

# Notes on statistical learning

Dávid Iván

December 31, 2020

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Overview of Supervised Learning</b>	<b>3</b>
2.1	Linear models and least squares . . . . .	3
2.2	Statistical decision theory . . . . .	4
2.2.1	application. Simple linear fit. . . . .	4
2.3	$E Y - c $ and the median . . . . .	7
2.3.1	discrete case . . . . .	7
2.3.2	continuous case . . . . .	9
2.4	Local methods in high dimensions . . . . .	11
2.4.1	deriving the prediction formula . . . . .	11
2.4.2	equation (2.47) on page 37 . . . . .	12
2.5	Solutions for the Exercises of chapter 2 . . . . .	13
2.5.1	Ex. 2.2 . . . . .	13
2.5.2	Ex. 2.3 . . . . .	14
2.5.3	Ex. 2.4 . . . . .	14
2.5.4	Ex. 2.5 . . . . .	15
2.5.5	Ex. 2.6 . . . . .	18
2.5.6	Ex. 2.7 . . . . .	18
2.5.7	Ex. 2.9 . . . . .	20
<b>3</b>	<b>Linear Methods for Regression</b>	<b>22</b>
3.1	equations on page 47 and 48 . . . . .	22
3.2	Equation (3.28) on page 54 . . . . .	26
3.3	Solutions for the Exercises of chapter 3 . . . . .	27
3.3.1	Ex. 3.1 . . . . .	27
3.3.2	Ex. 3.2 . . . . .	28
3.3.3	Ex. 3.3 . . . . .	29
3.3.4	Ex. 3.5 . . . . .	31
3.3.5	Ex. 3.6 . . . . .	32
3.3.6	Ex. 3.8 . . . . .	34
3.3.7	Ex. 3.9 <i>Forward stepwise regression</i> . . . . .	35

3.3.8	Ex. 3.10 Backward stepwise regression . . . . .	36
3.3.9	Ex. 3.11 . . . . .	36
<b>Appendices</b>		<b>38</b>
<b>A</b>	<b>differentiation</b>	<b>38</b>
A.1	differentiation w.r.t. a vector . . . . .	38
A.2	differentiation w.r.t. a matrix . . . . .	38
<b>B</b>	<b>variance and covariance properties</b>	<b>40</b>
B.1	scalar multiple . . . . .	40
B.2	vector multiple . . . . .	40
B.3	matrix multiple . . . . .	41
<b>C</b>	<b>distributions</b>	<b>41</b>
C.1	Student's t-distribution . . . . .	41
C.2	F-distribution . . . . .	41
C.3	Gaussian posterior distribution . . . . .	42
<b>D</b>	<b>projections</b>	<b>43</b>
D.1	sum of projections . . . . .	43
D.2	difference of projections . . . . .	44
D.3	The special vector $Xb$ . . . . .	44
<b>E</b>	<b>Matrices</b>	<b>45</b>
E.1	The matrix inversion lemma . . . . .	45

# 1 Introduction

## 2 Overview of Supervised Learning

### 2.1 Linear models and least squares

On page 12 we have that the residual sum of squares:

$$\text{RSS}(\beta) = \sum_{i=1}^N (y_i - x_i^T \beta)^2 = (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) \quad (1)$$

How can we differentiate with respect to  $\beta$ ?

$$\text{RSS}(\beta) = \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\beta - \beta^T \mathbf{X}^T \mathbf{y} + \beta^T \mathbf{X}^T \mathbf{X}\beta \quad (2)$$

Using (139, 140), we can differentiate RSS:

$$\frac{d}{d\beta} \text{RSS}(\beta) = 0 - (\mathbf{y}^T \mathbf{X})^T - \mathbf{X}^T \mathbf{y} + (\mathbf{X}^T \mathbf{X} + (\mathbf{X}^T \mathbf{X})^T) \beta \quad (3)$$

$$\frac{d}{d\beta} \text{RSS}(\beta) = -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X}\beta \quad (4)$$

Setting this to zero we get the normal equations:

$$\mathbf{X}^T \mathbf{y} = \mathbf{X}^T \mathbf{X}\beta \quad (5)$$

## 2.2 Statistical decision theory

On page 19 we have that

$$\beta = [E(XX^T)]^{-1}E(XY) \quad (6)$$

but how exactly do we get this equation? In general, we have the expected prediction error:

$$\text{EPE}(f) = E(Y - f(X))^2 \quad (7)$$

And we have that the prediction function is linear:

$$f(X) = X^T \beta \quad (8)$$

We seek a  $\beta$  for minimizing the expected prediction error.  $X$  and  $Y$  are random variables,  $X$  being a vector,  $Y$  being a scalar.

$$\begin{aligned} \frac{d}{d\beta} \text{EPE} &= \frac{d}{d\beta} E((Y - X^T \beta)^2) = E \left( \frac{d}{d\beta} (Y - X^T \beta)^2 \right) \\ &= E(2(Y - X^T \beta) \cdot (-X)) = -2E(YX) + 2E(X(X^T \beta)) \\ &= -2E(YX) + 2E((XX^T)\beta) = -2E(YX) + 2(E(XX^T))\beta \end{aligned} \quad (9)$$

We used the fact that the expected value is linear, and that  $\beta$  is not random, so we could factor out from the expected value. Setting this to zero we have that:

$$E(YX) = E(XX^T)\beta \quad (10)$$

which yields

$$\hat{\beta} = [E(XX^T)]^{-1}E(YX) \quad (11)$$

### 2.2.1 application. Simple linear fit.

Let's see an application for this equation. Let  $X = \begin{bmatrix} x \\ 1 \end{bmatrix}$ ,  $\beta = \begin{bmatrix} a \\ b \end{bmatrix}$ .

Now  $f(x) = a \cdot x + b$

$$XY = \begin{bmatrix} x \cdot y \\ y \end{bmatrix} \quad (12)$$

$$XX^T = \begin{bmatrix} x^2 & x \\ x & 1 \end{bmatrix} \quad (13)$$

If we have  $N$  datapoints  $\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$ , we can approximate the expectation values.

$$E(XY) \approx \frac{1}{N} \begin{bmatrix} \sum_i x_i \cdot y_i \\ \sum_i y_i \end{bmatrix} \quad (14)$$

$$E(XX^T) \approx \frac{1}{N} \begin{bmatrix} \sum_i x_i^2 & \sum_i x_i \\ \sum_i x_i & N \end{bmatrix} \quad (15)$$

Let's denote the followings:

$$\alpha_X = \sum_i x_i \quad (16)$$

$$\alpha_Y = \sum_i y_i \quad (17)$$

$$\alpha_{XY} = \sum_i x_i y_i \quad (18)$$

$$\alpha_{X^2} = \sum_i x_i^2 \quad (19)$$

With these notations:

$$E(XY) \approx \frac{1}{N} \begin{bmatrix} \alpha_{XY} \\ \alpha_Y \end{bmatrix} \quad (20)$$

$$E(XX^T) \approx \frac{1}{N} \begin{bmatrix} \alpha_{X^2} & \alpha_X \\ \alpha_X & N \end{bmatrix} \quad (21)$$

Inverting  $E(XX^T)$ :

$$[E(XX^T)]^{-1} \approx \frac{N}{N\alpha_{X^2} - \alpha_X^2} \begin{bmatrix} N & -\alpha_X \\ -\alpha_X & \alpha_{X^2} \end{bmatrix} \quad (22)$$

Plug these in to the equation:

$$\hat{\beta} = [E(XX^T)]^{-1} E(YX) \quad (23)$$

$$\hat{\beta} \approx \frac{N}{N\alpha_{X^2} - \alpha_X^2} \begin{bmatrix} N & -\alpha_X \\ -\alpha_X & \alpha_{X^2} \end{bmatrix} \frac{1}{N} \begin{bmatrix} \alpha_{XY} \\ \alpha_Y \end{bmatrix} \quad (24)$$

$$\hat{\beta} \approx \frac{1}{N\alpha_{X^2} - \alpha_X^2} \begin{bmatrix} N\alpha_{XY} - \alpha_X\alpha_Y \\ \alpha_{X^2}\alpha_Y - \alpha_X\alpha_{XY} \end{bmatrix} \quad (25)$$

From here we can get  $\hat{a}$  and  $\hat{b}$ , since  $\hat{\beta} = \begin{bmatrix} \hat{a} \\ \hat{b} \end{bmatrix}$

## 2.3 $E|Y - c|$ and the median

On page 20, it asks the question "What happens if we replace the  $L_2$  loss function with the  $L_1 : E|Y - f(X)|$  ?" Let's investigate this question.

### 2.3.1 discrete case

We can get rid of the conditional  $X = x$ , and just ask the question: What  $c$  will minimize  $E|Y - c|$ ? Denote this function with  $g$ , so  $g(c) = E|Y - c|$ . Let's look at two examples.

Example 1. The random variable  $Y$  takes 4 possible values with probabilities  $\frac{1}{7}, \frac{1}{7}, \frac{3}{7}, \frac{2}{7}$ . The figure below shows the probability mass function.

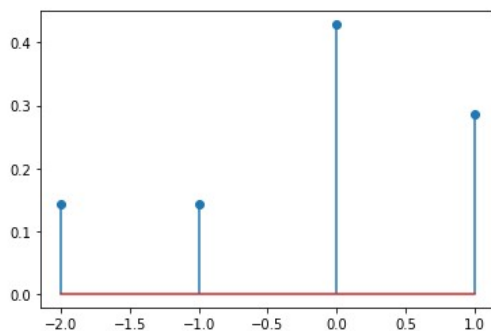


Figure 1: probability mass function of the first example random variable.

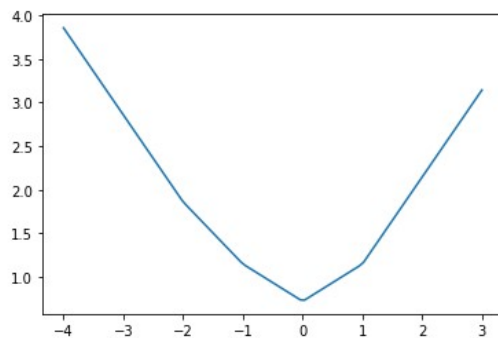


Figure 2:  $g(c)$  function. The horizontal axis is  $c$ .

Example 2. The random variable  $Y$  takes 4 possible values with probabilities 0.1, 0.4, 0.3, 0.2. The figure below shows the probability mass function.

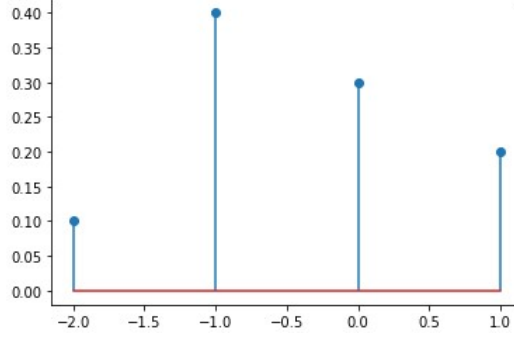


Figure 3: probability mass function of the second example random variable.

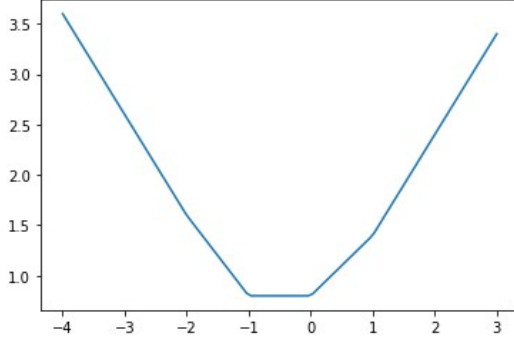


Figure 4:  $g(c)$  function. The horizontal axis is  $c$ .

We can see that  $g(c)$  is a piecewise linear function, and it has minimum, which is a point or a line segment. Let's say we have  $Y$  discrete random variable that takes values from  $S = \{x_1, x_2, \dots, x_n\}$ . The values are ordered:  $x_1 < x_2 < \dots < x_n$ .  $Y$  takes these values with corresponding probabilities  $p_1, p_2, \dots, p_n$ .

Let's calculate the equation of the piecewise linear function. Denote the interval  $I_k$  such that  $x \in I_k$  if and only if  $k$  values from  $S$  are smaller than  $x$ . So  $I_0 = (-\infty, x_1]$ ,  $I_1 = [x_1, x_2]$ , ...,  $I_n = [x_n, \inf)$ .

$$g(c) = E|Y - c| = \sum_{i=1}^n p_i \cdot |x_i - c| \quad (26)$$

If  $c \in I_k$ , then

$$g(c) = E|Y - c| = \sum_{i=1}^k p_i \cdot (c - x_i) + \sum_{i=k+1}^n p_i \cdot (x_i - c) \quad (27)$$

$$g(c) = c \cdot (\sum_{i=1}^k p_i - \sum_{i=k+1}^n p_i) + (\sum_{i=k+1}^n p_i x_i - \sum_{i=1}^k p_i x_i) \quad (28)$$



First we can show that this function is continuous. On one hand ( $c \in I_k = [x_k, x_k + 1]$ ):

$$g_1 = g(x_k) = x_k \cdot (\sum_{i=1}^k p_i - \sum_{i=k+1}^n p_i) + (\sum_{i=k+1}^n p_i x_i - \sum_{i=1}^k p_i x_i) \quad (29)$$

On the other hand we have ( $c \in I_{k-1} = [x_{k-1}, x_k]$ ):

$$g_2 = g(x_k) = x_k \cdot (\sum_{i=1}^{k-1} p_i - \sum_{i=k}^n p_i) + (\sum_{i=k}^n p_i x_i - \sum_{i=1}^{k-1} p_i x_i) \quad (30)$$

$$g_1 - g_2 = x_k \cdot (p_k + p_k) - p_k x_k - p_k x_k = 0 \quad (31)$$

Now that we showed that this function is continuous, let's find it's minimum. Since it is piecewise linear, its derivative is piecewise constant. Denote the derivative of  $g$  on the interval  $I_k$  with  $g'(I_k)$ .

$$\begin{aligned} g'(I_k) &= -1 \\ g'(I_1) &= -1 + 2p_1 \\ g'(I_2) &= -1 + 2p_1 + 2p_2 \\ &\dots \\ g'(I_n) &= -1 + 2p_1 + \dots + 2p_n = 1 \end{aligned} \quad (32)$$

So the derivative is increasing from  $-1$  to  $+1$ . We can distinguish two possibilities. First, assume that the derivative is never zero. In this case, we have a  $k$  where  $g'(I_{k-1}) < 0$  but  $g'(I_k) > 0$ , so the minimum is at  $x_k$ , the median. The second case is where there is an interval where the derivative is zero. In this case the whole interval is minimum, again, the median.

### 2.3.2 continuous case

Let's have the following function:

$$f(x) = \int_a^x g(x, t) dt \quad (33)$$

I state without proof that the derivative of this function is as follows:

$$f'(x) = \int_a^x \frac{\partial g(x, t)}{\partial x} dt + g(x, x) \quad (34)$$

Now we have that

$$g(c) = E(|Y - c| | X = x) = \int_{-\infty}^c (c - y) f_{Y|X}(y|x) dy + \int_c^{\infty} (y - c) f_{Y|X}(y|x) dy \quad (35)$$

$$g'(c) = \int_{-\infty}^c f_{Y|X}(y|x) dy + \int_c^{\infty} f_{Y|X}(y|x) dy \quad (36)$$

Setting this to zero, we get that

$$\int_{-\infty}^c f_{Y|X}(y|x) dy = \int_c^{\infty} f_{Y|X}(y|x) dy \quad (37)$$

$$P(Y < c | X = x) = P(Y > c | X = x) \quad (38)$$

Again, this means the minimum is at the median.

## 2.4 Local methods in high dimensions

### 2.4.1 deriving the prediction formula

On page 24 we see an example of a linear data with noise. At first I was confused how it gets  $\hat{y}_0 = x_0^T \beta + \sum_{i=1}^N l_i(x_0) \epsilon_i$ , where  $l_i(x_0)$  is the  $i$ th element of  $\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} x_0$

In this example we make  $N$  experiments, storing the  $X_i$  values in the rows of  $\mathbf{X}$ , and we have also  $\epsilon_i$  (elements of  $\vec{\epsilon}$ ) and  $Y_i = X_i^T \beta + \epsilon_i$  for some fixed  $\beta$ . For approximating  $\beta$ , we use the result:

$$\beta = [\mathbf{E}(X X^T)]^{-1} \mathbf{E}(Y X) \quad (39)$$

in this case it will be an approximation, since we have noise ( $\epsilon$ ).

Calculate first  $\mathbf{E}(Y X)$ :

$$\mathbf{E}(Y X) = \mathbf{E}((X^T \beta + \epsilon) X) = \mathbf{E}(X X^T) \beta + \mathbf{E}(\epsilon X) \quad (40)$$

Substitute this into the approximation of  $\beta$ :

$$\begin{aligned} \hat{\beta} &= [\mathbf{E}(X X^T)]^{-1} \mathbf{E}(Y X) \\ &= [\mathbf{E}(X X^T)]^{-1} (\mathbf{E}(X X^T) \beta + \mathbf{E}(\epsilon X)) \\ &= \beta + [\mathbf{E}(X X^T)]^{-1} \mathbf{E}(\epsilon X) \end{aligned} \quad (41)$$

We do not know of course the exact expectation values, but we have  $N$  data samples (training data). So how could we approximate the expectation values? Use the averages:

$$\mathbf{E}(\epsilon X)_i \approx \frac{1}{N} \sum_{k=1}^N \mathbf{X}_{ki} \epsilon_k \rightarrow \mathbf{E}(\epsilon X) \approx \frac{1}{N} \mathbf{X}^T \vec{\epsilon} \quad (42)$$

similarly,

$$[\mathbf{E}(X X^T)]^{-1} \approx N \cdot (\mathbf{X}^T \mathbf{X})^{-1} \quad (43)$$

putting these all together, we have:

$$\hat{y}_0 = x_0^T \hat{\beta} = x_0^T \beta + x_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \vec{\epsilon} = x_0^T \beta + \vec{\epsilon}^T \cdot \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} x_0 \quad (44)$$

and this is the formula that was to be explained.

### 2.4.2 equation (2.47) on page 37

$$\begin{aligned}\text{EPE}_k(x_0) &= \text{E} \left( (Y - \hat{f}_k(x_0))^2 | X = x_0 \right) \\ &= \text{E} \left( (f(x_0) + \epsilon_0 - \hat{f}_k(x_0))^2 \right)\end{aligned}\tag{45}$$

The data points are fixed:  $x_1, x_2, \dots, x_N$ . Denote the closest data point to  $x_0$  as  $x_{(1)}$ , the second closest  $x_{(2)}$ , etc. With this notation, the nearest neighbor estimate for  $f(x_0)$ :

$$\hat{f}_k(x_0) = \frac{1}{k} \sum_{l=1}^k (f(x_{(l)}) + \epsilon_l)\tag{46}$$

In this equation,  $f(x_{(l)})$  is fixed, and epsilons are iid random variables.

$$\begin{aligned}\text{EPE}_k(x_0) &= \text{E} \left( \left( f(x_0) + \epsilon_0 - \frac{1}{k} \sum_{l=1}^k (f(x_{(l)}) + \epsilon_l) \right)^2 \right) \\ &= \text{E} \left( \left[ \left( f(x_0) - \frac{1}{k} \sum_{l=1}^k f(x_{(l)}) \right) + \left( \epsilon_0 - \frac{1}{k} \sum_{l=1}^k \epsilon_l \right) \right]^2 \right) \\ &= \text{E} \left( [\hat{F} + \hat{E}]^2 \right) = \text{E} \left( \hat{F}^2 + 2 \cdot \hat{F} \hat{E} + \hat{E}^2 \right) = \hat{F}^2 + 2\hat{F} \cdot \text{E}(\hat{E}) + \text{E}(\hat{E}^2)\end{aligned}\tag{47}$$

where  $\hat{F}$  is nonrandom, and  $\hat{E}$  is random:

$$\begin{aligned}\hat{F} &\equiv f(x_0) - \frac{1}{k} \sum_{l=1}^k f(x_{(l)}), \\ \hat{E} &\equiv \epsilon_0 - \frac{1}{k} \sum_{l=1}^k \epsilon_l\end{aligned}\tag{48}$$

Let's calculate the expectation of  $\hat{E}$ :

$$\text{E}(\hat{E}) = \text{E} \left( \epsilon_0 - \frac{1}{k} \sum_{l=1}^k \epsilon_l \right) = \text{E}(\epsilon_0) - \frac{1}{k} \sum_{l=1}^k \text{E}(\epsilon_l) = 0 - \frac{1}{k} \sum_{l=1}^k 0 = 0\tag{49}$$

The expectation of  $\hat{E}^2$ :

$$\mathbb{E}(\hat{E}^2) = \mathbb{E} \left( \epsilon_0 - \frac{1}{k} \sum_{l=1}^k \epsilon_l \right)^2 = \mathbb{E} \left( \epsilon_0^2 + \sum_{l=1}^k \frac{\epsilon_l^2}{k^2} + \text{CrossProducts} \right) \quad (50)$$

The expectation of the cross products are zero, since epsilons are independent, so  $\mathbb{E}(\epsilon_i \epsilon_j) = \mathbb{E}\epsilon_i \cdot \mathbb{E}\epsilon_j = 0 \cdot 0 = 0$

$$\begin{aligned} \mathbb{E}(\hat{E}^2) &= \mathbb{E} \left( \epsilon_0^2 + \sum_{l=1}^k \frac{\epsilon_l^2}{k^2} \right) = \mathbb{E}\epsilon_0^2 + \sum_{l=1}^k \frac{\mathbb{E}\epsilon_l^2}{k^2} \\ &= \sigma^2 + \sum_{l=1}^k \frac{\sigma^2}{k^2} = \sigma^2 + \frac{\sigma^2}{k} \end{aligned} \quad (51)$$

We used the fact that the error has zero mean, so the variance is  $\sigma^2 = \text{Var}(\epsilon) = \mathbb{E}(\epsilon^2) - (\mathbb{E}\epsilon)^2 = \mathbb{E}(\epsilon^2)$ . So the final form is:

$$\begin{aligned} \mathbb{E}_k(x_0) &= \hat{F}^2 + 2\hat{F} \cdot \mathbb{E}(\hat{E}) + \mathbb{E}(\hat{E}^2) = \hat{F}^2 + \mathbb{E}(\hat{E}^2) \\ &= \left( f(x_0) - \frac{1}{k} \sum_{l=1}^k f(x_{(l)}) \right)^2 + \sigma^2 + \frac{\sigma^2}{k} \end{aligned} \quad (52)$$

## 2.5 Solutions for the Exercises of chapter 2

### 2.5.1 Ex. 2.2

We have  $X \in \mathbb{R}^p$  continuous and  $G$  discrete random variables. Assume we have  $K$  classes. Each class has its own distribution, let's say that class  $g$  has a pdf  $f_g(x)$  ( $x \in \mathbb{R}^p$ ). When generating points, we first choose a class with associated probabilities  $p_1, p_2, \dots, p_K$  ( $\sum p_i = 1$ ). When we have chosen the class, we generate a point with the appropriate distribution.

The Bayes classifier classifies each point  $x$  to the most probable class. So let's calculate the probability of class  $g$ , given the point. It should be noted that when I write  $P(x)$ , I mean "the probability that the chosen point is in the infinitesimal neighborhood of  $x$ ". So I should write  $P(X \in b_{dx}(x))$ , i.e., the probability that  $X$  is in the  $dx$ -volume ball around  $x$ . If the pdf was  $f(x)$ , this probability is  $f(x)dx$ . But instead, I'll write  $P(x) = f(x)$ . Likewise, when I write  $P(g)$ , I mean  $P(G = g)$ .

$$P(g|x) = \frac{P(g \cap x)}{P(x)} = \frac{P(g \cap x)}{P(x)} = \frac{P(x|g)P(g)}{\sum_{g'} P(x|g')P(g')} \quad (53)$$

The denominator is a normalizing constant, so the chosen class, for which  $P(g|x)$  is maximum:

$$\hat{g}(x) = \max_g P(x|g)P(g) \quad (54)$$

### 2.5.2 Ex. 2.3

Given a unit ball in  $p$ -dimension. We sample  $N$  data points from it uniformly. Let  $X$  be the distance from the origin. The pdf must be proportional to  $x^{p-1}$ , and integrating it from 0 to 1 gives 1, thus the pdf:

$$f(x) = p \cdot x^{p-1} \quad (55)$$

The probability that a random sample is at least  $x$  distant from the origin is:

$$P(X > x) = \int_x^1 f(x)dx = 1 - x^p \quad (56)$$

The probability that all  $N$  sample points are further from origin as  $x$ :

$$P(X_1 > x \cap X_2 > x \cap \dots \cap X_N > x) = (1 - x^p)^N \quad (57)$$

We seek and  $x$  for that this probability is a half (that will give us the median):

$$\begin{aligned} (1 - x^p)^N &= \frac{1}{2} \\ 1 - x^p &= \left(\frac{1}{2}\right)^{1/N} \\ \left[1 - \left(\frac{1}{2}\right)^{1/N}\right]^{1/p} &= x \end{aligned} \quad (58)$$

### 2.5.3 Ex. 2.4

If we choose  $a$  as the first unit base vector ( $a = [1, 0, 0, \dots, 0]^T$ ), then  $a^T \cdot x_i$  is the first coordinate of  $x_i$ . It is by definition (standard) normally distributed. Since the distribution is spherically symmetric, we can choose any direction  $a$ ,  $a^T \cdot x_i$  remains standard normal.

I created an experiment on this. Created 1000 sample points in  $p$  dimension, and rotated them into the first 2 dimension, so that we can visualize the distances. On the first image below we can see that the points get further and further away from the origin as the dimension increases.

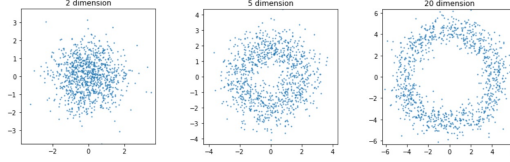


Figure 5: The sample points get further from the origin as we increase the dimension.

But this doesn't mean the points are close to each other. In the following experiment I took the random sample points, chose one of them and set it as the new origin. We can see that still the points are far from a random sample point as we increase the dimension.

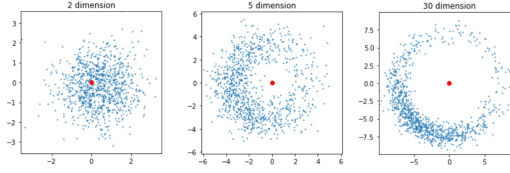


Figure 6: The points get further and further from each other (red dot is a randomly selected sample point) as we increase the dimension.

#### 2.5.4 Ex. 2.5

**equation (2.27) on page 26** I won't use indices at the expectation sign, it always confuses me. So this is the expected prediction error:

$$\text{EPE}(x_0) = \text{E}(y_0 - \hat{y}_0)^2 \quad (59)$$

Recall, that  $y_0 = x_0^T \beta + \epsilon$  is a random variable, since  $\epsilon \sim N(0, \sigma^2)$ . This is the label (the ground truth) for  $x_0$ . The prediction that we make for  $x_0$  is  $\hat{y}_0 = x_0^T \beta + \bar{\epsilon}^T \cdot \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} x_0$ .  $x_0$  is a  $p$ -vector,  $\beta$  is a  $p$ -vector,  $\bar{\epsilon}$  is an  $n$ -vector, and  $\mathbf{X}$  is a  $n$  by  $p$  matrix (each row is a training sample vector).  $\hat{y}_0 = x_0^T \beta + \bar{\epsilon}^T \cdot \mathbf{Z}^T x_0$ . Here I introduced the  $p$  by  $n$  matrix  $\mathbf{Z}$ :

$$\mathbf{Z} \equiv (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \quad (60)$$

In the expression of  $\hat{y}_0$ ,  $\vec{\epsilon}$  and  $\mathbf{Z}$  are the only random variables. The elements of  $\vec{\epsilon}$  are iid RVs.  $\vec{\epsilon}$  and  $\mathbf{Z}$  are independent. Let's calculate the expectation values of  $y_0$  and  $\hat{y}_0$ :

$$\mathbb{E}(y_0) = \mathbb{E}(x_0^T \beta + \epsilon) = x_0^T \beta \quad (61)$$

$$\begin{aligned} \mathbb{E}(\hat{y}_0) &= \mathbb{E}(x_0^T \beta + \vec{\epsilon}^T \cdot \mathbf{Z}^T x_0) \\ &= x_0^T \beta + \mathbb{E}(\vec{\epsilon}^T \cdot \mathbf{Z}^T x_0) \\ &= x_0^T \beta + \mathbb{E}(\vec{\epsilon}^T) \mathbb{E}(\mathbf{Z}^T x_0) \\ &= x_0^T \beta + \vec{0}^T \cdot \mathbb{E}(\mathbf{Z}^T x_0) \\ &= x_0^T \beta \end{aligned} \quad (62)$$

Here I used that  $\vec{\epsilon}$  and  $\mathbf{Z}$  are independent, so the expectation value of their product is the product of their expectations. For simplicity, denote  $\mu \equiv \mathbb{E}(y_0) = \mathbb{E}(\hat{y}_0) = x_0^T \beta$ . The expected prediction error:

$$\begin{aligned} \text{EPE}(x_0) &= \mathbb{E}(y_0 - \mu + \mu - \hat{y}_0)^2 \\ &= \mathbb{E}(y_0 - \mu)^2 - 2 \cdot \mathbb{E}((y_0 - \mu)(\hat{y}_0 - \mu)) + \mathbb{E}(\hat{y}_0 - \mu)^2 \\ &= \mathbb{E}(y_0 - \mu)^2 + 0 + \mathbb{E}(\hat{y}_0 - \mu)^2 \\ &= \text{Var}(y_0) + \text{Var}(\hat{y}_0) \end{aligned} \quad (63)$$

Note that  $y_0$  and  $\hat{y}_0$  are independent. The epsilon in  $y_0$  is a scalar and is nothing to do with the vector epsilon in  $\hat{y}_0$ . This is why  $\mathbb{E}((y_0 - \mu)(\hat{y}_0 - \mu))$  is zero. Now let's derive the variances:

$$\text{Var}(y_0) = \text{Var}(\mu + \epsilon) = \text{Var}(\epsilon) = \sigma^2 \quad (64)$$

Furthermore, we can write  $\text{Cov}(X, X) = \mathbb{E}((X - \mathbb{E}X) \cdot (X - \mathbb{E}X)^T) = \mathbb{E}(X \cdot X^T) - (\mathbb{E}X)(\mathbb{E}X^T)$ . Let's apply (146) and (147) to derive the variance of  $\hat{y}_0$ :

$$\begin{aligned} \text{Var}(\hat{y}_0) &= \text{Var}(\mu + x_0^T \cdot \mathbf{Z} \cdot \vec{\epsilon}) = \text{Var}(x_0^T \cdot \mathbf{Z} \cdot \vec{\epsilon}) \\ &= x_0^T \cdot \text{Cov}(\mathbf{Z} \cdot \vec{\epsilon}, \mathbf{Z} \cdot \vec{\epsilon}) \cdot x_0 \\ &= x_0^T \cdot \left( \mathbb{E}(\mathbf{Z} \vec{\epsilon} \cdot \vec{\epsilon}^T \mathbf{Z}^T) - \mathbb{E}(\mathbf{Z} \vec{\epsilon}) \cdot \mathbb{E}(\vec{\epsilon}^T \mathbf{Z}^T) \right) \cdot x_0 \\ &= x_0^T \cdot \mathbb{E}(\mathbf{Z} \vec{\epsilon} \cdot \vec{\epsilon}^T \mathbf{Z}^T) \cdot x_0 \end{aligned} \quad (65)$$

Here we used the fact that  $\mathbf{Z}$  and  $\vec{\epsilon}$  are independent, so  $\mathbb{E}(\mathbf{Z} \vec{\epsilon}) = \mathbb{E}\mathbf{Z} \cdot \mathbb{E}\vec{\epsilon} = \mathbb{E}\mathbf{Z} \cdot \vec{0} = \vec{0}$ , and  $\vec{0} \cdot \vec{0}^T = \mathbf{0}$ , zero matrix.



$$\begin{aligned}
(\mathbf{Z}\vec{\epsilon} \cdot \vec{\epsilon}^T \mathbf{Z}^T)_{i,j} &= \sum_{k,l} (\mathbf{Z})_{i,k} \cdot (\vec{\epsilon}\vec{\epsilon}^T)_{k,l} \cdot (\mathbf{Z}^T)_{l,j} \\
&\rightarrow \mathbb{E}(\mathbf{Z}\vec{\epsilon} \cdot \vec{\epsilon}^T \mathbf{Z}^T)_{i,j} = \sum_{k,l} \mathbb{E}(Z_{i,k} \cdot \epsilon_k \cdot \epsilon_l \cdot Z_{j,l}) \\
&= \sum_{k,l} \mathbb{E}(Z_{i,k} \cdot Z_{j,l}) \cdot \mathbb{E}(\epsilon_k \cdot \epsilon_l) = \sum_{k,l} \mathbb{E}(Z_{i,k} \cdot Z_{j,l}) \cdot \sigma^2 \delta_{k,l} \\
&= \sigma^2 \cdot \sum_k \mathbb{E}(Z_{i,k} \cdot Z_{j,k}) = \sigma^2 \cdot \mathbb{E}(\mathbf{Z}\mathbf{Z}^T)_{i,j} \\
&\rightarrow \mathbb{E}(\mathbf{Z}\vec{\epsilon} \cdot \vec{\epsilon}^T \mathbf{Z}^T) = \sigma^2 \cdot \mathbb{E}(\mathbf{Z}\mathbf{Z}^T)
\end{aligned} \tag{66}$$

Substituting this into (65):

$$\text{Var}(\hat{y}_0) = x_0^T \cdot \sigma^2 \mathbb{E}(\mathbf{Z}\mathbf{Z}^T) \cdot x_0 \tag{67}$$

$$\mathbb{E}(\mathbf{Z}\mathbf{Z}^T) = \mathbb{E}\left((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}\right) = \mathbb{E}\left((\mathbf{X}^T \mathbf{X})^{-1}\right) \tag{68}$$

Putting it all together:

$$\text{EPE}(x_0) = \text{Var}(y_0) + \text{Var}(\hat{y}_0) = \sigma^2 + \sigma^2 \cdot x_0^T \cdot \mathbb{E}\left((\mathbf{X}^T \mathbf{X})^{-1}\right) \cdot x_0 \tag{69}$$

And this is what we wanted to derive.

**equation (2.28) on page 26**

$$\mathbb{E}(x_0^T \text{Cov}(X)^{-1} x_0) = \mathbb{E} \sum_{i,j} x_{0,i} \text{Cov}(X)_{i,j}^{-1} x_{0,j} \tag{70}$$

Assuming that  $x_0 \sim X$ , i.e.,  $x_0$  (the test point) has the same distribution as  $X$  (the training data), and the expectation of it is the zero vector,  $\text{Cov}(x_0) = \text{Cov}(X) = \mathbb{E}(x_0 x_0^T)$ .

$$\begin{aligned}
\mathbb{E} \sum_{i,j} x_{0,i} \text{Cov}(X)_{i,j}^{-1} x_{0,j} &= \mathbb{E} \sum_{i,j} \text{Cov}(X)_{i,j}^{-1} \text{Cov}(x_0)_{i,j} \\
&= \mathbb{E} \sum_{i,j} \text{Cov}(X)_{i,j}^{-1} \text{Cov}(X)_{i,j} = \mathbb{E} \sum_i \left( \sum_j \text{Cov}(X)_{i,j}^{-1} \text{Cov}(X)_{j,i} \right) \\
&= \mathbb{E} \sum_i [\text{Cov}(X)^{-1} \text{Cov}(X)]_{i,i} = \mathbb{E}(\text{Trace}(\text{Cov}(X)^{-1} \text{Cov}(X))) \\
&= \mathbb{E}(\text{Trace} I_{p \times p}) = p
\end{aligned} \tag{71}$$

### 2.5.5 Ex. 2.6

Assume that we have  $n$  identical inputs  $x_1 = x_2 = \dots = x_n \equiv x$  with outputs  $y_1, y_2, \dots, y_n$ . The least squares formula:

$$RSS(\theta) = \sum_{i=1}^n (y_i - f_{\theta}(x))^2 \quad (72)$$

The weighted least squares formula:

$$RSS_w(\theta) = n \cdot \left( \frac{\sum_{i=1}^n y_i}{n} - f_{\theta}(x) \right)^2 \quad (73)$$

I claim that the two expressions differ by a constant term that doesn't depend on  $\theta$ , so both expressions lead to the same solution. This naturally extends to the case when we have groups of equal inputs.

Expanding  $RSS$ :

$$RSS(\theta) = \sum_{i=1}^n y_i^2 - 2f_{\theta}(x) \sum_{i=1}^n y_i + f_{\theta}^2(x) \quad (74)$$

Expanding  $RSS_w$ :

$$RSS_w(\theta) = \frac{(\sum_{i=1}^n y_i)^2}{n} - 2f_{\theta}(x) \sum_{i=1}^n y_i + f_{\theta}^2(x) \quad (75)$$

So the difference of the 2 expressions is a constant that doesn't depend on  $\theta$ . So when we derive wrt  $\theta$ , we get the same formulae.

Whenever we have observations with identical values  $x$ , we can always refactor the  $RSS$  for the groups according to (72)  $\rightarrow$  (73).

### 2.5.6 Ex. 2.7

Our estimator according to the problem statement:

$$\hat{f}(x_0) = \sum_{i=1}^N l_i(x_0; \mathcal{X}) y_i \quad (76)$$

a) For kNN, the weights are:

$$l_i(x_0; \mathcal{X}) = \frac{1}{k} \delta(x_i \in \text{kNN}(x_0)) \quad (77)$$

where  $\delta(x_i \in \text{kNN}(x_0))$  is 1 if  $x_i$  is in the set of k-nearest neighbors of  $x_0$ , and 0 otherwise. So in this case we average the  $y$ s of the k-nearest neighbors of  $x_0$ .

For linear regression we have

$$\hat{f}(x_0) = x_0^T \beta \quad (78)$$

Where  $\beta$  comes from the following equation (see Section 2.2):

$$\beta = \text{E}(XX^T)^{-1} \text{E}(XY) \quad (79)$$

Now let's calculate this expression. We estimate the expectation values with averages.

$$\begin{aligned} \text{E}(XX^T)^{-1} &\approx \left( \frac{1}{N} \sum_{j=1}^N x_j x_j^T \right)^{-1} \\ \text{E}(XY) &\approx \frac{1}{N} \sum_{i=1}^N x_i y_i \end{aligned} \quad (80)$$

With these, we can formulate  $\hat{f}(x_0)$  as follows:

$$\begin{aligned} \hat{f}(x_0) &= x_0^T \left( \frac{1}{N} \sum_{j=1}^N x_j x_j^T \right)^{-1} \cdot \frac{1}{N} \sum_{i=1}^N x_i y_i \\ &= \sum_{i=1}^N x_0^T \left( \sum_{j=1}^N x_j x_j^T \right)^{-1} x_i y_i \\ &\equiv \sum_{i=1}^N l_i(x_0; \mathcal{X}) y_i \end{aligned} \quad (81)$$

From this we get the weights:

$$l_i(x_0; \mathcal{X}) = x_0^T \left( \sum_{j=1}^N x_j x_j^T \right)^{-1} x_i \quad (82)$$

### 2.5.7 Ex. 2.9

The short answer is this:

$$ER_{tr}(\hat{\beta}) \leq ER_{tr}(E\hat{\beta}) = ER_{te}(E\hat{\beta}) \leq ER_{te}(\hat{\beta}) \quad (83)$$

Now I explain this in more details.

1. Proving the left inequality.  $\hat{\beta}$  comes from the following:

$$\hat{\beta} = \arg \min_{\beta'} R_{tr}(\beta') \quad (84)$$

This implies that for any fix  $\beta$ :

$$R_{tr}(\hat{\beta}) \leq R_{tr}(\beta) \quad (85)$$

Taking the expectation of both sides:

$$ER_{tr}(\hat{\beta}) \leq ER_{tr}(\beta) \quad (86)$$

$\hat{\beta}$  is a random variable (which depends on the training data), we can take the expectation, so we get  $E\hat{\beta}$  which is a fix, non-random vector. Substituting into the above inequality we get what we wanted to prove:

$$ER_{tr}(\hat{\beta}) \leq ER_{tr}(E\hat{\beta}) \quad (87)$$

2. Proving the equation in the middle. For any fix  $\beta$ :

$$ER_{tr}(\beta) = \frac{1}{N} \sum_{i=1}^N E(y_i - \beta^T x_i)^2 = E(Y - \beta^T X)^2 \quad (88)$$

$$ER_{te}(\beta) = \frac{1}{M} \sum_{i=1}^M E(\tilde{y}_i - \beta^T \tilde{x}_i)^2 = E(Y - \beta^T X)^2 \quad (89)$$

This is because both the train and the test data come from the same distribution. So for any fix  $\beta$ ,  $ER_{tr}(\beta) = ER_{te}(\beta)$ . Since  $E\hat{\beta}$  is a fix vector, we're done with this part.

3. Proving the right inequality. For this we use the fact that the training data and the test data are independent. Thus  $\hat{\beta}$  and the test data are also independent. For this part, just forget about the training data. Think of  $\hat{\beta}$  as a random vector independent from the (test) data.

$$\mathbb{E}R_{te}(\hat{\beta}) = \mathbb{E}(Y - \hat{\beta}^T X)^2 = \mathbb{E}\mathbb{E}\left((Y - \hat{\beta}^T X)^2 | X, Y\right) \quad (90)$$

$$\begin{aligned} \mathbb{E}\left((Y - \hat{\beta}^T X)^2 | X, Y\right) &= \mathbb{E}\left(Y^2 - 2Y\hat{\beta}^T X + (\hat{\beta}^T X)^2 | X, Y\right) \\ &= Y^2 - 2Y\mathbb{E}(\hat{\beta}^T)X + X^T\mathbb{E}(\hat{\beta}\hat{\beta}^T)X \\ &= Y^2 - 2Y\mathbb{E}(\hat{\beta}^T)X + X^T[\mathbb{E}\hat{\beta} \cdot \mathbb{E}\hat{\beta}^T + \text{Cov}(\hat{\beta})]X \\ &= Y^2 - 2Y\mathbb{E}(\hat{\beta}^T)X + (\mathbb{E}\hat{\beta}^T)XX^T(\mathbb{E}\hat{\beta}) + X^T\text{Cov}(\hat{\beta})X \end{aligned} \quad (91)$$

Since the covariance matrix is positive semi-definite,  $X^T\text{Cov}(\beta)X \geq 0$

$$\begin{aligned} \mathbb{E}\left((Y - \hat{\beta}^T X)^2 | X, Y\right) &\geq Y^2 - 2Y\mathbb{E}(\hat{\beta}^T)X + (\mathbb{E}\hat{\beta}^T)XX^T(\mathbb{E}\hat{\beta}) \\ \mathbb{E}\left((Y - \hat{\beta}^T X)^2 | X, Y\right) &\geq (Y - \mathbb{E}(\hat{\beta}^T)X)^2 \\ \mathbb{E}\mathbb{E}\left((Y - \hat{\beta}^T X)^2 | X, Y\right) &\geq \mathbb{E}(Y - \mathbb{E}(\hat{\beta}^T)X)^2 \\ \mathbb{E}(Y - \hat{\beta}^T X)^2 &\geq \mathbb{E}(Y - \mathbb{E}(\hat{\beta}^T)X)^2 \\ \mathbb{E}R_{te}(\hat{\beta}) &\geq \mathbb{E}R_{te}(\mathbb{E}\hat{\beta}) \end{aligned} \quad (92)$$

### 3 Linear Methods for Regression

#### 3.1 equations on page 47 and 48

**Variance of beta hat (page 47).** We know the formulae for  $\hat{\beta}$  (equation 3.6 on page 45):

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (93)$$

Using (146):

$$\begin{aligned} \text{Var}(\hat{\beta}) &\equiv \text{Cov}(\hat{\beta}) \\ &= \text{Cov}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}) \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \cdot \text{Cov}(\mathbf{y}) \cdot \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \cdot \sigma^2 \mathbf{I} \cdot \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \\ &= \sigma^2 \cdot (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \\ &= \sigma^2 \cdot (\mathbf{X}^T \mathbf{X})^{-1} \end{aligned} \quad (94)$$

**Sigma hat.** Deriving the expectation of  $\hat{\sigma}^2$ . We know that  $\mathbf{H}$ , that hat matrix is an orthogonal projection onto the column space of  $\mathbf{X}$ . This implies that  $\mathbf{H}^2 = \mathbf{H} = \mathbf{H}^T$ , and  $\text{Tr}(\mathbf{H}) = p + 1$  (the trace of an orthogonal projection is the dimension of the subspace it projects onto, that is, the rank of  $\mathbf{X}$ ). Another thing is that  $(\mathbf{I} - \mathbf{H})$  is also an orthogonal projection. It projects to the orthogonal complement of the column space of  $\mathbf{X}$ . So  $(\mathbf{I} - \mathbf{H})^2 = (\mathbf{I} - \mathbf{H}) = (\mathbf{I} - \mathbf{H})^T$ , and  $\text{Tr}(\mathbf{I} - \mathbf{H}) = N - p - 1$ .

$$\begin{aligned} \sum_{i=1}^N (y_i - \hat{y}_i)^2 &= (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}}) \\ &= (\mathbf{y} - \mathbf{H}\mathbf{y})^T (\mathbf{y} - \mathbf{H}\mathbf{y}) \\ &= ((\mathbf{I} - \mathbf{H})\mathbf{y})^T ((\mathbf{I} - \mathbf{H})\mathbf{y}) \\ &= \mathbf{y}^T (\mathbf{I} - \mathbf{H})^T (\mathbf{I} - \mathbf{H}) \mathbf{y} \\ &= \mathbf{y}^T (\mathbf{I} - \mathbf{H})^2 \mathbf{y} = \mathbf{y}^T (\mathbf{I} - \mathbf{H}) \mathbf{y} \\ &= \text{Tr}(\mathbf{y}^T (\mathbf{I} - \mathbf{H}) \mathbf{y}) = \text{Tr}((\mathbf{I} - \mathbf{H}) \mathbf{y} \mathbf{y}^T) \end{aligned} \quad (95)$$

We know that  $\text{Cov}(\mathbf{y}) = \sigma^2 \mathbf{I} = \text{E}(\mathbf{y} \mathbf{y}^T) - (\text{E} \mathbf{y}) \cdot (\text{E} \mathbf{y}^T)$ . Taking the expectation:

$$\begin{aligned}
\mathbb{E} \sum_{i=1}^N (y_i - \hat{y}_i)^2 &= \mathbb{E} (\text{Tr}((\mathbf{I} - \mathbf{H})\mathbf{y}\mathbf{y}^T)) \\
&= \text{Tr}((\mathbf{I} - \mathbf{H})\mathbb{E}(\mathbf{y}\mathbf{y}^T)) \\
&= \text{Tr}((\mathbf{I} - \mathbf{H}) \cdot (\sigma^2 \mathbf{I} + \mathbb{E}\mathbf{y} \cdot \mathbb{E}\mathbf{y}^T)) \\
&= \text{Tr}((\mathbf{I} - \mathbf{H})\sigma^2 + (\mathbf{I} - \mathbf{H}) \cdot \mathbb{E}\mathbf{y} \cdot \mathbb{E}\mathbf{y}^T) \quad (96) \\
&= \text{Tr}((\mathbf{I} - \mathbf{H})\sigma^2) + \text{Tr}((\mathbf{I} - \mathbf{H}) \cdot \mathbb{E}\mathbf{y} \cdot \mathbb{E}\mathbf{y}^T) \\
&= \text{Tr}(\mathbf{I} - \mathbf{H}) \cdot \sigma^2 + \text{Tr}(\mathbb{E}\mathbf{y}^T \cdot (\mathbf{I} - \mathbf{H}) \cdot \mathbb{E}\mathbf{y}) \\
&= (N - p - 1) \cdot \sigma^2 + \mathbb{E}\mathbf{y}^T \cdot (\mathbf{I} - \mathbf{H}) \cdot \mathbb{E}\mathbf{y} \\
&= (N - p - 1) \cdot \sigma^2
\end{aligned}$$

At the last step we had to assume that  $\mathbb{E}\mathbf{y}$  lies in the column space of  $\mathbf{X}$ , because it means that  $\mathbb{E}\mathbf{y}$  and  $(\mathbf{I} - \mathbf{H})\mathbb{E}\mathbf{y}$  are perpendicular to each other. This means that the response ( $y$ ) is linear in its inputs, plus a random variable with zero mean. Now we see that

$$\frac{1}{N - p - 1} \mathbb{E} \sum_{i=1}^N (y_i - \hat{y}_i)^2 = \sigma^2 \rightarrow \mathbb{E}\hat{\sigma}^2 = \sigma^2 \quad (97)$$

**Distribution of sigma hat.** (3.11) states that  $\hat{\sigma}^2$  is proportional to a Chi-square distribution with  $N - p - 1$  parameters. Now we use the assumption that  $\mathbf{y} = \mathbf{X}\beta + \epsilon$ , where  $\mathbf{X}$  and  $\beta$  are fixed, and  $\epsilon$  is a vector of iid normal random variables with zero mean and  $\sigma^2$  variance. According to (95) we can write that:

$$\sum_{i=1}^N (y_i - \hat{y}_i)^2 = \|(\mathbf{I} - \mathbf{H})\mathbf{y}\|^2 \quad (98)$$

Since  $\mathbf{H}$  is a projection to the column space of  $\mathbf{X}$ ,  $\mathbf{H}\mathbf{X} = \mathbf{X}$ , and  $(\mathbf{I} - \mathbf{H})\mathbf{X} = \mathbf{X} - \mathbf{H}\mathbf{X} = \mathbf{X} - \mathbf{X} = \mathbf{0}$ . So  $(\mathbf{I} - \mathbf{H})\mathbf{y} = (\mathbf{I} - \mathbf{H}) \cdot (\mathbf{X}\beta + \epsilon) = (\mathbf{I} - \mathbf{H}) \cdot \epsilon$ .

$$\sum_{i=1}^N (y_i - \hat{y}_i)^2 = \|(\mathbf{I} - \mathbf{H})\epsilon\|^2 \quad (99)$$

Now it is clear that this is  $\sigma^2 \cdot \chi_{N-p-1}^2$ , because we project the spherical normal distribution ( $\epsilon$ ) to a  $(N-p-1)$ -dimensional plane (subspace).

**The Z-score.** According to (3.12) we form the standardized coefficient or Z-score

$$z_j = \frac{\hat{\beta}_j}{\hat{\sigma}\sqrt{v_j}} \quad (100)$$

Why is this a t-distribution under the null hypothesis that  $\beta_j = 0$ ?

$$\hat{\beta} = \beta + \sigma \cdot (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon = \beta + \sigma \cdot (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{H} \epsilon \quad (101)$$

$$\hat{\sigma}^2 = \frac{1}{N - p - 1} \|(\mathbf{I} - \mathbf{H})\epsilon\|^2 \quad (102)$$

Now because  $\mathbf{H}\epsilon$  and  $(\mathbf{I} - \mathbf{H})\epsilon$  are independent,  $\hat{\beta}$  and  $\hat{\sigma}^2$  are also independent. Moreover,  $\hat{\beta}$  has a covariance  $\sigma^2(\mathbf{X}^T \mathbf{X})^{-1}$ ,  $\hat{\beta}_j$  has a variance  $\sigma^2 \cdot v_j$ , with  $v_j = [(\mathbf{X}^T \mathbf{X})^{-1}]_{jj}$ . So under the null hypothesis,

$$\frac{\beta_j}{\sigma\sqrt{v_j}} \sim N(0, 1) \quad (103)$$

Also,

$$\frac{\hat{\sigma}}{\sigma} \sim \sqrt{\frac{\chi_{N-p-1}^2}{N - p - 1}} \quad (104)$$

According to (149), the following has a Student's t-distribution with  $N - p - 1$  degrees of freedom

$$\frac{\frac{\beta_j}{\sigma\sqrt{v_j}}}{\frac{\hat{\sigma}}{\sigma}} = \frac{\beta_j}{\hat{\sigma}\sqrt{v_j}} = z_j \quad (105)$$

**F statistic.** According to (3.13) on page 48, we form the following statistic to decide whether we can drop groups of coefficients simultaneously.

$$F = \frac{(\text{RSS}_0 - \text{RSS}_1)/(p_1 - p_0)}{\text{RSS}_1/(N - p_1 - 1)} \quad (106)$$

I will show that under the null-hypothesis,  $\text{RSS}_0 - \text{RSS}_1$  is chi-squared with  $p_1 - p_0$  degrees of freedom,  $\text{RSS}_1$  is chi-squared with  $N - p_1 - 1$  degrees of freedom, and are independent. So according to appendix C.2,  $F$  has indeed an F-distribution with  $(p_1 - p_0), (N - p_1 - 1)$  parameters.



Let  $X$  be the  $N \times (p_0 + 1)$  data-matrix, while  $X_1 = [X|X']$  the extended  $N \times (p_1 + 1)$  data-matrix. We assume that the smaller model is true, i.e.,  $y = X\beta + \epsilon$ . We have two different estimates for  $y$ .  $\hat{y}_0 = H_0 y$  comes from the smaller model, while  $\hat{y}_1 = H_1 y$  comes from the bigger model.  $H_0$  and  $H_1$  are projections.  $H_1$  projects onto the column space of  $X_1$ , which we denote by  $W_1$ .  $H_0$  projects onto the column space of  $X$ , which we denote by  $W_0$ , this is actually a subspace of  $W_1$ . Let  $W_2$  be a subspace in  $W_1$  that is orthogonal to  $W_0$ .  $H_1 - H_0$  projects onto this subspace. Now calculate the residual sum of squares.

$$\begin{aligned}
\text{RSS}_0 &= \|y - \hat{y}_0\|^2 \\
&= \|y - H_0 y\|^2 \\
&= \|(I - H_0)y\|^2 \\
&= \|(I - H_0)(X\beta + \epsilon)\|^2 \\
&= \|(I - H_0)X\beta + (I - H_0)\epsilon\|^2 \\
&= \|0\beta + (I - H_0)\epsilon\|^2 \\
&= \|(I - H_0)\epsilon\|^2 \\
&= \epsilon^T (I - H_0)^T (I - H_0) \epsilon \\
&= \epsilon^T (I - H_0) (I - H_0) \epsilon \\
&= \epsilon^T (I - H_0) \epsilon
\end{aligned} \tag{107}$$

Similarly,

$$\text{RSS}_1 = \epsilon^T (I - H_1) \epsilon \tag{108}$$

Note that here we used the fact that the columns of  $X$  make up the subspace  $W_0$ .  $I - H_0$  projects onto  $W_0^\perp$ , so  $(I - H_0)X = 0$ . Similarly,  $(I - H_1)X = 0$ .

From these, we can calculate the difference of the residual sum of squares.

$$\begin{aligned}
\text{RSS}_0 - \text{RSS}_1 &= \epsilon^T (H_1 - H_0) \epsilon \\
&= \|(H_1 - H_0)\epsilon\|^2
\end{aligned} \tag{109}$$

So  $\text{RSS}_0 - \text{RSS}_1$  is a chi-squared random variable with  $p_1 - p_0$  degrees of freedom (the dimension of  $W_2$ ). And  $\text{RSS}_1$  is also chi-squared with  $N - p_1 - 1$  degrees of freedom (the dimension of  $W_1^\perp$ ).  $\text{RSS}_0 - \text{RSS}_1$  and  $\text{RSS}_1$  are independent, because  $I - H_1$  and  $H_1 - H_0$  project onto perpendicular subspaces.

### 3.2 Equation (3.28) on page 54

I'd like to confirm that in general,

$$\hat{\beta}_j = \frac{z_j^T y}{z_j^T z_j} \quad (110)$$

But first some notations and clarifications.  $y, z_j \in \mathbb{R}^N$ .  $x_j \in \mathbb{R}^N$  is the  $j$ th column vector of  $X$ , the data matrix.  $W_X$  is the column space of  $X$ ,  $W_{X(j)}$  is the subspace spanned by all the columns of  $X$  except the  $j$ th column.  $W_{X(j)}^\perp$  is a one-dimensional subspace that is orthogonal to  $W_{X(j)}$ , and

$$W_{X(j)} + W_{X(j)}^\perp = W_X \quad (111)$$

$z_j = x_j - P_j x_j$ , where  $P_j$  projects onto  $W_{X(j)}$ . We can also express it as

$$z_j = P_j^\perp x_j \quad (112)$$

where  $P_j^\perp$  projects onto  $W_{X(j)}^\perp$ . We know, that

$$\hat{\beta} = (X^T X)^{-1} X^T y \quad (113)$$

Denoting the  $j$ th column vector of  $(X^T X)^{-1}$  by  $b_j$  ( $(X^T X)^{-1}$  is symmetric, so  $b_j^T$  is the  $j$ th row vector), we can write:

$$\hat{\beta}_j = b_j^T X^T y \quad (114)$$

From appendix (D.3) we can construct  $P_j^\perp$ :

$$P_j^\perp = \frac{X b_j \cdot b_j^T X^T}{v_j} \quad (115)$$

Where  $v_j$  is the  $j$ th element of  $b_j$  ( $= [(X^T X)^{-1}]_{j,j}$ ). With this we can calculate  $z_j$  ( $A_{:,j}$  denotes the  $j$ th column vector of matrix  $A$ ):

$$\begin{aligned} z_j = P_j^\perp x_j &= (P_j^\perp X)_{:,j} = \left( \frac{X b_j \cdot b_j^T X^T}{v_j} X \right)_{:,j} \\ &= \left( \frac{X b_j b_j^T X^T X}{v_j} \right)_{:,j} = \left( \frac{X b_j \delta_j^T}{v_j} \right)_{:,j} = \frac{X b_j}{v_j} \end{aligned} \quad (116)$$

Here  $\delta_j = I_{:,j}$ , the  $j$ th column vector of the identity. We can plug this result into (110):

$$\hat{\beta}_j = \frac{\frac{b_j^T X^T}{v_j} y}{\frac{b_j^T X^T X b_j}{v_j^2}} = \frac{b_j^T X^T y}{\frac{b_j^T \delta_j}{v_j}} = \frac{b_j^T X^T y}{\frac{v_j}{v_j}} = b_j^T X^T y \quad (117)$$

It is indeed the same as we got in (114), so the proof is complete.

### 3.3 Solutions for the Exercises of chapter 3

#### 3.3.1 Ex. 3.1

According to Appendix (C.2), the F-statistics can be written in the form:

$$F \sim \frac{\chi_{d_1}^2/d_1}{\chi_{d_2}^2/d_2} \quad (118)$$

where in our case  $d_1 = p_1 - p_0$ ,  $d_2 = N - p_1 - 1$ . Dropping a single coefficient means that  $p_1 = p_0 + 1 \rightarrow p_1 - p_0 = 1$ , so

$$F \sim \frac{\chi_1^2/1}{\chi_{d_2}^2/d_2} \sim \frac{N(0,1)^2}{\chi_{d_2}^2/d_2} \sim \left( \frac{N(0,1)}{\sqrt{\frac{\chi_{N-p_1-1}^2}{(N-p_1-1)}}} \right)^2 \sim t_{N-p_1-1}^2 \quad (119)$$

The Z-score is t-distributed with  $N - p - 1$  parameters, so the F-statistics for dropping a single coefficient is indeed *distributed* as the square of the corresponding Z-score. Well, this doesn't prove that the square of the calculated  $Z$  is equal to the calculated  $F$ . So let's prove it. Without loss of generality we can assume that we test for the last coefficient,  $j = p + 1$ .

$$z_j^2 = \frac{\hat{\beta}_j^2}{\hat{\sigma}^2 v_j} \quad (120)$$

We have to show that this equals to the following  $F$ :

$$F = \frac{(\text{RSS}_0 - \text{RSS}_1)/(p_1 - p_0)}{\text{RSS}_1/(N - p_1 - 1)} \quad (121)$$

Now since  $(N - p_1 - 1)\hat{\sigma}^2 = \text{RSS}_1$ , and  $p_1 - p_0 = 1$ , we have to show that

$$\frac{\hat{\beta}_j^2}{v_j} = \text{RSS}_0 - \text{RSS}_1 \quad (122)$$

Denote the  $j$ th column (=  $j$ th row) of  $(X^T X)^{-1}$  as  $b_j$ . With this notation

$$\hat{\beta}_j = b_j^T X^T y = y^T X b_j \quad (123)$$

$$\frac{\hat{\beta}_j^2}{v_j} = y^T \frac{X b_j \cdot b_j^T X^T}{v_j} y \quad (124)$$

$$\text{RSS}_0 - \text{RSS}_1 = y^T (H_1 - H_0) y \quad (125)$$

Where  $H_1$  projects onto the column space of  $X$ ,  $H_0$  projects onto  $W_0$ .  $W_0$  is the column space of the matrix same as  $X$  but dropping the last column. According to Appendix (D.3):

$$H_1 - H_0 = \frac{X b_j \cdot b_j^T X^T}{v_j} \quad (126)$$

which concludes the proof.

### 3.3.2 Ex. 3.2

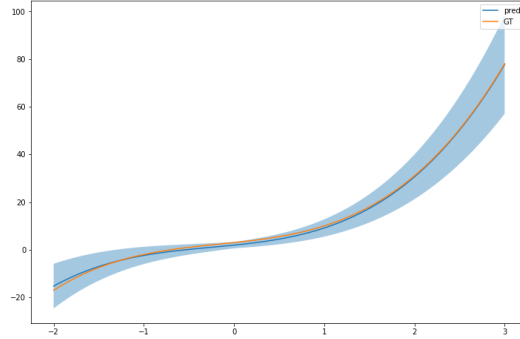


Figure 7: True function, predicted function, with 95% confidence band (point-wise).

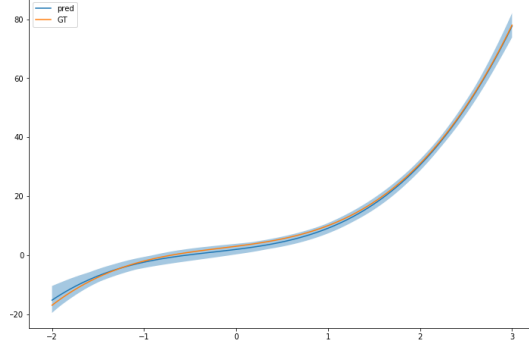


Figure 8: True function, predicted function, with 95% confidence band (from multivariate normal).

### 3.3.3 Ex. 3.3

**a.** We formulate an estimate of  $a^T \beta$  as  $c^T y$ . For the least squares estimate we have that

$$a^T \hat{\beta} = a^T (X^T X)^{-1} X^T y \quad (127)$$

From this,

$$c_0 = X(X^T X)^{-1} a \quad (128)$$

Any  $c$  can be written as

$$c = c_0 + c_1 = X(X^T X)^{-1} a + c_1 \quad (129)$$

The constraint is that  $E(c^T y) = a^T \beta$ .

$$E(c^T y) = E(c_0^T y) + E(c_1^T y) = a^T \beta \quad (130)$$

Since  $E(c_0^T y) = a^T \beta$ , we have that  $E(c_1^T y) = 0$ .

$$0 = E(c_1^T y) = c_1^T X \beta \quad (131)$$

Because  $\beta$  is unobservable, we conclude that

$$0 = c_1^T X \quad (132)$$

which means

$$c_1^T c_0 = c_1^T X(X^T X)^{-1} a = 0 \quad (133)$$

Now consider the variances.

$$\text{Var}(a^T \hat{\beta}) = a^T \text{Var}(\hat{\beta}) a = \sigma^2 a^T (X^T X)^{-1} a \quad (134)$$

Calculating the variance of a general unbiased estimate, using (133):

$$\begin{aligned} \text{Var}(c^T y) &= \sigma^2 c^T c \\ &= \sigma^2 (c_0^T + c_1^T)(c_0 + c_1) = \sigma^2 (c_0^T c_0 + 0 + 0 + c_1^T c_1) \\ &= \sigma^2 c_0^T c_0 + \sigma^2 c_1^T c_1 \\ &= \text{Var}(a^T \hat{\beta}) + \sigma^2 c_1^T c_1 \end{aligned} \quad (135)$$

Since  $\sigma^2 c_1^T c_1 \geq 0$ , we conclude that

$$\text{Var}(c^T y) \geq \text{Var}(a^T \hat{\beta}) \quad (136)$$

**b.** The solution is basically the same as for the previous one. Here we will use the fact that  $A^T A$  is a positive semidefinite matrix for any matrix  $A$ . A linear unbiased estimate for  $\beta$  can be expressed as  $\tilde{\beta} = C^T y$ , where  $C$  is a  $N \times (p+1)$  matrix. We can express  $C$  as  $C = X(X^T X)^{-1} + C_1 = C_0 + C_1$ . The estimates are unbiased, so  $E(C^T y) = \beta$ . From this we have

$$E(C^T y) = (C_0 + C_1)^T (X\beta) = \beta + C_1^T X\beta = \beta \rightarrow C_1^T X\beta = 0 \quad (137)$$

Because  $\beta$  is unobservable, we have that

$$C_1^T X = 0 \quad (138)$$

Now consider the variances.  $\hat{V} \equiv \text{Var}(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$ .  $\tilde{V} \equiv \text{Var}(\tilde{\beta}) = \text{Var}(C^T y) = \sigma^2 C^T C = \sigma^2 (C_0 + C_1)^T (C_0 + C_1)$ . Using (138), we can write that  $\tilde{V} = \sigma^2 C_0^T C_0 + \sigma^2 C_1^T C_1 = \hat{V} + \sigma^2 C_1^T C_1$ . From this:  $\tilde{V} - \hat{V} = \sigma^2 C_1^T C_1$  which is a positive semidefinite matrix. This concludes the proof.

### 3.3.4 Ex. 3.5

The ridge regression loss in vectorized form:

$$L_{\text{ridge}} = (\mathbf{y} - \beta_0 \mathbf{e} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \beta_0 \mathbf{e} - \mathbf{X}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^T \boldsymbol{\beta}$$

$\mathbf{e}$  is a vector of all ones in  $\mathbb{R}^N$ .  $\mathbf{X}$  is  $N \times p$  matrix, so it doesn't have  $\mathbf{e}$  as its first column. That's why we have a separate  $\beta_0$ , which is a scalar. All other coefficients are in the vector  $\boldsymbol{\beta}$ .

How to formulate  $L_{\text{ridge}}^c$ ? We can use  $L_{\text{ridge}}$ , but instead of  $\mathbf{X}$ , we write  $\mathbf{X} - \frac{1}{N} \mathbf{e} \mathbf{e}^T \mathbf{X}$ . It's because  $\frac{1}{N} \mathbf{e}^T \mathbf{X} = [\bar{x}_1, \bar{x}_2, \dots, \bar{x}_p]$ . Let's create the following notation:

$$\mathbf{A} \equiv \frac{\mathbf{e} \mathbf{e}^T}{\mathbf{e}^T \mathbf{e}} = \frac{\mathbf{e} \mathbf{e}^T}{N}$$

This is a projection onto the line defined by  $\mathbf{e}$ . So to get  $L_{\text{ridge}}^c$ , we replace  $\mathbf{X}$  in  $L_{\text{ridge}}$  with  $(\mathbf{I} - \mathbf{A})\mathbf{X}$ . So we only need to solve  $L_{\text{ridge}}$  for  $\beta_0$  and  $\boldsymbol{\beta}$ . Once we have these, we get  $\beta_0^c, \boldsymbol{\beta}^c$  by applying the change  $\mathbf{X} \rightarrow (\mathbf{I} - \mathbf{A})\mathbf{X}$ .

$$\begin{aligned} \frac{\partial L_{\text{ridge}}}{\partial \beta_0} &= -2 \mathbf{e}^T (\mathbf{y} - \beta_0 \mathbf{e} - \mathbf{X}\boldsymbol{\beta}) = 0 \\ \rightarrow \beta_0 &= \bar{y} - \frac{1}{N} \mathbf{e}^T \mathbf{X} \boldsymbol{\beta} \end{aligned}$$

Note that  $\mathbf{e}^T \mathbf{e} = N$ , and  $\mathbf{e}^T \mathbf{y} = N \bar{y}$ . Now we can get  $\beta_0^c$  by the above mentioned substitution. Note that  $\mathbf{e}^T \mathbf{A} = \frac{1}{N} \mathbf{e}^T \mathbf{e} \mathbf{e}^T = \frac{1}{N} N \mathbf{e}^T = \mathbf{e}^T$ .

$$\begin{aligned} \beta_0^c &= \bar{y} - \frac{1}{N} \mathbf{e}^T (\mathbf{I} - \mathbf{A}) \mathbf{X} \boldsymbol{\beta}^c \\ &= \bar{y} - \frac{1}{N} (\mathbf{e}^T \mathbf{I} - \mathbf{e}^T \mathbf{A}) \mathbf{X} \boldsymbol{\beta}^c \\ &= \bar{y} - \frac{1}{N} (\mathbf{e}^T - \mathbf{e}^T) \mathbf{X} \boldsymbol{\beta}^c \\ &= \bar{y} \end{aligned}$$

So we see that  $\beta_0^c$  does not depend on the features and coefficients, as  $\beta_0$  does. Now let's move on to  $\boldsymbol{\beta}$  and  $\boldsymbol{\beta}^c$ .

$$\frac{\partial L_{\text{ridge}}}{\partial \boldsymbol{\beta}} = -2 \mathbf{X}^T (\mathbf{y} - \beta_0 \mathbf{e} - \mathbf{X}\boldsymbol{\beta}) + 2 \lambda \boldsymbol{\beta}$$

Note that

$$\beta_0 \mathbf{e} = \mathbf{e} \bar{y} - \frac{1}{N} \mathbf{e} \mathbf{e}^T \mathbf{X} \boldsymbol{\beta} = \mathbf{A} \mathbf{y} - \mathbf{A} \mathbf{X} \boldsymbol{\beta}$$

So we have that

$$\begin{aligned} \frac{\partial L_{\text{ridge}}}{\partial \boldsymbol{\beta}} &= -2 \mathbf{X}^T (\mathbf{y} - \mathbf{A} \mathbf{y} + \mathbf{A} \mathbf{X} \boldsymbol{\beta} - \mathbf{X} \boldsymbol{\beta}) + 2 \lambda \boldsymbol{\beta} = 0 \\ \rightarrow \boldsymbol{\beta} &= (\mathbf{X}^T (\mathbf{I} - \mathbf{A}) \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T (\mathbf{I} - \mathbf{A}) \mathbf{y} \end{aligned}$$

From this we can get  $\boldsymbol{\beta}^c$  (by the substitution  $\mathbf{X} \rightarrow (\mathbf{I} - \mathbf{A}) \mathbf{X}$ ). Note that  $\mathbf{I} - \mathbf{A}$  is a projection, so its square is itself, as well as its transpose.

$$\begin{aligned} \boldsymbol{\beta}^c &= ((\mathbf{I} - \mathbf{A}) \mathbf{X})^T (\mathbf{I} - \mathbf{A}) (\mathbf{I} - \mathbf{A}) \mathbf{X} + \lambda \mathbf{I})^{-1} ((\mathbf{I} - \mathbf{A}) \mathbf{X})^T (\mathbf{I} - \mathbf{A}) \mathbf{y} \\ &= (\mathbf{X}^T (\mathbf{I} - \mathbf{A}) (\mathbf{I} - \mathbf{A}) (\mathbf{I} - \mathbf{A}) \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T (\mathbf{I} - \mathbf{A}) (\mathbf{I} - \mathbf{A}) \mathbf{y} \\ &= (\mathbf{X}^T (\mathbf{I} - \mathbf{A}) \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T (\mathbf{I} - \mathbf{A}) \mathbf{y} \\ &= \boldsymbol{\beta} \end{aligned}$$

So the conclusion is that by centering the inputs, we get the same result for  $\boldsymbol{\beta}$  and  $\boldsymbol{\beta}^c$ , while  $\beta_0$  gets replaced by  $\beta_0^c = \bar{y}$ .

### 3.3.5 Ex. 3.6

- $\mathbf{Y} = \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\epsilon}$
- $\boldsymbol{\beta} \sim N(\mathbf{0}, \tau^2 \mathbf{I}_p)$
- $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_N)$

**Solution 1.** Using Appendix (C.3),  $U = \mathbf{X} \boldsymbol{\beta}$ ,  $V = \boldsymbol{\epsilon}$ ,  $W = U + V = \mathbf{Y}$ .  $\mathbf{X} \boldsymbol{\beta} \sim N(\mathbf{0}, \tau^2 \mathbf{X} \mathbf{X}^T)$ . With the notations in Appendix (C.3),  $\boldsymbol{\mu}_U = \mathbf{0}$ ,  $\boldsymbol{\mu}_V = \mathbf{0}$ ,  $\boldsymbol{\Sigma}_U = \tau^2 \mathbf{X} \mathbf{X}^T$ ,  $\boldsymbol{\Sigma}_V = \sigma^2 \mathbf{I}_N$ .  $U|W$  is normally distributed with mean (and mode)

$$\boldsymbol{\mu}_{U|W} = \tau^2 \mathbf{X} \mathbf{X}^T (\tau^2 \mathbf{X} \mathbf{X}^T + \sigma^2 \mathbf{I}_N)^{-1} \mathbf{y}$$

Using the matrix inversion lemma described in Appendix (E.1), we can write that



$$\begin{aligned}
(\tau^2 \mathbf{X} \mathbf{X}^T + \sigma^2 \mathbf{I}_N)^{-1} &= \frac{1}{\sigma^2} \left( \frac{\tau^2}{\sigma^2} \mathbf{X} \mathbf{X}^T + \mathbf{I}_N \right)^{-1} \\
&= \frac{1}{\sigma^2} \left( \mathbf{I}_N - \frac{\tau}{\sigma} \mathbf{X} \left( \mathbf{I}_p + \frac{\tau^2}{\sigma^2} \mathbf{X}^T \mathbf{X} \right)^{-1} \frac{\tau}{\sigma} \mathbf{X}^T \right) \\
&= \frac{1}{\sigma^2} \mathbf{I}_N - \frac{1}{\sigma^2} \mathbf{X} \left( \frac{\sigma^2}{\tau^2} \mathbf{I}_p + \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \\
\boldsymbol{\mu}_{U|W} &= \frac{\tau^2}{\sigma^2} \mathbf{X} \mathbf{X}^T \mathbf{y} - \frac{\tau^2}{\sigma^2} \mathbf{X} \mathbf{X}^T \mathbf{X} \left( \frac{\sigma^2}{\tau^2} \mathbf{I}_p + \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y}
\end{aligned}$$

Now with this we can express  $\boldsymbol{\mu}_{\beta|y}$ :

$$\begin{aligned}
\boldsymbol{\mu}_{\beta|y} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\mu}_{U|W} \\
&= \frac{\tau^2}{\sigma^2} \mathbf{X}^T \mathbf{y} - \frac{\tau^2}{\sigma^2} \mathbf{X}^T \mathbf{X} \left( \frac{\sigma^2}{\tau^2} \mathbf{I}_p + \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y} \\
&= \left[ \frac{\tau^2}{\sigma^2} \left( \frac{\sigma^2}{\tau^2} \mathbf{I}_p + \mathbf{X}^T \mathbf{X} \right) - \frac{\tau^2}{\sigma^2} \mathbf{X}^T \mathbf{X} \right] \left( \frac{\sigma^2}{\tau^2} \mathbf{I}_p + \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y} \\
&= \left( \frac{\sigma^2}{\tau^2} \mathbf{I}_p + \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y}
\end{aligned}$$

So this is the mean (and mode, as it is normally distributed) of  $\beta$ , given  $\mathbf{y}$ . We got the same expression as the one for ridge regression, when we have

$$\lambda = \frac{\sigma^2}{\tau^2}$$

**Solution 2.** In Appendix (C.3), we assume both  $U$  and  $V$  have a non-singular covariance matrix. But in our case it is not the case, since  $\mathbf{X} \mathbf{X}^T$  is not invertible. But we can directly apply the Bayes-formula:

$$f_{\beta|y}(\boldsymbol{\beta}|\mathbf{y}) = \frac{f_{y|\beta}(\mathbf{y}|\boldsymbol{\beta})f_{\beta}(\boldsymbol{\beta})}{f_y(\mathbf{y})}$$

From this (again, after a lot of computation) we get that  $\beta|y$  is normally distributed and get that the mean is the one we got in solution 1.

### 3.3.6 Ex. 3.8

I assume that  $\mathbf{X}$  has full rank, i.e.,  $\text{rank}(\mathbf{X}) = p + 1$ . We can write that  $\mathbf{X} = [\mathbf{e}|\mathbf{X}_2]$ , where  $\mathbf{e}$  is the vector with all ones, and  $\mathbf{X}_2$  is  $N \times p$  matrix. Denote the column space of  $\mathbf{X}$  by  $W$ , the subspace spanned by  $\mathbf{e}$  is  $W_e$ , and so

$$W = W_e + W_e^\perp$$

By centering  $\mathbf{X}_2$ , we remove the projections of its columns onto  $\mathbf{e}$ , and so those columns will span  $W_e^\perp$ . the projection onto the line spanned by  $\mathbf{e}$  is

$$\mathbf{P}_e = \frac{1}{N}\mathbf{e}\mathbf{e}^T$$

With this we have that

$$\tilde{\mathbf{X}} = (\mathbf{I} - \mathbf{P}_e)\mathbf{X}_2$$

because this way the columns of  $\tilde{\mathbf{X}}$  are centered. The column space of  $\tilde{\mathbf{X}}$  is  $W_e^\perp$ .

Considering the QR decomposition of  $\mathbf{X}$ , we can write:

$$\mathbf{Q} = [\mathbf{e}/\sqrt{N}|\mathbf{Q}_2] = \frac{1}{\sqrt{N}}[\mathbf{e}|\sqrt{N}\mathbf{Q}_2]$$

The SVD of  $\tilde{\mathbf{X}}$ :

$$\tilde{\mathbf{X}} = \mathbf{U}\mathbf{D}\mathbf{V}^T$$

Denote the column space of (arbitrary)  $\mathbf{A}$  by  $R(\mathbf{A})$ . Then

$$R(\tilde{\mathbf{X}}) = R(\mathbf{U}) = W_e^\perp$$

$$R(\mathbf{Q}_2) = W_e^\perp$$

### 3.3.7 Ex. 3.9 Forward stepwise regression

Actually we don't need  $\mathbf{Q}$ , we only need  $\mathbf{H} = \mathbf{Q}\mathbf{Q}^T$ , and updating this matrix. Some notations first.  $W_1$  is the column space of  $\mathbf{X}_1$ .  $\mathbf{u}$  is an arbitrary column vector of  $\mathbf{X}_2$ .  $\mathbf{X}_{1u}$  is the matrix  $[\mathbf{X}_1 | \mathbf{u}]$ .  $W_{1u}$  is the column space of  $\mathbf{X}_{1u}$ .  $\mathbf{H}_1$  is the projection matrix that projects onto  $W_1$ .  $\mathbf{H}_{1u}$  projects onto  $W_{1u}$ .

We know the residual:

$$\mathbf{r} = \mathbf{y} - \hat{\mathbf{y}} = \mathbf{y} - \mathbf{H}_1\mathbf{y}$$

Now the new residual:

$$\mathbf{r}' = \mathbf{y} - \hat{\mathbf{y}}' = \mathbf{y} - \mathbf{H}_{1u}\mathbf{y}$$

Let's calculate  $\mathbf{H}_{1u}$  efficiently. It projects onto  $W_{1u}$ . We have an orthonormal base in  $W_1$ , namely, the columns of  $\mathbf{Q}$ . We need an orthonormal base in  $W_{1u}$  as well. So let's regress  $\mathbf{u}$  onto  $W_1$ .  $\mathbf{u}_p = \mathbf{H}_1\mathbf{u}$  is the projection of  $\mathbf{u}$  onto  $W_1$ , so  $\mathbf{u} - \mathbf{H}_1\mathbf{u}$  is orthogonal to  $W_1$ . Let's normalize this vector and append it to  $\mathbf{Q}$ .

$$\mathbf{q} = \frac{\mathbf{u} - \mathbf{H}_1\mathbf{u}}{\|\mathbf{u} - \mathbf{H}_1\mathbf{u}\|}$$

With this we can construct  $\mathbf{Q}' = [\mathbf{Q} | \mathbf{q}]$ , and the columns of this matrix are orthonormal. The projection matrix:

$$\mathbf{H}_{1u} = \mathbf{Q}'\mathbf{Q}'^T = \mathbf{Q}\mathbf{Q}^T + \mathbf{q}\mathbf{q}^T = \mathbf{H}_1 + \mathbf{q}\mathbf{q}^T$$

The new residual:

$$\mathbf{r}' = \mathbf{y} - \mathbf{H}_{1u}\mathbf{y} = \mathbf{y} - \mathbf{H}_1\mathbf{y} - \mathbf{q}\mathbf{q}^T\mathbf{y} = \mathbf{r} - \mathbf{q}\mathbf{q}^T\mathbf{y}$$

From this:

$$\text{RSS}' = \text{RSS} - 2\mathbf{r}^T\mathbf{q}\mathbf{q}^T\mathbf{y} + \mathbf{y}^T\mathbf{q}\mathbf{q}^T\mathbf{y}$$

The drop in RSS:

$$\text{RSS} - \text{RSS}' = (2\mathbf{r} - \mathbf{y})^T\mathbf{q}\mathbf{q}^T\mathbf{y}$$

So we iterate through the columns of  $\mathbf{X}_2$  and seek for the largest drop in RSS. Once we have the best column ( $\mathbf{u}^*$ ), we pass  $\mathbf{H}_{1u^*}$  to the next iteration. Consider the notebook "forward\_stepwise.ipynb" where I have implemented this algorithm.

### 3.3.8 Ex. 3.10 Backward stepwise regression

Now we seek for a column vector  $\mathbf{x}$  that we can drop from  $\mathbf{X}$ . Notations.  $\mathbf{x}_i$  is the  $i$ th column vector of  $\mathbf{X}$ .  $\mathbf{X}_i$  is the matrix obtained from  $\mathbf{X}$  by dropping the  $i$ th column.  $\mathbf{B} \equiv (\mathbf{X}^T \mathbf{X})^{-1}$ . The  $i$ th column vector of  $\mathbf{B}$  is  $\mathbf{b}_i$ .

Now,  $\mathbf{H} = \mathbf{X} \mathbf{B} \mathbf{X}^T$  projects onto the column space of  $\mathbf{X}$ . According to Appendix (D.3),  $\mathbf{H}_i = \mathbf{H} - \mathbf{P}_i$  projects onto the column space of  $\mathbf{X}_i$ , where

$$\mathbf{P}_i = \frac{\mathbf{X} \mathbf{b}_i \mathbf{b}_i^T \mathbf{X}^T}{v_i}$$

and  $v_i$  is the  $i$ th element in the diagonal of  $\mathbf{B}$ . The original residual:

$$\mathbf{r} = (\mathbf{I} - \mathbf{H})\mathbf{y}$$

The new residual (after dropping the  $i$ th column in  $\mathbf{X}$ ):

$$\mathbf{r}_i = (\mathbf{I} - \mathbf{H}_i)\mathbf{y} = \mathbf{r} + \mathbf{P}_i \mathbf{y}$$

From this we get the increment in RSS introduced by dropping the feature:

$$\Delta \text{RSS}_i = (2\mathbf{r} + \mathbf{y})^T \mathbf{P}_i \mathbf{y}$$

We seek for a column vector of  $\mathbf{X}$  for which  $\Delta \text{RSS}_i$  is minimal. We drop that feature. Consider the notebook "backward\_stepwise.ipynb" where I implement this algorithm.

### 3.3.9 Ex. 3.11

$$\begin{aligned} \text{RSS}(\mathbf{B}) &= \sum_{i=1}^N (y_i - f(x_i))^T \Sigma^{-1} (y_i - f(x_i)) \\ &= \text{Tr}((\mathbf{Y} - \mathbf{X} \mathbf{B}) \Sigma^{-1} (\mathbf{Y} - \mathbf{X} \mathbf{B})^T) \\ &= \text{Tr}(\mathbf{Y} \Sigma^{-1} \mathbf{Y}^T) - \text{Tr}(\mathbf{Y} \Sigma^{-1} \mathbf{B}^T \mathbf{X}^T) - \text{Tr}(\mathbf{X} \mathbf{B} \Sigma^{-1} \mathbf{Y}^T) + \text{Tr}(\mathbf{X} \mathbf{B} \Sigma^{-1} \mathbf{B}^T \mathbf{X}^T) \end{aligned}$$

Calculating the derivatives:

$$\frac{d}{d\mathbf{B}} \text{Tr}(\mathbf{Y} \Sigma^{-1} \mathbf{Y}^T) = \mathbf{0}$$

$$\frac{d}{d\mathbf{B}} \text{Tr}(\mathbf{Y} \Sigma^{-1} \mathbf{B}^T \mathbf{X}^T) = \frac{d}{d\mathbf{B}} \text{Tr}(\mathbf{X}^T \mathbf{Y} \Sigma^{-1} \mathbf{B}^T) = \mathbf{X}^T \mathbf{Y} \Sigma^{-1}$$

$$\frac{d}{d\mathbf{B}} \text{Tr}(\mathbf{X}\mathbf{B}\mathbf{\Sigma}^{-1}\mathbf{Y}^T) = \frac{d}{d\mathbf{B}} \text{Tr}(\mathbf{\Sigma}^{-1}\mathbf{Y}^T\mathbf{X}\mathbf{B}) = (\mathbf{\Sigma}^{-1}\mathbf{Y}^T\mathbf{X})^T = \mathbf{X}^T\mathbf{Y}\mathbf{\Sigma}^{-1}$$

$$\begin{aligned} \frac{d}{d\mathbf{B}} \text{Tr}(\mathbf{X}\mathbf{B}\mathbf{\Sigma}^{-1}\mathbf{B}^T\mathbf{X}^T) &= \frac{d}{d\mathbf{A}} \text{Tr}(\mathbf{X}\mathbf{A}\mathbf{\Sigma}^{-1}\mathbf{B}^T\mathbf{X}^T) + \frac{d}{d\mathbf{A}} \text{Tr}(\mathbf{X}\mathbf{B}\mathbf{\Sigma}^{-1}\mathbf{A}^T\mathbf{X}^T) \\ &= (\mathbf{\Sigma}^{-1}\mathbf{B}^T\mathbf{X}^T\mathbf{X})^T + \mathbf{X}^T\mathbf{X}\mathbf{B}\mathbf{\Sigma}^{-1} \\ &= 2\mathbf{X}^T\mathbf{X}\mathbf{B}\mathbf{\Sigma}^{-1} \end{aligned}$$

Now we can substitute:

$$\frac{d}{d\mathbf{B}} \text{RSS}(\mathbf{B}) = -2\mathbf{X}^T\mathbf{Y}\mathbf{\Sigma}^{-1} + 2\mathbf{X}^T\mathbf{X}\mathbf{B}\mathbf{\Sigma}^{-1}$$

Setting this to zero, we get:

$$\mathbf{B} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y}$$

What happens if the covariance matrices  $\mathbf{\Sigma}_i$  are different for each observation?  
The Residual Sum of Squares:

$$\text{RSS}(\mathbf{B}) = \sum_{i=1}^N (y_i - \mathbf{B}^T x_i)^T \mathbf{\Sigma}_i^{-1} (y_i - \mathbf{B}^T x_i)$$

Derivating this w.r.t  $\mathbf{B}$ , and setting to zero, we get:

$$\sum_{i=1}^N x_i x_i^T \mathbf{B} \mathbf{\Sigma}_i^{-1} = \sum_{i=1}^N x_i y_i^T \mathbf{\Sigma}_i^{-1}$$

# Appendices

Here I collected the useful mathematical knowledge required to understand some proofs.

## A differentiation

### A.1 differentiation w.r.t. a vector

1. Let  $\mathbf{a} \in \mathbb{R}^n$  be a constant vector,  $\mathbf{x} \in \mathbb{R}^n$ . Then

$$\frac{d}{d\mathbf{x}}(\mathbf{x}^T \mathbf{a}) = \frac{d}{d\mathbf{x}}(\mathbf{a}^T \mathbf{x}) = \mathbf{a} \quad (139)$$

2. Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be a constant matrix,  $\mathbf{x} \in \mathbb{R}^n$  Then

$$\frac{d}{d\mathbf{x}}(\mathbf{x}^T \mathbf{A} \mathbf{x}) = (\mathbf{A} + \mathbf{A}^T) \mathbf{x} \quad (140)$$

We can derive this as follows:

$$\begin{aligned} \frac{d}{d\mathbf{x}}(\mathbf{x}^T \mathbf{A} \mathbf{x}) &= \frac{d}{d\mathbf{y}}(\mathbf{y}^T \mathbf{A} \mathbf{x}) + \frac{d}{d\mathbf{y}}(\mathbf{x}^T \mathbf{A} \mathbf{y}) \\ &= \frac{d}{d\mathbf{y}}(\mathbf{y}^T \mathbf{A} \mathbf{x}) + \frac{d}{d\mathbf{y}}(\mathbf{y}^T \mathbf{A}^T \mathbf{x}) \\ &= \mathbf{A} \mathbf{x} + \mathbf{A}^T \mathbf{x} \end{aligned}$$

### A.2 differentiation w.r.t. a matrix

1. Let  $\mathbf{A} \in \mathbb{R}^{n \times m}$ ,  $\mathbf{X} \in \mathbb{R}^{m \times n}$ . Then

$$\frac{d}{d\mathbf{X}} \text{Tr}(\mathbf{A} \mathbf{X}) = \frac{d}{d\mathbf{X}} \text{Tr}(\mathbf{X} \mathbf{A}) = \mathbf{A}^T \quad (141)$$

2. Let  $\mathbf{A} \in \mathbb{R}^{n \times m}$ ,  $\mathbf{X} \in \mathbb{R}^{n \times m}$ . Then

$$\frac{d}{d\mathbf{X}} \text{Tr}(\mathbf{A} \mathbf{X}^T) = \frac{d}{d\mathbf{X}} \text{Tr}(\mathbf{X}^T \mathbf{A}) = \mathbf{A} \quad (142)$$

3. Let  $\mathbf{A} \in \mathbb{R}^{m \times m}$ ,  $\mathbf{X} \in \mathbb{R}^{m \times n}$ . Then

$$\frac{d}{d\mathbf{X}} \text{Tr}(\mathbf{X}^T \mathbf{A} \mathbf{X}) = (\mathbf{A} + \mathbf{A}^T) \mathbf{X} \quad (143)$$

We can derive it as follows:

$$\begin{aligned} \frac{d}{d\mathbf{X}} \text{Tr}(\mathbf{X}^T \mathbf{A} \mathbf{X}) &= \frac{d}{d\mathbf{Y}} \text{Tr}(\mathbf{Y}^T \mathbf{A} \mathbf{X}) + \frac{d}{d\mathbf{Y}} \text{Tr}(\mathbf{X}^T \mathbf{A} \mathbf{Y}) \\ &= \mathbf{A} \mathbf{X} + (\mathbf{X}^T \mathbf{A})^T \\ &= \mathbf{A} \mathbf{X} + \mathbf{A}^T \mathbf{X} = (\mathbf{A} + \mathbf{A}^T) \mathbf{X} \end{aligned}$$

**Example.** Consider now this example.

$$f(\mathbf{X}) = \text{Tr}(\mathbf{X}^T \mathbf{A} \mathbf{X} \mathbf{B} \mathbf{X}^T \mathbf{C})$$

where  $\mathbf{X} \in \mathbb{R}^{n \times m}$ ,  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{B} \in \mathbb{R}^{m \times m}$ ,  $\mathbf{C} \in \mathbb{R}^{n \times m}$ .

$$\begin{aligned} \frac{d}{d\mathbf{X}} f(\mathbf{X}) &= \frac{d}{d\mathbf{Y}} \text{Tr}(\mathbf{Y}^T \mathbf{A} \mathbf{X} \mathbf{B} \mathbf{X}^T \mathbf{C}) \\ &\quad + \frac{d}{d\mathbf{Y}} \text{Tr}(\mathbf{X}^T \mathbf{A} \mathbf{Y} \mathbf{B} \mathbf{X}^T \mathbf{C}) \\ &\quad + \frac{d}{d\mathbf{Y}} \text{Tr}(\mathbf{X}^T \mathbf{A} \mathbf{X} \mathbf{B} \mathbf{Y}^T \mathbf{C}) \end{aligned}$$

Calculating these:

$$\frac{d}{d\mathbf{Y}} \text{Tr}(\mathbf{Y}^T \mathbf{A} \mathbf{X} \mathbf{B} \mathbf{X}^T \mathbf{C}) = \mathbf{A} \mathbf{X} \mathbf{B} \mathbf{X}^T \mathbf{C}$$

$$\begin{aligned} \frac{d}{d\mathbf{Y}} \text{Tr}(\mathbf{X}^T \mathbf{A} \mathbf{Y} \mathbf{B} \mathbf{X}^T \mathbf{C}) &= \frac{d}{d\mathbf{Y}} \text{Tr}(\mathbf{Y} \mathbf{B} \mathbf{X}^T \mathbf{C} \mathbf{X}^T \mathbf{A}) \\ &= (\mathbf{B} \mathbf{X}^T \mathbf{C} \mathbf{X}^T \mathbf{A})^T \\ &= \mathbf{A}^T \mathbf{X} \mathbf{C}^T \mathbf{X} \mathbf{B}^T \end{aligned}$$

$$\frac{d}{d\mathbf{Y}} \text{Tr}(\mathbf{X}^T \mathbf{A} \mathbf{X} \mathbf{B} \mathbf{Y}^T \mathbf{C}) = \frac{d}{d\mathbf{Y}} \text{Tr}(\mathbf{Y}^T \mathbf{C} \mathbf{X}^T \mathbf{A} \mathbf{X} \mathbf{B}) = \mathbf{C} \mathbf{X}^T \mathbf{A} \mathbf{X} \mathbf{B}$$

So the result is:

$$\frac{d}{d\mathbf{X}} \text{Tr}(\mathbf{X}^T \mathbf{A} \mathbf{X} \mathbf{B} \mathbf{X}^T \mathbf{C}) = \mathbf{A} \mathbf{X} \mathbf{B} \mathbf{X}^T \mathbf{C} + \mathbf{A}^T \mathbf{X} \mathbf{C}^T \mathbf{X} \mathbf{B}^T + \mathbf{C} \mathbf{X}^T \mathbf{A} \mathbf{X} \mathbf{B}$$

## B variance and covariance properties

### B.1 scalar multiple

We know that if  $a \in \mathbb{R}$  is a constant scalar, and  $X$  is a scalar random variable, then

$$\text{Var}(aX) = a^2 \cdot \text{Var}(X) \quad (144)$$

### B.2 vector multiple

Now what if  $a \in \mathbb{R}^p$  is a (constant) vector, and  $X \in \mathbb{R}^p$  is a random vector, and we take the inner product  $a^T \cdot X$ ? What is the variance  $\text{Var}(a^T \cdot X)$ ?

$$\begin{aligned} \text{Var}(a^T \cdot X) &= \text{Var}(a_1 \cdot X_1 + a_2 \cdot X_2 + \cdots + a_n \cdot X_n) \\ &= \text{E} \left( \sum_i a_i \cdot (X_i - \text{E}X_i) \right)^2 \\ &= \text{E} \left( \sum_{i,j} a_i \cdot (X_i - \text{E}X_i) \cdot a_j \cdot (X_j - \text{E}X_j) \right) \\ &= \sum_{i,j} a_i \cdot a_j \cdot \text{E}((X_i - \text{E}X_i) \cdot (X_j - \text{E}X_j)) \\ &= \sum_{i,j} a_i \cdot a_j \cdot \text{Cov}(X_i, X_j) \\ &= a^T \cdot \text{Cov}(X, X) \cdot a = a^T \cdot \Sigma \cdot a \end{aligned} \quad (145)$$

Here  $\Sigma \equiv \text{Cov}(X, X)$  is the covariance matrix,  $\text{Cov}(X, X)_{i,j} = \text{Cov}(X_i, X_j)$ . Let's state our finding again.  $a \in \mathbb{R}^p$  is a constant vector,  $X \in \mathbb{R}^p$  is a random vector, then:

$$\text{Var}(a^T \cdot X) = a^T \cdot \text{Cov}(X, X) \cdot a \quad (146)$$

Furthermore, we can write for the covariance matrix:

$$\text{Cov}(X, X) = \text{E}((X - \text{E}X) \cdot (X - \text{E}X)^T) = \text{E}(X \cdot X^T) - (\text{E}X)(\text{E}X^T) \quad (147)$$



### B.3 matrix multiple

Let  $A \in \mathbb{R}^{n \times p}$  a constant matrix,  $X \in \mathbb{R}^p$  a random vector. The covariance matrix:

$$\begin{aligned}\text{Cov}(AX) &= E(AXX^T A^T) - E(AX)E(X^T A^T) \\ &= A \cdot E(XX^T) \cdot A^T - A \cdot E(X)E(X^T) \cdot A^T \\ &= A \cdot (E(XX^T) - E(X)E(X^T)) \cdot A^T \\ &= A \cdot \text{Cov}(X) \cdot A^T\end{aligned}\tag{148}$$

## C distributions

### C.1 Student's t-distribution

The t-distribution with  $\nu$  degrees of freedom can be expressed as

$$T = \frac{Z}{\sqrt{V/\nu}}\tag{149}$$

where

- $Z \sim N(0, 1)$
- $V \sim \chi_\nu^2$
- $Z$  and  $V$  are independent

### C.2 F-distribution

A random variate of the F-distribution with parameters  $d_1$  and  $d_2$  arises as the ratio of two appropriately scaled chi-squared variates:

$$X = \frac{U_1/d_1}{U_2/d_2}\tag{150}$$

where

- $U_1$  and  $U_2$  have chi-squared distributions with  $d_1$  and  $d_2$  degrees of freedom respectively, and
- $U_1$  and  $U_2$  are independent.

### C.3 Gaussian posterior distribution

- $U \sim N(\boldsymbol{\mu}_U \in \mathbb{R}^n, \boldsymbol{\Sigma}_U \in \mathbb{R}^{n \times n})$
- $V \sim N(\boldsymbol{\mu}_V \in \mathbb{R}^n, \boldsymbol{\Sigma}_V \in \mathbb{R}^{n \times n})$
- $U$  and  $V$  are independent

$W = U + V$  is also normal, as well as  $W|U$ :

- $W = U + V \sim N(\boldsymbol{\mu}_U + \boldsymbol{\mu}_V, \boldsymbol{\Sigma}_U + \boldsymbol{\Sigma}_V)$
- $W|U \sim N(U + \boldsymbol{\mu}_V, \boldsymbol{\Sigma}_V)$

The posterior distribution  $U|W$  is Gaussian, this is what I want to prove here.

- $U|W \sim N(\boldsymbol{\mu}_{U|W}, \boldsymbol{\Sigma}_{U|W})$
- $\boldsymbol{\mu}_{U|W} = \boldsymbol{\Sigma}_U(\boldsymbol{\Sigma}_U + \boldsymbol{\Sigma}_V)^{-1}(W - \boldsymbol{\mu}_V) + \boldsymbol{\Sigma}_V(\boldsymbol{\Sigma}_U + \boldsymbol{\Sigma}_V)^{-1}\boldsymbol{\mu}_U$
- $\boldsymbol{\Sigma}_{U|W} = \boldsymbol{\Sigma}_U(\boldsymbol{\Sigma}_U + \boldsymbol{\Sigma}_V)^{-1}\boldsymbol{\Sigma}_V = \boldsymbol{\Sigma}_V(\boldsymbol{\Sigma}_U + \boldsymbol{\Sigma}_V)^{-1}\boldsymbol{\Sigma}_U$

The proof is too long, I won't specify all the details here. First, assuming that both  $U$  and  $V$  have a pdf,

$$f_{U|W}(u|w) = \frac{f_{W|U}(w|u)f_U(u)}{f_W(w)}$$

where

$$f_{W|U}(w|u) = \frac{\exp\left(-\frac{1}{2}(w - u - \boldsymbol{\mu}_V)^T \boldsymbol{\Sigma}_V^{-1}(w - u - \boldsymbol{\mu}_V)\right)}{\sqrt{(2\pi)^n \det(\boldsymbol{\Sigma}_V)}}$$

$$f_U(u) = \frac{\exp\left(-\frac{1}{2}(u - \boldsymbol{\mu}_U)^T \boldsymbol{\Sigma}_U^{-1}(u - \boldsymbol{\mu}_U)\right)}{\sqrt{(2\pi)^n \det(\boldsymbol{\Sigma}_U)}}$$

$$f_W(w) = \frac{\exp\left(-\frac{1}{2}(w - \boldsymbol{\mu}_U - \boldsymbol{\mu}_V)^T (\boldsymbol{\Sigma}_U + \boldsymbol{\Sigma}_V)^{-1}(w - \boldsymbol{\mu}_U - \boldsymbol{\mu}_V)\right)}{\sqrt{(2\pi)^n \det(\boldsymbol{\Sigma}_U + \boldsymbol{\Sigma}_V)}}$$

Plugging these into the equation above, after a lot of computation, we get to

$$f_{U|W}(u|w) = \frac{\exp\left(-\frac{1}{2}(u - \boldsymbol{\mu}_{U|W})^T \boldsymbol{\Sigma}_{U|W}^{-1}(u - \boldsymbol{\mu}_{U|W})\right)}{\sqrt{(2\pi)^n \det(\boldsymbol{\Sigma}_{U|W})}}$$

## D projections

### D.1 sum of projections

Let  $V$  be a vector space,  $W_1 \subset V$  a subspace,  $W_2 \subset V$  a subspace such that  $W_1 \perp W_2$ .

Let  $P_1$  be an orthogonal projection onto the subspace  $W_1$ ,  $P_2$  be an orthogonal projection onto the subspace  $W_2$ . I claim that  $P_1 + P_2$  is an orthogonal projection onto  $W_1 + W_2$ .

*Proof.* Denote  $W_\perp$  the orthogonal complement of  $W_1 + W_2$ .

$$(W_1 + W_2) + W_\perp = V \quad (151)$$

and

$$(W_1 + W_2) \perp W_\perp \quad (152)$$

Any vector  $v \in V$  can be decomposed as

$$v = w_\perp + w_1 + w_2 \quad (153)$$

where  $w_\perp \in W_\perp$ ,  $w_1 \in W_1$ ,  $w_2 \in W_2$ . This decomposition is unique.

$$P_1 v = 0 + w_1 + 0 = w_1 \quad (154)$$

$$P_2 v = 0 + 0 + w_2 = w_2 \quad (155)$$

From these

$$(P_1 + P_2)v = w_1 + w_2 \quad (156)$$

So  $P_1 + P_2$  projects onto  $W_1 + W_2$ .

□

## D.2 difference of projections

Let  $V$  be a vector space,  $W_1 \subset V$  a subspace,  $W_2 \subset W_1$  a subspace,  $W_3 \subset W_1$  a subspace, such that  $W_2 \perp W_3$ , and  $W_2 + W_3 = W_1$ .

Let  $P_1$  be an orthogonal projection onto the subspace  $W_1$ ,  $P_2$  be an orthogonal projection onto the subspace  $W_2$ . I claim that  $P_1 - P_2$  is an orthogonal projection onto  $W_3$ .

*Proof.* Denote  $W_\perp$  the orthogonal complement of  $W_1$ :  $W_1 + W_\perp = V$ , and  $W_1 \perp W_\perp$ .

Any vector  $v \in V$  can be decomposed as

$$v = w_\perp + w_2 + w_3 \quad (157)$$

where  $w_\perp \in W_\perp$ ,  $w_2 \in W_2$ ,  $w_3 \in W_3$ . This decomposition is unique.

$$P_1 v = 0 + w_2 + w_3 = w_2 + w_3 \quad (158)$$

$$P_2 v = 0 + w_2 + 0 = w_2 \quad (159)$$

From these

$$(P_1 - P_2)v = (w_2 + w_3) - w_2 = w_3 \quad (160)$$

So  $P_1 - P_2$  projects onto  $W_3$ .

□

## D.3 The special vector $Xb$

(I couldn't find any better name for this subsection, sorry for this...) Let's begin with the  $N \times p$  matrix  $X$ , where we denote the column vectors by  $x_i$ . Assume that  $X$  has a full column-rank, so  $\text{rank}(X) = p$ . Denote the subspace  $W = \text{span}(x_1, x_2, \dots, x_{p-1})$ , which is the subspace generated by all the columns of  $X$ , except the last one. Let  $b$  be the last column vector of  $(X^T X)^{-1}$ . I claim that  $Xb$  is a vector that is perpendicular to  $W$ . Obviously,  $Xb$  is in the column space of  $X$ . We have that

$$X^T X (X^T X)^{-1} = I \quad (161)$$

Considering the  $i$ th row,  $j$ th column ( $q_j$  being the  $j$ th column vector of  $(X^T X)^{-1}$ , so  $q_p = b$ )

$$x_i^T X q_j = \delta_{i,j} \quad (162)$$

Choosing  $j = p$

$$x_i^T X b = \delta_{i,p} \quad (163)$$

This means that  $Xb$  is perpendicular to  $x_i$  ( $i \neq p$ ), which is what I wanted to prove. Now let's project onto the subspace  $W_p = \text{span}(Xb)$ :

$$P_p = \frac{Xb \cdot b^T X^T}{b^T X^T X b} = \frac{Xb \cdot b^T X^T}{v} \quad (164)$$

where  $v \equiv b_p$  is the last element of the vector  $b$ , that is,  $v \equiv [(X^T X)^{-1}]_{p,p}$ .

Now we can create the same projection according to Appendix D.2. Let  $P_X$  be the projection onto the column space of  $X$ , and  $P_W$  the projection onto  $W$ :

$$P_X = X(X^T X)^{-1} X^T \quad (165)$$

$$P_W = X_0(X_0^T X_0)^{-1} X_0^T \quad (166)$$

where we get  $X_0$  from  $X$  by dropping the last column. Now we see that

$$X(X^T X)^{-1} X^T - X_0(X_0^T X_0)^{-1} X_0^T = \frac{Xb \cdot b^T X^T}{v} \quad (167)$$

## E Matrices

### E.1 The matrix inversion lemma

A special form of the matrix inversion lemma:

$$(I_n + UV)^{-1} = I_n - U(I_m + VU)^{-1}V$$

where  $U \in \mathbb{R}^{n \times m}$ ,  $V \in \mathbb{R}^{m \times n}$ .

Proof.

$$(I_n + UV)^{-1} = I_n - (I_n + UV)^{-1}UV = I_n - U(I_m + VU)^{-1}V$$

At first we used the following:

$$(\mathbf{I} + \mathbf{P})^{-1} = \mathbf{I} - (\mathbf{I} + \mathbf{P})^{-1}\mathbf{P}$$

which comes from:

$$\begin{aligned} (\mathbf{I} + \mathbf{P})^{-1}(\mathbf{I} + \mathbf{P}) &= \mathbf{I} \\ (\mathbf{I} + \mathbf{P})^{-1} + (\mathbf{I} + \mathbf{P})^{-1}\mathbf{P} &= \mathbf{I} \\ (\mathbf{I} + \mathbf{P})^{-1} &= \mathbf{I} - (\mathbf{I} + \mathbf{P})^{-1}\mathbf{P} \end{aligned}$$

Secondly, we used:

$$(\mathbf{I}_n + \mathbf{UV})^{-1}\mathbf{U} = \mathbf{U}(\mathbf{I}_m + \mathbf{VU})^{-1}$$

which comes from:

$$\begin{aligned} \mathbf{U}(\mathbf{I}_m + \mathbf{VU}) &= (\mathbf{I}_n + \mathbf{UV})\mathbf{U} \\ (\mathbf{I}_n + \mathbf{UV})^{-1}\mathbf{U}(\mathbf{I}_m + \mathbf{VU}) &= \mathbf{U} \\ (\mathbf{I}_n + \mathbf{UV})^{-1}\mathbf{U} &= \mathbf{U}(\mathbf{I}_m + \mathbf{VU})^{-1} \end{aligned}$$

which concludes the proof.