

# Machine Learning Lecture Notes

Haoyu Zhen

June 20, 2022

# Contents

<b>1</b>	<b>Foundations</b>	<b>4</b>
1.1	Model evaluation: . . . . .	4
1.2	Performance . . . . .	4
1.3	Bias-Variance Decomposition . . . . .	4
<b>2</b>	<b>Regression</b>	<b>5</b>
2.1	Linear Regression . . . . .	5
2.2	Ridge Regression . . . . .	5
2.3	Lasso Regression . . . . .	6
2.4	Logistic Regression . . . . .	6
<b>3</b>	<b>Generalized Linear Models</b>	<b>7</b>
3.1	The Exponential Family . . . . .	7
3.2	Constructing GLMs . . . . .	7
<b>4</b>	<b>Kernel Method</b>	<b>8</b>
4.1	LMS with Features . . . . .	8
4.2	Properties of Kernels . . . . .	8
<b>5</b>	<b>Support Vector Machines</b>	<b>9</b>
5.1	Hard-SVM . . . . .	9
5.1.1	The Sample Complexity of Hard-SVM* . . . . .	9
5.2	Soft-SVM and Norm Regularization . . . . .	10
5.3	Duality . . . . .	10
5.4	SMO . . . . .	11
5.5	Implementing Soft-SVM Using SGD . . . . .	11
<b>6</b>	<b>Clustering</b>	<b>12</b>
<b>7</b>	<b>Decision Tree</b>	<b>13</b>
7.1	Random Forest . . . . .	14
<b>8</b>	<b>Dimensionality Reduction</b>	<b>15</b>
8.1	Principal Component Analysis . . . . .	15
8.2	Implementation . . . . .	16

## Preface

### Acknowledgement

These Notes contain material developed and copyright by:

- *CS229, Machine Learning*, © 2020 Summer by Andrew Ng, Stanford. Course website:
- *Understanding Machine Learning*, © 2014 by Shai Shalev-Shwartz and Shai Ben-David.
- *CS189, Introduction to Machine Learning*, © 2021 by Marvin Zhang, UCB.
- *The Elements of Statistical Learning*, © 2017 by Trevor Hastie, Robert Tibshirani and Jerome Friedman.
- *AI-2611, Machine Learning*, © 2022 by Bingbing Ni, Shanghai Jiao Tong University.

### Miscellaneous

Follow me or Contact me:

- Email: `anye[underscore]zhen[at]sjtu[dot]edu[dot]cn`.
- GitHub: [anyeZHY](#).

# 1 Foundations

## 1.1 Model evaluation:

Hold-out, cross validation and bootstrap.

For cross validation, we often let the numbers of the folds be 10. And in bootstrap, the equation  $\lim_{n \rightarrow \infty} (1 - 1/m)^m = 1/e$  is used to analyse the probability.

## 1.2 Performance

**Definition 1.1** (Sensitivity and FPR). Now we consider that

	prediction+	prediction-
Actual	1	0
1	TP	FP
0	FN	TN

$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}, \text{FPR} = \frac{\text{FP}}{\text{TN} + \text{FP}}.$$

**Remark 1.1.** ROC space and **AUC** are also useful to select models.

**Definition 1.2** (Precision and recall).

$$\text{precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}, \text{recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}.$$

$$F_\beta = \frac{(1 + \beta^2) \times P \times R}{\beta^2 \times P + R}.$$

$\beta$  depends on the preference of Precision and Recll.

## 1.3 Bias-Variance Decomposition

**Theorem 1.1.**

$$\begin{aligned} E(f; D) &= \text{bias}^2(x) + \text{var}(x) + \varepsilon^2 \\ &= (\bar{f}(x) - y)^2 + \mathbb{E}_D[f(x; D) - \bar{f}(x)] + \mathbb{E}_D[(y_D - y)^2] \end{aligned}$$

## 2 Regression

### 2.1 Linear Regression

The hypothesis class of linear regression predictors is simply the set of linear functions,

$$\mathcal{H}_{reg} = \{\mathbf{x} \mapsto \langle \mathbf{w}, \mathbf{x} \rangle + b : \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}\}.$$

Intuitively,

$$\mathcal{L}_S(h) = \frac{1}{m} \sum_{i=1}^m (h(\mathbf{x}_i) - y_i)^2, \forall h \in \mathcal{H}_{reg}.$$

To minimize the loss function, we need to solve  $A\mathbf{w} = \mathbf{b}$  where  $A \stackrel{def}{=} \sum \mathbf{x}_i \mathbf{x}_i^T = XX^T$  and  $\mathbf{b} \stackrel{def}{=} \sum y_i \mathbf{x}_i = X^T \mathbf{y}$ . If  $A$  is invertible then the solution is  $\mathbf{w} = A^{-1}\mathbf{b}$ .

#### Theorem 2.1.

$$\boldsymbol{\omega} = (X^T X)^{-1} X^T \mathbf{y}.$$

If the training instances do not span the entire space of  $\mathbb{R}^d$  then  $A$  is not invertible.

**Theorem 2.2.** Using  $A$ 's eigenvalue decomposition, we could write  $A$  as  $VD^+V^T$  where  $D$  is a diagonal matrix and  $V$  is an orthonormal matrix. Define  $D^+$  to be the diagonal matrix such that  $D_{i,i}^+ = 0$  if  $D_{i,i} = 0$  otherwise  $D_{i,i}^+ = 1/D_{i,i}$ . Then,

$$A\hat{\mathbf{w}} = \mathbf{b}$$

where  $\hat{\mathbf{w}} = VD^+V^T\mathbf{b}$

#### Proof.

$$A\hat{\mathbf{w}} = AA^+\mathbf{b} = VDV^TVD^+V^T\mathbf{b} = VDD^+V^T\mathbf{b} = \sum_{i:D_{i,i} \neq 0} \mathbf{v}_i \mathbf{v}_i^T \mathbf{b}.$$

That is,  $A\hat{\mathbf{w}}$  is the projection of  $\mathbf{b}$  onto the span of those vectors  $\mathbf{v}_i$  for which  $D_{i,i} \neq 0$ . Since the linear span of  $\mathbf{x}_1, \dots, \mathbf{x}_m$  is the same as the linear span of those  $\mathbf{v}_i$ , and  $\mathbf{b}$  is in the linear span of the  $\mathbf{x}_i$ , we obtain that  $A\hat{\mathbf{w}} = \mathbf{b}$ , which concludes our argument.  $\square$

**Remark 2.1.** Indeed we always use the **Gradient Descent** method to optimize the loss function.

Linear regression for polynomial regression tasks  $\mathcal{H}_{poly}^n = \{x \mapsto p(x)\}$  where  $\psi(x) = (1, x, x^2, \dots, x^n)$  and  $p(\psi(x)) = a_0 + a_1x + a_2x^2 + \dots + a_nx^n$ .

### 2.2 Ridge Regression

To ameliorate the effect of the invertible matrix, we could introduce the regularization.

**Definition 2.1** (Ridge Regularized Loss).

$$R(w) = \lambda \|w\|^2.$$

Now the loss function reads:

$$\mathcal{L} = \mathcal{L}_S(w) + R(w) = \frac{1}{m} \sum_{i=1}^m (h(\mathbf{x}_i) - \mathbf{y}_i)^2 + \lambda \|w\|^2.$$

Hence, the solution to ridge regression becomes

$$\mathbf{w} = (2\lambda m I + A)^{-1}.$$

**Theorem 2.3** (The stability of regularization). Let  $\mathcal{D}$  be a distribution over  $\mathcal{X} \times [-1, 1]$ , where  $\mathcal{X} = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| \leq 1\}$ . Let  $\mathcal{H} = \{\mathbf{w} \in \mathbb{R}^d : \|\mathbf{w}\| \leq B\}$ . For any  $\varepsilon \in (0, 1)$ , let  $m \geq 150B^2/\varepsilon^2$ . Then applying the ridge regression algorithm with parameter  $\lambda = \varepsilon/3B^2$  satisfies

$$\mathbb{E}_{S \sim \mathcal{D}^m} [L_D(A(S))] \leq \min_{\mathbf{w} \in \mathcal{H}} L(D) + \varepsilon.$$

## 2.3 Lasso Regression

**Definition 2.2** (Lasso Regularized Loss).

$$R(w) = \lambda \|w\|_1^2.$$

Under some assumptions on the distribution and the regularization parameter  $\lambda$ , the LASSO will find sparse solutions

## 2.4 Logistic Regression

The hypothesis class is:

$$H_{sig} = \{x \mapsto \text{sigmoid}(\mathbf{w}\mathbf{x}) : \mathbf{w} \in \mathbb{R}^d\}$$

where  $\text{sigmoid}(s) = 1/[1 + \exp(-s)]$ . The loss function is

$$\mathcal{L} = \frac{1}{m} \sum_{i=1}^m \log [1 + \exp(-y_i \mathbf{w}\mathbf{x}_i)].$$

**Remark 2.2.** Optimization in logistic regression

- The advantage of the logistic loss function is that it is a convex function with respect to  $\mathbf{w}$ .
- No close form solution.
- Identical to the problem of finding a Maximum Likelihood Estimator.

### 3 Generalized Linear Models

#### 3.1 The Exponential Family

**Definition 3.1.** We say that a class of distributions is in the exponential family if it can be written in the form

$$p(y; \eta) = b(y) \exp(\eta^T T(y) - a(\eta)).$$

Here,  $\eta$  is called the **natural parameter** (also called the canonical parameter) of the distribution;  $T(y)$  is the **sufficient statistic** (for the distributions we consider, it will often be the case that  $T(y) = y$ ); and  $a(\eta)$  is the log **partition function**. The quantity  $e^{-a(\eta)}$  essentially plays the role of a normalization constant, that makes sure the distribution  $p(y; \eta)$  sums/integrates over  $y$  to 1.

#### 3.2 Constructing GLMs

1.  $y | x; \theta \sim \text{ExponentialFamily}(\eta)$ . I.e., given  $x$  and  $\theta$ , the distribution of  $y$  follows some exponential family distribution, with parameter  $\eta$ .
2. Given  $x$ , our goal is to predict the expected value of  $T(y)$  given  $x$ . In most of our examples, we will have  $T(y) = y$ , so this means we would like the prediction  $h(x)$  output by our learned hypothesis  $h$  to satisfy  $h(x) = \mathbb{E}[y|x]$ . (Note that this assumption is satisfied in the choices for  $h_\theta(x)$  for both logistic regression and linear regression. For instance, in logistic regression, we had  $h_\theta(x) = p(y = 1|x; \theta) = 0 \cdot p(y = 0|x; \theta) + 1 \cdot p(y = 1|x; \theta) = E[y|x; \theta]$ .)
3. The natural parameter  $\eta$  and the inputs  $x$  are related linearly:  $\eta = \theta^T x$ . (Or, if  $\eta$  is vector-valued, then  $\eta_i = \theta_i^T x$ .)

**Example 3.1 (Logistic Regression).** Note that:  $y|x; \theta \sim \text{Bernoulli}(\phi)$ . Then we have  $\mathbb{E}[y|x; \theta] = \phi$ .

Thus

$$h_\theta(x) = \mathbb{E}[y|x; \theta] = \phi = \frac{1}{1 + e^{-\eta}} = \frac{1}{1 + e^{-\theta^T x}}.$$

If we have a training set of  $n$  examples  $\{(x^i, y^i); i = 1, \dots, n\}$  and would like to learn the parameters  $\theta_i$  of this model, we would begin by writing down the log-likelihood

$$\mathcal{L}(\theta) = \sum_{i=1}^n \log p(y^i | x^i; \theta) = \sum_{i=1}^n \log \left[ \left( \frac{1}{1 + e^{-\theta^T x}} \right)^{1\{y^i=1\}} \left( \frac{e^{-\theta^T x}}{1 + e^{-\theta^T x}} \right)^{1\{y^i=0\}} \right].$$

## 4 Kernel Method

Now we will introduce a function  $\phi(x) : \mathbb{R}^d \mapsto \mathbb{R}^p$  mapping the attributes to the features.

### 4.1 LMS with Features

Suppose that  $\theta = \sum_{i=1}^n \beta_i x^i$ . By updating rules of gradient descent,

$$\begin{aligned}\theta &:= \theta + \alpha \sum_{i=1}^n [y^i - \theta^T \phi(x^i)] \phi(x^i) \\ &= \sum_{i=1}^n \underbrace{\{\beta_i + \alpha [y^i - \theta^T \phi(x^i)]\}}_{\text{new } \beta} \phi(x^i)\end{aligned}$$

Then  $i \in \{1, \dots, n\}$ :

$$\beta_i := \beta_i + \alpha \left[ y^i - \sum_{j=1}^n \beta_j \phi(x^j)^T \phi(x^i) \right] = \beta_i + \alpha \left[ y^i - \sum_{j=1}^n \beta_j K(\phi(x^j), \phi(x^i)) \right]$$

where

$$K(x, z) \triangleq \langle \phi(x), \phi(z) \rangle.$$

**Remark 4.1.** Kernel is a corresponding to the feature map  $\phi$  as a function that maps  $\mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$ .

### 4.2 Properties of Kernels

**Definition 4.1** (Gaussian kernel).

$$K(x, z) = \exp\left(-\frac{\|x - z\|^2}{2\sigma^2}\right).$$

The gaussian kernel is corresponding to an **infinite** dimensional feature mapping  $\phi$ . Also,  $\phi$  lives in Hilbert space.

**Theorem 4.1.** The corresponding kernel matrix  $K \in \mathbb{R}^{n \times n}$  is symmetric positive semidefinite.

**Theorem 4.2** (Mercer Theorem). Let  $K : \mathbb{R}^d \times \mathbb{R}^d \mapsto \mathbb{R}$  be given. Then for  $K$  to be a valid Mercer Kernel, it is necessary and sufficient that for any  $\{x^1, \dots, x^n\}$ , ( $n < \infty$ ), the corresponding kernel matrix is symmetric positive semidefinite. Nota Bene: the generalized form involve  $L^2$  functions.



## 5 Support Vector Machines

SVMs are among the best (and many believe are indeed the best) off-the-shelf supervised learning algorithms. So, be **self-motivated** in this section.

### 5.1 Hard-SVM

**Hard-SVM** is the learning rule in which we return an ERM hyperplane that separates the training set with the largest possible margin. The Hard-SVM rule is

$$\arg \max_{(w,b): \|w\|=1} \min_{i \in [m]} |w^T x^i + b| \quad \text{s.t. } \forall i, y^i (w^T x^i + b) \geq 1.$$

Equivalently,

$$\arg \max_{(w,b): \|w\|=1} \min_{i \in [m]} y^i (w^T x^i + b) \quad (5.1)$$

Next, we give another equivalent formulation of the Hard-SVM rule as a quadratic optimization problem.<sup>1</sup>

Input:  $(x^1, y^1), \dots, (x^m, y^m)$

Solve:

$$(w_0, b_0) = \arg \min_{(w,b)} \frac{1}{2} \|w\|^2 \quad \text{s.t. } \forall i, y^i (w^T x^i + b) \geq 1. \quad (5.2)$$

Output:  $\hat{w} = w_0 / \|w_0\|, \hat{b} = b_0 / \|w_0\|$

**Lemma 5.1.** The output of Hard-SVM is a solution of Equation (5.1).

**Proof.** Let  $(w_1, b_1)$  be a solution of Equation (5.1) and  $\gamma_1 = \min_{i \in [m]} y_i (w_1^T x^i + b_1)$ . Then we have

$$y^i \left( \frac{w_1^T}{\gamma_1} x^i + \frac{b_1}{\gamma_1} \right) \geq 1.$$

Hence  $\|w_0\| \leq \|w_1 / \gamma_1\| = 1 / \gamma_1$ . It follows that for all  $i$ ,

$$y^i (\hat{w}^T x^i + \hat{b}) \geq \frac{1}{\|w_0\|} \geq \gamma_1.$$

Since  $\|\hat{w}\| = 1$  we obtain that  $(\hat{w}, \hat{b})$  is an optimal solution of Equation (5.1).  $\square$

#### 5.1.1 The Sample Complexity of Hard-SVM\*

**Definition 5.1 (Separability).** Let  $\mathcal{D}$  be a distribution over  $\mathbb{R}^d \times \{\pm 1\}$ . We say that  $\mathcal{D}$  is separable with a  $(\gamma, \rho)$ -margin if there exists  $(w^*, b^*)$  such that  $\|w^*\| = 1$  and such that with probability 1 over the choice of  $(x, y) \sim \mathcal{D}$  we have that  $y(w^{*T} x + b^*) \geq \gamma$  and  $\|x\| \leq \rho$ .

<sup>1</sup>A quadratic optimization problem is an optimization problem in which the objective is a convex quadratic function and the constraints are linear inequalities.

**Theorem 5.1.** Let  $\mathcal{D}$  be a distribution over  $\mathbb{R}^d \times \{\pm 1\}$  that satisfies the  $(\gamma, \rho)$ -separability with margin assumption using a homogenous halfspace. Then, with probability of at least  $1 - \delta$  over the choice of a training set of size  $m$ , the 0-1 error of the output of Hard-SVM is at most

$$\sqrt{\frac{4(\rho)/\gamma^2}{m}} + \sqrt{\frac{2 \log(2/\delta)}{m}}.$$

## 5.2 Soft-SVM and Norm Regularization

Input:  $(x^1, 1), \dots, (x^m, y^m)$

Parameter:  $\lambda > 0$

Solve:

$$\begin{aligned} \min_{w, b, \xi} & \left( \lambda \|w\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i \right) \\ \text{s.t. } & \forall i, y^i (w^T x^i + b) \geq 1 - \xi_i \text{ and } \xi_i \geq 0 \end{aligned}$$

Output:  $w, b$

**Definition 5.2** (hinge loss).

$$l^{\text{hinge}}((w, b), (x, y)) = \max \{0, 1 - yw^T x + b\}.$$

Now we just need to optimize  $\lambda \|w\|^2 + \mathcal{L}^{\text{hinge}}(w, b)$ .

## 5.3 Duality

The Lagrangian for EQ.5.2 is:

$$\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i \left[ y^{(i)} (w^T x^{(i)} + b) - 1 \right].$$

By the fact that  $\nabla \mathcal{L} = 0$

$$w = \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)} \quad \text{and} \quad \sum_{i=1}^n \alpha_i y^{(i)} = 0.$$

Plug them back we obtain

$$\begin{aligned} \max_{\alpha} \mathcal{L}(\alpha) &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle \\ \text{s.t. } & \alpha_i \geq 0 \\ & \sum_{i=1}^n \alpha_i y^{(i)} = 0 \end{aligned}$$

Note that

$$b = -\frac{\max_{y^{(i)}=-1} w^T x^{(i)} + \min_{y^{(i)}=1} w^T x^{(i)}}{2}.$$

## 5.4 SMO

Consider the dual form of Soft-SVM

$$\begin{aligned} \max_{\alpha} \quad & \mathcal{L}(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C \\ & \sum_{i=1}^n \alpha_i y^{(i)} = 0 \end{aligned}$$

The dual-complementarity conditions are

$$\begin{cases} \alpha_i = 0 & \implies y^{(i)}(w^T x^{(i)} + b) \geq 1 \\ \alpha_i = C & \implies y^{(i)}(w^T x^{(i)} + b) \leq 1 \\ 0 < \alpha_i < C & \implies y^{(i)}(w^T x^{(i)} + b) = 1 \end{cases}$$

Now we introduce SMO (sequential minimal optimization). Repeat till convergence:

1. Select some pair  $\alpha_i$  and  $\alpha_j$  to update next (using a heuristic that tries to pick the two that will allow us to make the biggest progress towards the global maximum).
2. Reoptimize  $W(\alpha)$  with respect to  $i$  and  $\alpha_j$ , while holding all the other  $\alpha_k$ 's ( $k \neq i, j$ ) fixed.

## 5.5 Implementing Soft-SVM Using SGD

---

**Algorithm 1** SGD for Solving Soft-SVM

---

```

 $\theta = 0$ 
for  $t = 1, \dots, T$  do
   $w^{(t)} = 1/\lambda t \times \theta$ 
  Choose  $i$  uniformly at random for  $[m]$ 
  if  $y_i w^T x^i < 1$  then
     $\theta \leftarrow \theta + y^i x^i$ 
  end if
end for
return  $\sum_{t=1}^T w^{(t)} / T$ 

```

---

## 6 Clustering

K-means and GMM algorithms are easy. Here I will not introduce them. These PDFs may be helpful:

- <https://cs229.stanford.edu/summer2020/cs229-notes7a.pdf>
- <https://cs229.stanford.edu/summer2020/cs229-notes7b.pdf>

## 7 Decision Tree

A general framework for growing a decision tree is as follows. We start with a tree with a single leaf (the root) and assign this leaf a label according to a majority vote among all labels over the training set. We now perform a series of iterations. On each iteration, we examine the effect of splitting a single leaf. We define some gain measure that quantifies the improvement due to this split. Then, among all possible splits, we either choose the one that maximizes the gain and perform it, or choose not to split the leaf at all.

---

### Algorithm 2 Iterative Dichotomizer 3

---

**Input:** training set  $S$ , feature subset  $A \subseteq [d]$

- 1: **if** all examples in  $S$  are labeled by 1 **then**
- 2:     **return** a leaf 1
- 3: **else if** all examples in  $S$  are labeled by 0 **then**
- 4:     **return** a leaf 0
- 5: **else if**  $A = \emptyset$  **then**
- 6:     **return** leaf whose value is the majority labels in  $S$
- 7: **end if**
- 8:  $j \leftarrow \arg \max_{i \in A} \text{Gain}(S, i)$
- 9:  $T_1 \leftarrow \text{ID3}(\{(x, y) \in S : x_j = 1\}, A/\{j\})$
- 10:  $T_2 \leftarrow \text{ID3}(\{(x, y) \in S : x_j = 0\}, A/\{j\})$
- 11: **return** a root  $r$  whose left subtree is  $T_1$  and the right tree is  $T_2$

---

**Definition 7.1 (Entropy).** The surprise of observing a discrete random variable  $Y$  takes on value  $k$  is  $-\log(\Pr(Y = k))$ . Then the entropy of  $Y$  is the expected surprise:

$$H(Y) = - \sum_k \Pr(Y = k) \log \Pr(Y = k).$$

When we choose a split-feature, we want to reduce entropy in some way. Thus we want to minimize the **conditional entropy**  $H(Y|X_j)$ :

$$\min H(Y|X_j) \triangleq \Pr(X_j = 1)H(Y|X_j = 1) + \Pr(X_j = 0)H(Y|X_j = 0)$$

which is equivalent to

$$\max I(X_j; Y) \triangleq H(Y) - H(Y|X_j).$$

The quantity  $I(X_j; Y)$  is known as the mutual information between  $X_j$  and  $Y$ .

**Definition 7.2 (Gini impurity/index).**

$$G(Y) = \sum_k \Pr(Y = k) \sum_{j \neq k} \Pr(Y = j) = 1 - \sum_k \Pr(Y = k)^2.$$

Similarly, we minimize the quantity  $G(Y|X_j) \triangleq \Pr(X_j = 1)G(Y|X_j = 1) + \Pr(X_j = 0)G(Y|X_j = 0)$ .

## 7.1 Random Forest

Random forests are a specific **ensemble method** where the individual models are decision trees trained in a randomized way so as to **reduce correlation** among them. Because the basic decision tree building algorithm is deterministic, it will produce the same tree every time if we give it the same dataset and use the same algorithm hyperparameters (stopping conditions, etc.).

Random forests are typically randomized in the following ways:

- Per-classifier **bagging** (short for **bootstrap aggregating**): sample some number  $m \leq n$  of datapoints uniformly with replacement, and use these as the training set.
- Per-split **feature randomization**: sample some number  $k < d$  of features as candidates to be considered for this split.

## 8 Dimensionality Reduction

Dimensionality reduction is the process of taking data in a high dimensional space and mapping it into a new space whose dimensionality is much smaller.

### 8.1 Principal Component Analysis

In PCA, we have a compressing matrix  $W \in \mathbb{R}^{n,d}$  and a recovering matrix  $U \in \mathbb{R}^{d,n}$ . For given data  $x_1, x_2, \dots, x_m$ , we aim at solving the problem:

$$\arg \min_{W,U} \sum_{i=1}^n \|x_i - UWx_i\|^2 \quad (8.1)$$

**Lemma 8.1.** Let  $(U, W)$  be a solution of Equation 8.1. Then  $U^T U = I$  and  $W = U^T$ . (The columns of  $U$  are orthonormal.)

**Proof.** Let  $R = \{UWx : x \in \mathbb{R}^d\}$  which is an  $n$  dimensional linear subspace of  $\mathbb{R}^d$ . Let  $V \in \mathbb{R}^{n,d}$  be a matrix satisfies the range of  $V$  is  $R$  and  $V^T V = I$ . Then  $\|x - Vy\|^2 = \|x\|^2 + \|y\|^2 - 2y^T V^T x$ . Minimizing this w.r.t.  $y$  gives that  $y = V^T x$ .  $\square$

By the fact that

$$\|x - UU^T x\|^2 = \|x\|^2 - \text{trace}(U^T x x^T U).$$

We could rewrite Equation 8.1 as follows:

$$\arg \max_{U \in \mathbb{R}^{d,n}: U^T U = I} \text{trace} \left[ U^T \left( \sum_{i=1}^m x_i x_i^T \right) U \right].$$

**Theorem 8.1.** Let  $x_1, \dots, x_m$  be arbitrary vectors in  $\mathbb{R}^d$ , let  $A = \sum_{i=1}^m x_i x_i^T$ , and let  $u_1, \dots, u_n$  be  $n$  eigenvectors of the matrix  $A$  corresponding to the largest  $n$  eigenvalues of  $A$ . Then, the solution to the PCA optimization problem given in Equation 8.1 is to set  $U$  to be the matrix whose columns are  $u_1, \dots, u_n$  and to set  $W = U^T$ .

**Proof.** Let  $VDV^T$  be the spectral decomposition of  $A$  (suppose that  $D_{1,1} \geq \dots \geq D_{d,d}$ ) and let  $B = V^T U$ . We have

$$\text{trace}(U^T A U) = \text{trace}(B^T D B) = \sum_{j=1}^d D_{j,j} \sum_{i=1}^n B_{j,i}^2 \leq \max_{\beta \in [0,1]^d: \|\beta\| \leq n} \sum_{j=1}^d D_{j,j} \beta_j = \sum_{j=1}^n D_{j,j}.$$

Nota Bene:  $B^T B = I$  which entails  $\sum_{j=1}^d \sum_{i=1}^n B_{j,i}^2 = n$ .  $\square$

## 8.2 Implementation

---

**Algorithm 3** PCA algorithm

---

**Input:** A matrix of  $m$  examples  $X \in R^{m,d}$  and number of components  $n$ .

```
1: if  $m > d$  then  
2:    $A = X^T X$   
3:   Let  $u_1, \dots, u_n$  be the eigenvectors of  $A$  with largest eigenvalues  
4: else  
5:    $B = X X^T$   
6:   Let  $v_1, \dots, v_n$  be the eigenvectors of  $B$  with largest eigenvalues  
7:    $\forall i, u_i = X^T v_i / \|X^T v_i\|$   
8: end if  
9: return  $u_1, \dots, u_n$ 
```

---

The algorithm use a more efficient method when  $d > m$ . The complexity is  $\mathcal{O}(m^2 d)$  under this case.