

# Modern Statistical Methods — Summary

Lucas Riedstra

November 16, 2020

## Contents

<b>1</b>	<b>Kernel machines</b>	<b>2</b>
1.1	Ridge regression	2
1.1.1	The SVD and PCA	3
1.2	$v$ -fold cross validation	4
1.3	The kernel trick	5
1.4	Kernels	6
1.4.1	Examples of kernels	7
1.4.2	Reproducing kernel Hilbert spaces	8
1.4.3	The representer theorem	9
1.5	Kernel ridge regression	10
1.6	Other kernel machines	12
1.6.1	The support vector machine	12
1.6.2	Logistic regression	13
1.7	Large scale kernel machines	14
<b>2</b>	<b>The Lasso</b>	<b>16</b>
2.1	Model selection	16
2.2	Lasso estimator	16
2.2.1	Slow prediction error rate	16
2.2.2	Concentration inequalities	17
2.2.3	Optimisation theory and convex analysis	19
2.2.4	Lasso solutions	20
2.2.5	Variable selection	20
2.2.6	Prediction and estimation	21
2.2.7	Computation	24
2.3	Extensions of the Lasso	25
2.3.1	Structural penalties	25
2.3.2	Debiasing the Lasso	26
<b>3</b>	<b>High-dimensional covariance estimation and PCA</b>	<b>27</b>
3.1	Covariance estimation	27
3.1.1	Maximum likelihood in multivariate normal model	27
3.1.2	Non-asymptotic error bounds	27

Classical models rely on so-called “large  $n$  asymptotics” (where  $n$  is the sample size). This course focuses on the scenario where  $p$ , the number of variables, is larger or about as large as  $n$ . In this case, the classical theory breaks down, so we need new methods.

## 1 Kernel machines

We represent data are pairs  $(Y_i, x_i) \in \mathbb{R} \times \mathbb{R}^p$  ( $i = 1, \dots, n$ ). The random variables  $Y_i$  are called the *responses*, and the (fixed) variables  $x_i$  are called *predictors*.

**Recap 1.1.** Let  $X = (X_1, \dots, X_n)^\top$  be a multivariate random variable. Its distribution function is given by

$$F_X: \mathbb{R}^n \rightarrow [0, 1]: \mathbf{x} \mapsto \mathbb{P}(X_1 \leq x_1, \dots, X_n \leq x_n).$$

Its expected value is given by

$$\mathbb{E}[X] := (\mathbb{E}[X_1], \dots, \mathbb{E}[X_n])^\top \in \mathbb{R}^n.$$

Its covariance matrix is given by

$$\text{Var}[X] = \mathbb{E}[(X - \mathbb{E}[X])(X - \mathbb{E}[X])^\top] = \mathbb{E}[XX^\top] - \mathbb{E}[X]\mathbb{E}[X]^\top.$$

The matrix  $\text{Var}[X]$  is symmetric positive semidefinite and satisfies  $(\text{Var}[X])_{ij} = \text{Cov}(X_i, X_j)$ .

**Definition 1.2.** In a *linear model*, we assume that

$$Y_i = x_i^\top \beta^0 + \varepsilon_i \quad (i = 1, \dots, n).$$

where  $\beta \in \mathbb{R}^p$  is unknown and the multivariate random variable  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)$  satisfies  $\mathbb{E}(\varepsilon) = 0$  and  $\text{Var}(\varepsilon) = \sigma^2 I$ .

**Definition 1.3.** For an estimator  $\tilde{\beta}$  of  $\beta^0$ , its *mean squared error* (MSE) is given by

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

Note that if  $\tilde{\beta}$  is unbiased, the second term will disappear and the MSE is simply the variance.

**Recap 1.4.** The maximum likelihood estimator (MLE) in this model is the ordinary least squares (OLS) estimator  $\hat{\beta}^{\text{OLS}} = (X^\top X)^{-1} X^\top Y$ , where the *design matrix*  $X \in \mathbb{R}^{n \times p}$  is the matrix whose rows are the vectors  $x_i$ . This estimator only exists if  $X$  has full column rank, so in particular, it is required that  $p \leq n$ .

The Cramér-Rao lower bound states that, out of all unbiased estimators, the MLE has the optimal variance *asymptotically* (i.e., for  $n \rightarrow \infty$ ).

### 1.1 Ridge regression

**Definition 1.5.** Let  $\lambda \geq 0$ , and let  $\mathbf{1} \in \mathbb{R}^n$  be the all-ones vector. Then we define the *ridge regression* estimators

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

where the used norm is the 2-norm. The parameter  $\lambda$  is called the *regularisation parameter*.

The parameter  $\lambda$  represents a penalty for large coefficients in the design matrix. The intercept is not penalised — this is because a shift in units should not affect the fitted values. However,  $X\hat{\beta}$  is not invariant under scale transformations, so it is common practice to centre the columns of  $X$  to have mean 0, and then scale them to have  $\ell_2$ -norm  $\sqrt{n}$ .

After that, we can compute  $\hat{\mu}_\lambda^R$  by taking the derivative:

$$\begin{aligned}\|Y - \mu\mathbf{1} - X\beta\|^2 &= \sum_i (Y_i - \mu - \sum_j X_{ij}\beta_j)^2. \\ \frac{\partial}{\partial \mu} \|Y - \mu\mathbf{1} - X\beta\|^2 &= -2 \sum_i \left( Y_i - \mu - \sum_j X_{ij}\beta_j \right).\end{aligned}$$

Setting this derivative equal to 0 yields

$$\begin{aligned}-2 \sum_i \left( Y_i - \mu - \sum_j X_{ij}\beta_j \right) &= 0 \\ \sum_i Y_i - n\mu - \sum_j \beta_j \left( \sum_i X_{ij} \right) &= 0 \\ \sum_i Y_i - n\mu &= 0 \\ \mu &= \frac{1}{n} \sum_i Y_i = \bar{Y}.\end{aligned}$$

Therefore we conclude  $\hat{\mu}_\lambda^R = \bar{Y}$ . After centering the responses (i.e. replacing  $Y_i$  by  $Y_i - \bar{Y}$ ), the problem can be reduced to

$$\hat{\beta}_\lambda^R = \arg \min_{\beta \in \mathbb{R}^p} Q(\beta) := \arg \min_{\beta \in \mathbb{R}^p} \left[ \|Y - X\beta\|^2 + \lambda \|\beta\|^2 \right].$$

Since  $Q(\beta)$  is convex quadratic, there is a unique root, and to find it we compute

$$\nabla_\beta Q(\beta) = 2X^\top(Y - X\beta) + 2\lambda\beta = 0 \iff \beta = (X^\top X + \lambda I)^{-1} X^\top Y.$$

We conclude that  $\hat{\beta}_\lambda^R = (X^\top X + \lambda I)^{-1} X^\top Y$ . Note that, even if  $X$  does not have full column rank, this estimator exists for all  $\lambda > 0$ . In fact, for  $\lambda$  sufficiently small, the ridge estimator outperforms the MLE in terms of mean squared error:

**Theorem 1.6.** Fix  $\beta^0, \sigma^2$ , and assume that  $\hat{\beta}^{\text{OLS}}$  exists (i.e.,  $X$  has full column rank). For some  $\lambda > 0$  sufficiently small, it holds that the MSE of  $\hat{\beta}_\lambda^{\text{OLS}}$  minus the MSE of  $\hat{\beta}_\lambda^R$  is positive definite.

*Proof.* This is simply writing out the MSE's. In the end, we find that the result holds for  $0 < \lambda < 2\sigma^2 / \|\beta^0\|^2$ .  $\square$

### 1.1.1 The SVD and PCA

**Recap 1.7.** Recall that any  $X \in \mathbb{R}^{n \times p}$  can be factorised as  $X = UDV^\top$ , where  $U, V$  are  $n \times n$  and  $p \times p$  orthogonal matrices respectively, and  $D \in \mathbb{R}^{n \times p}$  satisfies  $D_{11} \geq \dots \geq D_{mm} \geq 0$  where  $m := \min(n, p)$ , and all other entries of  $D$  are 0. This is called the *singular value decomposition* or SVD of  $X$ .

If  $n > p$ , we can replace  $U$  by its first  $p$  columns and  $D$  by its first  $p$  rows to produce the so-called *thin SVD* of  $X$ . Then  $U \in \mathbb{R}^{n \times p}$  has orthogonal columns (so  $U^\top U = I$ ) and  $D \in \mathbb{R}^{p \times p}$

is square and diagonal.

Suppose  $n \geq p$  and let  $X = UDV^\top$  be the thin SVD of our design matrix  $X$ . Then we can write the fitted values from the ridge regression as follows:

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

Note that for OLS ( $\lambda = 0$ ), this is simply the projection of  $Y$  onto the column space of  $X$  (if  $X$  has full column rank). If  $\lambda > 0$ ,  $Y$  is still projected onto the column space of  $X$ , but the projection is shrunk in the directions of the left singular vectors, and the lower the corresponding singular value, the higher the shrinkage.

**Principal component analysis** Consider  $v \in \mathbb{R}^p$  with norm 1, then since the columns of  $X$  have been centered, the sample mean of  $Xv$  is 0, and the sample variance is therefore

$$\frac{1}{n} \sum_i (Xv)_i^2 = \frac{1}{n} (Xv)^\top Xv = \frac{1}{n} v^\top X^\top Xv = \frac{1}{n} v^\top VD^2V^\top v.$$

Writing  $a = V^\top v$  (with  $\|a\| = 1$ ), we find

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

Therefore, we see that the above is maximised if  $a = \pm e_1$ , or equivalently  $v = \pm V_1$ . Therefore,  $V_1$  determines which combination of columns of  $X$  has the largest variance (subject to having norm 1), and  $XV_1 = D_{11}U_1$  is known as the *first principal component* of  $X$ . Analogously, it can be shown that  $D_{22}U_2, \dots, D_{pp}U_p$  have maximum variance  $D_{jj}^2/n$ , subject to being orthonormal to all earlier principal components.

We see that ridge regression shrinks  $Y$  most in the smaller principal components of  $X$ . Therefore it will work well if most of the information is in the larger principal components of  $X$ .

**A comment on computation** By analogous calculations as before, one can compute  $\hat{\beta}_\lambda^R = V(D^2 + \lambda I)^{-1}DU^\top Y$ . Since calculating the inverse of a diagonal matrix is trivial, we see that the complexity of computing  $\hat{\beta}_\lambda^R$  for any  $\lambda$  lies in  $O(np)$ . Of course, this is after computation of the SVD of  $X$ , which lies in  $O(np \min(n, p))$ .

## 1.2 $v$ -fold cross validation

Of course, we are still left with the problem of choosing  $\lambda$  in ridge regression. We consider one possible way of doing so, namely  $v$ -fold cross validation, which is a general way of selecting a good regression method from several competing methods. Here, we assume that our predictors are random, so that we have i.i.d. data pairs  $(x_i, Y_i)$  ( $i = 1, \dots, n$ ). Suppose  $(x^*, Y^*)$  is a new data pair, independent of

$(X, Y)$  and identically distributed. Ideally, we want to pick  $\lambda$  which minimises the prediction error (averaged over  $Y^*$  and  $x^*$ )

$$\mathbb{E} \left[ \left( Y^* - (x^*)^\top \hat{\beta}_\lambda^R(X, Y) \right)^2 \mid X, Y \right],$$

where the dependence of  $\hat{\beta}_\lambda^R$  on the training data  $(X, Y)$  is made explicit by denoting it  $\hat{\beta}_\lambda^R(X, Y)$ .

This is impossible to minimise, but it may be possible to minimise the expected prediction error (averaged over the training data)

$$\mathbb{E} \left\{ \mathbb{E} \left[ \left( Y^* - (x^*)^\top \hat{\beta}_\lambda^R(X, Y) \right)^2 \mid X, Y \right] \right\}. \quad (1)$$

This is still not possible to compute directly, but we estimate it using  $v$ -fold cross validation. Split the data into  $v$  groups or *folds* of roughly equal size  $(X^{(1)}, Y^{(1)}), \dots, (X^{(v)}, Y^{(v)})$  and let  $(X^{(-k)}, Y^{(-k)})$  denote all data except that in the  $k$ -th fold. Then we define

$$\text{CV}(\lambda) := \frac{1}{n} \sum_{i=1}^n \left[ Y_i - x_i^\top \hat{\beta}_\lambda^R(X^{(-\kappa(i))}, Y^{(-\kappa(i))}) \right]^2,$$

and choose the value of  $\lambda$  that minimises  $\text{CV}(\lambda)$ .

The function  $\text{CV}(\lambda)$  is called the *out-of-sample error*, since the training data does not include  $x_i$ .

**Recap 1.8.** The *tower rule* states that for random variables  $X, Y$  we have  $\mathbb{E}[X] = \mathbb{E}[\mathbb{E}[X \mid Y]]$ .

Note that for each  $i$ , we have

$$\mathbb{E} \left[ \left\{ Y_i - x_i^\top \hat{\beta}_\lambda^R(X^{(-\kappa(i))}, Y^{(-\kappa(i))}) \right\}^2 \right] = \mathbb{E} \left[ \mathbb{E} \left[ \left\{ Y_i - x_i^\top \hat{\beta}_\lambda^R(X^{(-\kappa(i))}, Y^{(-\kappa(i))}) \right\}^2 \mid X^{(-\kappa(i))}, Y^{(-\kappa(i))} \right] \right].$$

This equals the expected prediction error in eq. (1), except that the training data  $X, Y$  are replaced with a smaller data set.

We now have a bias-variance tradeoff in the size of the folds: if  $v = n$  (known as “leave-one-out” cross-validation), the estimation will be almost unbiased, but the averaged quantities in  $\text{CV}(\lambda)$  will be highly correlated which leads to high variance. Typical choices of  $v$  are 5 or 10.

Instead of finding the single best  $\lambda$ , we can also aim to find the best weighted combination of  $\lambda$ 's. For example, suppose  $\lambda$  is restricted to a grid  $\lambda_1 > \dots > \lambda_L$ . Then we can use any nonnegative least-squares optimization algorithm to minimise

$$\frac{1}{n} \sum_{i=1}^n \left[ Y_i - \sum_{\ell=1}^L w_\ell x_i^\top \hat{\beta}_{\lambda_\ell}^R(X^{(-\kappa(i))}, Y^{(-\kappa(i))}) \right]^2,$$

over all  $w \in \mathbb{R}_{\geq 0}^L$ . This procedure is known as *stacking* and often outperforms cross-validation.

### 1.3 The kernel trick

We note that

$$X^\top (X X^\top + \lambda I) = (X^\top X + \lambda I) X^\top,$$

and multiplying from the left with  $(X^\top X + \lambda I)^{-1}$  and from the right with  $(X X^\top + \lambda I)^{-1}$  gives

$$(X^\top X + \lambda I)^{-1} X^\top = X^\top (X X^\top + \lambda I)^{-1}.$$

Using this, we see that we can rewrite the fitted values from ridge regression as follows:

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

Two important remarks:

1. Computing the LHS of this equation takes roughly  $O(np^2 + p^3)$  operations, while computing the RHS takes  $O(n^2p + n^3)$  operations (this is because in the LHS we invert an  $p \times p$  matrix, while in the RHS we invert a  $n \times n$  matrix). Therefore, if  $p \gg n$ , the RHS can be much cheaper to compute.
2. The LHS depends only on the matrix  $K = XX^\top$  (this matrix is called the *kernel matrix*). Intuitively, since  $K_{ij} = \langle x_i, x_j \rangle$ , the entries of the kernel matrix show how ‘similar’ the corresponding predictors are.

**Example 1.9.** Suppose we have data  $(Y_i, z_i)_{i=1, \dots, n}$  with  $z_i = (z_{i1}, \dots, z_{id})^\top$ , and we believe the following quadratic relation holds:

$$Y_i = \sum_k \sqrt{2}\gamma_k z_{ik} + \sum_{k, \ell} \vartheta_{k\ell} z_{ik} z_{i\ell} + \varepsilon_i.$$

To compute fitted values using ridge regression, we can rewrite this as a linear model  $Y = X\beta + \varepsilon$  where

$$\beta = \begin{pmatrix} \gamma_1 \\ \vdots \\ \gamma_d \\ \vartheta_{11} \\ \vartheta_{12} \\ \vdots \\ \vartheta_{dd} \end{pmatrix}, \quad x_i = \begin{pmatrix} \sqrt{2}z_{i1} \\ \vdots \\ \sqrt{2}z_{id} \\ z_{i1}z_{i1} \\ z_{i1}z_{i2} \\ \vdots \\ z_{id}z_{id} \end{pmatrix}.$$

In this case, we have  $p = d^2 + d$  variables, which means computing  $(X^\top X + \lambda I)^{-1}$  takes  $O(d^6)$  operations. In this case, computing  $(XX^\top + \lambda I)^{-1}$  is probably easier.

We are still left with the problem of computing  $K := XX^\top$ , which can take  $O(n^2p) = O(n^2d^2)$  operations if done naively. However, observe that

$$K_{ij} = x_i^\top x_j = 2 \sum_k z_{ik} z_{jk} + \sum_{k, \ell} z_{ik} z_{i\ell} z_{jk} z_{j\ell} = \left(1 + \sum_k z_{ik} z_{jk}\right)^2 - 1 = (1 + z_i^\top z_j) - 1.$$

This quantity can be computed in  $O(d)$ , and therefore  $K$  can be computed in  $O(n^2d)$  operations: we have a factor  $d$  improvement.

The general point of the previous example is that we can bypass the features  $x_i$  entirely and instead think directly of  $K = XX^\top$  where an entry  $K_{ij}$  represents similarity between the inputs of the  $i$ -th and  $j$ -th samples. This leads to the notion of a kernel in general.

## 1.4 Kernels

We will assume our inputs  $x_1, \dots, x_n$  live in an abstract space  $\mathcal{X}$ .

**Definition 1.10.** A (*positive-definite*) *kernel* is a symmetric map  $k: \mathcal{X}^2 \rightarrow \mathbb{R}$  such that for all  $n \in \mathbb{N}$  and all  $x_1, \dots, x_n \in \mathcal{X}$ , the matrix  $K \in \mathbb{R}^{n \times n}$  with  $K_{ij} = k(x_i, x_j)$  is positive semi-definite.

**Proposition 1.11** (Cauchy-Schwarz for kernels). *Let  $k$  be a kernel and  $x, x' \in \mathcal{X}$ , then*

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

*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 its determinant must be non-negative.  $\square$

In our old models, the data points  $x_i$  were vectors in  $\mathbb{R}^p$ . Now we try to think of them as points in an abstract space with an associated *feature map*  $\varphi: \mathcal{X} \rightarrow \mathcal{H}$  (with  $\mathcal{H}$  an inner product space), and a kernel  $k(x, x')$  gives a measure of similarity between  $\varphi(x)$  and  $\varphi(x')$ . In this case, we have the following:

**Proposition 1.12.** *Let  $\mathcal{H}$  be an inner product space,  $\varphi: \mathcal{X} \rightarrow \mathcal{H}$  and define  $k(x, x') := \langle \varphi(x), \varphi(x') \rangle$ . Then  $k$  is a kernel.*

*Proof.* We have, for all  $x_1, \dots, x_n \in \mathcal{X}$  and  $\alpha \in \mathbb{R}^n$  that

$$\alpha^\top K \alpha = \sum_{i,j} K_{ij} \alpha_i \alpha_j = \sum_{i,j} \langle \varphi(x_i), \varphi(x_j) \rangle \alpha_i \alpha_j = \left\| \sum_i \alpha_i \varphi(x_i) \right\|^2 \geq 0.$$

□

### 1.4.1 Examples of kernels

The following proposition shows how to make new kernels from old:

**Proposition 1.13.** *Suppose  $k_1, k_2, \dots$  are kernels. Then:*

1. *If  $\alpha_1, \alpha_2 \geq 0$  then  $\alpha_1 k_1 + \alpha_2 k_2$  is a kernel.*
2. *The pointwise limit of a sequence of kernels is a kernel (if it exists).*
3. *The pointwise product  $k_1 k_2$  is a kernel.*

*Proof.* See Example Sheet 1.

□

**Example 1.14.** Let us consider some examples of kernels:

1. For  $\mathcal{X} = \mathbb{R}^p$  we have already seen the *linear kernel*  $k(x, x') = x^\top x'$ .
2. For  $\mathcal{X} = \mathbb{R}^p$ , the *polynomial kernel* is defined as  $k(x, x') = (1 + x^\top x')^d$ . This is a kernel since it is a power of a sum of two kernels.
3. For  $\mathcal{X} = \mathbb{R}^p$ , the *Gaussian kernel* is defined by

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

To show this is a kernel, write  $k$  as the pointwise product  $k_1 k_2$  where

$$k_1(x, x') = \exp \left( -\frac{\|x\|^2}{2\sigma^2} \right) \exp \left( -\frac{\|x'\|^2}{2\sigma^2} \right), \quad k_2(x, x') = \exp \left( \frac{x^\top x'}{\sigma^2} \right).$$

Clearly  $k_1$  is the kernel induced by the feature map  $\varphi(x) = \exp(-\|x\|^2/(2\sigma^2))$ , while  $k_2$  can be seen to be a kernel by using the Taylor expansion, which shows that  $k_2$  is a limit of nonnegative linear combinations of kernels.

4. For  $\mathcal{X} = [0, 1]$ , define the *Sobolev kernel*  $k(x, x') = \min(x, x')$ . The proof that this is a kernel is on example sheet 1.
5. For  $\mathcal{X} = \mathcal{P}(\{1, \dots, p\})$ , define the *Jaccard kernel*

$$k(x, x') = \frac{|x \cap x'|}{|x \cup x'|} \quad \text{where } 0/0 := 1.$$

The proof that this is a kernel is on example sheet 1.

### 1.4.2 Reproducing kernel Hilbert spaces

By proposition 1.12, we see that every feature map  $\varphi: \mathcal{X} \rightarrow \mathcal{H}$  gives rise to a kernel. In the next (important!) theorem, we will see that every kernel is in fact induced by a feature map.

**Theorem 1.15.** *Let  $k$  be a kernel, then there exists an inner product space  $\mathcal{H}$  and a feature map  $\varphi: \mathcal{X} \rightarrow \mathcal{H}$  such that*

$$k(x, x') = \langle \varphi(x), \varphi(x') \rangle \quad \text{for all } x, x' \in \mathcal{X}.$$

*Proof.* We will construct  $\mathcal{H}$  and  $\varphi$  explicitly. First we define the function space

$$\mathcal{H} = \left\{ \sum_{i=1}^n \alpha_i k(\cdot, x_i) \mid n \in \mathbb{N}, \alpha_i \in \mathbb{R}, x_i \in \mathcal{X} \right\}.$$

Let  $f = \sum_{i=1}^n \alpha_i k(\cdot, x_i)$  and  $g = \sum_{j=1}^m \beta_j k(\cdot, x'_j)$ , then the inner product on  $\mathcal{H}$  is given by

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

We define  $\varphi: \mathcal{X} \rightarrow \mathcal{H}$  as  $\varphi(x) = k(\cdot, x)$ .

We must check that the inner product does not depend on the choice of representation of  $f$  and  $g$ . For this, note that

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

which holds by symmetry of the kernel. Since  $\sum_i \alpha_i g(x_i)$  is independent of the representation of  $g$ , while  $\sum_j \beta_j f(x'_j)$  is independent of the representation of  $f$ , we conclude that the entire expression is independent of both representations.

Secondly, we must verify that the formula  $k(x, x') = \langle \varphi(x), \varphi(x') \rangle$  indeed holds. For any  $f \in \mathcal{H}, x \in \mathcal{X}$  we have

$$\langle k(\cdot, x), f \rangle = \sum_{i=1}^n \alpha_i k(x_i, x) = f(x), \quad (2)$$

i.e., evaluation of a function is a linear functional in  $\mathcal{H}$ .

In particular, we have

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

Finally, we must check that  $\langle \cdot, \cdot \rangle$  is indeed an inner product. Symmetry and bilinearity are clear. Furthermore, we have

$$\langle f, f \rangle = \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k(x_i, x_j) = \alpha^\top K \alpha \geq 0$$

by the fact that  $k$  is a kernel. We must now only show that  $f \neq 0 \implies \langle f, f \rangle > 0$ . For this, note that  $\langle \cdot, \cdot \rangle$  is a kernel on  $\mathcal{H}$ , so by proposition 1.11 (Cauchy-Schwarz) we have

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

and therefore if  $f$  is nonzero anywhere,  $\langle f, f \rangle$  must also be nonzero.  $\square$

While  $\mathcal{H}$  constructed in the proof is an inner product space, it is not necessarily a Hilbert space. Let  $(f_n) \subseteq \mathcal{H}$  be Cauchy, then by Cauchy-Schwarz for kernels we find

$$f_m(x) - f_n(x) = (f_m - f_n)(x) = \langle k(\cdot, x), f_m - f_n \rangle \leq \sqrt{k(x, x)} \|f_m - f_n\|.$$

We can do an analogous computation for  $f_n - f_m$  to conclude that  $|f_m(x) - f_n(x)| \leq \sqrt{k(x, x)} \|f_n - f_m\|$ , and therefore, if  $(f_n)$  is Cauchy, then it converges pointwise to some  $f^*: \mathcal{X} \rightarrow \mathbb{R}$ . We will not prove the following theorem:



**Theorem 1.16.** *The inner product space  $\mathcal{H}$  constructed in the proof of theorem 1.15 can be extended to a Hilbert space by adding all pointwise limits  $f^*$  of Cauchy sequences  $(f_n) \subseteq \mathcal{H}$ .*

The completion of  $\mathcal{H}$  is a special type of Hilbert space:

**Definition 1.17.** A Hilbert space  $\mathcal{B}$  of functions  $f: \mathcal{X} \rightarrow \mathbb{R}$  is called a *reproducing kernel Hilbert space* (RKHS) if for all  $x \in \mathcal{X}$ , there exists  $k_x \in \mathcal{B}$  such that

$$f(x) = \langle k_x, f \rangle,$$

i.e., evaluation of functions is a linear functional.

The function  $k(x, x') = \langle k_x, k_{x'} \rangle$  is known as the *reproducing kernel* of  $\mathcal{B}$  (induced by the feature map  $\varphi(x) = k_x$ ).

If we start with a kernel  $k$ , construct the corresponding RKHS  $\mathcal{B}$ , then it is easily checked that  $k$  is indeed the reproducing kernel of  $\mathcal{B}$ .

**Example 1.18** (Linear kernel). Let  $X = \mathbb{R}^p$  and  $k(x, x') = x^\top x'$ . Then we have

$$\mathcal{H} = \left\{ x \mapsto \sum_{i=1}^n \alpha_i x^\top x_i = x^\top \left( \sum_i \alpha_i x_i \right) \mid \alpha_i \in \mathbb{R}, x_i \in \mathbb{R}^p \right\} = \{x \mapsto x^\top \beta \mid \beta \in \mathbb{R}^p\},$$

and if  $f(x) = x^\top \beta$ ,  $g(x) = x^\top \beta'$ , then

$$\langle f, g \rangle = k(\beta, \beta') = \beta^\top \beta' \quad \text{so } \|f\|_{\mathcal{H}} = \|\beta\|_2.$$

### 1.4.3 The representer theorem

We started with data  $(Y_i, x_i)$  in a model where our fitted values depend only on  $XX^\top = K$ , where  $K$  is the kernel matrix derived from the kernel  $k(x_i, x_j) = \langle x_i, x_j \rangle = x_i^\top x_j$ . Suppose we change the kernel to some  $\tilde{k}$ , and replace  $K$  by  $\tilde{K}$  accordingly. Then by theorem 1.15 this is equivalent by replacing our predictors  $x_i$  by  $\varphi(x_i)$  for some feature map  $\varphi$ .

Recall that the ridge regression objective function is given by

$$\min_{\beta \in \mathbb{R}^p} \sum_{i=1}^n (Y_i - x_i^\top \beta)^2 + \lambda \|\beta\|^2.$$

Let  $\mathcal{H}$  be the RKHS of the linear kernel, then the above is equivalent to

$$\min_{f \in \mathcal{H}} \sum_{i=1}^n (Y_i - f(x_i))^2 + \lambda \|f\|^2. \tag{3}$$

We will show that in general, if we change our predictors  $x_i$  to  $\varphi(x_i)$ , and update the corresponding RKHS in eq. (3), then eq. (3) solves our new problem. Note that there are no references to the feature map  $\varphi$  in eq. (3).

**Theorem 1.19** (Representer theorem). *Let:*

- $(Y_i, x_i)_{i=1}^n$  be our data with  $x_i \in \mathcal{X}$ ;
- $c: \mathbb{R}^n \times \mathcal{X}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$  be an arbitrary loss function;
- $J: \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$  strictly increasing;
- $k$  a kernel with RKHS  $\mathcal{H}$  and kernel matrix  $K_{ij} = k(x_i, x_j)$ .

Then  $\hat{f} \in \mathcal{H}$  minimises

$$Q_1(f) := c(Y, x_1, \dots, x_n, f(x_1), \dots, f(x_n))$$

if and only if  $\hat{f}(\cdot) = \sum_{i=1}^n \hat{\alpha}_i k(\cdot, x_i)$ , where  $\hat{\alpha} \in \mathbb{R}^n$  minimises

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

*Proof.* Let  $f \in \mathcal{H}$  and  $U = \text{Span}\{k(\cdot, x_1), \dots, k(\cdot, x_n)\} \subseteq \mathcal{H}$ , then  $U$  is closed since it is finite-dimensional so we can write  $f = u + v$  with  $u \in U, v \in U^\perp$ . Note that

$$f(x_i) = \langle k(\cdot, x_i), u + v \rangle = \langle k(\cdot, x_i), u \rangle = u(x_i) \quad (i = 1, \dots, n),$$

so minimisation of  $Q_1$  is equivalent to minimisation of

$$c(Y, x_1, \dots, x_n, u(x_1), \dots, u(x_n)) + J(\|u\|^2 + \|v\|^2) \quad (4)$$

w.r.t.  $u \in U, v \in U^\perp$ . Now, since  $J$  is strictly increasing, any minimiser  $(u, v)$  of eq. (4) will satisfy  $v = 0$ .

Write  $u(\cdot) = \sum_{i=1}^n \alpha_i k(\cdot, x_i) \in U$ , then clearly

$$\begin{bmatrix} u(x_1) \\ \vdots \\ u(x_n) \end{bmatrix} = \begin{bmatrix} \sum_i \alpha_i k(x_1, x_i) \\ \vdots \\ \sum_i \alpha_i k(x_n, x_i) \end{bmatrix} = K\alpha \quad \text{and} \quad \|u\|^2 = \sum_{i,j} \alpha_i \alpha_j k(x_i, x_j) = \alpha^\top K\alpha.$$

This shows that minimisation of eq. (4) is equivalent to minimisation of

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

which completes the proof.  $\square$

**Example 1.20.** In ridge regression, the representer theorem tells us that minimising eq. (3) is equivalent to minimising  $\|Y - K\alpha\|^2 + \lambda \alpha^\top K\alpha$ , and indeed it can be shown (example sheet 1) that the minimiser satisfies  $K\hat{\alpha} = K(K + \lambda I)^{-1}Y = X\hat{\beta}_\lambda^R$ . Since  $\mathcal{H}$  may be infinite-dimensional, this gives a way to rewrite an infinite-dimensional optimization problem to a finite-dimensional optimisation problem.

## 1.5 Kernel ridge regression

We have now defined kernel ridge regression and shown how the estimator may be computed, but we have yet to assess its predictive performance. We consider the model

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

where we assume  $f^0 \in \mathcal{H}$  where  $\mathcal{H}$  is an RKHS with reproducing kernel  $k$ . By scaling the equation on both sides we may assume  $\|f^0\| \leq 1$  (note that this changes  $\text{Var}[\varepsilon]$ ). Let  $K$  be the kernel matrix with eigenvalues  $d_1 \geq \dots \geq d_n \geq 0$ , and let  $\hat{f} = \hat{f}_\lambda$  be the estimated regression function, so

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

**Theorem 1.21.** *The mean squared prediction error (MSPE) satisfies*

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

*Proof.* The representer theorem tells us that  $(\hat{f}(x_1), \dots, \hat{f}(x_n))^\top = K(K + \lambda I)^{-1}Y$ . By projecting  $f^0$  onto  $\text{Span } k(\cdot, x_1), \dots, k(\cdot, x_n)$  it is easily seen that there exists  $\alpha \in \mathbb{R}^n$  such that

$$(f^0(x_1), \dots, f^0(x_n))^\top = K\alpha \quad \text{and} \quad \|f^0\|^2 \geq \alpha^\top K\alpha.$$

Write  $K = UDU^\top$  where  $D_{ii} = d_i$  and define  $\vartheta := U^\top K\alpha = DU^\top \alpha$  (note  $U\vartheta = K\alpha$  and note that  $d_i = 0 \implies \vartheta_i = 0$ ). Then we have

$$\begin{aligned} \mathbb{E} \left[ \sum_{i=1}^n \left( f^0(x_i) - \hat{f}_\lambda(x_i) \right)^2 \right] &= \mathbb{E} \|K\alpha - K(K + \lambda I)^{-1}Y\|^2 \\ &= \mathbb{E} \|K(K + \lambda I)^{-1}(U\vartheta + \varepsilon) - U\vartheta\|^2 \\ &= \mathbb{E} \|DU^\top (UDU^\top + \lambda I)^{-1}(U\vartheta + \varepsilon) - \vartheta\|^2 \\ &= \mathbb{E} \|D(D + \lambda I)^{-1}(\vartheta + U^\top \varepsilon) - \vartheta\|^2 \\ &= \mathbb{E} \left\| [D(D + \lambda I)^{-1} - I]\vartheta + D(D + \lambda I)^{-1}U^\top \varepsilon \right\|^2 \\ &= \mathbb{E} \left\| [D(D + \lambda I)^{-1} - I]\vartheta \right\|^2 + \mathbb{E} \|D(D + \lambda I)^{-1}U^\top \varepsilon\|^2. \end{aligned}$$

Note that the cross-term in the final equality disappears since it is the expectation of a linear combination of  $\varepsilon$ , which is 0.

The first term is a deterministic quantity, which equals

$$\left\| [D(D + \lambda I)^{-1} - I]\vartheta \right\|^2 = \sum_{i=1}^n \left( \left( \frac{d_i}{d_i + \lambda} - 1 \right) \vartheta_i \right)^2 = \sum_{i=1}^n \frac{\lambda^2 \vartheta_i^2}{(d_i + \lambda)^2}.$$

Let  $D^+$  be the diagonal matrix with  $D_{ii}^+ = d_i^{-1}$  if  $d_i \neq 0$  and 0 else. Then we have

$$\sum_{i: d_i > 0} \frac{\vartheta_i^2}{d_i} = \vartheta^\top D^+ \vartheta = \alpha^\top K U D^+ U^\top K \alpha = \alpha^\top U D D^+ D U^\top \alpha = \alpha^\top U D U^\top \alpha = \alpha^\top K \alpha \geq 1,$$

Using this, we can bound the first term by

$$\sum_{i=1}^n \frac{\lambda^2 \vartheta_i^2}{(d_i + \lambda)^2} = \sum_{i: d_i \neq 0} \frac{\vartheta_i^2}{d_i} \frac{d_i \lambda^2}{(d_i + \lambda)^2} \leq \max_i \frac{d_i \lambda^2}{(d_i + \lambda)^2} \leq \max_i \frac{d_i \lambda^2}{4\lambda d_i} \leq \frac{\lambda}{4}.$$

To compute the second term, we use the *trace trick*:

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

This gives our first inequality  $\text{MSPE} \leq \frac{\sigma^2}{n} \sum_{i=1}^n \frac{d_i^2}{(d_i + \lambda)^2} + \frac{\lambda}{4n}$ . For the second inequality, note that  $\frac{d_i^2}{(d_i + \lambda)^2} \leq 1$  and also  $\frac{d_i^2}{(d_i + \lambda)^2} \leq \frac{d_i^2}{4d_i \lambda} = \frac{d_i}{4\lambda}$ . Therefore, we can write  $\frac{d_i^2}{(d_i + \lambda)^2} \leq \frac{1}{\lambda} \min(\lambda, \frac{d_i}{4})$  which proves the second inequality.  $\square$

We can now ask ourselves if kernel ridge regression is optimal in any sense. To this end, we define  $\hat{\mu}_i := d_i/n$  and  $\lambda_n := \lambda/n$ . Then we can rewrite the upper bound from our previous theorem as

$$\text{MSPE} \leq \frac{\sigma^2}{n\lambda_n} \sum_{i=1}^n \min\left(\frac{\hat{\mu}_i}{4}, \lambda_n\right) + \frac{\lambda_n}{4} =: \delta_n(\lambda_n).$$

Now, instead of taking  $x_1, \dots, x_n \in \mathcal{X}$  to be fixed, assume that they are i.i.d. random variables (then  $K$  is random, so the  $\hat{\mu}_i$  are random as well), and taking an expectation on both sides yields

$$\mathbb{E}[\text{MSPE}] \leq \mathbb{E}[\delta_n(\lambda_n)].$$

We want to bound  $\delta_n(\lambda_n)$  by some function of  $n$ . For this we will use Mercer's theorem:

**Theorem 1.22** (Mercer). *Let  $\mathcal{X} = [a, b]$ . If  $k$  is a continuous kernel on  $\mathcal{X}$ , then there exists an orthonormal basis  $(e_i)$  of  $L^2[a, b]$  such that*

$$k(x, x') = \sum_{j=1}^{\infty} \mu_j e_j(x) e_j(x')$$

as well as a sequence of nonnegative eigenvalues  $(\mu_j)$  such that

$$\int_a^b k(x, x') e_j(x) dx = \mu_j e_j(x).$$

Applying Mercer's theorem to  $k$ , it can be proved that for some  $C > 0$  we have

$$\mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n \min\left(\frac{\hat{\mu}_i}{4}, \lambda_n\right)\right] \leq \frac{C}{n} \sum_{i=1}^{\infty} \min\left(\frac{\mu_i}{4}, \lambda_n\right).$$

**Example 1.23.** Let  $k$  be the Sobolev kernel. Then the eigenvalues satisfy

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

and with some calculations it can be shown that  $\sum_{i=1}^{\infty} \min(\frac{\mu_i}{4}, \lambda_n) = O(\sqrt{\lambda_n})$  as  $\lambda_n \rightarrow 0$ . Combining this with the rest of the bound shows

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

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

It has been shown that this is in fact the optimal rate for estimating a function  $f^0 \in \mathcal{H}$  (up to multiplicative constants).

## 1.6 Other kernel machines

### 1.6.1 The support vector machine

Consider a *classification problem* with data  $x_1, \dots, x_n \in \mathbb{R}^p$  and binary response  $Y_i \in \{\pm 1\}$ . Furthermore assume that  $x_1, \dots, x_n$  are separated by a hyperplane through the origin, that is, for some  $\beta \in \mathbb{R}^p$  we have  $Y_i x_i^\top \beta > 0$  for all  $i$ .

To choose between different planes that separate the classes, we pick the hyperplane with the highest margin: that is, we compute

$$\max_{\beta \in \mathbb{R}^p, M \geq 0} M \quad \text{such that} \quad \frac{Y_i x_i^\top \beta}{\|\beta\|} \geq M \quad (i = 1, \dots, n).$$

When the classes are not separable, we can penalise according to the distance of a point “over the margin”. The penalty should be 0 if  $x$  is on the correct side and should equal the distance over the boundary otherwise, measured in units of  $M$ . In this case, by considering minimisation of  $1/M^2$  instead of maximisation of  $M$  we find

$$\arg \min_{M \geq 0, \beta \in \mathbb{R}^p} \frac{1}{M^2} + \lambda \sum_{i=1}^n \left(1 - \frac{Y_i x_i^\top \beta}{\|\beta\| M}\right)_+.$$

Since the above equation is independent of the norm of  $\beta$ , we may set  $\|\beta\| = 1/M$  and remove  $M$  from the equation entirely, thus obtaining

$$\arg \min_{\beta \in \mathbb{R}^p} \|\beta\|^2 + \lambda \sum_{i=1}^n (1 - Y_i x_i^\top \beta)_+.$$

Replacing  $\lambda$  by  $\frac{1}{\lambda}$  and multiplying the equation with  $\lambda$  we obtain

$$\arg \min_{\beta \in \mathbb{R}^p} \sum_{i=1}^n (1 - Y_i x_i^\top \beta)_+ + \lambda \|\beta\|^2.$$

Finally, if we lift the restriction that the hyperplane must go through the origin, we find the objective function

$$(\hat{\mu}, \hat{\beta}) = \arg \min_{\mu \in \mathbb{R}, \beta \in \mathbb{R}^p} \sum_{i=1}^n (1 - Y_i (x_i^\top \beta + \mu))_+ + \lambda \|\beta\|^2. \quad (5)$$

The solution to this problem is known as the *support vector classifier*.

**Introducing kernels** Now we introduce kernels to the problem. Letting  $k$  be the linear kernel and  $\mathcal{H}$  the corresponding RKHS, we can rewrite eq. (5) as

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

Suppose we change  $\mathcal{H}$  to a different RKHS with a different kernel  $k$ . By a variant of the representer theorem (Example Sheet 1, question 9) the solution is equivalent to

$$(\hat{\mu}, \hat{\alpha}) = \arg \min_{\mu \in \mathbb{R}, \alpha \in \mathbb{R}^n} \sum_{i=1}^n (1 - Y_i (K_i^\top \alpha + \mu))_+ + \lambda \alpha^\top K \alpha,$$

and the solution to the above problem is called the *support vector machine*. Predictions at a new point  $x^*$  are given by

$$\text{sgn}(\hat{\mu} + \sum_{i=1}^n \hat{\alpha}_i k(x^*, x_i)) \in \pm 1.$$

*Remark.* Note that the support vector machine can correspond to a nonlinear boundary: the boundary is a hyperplane *in the corresponding RKHS*, but can be nonlinear in the space  $\mathcal{X}$  where the data points live. This also makes it easy to overfit.

### 1.6.2 Logistic regression

**Recap 1.24.** The *standard logistic regression* model is motivated by assuming

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

The maximum-likelihood estimator  $\hat{\beta}$  is given by (Example sheet 2, question 4)

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

Just like in ridge regression, we will try to lower the variance of  $\hat{\beta}$  by penalising large values, i.e.,

$$\begin{aligned} \hat{\beta}_\lambda &:= \arg \min_{\beta \in \mathbb{R}^p} \sum_{i=1}^n \log(1 + \exp(-Y_i x_i^\top \beta)) + \lambda \|\beta\|^2 \\ &= \arg \min_{f \in \mathcal{H}} \sum_{i=1}^n \log(1 + \exp(-Y_i f(x_i))) + \lambda \|f\|^2, \end{aligned}$$

where  $\mathcal{H}$  is the RKHS corresponding to the linear kernel.

Note that this is exactly the same as the support vector machine, but with a different loss function. The loss function in the support vector machine is called “hinge loss”, while the loss function in logistic regression is called “logistic loss”. The main differences are:

1. The logistic loss function increases much quicker for incorrectly classified data.
2. Even when a point is correctly classified, the logistic loss function decreases the further the datapoint is from the boundary. The hinge loss function will always be 0 if a point is correctly classified and lies outside of the margin.

The second point also explains the name *support vector machine*: the only points that contribute to the loss function are those that either violate the boundary or are within the margin (the *support vectors*), and perturbing any of the other points does not change the boundary. Because of this property, SVM's tend to be more stable.

## 1.7 Large scale kernel machines

We have seen that kernels give very flexible regression and classification estimators, and that the parameter  $\lambda$  helps us with the bias-variance tradeoff. Also, when  $p \gg n$ , using the kernel matrix saves computational effort.

We now turn to the case  $n \gg p$ . In this case, both kernel ridge regression and the SVM are very expensive to compute: the former has cost  $O(n^3)$ , while the latter can be computed using an iterative algorithm where every iteration has cost  $O(n^2)$ .

One approach to speed this process up is that of *random feature maps*: we develop a random map  $\hat{\varphi}: \mathcal{X} \rightarrow \mathbb{R}^b$  (with  $b$  small) such that  $\mathbb{E}[\hat{\varphi}(x)^\top \hat{\varphi}(x')] = k(x, x')$  for all  $x, x' \in \mathcal{X}$ . Then we consider  $L$  i.i.d. copies of  $\hat{\varphi}$  (denotes  $\hat{\varphi}_1, \dots, \hat{\varphi}_L$ ), and we consider the feature map

$$\psi(x) = L^{-1/2}(\hat{\varphi}_1(x), \dots, \hat{\varphi}_L(x))^\top \in \mathbb{R}^{Lb}, \text{ so that } \psi(x)^\top \psi(x') = \frac{1}{L} \sum_{i=1}^L \hat{\varphi}_i(x)^\top \hat{\varphi}_i(x').$$

Letting  $\Phi$  be the matrix with rows  $\psi(x_1)^\top, \dots, \psi(x_n)^\top$ , we find that  $\mathbb{E}(\Phi \Phi^\top)_{ij} = k(x_i, x_j)$ , and that  $\text{Var}(\Phi \Phi^\top)_{ij}$  decreases with speed  $L^{-1}$ . So  $\Phi \Phi^\top$  is a good approximation for  $K$  when  $L$  is moderately large.

One specific example for shift-invariant kernels ( $k(x, x') = h(x - x')$  for all  $x, x' \in \mathcal{X}$ ) is based on the work of Rahimi and Recht.

**Theorem 1.25** (Bochner). *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, for some  $c > 0$  and some distribution  $F$  on  $\mathbb{R}^p$ , it holds for all  $W \sim F$  that*

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

To make use of this theorem, we let  $u \sim U[-\pi, \pi]$ , and it can then be computed (using  $u \stackrel{d}{=} -u$  and  $\cos(u) \stackrel{d}{=} \sin(u)$ ) that

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

Therefore, given a shift-invariant kernel  $k$  with associated distribution  $F$ , we define for  $W \sim F$  the feature map

$$\hat{\varphi}(x) = \sqrt{2c} \cos(W^\top x + u) \in \mathbb{R},$$

so that

$$\mathbb{E}\hat{\varphi}(x)\hat{\varphi}(x') = 2c\mathbb{E}[\mathbb{E}[\cos(W^\top x + u)\cos(W^\top x' + u) \mid W]] = c\mathbb{E}\cos((x - x')^\top W) = k(x, x').$$

**Example 1.26.** Let  $k$  be the Gaussian kernel  $k(x, x') = e^{-\|x-x'\|^2/(2\sigma^2)}$ . Let  $W \sim N(0, \sigma^{-2}I)$ , then using the characteristic function of a normal distribution we find that

$$\mathbb{E}\left(e^{i(x-x')^\top W}\right) = \exp\left(-\|x-x'\|^2/(2\sigma^2)\right) = k(x, x'),$$

so we can set  $c = 1$ , i.e., define the feature map  $\hat{\varphi}(x) = \sqrt{2} \cos(W^\top x + u)$ .

## 2 The Lasso

### 2.1 Model selection

We go back to the linear model  $Y = X\beta^0 + \varepsilon$  with  $\mathbb{E}[\varepsilon] = 0$  and  $\text{Var}(\varepsilon) = \sigma^2 I$ . Using the trace trick, one can easily compute that the MSPE of the OLS estimator is given by

$$\frac{1}{n} \mathbb{E} \left\| X\beta^0 - X\hat{\beta}^{\text{OLS}} \right\|_2^2 = \frac{\sigma^2 p}{n}.$$

Defining  $S = \{k \mid (\beta^0)_k \neq 0\}$ , there is often reason to assume that  $S$  is small, i.e.,  $s := |S| \ll p$ . If we could fit a model using only the variables in  $S$ , the MSPE would be much  $\frac{\sigma^2 s}{n} \ll \frac{\sigma^2 p}{n}$ .

**Best subset selection** A natural way to find  $S$  is to consider all possible subsets of  $\{1, \dots, p\}$ , and pick the best regression procedure using, for example, cross-validation. However, this can become computationally infeasible for moderately large  $p$  (say  $p \approx 10$ ).

**Forward selection** This is a greedy way of performing best subset regression. Given a target model size  $m$ , we first compute the intercept-only model  $M_0$ , and then one-by-one add the predictor variable that reduces the residual sum of squares the most, until we have a model with  $m$  variables.

### 2.2 Lasso estimator

The *Least absolute shrinkage and selection operator* or *Lasso* is given by

$$(\hat{\mu}_\lambda^L, \hat{\beta}_\lambda^L) := \arg \min_{(\mu, \beta) \in \mathbb{R} \times \mathbb{R}^p} \frac{1}{2n} \|Y - \mu \mathbf{1} - X\beta\|_2^2 + \lambda \|\beta\|_1.$$

As with ridge regression, we usually centre and scale the matrix  $X$ , as well as centre the responses  $Y$ , in which case we find

$$\hat{\beta}_\lambda^L = \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2n} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_1.$$

The main difference between the lasso estimator and the ridge regression estimator is that it is likely that the lasso estimator has some zero components. This means that the lasso estimator also estimates which variables are relevant.

#### 2.2.1 Slow prediction error rate

We have the following bound on the prediction error:

**Theorem 2.1** (Slow rate). *Assume  $X$  has centred and scaled columns, and assume that  $Y$  has been centred, so  $Y = X\beta^0 + \varepsilon - \bar{\varepsilon} \mathbf{1}$ . Let  $A > 0$  and suppose*

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

*Let  $\hat{\beta} = \hat{\beta}_\lambda^L$ , then with probability at least  $1 - 2p^{-(A^2/2-1)}$  we have that the MSPE satisfies*

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

*Proof.* By definition we have

$$\frac{1}{2n} \left\| Y - X\hat{\beta} \right\|_2^2 + \lambda \|\hat{\beta}\|_1 \leq \frac{1}{2n} \left\| Y - X\beta^0 \right\|_2^2 + \lambda \|\beta^0\|_1,$$



and rearranging the terms gives

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

By Hölder's inequality we have  $\left| \varepsilon^\top X(\hat{\beta} - \beta^0) \right| \leq \left\| X^\top \varepsilon \right\|_\infty \left\| \hat{\beta} - \beta^0 \right\|_1$ . Define the event  $\Omega = \{ \left\| X^\top \varepsilon \right\|_\infty / n \leq \lambda \}$ , then conditional on  $\Omega$  we find

$$\frac{1}{n} \left\| X(\beta^0 - \hat{\beta}) \right\|_2^2 \leq 2\lambda \left( \left\| \beta^0 - \hat{\beta} \right\|_1 + \left\| \beta^0 \right\|_1 - \left\| \hat{\beta} \right\|_1 \right) \leq 4\lambda \left\| \beta^0 \right\|_1.$$

In lemma 2.5, we will show that  $\mathbb{P}(\Omega) \geq 1 - 2p^{-(A^2/2-1)}$ , which completes the proof.  $\square$

### 2.2.2 Concentration inequalities

Let  $W$  be any random variable and  $\varphi: \mathbb{R} \rightarrow [0, \infty)$  strictly increasing. Then by Markov's inequality we have

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

Plugging in  $\varphi(x) = e^{\alpha x}$  (for some  $\alpha > 0$ ), we get

$$\mathbb{P}(W \geq t) \leq e^{-\alpha t} \mathbb{E}[e^{\alpha W}] = e^{-\alpha t} M_W(\alpha).$$

Now we can take the infimum over all  $\alpha$  on the right-hand side, and we get what is called the *Chernoff bound*:

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

**Definition 2.2.** A random variable  $W$  with mean  $\mu$  is called *sub-Gaussian with parameter  $\sigma > 0$*  or  $\sigma$ -sub-Gaussian if

$$M_{W-\mu} \leq M_{N(0, \sigma^2)} \quad \text{or equivalently} \quad \mathbb{E}[e^{\alpha(W-\mu)}] \leq e^{\alpha^2 \sigma^2 / 2} \text{ for all } \alpha \in \mathbb{R}.$$

We need the following lemma, which characterises an important class of sub-Gaussian random variables:

**Lemma 2.3** (Hoeffding). *If  $W$  is a mean-zero random variable which takes values in  $[a, b]$ , then  $W$  is sub-Gaussian with parameter  $(b - a)/2$ .*

By Chernoff bounding, we obtain for a  $\sigma$ -sub-Gaussian random variable  $W$  that

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

**Proposition 2.4.** *Let  $W_1, \dots, W_n$  be independent, mean-zero random variables where  $W_i$  is  $\sigma_i$ -sub-Gaussian. For any  $\gamma \in \mathbb{R}^n$ , the random variable  $\gamma^\top W$  is sub-Gaussian with parameter  $(\sum_i \sigma_i^2 \gamma_i^2)^{1/2}$ .*

*Proof.* Since the  $W_i$  are independent we have for all  $\alpha \in \mathbb{R}$  that

$$\mathbb{E}[e^{\alpha \sum_i \gamma_i W_i}] = \prod_{i=1}^n \mathbb{E}[e^{\alpha \gamma_i W_i}] \leq \prod_{i=1}^n e^{\alpha^2 \gamma_i^2 \sigma_i^2 / 2} = e^{\alpha^2 \sum_i \gamma_i^2 \sigma_i^2 / 2}.$$

$\square$

**Lemma 2.5.** *Suppose  $\varepsilon_1, \dots, \varepsilon_n$  are independent mean-zero  $\sigma$ -sub-Gaussian random variables and let  $\lambda := A\sigma\sqrt{\log(p)}/n$ . Then*

$$\mathbb{P}\left(\frac{\left\| X^\top \varepsilon \right\|_\infty}{n} \leq \lambda\right) \geq 1 - 2p \exp\left(\frac{-\lambda^2}{2(\sigma/\sqrt{n})^2}\right) = 1 - 2p^{-(A^2/2-1)}.$$

*Proof.* We have

$$\mathbb{P}\left(\frac{\|X^\top \varepsilon\|_\infty}{n} > \lambda\right) = \mathbb{P}\left(\bigcup_i \frac{|X_i^\top \varepsilon|}{n} > \lambda\right) \leq \sum_i \mathbb{P}\left(\frac{|X_i^\top \varepsilon|}{n} > \lambda\right).$$

By the previous proposition, both  $X_i^\top \varepsilon$  and  $-X_i^\top \varepsilon$  are sub-Gaussian with parameter  $\sigma/\sqrt{n}$ , so we have

$$\sum_i \mathbb{P}\left(\frac{|X_i^\top \varepsilon|}{n} > \lambda\right) \leq 2p \exp\left(\frac{-\lambda^2}{2(\sigma/\sqrt{n})^2}\right) = 2p \exp(-A^2 \log(p)/2) = 2p^{-(A^2/2-1)}.$$

□

We now move to finding tail bounds for products of sub-Gaussian random variables. To this end, we consider Bernstein's inequality.

**Definition 2.6.** A random variable  $W$  with mean  $\mu$  is said to satisfy *Bernstein's condition* with parameters  $\sigma > 0, b > 0$  if

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

**Proposition 2.7** (Bernstein's inequality). *Let  $W_1, W_2, \dots$  be independent random variables with mean  $\mu$ , where each  $W_i$  satisfies Bernstein's condition with parameter  $(\sigma, b)$ . Then*

$$1. \mathbb{E}[\exp(\alpha(W_i - \mu))] \leq \exp\left(\frac{\alpha^2 \sigma^2}{2(1-b|\alpha|)}\right) \text{ for all } |\alpha| < b^{-1};$$

$$2. \mathbb{P}(\overline{W} - \mu \geq t) \leq \exp\left(\frac{-nt^2}{2(\sigma^2 + bt)}\right) \text{ for all } t \geq 0.$$

*Proof.* See notes from topics in statistical theory. □

**Lemma 2.8.** *Let  $W, Z$  be mean-zero, sub-Gaussian random variables with parameters  $\sigma_W, \sigma_Z$ . Then  $WZ$  satisfies the Bernstein condition with parameters  $(8\sigma_W\sigma_Z, 4\sigma_W\sigma_Z)$ .*

*Proof.* Using the tail probability formula for expectations one can prove that any mean-zero sub-Gaussian random variable  $X$  satisfies  $\mathbb{E}[X^{2k}] \leq 2^{k+1} \sigma_X^{2k} k!$ .

Furthermore for any random variable  $Y$  we have

$$\begin{aligned} \mathbb{E}|Y - \mathbb{E}Y|^k &\leq \mathbb{E}(|Y| + |\mathbb{E}Y|)^k \\ &= \mathbb{E} \sum_{t=0}^k \binom{k}{t} |Y|^t |\mathbb{E}Y|^{k-t} \\ &\stackrel{*}{\leq} \sum_{t=0}^k \binom{k}{t} \mathbb{E}(|Y|^t) \mathbb{E}(|Y|^{k-t}) \\ &\stackrel{*}{\leq} \sum_{t=0}^k \binom{k}{t} \mathbb{E}(|Y|^k)^{t/k} \mathbb{E}(|Y|^k)^{(k-t)/k} \\ &= \mathbb{E}(|Y|^k) \sum_{t=0}^k \binom{k}{t} = 2^k \mathbb{E}|Y|^k, \end{aligned}$$

where both equations  $*$  are applications of Jensen's inequality.

From the above bound, we obtain

$$\begin{aligned} \mathbb{E}|WZ - \mathbb{E}WZ|^k &\leq 2^k \mathbb{E}(|W|^k |Z|^k) \stackrel{\text{CS}}{\leq} 2^k (\mathbb{E}W^{2k})^{1/2} (\mathbb{E}Z^{2k})^{1/2} \leq 2^{2k+1} \sigma_W^k \sigma_Z^k k! \\ &= \frac{k!}{2} (8\sigma_W\sigma_Z)^2 (4\sigma_W\sigma_Z)^{k-2}, \end{aligned}$$

which proves the claim. □

### 2.2.3 Optimisation theory and convex analysis

#### Lagrangian optimisation

**Definition 2.9.** For the optimisation problem

$$\inf \{f(x) \mid g(x) = 0, h(x) \leq 0\}, \quad (6)$$

its *Lagrangian* is given by

$$L(x, \lambda, \nu) := f(x) + \lambda h(x) + \nu g(x),$$

and its *Lagrange dual function* is

$$\tilde{f}(\lambda, \nu) = \inf_x L(x, \lambda, \nu). \quad (7)$$

Note that if  $p^*$  is the minimum value of (6), then for any feasible  $x$  (i.e.,  $x$  with  $g(x) = 0, h(x) \leq 0$ ) we have for  $\lambda \geq 0$  that

$$f(x) \geq L(x, \lambda, \nu) \geq \tilde{f}(\lambda, \nu).$$

By taking the infimum over  $x$  in the left-hand side and the supremum over  $\lambda, \nu$  in the right-hand side we find that

$$p^* \geq \sup_{\lambda \geq 0, \nu} \tilde{f}(\lambda, \nu),$$

Sometimes we actually have equality here, and in this case we say that the *duality gap* is zero. Note that  $\tilde{f}$  is a concave function, and is therefore easy to maximise.

A sufficient condition for a zero duality gap is that  $L(\cdot, \lambda, \nu)$  is convex and that the problem is feasible.

#### Convex analysis

Let  $\overline{\mathbb{R}} = \mathbb{R} \cup \{\infty\}$ , and define for any  $f: \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$  its (*effective*) *domain*  $\text{dom}(f) := \{x \mid f(x) < \infty\}$ .

**Proposition 2.10.** Let  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  be twice continuously differentiable. Then  $f$  is convex (resp. strictly convex) if and only if its Hessian is positive semi-definite (resp. positive definite) for all  $x \in \mathbb{R}^d$ .

**Definition 2.11.** Let  $f: \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$  be convex and  $x \in \mathbb{R}^d$ . Then  $v \in \mathbb{R}^d$  is called a *subgradient* of  $f$  at  $x$  if

$$f(y) \geq f(x) + v^\top (y - x) \quad \forall y \in \mathbb{R}^d.$$

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

**Proposition 2.12.** If  $f$  is convex and differentiable at  $x \in \text{Int}(\text{dom } f)$ , then  $\partial f(x) = \{\nabla f(x)\}$ .

**Proposition 2.13.** Let  $f, g: \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$  be convex with  $\text{Int}(\text{dom } f) \cap \text{Int}(\text{dom } g) \neq \emptyset$ ,  $x \in \mathbb{R}^d$  and  $\alpha > 0$ . Then  $\partial(\alpha f)(x) = \alpha \partial f(x)$  and  $\partial(f + g)(x) = \partial f(x) + \partial g(x)$ .

**Proposition 2.14** (KKT conditions). Let  $f: \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$  be convex. Then  $x^*$  is a global minimiser of  $f$  if and only if  $0 \in \partial f(x^*)$ .

*Proof.* Clearly  $x^*$  is a global minimiser of  $f$  if and only if  $f(y) \geq f(x^*) + 0^\top (y - x)$  for all  $y \in \mathbb{R}^d$ .  $\square$

We define for  $x \in \mathbb{R}$

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

and for  $x \in \mathbb{R}^d$  we define  $\text{sgn}(x) := (\text{sgn}(x_1), \dots, \text{sgn}(x_n))$ .

**Proposition 2.15.** *Let  $x \in \mathbb{R}^d$  and  $A := \{j \mid x_j \neq 0\}$ . Then*

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

*Proof.* Define  $g_j(x) = |x_j|$  so that  $\|\cdot\|_1 = \sum_j g_j$  and  $\partial\|x\|_1 = \sum_j \partial g_j(x)$ . Now, if  $x_j \neq 0$ , then  $g_j$  is differentiable at  $x$  so  $\partial g_j(x) = \{\text{sgn}(x_j)e_j\}$ . If  $x_j = 0$ , then it can be shown that

$$\partial g_j(x) = \{ce_j \mid c \in [-1, 1]\}.$$

Combining the above facts proves the claim.  $\square$

## 2.2.4 Lasso solutions

With the previous proposition, we can write characterise solutions to the Lasso: namely we have

$$\hat{\beta}_\lambda^L = \arg \min_{\beta} \frac{1}{2n} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_1 = \arg \min_{\beta} Q(\lambda)$$

if and only if  $0 \in \partial Q(\hat{\beta}_\lambda^L)$ . A simple computation gives, with  $A = \{k \mid \beta_k \neq 0\}$ ,

$$\begin{aligned} \partial Q(\beta) &= \partial\left(\frac{1}{2n} \|Y - X\beta\|_2^2\right) + \partial(\lambda \|\beta\|_1) \\ &= \left\{-\frac{1}{n} X^\top (Y - X\beta)\right\} + \{\lambda \nu : \|\nu\|_\infty \leq 1, \nu_A = \text{sgn}(\beta_A)\}. \end{aligned}$$

We conclude that  $\hat{\beta}$  is a Lasso solution if and only if there exists a vector  $\hat{\nu}$  with  $\|\hat{\nu}\|_\infty \leq 1$  and  $\hat{\nu}_A = \text{sgn}(\hat{\beta}_A)$  such that

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

Note that  $\hat{\beta}_\lambda^L$  need not be unique, however, the fitted values are unique:

**Proposition 2.16.**  *$X\hat{\beta}_\lambda^L$  is unique.*

*Proof.* Write  $\frac{1}{2n} \|Y - X\beta\|_2^2 = Q_1(X\beta)$  and  $\|\beta\|_1 = Q_2(\beta)$ . Note that  $Q_1$  is strictly convex. Now, for any  $t \in (0, 1)$  and any two lasso solutions  $\hat{\beta}, \hat{\gamma}$  with  $Q(\hat{\beta}) = Q(\hat{\gamma}) = c^*$  we have

$$\begin{aligned} c^* &\leq Q(t\hat{\beta} + (1-t)\hat{\gamma}) = Q_1(tX\hat{\beta} + (1-t)X\hat{\gamma}) + Q_2(t\hat{\beta} + (1-t)\hat{\gamma}) \\ &\stackrel{*}{\leq} tQ_1(X\hat{\beta}) + (1-t)Q_2(X\hat{\gamma}) + Q_2(t\hat{\beta} + (1-t)\hat{\gamma}) \\ &\leq tQ(\hat{\beta}) + (1-t)Q(\hat{\gamma}) = c^*, \end{aligned}$$

so all inequalities must be equalities. In particular, by  $\star$  we find  $Q_1(tX\hat{\beta} + (1-t)X\hat{\gamma}) = tQ_1(X\hat{\beta}) + (1-t)Q_1(X\hat{\gamma})$ , and since  $Q_1$  is strictly convex this implies  $X\hat{\beta} = X\hat{\gamma}$ .  $\square$

## 2.2.5 Variable selection

We consider the noiseless model  $Y = X\beta^0$ , where no randomness is involved. Let  $S = \{k \mid \beta_k \neq 0\}$ ,  $N = \{1, \dots, p\} \setminus S$ , and assume that  $S = \{1, \dots, s\}$  (so  $\beta_1, \dots, \beta_s$  are nonzero and  $\beta_{s+1}, \dots, \beta_p$  are zero).

Furthermore, write  $X = [X_S \ X_N]$  and assume that  $\text{rank}(X_S) = s$ , so that  $X_S^\top X_S$  is invertible (NB in particular this implies  $n \geq s$ ).

**Theorem 2.17.** *Let  $\lambda > 0$  and define  $\Delta = X_N^\top X_S (X_S^\top X_S)^{-1} \text{sgn}(\beta_S^0)$ .*

1. If  $\|\Delta\|_\infty \leq 1$  and for  $k = 1, \dots, s$  we have

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

then there exists a Lasso solution  $\hat{\beta}_\lambda^L$  with  $\text{sgn}(\hat{\beta}_\lambda^L) = \text{sgn}(\beta^0)$ .

2. If there exists a Lasso solution  $\hat{\beta}_\lambda^L$  with  $\text{sgn}(\hat{\beta}_\lambda^L) = \text{sgn}(\beta^0)$ , then  $\|\Delta\|_\infty = 1$ .

*Remark.* The condition  $\|\Delta\|_\infty \leq 1$  is called the *irrepresentable condition*, and it roughly states that none of the insignificant predictors align “too well” with the response. **TODO:** Elaborate on this (lecture 12).

*Proof.* Let  $\hat{\beta}$  be any lasso solution and write  $\hat{S} = \{k \mid \beta_k \neq 0\}$ . Using the fact that  $\beta_N^0 = 0$ , the KKT conditions for the Lasso can be expanded as

$$\frac{1}{n} \begin{pmatrix} X_S^\top X_S & X_S^\top X_N \\ X_N^\top X_S & X_N^\top X_N \end{pmatrix} \begin{pmatrix} \beta_S^0 - \hat{\beta}_S \\ -\hat{\beta}_N \end{pmatrix} = \lambda \begin{pmatrix} \hat{\nu}_S \\ \hat{\nu}_N \end{pmatrix} \quad (8)$$

for some  $\hat{\nu}$  with  $\|\hat{\nu}\|_\infty \leq 1$  and  $\hat{\nu}_{\hat{S}} = \text{sgn}(\hat{\beta}_{\hat{S}})$ .

1. It is easily checked that taking

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

is a Lasso solution with  $\text{sgn}(\hat{\beta}_S^0) = \text{sgn}(\hat{\beta}_S)$ .

2. If  $\text{sgn}(\hat{\beta}) = \text{sgn}(\beta^0)$ , then  $\hat{S} = S$ ,  $\hat{\nu}_S = \text{sgn}(\beta_S^0)$  and  $\hat{\beta}_N = 0$ . The top block of eq. (8) can be written as

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

Plugging that into the bottom block of the same equation, we get

$$\lambda \frac{1}{n} X_N^\top X_S \left[ \frac{1}{n} X_S^\top X_S \right]^{-1} \text{sgn}(\beta_S^0) = \lambda \hat{\nu}_N \implies \|\Delta\|_\infty \leq 1.$$

□

### 2.2.6 Prediction and estimation

**Definition 2.18.** Let  $X \in \mathbb{R}^{n \times p}$  be a design matrix and  $S \subseteq \{1, \dots, p\}$  where  $s := |S| > 0$ , then the *compatibility constant*  $\varphi^2 = \varphi^2(X, S)$  is defined as

$$\varphi^2 := \inf \left\{ \frac{\frac{1}{n} \|X\beta\|_2^2}{\frac{1}{s} \|\beta_S\|_1^2} : \beta \in \mathbb{R}^p, \beta_S \neq 0, \|\beta_N\|_1 \leq 3\|\beta_S\|_1 \right\}.$$

The *compatibility condition* is that  $\varphi^2 > 0$ .

Note that if  $X_S$  is rank-deficient, then there exists  $\beta_S \neq 0$  with  $X_S \beta_S = 0$ , and by setting  $\beta_N = 0$  we find that  $\varphi^2 = 0$ . Therefore,  $X_S$  must have full column rank for the compatibility condition to be satisfied, and in particular we must have  $n \geq s$ .

**Proposition 2.19.** *If  $X^\top X/n$  has minimal eigenvalue  $c_{\min}$  then  $\varphi^2 \geq c_{\min}$ . In particular, if  $X$  has full column rank (i.e.,  $c_{\min} > 0$ ), the compatibility condition is always satisfied.*

*Proof.* We have

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

Therefore, we find

$$\varphi^2 \geq \inf_{\beta \neq 0} \frac{\frac{1}{n} \|X\beta\|_2^2}{\frac{1}{s} \|\beta\|_1^2} \geq \frac{1}{n} \inf_{\beta \neq 0} \frac{\|X\beta\|_2^2}{\|\beta\|_2^2} = \inf_{\|\beta\|=1} \beta^\top \left( \frac{X^\top X}{n} \right) \beta = c_{\min}.$$

□

Note that the converse of this proposition does not necessarily hold:  $X$  does not need to have full column rank to satisfy the compatibility condition. In particular, it does not need to be the case that  $n < p$ .

Now we present the “fast” convergence rate of the lasso (compare with theorem 2.1):

**Theorem 2.20** (Fast rate). *Let the conditions of theorem 2.1 be satisfied and assume furthermore that the compatibility condition holds. Then with probability at least  $1 - 2p^{-(A^2/8-1)}$  we have*

$$\frac{1}{n} \|X(\beta^0 - \hat{\beta})\|_2^2 + \lambda \|\hat{\beta} - \beta^0\|_1 \leq \frac{16\lambda^2 s}{\varphi^2} = \frac{16A^2 \log(p) \sigma^2 s}{\varphi^2 n}.$$

*Proof.* As in the proof of theorem 2.1, we have

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

Define the event  $\Omega = \{2\|X^\top \varepsilon\|_\infty / n \leq \lambda\}$ , then by Hölder’s inequality we get, conditional on  $\Omega$ ,

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

By lemma 2.5, we find that  $\mathbb{P}(\Omega) \geq 1 - 2p^{-(A^2/8-1)}$ .

Dividing by  $\lambda$ , we have

$$\begin{aligned} \frac{1}{\lambda n} \|X(\hat{\beta} - \beta^0)\|_2^2 + 2\|\hat{\beta}\|_1 &\leq \|\hat{\beta} - \beta^0\|_1 + 2\|\beta^0\|_1 \\ \frac{1}{\lambda n} \|X(\hat{\beta} - \beta^0)\|_2^2 + 2(\|\hat{\beta}_S\|_1 + \|\hat{\beta}_N\|_1) &\leq \|\hat{\beta}_S - \beta_S^0\|_1 + \|\hat{\beta}_N\|_1 + 2\|\beta_S^0\|_1 \\ \frac{1}{\lambda n} \|X(\hat{\beta} - \beta^0)\|_2^2 + \|\hat{\beta}_N\|_1 &\leq \|\hat{\beta}_S - \beta_S^0\|_1 - 2\|\hat{\beta}_S\|_1 + 2\|\beta_S^0\|_1 \\ &\leq 3\|\hat{\beta}_S - \beta_S^0\|_1. \end{aligned}$$

Now we note that  $\|\hat{\beta}_N\|_1 = \|\hat{\beta}_N + \beta_N^0\|_1$ , and we add  $\|\hat{\beta}_S - \beta_S^0\|_1$  to both sides of the inequality to get

$$\frac{1}{\lambda n} \|X(\hat{\beta} - \beta^0)\|_2^2 + \|\hat{\beta} - \beta^0\|_1 \leq 4\|\beta_S - \beta_S^0\|_1. \quad (9)$$

Assume that  $\|\hat{\beta} - \beta^0\| \neq 0$  (otherwise the theorem is definitely true), and note that  $\|(\hat{\beta} - \beta^0)_N\|_1 \leq 3\|\hat{\beta}\|_1 \leq 3\|(\hat{\beta} - \beta^0)_S\|_1$  by our previous inequality, we find that the compatibility constant satisfies

$$\varphi^2 \leq \frac{\frac{1}{n} \|X(\hat{\beta} - \beta^0)\|_2^2}{\frac{1}{s} \|\hat{\beta} - \beta^0\|_1^2} \implies \frac{1}{\varphi} \geq \sqrt{\frac{n}{s}} \frac{\|\hat{\beta} - \beta^0\|_1}{\|X(\hat{\beta} - \beta^0)\|_2},$$

and plugging this into eq. (9) we get

$$\frac{1}{\lambda n} \left\| X(\hat{\beta} - \beta^0) \right\|_2^2 + \left\| \hat{\beta} - \beta^0 \right\|_1 \leq 4 \left\| \beta_S - \beta_S^0 \right\|_1 \leq \frac{4}{\varphi} \sqrt{\frac{s}{n}} \left\| X(\hat{\beta} - \beta^0) \right\|_2, \quad (10)$$

or equivalently

$$\frac{1}{\lambda n} \left\| X(\hat{\beta} - \beta^0) \right\|_2^2 - \frac{4}{\varphi} \sqrt{\frac{s}{n}} \left\| X(\hat{\beta} - \beta^0) \right\|_2 + \left\| \hat{\beta} - \beta^0 \right\|_1 \leq 0.$$

This is a quadratic inequality in  $\left\| X(\hat{\beta} - \beta^0) \right\|_2$ , and solving it yields

$$\frac{1}{\sqrt{n}} \left\| X(\hat{\beta} - \beta^0) \right\|_2 \leq \frac{4\lambda\sqrt{s}}{\varphi}.$$

Plugging that into the right hand side of eq. (10) yields the required result.  $\square$

### The compatibility condition and random design

We will show that if  $X$  is a random matrix, then the compatibility condition holds with high probability.

Define, for  $\Sigma \in \mathbb{R}^{p \times p}$  and  $S \subseteq \{1, \dots, p\}$ , the quantity

$$\varphi_\Sigma^2(S) = \inf_{\beta_S \neq 0, \|\beta_N\|_1 \leq 3\|\beta_S\|_1} \frac{\beta^\top \Sigma \beta}{\frac{1}{|S|} \|\beta_S\|_1^2},$$

so that the compatibility constant satisfies  $\varphi^2 = \varphi_\Sigma^2(S)$  with  $\hat{\Sigma} := X^\top X/n$  and  $S$  the support of  $\beta^0$ .

To prove our main result, we first need a lemma:

**Lemma 2.21.** *Suppose  $\varphi_\Sigma^2(S) > 0$  and  $\max_{jk} |\hat{\Sigma}_{jk} - \check{\Sigma}_{jk}| \leq \varphi_\Sigma^2(S)/(32|S|)$ . Then  $\varphi_{\hat{\Sigma}}^2(S) \geq \varphi_\Sigma^2(S)$ .*

*Proof.* Let  $s = |S|$  and  $t := \varphi_\Sigma^2/(32s)$ . Then by applying Hölder's inequality we obtain

$$\left| \beta^\top (\hat{\Sigma} - \check{\Sigma}) \beta \right| \leq \|\beta\|_1 \left\| (\hat{\Sigma} - \check{\Sigma}) \beta \right\|_\infty \leq \|\beta\|_1^2 \max_{j,k} |\hat{\Sigma}_{jk} - \check{\Sigma}_{jk}| \leq \|\beta\|_1^2 t.$$

Furthermore, if  $\|\beta_N\|_1 \leq 3\|\beta_S\|_1$ , then

$$\|\beta\|_1 = \|\beta_N\|_1 + \|\beta_S\|_1 \leq 4\|\beta_S\|_1 \leq \frac{4\sqrt{\beta^\top \hat{\Sigma} \beta}}{\varphi_{\hat{\Sigma}}/s} \implies \|\beta\|_1^2 \leq \frac{16\beta^\top \hat{\Sigma} \beta}{\varphi_{\hat{\Sigma}}^2/s^2},$$

and therefore

$$\frac{1}{2} \beta^\top \hat{\Sigma} \beta = \beta^\top \check{\Sigma} \beta - \frac{\varphi_{\hat{\Sigma}}^2}{32s} \cdot \frac{16\beta^\top \hat{\Sigma} \beta}{\varphi_{\hat{\Sigma}}^2/s^2} \leq \beta^\top \check{\Sigma} \beta - t \|\beta\|_1^2 \leq \beta^\top \check{\Sigma} \beta - \beta^\top (\hat{\Sigma} - \check{\Sigma}) \beta = \beta^\top \hat{\Sigma} \beta.$$

Taking the infimum over all  $\beta$  with  $\|\beta_N\|_1 \leq 3\|\beta_S\|_1$  and  $\beta_S \neq 0$ , we obtain the required result.  $\square$

Now we can state the main theorem. We consider a sequence of models with sample size  $n$  (where  $n \rightarrow \infty$ ), and with  $p$  and  $s$  increasing to  $\infty$  as functions of  $n$ .

**Theorem 2.22.** *Suppose the rows of  $X$  are i.i.d. and each entry is a mean-zero sub-Gaussian random variable with parameter  $v$ . Let  $\Sigma^0 := \mathbb{E}[X^\top X/n]$ , and suppose that  $s\sqrt{\log(p)/n} \rightarrow 0$  as  $n \rightarrow \infty$ . Define*

$$\varphi_{\Sigma^0, s}^2 := \min_{|S|=s} \varphi_{\Sigma^0}^2(S), \quad \varphi_{\Sigma^0, s}^2 := \min_{|S|=s} \varphi_{\Sigma^0}^2(S).$$

*If  $\varphi_{\Sigma^0, s}^2 > c > 0$  (for all  $n$ ), then  $\mathbb{P}(\varphi_{\Sigma^0, s}^2 \geq \varphi_{\Sigma^0, s}^2/2) \rightarrow 1$  as  $n \rightarrow \infty$ .*

*Proof.* By the previous lemma, it suffices to show that

$$(*) = \mathbb{P}(\max_{jk} |\hat{\Sigma}_{jk} - \Sigma_{jk}^0| \geq \varphi_{\Sigma^0, s}^2 / (32s)) \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

Letting  $t := \varphi_{\Sigma^0, s}^2 / (32s)$ , then by a union bound and Bernstein's theorem, we have that

$$\begin{aligned} (*) &\leq p^2 \max_{j,k} \mathbb{P}\left(\left|\sum_{i=1}^n X_{ij} X_{ik} / n - \Sigma_{jk}^0\right| \geq t\right) \\ &\stackrel{*}{\leq} p^2 2 \exp\left(-\frac{nt^2}{2(64v^4 + 4v^2t)}\right) \\ &\leq c_1 \exp\left(-c_2 \frac{n}{s^2} + c_3 \log(p)\right) \rightarrow 0 \quad \text{if } s\sqrt{\log(p)/n} \rightarrow 0 \text{ as } n \rightarrow \infty. \end{aligned}$$

where  $\star$  follows from the fact that  $X_{ij}X_{ik}$  satisfies Bernstein's condition with parameters  $(8v^2, 4v^2)$ .  $\square$

What does the previous theorem tell us? We have a weak growth condition: for example, the theorem holds if  $p = n^d$  and  $s = (\sqrt{n})^{1-\alpha}$  for any  $d, \alpha > 0$ .

**Example 2.23.** Suppose the rows of  $X$  are taken from a  $N_p(0, \Sigma^0)$  distribution, where  $\Sigma^0$  has bounded diagonal entries and the smallest eigenvalue of  $\Sigma^0$  is  $c_{\min} > c > 0$  (for all  $n$ ). Then  $\mathbb{P}(\varphi_{\Sigma, s}^2 \geq \frac{c_{\min}}{2}) \rightarrow 1$  as  $n \rightarrow \infty$  if the growth condition is satisfied. This holds for example if  $\Sigma^0 = I_p$ .

**Example 2.24.** Suppose  $X$  has i.i.d. entries  $\mathbb{P}(X_{ij} = 1) = \frac{1}{2} = \mathbb{P}(X_{ij} = -1)$ . Then by Hoeffding's lemma,  $X_{ij}$  is sub-Gaussian, and therefore the theorem holds under appropriate conditions.

**Example 2.25.** Let  $W$  be a matrix whose rows are vectors in a basis, and let  $X$  have rows which are random subsets of  $W$ . This does not satisfy the conditions of theorem 2.22, yet the lasso still works very well in this case.

Generally, the Lasso demonstrates a very low estimation if  $n/p$  is much larger than  $s/p$ , but a very high estimation error if  $s/p$  is much larger than  $n/p$ . Usually, there is a sharp “phase transition” here.

## 2.2.7 Computation

The *coordinate descent* algorithm is an algorithm for minimising functions of the form  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ , and works particularly well for functions of the form

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

where  $g$  is convex and differentiable and the  $h_j: \mathbb{R} \rightarrow \mathbb{R}$  are convex and continuous. Note that the Lasso objective function is clearly of this form. We start with an initial guess  $x^{(0)}$  for a minimiser, and then for  $m = 1, 2, \dots$  we choose, for  $i = 1, \dots, d$ ,

$$x_i^{(m)} = \arg \min_{x_i \in \mathbb{R}} f\left(x_1^{(m)}, \dots, x_{i-1}^{(m)}, x_i, x_{i+1}^{(m-1)}, \dots, x_d^{(m-1)}\right).$$

It can be shown that for the Lasso objective function, each step of the coordinate descent algorithm admits a closed-form solution.

We have the following:

**Theorem 2.26** (Tseng). *If  $A_0 := \{x \mid f(x) \leq f(x^0)\}$  is compact, then every converging subsequence of  $(x^{(m)})$  converges to a minimiser of  $f$ .*



Using this theorem, we can see the following:

**Corollary 2.27.** *Suppose  $A^0$  is compact. Then*

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

*Proof.* Since  $f$  is continuous, it attains a minimum  $f(x^*)$  on  $A_0$ . Suppose  $f(x^{(m)}) \not\rightarrow f(x^*)$ . Note that  $f(x^{(m)})$  is a monotone decreasing sequence, so it must converge. Furthermore,  $(x^{(m)}) \subseteq A_0$ , so there exists a subsequence of  $(f(x^{(m)}))$  converging to  $f(x^*)$  by theorem 2.26. This proves that  $f(x^{(m)}) \rightarrow f(x^*)$ .

Suppose that  $x^*$  is the unique minimiser of  $f$ , then by theorem 2.26 we know that every convergent subsequence of  $x^{(m)}$  converges to  $x^*$ . Since  $x^{(m)}$  is bounded, this proves that  $x^{(m)} \rightarrow x^*$  as well.  $\square$

Often, we want to solve the Lasso on a grid of values  $\lambda_0 > \dots > \lambda_L$ . To speed up the coordinate descent procedure, we can use the minimiser for  $\lambda_{\ell-1}$  as a starting point for coordinate descent on  $\lambda_\ell$ . An “active set” strategy can speed up computation even further:

1. Initialise  $A_\ell := \{k \mid \beta_{\ell-1,k}^L \neq 0\}$ ;
2. Perform coordinate descent with  $\lambda_\ell$  only on the coordinates in  $A_\ell$ , call this result  $\hat{\beta}$ .
3. Let  $V$  be the set of coordinates which violate the KKT conditions;
4. If  $V = \emptyset$ , set  $\hat{\beta}_{\lambda_\ell}^L = \hat{\beta}$ . Else, update  $A_\ell = A_\ell \cup V$  and return to step 2.

*Remark.* There are multiple possible algorithms for computing Lasso solutions, coordinate descent is simply one option.

## 2.3 Extensions of the Lasso

One of the ways to extend the Lasso is to change the loss function (for example, the logistic loss function can be used if  $Y_i \in \{-1, 1\}$ ). There are other extensions.

### 2.3.1 Structural penalties

**Example 2.28** (Group lasso). Let  $G_1, \dots, G_q$  be a partition of  $\{1, \dots, p\}$ , then instead of the Lasso penalty we can use

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

where  $m_j$  is typically chosen as  $\sqrt{|G_j|}$ . This estimator will generally ensure that either the entire group  $G_j$  is shrunk to 0 or not at all.

**Example 2.29** (Fused lasso). Suppose the  $(\beta_j)$  are ordered, for example if  $j$  fixes time. If we expect  $(\beta_j)$  to be piecewise constant (or at least that  $\beta_j^0$  is close to  $\beta_{j+1}^0$ ), we can use the penalty

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

where the term  $\lambda_2 \|\beta\|_1$  can be omitted if shrinkage is not desired.

### 2.3.2 Debiasing the Lasso

A useful property of the Lasso is that non-significant coefficients can be shrunk towards 0, but this same effect also causes many significant coefficients to be shrunk towards 0. A possible solution is to only use the Lasso to estimate which coefficients are 0 and then use a different estimator (i.e., OLS) for the other coefficients.

Another option is to use a non-convex penalty function which does not shrink large values of  $\hat{\beta}$  as much.

### 3 High-dimensional covariance estimation and PCA

Until now, we've focused on regression. In this chapter, we will assume the samples  $x_1, \dots, x_n$  are i.i.d. from some distribution in  $\mathbb{R}^d$ . We are interested in estimating the covariance matrix of  $x_i$ . There are two reasons for this:

1. Covariance can reveal underlying structure relating the variables.
2. When  $d$  is large, interpreting the covariance matrix directly is hard. However, its largest eigenvalues and corresponding eigenvectors can reveal the principal modes of variation in the data.

#### 3.1 Covariance estimation

##### 3.1.1 Maximum likelihood in multivariate normal model

Let  $x_1, \dots, x_n$  be i.i.d. samples from a  $N_d(\mu, \Sigma)$  distribution. We will assume that  $\Sigma$  is invertible with inverse  $\Omega$  (this inverse is sometimes called the *precision matrix*). Furthermore, we define the *sample covariance*

$$\hat{\Sigma} := \frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})(x_i - \bar{X})^\top.$$

It can be shown that the maximum likelihood estimate for  $\mu$  is  $\bar{X} := \frac{1}{n} \sum_{i=1}^n x_i$ , and that the maximum likelihood estimate  $\Omega$  is

$$\Omega = \min_{\Omega > 0} \left( -\log(\det \Omega) + \text{tr}(\hat{\Sigma}\Omega) \right).$$

Since this is the minimisation of a convex function over a convex set, we know that  $\Omega$  is a minimiser if and only if the derivative w.r.t. each element of  $\Omega$  is 0. From this, we can compute that  $(\hat{\Omega}^{\text{ML}})^{-1} = \hat{\Sigma}$ . Therefore, if  $X$  has full column rank so that  $\hat{\Sigma}$  is invertible, we find  $\hat{\Omega}^{\text{ML}} = \hat{\Sigma}^{-1}$ , and since matrix inversion is a bijective map on the set of positive definite matrices, we conclude  $\hat{\Sigma}^{\text{ML}} = \hat{\Sigma}$ .

##### 3.1.2 Non-asymptotic error bounds

In the classical setting, we know that the maximum likelihood estimator is asymptotically optimal if we let  $n \rightarrow \infty$  and keep all other parameters fixed. However, we are interested in *non-asymptotic error bounds*, which do not rely on  $n$  going to  $\infty$ . From these bounds, we can infer much more about the  $n \rightarrow \infty$  case, namely, how much the other parameters are allowed to grow with  $n$ .

We will assume for simplicity that  $\mu = 0$  and that  $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n x_i x_i^\top$ . Recall that the largest (resp. smallest) eigenvalue of a symmetric matrix  $A \in \mathbb{R}^{n \times n}$  is given by the maximum (resp. minimum) of  $v^\top A v$  over all  $v \in S^{n-1}$ . It follows that the operator norm of  $A$  is given by  $\|A\|_{\text{op}} = \max_{v \in S^{n-1}} |v^\top A v|$ .

**Theorem 3.1.** *Let  $X \in \mathbb{R}^{n \times d}$  have i.i.d. rows  $x_1, \dots, x_n$  with mean 0 and covariance  $\Sigma$ , so that*

$$\mathbb{E} \left[ e^{\lambda v^\top x_1} \right] \leq e^{\lambda^2 \sigma^2 / 2} \quad \text{for all } v \in S^{d-1}.$$

*Then there are universal constants  $c, C$  such that  $\hat{\Sigma}$  satisfies*

$$\mathbb{E} \left[ e^{\lambda \|\hat{\Sigma} - \Sigma\|_{\text{op}}} \right] \leq \exp \left( c \left[ \frac{\lambda^2 \sigma^4}{n} + d \right] \right) \quad \text{for all } |\lambda| \leq \frac{n}{16\sigma^2},$$

*and for all  $\delta > 0$ , we have*

$$\mathbb{P} \left( \frac{\|\hat{\Sigma} - \Sigma\|_{\text{op}}}{\sigma^2} \leq C \left( \frac{d + \delta}{n} \vee \sqrt{\frac{d + \delta}{n}} \right) \right) \geq 1 - e^{-\delta}.$$

This is a non-asymptotic error bound. Note that, as long as  $d/n \rightarrow 0$ , the estimator  $\hat{\Sigma}$  will be consistent.

*Proof.* The tail probability bound follows from Chernoff bounding, see lecture 17 for the details.

For the proof of (i), cover  $S^{d-1}$  with balls of radius  $1/8$  centred at points  $v_1, \dots, v_N \in S^{d-1}$ . It is a known geometric fact that this is possible with  $N \leq e^{4d}/2$ . We will prove that for any symmetric  $Q$  we have  $\|Q\|_{\text{op}} \leq 2 \max_i |v_i^\top Q v_i|$ . For this, let  $v \in S^{d-1}$  and choose  $i$  such that  $\|v - v_i\| \leq \frac{1}{8}$ , then

$$\begin{aligned} |v^\top Q v| &= |(v - v_i)^\top Q (v - v_i) + 2v_i^\top Q (v - v_i) - v_i^\top Q v_i| \\ &\leq |(v - v_i)^\top Q (v - v_i)| + 2|v_i^\top Q (v - v_i)| + |v_i^\top Q v_i| \\ &\leq \|v - v_i\|^2 \|Q\|_{\text{op}} + 2\|v_i\| \|Q\|_{\text{op}} \|v - v_i\| + |v_i^\top Q v_i| \\ &\leq \left(\frac{2}{8} + \frac{1}{64}\right) \|Q\|_{\text{op}} + |v_i^\top Q v_i| \leq \frac{1}{2} \|Q\|_{\text{op}} + |v_i^\top Q v_i| \\ &\leq \frac{1}{2} \|Q\|_{\text{op}} + \max_i |v_i^\top Q v_i|. \end{aligned}$$

Since the right-hand side is independent of  $v$  we can maximise the left-hand side over  $v$  and obtain

$$\|Q\|_{\text{op}} \leq \frac{1}{2} \|Q\|_{\text{op}} + \max_i |v_i^\top Q v_i| \implies \|Q\|_{\text{op}} \leq 2 \max_i |v_i^\top Q v_i|.$$

Now, we apply this to  $Q := \hat{\Sigma} - \Sigma$ . Then we have

$$\mathbb{E}[e^{\lambda \|Q\|_{\text{op}}}] \leq \mathbb{E}[e^{2\lambda \max_i |v_i^\top Q v_i|}] \leq \sum_{i=1}^N \mathbb{E}[e^{2\lambda v_i^\top Q v_i}] + \mathbb{E}[e^{-2\lambda v_i^\top Q v_i}], \quad (11)$$

where for the final step we use that  $\max_i |v_i^\top Q v_i| \in \{\pm v_i^\top Q v_i \mid i = 1, \dots, n\}$ .

Next, we claim that for fixed  $u \in S^{d-1}$ , we have  $\mathbb{E}[e^{tu^\top Q u}] \leq \exp\left(64 \frac{t^2 \sigma^4}{n}\right)$  for  $|t| \leq \frac{n}{8\sigma^2}$ . For this, observe that

$$u^\top Q u = \frac{1}{n} \sum_{i=1}^n (u^\top x_i)^2 - u^\top \Sigma u,$$

and note that  $u^\top \Sigma u = \text{Var}(u^\top x_i) = \mathbb{E}[(u^\top x_i)^2]$  since the  $x_i$  have mean zero. Using the fact that the  $x_i$  are independent we have

$$\mathbb{E}[\exp(tu^\top Q u)] = \mathbb{E}\left[\exp\left(\frac{t}{n} [(u^\top x_1)^2 - u^\top \Sigma u]\right)\right]^n = \mathbb{E}[\exp(Z^2 - \mathbb{E}[Z^2])]^n,$$

where  $Z := \sqrt{\frac{t}{n}} u^\top x_1$  is sub-Gaussian with parameter  $\sqrt{|t|\sigma^2/n}$ . Therefore,  $Z^2$  satisfies the Bernstein condition with parameters  $(8|t|\sigma^2/n, 4|t|\sigma^2/n)$ , and recalling that  $|t| \leq \frac{n}{8\sigma^2}$  we obtain

$$\begin{aligned} \mathbb{E}[\exp(Z^2 - \mathbb{E}[Z^2])] &\leq 1 + \sum_{k=2}^{\infty} \frac{\mathbb{E}[(Z^2 - \mathbb{E}[Z^2])^k]}{k!} \\ &\stackrel{*}{\leq} 1 + \sum_{k=2}^{\infty} 2^{2k+1} \sigma^{2k} \left(\frac{|t|}{n}\right)^k \\ &\stackrel{**}{\leq} 1 + \frac{2^5 \sigma^4 t^2}{n} \frac{1}{1 - 4\sigma^2 |t|/n} \\ &\leq 1 + \frac{64 \sigma^4 t^2}{n} \leq e^{64 \sigma^4 t^2 / n}, \end{aligned}$$

where  $\star$  follows from Bernstein's inequality and  $\star\star$  from a geometric series, which proves the claim.

Substituting this back into eq. (11) using  $t = \pm 2\lambda$  shows that for  $|\lambda| < \frac{n}{16\sigma^2}$  we have

$$\mathbb{E}[e^{\lambda\|Q\|_{\text{op}}}] \leq 2Ne^{256\lambda^2\sigma^4/n} \leq e^{256\lambda^2\sigma^4/n+4d} \leq e^{256(\lambda^2\sigma^4/n+d)}.$$

□

The upper bound in this theorem is tight when  $\Sigma = I$ . However, note that our bound in eq. (11) is quite crude, and we can obtain much better rates if we assume that the  $x_i$  are normally distributed (in this case, more is known about the supremum).

Furthermore, when the spectrum of  $\Sigma$  is concentrated on a few directions, the rate can be improved significantly. For this, we need a definition.

**Definition 3.2.** The *effective rank*  $r(\Sigma)$  of a nonzero symmetric positive semi-definite matrix  $\Sigma$  with eigenvalues  $0 \leq \gamma_d(\Sigma) \leq \dots \leq \gamma_1(\Sigma)$  is given by

$$r(\Sigma) = \frac{\text{Tr}(\Sigma)}{\|\Sigma\|_{\text{op}}} = \frac{\sum_{i=1}^d \gamma_i(\Sigma)}{\gamma_1(\Sigma)}.$$

It is easily seen that  $r(\Sigma) \leq \text{rank}(\Sigma)$ , with equality if and only if all nonzero eigenvalues of  $\Sigma$  are equal. The following theorem (which we will not prove) gives a tighter bound if we can assume that the rows of  $X$  are normally distributed and that  $r(\Sigma)$  is constrained:

**Theorem 3.3.** *Suppose  $X$  has i.i.d.  $N(0, \Sigma)$  rows with  $r(\Sigma) \leq cn$  (??). Then there exists  $C$  such that for all  $\delta > 0$  we have*

$$\mathbb{P}\left(\frac{\|\hat{\Sigma} - \Sigma\|_{\text{op}}}{\|\Sigma\|_{\text{op}}} \leq C\left(\frac{r(\Sigma) + \delta}{n} \vee \sqrt{\frac{r(\Sigma) + \delta}{n}}\right)\right) \geq 1 - e^{-\delta}.$$

In comparison with the previous bound, the dimension of the samples  $d$  is replaced by the effective rank  $r(\Sigma)$  (in fact,  $d$  does not occur at all in this theorem!). This means that  $d$  can grow rapidly with  $n$ , as long as  $r(\Sigma)/n \rightarrow 0$ . This can even be generalised to an infinite-dimensional setting with appropriate definitions.