

# Notes

STA314H1 - Fall 2020

Ziyue Yang

October 30, 2020

## Contents

<b>1</b>	<b>Week 5</b>	<b>2</b>
1.1	The Singular Value Decomposition . . . . .	2
1.2	Ridge Regression . . . . .	3
<b>2</b>	<b>Week 7</b>	<b>4</b>
2.1	Introduction to Lasso . . . . .	4

# 1 Week 5

## 1.1 The Singular Value Decomposition

**Remark 1.1.** We know that  $\|X\|_2$  is the square root of the largest eigenvalue of  $X^T X$ , hence at some point, whenever we see one of the inner product matrices, we should recall the **PCA**.

The SVD is a good way to understand exactly how PCA is working in terms of the feature matrix  $X$ .

**Theorem 1.1** (The Singular Value Decomposition). Assume that  $p < n$ .

Let  $X$  be an  $n \times p$  matrix. Then there exists an orthogonal  $p \times p$  matrix  $V$  (i.e.  $V^T V = V V^T = I$ ) and an orthogonal  $n \times n$  matrix  $U$  such that

$$U^T X V = D, \quad (1.1)$$

where  $D = \text{diag}(\sigma_1, \dots, \sigma_p)$ , and  $\sigma_1 \geq \dots \geq \sigma_p \geq 0$ .

*Proof.* Omitted for now. □

**Remark 1.2.** We can express the SVD of an  $n \times p$  matrix  $X$  in several equivalent ways:

1. As a singular tuple  $(\sigma, u, v)$  that satisfies  $Xv = \sigma u$  and  $X^T u = \sigma v$ .
2. As a matrix decomposition  $X = U D V^T$ , where  $V$  is a  $p \times p$  orthogonal matrix, and  $U$  is a  $n \times n$  orthogonal matrix.
3. As a way of representing the matrix as a sum

$$X = \sum_{j=1}^p \sigma_j u_j v_j^T. \quad (1.2)$$

**Remark 1.3** (The SVD and Principal Components). Recall that the factor loadings are the eigenvectors of  $X^T X$ .

If  $X = U D V^T$ , then  $X^T X = V^T D U^T U D V^T = V D^2 V^T$ .

- The  $V$  in the SVD is exactly the matrix of factor loadings.
- The eigenvalues of  $X^T X$  are the squares of the singular values.

Note that the score vectors were defined as  $t_j = X v_j$ , and We can use one of the representations of singular vectors to see that

$$t_j = X v_j = \sigma_j u_j. \quad (1.3)$$

**Remark 1.4.** The SVD makes it easy to solve the normal equations.

Recall that

$$\hat{\beta} = (X^T X)^{-1} X^T y \quad (1.4)$$

$$= V D^{-2} V^T V D U^T y \quad (1.5)$$

$$= V D^{-1} U^T y \quad (1.6)$$

$$= \sum_{j=1}^p \frac{u_j^T y}{\sigma_j} v_j. \quad (1.7)$$

PCR just snipes off the small eigenvectors:

$$\hat{\beta}_{\text{pcr}} = V_k D_k^{-2} V_k^T V D U^T = \sum_{j=1}^k \frac{u_j^T y}{\sigma_j} v_j. \quad (1.8)$$

## 1.2 Ridge Regression

What does overfitting look like?

*Skipped, TODO.*

## 2 Week 7

### 2.1 Introduction to Lasso

**Remark 2.1.** Ridge regression stabilizes the least-squares estimates by shrinking low-variance directions, which makes it like a *softer* version of **principal component regression**.

Can we use penalized regression to make a softer version of variable selection? Yes. But we need to use a different penalty.