

NOTES

IMPERIAL COLLEGE LONDON

DEPARTMENT OF COMPUTING

Mathematics For Machine Learning

Author:

W (CID: your college-id number)

Date: November 30, 2019

1 Linear Algebra

1.1 Vector Space

Given a set \mathcal{G} and an operation \otimes on \mathcal{G} , (\mathcal{G}, \otimes) is called a **group** if the following holds:

- $\forall x, y \in \mathcal{G} : x \otimes y \in \mathcal{G}$
- $\forall x, y, z \in \mathcal{G} : (x \otimes y) \otimes z = x \otimes (y \otimes z)$
- $\exists e \in \mathcal{G} \forall x \in \mathcal{G} : x \otimes e = e \otimes x = x$
- $\forall x \in \mathcal{G} \exists y \in \mathcal{G} : x \otimes y = y \otimes x = e$

The group is an **Abelian group** if $\forall x, y \in \mathcal{G} : x \otimes y = y \otimes x$.

A real value vector space is a set \mathcal{V} with two operations:

$$+ : \mathcal{V} \times \mathcal{V} \rightarrow \mathcal{V} \qquad \cdot : \mathbb{R} \times \mathcal{V} \rightarrow \mathcal{V}$$

Where $+$ is vector addition and \cdot is multiplication by scalar. $(\mathcal{V}, +)$ is an Abelian group with neutral element $e = \mathbf{0}$.

1.2 Basis and Rank

Consider a vector space $(\mathcal{V}, +, \cdot)$ and a set of vectors \mathcal{A} . \mathcal{A} is the generating set of \mathcal{V} if every $v \in \mathcal{V}$ can be expressed as a linear combination of vectors in \mathcal{A} . The **basis** of \mathcal{V} is the minimal generating set.

The **rank** of a matrix $A \in \mathbb{R}^{m \times n}$, denoted as $rk(A)$, equals the number of linearly independent columns or rows in matrix A . A matrix has **full rank** if $rk(A) = \min(m, n)$.

2 Analytic Geometry

2.1 Orthonormal Basis

Given the basis $\{\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_n\}$ of a vector space \mathcal{V} , the basis is called **orthogonal basis** if every basis vector is orthogonal to each other, i.e. their dot product is 0

$$\langle \mathbf{b}_i, \mathbf{b}_j \rangle = 0, \quad i \neq j$$

Additionally, the basis is called **orthonormal basis** if it is an orthogonal basis and each basis vector has length/norm 1.

$$\langle \mathbf{b}_i, \mathbf{b}_i \rangle = 1$$

A matrix A is an **orthogonal matrix** if its rows/columns form an orthonormal basis. A property of orthogonal matrix is that

$$AA^\top = A^\top A = I$$

It follows that orthogonal matrices are always invertible and $A^{-1} = A^\top$

3 Matrix Decomposition

3.1 Eigenvalues and Eigenvectors

Eigenvalues and Eigenvectors are only defined for square matrices. Given a square matrix $A \in \mathbb{R}^{n \times n}$, then $\lambda \in \mathbb{R}$ is an eigenvalue of A and $\mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$ is an eigenvector of A if $A\mathbf{x} = \lambda\mathbf{x}$. This means the transformation of A does not change the direction of eigenvectors \mathbf{x} but only their magnitude.

We can find the eigenvalues of a matrix A by finding the roots of the characteristic polynomial $p_A(\lambda)$ of A

$$p_A(\lambda) \doteq \det(A - \lambda I)$$

Once we find the eigenvalues, we can obtain the eigenvectors by solving the set of linear equations. If \mathbf{x} is an eigenvector of A associated with eigenvalue λ , then for any $c \in \mathbb{R} \setminus \{0\}$ $c\mathbf{x}$ is also an eigenvector of A

$$A(c\mathbf{x}) = cA\mathbf{x} = c\lambda\mathbf{x} = \lambda(c\mathbf{x})$$

3.1.1 Spectral Theorem

Given a **symmetric** matrix $A \in \mathbb{R}^{n \times n}$, then the eigenvalues of A are real and the eigenvectors of A forms an orthonormal basis

3.1.2 Determinants and Traces

The determinant of a matrix A is the product of its eigenvalues

$$\det(A) = \prod_{n=1}^N \lambda_n$$

The trace of a matrix A is the sum of its eigenvalues

$$\text{tr}(A) = \sum_{n=1}^N \lambda_n$$

3.2 Eigendecomposition and Diagonalisation

A diagonal matrix D is a matrix with all values 0 except for the main diagonal. Diagonal matrices are useful because they allow fast computations for certain operations:

- The determinant $\det(D)$ is the product of its diagonal entries
- The matrix power D^k is each diagonal element raised to power k
- The inverse D^{-1} is the reciprocal of each diagonal element

A squared matrix $A \in \mathbb{R}^{n \times n}$ is **diagonalisable** if there is an invertible matrix $P \in \mathbb{R}^{n \times n}$ such that $D = P^{-1}AP$. It turns out that the columns of P are the eigenvectors of A .

We define D as the diagonal matrix where the diagonal values are the eigenvalues of A , and we define $P = [p_1, \dots, p_n]$ where each p_i is an eigenvector of A . Then

$$\begin{aligned} AP &= A[p_1 \ \dots \ p_n] = [Ap_1 \ \dots \ Ap_n] \\ PD &= [p_1 \ \dots \ p_n] \begin{bmatrix} \lambda_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_n \end{bmatrix} = [\lambda_1 p_1 \ \dots \ \lambda_n p_n] \\ \therefore Ap_1 &= \lambda_1 p_1 \\ \therefore \vdots &= \vdots \\ \therefore Ap_n &= \lambda_n p_n \\ \therefore AP &= PD \end{aligned}$$

Eigendecomposition is factoring a square matrix $A \in \mathbb{R}^{n \times n}$ into

$$A = PDP^{-1}$$

In the case where A is symmetric, from the spectral theorem we know the eigenvectors of A forms an orthonormal basis, meaning $P^{-1} = P^\top$, hence

$$A = PDP^\top$$

3.3 Singular Value Decomposition (SVD)

Singular Value Decomposition can be applied to all matrices, not just squared matrices, and it always exists. Given a rectangular matrix A , SVD decomposes A into

$$A_{m \times n} = U_{m \times m} \Sigma_{m \times n} V_{n \times n}^\top$$

Where U and V are orthogonal matrices and Σ is a diagonal matrix.

Σ is the singular matrix and the diagonal entries are called the singular values. The singular matrix has the same shape as A hence needs additional zero padding.

$$\Sigma = \begin{bmatrix} \sigma_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_n \\ 0 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} \sigma_1 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_n & 0 & \dots & 0 \end{bmatrix}$$

From Eigendecomposition, we saw that a symmetric matrix S can be decomposed into $S = PDP^\top$. Also, we can construct a symmetric, positive semidefinite matrix

from a rectangular matrix A by AA^\top or $A^\top A$.

Therefore, constructing the SVD consists of find the eigenvectors of AA^\top and $A^\top A$. The eigenvectors of AA^\top make up the columns of U , and the eigenvectors of $A^\top A$ make up the columns of V^\top . The singular values in Σ are the square roots of eigenvalues from AA^\top or $A^\top A$.

$$\begin{aligned} A^\top A &= (U\Sigma V^\top)^\top (U\Sigma V^\top) = V\Sigma^\top U^\top U\Sigma V^\top = V\Sigma^\top \Sigma V^\top \\ AA^\top &= (U\Sigma V^\top)(U\Sigma V^\top)^\top = U\Sigma V^\top V\Sigma^\top U^\top = U\Sigma\Sigma^\top U^\top \end{aligned}$$

4 Vector Calculus

Many algorithms in machine learning optimise an objective function with respect to a set of desired model parameters, and so calculus plays an important role in machine learning.

A function f relates two quantities to each other. In machine learning, the input is a vector $\mathbf{x} \in \mathbb{R}^D$ with D features. Here, \mathbb{R}^D is the **domain** of f and the values $f(\mathbf{x})$ are **image** of f .

4.1 Differentiate of Univariate Functions

Given a variable $x \in \mathbb{R}$ and a function $f : \mathbb{R} \rightarrow \mathbb{R}$. The **derivative** for f with respect to x , for $h > 0$

$$\frac{df}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

4.2 Partial Differentiations

In the following, we consider the general case of f with input $\mathbf{x} \in \mathbb{R}^D$.

The gradient of f with respect to \mathbf{x} is a collection of partial derivatives. Each partial derivative is obtained differentiating one variable at a time and keeping the others constant.

The gradient of f is a row vector $\mathbb{R}^{1 \times D}$ is called the **Jacobian**.

$$\frac{df}{d\mathbf{x}} = \left[\frac{\partial f(\mathbf{x})}{\partial x_1} \quad \frac{\partial f(\mathbf{x})}{\partial x_2} \quad \dots \quad \frac{\partial f(\mathbf{x})}{\partial x_D} \right] \in \mathbb{R}^{1 \times D}$$

4.3 Vector-Valued Functions

In the following, we consider the general case where $f : \mathbb{R}^D \rightarrow \mathbb{R}^E$ with input $\mathbf{x} \in \mathbb{R}^D$.

$$f(\mathbf{x}) = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_E(\mathbf{x}) \end{bmatrix} \in \mathbb{R}^E$$

The derivative of f with respect to \mathbf{x} is a $E \times D$ matrix called the **Jacobian**.

$$\frac{df(\mathbf{x})}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \cdots & \frac{\partial f_1(\mathbf{x})}{\partial x_D} \\ \vdots & & \vdots \\ \frac{\partial f_E(\mathbf{x})}{\partial x_1} & \cdots & \frac{\partial f_E(\mathbf{x})}{\partial x_D} \end{bmatrix} \in \mathbb{R}^{E \times D}$$

Example

Find the derivatives $df/d\mathbf{x}$ and $df/d\mathbf{A}$.

$$\mathbf{f}(\mathbf{x}) = \mathbf{A}\mathbf{x} \quad \mathbf{f} \in \mathbb{R}^M \quad \mathbf{A} \in \mathbb{R}^{M \times N} \quad \mathbf{x} \in \mathbb{R}^N$$

Writing down each function $f_i(\mathbf{x})$ makes it easier to find the gradients.

$$f_i(\mathbf{x}) = \sum_{j=1}^N A_{ij}x_j$$

Since function f has dimensions M and the target \mathbf{x} has dimensions N . The derivative $df/d\mathbf{x}$ has the shape $M \times N$. Because

$$\frac{\partial f_i}{\partial x_j} = A_{ij} \Rightarrow \frac{df}{d\mathbf{x}} = \mathbf{A}$$

Therefore the derivative $df/d\mathbf{x} = \mathbf{A}$.

Since function f has dimensions M and the target \mathbf{A} has dimensions $M \times N$. The derivative $df/d\mathbf{A}$ has the shape $M \times M \times N$.

$$\frac{\partial f_i}{\partial A_{ij}} = x_j \Rightarrow \frac{\partial f_i}{\partial A_k} = \begin{cases} \mathbf{x}^T, & k = i \\ \mathbf{0}^T, & k \neq i \end{cases}$$

Therefore the derivative $df/d\mathbf{A}$ is

$$\frac{df}{d\mathbf{A}} = \begin{bmatrix} \frac{\partial f_1}{\partial \mathbf{A}} \\ \vdots \\ \frac{\partial f_M}{\partial \mathbf{A}} \end{bmatrix}$$

Each of the partial derivatives $\frac{\partial f_i}{\partial \mathbf{A}}$ is a $M \times N$ matrix where the k^{th} is \mathbf{x}^T if $k = i$, and $\mathbf{0}^T$ otherwise.

4.4 Higher Order Derivatives

Some methods of optimisation requires second derivatives. Given a function $f : \mathbb{R}^D \rightarrow \mathbb{R}$, the second derivative is represented by a $D \times D$ matrix called the **Hessian**. The Hessian matrix is a collection of second partial derivatives.

The second partial derivative means f is first differentiated with respect to x then with respect to y . Notice the order of differentiation does not matter.

$$\frac{\partial^2 f}{\partial x \partial y} = \frac{\partial f}{\partial x} \frac{\partial f}{\partial y} = \frac{\partial^2 f}{\partial y \partial x}$$

The point is a minimum point if the Hessian matrix is **positive definite**.

Given a function $f(x, y)$ that is twice differentiable. The Hessian matrix is given by

$$\mathbf{H} = \begin{bmatrix} \frac{\partial^2 f}{\partial x^2} & \frac{\partial^2 f}{\partial x \partial y} \\ \frac{\partial^2 f}{\partial y \partial x} & \frac{\partial^2 f}{\partial y^2} \end{bmatrix}$$

5 Probability and Distributions

In machine learning and statistics, there are two major interpretations of probabilities: the Bayesian and Frequentist interpretations. The Bayesian interpretation uses probability to specify the degree of uncertainty a user has about an event. The Frequentist interpretation considers the relative frequencies of events of interests to the total number of events.

5.1 Random Variables

The **sample space** Ω is the set of all possible outcomes of the experiment. The **event space** \mathcal{A} is obtained by considering the collection of subsets of Ω . Each event $A \in \mathcal{A}$ is associated with a **probability** $P(A)$ that measures the probability the event will occur.

In machine learning, we often avoid explicitly referring to the probability space, but instead refer to probabilities on quantity of interest. This quantity of interest is in the **target space** \mathcal{T} . A **random variable** $X : \Omega \rightarrow \mathcal{T}$ is a function that maps elements of Ω to a quantity of interest $x \in \mathcal{T}$.

It is possible to equate the probability of output of X and the probability of the samples in Ω . Given a set $\mathcal{S} \subseteq \mathcal{T}$, the probability of \mathcal{S}

$$P_X(\mathcal{S}) = P(X^{-1}(\mathcal{S})) = P(\{\omega \in \Omega : X(\omega) \in \mathcal{S}\})$$

Example

In the case of tossing two coins, the sample space $\Omega = \{HH, HT, TH, TT\}$. The quantity of interest is the number of heads in each outcome. The random variable X maps each element in the sample space to the target space $\mathcal{T} = \{0, 1, 2\}$.

$$X(HH) = 2 \quad X(HT) = 1 \quad X(TH) = 1 \quad X(TT) = 0$$

Say we want to find the probability of obtaining one head, $\mathcal{S} = \{1\}$. Then

$$P_X(\mathcal{S}) = P(X^{-1}(\mathcal{S})) = P(\{HT, TH\})$$

5.2 Discrete and Continuous Probabilities

Depending on whether the target space \mathcal{T} is discrete or continuous, the natural way to refer to distributions is very different.

Continuous probabilities uses a **Probability Density Function (pdf)** $f : \mathbb{R}^D \rightarrow \mathbb{R}$ that must satisfy two conditions:

1. $\forall \mathbf{x} \in \mathbb{R}^D : f(\mathbf{x}) \geq 0$
2. $\int f(\mathbf{x}) d\mathbf{x} = 1$

Discrete probabilities uses **Probability Mass Function (pmf)** that is the same as pdf but the integral is replace with sum in the second condition.

Continuous probabilities also uses a **Cumulative Distribution Function (cdf)** defined as

$$F_X(\mathbf{x}) = P(X_1 \leq x_1, \dots, X_D \leq x_D)$$

5.3 Sum Rule, Product Rule, and Bayes Theorem

The **joint probability** $p(\mathbf{x}, \mathbf{y})$ is the probability of intersection of both events.

The **sum rule** is used to find corresponding **marginal distributions** $p(\mathbf{x})$ and $p(\mathbf{y})$.

$$p(\mathbf{x}) = \int_{\mathcal{Y}} p(\mathbf{x}, \mathbf{y}) d\mathbf{y} \quad p(\mathbf{x}) = \sum_{\mathbf{y} \in \mathcal{Y}} p(\mathbf{x}, \mathbf{y})$$

The **product rule** relates the joint distribution to conditional distribution

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

Bayes' Theorem is used to draw conclusions about \mathbf{x} given observed values in \mathbf{y} .

$$p(\mathbf{x} | \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{x})p(\mathbf{x})}{p(\mathbf{y})}$$

The prior $p(\mathbf{x})$ encapsulates subjective prior knowledge \mathbf{x} . The likelihood $p(\mathbf{y} | \mathbf{x})$ describes how \mathbf{x} and \mathbf{y} are related. The evidence $p(\mathbf{y})$ can be seen as a normalising constant.

5.4 Summary Statistics

Mean and (co)variance are often useful to describe properties of probability distributions.

The **Expected Value** of a function $g : \mathbb{R} \rightarrow \mathbb{R}$ of a univariate random variable X

$$\mathbb{E}_X[g(x)] = \int_{\mathcal{X}} g(x)p(x)dx \quad \mathbb{E}_X[g(x)] = \sum_{x \in \mathcal{X}} g(x)p(x)$$

For multivariate random variables, the expected value is defined element wise

$$\mathbb{E}_X[g(\mathbf{x})] = \begin{bmatrix} \mathbb{E}_{X_1}[g(x_1)] \\ \dots \\ \mathbb{E}_{X_D}[g(x_D)] \end{bmatrix}$$

The **Mean** is a special case of expected value where the function g is an identity function.

The **Covariance** between two univariate random variable X, Y is the expected product of their deviations from their respective mean

$$\begin{aligned} \text{Cov}_{X,Y}[x, y] &= \mathbb{E}_{X,Y}[(x - \mathbb{E}_X[x])(y - \mathbb{E}_Y[y])] \\ &= \mathbb{E}[xy] - \mathbb{E}[x]\mathbb{E}[y] \end{aligned}$$

For two multivariate random variables X, Y , the covariance is defined as

$$\text{Cov}[\mathbf{x}, \mathbf{y}] = \mathbb{E}[\mathbf{x}\mathbf{y}^T] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{y}]^T = \text{Cov}[\mathbf{y}, \mathbf{x}]^T$$

The **Variance** of a random variable X is the covariance of the variable with itself

$$\begin{aligned} \mathbb{V}_X[\mathbf{x}] &= \text{Cov}_X[\mathbf{x}, \mathbf{x}] \\ &= \mathbb{E}_X[\mathbf{x}\mathbf{x}^T] - \mathbb{E}_X[\mathbf{x}]\mathbb{E}_X[\mathbf{x}]^T \end{aligned}$$

The covariance matrix in this case is symmetric and positive semidefinite.

5.4.1 Empirical Mean and Covariance

The definitions above are **population mean** and **population covariance**. In machine learning, we usually have to learn from empirical observations of data. Given $\mathbf{x} \in \mathbb{R}^D$:

The **empirical mean** vector is the arithmetic average of the observations for each variable.

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

The **empirical covariance** is a $D \times D$ matrix

$$\Sigma = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T$$

If we focus on a single random variable X , the **variance** can be defined as

$$\mathbb{V}_X[x] = \mathbb{E}_X[(x - \mu)^2], \text{ where } \mu = \mathbb{E}_X[x]$$

5.5 Statistical Independence

Two random variables X and Y are **statistically independent** if and only if

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y})$$

X and Y are **conditionally independent** given Z if and only if

$$\forall \mathbf{z} \in Z : p(\mathbf{x}, \mathbf{y} | \mathbf{z}) = p(\mathbf{x} | \mathbf{z})p(\mathbf{y} | \mathbf{z})$$

From the product rule

$$p(\mathbf{x} | \mathbf{y}) = \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{y})} \Rightarrow p(\mathbf{x} | \mathbf{y}, \mathbf{z}) = \frac{p(\mathbf{x}, \mathbf{y} | \mathbf{z})}{p(\mathbf{y} | \mathbf{z})}$$

If X and Y are conditionally independent given Z , then

$$\begin{aligned} p(\mathbf{x} | \mathbf{y}, \mathbf{z}) &= \frac{p(\mathbf{x} | \mathbf{z})p(\mathbf{y} | \mathbf{z})}{p(\mathbf{y} | \mathbf{z})} \\ &= p(\mathbf{x} | \mathbf{z}) \end{aligned}$$

5.6 Probability Distributions

5.6.1 Bernoulli Distribution

The Bernoulli distribution is a distribution for a single binary random variable X . It has a single parameter $\mu \in [0, 1]$ that represents the probability $P(X = 1)$. The Bernoulli distribution $Ber(\mu)$ is defined as

$$p(x | \mu) = \mu^x(1 - \mu)^{1-x}, \quad x \in \{0, 1\}$$

The mean of Bernoulli $\mathbb{E}[x] = \mu$ and the variance $\mathbb{V}[x] = \mu(1 - \mu)$

5.6.2 Binomial Distribution

The Binomial distribution can be used to describe the probability of observing m occurrences of $X = 1$ in a set of N samples from a Bernoulli distribution $Ber(\mu)$. The Binomial distribution $Bin(N, \mu)$ is defined as

$$p(m | N, \mu) = \binom{N}{m} \mu^m (1 - \mu)^{N-m}$$

The mean of Binomial $\mathbb{E}[m] = N\mu$ and the variance $\mathbb{V}[m] = N\mu(1 - \mu)$

5.6.3 Beta Distribution

The Beta distribution represents a probability distribution of probabilities, it is often used to represent the probability of μ in Bernoulli distribution. The Beta distribution $Beta(\alpha, \beta)$ is defined as

$$p(\mu | \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \mu^{\alpha-1} (1 - \mu)^{\beta-1}$$

The two parameters $\alpha > 0$, $\beta > 0$ governs the distribution. Intuitively, α moves the probability mass towards 1 and β moves the probability mass towards 0.

An interesting case is when $\alpha = \beta = 1$, in this case we obtain Uniform distribution.

5.7 Gaussian Distribution

The Gaussian Distribution, also known as normal distribution, is extensively used in machine learning.

Given a multivariate data $\mathbf{x} \in \mathbb{R}^D$. The **multivariate Gaussian distribution** is defined by two parameters: the mean vector $\boldsymbol{\mu}$ and a covariance matrix Σ .

$$p(\mathbf{x} | \boldsymbol{\mu}, \Sigma) = (2\pi)^{-\frac{D}{2}} |\Sigma|^{-\frac{1}{2}} e^{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})}$$

The Gaussian with $\boldsymbol{\mu} = \mathbf{0}$ and $\Sigma = \mathbf{I}$ is called the **standard normal distribution**.

5.7.1 Marginal and Conditional of Gaussians

Let X and Y be multivariate random variables. Their joint probability is a Gaussian distribution defined as

$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N}\left(\begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}\right)$$

Where $\Sigma_{xx} = \text{Cov}[\mathbf{x}, \mathbf{x}]$ and $\Sigma_{yy} = \text{Cov}[\mathbf{y}, \mathbf{y}]$. Σ_{xy} and Σ_{yx} are transpose of each other.

The **marginal distribution** $p(\mathbf{x})$ is a Gaussian distribution

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_x, \Sigma_{xx})$$

The **conditional probability** $p(\mathbf{x} | \mathbf{y})$ is a Gaussian distribution

$$p(\mathbf{x} | \mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_x + \Sigma_{xy} \Sigma_{yy}^{-1} (\mathbf{y} - \boldsymbol{\mu}_y), \Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx})$$

Notice that \mathbf{y} here is not a random variable but a value.

5.7.2 Product of Gaussian

The product of two Gaussians $\mathcal{N}(\mathbf{x} | \mathbf{a}, \mathbf{A})\mathcal{N}(\mathbf{x} | \mathbf{b}, \mathbf{B})$ is a Gaussian $z\mathcal{N}(\mathbf{x} | \mathbf{c}, \mathbf{C})$. Where

$$\begin{aligned}\mathbf{C} &= (\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1} \\ \mathbf{c} &= \mathbf{C}(\mathbf{A}^{-1}\mathbf{a} + \mathbf{B}^{-1}\mathbf{b}) \\ z &= (2\pi)^{-\frac{D}{2}} |\mathbf{A} + \mathbf{B}|^{-\frac{1}{2}} e^{-\frac{1}{2}(\mathbf{a}-\mathbf{b})^T(\mathbf{A}+\mathbf{B})^{-1}(\mathbf{a}-\mathbf{b})} \\ &= \mathcal{N}(\mathbf{a} | \mathbf{b}, \mathbf{A} + \mathbf{B}) = \mathcal{N}(\mathbf{b} | \mathbf{a}, \mathbf{A} + \mathbf{B})\end{aligned}$$

The scaling constant z can be written in the functional form of a Gaussian density. Meaning in $\mathcal{N}(\mathbf{a} | \mathbf{b}, \mathbf{A} + \mathbf{B})$ and $\mathcal{N}(\mathbf{b} | \mathbf{a}, \mathbf{A} + \mathbf{B})$, \mathbf{a} and \mathbf{b} are not random variables but actual values.

5.7.3 Linear Transformation

Consider a Gaussian distributed random sample $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Let \mathbf{Y} be a random variable that is a transformed version of \mathbf{x} by applying matrix \mathbf{A} and vector \mathbf{b} , $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b}$.

$$\begin{aligned}\mathbb{E}[\mathbf{y}] &= \mathbb{E}[\mathbf{A}\mathbf{x} + \mathbf{b}] = \mathbf{A}\mathbb{E}[\mathbf{x}] + \mathbf{b} = \mathbf{A}\boldsymbol{\mu} + \mathbf{b} \\ \mathbb{V}[\mathbf{y}] &= \mathbb{V}[\mathbf{A}\mathbf{x} + \mathbf{b}] = \mathbf{A}\mathbb{V}[\mathbf{x}]\mathbf{A}^T = \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T\end{aligned}$$

Therefore \mathbf{Y} is a Gaussian distribution defined as

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | \mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$$

Consider the reverse transformation. Given $\mathbf{y} = \mathbf{A}\mathbf{x}$ and we know the probability distribution of $p(\mathbf{y})$. We want to find the probability distribution of $p(\mathbf{x})$.

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | \mathbf{A}\mathbf{x}, \boldsymbol{\Sigma})$$

If \mathbf{A} is **invertible**. Then we can write $\mathbf{x} = \mathbf{A}^{-1}\mathbf{y}$ and apply the transformation as shown above.

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \mathbf{A}^{-1}\mathbf{A}\mathbf{x}, \mathbf{A}^{-1}\boldsymbol{\Sigma}(\mathbf{A}^{-1})^T)$$

If \mathbf{A} is **not invertible**. We multiply both sides with \mathbf{A}^T then invert $\mathbf{A}^T\mathbf{A}$. Then apply the transformation as shown above.

$$\begin{aligned}(\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{y} &= \mathbf{x} \\ p(\mathbf{x}) &= \mathcal{N}(\mathbf{x} | (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{A}\mathbf{x}, (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\boldsymbol{\Sigma}\mathbf{A}(\mathbf{A}^T\mathbf{A})^{-1})\end{aligned}$$

5.8 Conjugacy

From Bayes' Theorem, the posterior is proportional to the product of the prior and likelihood. A prior is **conjugate** for the likelihood function if the posterior is the same as the prior.

Conjugacy is particularly convenient because we can algebraically calculate our posterior distribution by updating the parameters of the prior distribution.

6 Continuous Optimisation

Training a machine learning model means finding a good set of parameters that minimise a predefined objective function. In this chapter we will look at optimisation using gradient descent.

6.1 Gradient Descent

Gradient descent is an optimisation algorithm. It finds the minimum of a function by iteratively taking a step proportional to the negative of the gradient of the function at the current point.

Given a multivariate function $f(\mathbf{x})$ starting at location \mathbf{x}_0 . In each step we update the location by

$$\mathbf{x}_{i+1} \leftarrow \mathbf{x}_i - \alpha_i ((\nabla f)(\mathbf{x}_i))^{\top}, \alpha_i \in (0, 1]$$

where α_i is the learning rate, or step-size.

6.1.1 Learning Rate

In standard gradient descent, the step-size α is constant throughout training. Choosing an appropriate step-size is important,

- A step-size that is too small takes a long time to converge
- A step-size that is too large may lead to overshooting or even diverge.

Ideally we want to have large step-size initially to make big steps towards the goal, but a small step-size for fine tuning when it is near the optimum. We can use **adaptive gradient methods** that rescales the step-size at each iteration.

6.1.2 Momentum

Gradient descent struggles in pathological curvatures (often described as valleys or trenches). Since the surface at the ridge curves is much steeper, the parameters will zigzag back and forth across the valley, making convergence very slow.

Momentum takes advantage of the knowledge from previous gradient to determine the direction to go.

$$\mathbf{x}_{i+1} \leftarrow \mathbf{x}_i - \alpha_i ((\nabla f)(\mathbf{x}_i))^{\top} - \beta \alpha_{i-1} ((\nabla f)(\mathbf{x}_{i-1}))^{\top}, \beta \in (0, 1]$$

6.1.3 Stochastic Gradient Descent (SGD)

Computing the gradient can be time consuming, stochastic gradient descent approximates the gradient of an objective function that is written as a sum of differentiable functions.

Given a dataset with N data, the objective function $L(\theta)$ usually sum of losses $L_n(\theta)$ of each individual data.

$$L(\theta) = \sum_{n=1}^N L_n(\theta)$$

In standard gradient descent, we calculate the gradient using the entire dataset in each iteration. However when the training set is enormous, evaluating the sum of gradients become expensive.

$$\theta_{i+1} \leftarrow \theta_i - \alpha_i \sum_{n=1}^N (\nabla L_n(\theta_i))^\top$$

Instead of using all L_n from $n = 1 \dots N$, we can randomly select a subset for batch gradient descent. In the extreme case, we can select a single L_n to estimate the gradient.

Taking only a subset of a data to approximate the gradient still works because for gradient descent to converge, we only require the gradient to be an unbiased estimate of the true gradient.

7 When Models Meet Data

7.1 Directed Graph Models

Directed Graph Models are used for representing conditional dependencies in a probabilistic model. In a graphical model, **nodes** are random variables. **Directed link** from a to b indicates conditional probability $p(b | a)$.

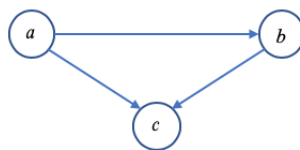
7.1.1 Constructing a Graph

Given a joint distribution, we can construct a graph in two steps:

1. Create a node for all random variables
2. For each conditional probability, add a directed link between the nodes

Example

The joint distribution $p(a, b, c) = p(c | a, b)p(b | a)p(a)$ is represented by the graph



7.1.2 Extracting Joint Probability

Given a graph, we can extract the joint probability in two steps

1. The joint distribution $p(x_1, x_2, \dots, x_K)$ should include all random variable appeared in the graph
2. Each conditional depends only on the parents of a node.

In general, the joint distribution is given by

$$p(x_1, x_2, \dots, x_K) = \prod_{k=1}^K p(x_k \mid \text{Parents}(x_k))$$

7.1.3 Conditional Independence

Recall the definition of conditional independence.

$$\begin{aligned} x \perp\!\!\!\perp y \mid z &\doteq p(x, y \mid z) = p(x \mid z)p(y \mid z) \\ &\doteq p(x \mid y, z) = p(x \mid z) \end{aligned}$$

Directed graph models allow us to find conditional independence using a concept called **d-separation**. Consider a directed graph with \mathcal{A} , \mathcal{B} , and \mathcal{C} which are nonintersecting sets of nodes.

We consider all paths (ignoring arrow direction) from any node in \mathcal{A} to any node in \mathcal{B} . The path is blocked if

- Links meet **head to tail** or **tail to tail** at a node in \mathcal{C} , or
- Links meet **head to head** at a node. The node and all its descendants are not in \mathcal{C} .

If all paths are blocked, then we say \mathcal{A} is d-separated from \mathcal{B} by \mathcal{C} , and $\mathcal{A} \perp\!\!\!\perp \mathcal{B} \mid \mathcal{C}$.

7.2 Model Selection

In machine learning, the choice of model influence the number and type of free parameters in the model and thereby also the flexibility and expressivity of the model. In this section, we look at mechanisms for assessing how good a model generalises to unseen test data.

7.2.1 Cross Validation

K -fold cross validation partitions the entire data set into K chunks. In each iteration k , it uses $K - 1$ chunks as the training set $\mathcal{R}^{(k)}$ to produce a predictor $f^{(k)}$, which is then applied to the validation set $\mathcal{V}^{(k)}$ to compute the loss $L(f^{(k)}, \mathcal{V}^{(k)})$. The process is

repeated K times using a different validation set in each iteration. The performance of the model is then averaged from all iterations:

$$\mathbb{E}_{\mathcal{V}}[L(f, \mathcal{V})] = \frac{1}{K} \sum_{k=1}^K L(f^{(k)}, \mathcal{V}^{(k)})$$

We repeat this procedure for all models and choose the model that performs best. Once the model is chosen, we can evaluate the final performance on the test set.

7.2.2 Bayesian Model Selection

Simpler models are less expressive but are also less prone to overfitting. Bayesian model selection aims to find the simplest model that explains the data reasonably well.

Given a finite number of models $M = \{M_1, M_2, \dots, M_K\}$ where each model M_k possesses parameters θ_k . Bayesian model selection places a prior on the set of model to create a generative process to generate data \mathcal{D} .

$$M_k \sim p(M) \quad \theta_k \sim p(\theta | M_k) \quad \mathcal{D} \sim p(\mathcal{D} | \theta_k)$$

We can compute the posterior distribution of the model using Bayes' Theorem. Where $p(\mathcal{D} | M_k)$ is the marginal likelihood. The model no longer depends on the parameters because they are integrated out.

$$p(M_k | \mathcal{D}) \propto p(\mathcal{D} | M_k) p(M_k)$$

$$p(\mathcal{D} | M_k) = \int p(\mathcal{D} | \theta_k) p(\theta_k | M_k) d\theta_k$$

The best model M^* is the one that maximises the posterior distribution. If the prior $p(M)$ is uniformly distributed, then maximising the posterior is same as maximising the marginal likelihood.

$$M^* = \operatorname{argmax}_{M_k} p(M_k | \mathcal{D})$$

8 Linear Regression

Regression task aims to find a function f that maps inputs $\mathbf{x} \in \mathbb{R}^D$ to corresponding noisy observed values $y \in \mathbb{R}$. The relationship between \mathbf{x} and y is given as

$$y = f(\mathbf{x}) + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2)$$

In linear regression, we only consider the special case where the parameters $\theta \in \mathbb{R}^D$ appears linearly in the model. The aim is to seek for parameters θ .

$$y = \mathbf{x}^\top \theta + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2)$$

This chapter will discuss maximum likelihood and maximum a posteriori (MAP) estimation to find optimal model parameters.

8.1 Maximum Likelihood Estimation

A widely used approach to finding the desired parameters θ_{ML} is to find one that maximise the likelihood, which in this case is the joint probability.

$$\theta_{ML} = \underset{\theta}{\operatorname{argmax}} p(\mathcal{Y} | \mathcal{X}, \theta)$$

Where $\mathcal{Y} = \{y_1, y_2, \dots, y_N\}$ and $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$.

Given a set of training data $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$ where all data points are independent and identically distributed, the joint probability is the product of all individual probability. Each individual probability is a Gaussian distribution because the noise ϵ is Gaussian distributed.

$$p(\mathcal{Y} | \mathcal{X}, \theta) = \prod_{n=1}^N p(y_n | \mathbf{x}_n, \theta) = \prod_{n=1}^N \mathcal{N}(y_n | \mathbf{x}_n^\top \theta, \sigma^2)$$

In practice, instead of maximising the likelihood function, we maximise the log of likelihood.

$$\begin{aligned} \theta_{ML} &= \underset{\theta}{\operatorname{argmax}} (\log(p(\mathcal{Y} | \mathcal{X}, \theta))) \\ &= \underset{\theta}{\operatorname{argmax}} \left(\log \left(\prod_{n=1}^N \mathcal{N}(y_n | \mathbf{x}_n^\top \theta, \sigma^2) \right) \right) \\ &= \underset{\theta}{\operatorname{argmax}} \left(\sum_{n=1}^N \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_n - \mathbf{x}_n^\top \theta)^2}{2\sigma^2}} \right) \right) \\ &= \underset{\theta}{\operatorname{argmax}} \left(\sum_{n=1}^N -\frac{(y_n - \mathbf{x}_n^\top \theta)^2}{2\sigma^2} + \text{const} \right) \\ &= \underset{\theta}{\operatorname{argmax}} \left(-\sum_{n=1}^N (y_n - \mathbf{x}_n^\top \theta)^2 \right) \end{aligned}$$

We can write the above equation in vectors. Notice that maximising the likelihood is same as minimising the sum of square errors.

$$\begin{aligned} \theta_{ML} &= \underset{\theta}{\operatorname{argmax}} \left(-(\mathbf{y} - \mathbf{X}\theta)^\top (\mathbf{y} - \mathbf{X}\theta) \right) \\ &= \underset{\theta}{\operatorname{argmin}} \left((\mathbf{y} - \mathbf{X}\theta)^\top (\mathbf{y} - \mathbf{X}\theta) \right) \end{aligned}$$

Where $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^\top$ is a $N \times D$ matrix and $\mathbf{y} = [y_1, y_2, \dots, y_N]$ is a column vector containing all the targets.

We can solve for θ_{ML} by finding the derivative and setting it to 0.

$$\frac{d}{d\theta} \left(-(\mathbf{y} - \mathbf{X}\theta)^\top (\mathbf{y} - \mathbf{X}\theta) \right) = \mathbf{0}^\top$$

$$\begin{aligned}
2(\mathbf{y} - \mathbf{X}\boldsymbol{\theta}_{ML})^\top \mathbf{X} &= \mathbf{0}^\top \\
\boldsymbol{\theta}_{ML}^\top \mathbf{X}^\top \mathbf{X} &= \mathbf{y}^\top \mathbf{X} \\
\boldsymbol{\theta}_{ML} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}
\end{aligned}$$

Notice that we require $\mathbf{X}^\top \mathbf{X}$ to be invertible, which is when all row in \mathbf{X} are linearly independent.

8.1.1 Maximum Likelihood with Features

The previous section only allows us to fit straight lines to data using maximum likelihood estimation. To overcome this problem, we can apply a nonlinear transformation $\phi(\mathbf{x})$ to the inputs $\mathbf{x} \in \mathbb{R}^D$. Therefore the relation between \mathbf{x} and observed value y becomes

$$y = \phi^\top(\mathbf{x})\boldsymbol{\theta} + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

Where $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^K$ is a nonlinear transformation to the input \mathbf{x} make up of $\phi_k : \mathbb{R}^D \rightarrow \mathbb{R}$.

$$\phi(\mathbf{x}) = \begin{bmatrix} \phi_0(\mathbf{x}) \\ \vdots \\ \phi_{K-1}(\mathbf{x}) \end{bmatrix} \in \mathbb{R}^K$$

We can group all the transformed inputs into one single matrix Φ .

$$\Phi = \begin{bmatrix} \phi^\top(\mathbf{x}_1) \\ \vdots \\ \phi^\top(\mathbf{x}_N) \end{bmatrix} = \begin{bmatrix} \phi_0(\mathbf{x}_1) & \dots & \phi_{K+1}(\mathbf{x}_1) \\ \vdots & & \vdots \\ \phi_0(\mathbf{x}_N) & \dots & \phi_{K+1}(\mathbf{x}_N) \end{bmatrix} \in \mathbb{R}^{N \times K}$$

Then find the parameters using maximum likelihood as above

$$\boldsymbol{\theta}_{ML} = (\Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y}$$

Similarly, we require $\Phi^\top \Phi$ to be invertible, meaning the rank of $\Phi = K$.

8.1.2 Overfitting

A problem with maximum likelihood estimation is that it can lead to severe overfitting if complex models are trained using data sets of limited size. This is because it is possible to increase the likelihood beyond any bound.

We can evaluate the quality of a model by computing the error incurred. In linear regression we often use the **Root Mean Squared Error (RMSE)**.

$$\sqrt{\frac{1}{N} \|\mathbf{y} - \Phi\boldsymbol{\theta}\|^2}$$

For model selection, we can use RMSE to determine the best degree of the polynomial by finding the polynomial degree M that minimises the loss function. Given a

dataset of size N , we can fit a polynomial of degree $0 \leq M \leq N - 1$.

We will notice that polynomials of low degree fit the data poorly. As the degree increase RMSE decreases. Until we reach even higher degrees the function passes through every data point and can not be used to generalise the data, this is when overfitting occurs.

8.2 Maximum a Posteriori Estimation

Maximum likelihood estimation is prone to overfitting. To mitigate the effect of large parameter values, we can place a prior distribution $p(\theta)$ on the parameters. Instead of maximising the likelihood, we find the parameters that maximise the posterior distribution.

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

We can set the prior to a Gaussian distribution to explicitly encode what parameter values are plausible.

$$p(\theta) = \mathcal{N}(\mathbf{0}, b^2 I)$$

To find the MAP estimate, we follow the steps similar to maximum likelihood estimation.

$$\begin{aligned} \theta_{MAP} &= \underset{\theta}{\operatorname{argmax}} (\log(p(\theta | \mathcal{X}, \mathcal{Y}))) \\ &= \underset{\theta}{\operatorname{argmax}} \left(\log \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_n - \mathbf{x}_n^\top \theta)^2}{2\sigma^2}} + \log \prod_{d=1}^D \frac{1}{\sqrt{2\pi b^2}} e^{-\frac{\theta^\top \theta}{2b^2}} \right) \\ &= \underset{\theta}{\operatorname{argmax}} \left(-\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \mathbf{x}_n^\top \theta)^2 - \frac{1}{2b^2} \theta^\top \theta \right) \\ &= \underset{\theta}{\operatorname{argmax}} \left(-(\mathbf{y} - \mathbf{X}\theta)^\top (\mathbf{y} - \mathbf{X}\theta) - \frac{\sigma^2}{b^2} \theta^\top \theta \right) \end{aligned}$$

Maximising the posterior is same as minimising the sum of square errors with L2 regularisation. We use the parameter λ to control the degree of regularisation.

$$\begin{aligned} \theta_{MAP} &= \underset{\theta}{\operatorname{argmax}} \left(-(\mathbf{y} - \mathbf{X}\theta)^\top (\mathbf{y} - \mathbf{X}\theta) - \frac{\sigma^2}{b^2} \theta^\top \theta \right) \\ &= \underset{\theta}{\operatorname{argmin}} \left((\mathbf{y} - \mathbf{X}\theta)^\top (\mathbf{y} - \mathbf{X}\theta) + \lambda \theta^\top \theta \right) \end{aligned}$$

We can solve for θ_{MAP} by finding the derivative and setting it to 0.

$$\frac{d}{d\theta} \left((\mathbf{y} - \mathbf{X}\theta)^\top (\mathbf{y} - \mathbf{X}\theta) + \frac{\sigma^2}{b^2} \theta^\top \theta \right) = \mathbf{0}^\top$$

$$\begin{aligned}
\frac{\sigma^2}{b^2} \boldsymbol{\theta}_{MAP}^\top &= (\mathbf{y} - \mathbf{X} \boldsymbol{\theta}_{MAP})^\top \mathbf{X} \\
\frac{\sigma^2}{b^2} \boldsymbol{\theta}_{MAP}^\top &= \mathbf{y}^\top \mathbf{X} - \boldsymbol{\theta}_{MAP}^\top \mathbf{X}^\top \mathbf{X} \\
\boldsymbol{\theta}_{MAP} &= \left(\mathbf{X}^\top \mathbf{X} + \frac{\sigma^2}{b^2} \mathbf{I} \right)^{-1} \mathbf{X}^\top \mathbf{y}
\end{aligned}$$

Comparing this with $\boldsymbol{\theta}_{ML}$, the only different is the additional term $\frac{\sigma^2}{b^2} \mathbf{I}$. If we set b^2 to be infinitely large, the term approaches 0 and we get the same result as maximum likelihood estimation. Intuitively, setting the variance b^2 of the prior to a large number removes the constraint on the parameters $\boldsymbol{\theta}$.

8.3 Bayesian Linear Regression

Unlike Maximum likelihood estimation and maximising a posteriori that aims to find a single best value of model parameter, Bayesian linear regression aims to determine the posterior distribution of the model parameters.

Not only is the response generated from a probability distribution, but the model parameters are assumed to come from a distribution as well. This means we do not fit any parameters, but compute the mean over all plausible parameters settings.

8.3.1 Prior Estimation

Given a prior distribution of the parameters $\boldsymbol{\theta}$, we can make a prediction of the noise corrupted target y_* at input \mathbf{x} by averaging all plausible parameter settings.

$$p(y_* | \mathbf{x}) = \int p(y_* | \mathbf{x}, \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta} = \mathbb{E}_{\boldsymbol{\theta}}[p(y_* | \mathbf{x})]$$

Since the likelihood is a Gaussian distribution, we choose a conjugate prior on $\boldsymbol{\theta}$ so that the above integration can be computed in close form.

$$\begin{aligned}
p(y | \mathbf{x}, \boldsymbol{\theta}) &= \mathcal{N}(y | \phi^\top(\mathbf{x})\boldsymbol{\theta}, \sigma^2) \\
p(\boldsymbol{\theta}) &= \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}_0, \mathbf{S}_0)
\end{aligned}$$

We can apply the marginalisation property of Gaussian to find $p(y_* | \mathbf{x})$.

$$\begin{aligned}
\mathbb{E}[y_*] &= \mathbb{E}[\phi^\top(\mathbf{x})\boldsymbol{\theta} + \epsilon] = \phi^\top(\mathbf{x})\mathbf{m}_0 \\
\mathbb{V}[y_*] &= \mathbb{V}[\phi^\top(\mathbf{x})\boldsymbol{\theta} + \epsilon] = \phi^\top(\mathbf{x})\mathbf{S}_0\phi(\mathbf{x}) + \sigma^2 \\
p(y_* | \mathbf{x}) &= \mathcal{N}(y_* | \phi^\top(\mathbf{x})\mathbf{m}_0, \phi^\top(\mathbf{x})\mathbf{S}_0\phi(\mathbf{x}) + \sigma^2)
\end{aligned}$$

The noise free target y is simply removing the noise variance σ^2 in the predictive variance.

$$p(y | \mathbf{x}) = \mathcal{N}(\phi^\top(\mathbf{x})\mathbf{m}_0, \phi^\top(\mathbf{x})\mathbf{S}_0\phi(\mathbf{x}))$$

8.3.2 Posterior Distribution

Given the distribution of the parameters θ , we can update the distribution using the training data to obtain a better set of plausible parameters. Given training set inputs \mathcal{X} and the corresponding observations \mathcal{Y} . We can compute the posterior distribution using Bayes' Theorem, which is a better representation of the parameter distribution.

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

The likelihood and prior are Gaussian distributions.

$$p(\mathcal{Y} | \mathcal{X}, \theta) = \mathcal{N}(\mathbf{y} | \Phi\theta, \sigma^2 \mathbf{I})$$

$$p(\theta) = \mathcal{N}(\theta | \mathbf{m}_0, \mathbf{S}_0)$$

The posterior distribution is in the form $\mathcal{N}(\theta | \mathbf{m}_N, \mathbf{S}_N)$. We can solve for \mathbf{m}_N and \mathbf{S}_N by finding the log posterior, log likelihood and log prior and comparing like terms.

The log posterior is

$$\begin{aligned} \log \mathcal{N}(\theta | \mathbf{m}_N, \mathbf{S}_N) &= -\frac{1}{2}(\theta - \mathbf{m}_N)^\top \mathbf{S}_N^{-1}(\theta - \mathbf{m}_N) \\ &= -\frac{1}{2}(\theta^\top \mathbf{S}_N^{-1} \theta - 2\mathbf{m}_N^\top \mathbf{S}_N^{-1} \theta) + \text{const} \end{aligned}$$

The sum of log likelihood is log prior is

$$\begin{aligned} \log \left(\frac{p(\mathcal{Y} | \mathcal{X}, \theta)p(\theta)}{p(\mathcal{Y} | \mathcal{X})} \right) &= \log \mathcal{N}(\mathbf{y} | \Phi\theta, \sigma^2 \mathbf{I}) + \log \mathcal{N}(\theta | \mathbf{m}_0, \mathbf{S}_0) + \text{const} \\ &= -\frac{1}{2}(\sigma^{-2}(\mathbf{y} - \Phi\theta)^\top (\mathbf{y} - \Phi\theta) + (\theta - \mathbf{m}_0)^\top \mathbf{S}_0^{-1}(\theta - \mathbf{m}_0)) + \text{const} \\ &= -\frac{1}{2}(\theta^\top (\sigma^{-2}\Phi^\top \Phi + \mathbf{S}_0^{-1})\theta - 2(\sigma^{-2}\Phi^\top \mathbf{y} + \mathbf{S}_0^{-1}\mathbf{m}_0)^\top \theta) + \text{const} \end{aligned}$$

By comparing like terms

$$\mathbf{S}_N = (\sigma^{-2}\Phi^\top \Phi + \mathbf{S}_0^{-1})^{-1}$$

$$\mathbf{m}_N = \mathbf{S}_N(\sigma^{-2}\Phi^\top \mathbf{y} + \mathbf{S}_0^{-1}\mathbf{m}_0)$$

After we find the posterior distribution of the parameters, we can use this distribution as the prior in prior estimation to predict samples.

8.3.3 Computing Marginal Likelihood

In the following section, we compute the marginal likelihood for Bayesian linear regression with a conjugate Gaussian prior on the parameters.

$$p(\theta) = \mathcal{N}(\theta | \mathbf{m}_0, \mathbf{S}_0) \quad \text{prior}$$

$$p(\mathcal{Y} | \mathcal{X}, \theta) = \mathcal{N}(\mathbf{y} | \mathbf{X}\theta, \sigma^2 \mathbf{I}) \quad \text{likelihood}$$

The marginal likelihood is given by

$$\begin{aligned} p(\mathcal{Y} | \mathcal{X}) &= \int p(\mathcal{Y} | \mathcal{X}, \theta) p(\theta) d\theta \\ &= \int \mathcal{N}(\mathbf{y} | \mathbf{X}\theta, \sigma^2 \mathbf{I}) \mathcal{N}(\theta | \mathbf{m}_0, \mathbf{S}_0) d\theta \end{aligned}$$

Recall from Section 2.7

1. The product of two Gaussians $\mathcal{N}(\mathbf{x} | \mathbf{a}, \mathbf{A}) \mathcal{N}(\mathbf{x} | \mathbf{b}, \mathbf{B})$ is a Gaussian in the form $z \mathcal{N}(\mathbf{x} | \mathbf{c}, \mathbf{C})$, where z is the functional form of a Gaussian density $\mathcal{N}(\mathbf{a} | \mathbf{b}, \mathbf{A} + \mathbf{B})$
2. We can apply a linear transformation to bring $\mathcal{N}(\theta | \mathbf{m}_0, \mathbf{S}_0)$ into the form $\mathcal{N}(\mathbf{y} | \mu, \Sigma)$ because $\mathbf{y} = \mathbf{X}\theta$.

$$\begin{aligned} p(\mathcal{Y} | \mathcal{X}) &= \int \mathcal{N}(\mathbf{y} | \mathbf{X}\theta, \sigma^2 \mathbf{I}) \mathcal{N}(\theta | \mathbf{m}_0, \mathbf{S}_0) d\theta \\ &= \int \mathcal{N}(\mathbf{y} | \mathbf{X}\theta, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{m}_0, \mathbf{X}\mathbf{S}_0\mathbf{X}^\top) d\theta \\ &= z \int \mathcal{N}(\mathbf{c}, \mathbf{C}) d\theta \\ &= \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{m}_0, \mathbf{X}\mathbf{S}_0\mathbf{X}^\top + \sigma^2 \mathbf{I}) \end{aligned}$$

9 Dimensionality Reduction with PCA

Storing high-dimensional data vectors is expensive, dimensionality reduction exploits structures and correlations and allow us to work with more compact representation of data. This chapter looks at an algorithm for linear dimensionality reduction called **Principle Component Analysis (PCA)**.

Consider a dataset $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ where $\mathbf{x}_n \in \mathbb{R}^D$, with mean $\mathbf{0}$ and covariance

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top$$

Furthermore, we assume there is a low-dimensional compressed representation of \mathbf{x}_n using the projection matrix $\mathbf{B} \in \mathbb{R}^M$, where $M < D$.

$$\begin{aligned} \mathbf{B} &= [\mathbf{b}_1 \quad \dots \quad \mathbf{b}_M] \in \mathbb{R}^{D \times M} \\ \mathbf{z}_n &= \mathbf{B}^\top \mathbf{x}_n \in \mathbb{R}^M \end{aligned}$$

The aim is to find matrix \mathbf{B} that retains as much information as possible when compressing the data. Retaining most information after data compression is equivalent to capturing the largest amount of variance in the low-dimensional code.

9.1 Direction with Maximum Variance

We will start by finding a single vector $\mathbf{b}_1 \in \mathbb{R}^D$ that maximises the variance of projected data, i.e. the first coordinate z_1 of \mathbf{z} .

$$\begin{aligned} \max \frac{1}{N} \sum_{n=1}^N z_{1n}^2 &= \operatorname{argmax}_{\mathbf{b}_1} \frac{1}{N} \sum_{n=1}^N (\mathbf{b}_1^\top \mathbf{x}_n)^2 \\ &= \operatorname{argmax}_{\mathbf{b}_1} \frac{1}{N} \sum_{n=1}^N \mathbf{b}_1^\top \mathbf{x}_n \mathbf{x}_n^\top \mathbf{b}_1 \\ &= \operatorname{argmax}_{\mathbf{b}_1} \mathbf{b}_1^\top \left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right) \mathbf{b}_1 \\ &= \operatorname{argmax}_{\mathbf{b}_1} \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 \end{aligned}$$

One problem is that \mathbf{b}_1 can become infinitely large to maximise the variance. Therefore we add a restriction $\|\mathbf{b}_1\| = 1$ and solved for the constrained optimisation problem using **Lagrangian**.

$$\mathcal{L}(\mathbf{b}_1, \lambda_1) = \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 + \lambda_1 (1 - \mathbf{b}_1^\top \mathbf{b}_1)$$

We can solve for the partial derivatives

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}_1} = 2\mathbf{S} \mathbf{b}_1 - 2\lambda_1 \mathbf{b}_1 \qquad \frac{\partial \mathcal{L}}{\partial \lambda_1} = 1 - \mathbf{b}_1^\top \mathbf{b}_1$$

Setting the derivatives to zero, we can obtain

$$\mathbf{S} \mathbf{b}_1 = \lambda_1 \mathbf{b}_1 \qquad \mathbf{b}_1^\top \mathbf{b}_1 = 1$$

Comparing this with the definition of Eigendecomposition, we see that \mathbf{b}_1 is an eigenvector of the covariance matrix \mathbf{S} , and λ_1 corresponds to the eigenvalue. Therefore, to maximise the variance of the low-dimensional code, we choose the basis vector associated with the largest eigenvalue of the data covariance matrix.

$$\max \mathbf{b}_1^\top \mathbf{S} \mathbf{b}_1 = \max \lambda_1 \mathbf{b}_1^\top \mathbf{b}_1 = \max \lambda_1$$

9.2 M-dimensional Maximum Variance

Now we consider the case of M dimensions, that is

$$\mathbf{z}_n = \begin{bmatrix} z_{1n} \\ \vdots \\ z_{Mn} \end{bmatrix} = \begin{bmatrix} \mathbf{b}_1^\top \mathbf{x}_n \\ \vdots \\ \mathbf{b}_M^\top \mathbf{x}_n \end{bmatrix} = \mathbf{B}^\top \mathbf{x}_n$$

Similar to above, we want to maximise the variance

$$\max \sum_{m=1}^M \sum_{n=1}^N z_{mn}^2 = \operatorname{argmax}_{\mathbf{B}} \frac{1}{N} \sum_{n=1}^N \sum_{m=1}^M (\mathbf{b}_m^\top \mathbf{x}_n)^2$$

$$\begin{aligned}
&= \operatorname{argmax}_{\mathbf{B}} \frac{1}{N} \sum_{n=1}^N \sum_{m=1}^M \mathbf{b}_m^\top \mathbf{x}_n \mathbf{x}_n^\top \mathbf{b}_m \\
&= \operatorname{argmax}_{\mathbf{B}} \sum_{m=1}^M \mathbf{b}_m^\top \left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right) \mathbf{b}_m \\
&= \operatorname{argmax}_{\mathbf{B}} \sum_{m=1}^M \mathbf{b}_m^\top \mathbf{S} \mathbf{b}_m \\
&= \operatorname{argmax}_{\mathbf{B}} \operatorname{tr}(\mathbf{B}^\top \mathbf{S} \mathbf{B})
\end{aligned}$$

Similar to above, we add a constraint such that \mathbf{B} is an orthogonal matrix. Hence, the Lagrangian we obtain becomes

$$\mathcal{L}(\mathbf{B}, \mathbf{\Lambda}) = \operatorname{tr}(\mathbf{B}^\top \mathbf{S} \mathbf{B}) + \operatorname{tr}(\mathbf{\Lambda}(\mathbf{I} - \mathbf{B}^\top \mathbf{B}))$$

Finding the derivative and setting them to $\mathbf{0}$, we obtain

$$\mathbf{S} \mathbf{B} = \mathbf{B} \mathbf{\Lambda} \Rightarrow \mathbf{S} \mathbf{b}_m = \lambda_m \mathbf{b}_m$$

This means the columns of \mathbf{B} corresponds to the eigenvectors of \mathbf{S} . To maximise the variance, it means finding the eigenvectors with the largest eigenvalues.

9.3 Computing PCA

In this section, we will look at how to compute PCA for **small sample size problems**, where the the number of features D is much larger than the number of samples N .

Lemma Let $\mathbf{B} = \mathbf{X} \mathbf{X}^\top$ and $\mathbf{C} = \mathbf{X}^\top \mathbf{X}$. Then

1. \mathbf{B} and \mathbf{C} have the same eigenvalues $\mathbf{\Lambda}$
2. The eigenvectors of \mathbf{B} , \mathbf{U} , and the eigenvectors of \mathbf{C} , \mathbf{V} , are related by $\mathbf{U} = \mathbf{X} \mathbf{V} \mathbf{\Lambda}^{-\frac{1}{2}}$

At the previous section, we saw that computing PCA requires computing the eigenvectors of $\mathbf{S} = \mathbf{X} \mathbf{X}^\top$, which has shape $D \times D$. When $D \gg N$, this is very inefficient. Using the Lemma, we can instead compute eigenvectors of \mathbf{S} from $\mathbf{X}^\top \mathbf{X}$, which is more efficient.

1. Centralised \mathbf{X} to mean $\mathbf{0}$. $\mathbf{X} = \{\mathbf{x}_1 - \boldsymbol{\mu}, \dots, \mathbf{x}_N - \boldsymbol{\mu}\}$
2. Compute dot product matrix $\mathbf{X}^\top \mathbf{X}$ and perform eigenanalysis $\mathbf{X}^\top \mathbf{X} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^\top$
3. Compute eigenvectors of $\mathbf{S} = \mathbf{X} \mathbf{X}^\top$ (We dropped $\frac{1}{N}$ for convenience) using Lemma, $\mathbf{U} = \mathbf{X} \mathbf{V} \mathbf{\Lambda}^{-\frac{1}{2}}$
4. Compute features $\mathbf{Y} = \mathbf{U}^\top \mathbf{X}$

9.4 Kernel PCA (KPCA)

PCA is a linear method. That is it can only be applied to datasets which are linearly separable. But, if we use it to non-linear datasets, we might get a result which may not be the optimal dimensionality reduction. Kernel PCA uses a kernel function to project dataset into a higher dimensional feature space, where it is linearly separable.

Given a dataset $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$. We can assume there is a non-linear mapping function that maps $\mathbf{x}_i \in \mathbb{R}^D$ to another feature space $\mathbf{X}_\Phi = [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_N)]$

$$\phi : \mathbb{R}^D \rightarrow \mathbb{R}^E$$

Usually this mapping function ϕ is unknown or is extremely expensive to compute. Instead, we only know the kernel functions which computes the dot product in feature space.

$$k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$$

We define the centralised feature space as

$$\bar{\mathbf{X}}_\Phi = [\phi(\mathbf{x}_1) - \boldsymbol{\mu}_\Phi, \dots, \phi(\mathbf{x}_N) - \boldsymbol{\mu}_\Phi] = \mathbf{X}_\Phi(\mathbf{I} - \mathbf{E}) \quad \boldsymbol{\mu}_\Phi = \frac{1}{N} \sum_{i=1}^N \phi(\mathbf{x}_i) \quad \mathbf{E} = \frac{1}{N} \mathbf{1}\mathbf{1}^\top$$

We define the kernel as

$$\mathbf{K} = [\phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)] = \mathbf{X}_\Phi^\top \mathbf{X}_\Phi$$

We also define a centralised kernel

$$\begin{aligned} \bar{\mathbf{K}} &= \left[(\phi(\mathbf{x}_i) - \boldsymbol{\mu}_\Phi) (\phi(\mathbf{x}_j) - \boldsymbol{\mu}_\Phi)^\top \right] \\ &= (\mathbf{I} - \mathbf{E}) \mathbf{X}_\Phi^\top \mathbf{X}_\Phi (\mathbf{I} - \mathbf{E}) \\ &= \bar{\mathbf{X}}_\Phi^\top \bar{\mathbf{X}}_\Phi \end{aligned}$$

We can now perform PCA on the feature space, solving the optimisation problem

$$\underset{\mathbf{B}}{\operatorname{argmax}} \operatorname{tr}(\mathbf{B}^\top \mathbf{S}_\Phi \mathbf{B}) \quad \text{Subject to } \mathbf{B}^\top \mathbf{B} = \mathbf{I}$$

We cannot find \mathbf{S}_Φ since we do not know the mapping function ϕ . However, we can use the Lemma defined in the previous section to compute \mathbf{B}

$$\begin{aligned} \bar{\mathbf{K}} &= \bar{\mathbf{X}}_\Phi^\top \bar{\mathbf{X}}_\Phi = \mathbf{V} \boldsymbol{\Lambda} \mathbf{V}^\top \\ \mathbf{S}_\Phi &= \bar{\mathbf{X}}_\Phi \bar{\mathbf{X}}_\Phi^\top = \mathbf{B} \boldsymbol{\Lambda} \mathbf{B}^\top \\ \mathbf{B} &= \bar{\mathbf{X}}_\Phi \mathbf{V} \boldsymbol{\Lambda}^{-\frac{1}{2}} \end{aligned}$$

We can use this \mathbf{B} to calculate the latent features of data in the feature space

$$\mathbf{z}_i = \mathbf{B}^\top (\phi(\mathbf{x}_i) - \boldsymbol{\mu}_\Phi)$$

$$\begin{aligned}
&= \mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{V}^\top \overline{\mathbf{X}}_\Phi^\top (\phi(\mathbf{x}_i) - \boldsymbol{\mu}_\Phi) \\
&= \mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{V}^\top (\mathbf{I} - \mathbf{E}) \mathbf{X}_\Phi^\top \left(\phi(\mathbf{x}_i) - \frac{1}{N} \mathbf{X}_\Phi \mathbf{1} \right) \\
&= \mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{V}^\top (\mathbf{I} - \mathbf{E}) \left(\mathbf{X}_\Phi^\top \phi(\mathbf{x}_i) - \frac{1}{N} \mathbf{K} \mathbf{1} \right) \\
&= \mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{V}^\top (\mathbf{I} - \mathbf{E}) \left(\begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_i) \\ \vdots \\ k(\mathbf{x}_N, \mathbf{x}_i) \end{bmatrix} - \frac{1}{N} \mathbf{K} \mathbf{1} \right)
\end{aligned}$$

10 Dimensionality Reduction with LDA

Both LDA and PCA are linear transformation techniques, while PCA is used for unsupervised approach, **Linear Discriminant Analysis (LDA)** is used in supervised classification problems.

Consider a dataset $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ each with a corresponding label l_i , where $l_i \in 1, \dots, C$ and C is the number of classes. The aim of LDA is the project \mathbf{x}_i to a low-dimensional space \mathbf{y}_i such that

1. The distance between mean of classes is maximised
2. The variance within each class is minimised

10.1 LDA with Two Classes

This section looks at LDA when there are only two classes, i.e. $C = 2$. We want to project \mathbf{x}_i to a latent space of \mathbf{y}_i such that the two classes have means $\mu_y(c_1)$, $\mu_y(c_2)$ and variances $\sigma_y^2(c_1)$, $\sigma_y^2(c_2)$. The goal is to

1. Maximise the distance between mean of classes $(\mu_y(c_1) - \mu_y(c_2))^2$
2. Minimise the variance within each class $\sigma_y^2(c_1) + \sigma_y^2(c_2)$

$$\begin{aligned}
\sigma_y^2(c_1) + \sigma_y^2(c_2) &= \frac{1}{N_{c_1}} \sum_{\mathbf{x}_i \in c_1} (\mathbf{y}_i - \mu_y(c_1))^2 + \frac{1}{N_{c_2}} \sum_{\mathbf{x}_i \in c_2} (\mathbf{y}_i - \mu_y(c_2))^2 \\
&= \frac{1}{N_{c_1}} \sum_{\mathbf{x}_i \in c_1} (\mathbf{w}^\top (\mathbf{x}_i - \mu_x(c_1)))^2 + \frac{1}{N_{c_2}} \sum_{\mathbf{x}_i \in c_2} (\mathbf{w}^\top (\mathbf{y}_i - \mu_x(c_2)))^2 \\
&= \mathbf{w}^\top \frac{1}{N_{c_1}} \sum_{\mathbf{x}_i \in c_1} ((\mathbf{x}_i - \mu_x(c_1))(\mathbf{x}_i - \mu_x(c_1))^\top) \mathbf{w} + \\
&\quad \mathbf{w}^\top \frac{1}{N_{c_2}} \sum_{\mathbf{x}_i \in c_2} ((\mathbf{x}_i - \mu_x(c_2))(\mathbf{x}_i - \mu_x(c_2))^\top) \mathbf{w} \\
&= \mathbf{w}^\top \mathbf{S}_1 \mathbf{w} + \mathbf{w}^\top \mathbf{S}_2 \mathbf{w}
\end{aligned}$$

$$\begin{aligned}
&= \mathbf{w}^\top (\mathbf{S}_1 + \mathbf{S}_2) \mathbf{w} \\
&= \mathbf{w}^\top \mathbf{S}_w \mathbf{w}
\end{aligned}$$

We denote the within class variance $\mathbf{S}_1 + \mathbf{S}_2$ as \mathbf{S}_w

$$(\mu_y(c_1) - \mu_y(c_2))^2 = \mathbf{w}^\top (\mu_x(c_1) - \mu_x(c_2)) (\mu_x(c_1) - \mu_x(c_2))^\top \mathbf{w} = \mathbf{w}^\top \mathbf{S}_b \mathbf{w}$$

We denote the between class variance $(\mu_x(c_1) - \mu_x(c_2)) (\mu_x(c_1) - \mu_x(c_2))^\top$ as \mathbf{S}_b .

The quantity we want to maximise becomes

$$\frac{\mathbf{w}^\top \mathbf{S}_b \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_w \mathbf{w}}$$

This is equivalent to solving for the constrained optimisation problem using Lagrangian

$$\mathcal{L}(\mathbf{w}, \lambda) = \mathbf{w}^\top \mathbf{S}_b \mathbf{w} + \lambda(1 - \mathbf{w}^\top \mathbf{S}_w \mathbf{w})$$

Solving for the partial derivatives

$$\frac{\partial}{\partial \mathbf{w}} \mathcal{L}(\mathbf{w}, \lambda) = 2\mathbf{S}_b \mathbf{w} - 2\lambda \mathbf{S}_w \mathbf{w}$$

Setting the derivatives to 0, we can obtain

$$\mathbf{S}_w^{-1} \mathbf{S}_b \mathbf{w} = \lambda \mathbf{w}$$

Therefore, \mathbf{w} is the eigenvector that corresponds to the largest eigenvalue of $\mathbf{S}_w^{-1} \mathbf{S}_b$

10.2 LDA with Multi-class

We can now extend LDA to classify multiple classes, $C > 2$. The within class variance matrix is defined as

$$\mathbf{S}_w = \sum_{j=1}^C \frac{1}{N_{c_j}} \sum_{\mathbf{x}_i \in c_j} (\mathbf{x}_i - \mu_x(c_j)) (\mathbf{x}_i - \mu_x(c_j))^\top = \sum_{j=1}^C \mathbf{S}_j$$

The between class matrix is defined as

$$\mathbf{S}_b = \sum_{j=1}^C N_{c_j} (\mu_x(c_j) - \boldsymbol{\mu}) (\mu_x(c_j) - \boldsymbol{\mu})^\top$$

Where $\boldsymbol{\mu}$ is the mean of all the classes center.

The goal is to find the projection matrix \mathbf{W} by solving the Lagrangian

$$\mathcal{L}(\mathbf{W}, \boldsymbol{\Lambda}) = \text{tr}(\mathbf{W}^\top \mathbf{S}_b \mathbf{W}) + \text{tr}(\boldsymbol{\Lambda}(\mathbf{I} - \mathbf{W}^\top \mathbf{S}_w \mathbf{W}))$$

Finding the derivatives and setting it to zero, we obtain

$$\mathbf{S}_w^{-1} \mathbf{S}_b \mathbf{W} = \mathbf{W} \boldsymbol{\Lambda}$$

Therefore the columns of \mathbf{W} are the eigenvectors that corresponds to the largest eigenvalues of $\mathbf{S}_w^{-1} \mathbf{S}_b$.

10.3 Computing LDA

In this section, we will look at how to compute LDA for **small sample size problems**, where the number of features D is much larger than the number of samples N , and without assuming S_w is invertible.

10.3.1 Scatter Matrix

The scatter matrix is defined as

$$\mathbf{M} = \begin{bmatrix} \mathbf{E}_1 & 0 & \dots & 0 \\ 0 & \mathbf{E}_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \mathbf{E}_C \end{bmatrix} \quad \mathbf{E}_i = \begin{bmatrix} \frac{1}{N_{c_i}} & \dots & \frac{1}{N_{c_i}} \\ \vdots & \ddots & \vdots \\ \frac{1}{N_{c_i}} & \dots & \frac{1}{N_{c_i}} \end{bmatrix}_{N_{c_i} \times N_{c_i}} = \frac{1}{N_{c_i}} \mathbf{1}\mathbf{1}^\top$$

For example, given a dataset with 3 samples from class 1 and 2 samples from class 2, the scatter matrix is defined as

$$\mathbf{M} = \begin{bmatrix} \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 & 0 \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 & 0 \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 & 0 \\ 0 & 0 & 0 & \frac{1}{2} & \frac{1}{2} \\ 0 & 0 & 0 & \frac{1}{2} & \frac{1}{2} \end{bmatrix}$$

One interesting property of scatter matrix is that it is **idempotent**, meaning $\mathbf{M}\mathbf{M} = \mathbf{M}$.

We can express S_b and S_w in terms of the scatter matrix.

$$S_b = \sum_{j=1}^C N_{c_j} \mu_x(c_j) \mu_x(c_j)^\top \quad \text{Assume the center points are centered}$$

$$= \mathbf{X}\mathbf{M}\mathbf{M}\mathbf{X}^\top$$

$$S_w = \sum_{j=1}^C \sum_{x_i \in c_j} (\mathbf{x}_i - \mu_x(c_j)) (\mathbf{x}_i - \mu_x(c_j))^\top \quad \text{Drop } \frac{1}{N_{c_j}} \text{ for convenience}$$

$$= \sum_{j=1}^C \sum_{x_i \in c_j} (\mathbf{x}_i \mathbf{x}_i^\top - \mu_x(c_j) \mathbf{x}_i^\top - \mathbf{x}_i \mu_x(c_j)^\top + \mu_x(c_j) \mu_x(c_j)^\top)$$

$$= \sum_{j=1}^C \left(\sum_{x_i \in c_j} (\mathbf{x}_i \mathbf{x}_i^\top) - N_{c_j} \mu_x(c_j) \mu_x(c_j)^\top \right)$$

$$= \sum_{j=1}^C \sum_{x_i \in c_j} \mathbf{x}_i \mathbf{x}_i^\top - \sum_{j=1}^C N_{c_j} \mu_x(c_j) \mu_x(c_j)^\top$$

$$= \mathbf{X}\mathbf{X}^\top - \mathbf{X}\mathbf{M}\mathbf{X}^\top$$

$$= \mathbf{X}(\mathbf{I} - \mathbf{M})(\mathbf{I} - \mathbf{M})\mathbf{X}^\top$$

$(\mathbf{I} - \mathbf{M})$ is also idempotent

10.3.2 Simultaneous Diagonalisation

Recall in the original LDA problem we want to optimise

$$\operatorname{argmax}_W \operatorname{tr}(W^\top S_b W) \quad \text{Subject to } W^\top S_w W = I$$

We can now rewrite the optimisation problem expressing S_b and S_w in terms of M .

$$\operatorname{argmax}_W \operatorname{tr}(W^\top X M M X^\top W) \quad \text{Subject to } W^\top X(I - M)(I - M)X^\top W = I$$

Assume that we can express $W = UQ$. The above optimisation problem becomes

$$\operatorname{argmax}_Q \operatorname{tr}(Q^\top U^\top X M M X^\top U Q) \quad \text{Subject to } Q^\top U^\top X(I - M)(I - M)X^\top U Q = I$$

The next step is to find a U such that $U^\top X(I - M)(I - M)X^\top U = I$. The purpose of this step is so that we can rewrite the above optimisation into

$$\operatorname{argmax}_Q \operatorname{tr}(Q^\top U^\top X M M X^\top U Q) \quad \text{Subject to } Q^\top Q = I$$

This optimisation problem has the same form as Principle Component Analysis, where the columns of Q are the eigenvectors that corresponds to the largest eigenvalues of $U^\top X M M X^\top U$.

To solve for U , we perform eigenanalysis and whitening to $X(I - M)(I - M)X^\top$. We will use the Lemma in PCA section to compute eigenvectors more efficiently

1. Perform eigenanalysis $(I - M)X^\top X(I - M) = V \Lambda V^\top$
2. Use Lemma $U = X(I - M)V \Lambda^{-\frac{1}{2}}$
3. Perform whitening so all eigenvalues are normalised to 1, this is to create identity matrix. $U = U \Lambda^{-\frac{1}{2}} = X(I - M)V \Lambda^{-1}$

Obtaining U , we can solve for Q by performing eigenanalysis on $U^\top X M M X^\top U$.

We can now solve for $W = UQ$

11 Support Vector Machine

Support Vector Machines (SVM) is a supervised machine learning algorithm which is mostly used in classification problems. In this algorithm, each data is a vector $x_i \in \mathbb{R}^D$ that belongs to class $y_i \in \{-1, 1\}$. The goal of SVM is to find a separating hyperplane that separates the two classes.

More formally, SVM aims to find parameters w and b such that the distance between the lines $w^\top x + b = -1$ and $w^\top x + b = 1$ is maximised.

Let \tilde{x} be a point on $w^\top x + b = -1$. We know w is a vector perpendicular to the line. Starting from \tilde{x} , we move along direction w until we reach the line $w^\top x + b = 1$. This gives us the equations

$$\begin{cases} w^\top(\tilde{x} + tw) + b = 1 \\ w^\top \tilde{x} + b = -1 \end{cases} \Rightarrow \begin{cases} w^\top \tilde{x} + b + tw^\top w = 1 \\ w^\top \tilde{x} + b = -1 \end{cases} \Rightarrow tw^\top w = 2$$

This goal is to maximise the distance between the two lines

$$\begin{aligned} \max_{w,b} \|tw\| &= \max_{w,b} \frac{2}{w^\top w} \|w\| \\ &= \max_{w,b} \frac{2}{\|w\|} \\ &= \min_{w,b} \frac{1}{2} w^\top w \end{aligned}$$

We want to introduce additional constraints such that data in class 1 must be on the right of $w^\top x + b = 0$, whereas data in class 2 must be on the left.

$$\begin{aligned} w^\top x_i + b &\geq 1 && \text{if } y_i = 1 \\ w^\top x_i + b &\leq -1 && \text{if } y_i = -1 \end{aligned}$$

These two constraints can be combined and written as a single constraint

$$y_i(w^\top x_i + b) \geq 1, \quad i \in 1 \dots N$$

Hence, the problem can be written as a constrained optimisation problem

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} w^\top w \\ \text{Subject to} \quad & y_i(w^\top x_i + b) \geq 1, \quad i \in 1 \dots N \end{aligned}$$

11.1 Lagrangian Duality

We can use the method of Lagrange Multiplier to solve the constrained optimisation problem. Given an optimisation of the form

$$\begin{aligned} \min_w \quad & f(w) \\ \text{Subject to} \quad & g(w) \leq 0 \end{aligned}$$

We first define the Lagrangian to be the original objective function added to a weighted combination of the constraints.

$$\mathcal{L}(w, a) = f(w) + ag(w)$$

Then rewrite the original optimisation problem in the following form where p is the **primal solution**.

$$p = \min_w \max_{a \geq 0} \mathcal{L}(w, a)$$

Optimisation problems can be viewed from a different perspective as a dual problem, d is the **dual solution**.

$$d = \max_{a \geq 0} \min_w \mathcal{L}(w, a)$$

The solution to the dual problem provides a lower bound to the solution of the primal problem. The two solutions are equal if the following conditions are met

- $f(w)$ is convex
- $g(w)$ is affine

For SVM problems, these conditions hold and therefore we can solve for the optimisation problem using the dual.

11.1.1 Karush-Kuhn-Tucker Conditions

The method of Lagrange Multipliers is used to find the solution for optimisation problems constrained to one or more equalities. When our constraints also have inequalities, we need to extend the method to the KKT conditions:

- $g(w) \leq 0$ and $a \geq 0$
- $ag(w) = 0$

The KKT conditions hold at the optimal solution.

11.2 Solving SVM

Going back to the original problem

$$\min_{w, b} \frac{1}{2} w^\top w$$

$$\text{Subject to } y_i(w^\top x_i + b) \geq 1, \quad i \in 1 \dots N$$

We can now define the Lagrangian as

$$\mathcal{L}(w, b, a) = \frac{1}{2} w^\top w - \sum_{i=1}^N a_i (y_i(w^\top x_i + b) - 1)$$

And solve for the dual problem

$$\max_{a_i \geq 0} \min_{w, b} \mathcal{L}(w, b, a)$$

To solve for the dual, we first minimise $\mathcal{L}(w, b, a)$ with respect to w and b for a fixed value a .

$$\nabla_w \mathcal{L}(w, b, a) = w - \sum_{i=1}^N a_i y_i x_i \qquad \nabla_b \mathcal{L}(w, b, a) = - \sum_{i=1}^N a_i y_i$$

$$\nabla_w \mathcal{L}(w, b, a) = 0 \Rightarrow w = \sum_{i=1}^N a_i y_i x_i \quad \nabla_b \mathcal{L}(w, b, a) = 0 \Rightarrow \sum_{i=1}^N a_i y_i = 0$$

Putting this back into the original equation

$$\begin{aligned} \mathcal{L}(w, b, a) &= \frac{1}{2} w^\top w - \sum_{i=1}^N a_i (y_i (w^\top x_i + b) - 1) \\ \mathcal{L}(a) &= \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N a_i a_j y_i y_j x_i^\top x_j - \sum_{i=1}^N \sum_{j=1}^N a_i a_j y_i y_j x_i^\top x_j - b \sum_{i=1}^N a_i y_i + \sum_{i=1}^N a_i \\ &= -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N a_i a_j y_i y_j x_i^\top x_j - b(0) + \sum_{i=1}^N a_i \\ &= \sum_{i=1}^N a_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N a_i a_j y_i y_j x_i^\top x_j \end{aligned}$$

Putting this together with the constraints $a_i \geq 0$ and $\sum a_i y_i = 0$. We obtain the following optimisation problem

$$\begin{aligned} \max_a \quad & \mathbf{1}^\top a - \frac{1}{2} a^\top K a & K &= y_i y_j x_i^\top x_j \\ \text{Subject to} \quad & a_i \geq 0, \quad i \in 1 \dots N \\ \text{Subject to} \quad & a^\top y = 0 \end{aligned}$$

11.3 Slack Variables

Sometimes the training data is not linearly separable, we can introduce slack variables, ξ , to each data to allow misclassification. We define ξ for each data to be

$$\xi = \begin{cases} 0, & \text{Point is on the correct side of the margin} \\ 0 < \xi \leq 1, & \text{Point is inside margin but correct side of hyperplane} \\ > 1, & \text{Point is at wrong side of hyperplane} \end{cases}$$

Introducing slack variables, the optimisation problem becomes

$$\begin{aligned} \min_{w, b} \quad & \frac{1}{2} w^\top w + C \sum_{i=1}^N \xi_i \\ \text{Subject to} \quad & y_i (w^\top x_i + b) \geq 1 - \xi_i, \quad i \in 1 \dots N \\ \text{Subject to} \quad & \xi_i \geq 0, \quad i \in 1 \dots N \end{aligned}$$

Where C is a hyperparameter called the regularisation term. A smaller C leads to larger margin and is able to tolerate more mistakes, whereas a larger C will make constraints hard to be ignored. The Lagrangian is defined as

$$\mathcal{L}(w, b, \xi, a, r) = \frac{1}{2} w^\top w + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N a_i (y_i (w^\top x_i + b) - 1 + \xi_i) - \sum_{i=1}^N r_i \xi_i$$

Same as the previous section, we first minimise $\mathcal{L}(w, b, \xi, a, r)$ with respect to w , b and ξ .

$$\begin{aligned}\nabla_{\xi_i} \mathcal{L}(w, b, \xi, a, r) &= C - a_i - r_i \\ \nabla_{\xi_i} \mathcal{L}(w, b, \xi, a, r) &= 0 \Rightarrow C - a_i - r_i = 0\end{aligned}$$

Putting this back into the original equation

$$\begin{aligned}\mathcal{L}(w, b, \xi, a, r) &= \frac{1}{2} w^\top w + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N a_i (y_i (w^\top x_i + b) - 1 + \xi_i) - \sum_{i=1}^N r_i \xi_i \\ \mathcal{L}(a) &= \sum_{i=1}^N a_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N a_i a_j y_i y_j x_i^\top x_j + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N a_i \xi_i - \sum_{i=1}^N r_i \xi_i \\ &= \sum_{i=1}^N a_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N a_i a_j y_i y_j x_i^\top x_j + \sum_{i=1}^N \xi_i (C - a_i - r_i) \\ &= \sum_{i=1}^N a_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N a_i a_j y_i y_j x_i^\top x_j\end{aligned}$$

Putting this together with the new constraint $C - a_i - r_i = 0 \Rightarrow a_i \leq C$. We obtain the following optimisation problem

$$\begin{aligned}\max_a \quad & \mathbf{1}^\top a - \frac{1}{2} a^\top K a & K &= y_i y_j x_i^\top x_j \\ \text{Subject to} \quad & 0 \leq a_i \leq C, \quad i \in 1 \dots N \\ \text{Subject to} \quad & a^\top y = 0\end{aligned}$$