

Part 1.2 Linear Algebra Applications

Scriber: Yingyu Liang

1 Linear Regression

Suppose we are given a set of data points $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ where $\mathbf{x}_i \in \mathbb{R}^d$ and $y \in \mathbb{R}$. We would like to learn a linear regression model to fit the data:

$$y = \boldsymbol{\theta}_x^\top \mathbf{x} + \theta_0$$

where $\boldsymbol{\theta}_x \in \mathbb{R}^d$, $\theta \in \mathbb{R}$.

First, we conveniently rewrite the problem. Let $\bar{\mathbf{x}}_i \in \mathbb{R}^{d+1}$ be the concatenation of \mathbf{x}_i and 1, and $\boldsymbol{\theta} \in \mathbb{R}^{d+1}$ be the concatenation of $\boldsymbol{\theta}_x$ and θ_0 . Then the model becomes

$$y = \boldsymbol{\theta}^\top \bar{\mathbf{x}}.$$

Next, we consider learning a model by optimizing the mean square error (MSE):

$$\hat{\boldsymbol{\theta}} = \operatorname{argmin}_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^n (y_i - \boldsymbol{\theta}^\top \bar{\mathbf{x}}_i)^2.$$

Let $\mathbf{X} \in \mathbb{R}^{n \times (d+1)}$ be a matrix with $\bar{\mathbf{x}}_i$ as the i -th row, and let $\mathbf{y} \in \mathbb{R}^n$ be a vector with y_i as the i -th entry. Then the MSE objective can be rewritten in the matrix form:

$$\hat{\boldsymbol{\theta}} = \operatorname{argmin}_{\boldsymbol{\theta}} \frac{1}{n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2 = \operatorname{argmin}_{\boldsymbol{\theta}} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})$$

where $\|\mathbf{v}\|^2 = \mathbf{v}^\top \mathbf{v}$ denotes the Euclidean norm.

Now we do the optimization. Rewrite the objective as:

$$(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) = \mathbf{y}^\top \mathbf{y} - 2\boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{y} + \boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\theta}.$$

Take the gradient with respect to $\boldsymbol{\theta}$ and set it to the zero vector:

$$\begin{aligned} -2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X}\boldsymbol{\theta} &= 0. \\ \mathbf{X}^\top \mathbf{X}\boldsymbol{\theta} &= \mathbf{X}^\top \mathbf{y}. \end{aligned}$$

Case 1. Assuming $\mathbf{X}^\top \mathbf{X}$ is invertible, we obtain the solution

$$\hat{\boldsymbol{\theta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

Case 2. If $\mathbf{X}^\top \mathbf{X}$ is not invertible, i.e., $\text{rank}(\mathbf{X}^\top \mathbf{X}) < d$, we can use its pseudo-inverse. Suppose $\text{rank}(\mathbf{X}) = r < d$. Consider the (truncated) SVD $\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top$, where $\Sigma \in \mathbb{R}^{r \times r}$ is the matrix removing the 0 singular values, and $\mathbf{U} \in \mathbb{R}^{n \times r}, \mathbf{V} \in \mathbb{R}^{d \times r}$ are those matrices removing the columns associated with 0 singular values. Then the pseudo-inverse $\mathbf{X}^\dagger = \mathbf{U}\Sigma^{-1}\mathbf{V}^\top$. And $\mathbf{X}^\top \mathbf{X} = \mathbf{V}\Sigma^2\mathbf{V}^\top$ and the pseudo-inverse $(\mathbf{X}^\top \mathbf{X})^\dagger = \mathbf{V}(\Sigma^2)^{-1}\mathbf{V}^\top$. Then we can obtain

$$\hat{\boldsymbol{\theta}} = (\mathbf{X}^\top \mathbf{X})^\dagger \mathbf{X}^\top \mathbf{y}.$$

It turns out $(\mathbf{X}^\top \mathbf{X})^\dagger \mathbf{X}^\top = \mathbf{X}^\dagger$, since

$$\begin{aligned} (\mathbf{X}^\top \mathbf{X})^\dagger \mathbf{X}^\top &= \mathbf{V}(\Sigma^2)^{-1}\mathbf{V}^\top \mathbf{X}^\top \\ &= \mathbf{V}(\Sigma^2)^{-1}\mathbf{V}^\top \mathbf{V}\Sigma\mathbf{U}^\top \\ &= \mathbf{V}(\Sigma^2)^{-1}\mathbf{V}^\top \mathbf{V}\Sigma\mathbf{U}^\top \\ &= \mathbf{V}\Sigma^{-1}\mathbf{U}^\top \\ &= \mathbf{X}^\dagger. \end{aligned}$$

So when $\mathbf{X}^\top \mathbf{X}$ is not invertible, we can use

$$\hat{\boldsymbol{\theta}} = \mathbf{X}^\dagger \mathbf{y}.$$

Ridge Regression. However, in practice, we typically consider the regularized version of linear regression (for statistical reasons we will learn later):

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \| \mathbf{y} - \mathbf{X}\boldsymbol{\theta} \|^2 + \lambda \|\boldsymbol{\theta}\|^2$$

where $\lambda > 0$ is the regularization coefficient. It turns out that this also helps mitigate the invertibility issue.

Taking the derivative again and setting it to 0:

$$\begin{aligned} -2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X}\boldsymbol{\theta} + 2\lambda\boldsymbol{\theta} &= 0 \\ (\mathbf{X}^\top \mathbf{X} + \lambda\mathbf{I})\boldsymbol{\theta} &= \mathbf{X}^\top \mathbf{y} \\ \hat{\boldsymbol{\theta}} &= (\mathbf{X}^\top \mathbf{X} + \lambda\mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}. \end{aligned}$$

Note that $\mathbf{X}^\top \mathbf{X} + \lambda\mathbf{I}$ is always invertible. To see this, consider the (full) SVD $\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top$, where $\mathbf{U} \in \mathbb{R}^{n \times n}, \Sigma \in \mathbb{R}^{n \times d}, \mathbf{V} \in \mathbb{R}^{d \times d}$. Then

$$\begin{aligned} \mathbf{X}^\top \mathbf{X} + \lambda\mathbf{I} &= (\mathbf{U}\Sigma\mathbf{V}^\top)^\top (\mathbf{U}\Sigma\mathbf{V}^\top) + \lambda\mathbf{I} \\ &= \mathbf{V}\Sigma^\top \Sigma\mathbf{V}^\top + \lambda\mathbf{I} \\ &= \mathbf{V}(\Sigma^\top \Sigma + \lambda\mathbf{I})\mathbf{V}^\top. \end{aligned}$$

Here $\Sigma^\top \Sigma + \lambda\mathbf{I}$ is a $d \times d$ diagonal matrix with diagonal entries $\sigma_i^2 + \lambda > 0$, so it is invertible.

We can also see the form of the solution:

$$\begin{aligned}\hat{\theta} &= (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y} \\ &= \mathbf{V}(\Sigma^\top \Sigma + \lambda \mathbf{I})^{-1} \mathbf{V}^\top (\mathbf{U} \Sigma \mathbf{V}^\top)^\top \mathbf{y} \\ &= \mathbf{V}(\Sigma^\top \Sigma + \lambda \mathbf{I})^{-1} \Sigma^\top \mathbf{U}^\top \mathbf{y}.\end{aligned}$$

Here $(\Sigma^\top \Sigma + \lambda \mathbf{I})^{-1} \Sigma^\top \in \mathbb{R}^{d \times n}$ has diagonal entries $\frac{\sigma_i}{\sigma_i^2 + \lambda}$.

We can also use gradient descent or its variants to solve the optimization, which will be discussed in the later lectures.

2 Principal Component Analysis

Suppose we are given a set of data points $\{\mathbf{x}_i\}_{i=1}^n$ where $\mathbf{x}_i \in \mathbb{R}^d$. We would like to learn a lower-dimensional subspace (of dimension r) and represent the data points as the projections in the subspace (i.e., a point in \mathbb{R}^r). For convenience, we shall *assume the data points are centered* for the rest of this section. That is,

- **(Assumption)** The data points have zero mean: $\mu := \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i = 0$.

If the mean $\mu \neq 0$, then we can centralize the data by subtracting the mean: $\mathbf{x}_i \leftarrow \mathbf{x}_i - \mu$.

Algorithm 1 Principal Component Analysis (PCA)

Require: $\{\mathbf{x}_i\}_{i=1}^n \subseteq \mathbb{R}^d$ such that $\sum_{i=1}^n \mathbf{x}_i = 0$, number of principal components $r \leq d$

Form the sample covariance matrix $\mathbf{S} := \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$

Do eigendecomposition: $\mathbf{S} = \mathbf{U} \Lambda \mathbf{U}^\top$, where Λ is a diagonal matrix with eigenvalues on the diagonal, and the columns of $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d]$ are the corresponding eigenvectors.

Suppose the eigenvalues are sorted in decreasing order: $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$.

Form $\mathbf{U}_r = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r]$, and let $\tilde{\mathbf{x}}_i = \mathbf{U}_r^\top \mathbf{x}_i = [\mathbf{u}_1^\top \mathbf{x}_i, \mathbf{u}_2^\top \mathbf{x}_i, \dots, \mathbf{u}_r^\top \mathbf{x}_i]^\top$

Ensure: the principal components \mathbf{U}_r , and the new representations $\{\tilde{\mathbf{x}}_i\}_{i=1}^n$

2.1 The Variance Preservation View

PCA can be justified in several ways. Here we consider the goal of finding a subspace to preserve the variance as much as possible.

Consider $r = 1$. Our goal is to find a 1-dimensional subspace (a line going through the origin) such that the projections of the data points on this line have as large a variance as possible, in the hope that this maximally preserves the distinction among the points.

A line going through the origin can be represented by a vector $\mathbf{w} \in \mathbb{R}^d$. The projection of a point \mathbf{x} onto this line is then:

$$\mathbf{p} = \frac{\mathbf{w}^\top \mathbf{x}}{\|\mathbf{w}\|}$$

where $\|\mathbf{w}\| = \sqrt{\mathbf{w}^\top \mathbf{w}}$ is the Euclidean norm.

For simplicity, let us consider \mathbf{w} with unit length. Then the projections of the data points are $\mathbf{w}^\top \mathbf{x}_i$. Since their mean is $\frac{1}{n} \sum_{i=1}^n \mathbf{w}^\top \mathbf{x}_i = \frac{1}{n} \mathbf{w}^\top (\sum_{i=1}^n \mathbf{x}_i) = 0$, their variance is

$$\frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i)^2 = \frac{1}{n} \sum_{i=1}^n \mathbf{w}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{w} = \mathbf{w}^\top \mathbf{S} \mathbf{w}$$

where $\mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$ is the sample covariance.

Recall that our goal is to maximize the variance of the projections:

$$\begin{aligned} & \max_{\mathbf{w}} \mathbf{w}^\top \mathbf{S} \mathbf{w} \\ & \text{s.t. } \|\mathbf{w}\| = 1. \end{aligned}$$

This is exactly finding the first eigenvector of the sample covariance \mathbf{S} (recall the proof of the eigendecomposition). We can also show it by solving the optimization via the Lagrangian:

$$\mathbf{w}^\top \mathbf{S} \mathbf{w} + \lambda(1 - \mathbf{w}^\top \mathbf{w}).$$

The gradient w.r.t. \mathbf{w} is

$$2\mathbf{S} \mathbf{w} - 2\lambda \mathbf{w}.$$

The optimal solution satisfies that the gradient equals 0, so

$$\mathbf{S} \mathbf{w} = \lambda \mathbf{w}$$

showing that it is the first eigenvector of the sample covariance \mathbf{S} . Furthermore, the projected variance for this optimal solution is

$$\mathbf{w}^\top \mathbf{S} \mathbf{w} = \mathbf{w}^\top \lambda \mathbf{w} = \lambda$$

showing that the projected variance is just the first eigenvalue of the sample covariance \mathbf{S} .

Consider $r > 1$. After finding the first principal component \mathbf{u}_1 , we can recursively find the second principal component \mathbf{w} : find the direction that is orthogonal to the first principal component and maximizes the projected variance. Again, it can be shown that the second principal component is just the second eigenvector of the sample covariance, and the projected variance is just the second eigenvalue. Similarly for the other principal components.

In summary, PCA to dimension r is just finding the top r eigenvectors of the sample covariance, and projecting the original points to this subspace.

2.2 The Minimum Reconstruction Error View

Using any orthonormal basis $\mathbf{w}_1, \dots, \mathbf{w}_d$, a point \mathbf{x}_i can be written as

$$\mathbf{x}_i = \sum_{j=1}^d \alpha_{ij} \mathbf{w}_j = \sum_{j=1}^d (\mathbf{w}_j^\top \mathbf{x}_i) \mathbf{w}_j = \sum_{j=1}^d \mathbf{w}_j \mathbf{w}_j^\top \mathbf{x}_i$$

where

$$\alpha_{ij} = \mathbf{w}_j^\top \mathbf{x}_i.$$

Consider the r -term approximation to \mathbf{x}_i :

$$\hat{\mathbf{x}}_i = \sum_{j=1}^r \alpha_{ij} \mathbf{w}_j.$$

We want the approximation error to be small for all training points:

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \|\hat{\mathbf{x}}_i - \mathbf{x}_i\|^2 &= \frac{1}{n} \sum_{i=1}^n \left\| \sum_{j=r+1}^d \alpha_{ij} \mathbf{w}_j \right\|^2 \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{j=r+1}^d \alpha_{ij}^2 \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{j=r+1}^d \mathbf{w}_j^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{w}_j \\ &= \sum_{j=r+1}^d \mathbf{w}_j^\top \mathbf{S} \mathbf{w}_j. \end{aligned}$$

If $r = d - 1$, i.e., we need to remove a single dimension, it is easy to see that \mathbf{w}_d should be the least eigenvector \mathbf{u}_d , because $\mathbf{u}_d^\top \mathbf{S} \mathbf{u}_d$ is the smallest among all unit vectors. Similarly, the other dimensions to remove are subsequently the eigenvectors corresponding to the least eigenvalues. In other words, we should use the top r eigenvectors for the r -term approximation.

2.3 Equivalence between the Two Views

The optimization objectives of the two views are equivalent, since the r -term approximation $\hat{\mathbf{x}}_i$ is just the coordinates of the projection in the original coordinate system.

Consider a subspace spanned by r orthogonal basis vectors $\mathbf{w}_1, \dots, \mathbf{w}_r$. The projection of \mathbf{x}_i into this subspace is

$$\mathbf{p}_i = [\alpha_{i1}, \dots, \alpha_{id}]^\top = [\mathbf{w}_1^\top \mathbf{x}_i, \dots, \mathbf{w}_r^\top \mathbf{x}_i]^\top.$$

Its contribution to the projected covariance is $\mathbf{p}_i^\top \mathbf{p}_i$.

Geometrically, $\mathbf{p}_i^\top \mathbf{p}_i = \|\mathbf{p}_i\|^2$ is just the square of the length of the projection and equals $\|\hat{\mathbf{x}}_i\|^2$. On the other hand, $\|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2$ is the distance between \mathbf{x}_i and its projection, which is its distance to the subspace spanned by $\mathbf{w}_1, \dots, \mathbf{w}_r$. Clearly, $\|\hat{\mathbf{x}}_i\|^2 + \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2 = \|\mathbf{x}_i\|^2$ (recall the Pythagorean Theorem!). This is a constant. Therefore, maximizing the projected variance $\|\mathbf{p}_i\|^2 = \|\hat{\mathbf{x}}_i\|^2$ is equivalent to minimizing the reconstruction error $\|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2$.

Why $\|\mathbf{p}_i\|^2 = \|\widehat{\mathbf{x}}_i\|^2$? This is because \mathbf{p}_i is the coordinate vector in the basis $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_r]$, and back to the original coordinate system of \mathbf{x}_i , the coordinate vector is just:

$$\mathbf{W}\mathbf{p}_i = \sum_{j=1}^r \alpha_{ij} \mathbf{w}_j = \widehat{\mathbf{x}}_i.$$

Therefore, their square lengths are the same:

$$\|\widehat{\mathbf{x}}_i\|^2 = \|\mathbf{W}\mathbf{p}_i\|^2 = (\mathbf{W}\mathbf{p}_i)^\top \mathbf{W}\mathbf{p}_i = \mathbf{p}_i^\top \mathbf{W}^\top \mathbf{W}\mathbf{p}_i = \mathbf{p}_i^\top \mathbf{p}_i = \|\mathbf{p}_i\|^2.$$

2.4 Solved by Singular Value Decomposition

PCA is just finding the top eigenvectors of the sample covariance \mathbf{S} . Note that

$$\mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top = \frac{1}{n} \mathbf{X}^\top \mathbf{X}$$

where $\mathbf{X} \in \mathbb{R}^{n \times d}$ is the data matrix with \mathbf{x}_i as the i -th row. By the connection between the eigendecomposition of $\mathbf{X}^\top \mathbf{X}$ and the SVD of \mathbf{X} , we know that PCA is just finding the top right singular vectors of the data matrix \mathbf{X} .