

Statistical Machine Learning

Exercise sheet 3

Exercise 3.1 (Singular Value Decomposition.) Let \mathbf{X} be a $n \times p$ matrix of rank r .

- (a) Show that $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^\top$, where \mathbf{U} and \mathbf{V} are $n \times r$ and $p \times r$ semi-orthogonal matrices respectively, in that $\mathbf{U}^\top\mathbf{U} = \mathbf{V}^\top\mathbf{V} = \mathbf{I}_r$ (while $\mathbf{U}\mathbf{U}^\top$ and $\mathbf{V}\mathbf{V}^\top$ are projections but need not be identity matrices unless \mathbf{U} or \mathbf{V} are square matrices, in which case the definition is equivalent to that of an orthogonal matrix.) and \mathbf{D} is a $r \times r$ diagonal matrix with r non-negative entries which we refer to as the singular values of \mathbf{X} . This is sometimes called a *compact* SVD.

Hint: For the special case $r = p$, consider the eigendecomposition of $\mathbf{X}^\top\mathbf{X} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^\top$ and let $\mathbf{U} = \mathbf{X}\mathbf{V}\mathbf{\Lambda}^{-1/2}$.

- (b) Show that $\mathbf{X} = \tilde{\mathbf{U}}\tilde{\mathbf{D}}\tilde{\mathbf{V}}^\top$ where $\tilde{\mathbf{U}}$ and $\tilde{\mathbf{V}}$ are $n \times n$ and $p \times p$ orthogonal matrices respectively and $\tilde{\mathbf{D}}$ is a $n \times p$ diagonal matrix with r non-negative entries.
- (c) To interpret the result from (b) we need to understand the orthonormal basis $(\mathbf{v}_1, \dots, \mathbf{v}_p)$ of \mathbb{R}^p . For each $j = 1, \dots, p$, calculate the norm of $\mathbf{X}\mathbf{v}_j$ and discuss the implication of (b); note that $\mathbf{X}\mathbf{v}_j$ are the j th coordinates of each row vector of \mathbf{X} in the basis $(\mathbf{v}_1, \dots, \mathbf{v}_p)$.

*Note: The vectors $\mathbf{v}_1, \dots, \mathbf{v}_p \in \mathbb{R}^p$ are called the *principal components directions* of \mathbf{X} . The vectors $\mathbf{u}_1, \dots, \mathbf{u}_p \in \mathbb{R}^n$ are the *normalized principal component scores*. We will go back to this in this class when discussing unsupervised learning (PCA).*

- (d) Use the compact singular value decomposition $\mathbf{U}\mathbf{D}\mathbf{V}^\top$ of \mathbf{X} (full rank) to show that applying the hat matrix $\mathbf{H} = \mathbf{X}[\mathbf{X}^\top\mathbf{X}]^{-1}\mathbf{X}^\top$ to a vector \mathbf{y} projects it onto the subspace spanned by the columns $\{\mathbf{u}_j\}$ of \mathbf{U} as in (a). In other words, \mathbf{H} is a projection matrix. Using this result, show that the fitted value of \mathbf{y} from ordinary least squares can be written in the following way:

$$\hat{\mathbf{y}}^{\text{ols}} = \sum_{j=1}^r \mathbf{u}_j \mathbf{u}_j^\top \mathbf{y}.$$

Exercise 3.2 § (Geometric interpretation of linear ridge regression) In this question, we will assume that the design matrix \mathbf{X} is an $n \times p$ full-rank matrix.

- (a) Show that

$$\begin{aligned} \hat{\boldsymbol{\beta}}^{\text{ridge}} &= \underset{\boldsymbol{\beta} \in \mathbb{R}^p}{\operatorname{argmin}} \frac{1}{n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|_2^2 \\ &= \mathbf{V}(\mathbf{D}^2 + n\lambda\mathbf{I})^{-1}\mathbf{D}\mathbf{U}^\top \mathbf{y}. \end{aligned}$$

(b) Show that the fitted values from ridge regression can be written as

$$\hat{\mathbf{y}}^{\text{ridge}}(\lambda) = \sum_{j=1}^p \frac{d_j^2}{d_j^2 + n\lambda} \mathbf{u}_j \mathbf{u}_j^\top \mathbf{y},$$

where the d_j are the diagonal elements of the matrix \mathbf{D} from (b). Discuss.

Exercise 3.3 (Artificial feature noising in linear regression) This exercise is inspired by the paper *Dropout Training as Adaptive Regularization*, by S. Wager, S. Wang, and P. Liang, published in Advances in Neural Information Processing Systems (NIPS), in 2013. In this exercise we focus on the case of feature noising applied to linear regression.

Consider linear regression with no intercept (recall from Exercise 2.4(a) that the intercept can always be estimated in a first step, so there is no loss of generality of assuming it is zero). Let $\mathbf{y} \in \mathbb{R}^n$ denote the response vector, and $\mathbf{X} \in \mathbb{R}^{n \times p}$ the design matrix with centered columns. Recall the least-squares estimator of the regression parameter

$$\hat{\boldsymbol{\beta}}^{\text{ols}} = \underset{\boldsymbol{\beta} \in \mathbb{R}^p}{\operatorname{argmin}} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2.$$

Consider the following artificial feature noising scheme: the features \mathbf{x}_i are replaced by the artificial features $\tilde{\mathbf{x}}_i = \nu(\mathbf{x}_i, \boldsymbol{\xi}_i)$ where ν is a noising function and the $\boldsymbol{\xi}_i$ are independent random vectors from some distribution $P_{\boldsymbol{\xi}}$. In practice, artificial features noising is implemented as follows:

1. For $l = 1, \dots, n$, generate a large number L of independent $\tilde{\mathbf{x}}_{il} = \nu(\mathbf{x}_i, \boldsymbol{\xi}_{il})$, where $\boldsymbol{\xi}_{il} \sim P_{\boldsymbol{\xi}}$ are independent copies of $\boldsymbol{\xi}_i$.
2. Estimate the regression parameter by

$$\hat{\boldsymbol{\beta}}^L = \underset{\boldsymbol{\beta} \in \mathbb{R}^p}{\operatorname{argmin}} \sum_{i=1}^n \sum_{l=1}^L (y_i - \tilde{\mathbf{x}}_{il}^\top \boldsymbol{\beta})^2. \quad (1)$$

The goal of this exercise is to study the link between feature noising and regularization. In theory, we study the limiting behavior of $\hat{\boldsymbol{\beta}}^L$ as $L \rightarrow \infty$. We will assume that

$$\hat{\boldsymbol{\beta}}^\infty = \underset{\boldsymbol{\beta} \in \mathbb{R}^p}{\operatorname{argmin}} \sum_{i=1}^n \mathbb{E}_{\boldsymbol{\xi}_i} \left\{ (y_i - \tilde{\mathbf{x}}_i^\top \boldsymbol{\beta})^2 \right\}.$$

- (a) Consider the following additive feature noising scheme: $\tilde{\mathbf{x}}_i = \nu(\mathbf{x}_i, \boldsymbol{\xi}_i) = \mathbf{x}_i + \boldsymbol{\xi}_i$ where $\boldsymbol{\xi}_i \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ for some $\sigma^2 > 0$. Show that the resulting estimator $\hat{\boldsymbol{\beta}}^\infty$ is a ridge estimator $\hat{\boldsymbol{\beta}}^{\text{ridge}}(\lambda)$ for some value of λ in terms of σ^2 .
- (b) Consider the following dropout noising scheme: $\tilde{\mathbf{x}}_i = \nu(\mathbf{x}_i, \boldsymbol{\xi}_i) = \mathbf{x}_i \cdot \boldsymbol{\xi}_i$, where the operator \cdot denotes the element-wise product of two vectors, and the components ξ_{ij} of the vectors $\boldsymbol{\xi}_i$ are independent random variables with

$$\xi_{ij} = \begin{cases} 0, & \text{with probability } \delta, \\ \frac{1}{1-\delta}, & \text{with probability } 1-\delta, \end{cases}$$

i.e., the ξ_{ij} follow independent scaled Bernoulli distributions, for some parameter $\delta \in (0, 1)$. Show that the resulting estimator $\hat{\beta}^\infty$ is a penalized least-squares estimator, and express the regularization term in terms of δ . What happens when the features are normalized so that $\sum_{i=1}^n x_{ij}^2 = 1$ for $j = 1, \dots, p$?

Practical exercise

Exercise 3.4 (Reading for next week's practical) Read §§6.5,6.6 from ISL. This is to familiarize yourself with the packages and the functions in R you will need for the purpose of next week's exercises.