

# Modern Statistical Methods

Rajen D. Shah

`r.shah@statslab.cam.ac.uk`

Course webpage:

[http://www.statslab.cam.ac.uk/~rds37/modern\\_stat\\_methods.html](http://www.statslab.cam.ac.uk/~rds37/modern_stat_methods.html)

The field of statistics has undergone profound changes in recent decades. Firstly, the types of datasets that statisticians are asked to analyse have transformed dramatically. In the past, we typically dealt with datasets containing many observations and a modest number of carefully chosen variables. Today, by contrast, it is common to encounter datasets with thousands of variables—sometimes even far exceeding the number of observations. For instance, in genomics, we might measure the expression levels of several thousand genes but only across a few hundred tissue samples. Classical statistical methods are often simply not applicable in these “high-dimensional” settings. As the scale of data collection has expanded, so too has the scope of the questions we seek to answer. Whereas statistics was once primarily concerned with uncovering associations between variables, we are now increasingly interested in understanding the causal structure of data. And rather than focusing solely on prediction, we often aim to predict the effects of interventions. At the same time, the rapid rise of machine learning has provided us with powerful new tools. In this course, we will explore how these advances can be harnessed to tackle some of the modern statistical challenges outlined above. The selection of material is heavily biased towards my own interests, but I hope it will nevertheless give you a flavour of some of the most important recent methodological developments in statistics.

The course is divided into 4 chapters (of unequal size). Our **first chapter** will start by introducing ridge regression, a simple generalisation of ordinary least squares. Our study of this will lead us to some beautiful connections with functional analysis and ultimately one of the most successful and flexible classes of learning algorithms: kernel machines. The **second chapter** concerns the Lasso and its extensions. The Lasso has been at the centre of much of the developments that have occurred in high-dimensional statistics, and will allow us to perform regression in the seemingly hopeless situation when the number of parameters we are trying to estimate is larger than the number of observations. In the **third chapter** we will study graphical modelling and provide an introduction to the exciting field of causal inference. Where the previous chapters consider methods for relating a particular response variable to a potentially large collection of (explanatory) variables, in the third chapter, we will study how to infer relationships between the variables themselves and answer causal questions using so-called *double/debiased machine learning* approaches. In the **final chapter**, we will turn to the problem of multiple testing which concerns handling settings where we may be performing thousands of hypothesis tests at the same time.

Before we begin the main content of the course, we will briefly review two key classical statistical methods: ordinary least squares and maximum likelihood estimation. This will help to set the scene and provide a warm-up for the modern methods to come later.

# Classical statistics

## Ordinary least squares

Imagine data are available in the form of observations  $(Y_i, x_i) \in \mathbb{R} \times \mathbb{R}^p$ ,  $i = 1, \dots, n$ , and the aim is to infer a simple *regression function* relating the average value of a *response*,  $Y_i$ , and a collection of *predictors* or *variables*,  $x_i$ . This is an example of regression analysis, one of the most important tasks in statistics.

A *linear model* for the data assumes that it is generated according to

$$Y = X\beta^0 + \varepsilon, \quad (1)$$

where  $Y \in \mathbb{R}^n$  is the vector of responses;  $X \in \mathbb{R}^{n \times p}$  is the predictor matrix (or design matrix) with  $i$ th row  $x_i^\top$ ;  $\varepsilon \in \mathbb{R}^n$  represents random error; and  $\beta^0 \in \mathbb{R}^p$  is the unknown vector of coefficients.

Provided  $p \ll n$ , a sensible way to estimate  $\beta$  is by ordinary least squares (OLS). This yields an estimator  $\hat{\beta}^{\text{OLS}}$  with

$$\hat{\beta}^{\text{OLS}} := \arg \min_{\beta \in \mathbb{R}^p} \|Y - X\beta\|_2^2 = (X^\top X)^{-1} X^\top Y, \quad (2)$$

provided  $X$  has full column rank.

Under the assumptions that (i)  $\mathbb{E}(\varepsilon_i) = 0$  and (ii)  $\text{Var}(\varepsilon) = \sigma^2 I$ , we have that:

- $\mathbb{E}_{\beta^0, \sigma^2}(\hat{\beta}^{\text{OLS}}) = \mathbb{E}\{(X^\top X)^{-1} X^\top (X\beta^0 + \varepsilon)\} = \beta^0$ .
- $\text{Var}_{\beta^0, \sigma^2}(\hat{\beta}^{\text{OLS}}) = (X^\top X)^{-1} X^\top \text{Var}(\varepsilon) X (X^\top X)^{-1} = \sigma^2 (X^\top X)^{-1}$ .

The Gauss–Markov theorem states that OLS is the best linear unbiased estimator in our setting: for any other estimator  $\tilde{\beta}$  that is linear in  $Y$  (so  $\tilde{\beta} = AY$  for some fixed matrix  $A$ ), we have

$$\text{Var}_{\beta^0, \sigma^2}(\tilde{\beta}) - \text{Var}_{\beta^0, \sigma^2}(\hat{\beta}^{\text{OLS}})$$

is positive semi-definite.

## Maximum likelihood estimation

The method of least squares is just one way to construct an estimator. A more general technique is that of maximum likelihood estimation. Here given data  $y \in \mathbb{R}^n$  that we take as a realisation of a random variable  $Y$ , we specify its density  $f(y; \theta)$  up to some unknown vector of parameters  $\theta \in \Theta \subseteq \mathbb{R}^d$ , where  $\Theta$  is the parameter space. The likelihood function is a function of  $\theta$  for each fixed  $y$  given by

$$L(\theta) := L(\theta; y) = c(y)f(y; \theta),$$

where  $c(y)$  is an arbitrary constant of proportionality. The maximum likelihood estimate of  $\theta$  maximises the likelihood, or equivalently it maximises the log-likelihood

$$\ell(\theta) := \ell(\theta; y) = \log f(y; \theta) + \log(c(y)).$$

A key quantity in the context of maximum likelihood estimation is the *Fisher information* matrix  $i(\theta) := \text{Cov}_\theta(\nabla \ell(\theta))$ . It can be thought of as a measure of how hard it is to estimate  $\theta$  when it is the true parameter value. The Cramér–Rao lower bound states that if  $\tilde{\theta}$  is an unbiased estimator of  $\theta$ , then under regularity conditions,

$$\text{Var}_\theta(\tilde{\theta}) - i^{-1}(\theta)$$

is positive semi-definite.

A remarkable fact about maximum likelihood estimators (MLEs) is that (under quite general conditions) they are asymptotically normally distributed, asymptotically unbiased and asymptotically achieve the Cramér–Rao lower bound.

Assume that the Fisher information matrix when there are  $n$  observations,  $i^{(n)}(\theta)$  (where we have made the dependence on  $n$  explicit) satisfies  $i^{(n)}(\theta)/n \rightarrow I(\theta)$  for some positive definite matrix  $I$ . Then denoting the maximum likelihood estimator of  $\theta$  when there are  $n$  observations by  $\hat{\theta}^{(n)}$ , under regularity conditions, as the number of observations  $n \rightarrow \infty$  we have

$$\sqrt{n}(\hat{\theta}^{(n)} - \theta) \xrightarrow{d} N_d(0, I^{-1}(\theta)).$$

Returning to our linear model, if we assume in addition that  $\varepsilon_i \stackrel{\text{i.i.d.}}{\sim} N(0, \sigma^2)$ , then the log-likelihood for  $(\beta, \sigma^2)$  is

$$\ell(\beta, \sigma^2) = -\frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - x_i^\top \beta)^2.$$

We see that the maximum likelihood estimate of  $\beta$  and OLS coincide. It is easy to check that

$$i(\beta, \sigma^2) = \begin{pmatrix} \sigma^{-2} X^\top X & 0 \\ 0 & n\sigma^{-4}/2 \end{pmatrix}.$$

The general theory for MLEs would suggest that approximately  $\sqrt{n}(\hat{\beta} - \beta) \sim N_p(0, \sigma^2(n^{-1} X^\top X)^{-1})$ ; in fact it is straightforward to show that this distributional result is exact.

# Contents

<b>1</b>	<b>Kernel machines</b>	<b>1</b>
1.1	Ridge regression . . . . .	1
1.1.1	Connection to principal components analysis . . . . .	3
1.2	The kernel trick . . . . .	4
1.3	Kernels . . . . .	6
1.3.1	Examples of kernels . . . . .	7
1.3.2	Reproducing kernel Hilbert spaces . . . . .	9
1.3.3	The representer theorem . . . . .	12
1.4	Kernel ridge regression . . . . .	14
1.5	Large-scale kernel machines . . . . .	18
<b>2</b>	<b>The Lasso and extensions</b>	<b>20</b>
2.1	Model selection . . . . .	20
2.2	The Lasso estimator . . . . .	21
2.2.1	Prediction error of the Lasso (slow rate) . . . . .	22
2.2.2	Concentration inequalities I . . . . .	22
2.2.3	Some facts from optimisation theory and convex analysis . . . . .	24
2.2.4	Lasso solutions . . . . .	27
2.2.5	Variable selection . . . . .	28
2.2.6	Prediction and estimation . . . . .	29
2.2.7	The compatibility condition . . . . .	31
2.2.8	Concentration inequalities II . . . . .	32
2.2.9	Random design . . . . .	34
2.2.10	Computation . . . . .	35
2.3	Extensions of the Lasso . . . . .	36
2.3.1	The square-root Lasso . . . . .	37
2.3.2	Structural penalties . . . . .	37
2.3.3	Correlated predictors . . . . .	38
2.3.4	Reducing the bias of the Lasso . . . . .	38
<b>3</b>	<b>Graphical modelling and causal inference</b>	<b>40</b>
3.1	Conditional independence . . . . .	40
3.2	Graphs . . . . .	41

3.3	Undirected graphical models . . . . .	43
3.4	Directed graphical models and causality . . . . .	44
3.4.1	Structural causal models . . . . .	45
3.4.2	Interventions . . . . .	45
3.4.3	Markov properties for DAGs . . . . .	47
3.4.4	Causal structure learning . . . . .	48
3.5	Gaussian graphical models . . . . .	48
3.5.1	Nodewise regression . . . . .	49
3.5.2	The Graphical Lasso . . . . .	49
3.6	Basic asymptotic statistics . . . . .	50
3.7	Conditional independence testing . . . . .	51
3.8	Average treatment effect estimation . . . . .	54
<b>4</b>	<b>Multiple testing</b>	<b>57</b>
4.1	The closed testing procedure . . . . .	58
4.2	The False Discovery Rate . . . . .	59

# Chapter 1

## Kernel machines

Let us revisit the linear model with

$$Y_i = \mu^0 + x_i^\top \beta^0 + \varepsilon_i.$$

Note that here we have included an explicit intercept term, for reasons that will become clear later. However, our interest will continue to centre on  $\beta^0$ , which quantifies the contributions of each of the predictors to the regression function. For unbiased estimators of  $\beta^0$ , their variance gives a way of comparing their quality in terms of squared error loss. For a potentially biased estimator,  $\tilde{\beta}$ , the relevant quantity is the mean-squared error (MSE),

$$\begin{aligned}\mathbb{E}_{\beta^0, \sigma^2}\{(\tilde{\beta} - \beta^0)(\tilde{\beta} - \beta^0)^\top\} &= \mathbb{E}\{[\tilde{\beta} - \mathbb{E}(\tilde{\beta}) + \mathbb{E}(\tilde{\beta}) - \beta^0][\tilde{\beta} - \mathbb{E}(\tilde{\beta}) + \mathbb{E}(\tilde{\beta}) - \beta^0]^\top\} \\ &= \text{Var}(\tilde{\beta}) + \{\mathbb{E}(\tilde{\beta} - \beta^0)\}\{\mathbb{E}(\tilde{\beta} - \beta^0)\}^\top,\end{aligned}$$

a sum of squared bias and variance terms. A crucial part of the optimality arguments for OLS and MLEs was *unbiasedness*. Do there exist biased methods whose variance is reduced compared to OLS such that their overall prediction error is lower? Yes—in fact the use of biased estimators is essential in dealing with settings where the number of parameters to be estimated is large compared to the number of observations. In the first two chapters we will explore two important methods for variance reduction based on different forms of penalisation: rather than forming estimators via optimising a least squares or log-likelihood term, we will introduce an additional penalty term that encourages estimates to be shrunk towards 0 in some sense. This will allow us to produce reliable estimators that work well when classical MLEs are infeasible, and in other situations can greatly outperform the classical approaches.

### 1.1 Ridge regression

One way to reduce the variance of  $\hat{\beta}^{\text{OLS}}$  is to shrink the estimated coefficients towards 0. *Ridge regression* [Hoerl and Kennard, 1970] does this by solving the following optimisation

problem

$$(\hat{\mu}_\lambda^R, \hat{\beta}_\lambda^R) = \arg \min_{(\mu, \beta) \in \mathbb{R} \times \mathbb{R}^p} \{\|Y - \mu \mathbf{1} - X\beta\|_2^2 + \lambda \|\beta\|_2^2\}.$$

Here  $\mathbf{1}$  is an  $n$ -vector of 1's. We see that the usual OLS objective is penalised by an additional term proportional to  $\|\beta\|_2^2$ . The parameter  $\lambda \geq 0$ , which controls the severity of the penalty and therefore the degree of the shrinkage towards 0, is known as a *regularisation parameter* or *tuning parameter*. We have explicitly included an intercept term which is not penalised. The reason for this is that were the variables to have their origins shifted so e.g. a variable representing temperature is given in units of Kelvin rather than Celsius, the fitted values would not change. However,  $X\hat{\beta}$  is not invariant under scale transformations of the variables so it is standard practice to centre each column of  $X$  (hence making them orthogonal to the intercept term) and then scale them to have  $\ell_2$ -norm  $\sqrt{n}$ .

It is straightforward to show that after this standardisation of  $X$ ,  $\hat{\mu}_\lambda^R = \bar{Y} := \sum_{i=1}^n Y_i/n$ , and

$$\hat{\beta}_\lambda^R = (X^\top X + \lambda I)^{-1} X^\top Y.$$

In this form, we can see how the addition of the  $\lambda I$  term helps to stabilise the estimator. Note that when  $X$  does not have full column rank (such as in high-dimensional situations), we can still compute this estimator. On the other hand, when  $X$  does have full column rank, we have the following theorem.

**Theorem 1.** *For  $\lambda$  sufficiently small (depending on  $\beta^0$  and  $\sigma^2$ ),*

$$\mathbb{E}(\hat{\beta}^{\text{OLS}} - \beta^0)(\hat{\beta}^{\text{OLS}} - \beta^0)^\top - \mathbb{E}(\hat{\beta}_\lambda^R - \beta^0)(\hat{\beta}_\lambda^R - \beta^0)^\top$$

*is positive definite.*

*Proof.* First we compute the bias of  $\hat{\beta}_\lambda^R$ . We drop the subscript  $\lambda$  and superscript  $R$  for convenience.

$$\begin{aligned} \mathbb{E}(\hat{\beta}) - \beta^0 &= (X^\top X + \lambda I)^{-1} X^\top X \beta^0 - \beta^0 \\ &= (X^\top X + \lambda I)^{-1} (X^\top X + \lambda I - \lambda I) \beta^0 - \beta^0 \\ &= -\lambda (X^\top X + \lambda I)^{-1} \beta^0. \end{aligned}$$

Now we look at the variance of  $\hat{\beta}$ .

$$\begin{aligned} \text{Var}(\hat{\beta}) &= \mathbb{E}\{(X^\top X + \lambda I)^{-1} X^\top \varepsilon\} \{(X^\top X + \lambda I)^{-1} X^\top \varepsilon\}^\top \\ &= \sigma^2 (X^\top X + \lambda I)^{-1} X^\top X (X^\top X + \lambda I)^{-1}. \end{aligned}$$

Thus  $\mathbb{E}(\hat{\beta}^{\text{OLS}} - \beta^0)(\hat{\beta}^{\text{OLS}} - \beta^0)^\top - \mathbb{E}(\hat{\beta} - \beta^0)(\hat{\beta} - \beta^0)^\top$  is equal to

$$\sigma^2 (X^\top X)^{-1} - \sigma^2 (X^\top X + \lambda I)^{-1} X^\top X (X^\top X + \lambda I)^{-1} - \lambda^2 (X^\top X + \lambda I)^{-1} \beta^0 \beta^{0\top} (X^\top X + \lambda I)^{-1}.$$

After some simplification, we see that this is equal to

$$\lambda (X^\top X + \lambda I)^{-1} [\sigma^2 \{2I + \lambda (X^\top X)^{-1}\} - \lambda \beta^0 \beta^{0\top}] (X^\top X + \lambda I)^{-1}.$$

Thus  $\mathbb{E}(\hat{\beta}^{\text{OLS}} - \beta^0)(\hat{\beta}^{\text{OLS}} - \beta^0)^\top - \mathbb{E}(\hat{\beta} - \beta^0)(\hat{\beta} - \beta^0)^\top$  is positive definite for  $\lambda > 0$  if and only if

$$\sigma^2\{2I + \lambda(X^\top X)^{-1}\} - \lambda\beta^0\beta^{0\top}$$

is positive definite, which is true for  $\lambda > 0$  sufficiently small (we can take  $0 < \lambda < 2\sigma^2/\|\beta^0\|_2^2$ ).  $\square$

The theorem says that  $\hat{\beta}_\lambda^{\text{R}}$  outperforms  $\hat{\beta}^{\text{OLS}}$  provided  $\lambda$  is chosen appropriately. To be able to use ridge regression effectively, we need a way of selecting a good  $\lambda$ —we will come to this very shortly. What the theorem doesn't really tell us is in what situations we expect ridge regression to perform well. To understand that, we will turn to one of the key matrix decompositions used in statistics, the singular value decomposition (SVD).

### 1.1.1 Connection to principal components analysis

The singular value decomposition (SVD) is a generalisation of an eigendecomposition of a square matrix. We can factorise any  $X \in \mathbb{R}^{n \times p}$  into its SVD

$$X = UDV^\top.$$

Here the  $U \in \mathbb{R}^{n \times n}$  and  $V \in \mathbb{R}^{p \times p}$  are orthogonal matrices and  $D \in \mathbb{R}^{n \times p}$  has  $D_{11} \geq D_{22} \geq \dots \geq D_{mm} \geq 0$  where  $m := \min(n, p)$  and all other entries of  $D$  are zero. To compute such a decomposition requires  $O(np \min(n, p))$  operations. The  $r$ th columns of  $U$  and  $V$  are known as the  $r$ th left and right singular vectors of  $X$  respectively, and  $D_{rr}$  is the  $r$ th singular value.

When  $n > p$ , we can replace  $U$  by its first  $p$  columns and  $D$  by its first  $p$  rows to produce another version of the SVD (sometimes known as the thin SVD). Then  $X = UDV^\top$  where  $U \in \mathbb{R}^{n \times p}$  has orthonormal columns (but is no longer square) and  $D$  is square and diagonal. There is an analogous version for when  $p > n$ .

Let us take  $X \in \mathbb{R}^{n \times p}$  as our matrix of predictors and suppose  $n \geq p$ . Using the (thin) SVD we may write the fitted values<sup>1</sup> from ridge regression as follows.

$$\begin{aligned} X\hat{\beta}_\lambda^{\text{R}} &= X(X^\top X + \lambda I)^{-1}X^\top Y \\ &= UDV^\top(VD^2V^\top + \lambda I)^{-1}VDU^\top Y \\ &= UD(D^2 + \lambda I)^{-1}DU^\top Y \\ &= \sum_{j=1}^p U_j \frac{D_{jj}^2}{D_{jj}^2 + \lambda} U_j^\top Y. \end{aligned}$$

Here we have used the notation (that we shall use throughout the course) that  $U_j$  is the  $j$ th column of  $U$ . For comparison, the fitted values from OLS (when  $X$  has full column rank) are

$$X\hat{\beta}^{\text{OLS}} = X(X^\top X)^{-1}X^\top Y = UU^\top Y.$$

---

<sup>1</sup>There is a slight abuse of terminology here as we are ignoring the contribution of  $\hat{\mu}_\lambda^{\text{R}}$ .



Both OLS and ridge regression compute the coordinates  $(U_j^\top Y)_{j=1}^p$  of  $Y$  with respect to the basis of the column space of  $X$  given by the columns of  $U$ . Ridge regression then shrinks these coordinates by the factors  $D_{jj}^2/(D_{jj}^2 + \lambda)$ ; if  $D_{jj}$  is small, the amount of shrinkage will be larger.

To interpret this further, note that the SVD is intimately connected with Principal Components Analysis (PCA). Consider  $v \in \mathbb{R}^p$  with  $\|v\|_2 = 1$ . Since the columns of  $X$  have had their means subtracted, the sample variance of  $Xv \in \mathbb{R}^n$ , is

$$\frac{1}{n} v^\top X^\top X v = \frac{1}{n} v^\top V D^2 V^\top v.$$

Writing  $a = V^\top v$ , so  $\|a\|_2 = 1$ , we have

$$\frac{1}{n} v^\top V D^2 V^\top v = \frac{1}{n} a^\top D^2 a = \frac{1}{n} \sum_j a_j^2 D_{jj}^2 \leq \frac{1}{n} D_{11} \sum_j a_j^2 = \frac{1}{n} D_{11}^2.$$

As  $\|XV_1\|_2^2/n = D_{11}^2/n$ ,  $V_1$  determines the linear combination of the columns of  $X$  which has the largest sample variance, when the coefficients of the linear combination are constrained to have  $\ell_2$ -norm 1.  $XV_1 = D_{11}U_1$  is known as the first principal component of  $X$ . Subsequent principal components  $D_{22}U_2, \dots, D_{pp}U_p$  have maximum variance  $D_{jj}^2/n$ , subject to being orthogonal to all earlier ones—see example sheet 1 for details.

Returning to ridge regression, we see that it shrinks  $Y$  most in the smaller principal components of  $X$ . Thus it will work well when most of the signal is spanned by the large principal components of  $X$ . We now turn to the problem of choosing  $\lambda$ .

## 1.2 The kernel trick

An alternative expression for the ridge regression solution is given by the following

$$\begin{aligned} X^\top (XX^\top + \lambda I_n) &= (X^\top X + \lambda I_p) X^\top \\ (X^\top X + \lambda I_p)^{-1} X^\top &= X^\top (XX^\top + \lambda I_n)^{-1} \\ (X^\top X + \lambda I_p)^{-1} X^\top Y &= X^\top (XX^\top + \lambda I_n)^{-1} Y = \hat{\beta}_\lambda^R. \end{aligned} \tag{1.1}$$

Two remarks are in order:

- Note while  $X^\top X$  is  $p \times p$ ,  $XX^\top$  is  $n \times n$ . Computing fitted values via the the LHS of (1.1) would require roughly  $O(np^2 + p^3)$  operations. If  $p \gg n$  this could be extremely costly. However, the alternative formulation

$$X \hat{\beta}_\lambda^R = XX^\top (XX^\top + \lambda I)^{-1} Y$$

would only require roughly  $O(n^2p + n^3)$  operations, which could be substantially smaller.

- We see that the fitted values of ridge regression depend only on inner products  $K = XX^\top$  between observations (note  $K_{ij} = x_i^\top x_j$ ).

Now suppose that we believe the signal depends quadratically on the predictors:

$$Y_i = x_i^\top \beta + \sum_{k,l} x_{ik} x_{il} \theta_{kl} + \varepsilon_i.$$

We can still use ridge regression provided we work with an enlarged set of predictors

$$x_{i1}, \dots, x_{ip}, x_{i1}x_{i1}, \dots, x_{i1}x_{ip}, x_{i2}x_{i1}, \dots, x_{i2}x_{ip}, \dots, x_{ip}x_{ip}. \quad (1.2)$$

This will give us  $O(p^2)$  predictors. Our new approach to computing fitted values would therefore have complexity  $O(n^2p^2 + n^3)$ , which could be rather costly if  $p$  is large.

However, rather than first creating all the additional predictors and then computing the new  $K$  matrix, we can attempt to directly compute  $K$ . To this end consider

$$\begin{aligned} (1/2 + x_i^\top x_j)^2 - 1/4 &= \left( \frac{1}{2} + \sum_k x_{ik} x_{jk} \right)^2 - \frac{1}{4} \\ &= \sum_k x_{ik} x_{jk} + \sum_{k,l} x_{ik} x_{il} x_{jk} x_{jl}. \end{aligned}$$

Observe this amounts to an inner product between vectors of the form (1.2). Thus if we set

$$K_{ij} = (1/2 + x_i^\top x_j)^2 - 1/4 \quad (1.3)$$

and plug this into the formula for the fitted values, it is *exactly* as if we had performed ridge regression with an enlarged predictor matrix

$$\Phi := \begin{pmatrix} \phi(x_1)^\top \\ \vdots \\ \phi(x_n)^\top \end{pmatrix}.$$

Now computing  $K$  using (1.3) would require only  $O(p)$  operations per entry, so  $O(n^2p)$  operations in total, compared to  $O(n^2p^2)$  for our earlier approach.

Predictions at a new  $x \in \mathbb{R}^p$  may be computed similarly. From (1.1), we have

$$\phi(x)^\top \hat{\beta}_\lambda^R = \phi(x) \Phi^\top (\Phi \Phi^\top + \lambda I)^{-1} Y = \sum_{i=1}^n k(x, x_i) \hat{\alpha}_i$$

where  $\hat{\alpha} := (K + \lambda I)^{-1} Y \in \mathbb{R}^n$ .

This is a nice computational trick, but more importantly for us it serves to illustrate some general points.

- Since ridge regression only depends on inner products between observations, rather than fitting non-linear models by first mapping the original data  $x_i \in \mathbb{R}^p$  to  $\phi(x_i) \in \mathbb{R}^d$  (say) using some *feature map*  $\phi$  (which could, for example introduce quadratic effects), we can instead try to directly compute  $k(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$ .
- In fact rather than thinking in terms of feature maps, we can instead try to think about an appropriate measure of similarity  $k(x_i, x_j)$  between observations. Modelling in this fashion is sometimes much easier.

We will now formalise and extend what we have learnt with this example.

## 1.3 Kernels

We have seen how a model with quadratic effects can be fitted very efficiently by replacing the inner product matrix (known as the *Gram matrix*)  $XX^\top$  in (1.1) with the matrix in (1.3). It is then natural to ask what other non-linear models can be fitted efficiently using this sort of approach.

We won't answer this question directly, but instead we will try to understand the sorts of similarity measures  $k$  that can be represented as inner products between transformations of the original data.

That is, we will study the similarity measures  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  from the input space  $\mathcal{X}$  to  $\mathbb{R}$  for which there exists a *feature map*  $\phi : \mathcal{X} \rightarrow \mathcal{H}$  where  $\mathcal{H}$  is some (real) inner product space with

$$k(x, x') = \langle \phi(x), \phi(x') \rangle. \quad (1.4)$$

Recall that an inner product space is a real vector space  $\mathcal{H}$  endowed with a map  $\langle \cdot, \cdot \rangle : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$  that obeys the following properties.

- (i) Symmetry:  $\langle u, v \rangle = \langle v, u \rangle$ .
- (ii) Linearity: for  $a, b \in \mathbb{R}$   $\langle au + bw, v \rangle = a\langle u, v \rangle + b\langle w, v \rangle$ .
- (iii) Positive-definiteness:  $\langle u, u \rangle \geq 0$  with equality if and only if  $u = 0$ .

**Definition 1.** A *positive definite kernel* or more simply a *kernel* (for brevity)  $k$  is a symmetric map  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  for which for all  $n \in \mathbb{N}$  and all  $x_1, \dots, x_n \in \mathcal{X}$ , the matrix  $K$  with entries

$$K_{ij} = k(x_i, x_j)$$

is positive semi-definite.

A kernel is a little like an inner product, but need not be bilinear in general. However, a form of the Cauchy–Schwarz inequality does hold for kernels.

**Proposition 2.**

$$k(x, x')^2 \leq k(x, x)k(x', x').$$

*Proof.* The matrix

$$\begin{pmatrix} k(x, x) & k(x, x') \\ k(x', x) & k(x', x') \end{pmatrix}$$

must be positive semi-definite so in particular its determinant must be non-negative.  $\square$

First we show that any inner product of feature maps will give rise to a kernel.

**Proposition 3.**  *$k$  defined by  $k(x, x') = \langle \phi(x), \phi(x') \rangle$  is a kernel.*

*Proof.* Let  $x_1, \dots, x_n \in \mathcal{X}$ ,  $\alpha_1, \dots, \alpha_n \in \mathbb{R}$  and consider

$$\begin{aligned} \sum_{i,j} \alpha_i k(x_i, x_j) \alpha_j &= \sum_{i,j} \alpha_i \langle \phi(x_i), \phi(x_j) \rangle \alpha_j \\ &= \left\langle \sum_i \alpha_i \phi(x_i), \sum_j \alpha_j \phi(x_j) \right\rangle \geq 0. \end{aligned} \quad \square$$

Showing that every kernel admits a representation of the form (1.4) is slightly more involved, and we delay this until after we have studied some examples.

### 1.3.1 Examples of kernels

**Proposition 4.** *Suppose  $k_1, k_2, \dots : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  are kernels.*

- (i) *If  $a_1, a_2 \geq 0$  then  $a_1 k_1 + a_2 k_2$  is a kernel.*
- (ii) *If  $\lim_{m \rightarrow \infty} k_m(x, x') =: k(x, x')$  exists for all  $x, x' \in \mathcal{X}$ , then  $k$  is a kernel.*
- (iii) *The pointwise product  $k$  given by  $k(x, x') := k_1(x, x') k_2(x, x')$  is a kernel.*

*Proof.* Let  $K, K_1, K_2, \dots$  be the corresponding kernel matrices and take  $\alpha \in \mathbb{R}^n$ .

- (i)  $\alpha^\top K \alpha = a_1 \alpha^\top K_1 \alpha + a_2 \alpha^\top K_2 \alpha \geq 0$ .
- (ii)  $\alpha^\top K \alpha = \alpha^\top \lim_{m \rightarrow \infty} K_m \alpha = \lim_{m \rightarrow \infty} \alpha^\top K_m \alpha \geq 0$ .
- (iii) Let  $X$  and  $Y$  be independent random vectors with  $\text{Var}(X) = K_1$ ,  $\text{Var}(Y) = K_2$ . The entrywise (Hadamard) product  $K = K_1 \odot K_2$  satisfies

$$K_{ij} = \mathbb{E}(X_i X_j) \mathbb{E}(Y_i Y_j) = \mathbb{E}(X_i Y_i X_j Y_j) = (\text{Var}(X \odot Y))_{ij},$$

and  $\text{Var}(X \odot Y)$  is positive semi-definite as a covariance matrix.  $\square$

**Linear kernel.**  $k(x, x') = x^\top x'$ .

**Polynomial kernel.**  $k(x, x') = (1 + x^\top x')^d$ . To show this is a kernel, we can simply note that  $1 + x^\top x'$  gives a kernel owing to the fact that 1 is a kernel and (i) of Proposition 4. Next (ii) and induction shows that  $k$  as defined above is a kernel.

**Gaussian kernel.** The highly popular Gaussian kernel is defined by

$$k(x, x') = \exp\left(-\frac{\|x - x'\|_2^2}{2h^2}\right).$$

For  $x$  close to  $x'$  it is large whilst for  $x$  far from  $x'$  the kernel quickly decays towards 0. The additional parameter  $h > 0$  known as the *bandwidth* controls the speed of the decay to zero. Note it is less clear how one might find a corresponding feature map and indeed any feature map that represents this must be infinite dimensional.

To show that it is a kernel first decompose  $\|x - x'\|_2^2 = \|x\|_2^2 + \|x'\|_2^2 - 2x^\top x'$ . Note that by Proposition 3,

$$k_1(x, x') = \exp\left(-\frac{\|x\|_2^2}{2h^2}\right) \exp\left(-\frac{\|x'\|_2^2}{2h^2}\right)$$

is a kernel. Next writing

$$k_2(x, x') = \exp(x^\top x' / h^2) = \sum_{r=0}^{\infty} \frac{(x^\top x' / h^2)^r}{r!}$$

and using (i) of Proposition 4 shows that  $k_2$  is a kernel. Finally observing that  $k = k_1 k_2$  and using (ii) shows that the Gaussian kernel is indeed a kernel.

**First order Sobolev kernel.** Take  $\mathcal{X}$  to be  $[0, 1]$  and let  $k(x, x') := x \wedge x' = \min(x, x')$ . We have

$$k(x, x') = \int_0^1 \mathbb{1}_{[0, x]}(u) \mathbb{1}_{[0, x']}(u) du = \langle \mathbb{1}_{[0, x]}, \mathbb{1}_{[0, x']} \rangle$$

so  $k$  is a kernel by Proposition 3.

**Second order Sobolev kernel.** Take  $\mathcal{X}$  to be  $[0, 1]$  and let

$$k(x, x') := \int_0^{x \wedge x'} \int_0^{x \wedge y} (x - u)(y - u) du$$

**Jaccard similarity kernel.** Take  $\mathcal{X}$  to be the set of all subsets of  $\{1, \dots, p\}$ . For  $x, x' \in \mathcal{X}$  with  $x \cup x' \neq \emptyset$  define

$$k(x, x') = \frac{|x \cap x'|}{|x \cup x'|}$$

and if  $x \cup x' = \emptyset$  then set  $k(x, x') = 1$ . Showing that this is a kernel is left to the example sheet.

### 1.3.2 Reproducing kernel Hilbert spaces

Recall that we wish to show that each kernel  $k$  admits a representation for the form  $k(x, x') = \langle \phi(x), \phi(x') \rangle$  for some feature map  $\phi : \mathcal{X} \rightarrow \mathcal{H}$  where  $\mathcal{H}$  is an inner product space. Before showing this, let us first consider the case where  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  for finite  $\mathcal{X}$ , so without loss of generality,  $\mathcal{X} = \{1, \dots, N\}$  for some  $N \in \mathbb{N}$ . Then  $K \in \mathbb{R}^{N \times N}$  with  $K_{ij} = k(i, j)$  contains all the information about  $k$ . The eigendecomposition  $K = P^\top D P$  then readily confirms what we want to show: taking  $\phi(i) = D^{1/2} P_i$ , we have

$$k(i, j) = K_{ij} = P_i^\top D P_j = (D^{1/2} P_i)^\top (D^{1/2} P_j) = \phi(i)^\top \phi(j).$$

Note however that representation of  $k$  through an inner product is not unique. Consider  $\phi(i) = K_i$  and the weighted Euclidean inner product on the column space  $\mathcal{H}_c$  of  $K$  given by  $\langle u, v \rangle = u^\top P^\top D^+ P u$ , where  $D^+$  has  $(D^+)_{ij} = D_{ij}^{-1}$  when  $D_{ij} > 0$  and 0 otherwise. Then for  $\alpha \in \mathbb{R}^N$ ,

$$\langle \phi(i), K\alpha \rangle = (P^\top D P e_i)^\top P^\top D^+ P P^\top D P \alpha = e_i^\top K \alpha = (K\alpha)_i,$$

so

$$\langle \phi(i), \phi(j) \rangle = (K_j)_i = K_{ij}$$

as required. Note also the interesting property that the inner product of  $\phi(i)$  and  $K\alpha$  extracts the  $i$ th component of  $K\alpha$ . It is this second representation that generalises most fruitfully to the case where  $\mathcal{X}$  may be infinite.

Consider now taking  $\mathcal{H}_0$  to be the linear span of the functions  $\{k(\cdot, x) : x \in \mathcal{X}\}$ , i.e. the vector space of functions of the form

$$f(\cdot) = \sum_{i=1}^n \alpha_i k(\cdot, x_i), \tag{1.5}$$

where  $n \in \mathbb{N}$ ,  $x_i \in \mathcal{X}$  and  $\alpha_i \in \mathbb{R}$ . Given another function

$$g(\cdot) = \sum_{j=1}^m \beta_j k(\cdot, x'_j) \tag{1.6}$$

we define their inner product to be

$$\langle f, g \rangle = \sum_{i=1}^n \sum_{j=1}^m \alpha_i \beta_j k(x_i, x'_j). \tag{1.7}$$

We need to check this is well-defined as the representations of  $f$  and  $g$  in (1.5) and (1.6) need not be unique. To this end, note that

$$\sum_{i=1}^n \sum_{j=1}^m \alpha_i \beta_j k(x_i, x'_j) = \sum_{i=1}^n \alpha_i g(x_i) = \sum_{j=1}^m \beta_j f(x'_j). \tag{1.8}$$

The first equality shows that the inner product does not depend on the particular expansion of  $g$  whilst the second equality shows that it also does not depend on the expansion of  $f$ . Thus the inner product is well-defined. We define our feature map  $\phi : \mathcal{X} \rightarrow \mathcal{H}_0$  to be

$$\phi(x) = k(\cdot, x). \quad (1.9)$$

**Theorem 5.** *For every kernel  $k$ , the space  $\mathcal{H}_0$  above is an inner product space and the feature map (1.9) satisfies*

$$k(x, x') = \langle \phi(x), \phi(x') \rangle. \quad (1.10)$$

*Proof.* First we check that with  $\phi$  defined as in (1.9) we do have relationship (1.10). Observe that

$$\langle k(\cdot, x), f \rangle = \sum_{i=1}^n \alpha_i k(x_i, x) = f(x), \quad (1.11)$$

so in particular we have

$$\langle \phi(x), \phi(x') \rangle = \langle k(\cdot, x), k(\cdot, x') \rangle = k(x, x').$$

It remains to show that (1.7) is indeed an inner product. It is clearly symmetric and (1.8) shows linearity. We now need to show positive definiteness.

First note that

$$\langle f, f \rangle = \sum_{i,j} \alpha_i k(x_i, x_j) \alpha_j \geq 0 \quad (1.12)$$

by positive definiteness of the kernel. Now from (1.11),

$$f(x)^2 = (\langle k(\cdot, x), f \rangle)^2.$$

If we could use the Cauchy–Schwarz inequality on the right-hand side, we would have

$$f(x)^2 \leq \langle k(\cdot, x), k(\cdot, x) \rangle \langle f, f \rangle, \quad (1.13)$$

which would show that if  $\langle f, f \rangle = 0$  then necessarily  $f = 0$ ; the final property we need to show that  $\langle \cdot, \cdot \rangle$  is an inner product. However, in order to use the traditional Cauchy–Schwarz inequality we need to first know we’re dealing with an inner product, which is precisely what we’re trying to show!

Although we haven’t yet shown that  $\langle \cdot, \cdot \rangle$  is an inner product, we do have enough information to show that it is itself a kernel. We may then appeal to Proposition 2 to obtain (1.13). With this in mind, we argue as follows. Given functions  $f_1, \dots, f_m$  and coefficients  $\gamma_1, \dots, \gamma_m \in \mathbb{R}$ , we have

$$\sum_{i,j} \gamma_i \langle f_i, f_j \rangle \gamma_j = \left\langle \sum_i \gamma_i f_i, \sum_j \gamma_j f_j \right\rangle \geq 0$$

where we have used linearity and (1.12), showing that it is a kernel.  $\square$

To further discuss the space  $\mathcal{H}_0$  we recall some facts from analysis. Any inner product space  $\mathcal{B}$  is also a normed space: for  $f \in \mathcal{B}$  we may define  $\|f\|_{\mathcal{B}}^2 := \langle f, f \rangle_{\mathcal{B}}$ . Recall that a Cauchy sequence  $(f_m)_{m=1}^{\infty}$  in  $\mathcal{B}$  has  $\|f_m - f_n\|_{\mathcal{B}} \rightarrow 0$  as  $n, m \rightarrow \infty$ . A normed space where every Cauchy sequence has a limit (in the space) is called *complete*, and a complete inner product space is called a *Hilbert space*.

Hilbert spaces may be thought of as the (potentially) infinite-dimensional analogues of finite-dimensional Euclidean spaces. For later use we note that if  $V$  is a closed subspace of a Hilbert space  $\mathcal{B}$ , then any  $f \in \mathcal{B}$  has a decomposition  $f = u + v$  with  $u \in V$  and

$$v \in V^{\perp} := \{v \in \mathcal{B} : \langle v, u \rangle_{\mathcal{B}} = 0 \text{ for all } u \in V\}.$$

Moreover, if  $V$  is finite-dimensional, then it is closed.

By adding the limits of Cauchy sequences to  $\mathcal{H}_0$  (from Theorem 5) we can make create a Hilbert space. If  $(f_m)_{m=1}^{\infty} \in \mathcal{H}$  is Cauchy, then since by (1.13) we have

$$|f_m(x) - f_n(x)| \leq \sqrt{k(x, x)} \|f_m - f_n\|_{\mathcal{H}},$$

we may define function  $f^* : \mathcal{X} \rightarrow \mathbb{R}$  by  $f^*(x) = \lim_{m \rightarrow \infty} f_m(x)$ . We can check that all such  $f^*$  can be added to  $\mathcal{H}_0$  to create a Hilbert space.

In fact, the completion of  $\mathcal{H}_0$  is a special type of Hilbert space known as a *reproducing kernel Hilbert space* (RKHS).

**Definition 2.** A Hilbert space  $\mathcal{H}$  of functions  $f : \mathcal{X} \rightarrow \mathbb{R}$  is a *reproducing kernel Hilbert space* (RKHS) if for all  $x \in \mathcal{X}$ , there exists  $k_x \in \mathcal{B}$  such that it satisfies the reproducing property

$$f(x) = \langle k_x, f \rangle \quad \text{for all } f \in \mathcal{B}.$$

The function

$$\begin{aligned} k : \mathcal{X} \times \mathcal{X} &\rightarrow \mathbb{R} \\ (x, x') &\mapsto \langle k_x, k_{x'} \rangle = k_{x'}(x) \end{aligned}$$

is known as the *reproducing kernel* of  $\mathcal{H}$ .

Note that the reproducing kernel is well-defined: if  $k_x$  and  $h_x$  satisfy the reproducing property above, then

$$\|k_x - h_x\|_{\mathcal{H}}^2 = \langle k_x, k_x - h_x \rangle - \langle h_x, k_x - h_x \rangle = (k_x - h_x)(x) - (k_x - h_x)(x) = 0.$$

To summarise what we have learnt, by Proposition 3, the reproducing kernel of any RKHS is a (positive definite) kernel, and Theorem 5 shows that to any kernel  $k$  is associated an RKHS that has reproducing kernel  $k$ . One can further show that this is the unique RKHS with  $k$  as its reproducing kernel.

## Examples

**Linear kernel.** Here  $\mathcal{H} = \{f : f(x) = \beta^{\top} x, \beta \in \mathbb{R}^p\}$  and if  $f(x) = \beta^{\top} x$  then  $\|f\|_{\mathcal{H}}^2 = \|\beta\|_2^2$ .



**First-order Sobolev kernel.** Take  $\mathcal{H}$  to be the class of almost everywhere differentiable functions  $f : [0, 1] \rightarrow \mathbb{R}$ , with  $f(0) = 0$  and  $\int_0^1 (f'(u))^2 du < \infty$ . This is an RKHS with reproducing kernel  $k(x, y) = x \wedge y$  and inner product

$$\langle f, g \rangle := \int_0^1 f'(u)g'(u) du.$$

We can check

$$\langle f, k(\cdot, x) \rangle = \int_0^1 f'(u) \mathbb{1}_{[0, x]}(u) du = \int_0^x f'(u) du = f(x).$$

**Second-order Sobolev kernel.** Take  $\mathcal{H}$  to be the set of differentiable functions  $f : [0, 1] \rightarrow \mathbb{R}$  with  $f(0) = 0, f'(0) = 0$  and where  $f'$  is almost everywhere differentiable with  $\int_0^1 (f''(u))^2 du < \infty$ . For  $f, g \in \mathcal{H}$  define

$$\langle f, g \rangle := \int_0^1 f''(u)g''(u) du.$$

Recall (see Ex. sheet) that  $k_x(y) := k(x, y) := \int_0^{x \wedge y} (x-u)(y-u) du$  satisfies  $k_x''(y) = (x-y)_+$  and for  $f \in \mathcal{H}$ ,

$$\begin{aligned} \langle f, k_x \rangle &= \int_0^1 f''(u)(x-u)_+ du \\ &= [f'(u)(x-u)_+]_0^1 + \int_0^1 f'(u) \mathbb{1}_{[0, x]}(u) du \\ &= f(x). \end{aligned}$$

### 1.3.3 The representer theorem

To recap, what we have shown so far is that replacing the matrix  $XX^\top$  in the definition of an algorithm by  $K$  derived from a positive definite kernel is essentially equivalent to running the same algorithm on some mapping of the original data, though with the modification that instances of  $x_i^\top x_j$  become  $\langle \phi(x_i), \phi(x_j) \rangle$ . This corresponds to ridge regression on a predictor matrix with  $i$ th row  $\phi(x_i)$  in the case where  $\phi : \mathcal{X} \rightarrow \mathcal{H}$  maps into a Euclidean space, but the question remains as to how to interpret this when  $\mathcal{H}$  is infinite-dimensional.

If  $\mathcal{H}$  denotes the RKHS of the linear kernel, then

$$\hat{f} := \arg \min_{f \in \mathcal{H}} \left\{ \sum_{i=1}^n \{Y_i - f(x_i)\}^2 + \lambda \|f\|_{\mathcal{H}}^2 \right\} \quad (1.14)$$

is the usual fitted regression function from ridge regression. The following theorem shows in particular that kernel ridge regression (i.e. ridge regression replacing  $XX^\top$  with  $K$ ) with kernel  $k$  is equivalent to the above with  $\mathcal{H}$  now being the RKHS corresponding to  $k$ .

**Theorem 6** (Representer theorem, [Kimeldorf and Wahba, 1970, Schölkopf et al., 2001]).  
Let  $c : \mathbb{R}^n \times \mathcal{X}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$  be an arbitrary loss function, and let  $J : [0, \infty) \rightarrow \mathbb{R}$  be strictly increasing. Let  $x_1, \dots, x_n \in \mathcal{X}$ ,  $Y \in \mathbb{R}^n$ . Finally, let  $f \in \mathcal{H}$  where  $\mathcal{H}$  is an RKHS with reproducing kernel  $k$ , and let  $K_{ij} = k(x_i, x_j)$   $i, j = 1, \dots, n$ . Then  $\hat{f}$  minimises

$$Q_1(f) := c(Y, x_1, \dots, x_n, f(x_1), \dots, f(x_n)) + J(\|f\|_{\mathcal{H}}^2)$$

over  $f \in \mathcal{H}$  iff.  $\hat{f}(\cdot) = \sum_{i=1}^n \hat{\alpha}_i k(\cdot, x_i)$  and  $\hat{\alpha} \in \mathbb{R}^n$  minimises  $Q_2$  over  $\alpha \in \mathbb{R}^n$  where

$$Q_2(\alpha) = c(Y, x_1, \dots, x_n, K\alpha) + J(\alpha^\top K\alpha).$$

*Proof.* Let  $U$  be the linear span of  $\{k(\cdot, x_i), i = 1, \dots, n\}$ . As  $U$  is finite-dimensional and hence closed, we can decompose any  $f \in \mathcal{H}$  as  $f = u + v$  with  $u \in U$  and  $v \in U^\perp$ . Then

$$f(x_i) = \langle f, k(\cdot, x_i) \rangle = \langle u + v, k(\cdot, x_i) \rangle = \langle u, k(\cdot, x_i) \rangle = u(x_i).$$

Meanwhile,

$$J(\|f\|_{\mathcal{H}}^2) = J(\|u\|_{\mathcal{H}}^2 + \|v\|_{\mathcal{H}}^2) \geq J(\|u\|_{\mathcal{H}}^2),$$

with equality if and only  $v = 0$ . Thus in minimising  $Q_1$ , we may restrict attention to those  $f \in U$ , i.e., those

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

for some  $\alpha \in \mathbb{R}^n$ . But for such  $f$ ,

$$\|f\|_{\mathcal{H}}^2 = \left\langle \sum_{i=1}^n \alpha_i k(\cdot, x_i), \sum_{i=1}^n \alpha_i k(\cdot, x_i) \right\rangle = \alpha^\top K\alpha,$$

and  $f(x_i) = K_i^\top \alpha$  so  $(f(x_i))_{i=1}^n = K\alpha$ . Thus  $Q_1(f) = Q_2(\alpha)$ .  $\square$

Consider the result specialised the ridge regression objective. We see that (1.14) is essentially equivalent to minimising

$$\|Y - K\alpha\|_2^2 + \lambda \alpha^\top K\alpha,$$

and you may check (see example sheet 1) that the minimiser  $\hat{\alpha}$  satisfies  $K\hat{\alpha} = K(K + \lambda I)^{-1}Y$ . Thus (1.14) is indeed an alternative way of expressing kernel ridge regression. The result also tells us how to form predictions: given a new observation  $x$ , our prediction for  $f(x)$  is

$$\hat{f}(x) = \sum_{i=1}^n \hat{\alpha}_i k(x, x_i).$$

The application of the result is however not limited to ridge regression and shows that a whole host of algorithms can be ‘kernelised’. For example, recall that in the classification setting where  $Y_i \in \{-1, 1\}$ , standard logistic regression may be motivated by assuming

$$\log \left( \frac{\mathbb{P}(Y_i = 1)}{\mathbb{P}(Y_i = -1)} \right) = \mu^0 + x_i^\top \beta^0$$

and picking  $(\hat{\mu}, \hat{\beta})$  to maximise the log-likelihood. This leads to the following optimisation problem:

$$\arg \min_{(\mu, \beta) \in \mathbb{R}^p} \sum_{i=1}^n \log\{1 + \exp(-Y_i(\mu + x_i^\top \beta))\}.$$

The kernelised version is then given by

$$\arg \min_{\mu \in \mathbb{R}, f \in \mathcal{H}} \left\{ \sum_{i=1}^n \log[1 + \exp\{-Y_i(\mu + f(x_i))\}] + \lambda \|f\|_{\mathcal{H}}^2 \right\},$$

where  $\mathcal{H}$  is an RKHS (note that here we have included an unpenalised intercept term).

## 1.4 Kernel ridge regression

We have seen how the kernel trick allows us to solve a potentially infinite-dimensional version of ridge regression. This may seem impressive, but ultimately we should judge kernel ridge regression on its statistical properties e.g. predictive performance. Consider a setting where

$$Y_i = f^0(x_i) + \varepsilon_i, \quad \mathbb{E}(\varepsilon) = 0, \quad \text{Var}(\varepsilon) = \sigma^2 I.$$

We shall assume that  $f^0 \in \mathcal{H}$  where  $\mathcal{H}$  is an RKHS with reproducing kernel  $k$ . Let  $K$  be the kernel matrix  $K_{ij} = k(x_i, x_j)$  with eigenvalues  $d_1 \geq d_2 \geq \dots \geq d_n \geq 0$ . We will see that the predictive performance depends delicately on these eigenvalues.

Let  $\hat{f}_\lambda$  be the estimated regression function from kernel ridge regression with kernel  $k$ :

$$\hat{f}_\lambda = \arg \min_{f \in \mathcal{H}} \left\{ \sum_{i=1}^n \{Y_i - f(x_i)\}^2 + \lambda \|f\|_{\mathcal{H}}^2 \right\}.$$

**Theorem 7.** *The mean squared prediction error (MSPE) may be bounded above in the following way:*

$$\begin{aligned} \frac{1}{n} \mathbb{E} \left\{ \sum_{i=1}^n \{f^0(x_i) - \hat{f}_\lambda(x_i)\}^2 \right\} &\leq \frac{\sigma^2}{n} \sum_{i=1}^n \frac{d_i^2}{(d_i + \lambda)^2} + \frac{\lambda \|f^0\|_{\mathcal{H}}^2}{4n} \\ &\leq \frac{\sigma^2}{n} \frac{1}{\lambda} \sum_{i=1}^n \min(d_i/4, \lambda) + \frac{\lambda \|f^0\|_{\mathcal{H}}^2}{4n}. \end{aligned} \tag{1.15}$$

*Proof.* We know from the representer theorem that

$$\left( \hat{f}_\lambda(x_1), \dots, \hat{f}_\lambda(x_n) \right)^\top = K(K + \lambda I)^{-1} Y.$$

You will show on the example sheet that

$$\left( f^0(x_1), \dots, f^0(x_n) \right)^\top = K\alpha,$$

for some  $\alpha \in \mathbb{R}^n$ , and moreover that  $\|f^0\|_{\mathcal{H}}^2 \geq \alpha^\top K \alpha$ . Let the eigendecomposition of  $K$  be given by  $K = UDU^\top$  with  $D_{ii} = d_i$  and define  $\theta = U^\top K \alpha$ . We see that  $n$  times the LHS of (1.15) is

$$\begin{aligned} \mathbb{E}\|K(K + \lambda I)^{-1}(U\theta + \varepsilon) - U\theta\|_2^2 &= \mathbb{E}\|DU^\top(UDU^\top + \lambda I)^{-1}(U\theta + \varepsilon) - \theta\|_2^2 \\ &= \mathbb{E}\|D(D + \lambda I)^{-1}(\theta + U^\top \varepsilon) - \theta\|_2^2 \\ &= \|\{D(D + \lambda I)^{-1} - I\}\theta\|_2^2 + \mathbb{E}\|D(D + \lambda I)^{-1}U^\top \varepsilon\|_2^2. \end{aligned}$$

To compute the second term, we use the ‘trace trick’:

$$\begin{aligned} \mathbb{E}\|D(D + \lambda I)^{-1}U^\top \varepsilon\|_2^2 &= \mathbb{E}[\{D(D + \lambda I)^{-1}U^\top \varepsilon\}^\top D(D + \lambda I)^{-1}U^\top \varepsilon] \\ &= \mathbb{E}[\text{tr}\{D(D + \lambda I)^{-1}U^\top \varepsilon \varepsilon^\top U D(D + \lambda I)^{-1}\}] \\ &= \sigma^2 \text{tr}\{D(D + \lambda I)^{-1}D(D + \lambda I)^{-1}\} \\ &= \sigma^2 \sum_{i=1}^n \frac{d_i^2}{(d_i + \lambda)^2}. \end{aligned}$$

For the first term, we have

$$\|\{D(D + \lambda I)^{-1} - I\}\theta\|_2^2 = \sum_{i=1}^n \frac{\lambda^2 \theta_i^2}{(d_i + \lambda)^2}.$$

Now as  $\theta = DU^\top \alpha$ , note that  $\theta_i = 0$  when  $d_i = 0$ . Let  $D^+$  be the diagonal matrix with  $i$ th diagonal entry equal to  $D_{ii}^{-1}$  if  $D_{ii} > 0$  and 0 otherwise. Then

$$\sum_{i:d_i>0} \frac{\theta_i^2}{d_i} = \|\sqrt{D^+}\theta\|_2^2 = \alpha^\top K U D^+ U^\top K \alpha = \alpha^\top U D D^+ D U^\top \alpha = \alpha^\top K \alpha \leq \|f^0\|_{\mathcal{H}}^2.$$

By Hölder’s inequality we have

$$\sum_{i=1}^n \frac{\lambda^2 \theta_i^2}{(d_i + \lambda)^2} = \sum_{i:d_i>0} \frac{\theta_i^2}{d_i} \frac{d_i \lambda^2}{(d_i + \lambda)^2} \leq \max_{i=1,\dots,n} \frac{d_i \lambda^2}{(d_i + \lambda)^2} \leq \|f^0\|_{\mathcal{H}}^2 \lambda / 4,$$

using the inequality  $(a + b)^2 \geq 4ab$  in the final line. Finally note that

$$\frac{d_i^2}{(d_i + \lambda)^2} \leq \min\{1, d_i^2/(4d_i \lambda)\} = \min(\lambda, d_i/4)/\lambda. \quad \square$$

To interpret this result further, it will be helpful to express it in terms of  $\hat{\mu}_i := d_i/n$  (the eigenvalues of  $K/n$ ) and  $\gamma_n := \lambda/n$ . We have

$$\frac{1}{n} \mathbb{E} \left\{ \sum_{i=1}^n \{f^0(x_i) - \hat{f}_{n\gamma_n}(x_i)\}^2 \right\} \leq \frac{\sigma^2}{\gamma_n} \frac{1}{n} \sum_{i=1}^n \min(\hat{\mu}_i/4, \gamma_n) + \|f^0\|_{\mathcal{H}}^2 \gamma_n / 4 =: \delta_n(\gamma). \quad (1.16)$$

Here we have treated the  $x_i$  as fixed, but we could equally well think of them as random. Consider a setup where the  $x_i$  are i.i.d. and independent of  $\varepsilon$ . If we take a further expectation on the RHS of (1.16), our result still holds true (the  $\hat{\mu}_i$  are random in this setting). Ideally we would like to then replace  $\mathbb{E} \min(\hat{\mu}_i/4, \gamma)$  with a quantity more directly related to the kernel  $k$ .

Mercer's theorem is helpful in this regard. This guarantees (under some mild conditions) an eigendecomposition for kernels, which recall are somewhat like infinite-dimensional analogues of symmetric positive semi-definite matrices.

Given a random variable  $X$  taking values in  $\mathcal{X}$ , we say a non-zero function  $e \in \mathcal{H}$  is an *eigenfunction* with *eigenvalue*  $\mu \in \mathbb{R}$  if

$$\mu e(x) = \mathbb{E} k(x, X) e(X).$$

Mercer's theorem states the following under mild conditions, including that  $\mathbb{E} k(X, X) < \infty$ :

- The set of positive eigenvalues is at most countable.
- The subspace spanned by the eigenfunctions corresponding to each positive eigenvalue has a finite dimension known as the *multiplicity* of the eigenvalue.
- Writing  $(\mu_j)_{j \in J}$  (where  $J = \{1, \dots, m\}$ , some  $m$  or  $J \in \mathbb{N}$ ) for the eigenvalues counted with multiplicity, there exist corresponding eigenfunctions  $(e_j)_{j \in J}$  that are orthonormal in the sense that

$$\mathbb{E} e_j(X) e_k(X) = \mathbb{1}_{\{j=k\}}$$

and satisfy

$$k(x, y) = \sum_{j \in J} \mu_j e_j(x) e_j(y). \quad (1.17)$$

**Lemma 8.** *When (1.17) holds, we have for  $\gamma > 0$ ,*

$$\mathbb{E} \left( \frac{1}{n} \sum_{i=1}^n \min(\hat{\mu}_i/4, \gamma) \right) \leq \frac{1}{n} \sum_{j \in J} \min(\mu_j/4, \gamma).$$

**Theorem 9.** *Provided the eigendecomposition (1.17) holds, there exists  $\gamma_n$  such that for fixed  $\sigma^2 > 0$ ,*

$$\frac{1}{n} \mathbb{E} \left\{ \sum_{i=1}^n \{f^0(x_i) - \hat{f}_{\gamma_n}(x_i)\}^2 \right\} = o(n^{-1/2}).$$

*Proof.* Let  $\phi : [0, \infty) \rightarrow [0, \infty)$  be given by

$$\phi(\gamma) := \sum_{j \in J} \min(\mu_j, \gamma).$$

Observe that  $\phi$  is increasing and as  $\sum_{j \in J} \mu_j < \infty$ ,  $\lim_{\gamma \downarrow 0} \phi(\gamma) = 0$  (this is clear when  $J$  is finite; otherwise note that given an arbitrary  $\epsilon > 0$ , there exists  $M$  such that  $\sum_{j=M}^{\infty} \mu_j \leq \epsilon$ , but then  $\phi(\gamma) \leq M\gamma + \epsilon \rightarrow \epsilon$  as  $\gamma \downarrow 0$ ). Let  $\gamma_n = n^{-1/2} \sqrt{\phi(n^{-1/2})}$  so  $\gamma_n = o(n^{-1/2})$ . Thus for  $n$  sufficiently large  $\phi(\gamma_n) \leq \phi(n^{-1/2})$ , whence for such  $n$  we have

$$\begin{aligned} \inf_{\gamma > 0} \{\phi(\gamma)/(n\gamma) + \gamma\} &\leq \frac{\phi(\gamma_n)}{n\gamma_n} + \gamma_n \\ &\leq 2\sqrt{\phi(n^{-1/2})/\sqrt{n}} = o(n^{-1/2}). \end{aligned} \quad \square$$

In specific cases, we can get faster rates.

**First-order Sobolev kernel.** When  $k$  is the Sobolev kernel, and considering a uniform distribution on  $\mathcal{X} = [0, 1]$ , an eigenvalue–eigenfunction pair  $(\mu, e)$  must satisfy

$$\mu e(x) = \int_0^1 \min(x, y) e(y) dy = \int_0^x y e(y) dy + x \int_x^1 e(y) dy,$$

so

$$\mu e'(x) = x e(x) + \int_x^1 e(y) dy - x e(x) = \int_x^1 e(y) dy, \quad (1.18)$$

hence

$$\mu e''(x) = -e(x).$$

This ODE has general solution  $e(x) = a_\mu \sin(x/\sqrt{\mu}) + b_\mu \cos(x/\sqrt{\mu})$  and the boundary condition  $e(0) = 0$  gives  $b_\mu = 0$ . Also from (1.18), we see that  $\mu e'(1) = 0$ , so  $1/\sqrt{\mu} = \pi/2 + k\pi$  for some  $k = 0, 1, 2, \dots$ . Thus the  $j$ th eigenvalue satisfies

$$\mu_j/4 = \frac{1}{\pi^2(2j-1)^2}.$$

We therefore have

$$\begin{aligned} \sum_{i=1}^{\infty} \min(\mu_i/4, \gamma_n) &\leq \frac{\gamma_n}{2} \left( \frac{1}{\sqrt{\pi^2 \gamma_n}} + 1 \right) + \frac{1}{\pi^2} \int_{\{(\pi^2 \gamma_n)^{-1/2} + 1\}/2}^{\infty} \frac{1}{(2x-1)^2} dx \\ &= \sqrt{\gamma_n}/\pi + \gamma_n/2 = O(\sqrt{\gamma_n}) \end{aligned}$$

as  $\gamma_n \rightarrow 0$ . Putting things together, we see that

$$\mathbb{E} \delta_n(\gamma_n) = O\left( \frac{\sigma^2}{n\gamma_n^{1/2}} + \gamma_n \right).$$

Thus an optimal  $\gamma_n \sim (\sigma^2/n)^{2/3}$  gives an error rate of order  $(\sigma^2/n)^{2/3}$ .

## 1.5 Large-scale kernel machines

We introduced the kernel trick as a computational device that avoided performing calculations in a high or infinite dimensional feature space and, in the case of kernel ridge regression reduced computation down to forming the  $n \times n$  matrix  $K$  and then inverting  $K + \lambda I$ . This can be a huge saving, but when  $n$  is very large, this can present serious computational difficulties. Even if  $p$  is small, the  $O(n^3)$  cost of inverting  $K + \lambda I$  may cause problems. What's worse, the fitted regression function is a sum over  $n$  terms:

$$\hat{f}(\cdot) = \sum_{i=1}^n \hat{\alpha}_i k(x_i, \cdot).$$

Even to evaluate a prediction at a single new observation requires  $O(n)$  computations unless  $\hat{\alpha}$  is sparse.

In recent years, there has been great interest in speeding up computations for kernel machines. We will discuss one exciting approach based on random feature expansions. Given a kernel  $k$ , the key idea is to develop a random map

$$\hat{\phi} : \mathcal{X} \rightarrow \mathbb{R}^b$$

with  $b$  small such that  $\mathbb{E}\{\hat{\phi}(x)^\top \hat{\phi}(x')\} = k(x, x')$ . In a sense we are trying to reverse the kernel trick by approximating the kernel using a random feature map. To increase the quality of the approximation of the kernel, we can consider

$$x \mapsto \frac{1}{\sqrt{L}}(\hat{\phi}_1(x), \dots, \hat{\phi}_L(x)) \in \mathbb{R}^{Lb}$$

with each  $(\hat{\phi}_l(x))_{l=1}^L$  being i.i.d. for each  $x$ . Let  $\Phi$  be the matrix with  $i$ th row given by  $(\hat{\phi}_1(x_i), \dots, \hat{\phi}_L(x_i))/\sqrt{L}$ . We may then run our learning algorithm replacing the initial matrix of predictors  $X$  with  $\Phi$ . For example, when performing ridge regression, we can compute

$$(\Phi^\top \Phi + \lambda I)^{-1} \Phi^\top Y,$$

which would require  $O(nL^2b^2 + L^3b^3)$  operations: a cost linear in  $n$ . Predicting a new observation would cost  $O(Lb)$ .

The work of Rahimi and Recht [2007] proposes a construction of such a random mapping  $\hat{\phi}$  for shift-invariant kernels, that is kernels for which there exists a function  $g$  with  $k(x, x') = g(x - x')$  for all  $x, x' \in \mathcal{X} = \mathbb{R}^p$ . A useful property of such kernels is given by Bochner's theorem.

**Theorem 10** (Bochner's theorem). *Let  $k : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$  be a continuous kernel. Then  $k$  is shift-invariant if and only if there exists some  $c > 0$  and distribution  $F$  on  $\mathbb{R}^p$  such that when  $W \sim F$*

$$k(x, x') = c\mathbb{E}e^{i(x-x')^\top W} = c\mathbb{E}\cos((x-x')^\top W).$$

To make use of this theorem, first observe the following. Let  $u \sim U[-\pi, \pi]$ ,  $x, y \in \mathbb{R}$ . Then

$$2\mathbb{E} \cos(x+u) \cos(y+u) = 2\mathbb{E}\{(\cos x \cos u - \sin x \sin u)(\cos y \cos u - \sin y \sin u)\}.$$

Now as  $u \stackrel{d}{=} -u$ ,  $\mathbb{E} \cos u \sin u = \mathbb{E} \cos(-u) \sin(-u) = -\mathbb{E} \cos u \sin u = 0$ . Also of course  $\cos^2 u + \sin^2 u = 1$  so  $\mathbb{E} \cos^2 u = \mathbb{E} \sin^2 u = 1/2$ . Thus

$$2\mathbb{E} \cos(x+u) \cos(y+u) = \cos x \cos y + \sin x \sin y = \cos(x-y).$$

Given a shift-invariant kernel  $k$  with associated distribution  $F$ , suppose  $W \sim F$  and let  $u \sim U[-\pi, \pi]$  independently. Define

$$\hat{\phi}(x) = \sqrt{2c} \cos(W^\top x + u).$$

Then

$$\begin{aligned} \mathbb{E} \hat{\phi}(x) \hat{\phi}(x') &= 2c \mathbb{E}[\mathbb{E}\{\cos(W^\top x + u) \cos(W^\top x' + u) | W\}] \\ &= c \mathbb{E} \cos((x - x')^\top W) = k(x, x'). \end{aligned}$$

As a concrete example of this approach, let us take the Gaussian kernel  $k(x, x') = \exp\{-\|x - x'\|_2^2 / (2h^2)\}$ . Note that if  $W \sim N(0, h^{-2}I)$ , it has characteristic function  $\mathbb{E}(e^{it^\top W}) = e^{-\|t\|_2^2 / (2h^2)}$  so we may take  $\hat{\phi}(x) = \sqrt{2} \cos(W^\top x + u)$ .



# Chapter 2

## The Lasso and extensions

### 2.1 Model selection

Let us revisit the linear model  $Y = X\beta^0 + \varepsilon$  where  $\mathbb{E}(\varepsilon) = 0$ ,  $\text{Var}(\varepsilon) = \sigma^2 I$ . In many modern datasets, there are reasons to believe there are many more variables present than are necessary to explain the response. Let  $S$  be the set  $S = \{k : \beta_k^0 \neq 0\}$  and suppose  $s := |S| \ll p$ .

The MSPE of OLS is

$$\begin{aligned} \frac{1}{n} \mathbb{E} \|X\beta^0 - X\hat{\beta}^{\text{OLS}}\|_2^2 &= \frac{1}{n} \mathbb{E} \{(\beta^0 - \hat{\beta}^{\text{OLS}})^\top X^\top X (\beta^0 - \hat{\beta}^{\text{OLS}})\} \\ &= \frac{1}{n} \mathbb{E} [\text{tr}\{(\beta^0 - \hat{\beta}^{\text{OLS}})(\beta^0 - \hat{\beta}^{\text{OLS}})^\top X^\top X\}] \\ &= \frac{1}{n} \text{tr}[\mathbb{E}\{(\beta^0 - \hat{\beta}^{\text{OLS}})(\beta^0 - \hat{\beta}^{\text{OLS}})^\top\} X^\top X] \\ &= \frac{1}{n} \text{tr}(\text{Var}(\hat{\beta}^{\text{OLS}}) X^\top X) = \frac{p}{n} \sigma^2. \end{aligned}$$

If we could identify  $S$  and then fit a linear model using just these variables, we'd obtain an MSPE of  $\sigma^2 s/n$  which could be substantially smaller than  $\sigma^2 p/n$ . Furthermore, it can be shown that parameter estimates from the reduced model are more accurate. The smaller model would also be easier to interpret.

We now briefly review some classical model selection strategies.

#### Best subset regression

A natural approach to finding  $S$  is to consider all  $2^p$  possible regression procedures each involving regressing the response on a different sets of explanatory variables  $X_M$  where  $M$  is a subset of  $\{1, \dots, p\}$ . We can then pick the best regression procedure using cross-validation (say). For general design matrices, this involves an exhaustive search over all subsets, so this is not really feasible for  $p > 50$ .

## Forward selection

This can be seen as a greedy way of performing best subsets regression. Given a target model size  $m$  (the tuning parameter), this works as follows.

1. Start by fitting an intercept only model.
2. Add to the current model the predictor variable that reduces the residual sum of squares the most.
3. Continue step 2 until  $m$  predictor variables have been selected.

## 2.2 The Lasso estimator

The *Least absolute shrinkage and selection operator (Lasso)* [Tibshirani, 1996] estimates  $\beta^0$  by  $\hat{\beta}_\lambda^L$ , where  $(\hat{\mu}_\lambda^L, \hat{\beta}_\lambda^L)$  minimise

$$\frac{1}{2n} \|Y - \mu \mathbf{1} - X\beta\|_2^2 + \lambda \|\beta\|_1 \quad (2.1)$$

over  $(\mu, \beta) \in \mathbb{R} \times \mathbb{R}^p$ . Here  $\|\beta\|_1$  is the  $\ell_1$ -norm of  $\beta$ :  $\|\beta\|_1 = \sum_{k=1}^p |\beta_k|$ .

Like ridge regression,  $\hat{\beta}_\lambda^L$  shrinks the OLS estimate towards the origin, but there is an important difference. The  $\ell_1$  penalty can force some of the estimated coefficients to be exactly 0. In this way the Lasso can perform simultaneous variable selection and parameter estimation. As we did with ridge regression, we can centre and scale the  $X$  matrix, so then  $\hat{\mu}_\lambda^L = \bar{Y}$ . As before, our target of interest is  $\beta^0$ . Define

$$Q_\lambda(\beta) = \frac{1}{2n} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_1. \quad (2.2)$$

Now the minimiser(s) of  $Q_\lambda(\beta)$  will also be the minimiser(s) of

$$\|Y - X\beta\|_2^2 \text{ subject to } \|\beta\|_1 \leq \|\hat{\beta}_\lambda^L\|_1.$$

Similarly, with the Ridge regression objective, we know that  $\hat{\beta}_\lambda^R$  minimises  $\|Y - X\beta\|_2^2$  subject to  $\|\beta\|_2 \leq \|\hat{\beta}_\lambda^R\|_2$ .

Now the contours of the OLS objective  $\|Y - X\beta\|_2^2$  are ellipsoids centred at  $\hat{\beta}^{\text{OLS}}$ , while the contours of  $\|\beta\|_2^2$  are spheres centred at the origin, and the contours of  $\|\beta\|_1$  are ‘diamonds’ centred at 0.

The important point to note is that the  $\ell_1$  ball  $\{\beta \in \mathbb{R}^p : \|\beta\|_1 \leq \|\hat{\beta}_\lambda^L\|_1\}$  has corners where some of the components are zero, and it is likely that the OLS contours will intersect the  $\ell_1$  ball at such a corner.

### 2.2.1 Prediction error of the Lasso (slow rate)

A remarkable property of the Lasso is that even when  $p \gg n$ , it can still perform well in terms of prediction error. Suppose the columns of  $X$  have been centred and scaled (as we will always assume from now on unless stated otherwise) and assume the normal linear model

$$Y = \mu \mathbf{1} + X\beta^0 + \varepsilon \quad (2.3)$$

where  $\varepsilon \sim N_n(0, \sigma^2 I)$ .

**Theorem 11.** *Let  $\hat{\beta}$  be any Lasso solution when*

$$\lambda = A\sigma\sqrt{\frac{\log(p)}{n}}.$$

*With probability at least  $1 - 2p^{-(A^2/2-1)}$*

$$\frac{1}{n}\|X(\beta^0 - \hat{\beta})\|_2^2 \leq 4A\sigma\sqrt{\frac{\log(p)}{n}}\|\beta^0\|_1.$$

*Proof.* From the definition of  $\hat{\beta}$  we have

$$\frac{1}{2n}\|Y - \bar{Y}\mathbf{1} - X\hat{\beta}\|_2^2 + \lambda\|\hat{\beta}\|_1 \leq \frac{1}{2n}\|Y - \bar{Y}\mathbf{1} - X\beta^0\|_2^2 + \lambda\|\beta^0\|_1.$$

Rearranging,

$$\frac{1}{2n}\|X(\beta^0 - \hat{\beta})\|_2^2 \leq \frac{1}{n}\varepsilon^\top X(\hat{\beta} - \beta^0) + \lambda\|\beta^0\|_1 - \lambda\|\hat{\beta}\|_1.$$

Now  $|\varepsilon^\top X(\hat{\beta} - \beta^0)| \leq \|X^\top \varepsilon\|_\infty \|\hat{\beta} - \beta^0\|_1$ . Let  $\Omega = \{\|X^\top \varepsilon\|_\infty/n \leq \lambda\}$ . Lemma 15 below shows that  $\mathbb{P}(\Omega) \geq 1 - 2p^{-(A^2/2-1)}$ . Working on the event  $\Omega$ , we obtain

$$\begin{aligned} \frac{1}{2n}\|X(\beta^0 - \hat{\beta})\|_2^2 &\leq \lambda\|\beta^0 - \hat{\beta}\|_1 + \lambda\|\beta^0\|_1 - \lambda\|\hat{\beta}\|_1, \\ \frac{1}{n}\|X(\beta^0 - \hat{\beta})\|_2^2 &\leq 4\lambda\|\beta^0\|_1, \quad \text{by the triangle inequality.} \end{aligned} \quad \square$$

### 2.2.2 Concentration inequalities I

The proof of Theorem 11 relies on a lower bound for the probability of the event  $\Omega$ . A union bound gives

$$\begin{aligned} \mathbb{P}(\|X^\top \varepsilon\|_\infty/n > \lambda) &= \mathbb{P}(\cup_{j=1}^p |X_j^\top \varepsilon|/n > \lambda) \\ &\leq \sum_{j=1}^p \mathbb{P}(|X_j^\top \varepsilon|/n > \lambda). \end{aligned}$$

Now  $X_j^\top \varepsilon/n \sim N(0, \sigma^2/n)$ , so if we obtain a bound on the tail probabilities of normal distributions, the argument above will give a bound for  $\mathbb{P}(\Omega)$ .

Motivated by the need to bound normal tail probabilities, we will briefly discuss the topic of *concentration inequalities* that provide such bounds for much wider classes of random variables. Concentration inequalities are vital for the study of many modern algorithms and in our case here, they will reveal that the attractive properties of the Lasso presented in Theorem 11 hold true for a variety of non-normal errors.

We begin our discussion with the simplest tail bound, *Markov's inequality*, which states that given a non-negative random variable  $W$ ,

$$\mathbb{P}(W \geq t) \leq \frac{\mathbb{E}(W)}{t}.$$

This immediately implies that given a strictly increasing function  $\varphi : \mathbb{R} \rightarrow [0, \infty)$  and any random variable  $W$ ,

$$\mathbb{P}(W \geq t) = \mathbb{P}\{\varphi(W) \geq \varphi(t)\} \leq \frac{\mathbb{E}(\varphi(W))}{\varphi(t)}.$$

Applying this with  $\varphi(t) = e^{\alpha t}$  ( $\alpha > 0$ ) yields the so-called *Chernoff bound*:

$$\mathbb{P}(W \geq t) \leq \inf_{\alpha > 0} e^{-\alpha t} \mathbb{E}e^{\alpha W}.$$

Consider the case when  $W \sim N(0, \sigma^2)$ . Recall that

$$\mathbb{E}e^{\alpha W} = e^{\alpha^2 \sigma^2 / 2}. \tag{2.4}$$

Thus

$$\mathbb{P}(W \geq t) \leq \inf_{\alpha > 0} e^{\alpha^2 \sigma^2 / 2 - \alpha t} = e^{-t^2 / (2\sigma^2)}.$$

Note that to arrive at this bound, all we required was (an upper bound on) the moment generating function (mgf) of  $W$  (2.4).

### Sub-Gaussian variables

**Definition 3.** We say a random variable  $W$  is *sub-Gaussian* if there exists  $\sigma > 0$  such that

$$\mathbb{E}e^{\alpha(W - \mathbb{E}W)} \leq e^{\alpha^2 \sigma^2 / 2}$$

for all  $\alpha \in \mathbb{R}$ . We then say that  $W$  is *sub-Gaussian with parameter  $\sigma$* .

**Proposition 12** (Sub-Gaussian tail bound). *If  $W$  is sub-Gaussian with parameter  $\sigma$  then*

$$\mathbb{P}(W - \mathbb{E}W \geq t) \leq e^{-t^2 / (2\sigma^2)}.$$

As well as Gaussian random variables, the sub-Gaussian class includes bounded random variables.

**Lemma 13** (Hoeffding's lemma). *If  $W$  takes values in  $[a, b]$ , then  $W$  is sub-Gaussian with parameter  $(b - a)/2$ .*

The following proposition shows that analogously to how a linear combination of jointly Gaussian random variables is Gaussian, a linear combination of sub-Gaussian random variables is also sub-Gaussian.

**Proposition 14.** *Let  $(W_i)_{i=1}^n$  be a sequence of independent sub-Gaussian random variables with parameters  $(\sigma_i)_{i=1}^n$  and let  $\gamma \in \mathbb{R}^n$ . Then  $\gamma^\top W$  is sub-Gaussian with parameter  $\left(\sum_i \gamma_i^2 \sigma_i^2\right)^{1/2}$ .*

*Proof.* Wlog, we may assume  $\mathbb{E}W_i = 0$  for all  $i$ . We have

$$\begin{aligned} \mathbb{E} \exp\left(\alpha \sum_{i=1}^n \gamma_i W_i\right) &= \prod_{i=1}^n \mathbb{E} \exp(\alpha \gamma_i W_i) \\ &\leq \prod_{i=1}^n \exp(\alpha^2 \gamma_i^2 \sigma_i^2 / 2) \\ &= \exp\left(\alpha^2 \sum_{i=1}^n \gamma_i^2 \sigma_i^2 / 2\right). \end{aligned} \quad \square$$

We can now prove a more general version of the probability bound required for Theorem 11.

**Lemma 15.** *Suppose  $(\varepsilon_i)_{i=1}^n$  are independent, mean-zero and sub-Gaussian with common parameter  $\sigma$ . Note that this includes  $\varepsilon \sim N_n(0, \sigma^2 I)$ . Let  $\lambda = A\sigma\sqrt{\log(p)/n}$ . Then*

$$\mathbb{P}(\|X^\top \varepsilon\|_\infty / n \leq \lambda) \geq 1 - 2p^{-(A^2/2-1)}.$$

*Proof.*

$$\mathbb{P}(\|X^\top \varepsilon\|_\infty / n > \lambda) \leq \sum_{j=1}^p \mathbb{P}(|X_j^\top \varepsilon| / n > \lambda).$$

But  $\pm X_j^\top \varepsilon / n$  are both sub-Gaussian with parameter  $(\sigma^2 \|X_j\|_2^2 / n^2)^{1/2} = \sigma / \sqrt{n}$ . Thus the RHS is at most

$$2p \exp(-A^2 \log(p)/2) = 2p^{1-A^2/2}. \quad \square$$

### 2.2.3 Some facts from optimisation theory and convex analysis

In order to study the Lasso in detail, it will be helpful to review some basic facts from optimisation and convex analysis.

## Convexity

A set  $C \subseteq \mathbb{R}^d$  is *convex* if

$$x, y \in C \Rightarrow (1 - t)x + ty \in C \quad \text{for all } t \in (0, 1).$$

Given  $C \subseteq \mathbb{R}^d$ , we say a function  $f : C \rightarrow \mathbb{R}$  is *convex* if  $C$  is convex and

$$f((1 - t)x + ty) \leq (1 - t)f(x) + tf(y)$$

for all  $x, y \in C$  and  $t \in (0, 1)$ . It is *strictly convex* if the inequality is strict for all  $x, y \in C$  with  $x \neq y$ . If a strictly convex function has a minimiser, it must be unique. In the following,  $C \subseteq \mathbb{R}^d$  is a convex set.

**Proposition 16.** (i) Let  $f_1, \dots, f_m : C \rightarrow \mathbb{R}$  be convex functions. Then if  $c_1, \dots, c_m \geq 0$ ,  $c_1 f_1 + \dots + c_m f_m : C \rightarrow \mathbb{R}$  is a convex function.

(ii) If  $f : C \rightarrow \mathbb{R}$ , and  $A : \mathbb{R}^m \rightarrow \mathbb{R}^d$  is an affine function (so  $A(x) = Mx + b$  for  $M \in \mathbb{R}^{d \times m}$  and  $b \in \mathbb{R}^d$ ) then  $g : D \rightarrow \mathbb{R}$ , where  $D = \{x \in \mathbb{R}^m : A(x) \in C\}$  given by  $g(x) = f(A(x))$  is convex.

(iii) If  $f : C \rightarrow \mathbb{R}$  is convex with  $C$  open and  $f$  is twice continuously differentiable on  $C$ , then

- (a)  $f$  is convex iff. its Hessian  $H(x)$  is positive semi-definite for all  $x \in C$ ,
- (b)  $f$  is strictly convex if  $H(x)$  is positive definite for all  $x \in C$ .

## The Lagrangian method

Consider an optimisation problem of the form

$$\text{minimise } f(x), \text{ subject to } g(x) = 0, \quad x \in C \subseteq \mathbb{R}^d, \quad (2.5)$$

where  $g : C \rightarrow \mathbb{R}^b$ . Suppose the optimal value is  $c^* \in \mathbb{R}$ . The Lagrangian for this problem is defined as

$$L(x, \theta) = f(x) + \theta^\top g(x)$$

where  $\theta \in \mathbb{R}^b$ . Note that

$$\inf_{x \in C} L(x, \theta) \leq \inf_{x \in C: g(x)=0} L(x, \theta) = c^*$$

for all  $\theta$ . The Lagrangian method involves finding a  $\theta = \theta^*$  such that the minimising  $x = x^*$  on the LHS satisfies  $g(x^*) = 0$ . This  $x^*$  must then be a minimiser in the original problem (2.5).

## Subgradients

**Definition 4.** Given convex  $C \subseteq \mathbb{R}^d$ , a vector  $v \in \mathbb{R}^d$  is a *subgradient* of a convex function  $f : C \rightarrow \mathbb{R}$  at  $x$  if

$$f(y) \geq f(x) + v^\top(y - x) \quad \text{for all } y \in C.$$

The set of subgradients of  $f$  at  $x$  is called the *subdifferential* of  $f$  at  $x$  and denoted  $\partial f(x)$ .

In order to make use of subgradients, we will require the following two facts:

**Proposition 17.** Let  $f : C \rightarrow \mathbb{R}$  be convex, and suppose  $f$  is differentiable at  $x \in \text{int}(C)$ . Then  $\partial f(x) = \{\nabla f(x)\}$ .

**Proposition 18** (Subgradient calculus). Let  $f, f_1, f_2 : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex. Then

$$(i) \quad \partial(\alpha f)(x) = \{\alpha g : g \in \partial f(x)\} \text{ for } \alpha > 0,$$

$$(ii) \quad \partial(f_1 + f_2)(x) = \{g_1 + g_2 : g_1 \in \partial f_1(x), g_2 \in \partial f_2(x)\}.$$

Also if  $h : \mathbb{R}^m \rightarrow \mathbb{R}$  is given by  $h(x) = f(Ax + b)$  where  $A \in \mathbb{R}^{d \times m}$  and  $b \in \mathbb{R}^d$ , then

$$(iii) \quad \partial h(x) = \{A^\top g : g \in \partial f(Ax + b)\}.$$

The following easy (but key) result is often referred to in the statistical literature as the Karush–Kuhn–Tucker (KKT) conditions, though it is actually a much simplified version of them.

**Proposition 19.** Given convex  $f : C \rightarrow \mathbb{R}$ ,  $x^* \in \arg \min_{x \in C} f(x)$  if and only if  $0 \in \partial f(x^*)$ .

*Proof.*

$$\begin{aligned} f(y) \geq f(x^*) \quad \text{for all } y \in C &\Leftrightarrow f(y) \geq f(x^*) + 0^\top(y - x) \quad \text{for all } y \in C \\ &\Leftrightarrow 0 \in \partial f(x^*). \end{aligned} \quad \square$$

Let us now compute the subdifferential of the  $\ell_1$ -norm. First note that  $\|\cdot\|_1 : \mathbb{R}^d \rightarrow \mathbb{R}$  is convex. Indeed it is a norm so the triangle inequality gives  $\|tx + (1-t)y\|_1 \leq t\|x\|_1 + (1-t)\|y\|_1$ . We introduce some notation that will be helpful here and throughout the rest of the course.

For  $x \in \mathbb{R}^d$  and  $A = \{k_1, \dots, k_m\} \subseteq \{1, \dots, d\}$  with  $k_1 < \dots < k_m$ , by  $x_A$  we will mean  $(x_{k_1}, \dots, x_{k_m})^\top$ . Similarly if  $X$  has  $d$  columns we will write  $X_A$  for the matrix

$$X_A = (X_{k_1} \cdots X_{k_m}).$$

Further in this context, by  $A^c$ , we will mean  $\{1, \dots, d\} \setminus A$ . Additionally, when in subscripts we will use the shorthand  $-j = \{j\}^c$  and  $-jk = \{j, k\}^c$ . Note these column and component

extraction operations will always be considered to have taken place first before any further operations on the matrix, so for example  $X_A^\top = (X_A)^\top$ . Finally, define

$$\text{sgn}(x_1) = \begin{cases} -1 & \text{if } x_1 < 0 \\ 0 & \text{if } x_1 = 0 \\ 1 & \text{if } x_1 > 0, \end{cases}$$

and

$$\text{sgn}(x) = (\text{sgn}(x_1), \dots, \text{sgn}(x_d))^\top.$$

**Proposition 20.** *For  $x \in \mathbb{R}^d$  let  $A = \{j : x_j \neq 0\}$ . Then*

$$\partial\|x\|_1 = \{v \in \mathbb{R}^d : \|v\|_\infty \leq 1 \text{ and } v_A = \text{sgn}(x_A)\}$$

*Proof.* We can write  $\|x\|_1 = \sum_{j=1}^d |e_j^\top x|$  where  $e_j$  is the  $j$ th standard basis vector. The subdifferential of the function  $u \mapsto |u|$  is  $[-1, 1]$  if  $u = 0$  and  $\{\text{sgn}(u)\}$  otherwise. Thus from Proposition 18(c) we see that the subdifferential of  $g_j(x) =: |e_j^\top x|$  is

$$\partial g_j(x) = \begin{cases} \{\text{sgn}(x_j)e_j\} & \text{if } x_j \neq 0 \\ \{te_j : t \in [-1, 1]\} & \text{otherwise.} \end{cases}$$

Proposition 18(b) then gives the result. □

## 2.2.4 Lasso solutions

Equipped with these tools from convex analysis, we can now fully characterise the solutions to the Lasso. We have that  $\hat{\beta}_\lambda^L$  is a Lasso solution if and only if  $0 \in \partial Q_\lambda(\hat{\beta}_\lambda^L)$ , which is equivalent to

$$\frac{1}{n}X^\top(Y - X\hat{\beta}_\lambda^L) = \lambda\hat{\nu},$$

for  $\hat{\nu}$  with  $\|\hat{\nu}\|_\infty \leq 1$  and writing  $\hat{S}_\lambda = \{k : \hat{\beta}_{\lambda,k}^L \neq 0\}$ ,  $\hat{\nu}_{\hat{S}_\lambda} = \text{sgn}(\hat{\beta}_{\lambda,\hat{S}_\lambda}^L)$ .

Lasso solutions need not be unique (e.g. if  $X$  has duplicate columns), though for most reasonable design matrices, Lasso solutions will be unique. We will often tacitly assume Lasso solutions are unique in the statement of our theoretical results. It is however straightforward to show that the Lasso fitted values are unique.

**Proposition 21.** *Fix  $\lambda \geq 0$  and suppose  $\beta^{(1)}$  and  $\beta^{(2)}$  are two Lasso solutions. Then  $X\beta^{(1)} = X\beta^{(2)}$ .*

*Proof.* Suppose  $\beta^{(1)}$  and  $\beta^{(2)}$  both give an optimal objective value of  $c^*$ . Now by strict convexity of  $\|\cdot\|_2^2$ ,

$$\|Y - X\beta^{(1)}/2 - X\beta^{(2)}/2\|_2^2 \leq \|Y - X\beta^{(1)}\|_2^2/2 + \|Y - X\beta^{(2)}\|_2^2/2,$$



with equality if and only if  $X\beta^{(1)} = X\beta^{(2)}$ . Since  $\|\cdot\|_1$  is also convex, we see that

$$\begin{aligned}
c^* &\leq Q_\lambda(\beta^{(1)}/2 + \beta^{(2)}/2) \\
&= \|Y - X\beta^{(1)}/2 - X\beta^{(2)}/2\|_2^2/(2n) + \lambda\|\beta^{(1)}/2 + \beta^{(2)}/2\|_1 \\
&\leq \|Y - X\beta^{(1)}\|_2^2/(4n) + \|Y - X\beta^{(2)}\|_2^2/(4n) + \lambda\|\beta^{(1)}/2 + \beta^{(2)}/2\|_1 \\
&\leq \{\|Y - X\beta^{(1)}\|_2^2/(4n) + \lambda\|\beta^{(1)}\|_1/2\} + \{\|Y - X\beta^{(2)}\|_2^2/(4n) + \lambda\|\beta^{(2)}\|_1/2\} \\
&= Q(\beta^{(1)})/2 + Q(\beta^{(2)})/2 = c^*.
\end{aligned}$$

Equality must prevail throughout this chain of inequalities, so  $X\beta^{(1)} = X\beta^{(2)}$ .  $\square$

Define the *equicorrelation set*  $\hat{E}_\lambda$  to be the set of  $k$  such that

$$\frac{1}{n}|X_k^\top(Y - X\hat{\beta}_\lambda^L)| = \lambda.$$

Note that  $\hat{E}_\lambda$  is well-defined since it only depends on the fitted values, which (as we have just shown) are unique. By the KKT conditions, the equicorrelation set contains the set of non-zeroes of all Lasso solutions. Note that if  $\text{rank}(X_{\hat{E}_\lambda}) = |\hat{E}_\lambda|$  then the Lasso solution must be unique: indeed if  $\beta^{(1)}$  and  $\beta^{(2)}$  are two Lasso solutions, then as

$$X_{\hat{E}_\lambda}(\beta_{\hat{E}_\lambda}^{(1)} - \beta_{\hat{E}_\lambda}^{(2)}) = 0,$$

by linear independence of the columns of  $X_{\hat{E}_\lambda}$ , we have  $\beta_{\hat{E}_\lambda}^{(1)} = \beta_{\hat{E}_\lambda}^{(2)}$ .

### 2.2.5 Variable selection

Consider now the “noiseless” version of the high-dimensional linear model (2.3),  $Y = X\beta^0$ . The case with noise can be dealt with by similar arguments to those we will use below when we work on an event that  $\|X^\top \varepsilon\|_\infty/n$  is small (see example sheet).

Let  $S = \{k : \beta_k^0 \neq 0\}$ ,  $N = \{1, \dots, p\} \setminus S$  and assume wlog that  $S = \{1, \dots, s\}$ , and also that  $\text{rank}(X_S) = s$ .

**Theorem 22.** *Let  $\lambda > 0$  and define  $\Delta = X_N^\top X_S (X_S^\top X_S)^{-1} \text{sgn}(\beta_S^0)$ . If  $\|\Delta\|_\infty \leq 1$  and for  $k \in S$ ,*

$$|\beta_k^0| > \lambda |\text{sgn}(\beta_S^0)^\top [\{\frac{1}{n} X_S^\top X_S\}^{-1}]_k|, \quad (2.6)$$

*then there exists a Lasso solution  $\hat{\beta}_\lambda^L$  with  $\text{sgn}(\hat{\beta}_\lambda^L) = \text{sgn}(\beta^0)$ . As a partial converse, if there exists a Lasso solution  $\hat{\beta}_\lambda^L$  with  $\text{sgn}(\hat{\beta}_\lambda^L) = \text{sgn}(\beta^0)$ , then  $\|\Delta\|_\infty \leq 1$ .*

*Remark 1.* We can interpret  $\|\Delta\|_\infty$  as the maximum in absolute value over  $k \in N$  of the dot product of  $\text{sgn}(\beta_S^0)$  and  $(X_S^\top X_S)^{-1} X_S^\top X_k$ , the coefficient vector obtained by regressing  $X_k$  on  $X_S$ . The condition  $\|\Delta\|_\infty \leq 1$  is known as the irrepresentable condition.

*Proof.* Fix  $\lambda > 0$  and write  $\hat{\beta} = \hat{\beta}_\lambda^L$  and  $\hat{S} = \{k : \hat{\beta}_k \neq 0\}$  for convenience. The KKT conditions for the Lasso give

$$\frac{1}{n} X^\top X (\beta^0 - \hat{\beta}) = \lambda \hat{\nu}$$

where  $\|\hat{\nu}\|_\infty \leq 1$  and  $\hat{\nu}_{\hat{S}} = \text{sgn}(\hat{\beta}_{\hat{S}})$ . We can expand this into

$$\frac{1}{n} \begin{pmatrix} X_S^\top X_S & X_S^\top X_N \\ X_N^\top X_S & X_N^\top X_N \end{pmatrix} \begin{pmatrix} \beta_S^0 - \hat{\beta}_S \\ -\hat{\beta}_N \end{pmatrix} = \lambda \begin{pmatrix} \hat{\nu}_S \\ \hat{\nu}_N \end{pmatrix}. \quad (2.7)$$

We prove the converse first. If  $\text{sgn}(\hat{\beta}) = \text{sgn}(\beta^0)$  then  $\hat{\nu}_S = \text{sgn}(\beta_S^0)$  and  $\hat{\beta}_N = 0$ . The top block of (2.7) gives

$$\beta_S^0 - \hat{\beta}_S = \lambda \left( \frac{1}{n} X_S^\top X_S \right)^{-1} \text{sgn}(\beta_S^0).$$

Substituting this into the bottom block, we get

$$\lambda \frac{1}{n} X_N^\top X_S \left( \frac{1}{n} X_S^\top X_S \right)^{-1} \text{sgn}(\beta_S^0) = \lambda \hat{\nu}_N.$$

Thus as  $\|\hat{\nu}_N\|_\infty \leq 1$ , we have  $\|\Delta\|_\infty \leq 1$ .

For the positive statement, we need to find a  $\hat{\beta}$  and  $\hat{\nu}$  such that  $\text{sgn}(\hat{\beta}_S) = \text{sgn}(\beta_S^0)$  and  $\hat{\beta}_N = 0$ , for which the KKT conditions hold. We claim that taking

$$\begin{aligned} (\hat{\beta}_S, \hat{\beta}_N) &= (\beta_S^0 - \lambda \left( \frac{1}{n} X_S^\top X_S \right)^{-1} \text{sgn}(\beta_S^0), 0) \\ (\hat{\nu}_S, \hat{\nu}_N) &= (\text{sgn}(\beta_S^0), \Delta) \end{aligned}$$

satisfies (2.7). We only need to check that  $\text{sgn}(\beta_S^0) = \text{sgn}(\hat{\beta}_S)$ , but this follows from (2.6).  $\square$

## 2.2.6 Prediction and estimation

Consider once more the model  $Y = \mu^0 \mathbf{1} + X\beta^0 + \varepsilon$  where the components of  $\varepsilon$  are independent mean-zero sub-Gaussian random variables with common parameter  $\sigma$ . Let  $S$ ,  $s$  and  $N$  be defined as in the previous section. As we have noted before, in an artificial situation where  $S$  is known, we could apply OLS on  $X_S$  and have an MSPE of  $\sigma^2 s/n$ . Under a so-called *compatibility condition* on the design matrix, we can obtain a similar MSPE for the Lasso.

**Definition 5.** Given a matrix of predictors  $X \in \mathbb{R}^{n \times p}$  and support set  $S \neq \emptyset$ , define the *compatibility factor*

$$\phi^2 = \inf_{\delta \in \mathbb{R}^p: \delta_S \neq 0, \|\delta_N\|_1 \leq 3\|\delta_S\|_1} \frac{\frac{1}{n} \|X\delta\|_2^2}{\frac{1}{s} \|\delta_S\|_1^2},$$

where  $s = |S|$  and we take  $\phi \geq 0$ . The *compatibility condition* is that  $\phi^2 > 0$ .

Note that if  $X^\top X/n$  has minimum eigenvalue  $c_{\min} > 0$  (so necessarily  $p \leq n$ ), then  $\phi^2 > c_{\min}$ . Indeed by the Cauchy–Schwarz inequality,

$$\|\delta_S\|_1 = \text{sgn}(\delta_S)^\top \delta_S \leq \sqrt{s} \|\delta_S\|_2 \leq \sqrt{s} \|\delta\|_2.$$

Thus

$$\phi^2 \geq \inf_{\delta \neq 0} \frac{\frac{1}{n} \|X\delta\|_2^2}{\|\delta\|_2^2} = c_{\min}.$$

Although in the high-dimensional setting we would have  $c_{\min} = 0$ , the fact that the infimum in the definition of  $\phi^2$  is over a restricted set of  $\delta$  can still allow  $\phi^2$  to be positive even in this case, as we discuss after the presentation of the theorem.

**Theorem 23.** *Suppose the compatibility condition holds and let  $\hat{\beta}$  be a Lasso solution with tuning parameter  $\lambda = A\sigma\sqrt{\log(p)/n}$  for  $A > 2\sqrt{2}$ . Then with probability at least  $1 - 2p^{-(A^2/8-1)}$ , we have that for all  $\lambda \geq \lambda^*$ ,*

$$\frac{1}{n} \|X(\beta^0 - \hat{\beta})\|_2^2 + \lambda \|\beta^0 - \hat{\beta}_\lambda^L\|_1 \leq \frac{12\lambda^2 s}{\phi^2}.$$

In particular,

$$\frac{1}{n} \|X(\beta^0 - \tilde{\beta})\|_2^2 \leq \frac{12A^2 \log(p)}{\phi^2} \frac{\sigma^2 s}{n}, \quad \text{and} \quad \|\beta^0 - \tilde{\beta}\|_1 \leq \frac{12A\sigma s}{\phi^2} \sqrt{\frac{\log p}{n}}.$$

*Proof.* As in Theorem 11 we start with the “basic inequality”:

$$\frac{1}{2n} \|X(\hat{\beta} - \beta^0)\|_2^2 + \lambda \|\hat{\beta}\|_1 \leq \frac{1}{n} \varepsilon^\top X(\hat{\beta} - \beta^0) + \lambda \|\beta^0\|_1.$$

We work on the event  $\Omega = \{2\|X^\top \varepsilon\|_\infty/n \leq \lambda\}$  where after applying Hölder’s inequality, we get

$$\frac{1}{n} \|X(\hat{\beta} - \beta^0)\|_2^2 + 2\lambda \|\hat{\beta}\|_1 \leq \lambda \|\hat{\beta} - \beta^0\|_1 + 2\lambda \|\beta^0\|_1. \quad (2.8)$$

Lemma 15 shows that  $\mathbb{P}(\Omega) \geq 1 - 2p^{-(A^2/8-1)}$ .

To motivate the rest of the proof, consider the following idea. We know

$$\frac{1}{n} \|X(\hat{\beta} - \beta^0)\|_2^2 \leq 3\lambda \|\hat{\beta} - \beta^0\|_1.$$

If we could obtain

$$3\lambda \|\hat{\beta} - \beta^0\|_1 \leq \frac{c\lambda}{\sqrt{n}} \|X(\hat{\beta} - \beta^0)\|_2$$

for some constant  $c > 0$ , then we would have that  $\|X(\hat{\beta} - \beta^0)\|_2^2/n \leq c^2\lambda^2$  and also  $3\lambda \|\beta^0 - \hat{\beta}\|_1 \leq c^2\lambda^2$ .

Returning to the actual proof, write  $a = \|X(\hat{\beta} - \beta^0)\|_2^2/(n\lambda)$ . Then from (2.8) we can derive the following string of inequalities:

$$\begin{aligned}
a + 2(\|\hat{\beta}_N\|_1 + \|\hat{\beta}_S\|_1) &\leq \|\hat{\beta}_S - \beta_S^0\|_1 + \|\hat{\beta}_N\|_1 + 2\|\beta_S^0\|_1 \\
a + \|\hat{\beta}_N\|_1 &\leq \|\hat{\beta}_S - \beta_S^0\|_1 + 2\|\beta_S^0\|_1 - 2\|\hat{\beta}_S\|_1 \\
a + \|\hat{\beta}_N - \beta_N^0\|_1 &\leq 3\|\beta_S^0 - \hat{\beta}_S\|_1 \\
&\leq \frac{3}{\phi} \sqrt{\frac{s}{n}} \|X(\hat{\beta} - \beta^0)\|_2,
\end{aligned} \tag{2.9}$$

the final inequality coming from using the compatibility condition with  $\delta = \hat{\beta} - \beta^0$ . Thus

$$\frac{1}{n} \|X(\hat{\beta} - \beta^0)\|_2^2 \leq \frac{3\lambda}{\phi} \|X(\hat{\beta} - \beta^0)\|_2,$$

so

$$\frac{1}{\sqrt{n}} \|X(\hat{\beta} - \beta^0)\|_2 \leq \frac{3\lambda\sqrt{s}}{\phi}.$$

Substituting this into (2.9), we obtain

$$\frac{1}{n} \|X(\hat{\beta} - \beta^0)\|_2^2 + \lambda \|\hat{\beta}_N - \beta_N^0\|_1 \leq \frac{9\lambda^2 s}{\phi^2}.$$

Finally adding the inequality  $\lambda \|\hat{\beta}_S - \beta_S^0\|_1 \leq 3\lambda^2 s/\phi^2$ , we get the result.  $\square$

## 2.2.7 The compatibility condition

How strong is the compatibility condition? In order to answer this question, we shall think of  $X$  as random and try to understand what conditions on the population covariance matrix  $\Sigma^0 := \mathbb{E}(X^\top X/n)$  imply that  $X$  satisfies a compatibility condition with high probability. To this end let us define

$$\phi_\Sigma^2 = \inf_{\delta: \|\delta_S\|_1 \neq 0, \|\delta_N\|_1 \leq 3\|\delta_S\|_1} \frac{\delta^\top \Sigma \delta}{\|\delta_S\|_1^2/s} = s \inf_{\delta: \|\delta_S\|_1 = 1, \|\delta_N\|_1 \leq 3} \delta^\top \Sigma \delta,$$

where  $\Sigma \in \mathbb{R}^{p \times p}$ . Note then our  $\phi^2 = \phi_\Sigma^2$  where  $\hat{\Sigma} := X^\top X/n$ . The following result shows that if  $\Sigma$  is close to a matrix  $\Theta$  for which  $\phi_\Theta^2 > 0$ , then also  $\phi_\Sigma^2 > 0$ .

**Lemma 24.** *Suppose  $\phi_\Theta^2 > 0$  and  $\max_{jk} |\Sigma_{jk} - \Theta_{jk}| \leq \phi_\Theta^2/(32s)$ . Then  $\phi_\Sigma^2 \geq \phi_\Theta^2/2$ .*

*Proof.* Let  $\mathcal{B} := \{\delta : \|\delta_S\|_1 = 1, \|\delta_N\|_1 \leq 3\}$ . Take  $\delta \in \mathcal{B}$ . Then we have

$$s\delta^\top \Sigma \delta = s\delta^\top \Theta \delta - s\delta^\top (\Theta - \Sigma) \delta \geq \phi_\Theta^2 - s|\delta^\top (\Sigma - \Theta) \delta|.$$

Furthermore,

$$\begin{aligned} |\delta^\top (\Theta - \Sigma) \delta| &\leq \|\delta\|_1 \|(\Theta - \Sigma) \delta\|_\infty \quad (\text{Hölder}) \\ &\leq \frac{\phi_\Theta^2}{32s} \|\delta\|_1^2 \quad (\text{Hölder again}) \end{aligned}$$

and  $\|\delta\|_1 = \|\delta_N\|_1 + \|\delta_S\|_1 \leq 4$ . Thus

$$s\delta^\top \Sigma \delta \geq \phi_\Theta^2 - \phi_\Theta^2/2 = \phi_\Theta^2/2.$$

Taking the infimum over  $\delta \in \mathcal{B}$  gives the result.  $\square$

We would like to apply the result above with  $\Theta = \Sigma^0$ , and use it to argue that if  $\Sigma^0$  satisfies the compatibility condition, then so will  $\hat{\Sigma}$  with high probability. In order to do this, we need to argue that the event that  $\max_{jk} |\hat{\Sigma}_{jk} - \Sigma_{jk}^0|$  is small occurs with high probability. We can obtain such a result with the aid of concentration inequalities.

### 2.2.8 Concentration inequalities II

When trying to understand the concentration properties of  $\hat{\Sigma}_{jk}$ , it will be helpful to have a tail bound for a product of sub-Gaussian random variables. Bernstein's inequality, which applies to random variables satisfying the condition below, is helpful in this regard.

**Definition 6** (Bernstein's condition). We say that the random variable  $W$  satisfies Bernstein's condition with parameter  $(\sigma, b)$  where  $\sigma, b > 0$  if

$$\mathbb{E}(|W - \mathbb{E}W|^k) \leq \frac{1}{2} k! \sigma^2 b^{k-2} \quad \text{for } k = 2, 3, \dots$$

**Proposition 25** (Bernstein's inequality). *Let  $W_1, W_2, \dots$  be independent random variables with  $\mathbb{E}(W_i) = \mu$ . Suppose each  $W_i$  satisfies Bernstein's condition with parameter  $(\sigma, b)$ . Then*

$$\begin{aligned} \mathbb{E}(e^{\alpha(W_i - \mu)}) &\leq \exp\left(\frac{\alpha^2 \sigma^2 / 2}{1 - b|\alpha|}\right) \quad \text{for all } |\alpha| < 1/b \\ \mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n W_i - \mu \geq t\right) &\leq \exp\left(-\frac{nt^2}{2(\sigma^2 + bt)}\right) \quad \text{for all } t > 0. \end{aligned}$$

*Proof.* Fix  $i$  and let  $W = W_i$ . We have

$$\begin{aligned}
\mathbb{E}(e^{\alpha(W-\mu)}) &= \mathbb{E}\left(1 + \alpha(W - \mu) + \sum_{k=2}^{\infty} \frac{\alpha^k (W - \mu)^k}{k!}\right) \\
&\leq \mathbb{E}\left(1 + \sum_{k=2}^{\infty} \frac{|\alpha|^k |W - \mu|^k}{k!}\right) \\
&= 1 + \sum_{k=2}^{\infty} |\alpha|^k \frac{\mathbb{E}(|W - \mu|^k)}{k!} \\
&\leq 1 + \frac{\sigma^2 \alpha^2}{2} \sum_{k=2}^{\infty} |\alpha|^{k-2} b^{k-2} \\
&= 1 + \frac{\sigma^2 \alpha^2}{2} \frac{1}{1 - |\alpha|b} \leq \exp\left(\frac{\alpha^2 \sigma^2 / 2}{1 - b|\alpha|}\right),
\end{aligned}$$

provided  $|\alpha| < 1/b$  and using the inequality  $e^u \geq 1 + u$  in the final line. For the probability bound, first note that

$$\begin{aligned}
\mathbb{E} \exp\left(\sum_{i=1}^n \alpha(W_i - \mu)/n\right) &= \prod_{i=1}^n \mathbb{E} \exp\{\alpha(W_i - \mu)/n\} \\
&\leq \exp\left(n \frac{(\alpha/n)^2 \sigma^2 / 2}{1 - b|\alpha/n|}\right)
\end{aligned}$$

for  $|\alpha|/n < 1/b$ . Then we use the Chernoff method, though without minimising over  $\alpha > 0$ : instead we set  $\alpha/n = t/(bt + \sigma^2) \in (0, 1/b)$ .  $\square$

**Lemma 26.** *Let  $W, Z$  be mean-zero and sub-Gaussian with parameters  $\sigma_W$  and  $\sigma_Z$  respectively. Then the product  $WZ$  satisfies Bernstein's condition with parameter  $(8\sigma_W\sigma_Z, 4\sigma_W\sigma_Z)$ .*

*Proof.* In order to use Bernstein's inequality (Proposition 25) we first obtain bounds on the moments of  $W$  and  $Z$ . Note that  $W^{2k} = \int_0^\infty \mathbb{1}_{\{x < W^{2k}\}} dx$ . Thus by Fubini's theorem

$$\begin{aligned}
\mathbb{E}(W^{2k}) &= \int_0^\infty \mathbb{P}(W^{2k} > x) dx \\
&= 2k \int_0^\infty t^{2k-1} \mathbb{P}(|W| > t) dt \quad \text{substituting } t^{2k} = x \\
&\leq 4k \int_0^\infty t^{2k-1} \exp\{-t^2/(2\sigma_W^2)\} dt \quad \text{by Proposition 12} \\
&= 4k\sigma_W^2 \int_0^\infty (2\sigma_W^2 x)^{k-1} e^{-x} dx \quad \text{substituting } t^2/(2\sigma_W^2) = x \\
&= 2^{k+1} \sigma_W^{2k} k!.
\end{aligned}$$

Next note that for any random variable  $Y$ ,

$$\begin{aligned}\mathbb{E}|Y - \mathbb{E}Y|^k &= 2^k \mathbb{E}|Y/2 - \mathbb{E}Y/2|^k \\ &\leq 2^{k-1}(\mathbb{E}|Y|^k + |\mathbb{E}Y|^k) \quad \text{by Jensen's inequality applied to } t \mapsto |t|^k, \\ &\leq 2^k \mathbb{E}|Y|^k.\end{aligned}$$

Therefore

$$\begin{aligned}\mathbb{E}(|WZ - \mathbb{E}WZ|^k) &\leq 2^k \mathbb{E}|WZ|^k \\ &\leq 2^k (\mathbb{E}W^{2k})^{1/2} (\mathbb{E}Z^{2k})^{1/2} \quad \text{by Cauchy-Schwarz} \\ &\leq 2^k 2^{k+1} \sigma_W^k \sigma_Z^k k! \\ &= \frac{k!}{2} (8\sigma_W \sigma_Z)^2 (4\sigma_W \sigma_Z)^{k-2}.\end{aligned} \quad \square$$

### 2.2.9 Random design

We now show that we can expect the compatibility condition to hold with high probability. To make the result more readily interpretable, we shall state it in an asymptotic framework. Imagine a sequence of design matrices with  $n$ ,  $p$  and  $s$  (the size of the set  $S$ ) growing, each with their own compatibility condition. We will however suppress the asymptotic regime in the notation.

**Theorem 27.** *Suppose the rows of  $X$  are i.i.d. and each entry of  $X$  is mean-zero sub-Gaussian with parameter  $v$ . Let  $\hat{\Sigma} := X^\top X/n$  and  $\Sigma^0 := \mathbb{E}(\hat{\Sigma})$ . Suppose  $s\sqrt{\log(p)/n} \rightarrow 0$  (and  $s, p, n > 1$ ) as  $n \rightarrow \infty$ . Suppose  $\phi_{\Sigma^0}^2 > c$  for a constant  $c > 0$ . Then  $\mathbb{P}(\phi_{\hat{\Sigma}}^2 \geq \phi_{\Sigma^0}^2/2) \rightarrow 1$  as  $n \rightarrow \infty$ .*

*Proof.* In view of Lemma 24, we need only show that

$$\mathbb{P}(\max_{jk} |\hat{\Sigma}_{jk} - \Sigma_{jk}^0| \geq \phi_{\Sigma^0}^2/(32s)) \rightarrow 0.$$

Let  $t := \phi_{\Sigma^0, s}^2/(32s)$ . By a union bound and then Lemma 26 we have

$$\begin{aligned}\mathbb{P}(\max_{jk} |\hat{\Sigma}_{jk} - \Sigma_{jk}^0| \geq t) &< p^2 \max_{jk} \mathbb{P}\left(\left|\sum_{i=1}^n X_{ij} X_{ik}/n - \Sigma_{jk}^0\right| \geq t\right) \\ &\leq 2 \exp\left(-\frac{nt^2}{2(64v^4 + 4v^2t)} + 2\log p\right) \\ &\leq 2 \exp(-c'n/s^2 + 2\log p)\end{aligned} \quad (2.10)$$

for a constant  $c' > 0$ . To justify the last line, observe that any constant  $C > 0$

$$\frac{(\phi_{\Sigma^0}^2)^2}{C + \phi_{\Sigma^0}^2/s} \geq \min_{u \geq c} \frac{u^2}{C + u} = \min_{u \geq c} u \left(1 - \frac{C}{u + C}\right) = c \left(1 - \frac{C}{c + C}\right) > 0.$$

Thus returning to (2.10), we see that this tends to zero as  $\log p = o(n/s^2)$ , which completes the proof.  $\square$

### 2.2.10 Computation

One of the most efficient ways of computing Lasso solutions is to use a optimisation technique called *coordinate descent*. This is a quite general way of minimising a function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  and works particularly well for functions of the form

$$f(x) = g(x) + \sum_{j=1}^d h_j(x_j)$$

where  $g$  is convex and differentiable and each  $h_j : \mathbb{R} \rightarrow \mathbb{R}$  is convex (and so continuous). We start with an initial guess of the minimiser  $x^{(0)}$  (e.g.  $x^{(0)} = 0$ ) and repeat for  $m = 1, 2, \dots$

$$\begin{aligned} x_1^{(m)} &= \arg \min_{x_1 \in \mathbb{R}} f(x_1, x_2^{(m-1)}, \dots, x_d^{(m-1)}) \\ x_2^{(m)} &= \arg \min_{x_2 \in \mathbb{R}} f(x_1^{(m)}, x_2, x_3^{(m-1)}, \dots, x_d^{(m-1)}) \\ &\vdots \\ x_d^{(m)} &= \arg \min_{x_d \in \mathbb{R}} f(x_1^{(m)}, x_2^{(m)}, \dots, x_{d-1}^{(m)}, x_d). \end{aligned}$$

Tseng [2001] proves that provided  $A_0 = \{x : f(x) \leq f(x^{(0)})\}$  is compact, then every converging subsequence of  $x^{(m)}$  will converge to a minimiser of  $f$ .

**Theorem 28.** *Suppose  $A^0$  is compact. Then*

- (i) *There exists a minimiser of  $f$ ,  $x^*$  and  $f(x^{(m)}) \rightarrow f(x^*)$ .*
- (ii) *If  $x^*$  is the unique minimiser of  $f$  then  $x^{(m)} \rightarrow x^*$ .*

*\*Proof\*.* Function  $f$  is continuous and so attains its infimum on the compact set  $A_0$ . Suppose  $f(x^{(m)}) \not\rightarrow f(x^*)$ . Then there exists  $\epsilon > 0$  and a subsequence  $(x^{(m_j)})_{j=0}^\infty$  such that  $f(x^{(m_j)}) \geq f(x^*) + \epsilon$  for all  $j$ . Note that since  $f(x^{(m)}) \leq f(x^{(m-1)})$ , we know that  $x^{(m)} \in A_0$  for all  $m$ . Thus if  $A_0$  is compact then any subsequence of  $(x^{(m)})_{m=0}^\infty$  has a further subsequence that converges by the Bolzano–Weierstrass theorem. Let  $\tilde{x}$  be the limit of the converging subsequence of  $(x^{(m_j)})_{j=0}^\infty$ . Then  $f(\tilde{x}) \geq f(x^*) + \epsilon$ , contradicting the result of Tseng [2001]. Thus (i) holds. The proof of (ii) is similar.  $\square$

We can replace individual coordinates by blocks of coordinates and the same result holds. That is if  $x = (x_1, \dots, x_B)$  where now  $x_b \in \mathbb{R}^{d_b}$  and

$$f(x) = g(x) + \sum_{b=1}^B h_b(x_b)$$

with  $g$  convex and differentiable and each  $h_b : \mathbb{R}^{d_b} \rightarrow \mathbb{R}$  convex, then block coordinate descent can be used.



One of the reasons that coordinate descent is so effective for solving the Lasso is that the coordinatewise optimisations are very simple and have closed form solutions. To see this, suppose at the  $m$ th iteration we are optimising for variable  $k$ . Let us write

$$R := Y - \sum_{j=1}^{k-1} X_j \hat{\beta}_j^{(m)} - \sum_{j=k+1}^p X_j \hat{\beta}_j^{(m-1)}.$$

We have that

$$\hat{\beta}_k^{(m)} = \arg \min_{\beta \in \mathbb{R}} \left\{ \frac{1}{2n} \|R - \beta X_k\|_2^2 + \lambda |\beta| \right\}$$

A minimiser  $\hat{\beta}_k^{(m)}$  is characterised by the subgradient optimality condition:

$$-\frac{1}{n} X_k^\top R + \hat{\beta}_k^{(m)} + \lambda \hat{\nu} = 0,$$

where  $\hat{\nu} \in [-1, 1]$  and if  $\hat{\beta}_k^{(m)} \neq 0$ ,  $\hat{\nu} = \text{sgn}(\hat{\beta}_k^{(m)})$ . Rearranging, we have

$$\hat{\beta}_k^{(m)} = \frac{1}{n} X_k^\top R - \lambda \hat{\nu},$$

and we may check that this is satisfied by

$$\hat{\beta}_k^{(m)} = S_\lambda(X_k^\top R/n),$$

where  $S_t(u) := \text{sgn}(u)(|u| - t)_+$  is the *soft-thresholding* operator. Note that  $\hat{\beta}_k^{(m)}$  is the unique minimiser as the coordinatewise objective is strictly convex.

We often want to solve the Lasso on a grid of  $\lambda$  values  $\lambda_0 > \dots > \lambda_L$  (for the purposes of cross-validation for example). To do this, we can first solve for  $\lambda_0$ , and then solve at subsequent grid points by using the solution at the previous grid points as an initial guess (known as a *warm start*). An active set strategy can further speed up computation. This works as follows: For  $l = 1, \dots, L$

1. Initialise  $A_l = \{k : \hat{\beta}_{\lambda_{l-1}, k}^L \neq 0\}$ .
2. Perform coordinate descent only on coordinates in  $A_l$  obtaining a solution  $\hat{\beta}$  (all components  $\hat{\beta}_k$  with  $k \notin A_l$  are set to zero).
3. Let  $V = \{k : |X_k^\top (Y - X\hat{\beta})|/n > \lambda_l\}$ , the set of coordinates which violate the KKT conditions when  $\hat{\beta}$  is taken as a candidate solution.
4. If  $V$  is empty, we set  $\hat{\beta}_{\lambda_l}^L = \hat{\beta}$ . Else we update  $A_l = A_l \cup V$  and return to 2.

## 2.3 Extensions of the Lasso

We can add an  $\ell_1$  penalty to many other log-likelihoods, or more generally other loss functions besides the squared-error loss that arises from the normal linear model. For Lasso-penalised generalised linear models, such as logistic regression, similar theoretical results to those we have obtained are available and computations can proceed in a similar fashion to above.

### 2.3.1 The square-root Lasso

Consider the normal linear model

$$Y = \mu^0 \mathbf{1} + X\beta^0 + \varepsilon, \quad (2.11)$$

where  $\varepsilon \sim N_n(0, \sigma^2 I)$ . A misgiving one might have about the theoretical results on the Lasso is that they rely on knowledge of the unknown  $\sigma$  in that the  $\lambda$  concerned takes the form  $A\sigma\sqrt{\log(p)/n}$ .

Now given a Lasso estimate  $\hat{\beta}_\lambda^L$  for  $\beta^0$ , a sensible estimate of  $\sigma$  is given by

$$\hat{\sigma}_\lambda^L := \frac{1}{\sqrt{n}} \|Y - \bar{Y}\mathbf{1} - X\hat{\beta}_\lambda^L\|_2.$$

Given this estimate, a sensible tuning parameter to choose for the Lasso would be  $A\hat{\sigma}_\lambda^L\sqrt{\log(p)/n}$ , which would lead to a new estimate of  $\beta^0$ , which could then in turn give a new estimate of  $\sigma$ . Iterating this process amounts to performing a block coordinate descent optimisation (alternating between optimising over  $\beta$  and  $\sigma$ ) of the following convex objective function,

$$Q_\gamma^{\text{sq}}(\beta, \sigma) := \frac{1}{2n\sigma} \|Y - \bar{Y}\mathbf{1} - X\beta\|_2^2 + \frac{\sigma}{2} + \gamma\|\beta\|_1,$$

with  $\gamma = A\sqrt{\log(p)/n}$  and initial value for  $\sigma$  given by  $\hat{\sigma}_\lambda^L$ . Theorem 28 indicates that this will lead to a minimiser of  $Q_\gamma^{\text{sq}}$ . However, a more direct route to the minimiser is offered by the so-called *square-root Lasso*  $\hat{\beta}_\gamma^{\text{sq}}$  [Belloni et al., 2011, Sun and Zhang, 2012], which minimises

$$\frac{1}{\sqrt{n}} \|Y - \bar{Y}\mathbf{1} - X\beta\|_2 + \gamma\|\beta\|_1. \quad (2.12)$$

Note that the display above is equal to  $\min_{\sigma>0} Q_\gamma^{\text{sq}}(\beta, \sigma)$  provided  $Y \neq X\beta$ .

### 2.3.2 Structural penalties

The Lasso penalty encourages the estimated coefficients to be shrunk towards 0 and sometimes exactly to 0. Other penalty functions can be constructed to encourage different types of sparsity.

#### Group Lasso

Suppose we have a partition  $G_1, \dots, G_q$  of  $\{1, \dots, p\}$  (so  $\cup_{k=1}^q G_k = \{1, \dots, p\}$ ,  $G_j \cap G_k = \emptyset$  for  $j \neq k$ ). The *group Lasso* penalty [Yuan and Lin, 2006] is given by

$$\lambda \sum_{j=1}^q m_j \|\beta_{G_j}\|_2.$$

The multipliers  $m_j > 0$  serve to balance cases where the groups are of very different sizes; typically we choose  $m_j = \sqrt{|G_j|}$ . This penalty encourages either an entire group  $G$  to have  $\hat{\beta}_G = 0$  or  $\hat{\beta}_k \neq 0$  for all  $k \in G$ . Such a property is useful when groups occur through coding for categorical predictors or when expanding predictors using basis functions.

## Fused Lasso

If there is a sense in which the coefficients are ordered, so  $\beta_j^0$  is expected to be close to  $\beta_{j+1}^0$ , a *fused Lasso* penalty [Tibshirani et al., 2005] may be appropriate. This takes the form

$$\lambda_1 \sum_{j=1}^{p-1} |\beta_j - \beta_{j+1}| + \lambda_2 \|\beta\|_1,$$

where the second term may be omitted depending on whether shrinkage towards 0 is desired. As an example, consider the simple setting where  $Y_i = \mu_i^0 + \varepsilon_i$ , and it is thought that the  $(\mu_i^0)_{i=1}^n$  form a piecewise constant sequence. Then one option is to minimise over  $\mu \in \mathbb{R}^n$ , the following objective

$$\frac{1}{n} \|Y - \mu\|_2^2 + \lambda \sum_{i=1}^{n-1} |\mu_i - \mu_{i+1}|.$$

### 2.3.3 Correlated predictors

When the matrix of predictors  $X$  has columns that are highly correlated, it can be difficult to distinguish the contributions of individual predictors to the response. The Lasso will tend to pick at most one the associated estimated coefficients to be non-zero, which can be undesirable when signal variables are highly correlated with noise variables.

To mitigate this, one option is first to cluster variables into groups of highly correlated variables, and then to apply a group Lasso, with groups given by these clusters. A simpler solution is to use the *elastic net*, which is a hybrid of ridge regression and the Lasso defined by a penalty of the form

$$\lambda \left\{ \alpha \|\beta\|_1 + \frac{1-\alpha}{2} \|\beta\|_2^2 \right\}.$$

Here  $\alpha \in [0, 1]$  is an additional tuning parameter controlling the relative contributions of the Lasso and ridge penalties. With  $\alpha \in (0, 1)$ , the elastic net can still deliver sparse estimates similarly to the Lasso, but will for example always give equal coefficient estimates to duplicate columns due to the presence of the (squared)  $\ell_2$ -norm penalty.

### 2.3.4 Reducing the bias of the Lasso

One potential drawback of the Lasso is that the same shrinkage effect that sets many estimated coefficients exactly to zero also shrinks all non-zero estimated coefficients towards zero. One possible solution is to take  $\hat{S}_\lambda = \{k : \hat{\beta}_{\lambda,k}^L \neq 0\}$  and then re-estimate  $\beta_{\hat{S}_\lambda}^0$  by OLS regression on  $X_{\hat{S}_\lambda}$ .

Another option is to re-estimate using the Lasso on  $X_{\hat{S}_\lambda}$ ; this procedure is known as the *relaxed Lasso* [Meinshausen, 2007]. The *adaptive Lasso* [Zou, 2006] takes an initial

estimate of  $\beta^0$ ,  $\hat{\beta}^{\text{init}}$  (e.g. from the Lasso) and then performs weighted Lasso regression:

$$\hat{\beta}_\lambda^{\text{adapt}} = \arg \min_{\beta \in \mathbb{R}^p: \beta_{\hat{S}_{\text{init}}^c} = 0} \left\{ \frac{1}{2n} \|Y - X\beta\|_2^2 + \lambda \sum_{k \in \hat{S}_{\text{init}}} \frac{|\beta_k|}{|\hat{\beta}_k^{\text{init}}|} \right\},$$

where  $\hat{S}_{\text{init}} = \{k : \hat{\beta}_k^{\text{init}} \neq 0\}$ .

Yet another approach involves using a family of non-convex penalty functions  $p_{\lambda, \gamma} : [0, \infty) \rightarrow [0, \infty)$  and attempting to minimise

$$\frac{1}{2n} \|Y - X\beta\|_2^2 + \sum_{k=1}^p p_{\lambda, \gamma}(|\beta_k|).$$

A prominent example is the *minimax concave penalty* (MCP) [Zhang, 2010] which takes

$$p'_\lambda(u) = \left( \lambda - \frac{u}{\gamma} \right)_+.$$

One disadvantage of using a non-convex penalty is that there may be multiple local minima which can make optimisation problematic. However, typically if the non-convexity is not too severe, coordinate descent can produce reasonable results.

## Chapter 3

# Graphical modelling and causal inference

So far, we have mainly considered the problem of relating a particular response to a potentially large collection of explanatory variables. In some settings, however, we do not have a distinguished response variable and instead we would like to better understand relationships between all measured variables. One simple way to formalise the relatedness between variables is to measure their correlation, or test their independence. However, this does not always lead to the most interpretable results as many pairs of variables may exhibit dependence without a very meaningful relationship between them. For example, height and literacy levels are highly positively correlated. While this may at first appear interesting or alarming, a little thought reveals that this fact is an expected consequence of babies not knowing how to read! If we were to look only at the literacy levels of those individuals of a given age  $a$ , then we would not expect to see such a relationship. The statistical property of *conditional independence*, defined below, captures this idea.

### 3.1 Conditional independence

**Definition 7.** Let  $X, Y, Z$  be random vectors. We say that  $X$  and  $Y$  are *conditionally independent* given  $Z$ , and write

$$X \perp\!\!\!\perp Y \mid Z,$$

if for all measurable sets  $A, B$ , we have

$$\mathbb{P}(X \in A, Y \in B \mid Z) = \mathbb{P}(X \in A \mid Z) \cdot \mathbb{P}(Y \in B \mid Z) \quad (3.1)$$

Otherwise, they are *conditionally dependent*, which we denote by  $X \not\perp\!\!\!\perp Y \mid Z$ .

**Proposition 29.** We have that  $X \perp\!\!\!\perp Y \mid Z$  if and only if for all measurable sets  $A$ ,

$$\mathbb{P}(X \in A \mid Y, Z) = \mathbb{P}(X \in A \mid Z).$$

The interpretation of  $X \perp\!\!\!\perp Y \mid Z$  is that ‘knowing  $Z$  renders  $X$  irrelevant for learning about  $Y$ ’. In the case where  $Z$  is deterministic, the definition of conditional independence reduces to that of (unconditional) independence.

The following useful rules concerning conditional independence may be deduced from its definition.

- (a) *Weak union*: If  $X \perp\!\!\!\perp Y, Z$ , then  $X \perp\!\!\!\perp Y \mid Z$  and  $X \perp\!\!\!\perp Z \mid Y$ .
- (b) *Contraction*: If  $X \perp\!\!\!\perp Z$  and  $X \perp\!\!\!\perp Y \mid Z$ , then  $X \perp\!\!\!\perp Y, Z$ .

These properties also hold when conditioning everywhere on an additional random variable  $W$ , so for example, the weak union property becomes

$$X \perp\!\!\!\perp Y, Z \mid W \implies \begin{cases} X \perp\!\!\!\perp Y \mid Z, W \\ X \perp\!\!\!\perp Z \mid Y, W. \end{cases}$$

A converse to the weak union property, known as the *intersection property* holds when  $X, Y, Z, W$  have a joint density with respect to a product measure, and  $f_{WYZ}(w, y, z) > 0$  for all  $w, y, z$ :

$$\left. \begin{array}{l} X \perp\!\!\!\perp Y \mid Z, W \\ X \perp\!\!\!\perp Z \mid Y, W \end{array} \right\} \implies X \perp\!\!\!\perp Y, Z \mid W.$$

Note also we always have  $X \perp\!\!\!\perp Y \mid Z \implies f(X, Z) \perp\!\!\!\perp g(Y, Z) \mid Z$  for all appropriate functions  $f$  and  $g$ .

## 3.2 Graphs

The properties above give several ways of deducing new conditional independencies from old ones, but can be cumbersome to use when many variables are involved. It turns out that a convenient way of expressing conditional independence relationships, such that many of these deductions can be made immediately, is through graphs.

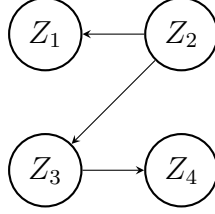
**Definition 8.** A *graph* is a pair  $\mathcal{G} = (V, E)$  where  $V$  is a set of *vertices* or *nodes* and  $E \subseteq V \times V$  with  $(v, v) \notin E$  for any  $v \in V$  is a set of *edges*.

Let  $Z = (Z_1, \dots, Z_p)^\top$  be a collection of random variables. The graphs we will consider will always have  $V = \{1, \dots, p\} =: [p]$ , so  $V$  indexes the random variables.

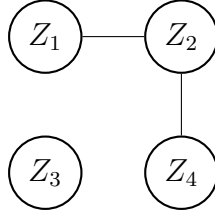
Let  $j, k \in V$ .

- We say there is an *edge* between  $j$  and  $k$  and that  $j$  and  $k$  are *adjacent* if either  $(j, k) \in E$  or  $(k, j) \in E$ .
- An edge  $(j, k)$  is *undirected* if also  $(k, j) \in E$ ; otherwise it is *directed* and we may write  $j \rightarrow k$  to represent this.

- If all edges in the graph are (un)directed we call it an *(un)directed graph*. We can represent graphs as pictures: for example, we can draw the graph when  $p = 4$  and  $E = \{(2, 1), (3, 4), (2, 3)\}$  as



If instead we have  $E = \{(1, 2), (2, 1), (2, 4), (4, 2)\}$  we get the undirected graph



- Say  $j$  is a *parent* of  $k$  and  $k$  is a *child* of  $j$  if  $j \rightarrow k$ . The sets of parents and children of  $k$  will be denoted  $\text{pa}(k)$  and  $\text{ch}(k)$  respectively.
- A *path* from  $j$  to  $k$  is a sequence  $j = j_1, j_2, \dots, j_m = k$  of (at least two) distinct vertices such that  $j_l$  and  $j_{l+1}$  are adjacent. Such a path is a *directed path* if  $j_l \rightarrow j_{l+1}$  for all  $l$ . We then call  $k$  a *descendant* of  $j$ . The set of descendants of  $j$  will be denoted  $\text{de}(j)$ .
- A *directed cycle* (or, for brevity, a ‘cycle’) is (almost) a directed path but with the start and end points the same. A *directed acyclic graph (DAG)* is a directed graph containing no directed cycles.

**Definition 9.** Given a DAG  $\mathcal{G} = ([p], E)$ , we say that a permutation  $\pi$  of  $[p]$  is a *topological ordering* if it satisfies

$$\pi(j) < \pi(k) \quad \text{whenever } k \in \text{de}(j).$$

**Proposition 30.** *Every DAG with  $p$  vertices has a topological ordering.*

*Proof.* We use induction on the number of nodes  $p$ . Clearly the result is true when  $p = 1$ .

Now we show that in any DAG, we can find a node with no parents. Pick any node and move to one of its parents, if possible. Then move to one of the new node’s parents, and continue in this fashion. This procedure must terminate since no node can be visited twice, or we would have found a cycle. The final node we visit must therefore have no parents, which we call a source node.

Suppose then that  $p \geq 2$ , and we know that all DAGs with  $p-1$  nodes have a topological ordering. Find a source  $s$  (wlog  $s = p$ ) and form a new DAG  $\tilde{\mathcal{G}}$  with  $p-1$  nodes by removing the source (and all edges emanating from it). Note we keep the labelling of the nodes in this new DAG the same. This smaller DAG must have a topological order  $\tilde{\pi}$ . A topological ordering  $\pi$  for our original DAG is then given by  $\pi(s) = 1$  and  $\pi(k) = \tilde{\pi}(k) + 1$  for  $k \neq s$ .  $\square$

### 3.3 Undirected graphical models

**Definition 10.** The *conditional independence graph* of a distribution  $P$  on  $\mathbb{R}^p$  with  $p \geq 2$  is the undirected graph  $\mathcal{G} = ([p], E)$  where given  $Z \sim P$ , we have for all  $j, k \in [p]$  that<sup>1</sup>

$$(j, k), (k, j) \in E \text{ if and only if } Z_j \not\perp\!\!\!\perp Z_k \mid Z_{-jk}.$$

Thus, undirected edges occur in a conditional independence graph if and only if the corresponding variables are conditionally dependent, given all other variables.

A second approach relating graphs and conditional independencies asks for the distribution to reflect further aspects of the structure of the graph through the notion of graph separation: given disjoint subsets  $A, B$  and  $S$  of vertices, we say  $S$  *separates*  $A$  from  $B$  if every path between a node in  $A$  and a node in  $B$  contains a node in  $S$ .

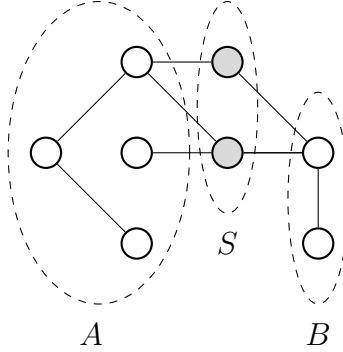


Figure 3.1: Illustration of separation in an undirected graph.

**Definition 11.** A distribution  $P$  on  $\mathbb{R}^p$  is *global Markov* with respect to an undirected graph  $\mathcal{G}$  with  $p$  vertices if whenever  $Z \sim P$  and  $A, B$  and  $S$  are disjoint sets of vertices such that  $S$  separates  $A$  from  $B$ , we have  $Z_A \perp\!\!\!\perp Z_B \mid Z_S$ .

**Theorem 31.** *If a distribution  $P$  on  $\mathbb{R}^p$  has a density with respect to a product measure that is positive everywhere, then it is global Markov with respect to its conditional independence graph.*

<sup>1</sup>Here and throughout, when  $A, B \subseteq [p]$  are non-empty and disjoint, and  $S = \emptyset$ , we interpret  $Z_A \perp\!\!\!\perp Z_B \mid Z_S$  as the (unconditional) independence relationship  $Z_A \perp\!\!\!\perp Z_B$ . Moreover, we adopt the convention that if either  $A$  or  $B$  is the empty set, then for every random vector  $Z$  and every  $S \subseteq [p]$ , the statement  $Z_A \perp\!\!\!\perp Z_B \mid Z_S$  holds.



*Proof.* We use backwards induction on the cardinality of the separating set. When  $|S| = p - 2$ , we must have  $S = [p] \setminus \{j, k\}$  for some distinct  $j, k \in [p]$  with  $(j, k), (k, j) \notin E$ . But then  $Z_j \perp\!\!\!\perp Z_k \mid Z_{-jk}$  by definition of the conditional independence graph. Suppose now that  $m \in [p - 2]$  and that the claim is true whenever the separating set  $S \subseteq [p]$  satisfies  $|S| = m$ . Take disjoint subsets of vertices  $A, B, S$  with  $|S| = m - 1$  such that  $S$  separates  $A$  from  $B$ . Any vertex not in  $A, B$  or  $S$  can be added to one of  $A$  or  $B$  while maintaining this separation, as otherwise  $A$  and  $B$  would not have been separated by  $S$  initially. We may therefore assume without loss of generality that  $|A| \geq 2$ . Fix  $j \in A$  and set  $A_- := A \setminus \{j\}$ . Since  $\{j\} \cup S$  separates  $A_-$  from  $B$  and  $A_- \cup S$  separates  $\{j\}$  from  $B$ , we have by the induction hypothesis that  $Z$  satisfies

$$Z_{A_-} \perp\!\!\!\perp Z_B \mid Z_S, Z_j \quad \text{and} \quad Z_j \perp\!\!\!\perp Z_B \mid Z_S, Z_{A_-}.$$

But then the intersection property yields that  $Z_A \perp\!\!\!\perp Z_B \mid Z_S$ , as required.  $\square$

An important application of the above ideas concerns a model-free notion of variable significance in regression settings. Suppose that  $X$  is a  $p$ -dimensional covariate vector and  $Y$  is a real-valued response. A *Markov blanket* for  $Y$  is a collection of predictors  $X_S$  such that

$$Y \perp\!\!\!\perp X_{S^c} \mid X_S.$$

The interpretation is that  $X_S$  contains all the information in  $X$  that is relevant for learning about  $Y$ . In general there are many sets  $S$  satisfying this requirement, with one example being  $S = [p]$ ; however, they must all contain

$$S^* := \{j \in [p] : Y \not\perp\!\!\!\perp X_j \mid X_{-j}\}.$$

Indeed, if  $X_S$  is a Markov blanket for  $Y$  and  $j \notin S$ , then by the weak union property,  $Y \perp\!\!\!\perp X_j \mid X_{-j}$ , so  $j \notin S^*$ . Moreover, if  $(X, Y)$  satisfies the intersection property, then since  $X_{S^*}$  separates  $Y$  and  $X_{S^{*,c}}$  in the conditional independence graph of  $(X, Y)$ , we see by Theorem 31 that  $Y \perp\!\!\!\perp X_{S^{*,c}} \mid X_{S^*}$ . Thus  $X_{S^*}$  is itself a Markov blanket, and is therefore, by the above argument, the unique minimal Markov blanket, where the minimality is with respect to set inclusion.

### 3.4 Directed graphical models and causality

Conditional independence graphs give us some understanding of relationships between variables. However they do not help us address questions such as ‘If we were to set the  $j$ th variable to a particular value, say 0.5, then how would the distribution of the other variables be altered?’. Yet this is often the sort of causal question that we would like to answer. The notion of a structural causal model can help us with this more ambitious goal.

### 3.4.1 Structural causal models

**Definition 12.** A *structural causal model*<sup>2</sup> (SCM) is a collection  $\mathcal{S} = (P_j, h_j, Q_j)_{j \in [p]}$ , where

- $P_j \subseteq [p] \setminus \{j\}$  for each  $j \in [p]$ , and these subsets are such that the graph with vertices  $[p]$  and edges satisfying  $\text{pa}(j) = P_j$  for each  $j \in [p]$  is a DAG;
- $h_j : \mathbb{R}^{|P_j|} \times \mathbb{R} \rightarrow \mathbb{R}$  for each  $j \in [p]$ ;
- $Q_j$  is a probability distribution on  $\mathbb{R}$ , for each  $j \in [p]$ .

We sometimes refer to the DAG defined by the SCM as a *causal DAG*. Associated with an SCM is a system of structural equations of the form

$$Z_j := h_j(Z_{P_j}, \varepsilon_j) \quad \text{for } j \in [p].$$

Here, the  $\varepsilon_1, \dots, \varepsilon_p$  are independent (noise) random variables with  $\varepsilon_j \sim Q_j$ , and when  $P_j = \emptyset$ , we interpret  $h_j$  to be a function of  $\varepsilon_j$  alone.

An SCM may be thought of as a recipe for how to generate a random vector. Indeed, using a topological ordering  $\pi$  for the associated DAG, we can write each  $Z_j$  as a function of  $\varepsilon_{\pi^{-1}(1)}, \varepsilon_{\pi^{-1}(2)}, \dots, \varepsilon_{\pi^{-1}(\pi(j))}$ . Importantly though, it also offers a framework for reasoning about how the distribution of the random vector will behave after changes to the system, as we now explain.

### 3.4.2 Interventions

Given an SCM  $\mathcal{S}$ , we can modify it to replace one of the structural assignments  $Z_j := h_j(Z_{P_j}, \varepsilon_j)$  with  $Z_j := \tilde{h}_j(Z_{\tilde{P}_j}, \varepsilon_j)$  where now  $\varepsilon_j \sim \tilde{Q}_j$ . This in turn results in a new distribution for  $Z$ . We say that node  $j$  has been *intervened* on. In the case where  $Z_j := a$ , we call this a *perfect intervention*. Expectations and probabilities are written by appending ‘ $|\text{do}(Z_j = a)$ ’, e.g.  $\mathbb{E}(Z_k | \text{do}(Z_j = a))$ . Note that in general this will be different from  $\mathbb{E}(Z_j | Z_k = a)$ .

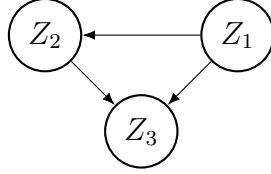
**Example 3.4.1.** The following SCM, involving three binary variables, is inspired by a real medical study comparing the effectiveness of two different treatments for kidney stones. These stones are classified as either large ( $Z_1 = 1$ ) or small ( $Z_1 = 0$ ) and the treatments given are open surgery ( $Z_2 = 1$ ), involving a large incision, or percutaneous nephrolithotomy ( $Z_2 = 0$ ), where only a small incision is made. The outcome variable  $Z_3$  being 1 indicates that the treatment has been successful. The relationships between these variables are modelled via the following system of structural equations:

$$\begin{aligned} Z_1 &= \varepsilon_1, & \varepsilon_1 &\sim \text{Bern}(1/2), \\ Z_2 &= \mathbb{1}_{\{\varepsilon_2(1+2Z_1) > 3/4\}}, & \varepsilon_2 &\sim U[0, 1], \\ Z_3 &= \mathbb{1}_{\{\varepsilon_3(2-Z_1+7Z_2/4-3Z_1Z_2/2) > 1/4\}}, & \varepsilon_3 &\sim U[0, 1]. \end{aligned}$$

---

<sup>2</sup>Some authors refer to structural causal models as (nonparametric) structural equation models.

The corresponding DAG is



Consider now the intervention  $\text{do}(Z_2 = 1)$ , i.e. the treatment is open surgery. This gives a new system of structural equations:

$$\begin{aligned}
 Z_1 &= \varepsilon_1, & \varepsilon_1 &\sim \text{Bern}(1/2), \\
 Z_2 &= 1, \\
 Z_3 &= \mathbb{1}_{\{\varepsilon_3(15/4 - 5Z_1/2) > 1/4\}}, & \varepsilon_3 &\sim U[0, 1].
 \end{aligned}$$

Thus  $\mathbb{P}(Z_3 = 1 \mid \text{do}(Z_2 = 1)) = \frac{1}{2}(\frac{4}{5} + \frac{14}{15}) = \frac{13}{15}$ . Similarly one may compute the probability of success when the treatment is percutaneous nephrolithotomy, and find that  $\mathbb{P}(Z_3 = 1 \mid \text{do}(Z_2 = 0)) = \frac{13}{16}$ , thus showing that open surgery is to be preferred. On the other hand,

$$\begin{aligned}
 \mathbb{P}(Z_3 = 1 \mid Z_2 = 1) &= \sum_{j=0}^1 \mathbb{P}(Z_3 = 1 \mid Z_2 = 1, Z_1 = j) \mathbb{P}(Z_1 = j \mid Z_2 = 1) \\
 &= \sum_{j=0}^1 \mathbb{P}(Z_3 = 1 \mid Z_2 = 1, Z_1 = j) \frac{\mathbb{P}(Z_2 = 1 \mid Z_1 = j) \mathbb{P}(Z_1 = j)}{\mathbb{P}(Z_2 = 1)} \\
 &= \frac{1}{\frac{1}{2}(\frac{1}{4} + \frac{3}{4})} \left( \frac{14}{15} \cdot \frac{1}{4} + \frac{4}{5} \cdot \frac{3}{4} \right) \cdot \frac{1}{2} = \frac{5}{6}.
 \end{aligned}$$

Similarly we may compute

$$\mathbb{P}(Z_3 = 1 \mid Z_2 = 0) = \frac{7}{8} \cdot \frac{3}{4} + \frac{3}{4} \cdot \frac{1}{4} = \frac{27}{32} > \frac{5}{6}.$$

Thus naively equating the conditional probability of  $Z_3 = 1$  given  $Z_2$  with success of treatment  $Z_2$  would lead us to conclude incorrectly that percutaneous nephrolithotomy is the preferred treatment. The reason for the discrepancy here is that as indicated in the original SCM, the presence of large kidney stones increases the chance of open surgery being performed, but also decreases the chance of success of either treatment.

The example above illustrates that the conclusions that may be drawn from an SCM go far beyond those of its associated joint distribution. The crucial additional piece of information included in the SCM is the causal structure encoded by the associated DAG. As we shall see below, a causal DAG can carry with it implications for the conditional independence structure. In principle, this can allow for a postulated DAG to be falsified if an implied conditional independence is not reflected in the data.

### 3.4.3 Markov properties for DAGs

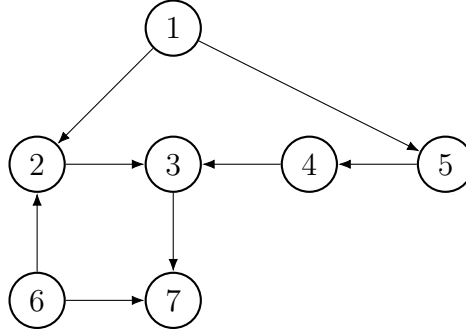
**Definition 13.** Given a DAG  $\mathcal{G}$  and a (not necessarily directed) path  $j_1, \dots, j_m$  with  $m \geq 3$ , we say  $j_\ell$  is a *collider* (relative to the path) if  $j_{\ell-1} \rightarrow j_\ell \leftarrow j_{\ell+1}$ . A path is *blocked* by a set of nodes  $S$  if there exists a node  $v$  on the path such that either:

- (i)  $v$  is not a collider and  $v \in S$ , or
- (ii)  $v$  is a collider and neither  $v$  nor any of its descendants are in  $S$ .

Thus, the path is not blocked if (i) every non-collider is not in  $S$  and (ii) every collider is either in  $S$  or has a descendant in  $S$ . Given disjoint subsets of nodes  $A$  and  $B$ , we say  $A$  and  $B$  are *d-separated*<sup>3</sup> by a subset of nodes  $S$ , and write  $A \perp\!\!\!\perp_{\mathcal{G}} B \mid S$ , if every path between  $A$  and  $B$  is blocked by  $S$ .

In particular, if there is no path between  $A$  and  $B$ , then  $A \perp\!\!\!\perp_{\mathcal{G}} B \mid S$  for every subset  $S$  of nodes including  $S = \emptyset$ .

**Example 3.4.2.** Consider the following DAG  $\mathcal{G}$ :



We claim that  $2 \perp\!\!\!\perp_{\mathcal{G}} 5 \mid 1, 6$ . To see this, observe that:

- the path  $2 \leftarrow 1 \rightarrow 5$  contains the node 1, which is a non-collider;
- the path  $2 \rightarrow 3 \leftarrow 4 \leftarrow 5$  contains the collider 3, and neither 3 nor its sole descendant 7 lie in  $\{1, 6\}$ ;
- the path  $2 \leftarrow 6 \rightarrow 7 \leftarrow 3 \leftarrow 4 \leftarrow 5$  between 2 and 5 via 6 contains for example  $2 \leftarrow 6 \rightarrow 7$ .

In fact, we also have  $2 \perp\!\!\!\perp_{\mathcal{G}} 5 \mid 1$ , because the path between 2 and 5 via 6 contains the collider 7, which has no descendants.

**Definition 14.** We say a distribution  $P$  on  $\mathbb{R}^p$  is *global Markov* with respect to a DAG  $\mathcal{G}$  on  $p$  vertices if whenever  $Z \sim P$ , and  $A, B$  and  $S$  are disjoint sets of vertices such that  $A \perp\!\!\!\perp_{\mathcal{G}} B \mid S$ , we have  $Z_A \perp\!\!\!\perp Z_B \mid Z_S$ .

**Theorem 32.** If a distribution  $P$  on  $\mathbb{R}^p$  is generated by an SCM with DAG  $\mathcal{G}$  on  $p$  vertices, then  $P$  is global Markov with respect to  $\mathcal{G}$ .

<sup>3</sup>d-separation is short for *directional separation*.

### 3.4.4 Causal structure learning

We have seen how a causal DAG implies certain conditional independencies. It follows that the presence of conditional *dependencies* in a distribution can be used to rule out potential causal DAGs, and thus reveal underlying causal structure. For example, the edges in a conditional independence graph can be given a causal interpretation:

**Proposition 33.** *Let  $p \geq 2$  and let  $Z \sim P$ . If  $Z_j \not\perp\!\!\!\perp Z_k | Z_{-[j,k]}$  for  $j, k \in [p]$ , then any causal DAG that generates  $P$  must either have  $j \rightarrow k$ ,  $j \leftarrow k$  or  $j \rightarrow \ell \leftarrow k$  for some  $\ell \in [p] \setminus \{j, k\}$ .*

*Proof.* Suppose that  $P$  is generated by a structural causal model with DAG  $\mathcal{G}$ . Then by Theorem 32,  $P$  satisfies the global Markov property with respect to  $\mathcal{G}$ . Thus, if  $Z_j \not\perp\!\!\!\perp Z_k | Z_{-[j,k]}$ , then the nodes  $j$  and  $k$  cannot be  $d$ -separated by  $S := [p] \setminus \{j, k\}$  in  $\mathcal{G}$ . But any path between  $j$  and  $k$  other than  $j \rightarrow k$ ,  $k \rightarrow j$  and  $j \rightarrow \ell \leftarrow k$  for some  $\ell \in [p] \setminus \{j, k\}$  must have a non-collider in  $S$  so is blocked by  $S$ . It follows that one of these three types of path must exist.  $\square$

In fact, Proposition 34 below reveals that under the stronger hypothesis that two variables are conditionally dependent when we condition on any subset of the remaining variables, we can obtain the stronger conclusion that they must be adjacent in any causal DAG; in other words, there is a direct causal link between the variables.

**Proposition 34.** *If nodes  $j$  and  $k$  in a DAG with  $p$  vertices are not adjacent and  $\pi$  is a topological order on  $[p]$  with  $\pi(j) < \pi(k)$ , then they are  $d$ -separated by  $\text{pa}(k)$ .*

*Proof.* Consider a path with vertices  $j = j_1, \dots, j_m = k$ ; we aim to show that it is blocked by  $\text{pa}(k)$ . If  $j_{m-1} \rightarrow k$  then it is blocked because  $j_{m-1}$  is a non-collider that belongs to  $\text{pa}(k)$ . If instead  $j_{m-1} \leftarrow k$ , then consider the maximal  $\ell \in \{2, 3, \dots, m-1\}$  such that  $j_{\ell-1} \rightarrow j_\ell \leftarrow j_{\ell+1}$ ; such an  $\ell$  must exist as otherwise we would have a directed path from  $k$  to  $j$ , contradicting the topological ordering. Since  $j_\ell$  is a collider relative to the path, and does not belong to  $\text{pa}(k)$ , the only way for the path not to be blocked by  $\text{pa}(k)$  would be for  $j_\ell$  to have a descendant in  $\text{pa}(k)$ . But this cannot occur as it would introduce a cycle. Thus the path must be blocked.  $\square$

## 3.5 Gaussian graphical models

We have seen that conditional dependence is a compelling way to assess the relatedness of variables in the presence of others. The problem of testing the null hypothesis that  $X \perp\!\!\!\perp Y | Z$  is thus of great interest. This however is a fundamentally difficult task unless  $Z$  is discrete [Shah and Peters, 2020], and we need to place certain assumptions on the form of the joint distribution of  $X, Y, Z$  to make progress. One assumption that offers a significant simplification of the problem is that of joint Gaussianity. Approaches that use this assumption are based on the following key result.

**Proposition 35.** Let  $Z \sim N_p(\mu, \Sigma)$  with  $\Sigma$  positive definite. Then for  $A, B \subseteq [p]$ ,

$$Z_A | Z_B = z_B \sim N_{|A|}(\mu_A + \Sigma_{A,B} \Sigma_{B,B}^{-1} (z_B - \mu_B), \Sigma_{A,A} - \Sigma_{A,B} \Sigma_{B,B}^{-1} \Sigma_{B,A}).$$

### 3.5.1 Nodewise regression

Specialising Proposition 35 to the case where  $A = \{k\}$  and  $B = A^c$  we see that when conditioning on  $Z_{-k} = z_{-k}$ , we may write

$$Z_k = m_k + z_{-k}^\top \Sigma_{-k,-k}^{-1} \Sigma_{-k,k} + \varepsilon_k,$$

where

$$m_k = \mu_k - \Sigma_{k,-k} \Sigma_{-k,-k}^{-1} \mu_{-k}$$

$$\varepsilon_k | Z_{-k} = z_{-k} \sim N(0, \Sigma_{k,k} - \Sigma_{k,-k} \Sigma_{-k,-k}^{-1} \Sigma_{-k,k}).$$

Note that if the  $j$ th element of the vector of coefficients  $\Sigma_{-k,-k}^{-1} \Sigma_{-k,k}$  is zero, then the distribution of  $Z_k$  conditional on  $Z_{-k}$  will not depend at all on the  $j$ th component of  $Z_{-k}$ . Then if that  $j$ th component was  $Z_{j'}$ , we would have that  $Z_k | Z_{-k} = z_{-k}$  has the same distribution as  $Z_k | Z_{-j'k} = z_{-j'k}$ , so  $Z_k \perp\!\!\!\perp Z_j | Z_{-j'k}$ .

Thus given  $x_1, \dots, x_n \stackrel{\text{i.i.d.}}{\sim} Z$  and writing

$$X = \begin{pmatrix} x_1^\top \\ \vdots \\ x_n^\top \end{pmatrix},$$

we may estimate the coefficient vector  $\Sigma_{-k,-k}^{-1} \Sigma_{-k,k}$  by regressing  $X_k$  on  $X_{\{k\}^c}$  and including an intercept term.

The technique of *neighbourhood selection* [Meinshausen and Bühlmann, 2006] involves performing such a regression for each variable, using the Lasso. There are two options for populating our estimate of the CIG with edges based on the Lasso estimates. Writing  $\hat{S}_k$  for the selected set of variables when regressing  $X_k$  on  $X_{\{k\}^c}$ , we can use the “OR” rule and put an edge between vertices  $j$  and  $k$  if and only if  $k \in \hat{S}_j$  or  $j \in \hat{S}_k$ . An alternative is the “AND” rule where we put an edge between  $j$  and  $k$  if and only if  $k \in \hat{S}_j$  and  $j \in \hat{S}_k$ .

### 3.5.2 The Graphical Lasso

Another popular approach to estimating the CIG works by first directly estimating the precision matrix  $\Omega := \Sigma^{-1}$ . To see how this works, the following fact concerning blockwise inversion of matrices is helpful.

**Proposition 36.** Let  $M \in \mathbb{R}^{p \times p}$  be a symmetric positive definite matrix and suppose

$$M = \begin{pmatrix} P & Q^\top \\ Q & R \end{pmatrix}$$

with  $P$  and  $R$  square matrices. The Schur complement of  $R$  is  $P - Q^\top R^{-1}Q =: S$ . We have that  $S$  is positive definite and

$$M^{-1} = \begin{pmatrix} S^{-1} & -S^{-1}Q^\top R^{-1} \\ -R^{-1}QS^{-1} & R^{-1} + R^{-1}QS^{-1}Q^\top R^{-1} \end{pmatrix}.$$

Furthermore  $\det(M) = \det(S)\det(R)$ .

We thus see that  $\Sigma_{-k,-k}^{-1}\Sigma_{-k,k} = -\Omega_{kk}^{-1}\Omega_{-k,k}$ , so

$$(\Sigma_{-k,-k}^{-1}\Sigma_{-k,k})_j = 0 \Leftrightarrow \begin{cases} \Omega_{j,k} = 0 & \text{for } j < k \\ \Omega_{j+1,k} = 0 & \text{for } j \geq k. \end{cases}$$

Thus

$$Z_k \perp\!\!\!\perp Z_j \mid Z_{-jk} \Leftrightarrow \Omega_{jk} = 0.$$

The *graphical Lasso* [Yuan and Lin, 2007] produces sparse estimates of  $\Omega$  given  $x_1, \dots, x_n \stackrel{\text{i.i.d.}}{\sim} N_p(\mu, \Omega^{-1})$  by minimising a sum of the negative log-likelihood  $-\ell(\mu, \Omega)$  and an  $\ell_1$ -penalty. We have

$$\max_{\mu \in \mathbb{R}^p} \ell(\mu, \Omega) = -\frac{n}{2} \{\text{tr}(S\Omega) - \log \det(\Omega)\};$$

see example sheet. Writing  $S := \frac{1}{n} \sum_{i=1}^n x_i x_i^\top - \bar{X} \bar{X}^\top$  where  $\bar{X} := \frac{1}{n} \sum_{i=1}^n x_i$ , the graphical Lasso minimises

$$-\log \det(\Omega) + \text{tr}(S\Omega) + \lambda \sum_{j,k} |\Omega_{jk}|,$$

over symmetric positive definite matrices  $\Omega \in \mathbb{R}^{p \times p}$ . Often the penalty term is modified such that the diagonal elements are not penalised.

## 3.6 Basic asymptotic statistics

**Definition 15.** Let  $W_1, W_2, \dots$  and  $W$  be real-valued random variables.

- We say the  $W_n$  converge in distribution to  $W$  with distribution function  $F$ , and write  $W_n \xrightarrow{d} W$ , if for all  $t \in \mathbb{R}$  at which  $F$  is continuous,

$$\mathbb{P}(W_n \leq t) \rightarrow F(t) \quad \text{as } n \rightarrow \infty.$$

- We say the  $W_n$  converge in probability to  $W$  and write  $W_n \xrightarrow{p} W$  if for all  $\epsilon > 0$ ,

$$\mathbb{P}(|W_n - W| > \epsilon) \rightarrow 0.$$

One can show that if  $W_n \xrightarrow{p} W$ , then  $W_n \xrightarrow{d} W$  (so convergence in probability is a stronger notion of convergence). However, they coincide if  $W = c \in \mathbb{R}$  is deterministic i.e. if  $W_n \xrightarrow{d} c$ , then  $W_n \xrightarrow{p} c$ .

**Lemma 37** (Weak law of large numbers (WLLN)). *If  $W_1, W_2, \dots$  are i.i.d. real-valued random variables and  $\mathbb{E}(W_1) = \mu < \infty$ , then as  $n \rightarrow \infty$ ,*

$$\frac{1}{n} \sum_{i=1}^n W_i \xrightarrow{p} \mu.$$

**Theorem 38** (Continuous mapping theorem (CMT)). *Suppose the sequence of random variables  $(W_n)_{n=1}^\infty$  is such that  $W_n \xrightarrow{p} W$ . Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be continuous at every point in a set  $C$  with  $\mathbb{P}(W \in C) = 1$ . Then  $f(W_n) \xrightarrow{p} f(W)$ .*

**Lemma 39** (Slutsky's lemma). *Let  $(U_n)_{n=1}^\infty$  and  $(W_n)_{n=1}^\infty$  be sequences of random variables where  $U_n \xrightarrow{d} U$  and  $W_n \xrightarrow{p} c$  for random variable  $U \in \mathbb{R}$  and deterministic  $c \in \mathbb{R}$ . Then*

1.  $U_n + W_n \xrightarrow{d} U + c$ ,
2.  $U_n W_n \xrightarrow{d} U c$ ,
3.  $U_n / W_n \xrightarrow{d} U / c$  if  $c \neq 0$ .

## 3.7 Conditional independence testing

While a Gaussian assumption offers simplicity, particularly in settings with larger sample sizes, it can be hard to defend. An alternative approach aims to leverage the predictive power of modern machine learning methods.

Consider the setting where our data  $(X_1, Y_1, Z_1), \dots, (X_n, Y_n, Z_n)$  consist of i.i.d. copies of the triple  $(X, Y, Z) \in \mathbb{R} \times \mathbb{R} \times \mathbb{R}^p$ . We collect together  $\mathcal{X} := (X_i)_{i=1}^n$ <sup>4</sup>, and similarly for  $\mathcal{Y}$  and  $\mathcal{Z}$ . Our goal is to test  $X \perp\!\!\!\perp Y \mid Z$ .

One way of connecting prediction (i.e. regression) with the problem of conditional independence testing is via the following observation:

$$X \perp\!\!\!\perp Y \mid Z \implies \mathbb{E}[\{X - \mathbb{E}(X \mid Z)\}\{Y - \mathbb{E}(Y \mid Z)\}] = 0 \quad (3.2)$$

provided  $\mathbb{E}X^2, \mathbb{E}Y^2 < \infty$ . In other words, the population level residuals from regressing each of  $X$  and  $Y$  on  $Z$  are uncorrelated. To see this, observe that in fact as  $\mathbb{E}(XY \mid Z) = \mathbb{E}(X \mid Z)\mathbb{E}(Y \mid Z)$ , under conditional independence, we have

$$\mathbb{E}[\{X - \mathbb{E}(X \mid Z)\}\{Y - \mathbb{E}(Y \mid Z)\} \mid Z] = 0,$$

so (3.2) follows from the tower property. The relationship in (3.2) suggests regressing each of  $X$  and  $Y$  on  $Z$ , and constructing a test statistic based on the products of the residuals. Let us write

$$X = f(Z) + \varepsilon \quad \text{and} \quad Y = g(Z) + \xi, \quad (3.3)$$

---

<sup>4</sup>Note we are departing from our previous notation where  $\mathcal{X}$  was the input space.



where  $f(z) := \mathbb{E}(X | Z = z)$ ,  $g(z) := \mathbb{E}(Y | Z = z)$ , and so under conditional independence  $\mathbb{E}(\varepsilon\xi) = 0$ .

Given fitted regression functions  $\hat{f}$  and  $\hat{g}$  from regressing  $X$  and  $Y$  respectively on  $Z$ , consider

$$\tau_N := \frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_i \hat{\xi}_i.$$

where  $\hat{\varepsilon}_i := X_i - \hat{f}(Z_i)$  and  $\hat{\xi}_i := Y_i - \hat{g}(Z_i)$ . Define  $\varepsilon_i := X_i - f(Z_i)$  and  $\xi_i := Y_i - g(Z_i)$ . If  $f$  and  $g$  are estimated sufficiently well by  $\hat{f}$  and  $\hat{g}$ , then the  $i$ th summand above will be close to  $\varepsilon_i \xi_i$ . Under the null these quantities are mean-zero and i.i.d., so the CLT suggests we can expect  $\sqrt{n}\tau_N$  to have a Gaussian distribution. To obtain a standard normal limit, we should normalise by the square root of

$$\tau_D^2 := \frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_i^2 \hat{\xi}_i^2 - \tau_N^2,$$

and so take as our final test statistic

$$T := \sqrt{n} \frac{\tau_N}{\tau_D},$$

known as the *Generalised Covariance Measure* (GCM) [Shah and Peters, 2020]. The result below formalises our intuition that provided

$$\mathcal{E}_f := \mathbb{E} \left( \frac{1}{n} \sum_{i=1}^n \{f(Z_i) - \hat{f}(Z_i)\}^2 \right) \quad \text{and} \quad \mathcal{E}_g := \mathbb{E} \left( \frac{1}{n} \sum_{i=1}^n \{g(Z_i) - \hat{g}(Z_i)\}^2 \right)$$

decay sufficiently fast, the test statistic will have an asymptotic standard normal distribution under the null.

**Theorem 40.** *Suppose there exists  $c > 0$  such that  $\text{Var}(\varepsilon | Z), \text{Var}(\xi | Z) < c$  almost surely and  $\text{Var}(\varepsilon\xi) > 0$ . Suppose also  $\mathcal{E}_f \rightarrow 0$ ,  $\mathcal{E}_g \rightarrow 0$  and  $\mathcal{E}_f \mathcal{E}_g = o(n^{-1})$ . Then under the null that  $X \perp\!\!\!\perp Y | Z$ , we have  $T \xrightarrow{d} N(0, 1)$ .*

*Proof.* Let us consider  $\tau_N$  first. Substituting  $X_i = f(Z_i) + \varepsilon_i$  and similarly for  $Y_i$ , we have

$$\begin{aligned} n\tau_N &= \sum_{i=1}^n \{f(Z_i) - \hat{f}(Z_i) + \varepsilon_i\} \{g(Z_i) - \hat{g}(Z_i) + \xi_i\} \\ &= \sum_{i=1}^n \{f(Z_i) - \hat{f}(Z_i)\} \{g(Z_i) - \hat{g}(Z_i)\} + \sum_{i=1}^n \varepsilon_i \{g(Z_i) - \hat{g}(Z_i)\} \\ &\quad + \sum_{i=1}^n \xi_i \{f(Z_i) - \hat{f}(Z_i)\} + \sum_{i=1}^n \varepsilon_i \xi_i \\ &=: A_1 + A_2 + A_3 + A_4. \end{aligned}$$

By the CLT,

$$\frac{A_4}{\sqrt{n}} \xrightarrow{d} N(0, \text{Var}(\varepsilon\xi)).$$

Now by the triangle inequality and the Cauchy–Schwarz inequality,

$$\begin{aligned} \mathbb{E}|A_1| &\leq \sum_{i=1}^n \mathbb{E}\{|f(Z_i) - \hat{f}(Z_i)| |g(Z_i) - \hat{g}(Z_i)|\} \\ &\leq \sum_{i=1}^n [\mathbb{E}\{f(Z_i) - \hat{f}(Z_i)\}^2]^{1/2} [\mathbb{E}\{g(Z_i) - \hat{g}(Z_i)\}^2]^{1/2} \\ &\leq \left( \sum_{i=1}^n \mathbb{E}\{f(Z_i) - \hat{f}(Z_i)\}^2 \right)^{1/2} \left( \sum_{i=1}^n \mathbb{E}\{g(Z_i) - \hat{g}(Z_i)\}^2 \right)^{1/2}. \end{aligned}$$

Thus, by Markov's inequality, given  $\delta > 0$ ,

$$\mathbb{P}(|A_1|/\sqrt{n} > \delta) \leq \frac{\delta^{-1}}{\sqrt{n}} \mathbb{E}|A_1| \leq \sqrt{n\mathcal{E}_f\mathcal{E}_g} \rightarrow 0,$$

by assumption, so  $A_1/\sqrt{n} \xrightarrow{p} 0$ . Turning to  $A_2$ , observe that

$$\mathbb{E}(\varepsilon_i \varepsilon_j \{g(Z_i) - \hat{g}(Z_i)\} \{g(Z_j) - \hat{g}(Z_j)\} \mid \mathcal{Y}, \mathcal{Z}) = \{g(Z_i) - \hat{g}(Z_i)\} \{g(Z_j) - \hat{g}(Z_j)\} \mathbb{E}(\varepsilon_i \varepsilon_j \mid \mathcal{Y}, \mathcal{Z}). \quad (3.4)$$

Now  $\mathbb{E}(\varepsilon_i \varepsilon_j \mid \mathcal{Y}, \mathcal{Z}) = \mathbb{E}(\varepsilon_i \varepsilon_j \mid \mathcal{Z})$  as  $\mathcal{X} \perp\!\!\!\perp \mathcal{Y} \mid \mathcal{Z}$  (see example sheet). Moreover,  $\mathbb{E}(\varepsilon_i \varepsilon_j \mid \mathcal{Z}) = \mathbb{E}(\varepsilon_i \mathbb{E}(\varepsilon_j \mid \mathcal{Z}, X_i) \mid \mathcal{Z})$ , but by the weak union property, when  $i \neq j$ ,  $X_j, Z_j \perp\!\!\!\perp \mathcal{Z}_{-j, [p]}, X_i \mid Z_j$ , so  $\mathbb{E}(\varepsilon_j \mid \mathcal{Z}, X_i) = \mathbb{E}(\varepsilon_j \mid Z_j) = 0$ . Similarly  $\mathbb{E}(\varepsilon_i^2 \mid \mathcal{Y}, \mathcal{Z}) = \mathbb{E}(\varepsilon_i^2 \mid Z_i)$ .

Thus, (3.4) is 0 if  $i \neq j$ . Therefore given  $\delta > 0$ , by Markov's inequality and the above,

$$\begin{aligned} \mathbb{P}(|A_2|/\sqrt{n} > \delta) &= \mathbb{P}(A_2^2/n > \delta^2) \leq \frac{\delta^{-2}}{n} \mathbb{E}(A_2^2) \\ &= \frac{\delta^{-2}}{n} \sum_{i=1}^n \mathbb{E}[\mathbb{E}(\varepsilon_i^2 \mid \mathcal{Y}, \mathcal{Z}) \{g(Z_i) - \hat{g}(Z_i)\}^2] \\ &\leq \delta^{-2} c\mathcal{E}_g \rightarrow 0, \end{aligned}$$

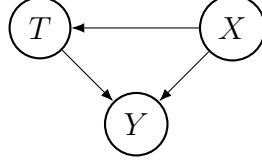
by assumption, so  $A_2/\sqrt{n} \xrightarrow{p} 0$ . Similarly  $A_3/\sqrt{n} \xrightarrow{p} 0$ . Thus by Slutsky's lemma,  $\sqrt{n}\tau_N \xrightarrow{d} N(0, \text{Var}(\varepsilon\xi))$ .

We also have  $\tau_D \xrightarrow{p} \sqrt{\text{Var}(\varepsilon\xi)}$  (see example sheet), so applying Slutsky's lemma again, we get  $T \xrightarrow{d} N(0, 1)$  as required.  $\square$

The key assumption  $\mathcal{E}_f\mathcal{E}_g = o(n^{-1})$  is relatively mild. Indeed, recall from Theorem 9 that if  $f$  and  $g$  lie in RKHS's that satisfy Mercer decompositions, then we can that when kernel ridge regression is used to produce estimates  $\hat{f}$  and  $\hat{g}$ , then we can expect  $\mathcal{E}_f = o(n^{-1/2})$ ,  $\mathcal{E}_g = o(n^{-1/2})$ .

### 3.8 Average treatment effect estimation

Consider now the setting where  $T$  is a binary treatment indicator,  $Y$  is an outcome measure and  $X$  is a vector of pre-treatment covariates. We assume we have an SCM with causal DAG



Note that there may be edges between components of the *confounder* (common cause)  $X$ ; we have not depicted these here and make no assumption about them. A key assumption however is that there are no unobserved confounders. Our goal is to estimate the *average treatment effect*

$$\tau := \mathbb{E}(Y \mid do(T = 1)) - \mathbb{E}(Y \mid do(T = 0)).$$

For simplicity, we will assume below that  $X$  is discrete, though if for example  $X$  is continuous, the same arguments will run through with sums replaced by integrals. We also assume throughout that  $\mathbb{E}Y^2 < \infty$ .

Let  $\mathcal{X} := \{x : \mathbb{P}(X = x) > 0\}$ . We also make the *overlap* assumption that the *propensity score*  $\pi(x) := \mathbb{P}(T = 1 \mid X = x)$  satisfies  $0 < \pi(x) < 1$  for all  $x \in \mathcal{X}$ . From the causal DAG, we see that

$$\mathbb{P}(Y \in A \mid do(T = 1), X = x) = \mathbb{P}(Y \in A \mid T = 1, X = x)$$

for all  $x \in \mathcal{X}$  and all measurable  $A \subseteq \mathbb{R}$ . Thus

$$\begin{aligned} \mathbb{E}(Y \mid do(T = 1)) &= \sum_{x \in \mathcal{X}} \underbrace{\mathbb{E}(Y \mid do(T = 1), X = x)}_{\mathbb{E}(Y \mid T=1, X=x)} \underbrace{\mathbb{P}(X = x \mid do(T = 1))}_{\mathbb{P}(X=x)} \\ &= \mathbb{E}[\mathbb{E}(Y \mid T = 1, X)] \end{aligned}$$

and similarly  $\mathbb{E}(Y \mid do(T = 0)) = \mathbb{E}[\mathbb{E}(Y \mid T = 0, X)]$ . This suggests an estimate of  $\tau$  formed by first obtaining estimates  $\hat{\mu}_j(x)$  of  $\mu_j(x) := \mathbb{E}(Y \mid T = j, X = x)$  for  $j = 0, 1$ . Next, given data  $\mathcal{D}$  formed of i.i.d. copies  $(X_1, Y_1, T_1), \dots, (X_n, Y_n, T_n)$  of  $(X, Y, T)$ , we can estimate  $\tau$  via the regression-based estimator

$$\frac{1}{n} \sum_{i=1}^n \{\hat{\mu}_1(X_i) - \hat{\mu}_0(X_i)\}.$$

An issue is that unless  $\mu_j$  is assumed to take a parametric form, the  $\hat{\mu}_j$  will typically have a substantial bias that decays slower than  $1/\sqrt{n}$  and will propagate to the estimator.

Another approach to estimating  $\tau$  involves re-weighting the data. Observe that  $\mathbb{E}(Y \mid T, X) = T\mu_1(X) + (1 - T)\mu_0(X)$ . Thus

$$\mathbb{E}\left(\frac{YT}{\pi(X)}\right) = \mathbb{E}\left(\frac{\mu_1(X)T}{\pi(X)}\right) = \mathbb{E}[\mu_1(X)],$$

so (by symmetry)

$$\tau = \mathbb{E} \left( \frac{YT}{\pi(X)} - \frac{Y(1-T)}{1-\pi(X)} \right).$$

This suggests an *inverse propensity weighting* (IPW) estimator of the form

$$\frac{1}{n} \sum_{i=1}^n \left( \frac{Y_i T_i}{\hat{\pi}(X_i)} - \frac{Y_i(1-T_i)}{1-\hat{\pi}(X_i)} \right)$$

where  $\hat{\pi}$  is an estimate of the propensity score  $\pi$ . Similarly to the regression estimator however, in a nonparametric setting,  $\hat{\pi}$  and hence IPW may suffer from undesirable bias.

The *augmented inverse propensity weighted* (AIPW) estimator [Robins et al., 1994] aims to combine the strengths of the two approaches and takes the form

$$\hat{\tau} := \hat{\tau}_1 - \hat{\tau}_0$$

where

$$\begin{aligned} \hat{\tau}_1 &:= \frac{1}{n} \sum_{i=1}^n \frac{(Y_i - \hat{\mu}_1(X_i))T_i}{\hat{\pi}(X_i)} + \hat{\mu}_1(X_i) \\ \hat{\tau}_0 &:= \frac{1}{n} \sum_{i=1}^n \frac{(Y_i - \hat{\mu}_0(X_i))(1-T_i)}{1-\hat{\pi}(X_i)} + \hat{\mu}_0(X_i). \end{aligned}$$

To understand the advantages this approach might afford, it is helpful to compare it to an oracle version  $\tau^* := \tau_1^* - \tau_0^*$  that replaces each of the estimates  $\hat{\mu}_1, \hat{\mu}_0$  and  $\hat{\pi}$  with their targets  $\mu_1, \mu_0$  and  $\pi$  respectively.

Now  $\tau^*$  is an average of i.i.d. copies of

$$\frac{\{Y - \mu_1(X)\}T}{\pi(X)} + \mu_1(X) - \frac{\{Y - \mu_0(X)\}(1-T)}{1-\pi(X)} - \mu_0(X)$$

which has mean  $\tau$ . By the central limit theorem, we thus have  $\sqrt{n}(\tau^* - \tau) \xrightarrow{d} N(0, v)$  where  $v$  is the variance of the display above.

If we could show  $\sqrt{n}(\hat{\tau} - \tau^*) \xrightarrow{p} 0$ , then by Slutsky, we would have that  $\sqrt{n}(\hat{\tau} - \tau) \xrightarrow{d} N(0, v)$ . To this end, observe that

$$\begin{aligned} n\{\hat{\tau}_1 - \tau_1^*\} &= \sum_{i=1}^n \left[ T_i \{\mu_1(X_i) - \hat{\mu}_1(X_i)\} \left( \frac{1}{\hat{\pi}(X_i)} - \frac{1}{\pi(X_i)} \right) \right. \\ &\quad + \{Y_i - \mu_1(X_i)\} T_i \left( \frac{1}{\hat{\pi}(X_i)} - \frac{1}{\pi(X_i)} \right) \\ &\quad \left. + \{\hat{\mu}_1(X_i) - \mu_1(X_i)\} \left( 1 - \frac{T_i}{\pi(X_i)} \right) \right] =: A_1 + A_2 + A_3. \end{aligned}$$

The term  $A_1$  can be controlled via Cauchy–Schwarz similarly to the term  $A_1$  in the proof of Theorem 40. The terms  $A_2$  and  $A_3$  could also be controlled similarly to their counterparts in the proof of Theorem 40 if  $\hat{\pi}$  and  $\hat{\mu}_1$  were deterministic functions. In this case, each of  $A_2$  and  $A_3$  would be sums of mean-zero i.i.d. random variables. To mimic this setting, we can estimate them on some held-out data, and then argue conditionally. However, this is wasteful; instead we can use an approach known as *cross-fitting* [Chernozhukov et al., 2018]. This involves splitting the data into  $K$  roughly equal parts. Let  $I_k$  denote the set of indices corresponding to the  $k$ th part or *fold* and for each  $k$ , form estimates  $\hat{\pi}^{(k)}$  and  $\hat{\mu}_j^{(k)}$  using only the data in  $I_k^c$ . We then take

$$\hat{\tau}_1 := \frac{1}{n} \sum_{k=1}^K \sum_{i \in I_k} \frac{(Y_i - \hat{\mu}_1^{(k)}(X_i))T_i}{\hat{\pi}^{(k)}(X_i)} + \hat{\mu}_1^{(k)}(X_i),$$

and forming  $\hat{\tau}_0$  similarly, and set  $\hat{\tau} := \hat{\tau}_1 - \hat{\tau}_0$  as before. Similarly to the case with the GCM, the quality of  $\hat{\tau}$  as an estimator depends on the performances of  $\hat{\pi}$  and the  $\hat{\mu}_j$ . To this end, let

$$\begin{aligned} \mathcal{E}_{\mu,0,n} &:= \mathbb{E}[\{\hat{\mu}_0(X) - \mu_0(X)\}^2], & \mathcal{E}_{\pi,0,n} &:= \mathbb{E}\left[\left\{\frac{1}{1 - \hat{\pi}(X)} - \frac{1}{1 - \pi(X)}\right\}^2\right], \\ \mathcal{E}_{\mu,1,n} &:= \mathbb{E}[\{\hat{\mu}_1(X) - \mu_1(X)\}^2], & \mathcal{E}_{\pi,1,n} &:= \mathbb{E}\left[\left\{\frac{1}{\hat{\pi}(X)} - \frac{1}{\pi(X)}\right\}^2\right]. \end{aligned}$$

Here, for example  $\hat{\pi}$  has been trained on data  $\mathcal{D}$  and is being evaluated at an independent  $X$ , and the expectations are over both of these.

**Theorem 41.** *Suppose there exists  $\eta > 0$  such that  $\eta < \min\{\pi(x), 1 - \pi(x)\}$  for all  $x \in \mathcal{X}$  and  $c > 0$  such that  $\text{Var}(Y | T, X) < c$  almost surely. Assume that  $\mathcal{E}_{\mu,j} \rightarrow 0$ ,  $\mathcal{E}_{\pi,j} \rightarrow 0$  and  $\mathcal{E}_{\mu,j}\mathcal{E}_{\pi,j} = o(n^{-1})$  for  $j = 0, 1$ . Then the cross-fitted AIPW estimator  $\hat{\tau}$  satisfies*

$$\sqrt{n}(\hat{\tau} - \tau) \xrightarrow{d} N(0, v).$$

# Chapter 4

## Multiple testing

In many modern applications, we may be interested in testing many hypotheses simultaneously. Suppose we are interested in testing null hypotheses  $H_1, \dots, H_m$  of which  $m_0$  are true and  $m - m_0$  are not (we do not mention the alternative hypotheses explicitly). Consider the following contingency table:

	Claimed non-significant	Claimed significant (reject)	Total
True null hypotheses	$N_{00}$	$N_{01}$	$m_0$
False null hypotheses	$N_{10}$	$N_{11}$	$m - m_0$
Total	$m - R$	$R$	$m$

The  $N_{jj}$  are unobserved random variables;  $R$  is observed.

Suppose we have  $p$ -values  $p_1, \dots, p_m$  associated with  $H_1, \dots, H_m$  and  $H_i, i \in I_0$  are the true null hypotheses, so

$$\mathbb{P}(p_i \leq \alpha) \leq \alpha$$

for all  $\alpha \in [0, 1]$ ,  $i \in I_0$ <sup>1</sup>. Traditional approaches to multiple testing have sought to control the familywise error rate (FWER) defined by

$$\text{FWER} = \mathbb{P}(N_{01} \geq 1)$$

at a prescribed level  $\alpha$ ; i.e. find procedures for which  $\text{FWER} \leq \alpha$ . The simplest such procedure is the *Bonferroni correction*, which rejects  $H_i$  if  $p_i \leq \alpha/m$ .

**Proposition 42.** *Using Bonferroni correction,*

$$\mathbb{P}(N_{01} \geq 1) \leq \mathbb{E}(N_{01}) \leq \frac{m_0 \alpha}{m} \leq \alpha.$$

---

<sup>1</sup>Note that the probability  $\mathbb{P}(p_i \leq \alpha)$  may depend on where in the null hypothesis we are; we have suppressed this dependence throughout.

*Proof.* The first inequality comes from Markov's inequality. Next

$$\begin{aligned}\mathbb{E}(N_{01}) &= \mathbb{E}\left(\sum_{i \in I_0} \mathbb{1}_{\{p_i \leq \alpha/m\}}\right) \\ &= \sum_{i \in I_0} \mathbb{P}(p_i \leq \alpha/m) \\ &\leq \frac{m_0 \alpha}{m}.\end{aligned}\quad \square$$

A more sophisticated approach is the closed testing procedure.

## 4.1 The closed testing procedure

Given our family of hypotheses  $\{H_i\}_{i=1}^m$ , define the *closure* of this family to be

$$\{H_I : I \subseteq \{1, \dots, m\}, I \neq \emptyset\}$$

where  $H_I = \cap_{i \in I} H_i$  is known as an *intersection hypothesis* ( $H_I$  is the hypothesis that all  $H_i$   $i \in I$  are true).

Suppose that for each  $I$ , we have an  $\alpha$ -level test  $\phi_I$  taking values in  $\{0, 1\}$  for testing  $H_I$  (we reject if  $\phi_I = 1$ ), so under  $H_I$ ,

$$\mathbb{P}_{H_I}(\phi_I = 1) \leq \alpha.$$

The  $\phi_I$  are known as *local tests*.

The *closed testing procedure* [Marcus et al., 1976] is defined as follows:

Reject  $H_I$  if and only if for all  $J \supseteq I$ ,  
 $H_J$  is rejected by the local test  $\phi_J$ .

Typically we only make use of the individual hypotheses that are rejected by the procedure i.e. those rejected  $H_I$  where  $I$  is a singleton.

We consider the case of 4 hypotheses as an example. Suppose the underlined hypotheses are rejected by the local tests.

$$\begin{array}{c} \underline{H_{1234}} \\ \underline{H_{123}} \ \underline{H_{124}} \ \underline{H_{134}} \ \underline{H_{234}} \\ \underline{H_{12}} \ \underline{H_{13}} \ \underline{H_{14}} \ \underline{H_{23}} \ \underline{H_{24}} \ \underline{H_{34}} \\ \underline{H_1} \ \underline{H_2} \ H_3 \ H_4 \end{array}$$

- Here  $H_1$  is rejected by the closed testing procedure.
- $H_2$  is not rejected by the closed testing procedure as  $H_{24}$  is not rejected by the local test.

- $H_{23}$  is rejected by the closed testing procedure.

**Theorem 43.** *The closed testing procedure makes no false rejections with probability  $1 - \alpha$ . In particular it controls the FWER at level  $\alpha$ .*

*Proof.* Assume  $I_0$  is not empty (as otherwise no rejection can be false anyway). Define the events

$$\begin{aligned} A &= \{\text{at least one false rejection}\} \supseteq \{N_{01} \geq 1\}, \\ B &= \{\text{reject } H_{I_0} \text{ with the local test}\} = \{\phi_{I_0} = 1\}. \end{aligned}$$

In order for there to be a false rejection, we must have rejected  $H_{I_0}$  with the local test. Thus  $B \supseteq A$ , so

$$\text{FWER} \leq \mathbb{P}(A) \leq \mathbb{P}(\phi_{I_0} = 1) \leq \alpha. \quad \square$$

Different choices for the local tests give rise to different testing procedures. *Holm's procedure* takes  $\phi_I$  to be the Bonferroni test i.e.

$$\phi_I = \begin{cases} 1 & \text{if } \min_{i \in I} p_i \leq \frac{\alpha}{|I|} \\ 0 & \text{otherwise.} \end{cases}$$

It can be shown (see example sheet) that Holm's procedure amounts to ordering the  $p$ -values  $p_1, \dots, p_m$  as  $p_{(1)} \leq \dots \leq p_{(m)}$  with corresponding hypothesis tests  $H_{(1)}, \dots, H_{(m)}$ , so  $(i)$  is the index of the  $i$ th smallest  $p$ -value, and then performing the following.

Step 1. If  $p_{(1)} \leq \alpha/m$  reject  $H_{(1)}$ , and go to step 2. Otherwise accept  $H_{(1)}, \dots, H_{(m)}$  and stop.

Step  $i$ . If  $p_{(i)} \leq \alpha/(m-i+1)$ , reject  $H_{(i)}$  and go to step  $i+1$ . Otherwise accept  $H_{(i)}, \dots, H_{(m)}$ .

Step  $m$ . If  $p_{(m)} \leq \alpha$ , reject  $H_{(m)}$ . Otherwise accept  $H_{(m)}$ .

The  $p$ -values are visited in ascending order and rejected until the first time a  $p$ -value exceeds a given critical value. This sort of approach is known (slightly confusingly) as a *step-down* procedure.

## 4.2 The False Discovery Rate

A different approach to multiple testing does not try to control the FWER, but instead attempts to control the *false discovery rate* (FDR) defined by

$$\begin{aligned} \text{FDR} &= \mathbb{E}(\text{FDP}) \\ \text{FDP} &= \frac{N_{01}}{\max(R, 1)}, \end{aligned}$$



where FDP is the *false discovery proportion*. Note the maximum in the denominator is to ensure division by zero does not occur. The FDR was introduced in Benjamini and Hochberg [1995], and it is now widely used across science, particularly biostatistics.

The *Benjamini–Hochberg procedure* attempts to control the FDR at level  $\alpha$  and works as follows. Let

$$\hat{k} = \max \left\{ i : p_{(i)} \leq \frac{i\alpha}{m} \right\}.$$

Reject  $H_{(1)}, \dots, H_{(\hat{k})}$  (or perform no rejections if  $\hat{k}$  is not defined).

**Theorem 44.** *Suppose that for each  $i \in I_0$ ,  $p_i$  is independent of  $\{p_j : j \neq i\}$ . Then the Benjamini–Hochberg procedure controls the FDR at level  $\alpha$ ; in fact  $FDR \leq \alpha m_0/m$ .*

*Proof.* For each  $i \in I_0$ , let  $R_i$  denote the number of rejections we get by applying a modified Benjamini–Hochberg procedure to

$$p_{-i} := \{p_1, p_2, \dots, p_{i-1}, p_{i+1}, \dots, p_m\}$$

with cutoff

$$\hat{k}_i = \max \left\{ j : p_{-i,(j)} \leq \frac{\alpha(j+1)}{m} \right\},$$

where  $p_{-i,(j)}$  is the  $j$ th smallest  $p$ -value among  $p_{-i}$ .

For  $r = 1, \dots, m$  and  $i \in I_0$ , note that

$$\begin{aligned} \left\{ p_i \leq \frac{\alpha r}{m}, R = r \right\} &= \left\{ p_i \leq \frac{\alpha r}{m}, p_{(r)} \leq \frac{\alpha r}{m}, p_{(s)} > \frac{\alpha s}{m} \text{ for all } s > r \right\} \\ &= \left\{ p_i \leq \frac{\alpha r}{m}, p_{-i,(r-1)} \leq \frac{\alpha r}{m}, p_{-i,(s-1)} > \frac{\alpha s}{m} \text{ for all } s > r \right\} \\ &= \left\{ p_i \leq \frac{\alpha r}{m}, R_i = r - 1 \right\}. \end{aligned}$$

Thus

$$\begin{aligned}
\text{FDR} &= \mathbb{E}\left(\frac{N_{01}}{\max(R, 1)}\right) \\
&= \sum_{r=1}^m \mathbb{E}\left(\frac{N_{01}}{r} \mathbb{1}_{\{R=r\}}\right) \\
&= \sum_{r=1}^m \frac{1}{r} \mathbb{E}\left(\sum_{i \in I_0} \mathbb{1}_{\{p_i \leq \alpha r/m\}} \mathbb{1}_{\{R=r\}}\right) \\
&= \sum_{r=1}^m \frac{1}{r} \sum_{i \in I_0} \mathbb{P}(p_i \leq \alpha r/m, R = r) \\
&= \sum_{r=1}^m \frac{1}{r} \sum_{i \in I_0} \mathbb{P}(p_i \leq \alpha r/m) \mathbb{P}(R_i = r - 1) \\
&\leq \frac{\alpha}{m} \sum_{i \in I_0} \sum_{r=1}^m \mathbb{P}(R_i = r - 1) \\
&= \frac{\alpha m_0}{m}.
\end{aligned}$$

□

# Bibliography

- A. Belloni, V. Chernozhukov, and L. Wang. Square-root lasso: pivotal recovery of sparse signals via conic programming. *Biometrika*, 98(4):791–806, 2011.
- Y. Benjamini and Y. Hochberg. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society, Series B*, 57:289–300, 1995.
- V. Chernozhukov, D. Chetverikov, M. Demirer, E. Duflo, C. Hansen, W. Newey, and J. Robins. Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, pages C1–C68, 2018.
- A. Hoerl and R. Kennard. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, pages 55–67, 1970.
- G. S. Kimeldorf and G. Wahba. A correspondence between bayesian estimation on stochastic processes and smoothing by splines. *The Annals of Mathematical Statistics*, 41(2):495–502, 1970.
- R. Marcus, P. Eric, and K. R. Gabriel. On closed testing procedures with special reference to ordered analysis of variance. *Biometrika*, 63(3):655–660, 1976.
- N. Meinshausen. Relaxed lasso. *Computational Statistics and Data Analysis*, 52:374–393, 2007.
- N. Meinshausen and P. Bühlmann. High dimensional graphs and variable selection with the lasso. *Annals of Statistics*, 34:1436–1462, 2006.
- A. Rahimi and B. Recht. Random features for large-scale kernel machines. In *Advances in Neural Information Processing Systems*, pages 1177–1184, 2007.
- J. M. Robins, A. Rotnitzky, and L. P. Zhao. Estimation of regression coefficients when some regressors are not always observed. *Journal of the American statistical Association*, 89(427):846–866, 1994.
- B. Schölkopf, R. Herbrich, and A. J. Smola. A generalized representer theorem. In *International Conference on Computational Learning Theory*, pages 416–426. Springer, 2001.

- R. D. Shah and J. Peters. The hardness of conditional independence testing and the generalised covariance measure. *Annals of Statistics*, 48:1514–1538, 2020.
- T. Sun and C.-H. Zhang. Scaled sparse linear regression. *Biometrika*, 99(4):879–898, 2012.
- R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B*, 58:267–288, 1996.
- R. Tibshirani, M. Saunders, S. Rosset, J. Zhu, and K. Knight. Sparsity and smoothness via the fused lasso. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(1):91–108, 2005.
- P. Tseng. Convergence of a block coordinate descent method for nondifferentiable minimization. *Journal of optimization theory and applications*, 109(3):475–494, 2001.
- S. Van de Geer, P. Bühlmann, Y. Ritov, R. Dezeure, et al. On asymptotically optimal confidence regions and tests for high-dimensional models. *The Annals of Statistics*, 42(3):1166–1202, 2014.
- M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(1):49–67, 2006.
- M. Yuan and Y. Lin. Model selection and estimation in the gaussian graphical model. *Biometrika*, 94(1):19–35, 2007.
- C.-H. Zhang. Nearly unbiased variable selection under minimax concave penalty. *The Annals of statistics*, pages 894–942, 2010.
- C.-H. Zhang and S. S. Zhang. Confidence intervals for low dimensional parameters in high dimensional linear models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(1):217–242, 2014.
- H. Zou. The adaptive lasso and its oracle properties. *Journal of the American Statistical Association*, 101:1418–1429, 2006.