 463–470. MIT Press, 2003.

[14] Brian Kulis and Michael I. Jordan. Revisiting k-means: New Algorithms via Bayesian Nonparametrics. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 513–520, 2012.

[15] Stephen Leslie, Bruce Winney, Garrett Hellenthal, Dan Davison, Abdelhamid Boumertit, Tammy Day, Katarzyna Hutnik, Ellen C Royrvik, Barry Cunliffe, Daniel J Lawson, et al. The fine-scale genetic structure of the british population. *Nature*, 519(7543):309–314, 2015.

[16] A. Rahimi and B. Recht. Random features for large-scale kernel machines. In *Advances in Neural Information Processing Systems*, 2007.

[17] C.E. Rasmussen and C. Williams. *Gaussian Processes for Machine Learning*. MIT Press, 2006.

[18] Danilo Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *International Conference on Machine Learning*, pages 1278–1286, 2014.

[19] Sam Roweis and Zoubin Ghahramani. A unifying review of linear gaussian models. *Neural Comput.*, 11(2):305–345, February 1999.

[20] B. Schölkopf, A. J. Smola, and K.-R. Müller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural Comput.*, 10:1299–1319, 1998.

[21] Shai Shalev-Shwartz and Shai Ben-David. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, New York, NY, USA, 2014.

[22] M. E. Tipping and Christopher Bishop. Probabilistic principal component analysis. *Journal of the Royal Statistical Society, Series B*, 21/3:611–622, January 1999.

[23] C. Zhang, S. Bengio, M. Herdt, B. Recht, and O Vinyals. Understanding deep learning requires rethinking generalization. In *ICLR*, 2017.

# Index

*Index*

Variational free energy, 52

weak duality, 20
within-cluster deviance, 9
witness function , 40