

Homework Set 2, CPSC 8420, Fall 2024

Your Name

Due 10/28/2024, Monday, 11:59PM EST

Problem 1

For Principle Component Analysis (PCA), from the perspective of maximizing variance (assume the data is already self-centered)

- show that the first column of \mathbf{U} , where $[\mathbf{U}, \mathbf{S}] = \text{svd}(\mathbf{X}^T \mathbf{X})$ will maximize $\|\mathbf{X}\phi\|_2^2$, s.t. $\|\phi\|_2 = 1$. (Note: you need prove why it is optimal than any other reasonable combinations of \mathbf{U}_i , say $\hat{\phi} = 0.8 * \mathbf{U}(:, 1) + 0.6 * \mathbf{U}(:, 2)$ which also satisfies $\|\hat{\phi}\|_2 = 1$.)
- show that the solution is not unique, say if ϕ is the optimal solution, so is $-\phi$.
- show that first r columns of \mathbf{U} , where $[\mathbf{U}, \mathbf{S}] = \text{svd}(\mathbf{X}^T \mathbf{X})$ maximize $\|\mathbf{X}\mathbf{W}\|_F^2$, s.t. $\mathbf{W}^T \mathbf{W} = \mathbf{I}_r$.
- Assume the singular values are all different in \mathbf{S} , then how many possible different \mathbf{W} 's will maximize the objective above?

Solution

Part (a)

Show that the first column of \mathbf{U} , where $[\mathbf{U}, \mathbf{S}] = \text{svd}(\mathbf{X}^T \mathbf{X})$, maximizes $\|\mathbf{X}\phi\|_2^2$ subject to $\|\phi\|_2 = 1$.

Proof:

The variance of the data when projected onto a unit vector ϕ is

$$\|\mathbf{X}\phi\|_2^2 = \phi^T \mathbf{X}^T \mathbf{X} \phi \quad \text{subject to} \quad \|\phi\|_2 = 1$$

$\mathbf{X}^T \mathbf{X}$ is a symmetric matrix, decompose it as

$$\mathbf{X}^T \mathbf{X} = \mathbf{U} \mathbf{S} \mathbf{U}^T,$$

where \mathbf{U} is an orthogonal matrix and \mathbf{S} is a diagonal matrix containing the singular values

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0.$$

Substituting into the objective function

$$\phi^T \mathbf{X}^T \mathbf{X} \phi = \phi^T \mathbf{U} \mathbf{S} \mathbf{U}^T \phi.$$

Let $\psi = \mathbf{U}^T \phi$. Since \mathbf{U} is orthogonal, $\|\psi\|_2 = \|\phi\|_2 = 1$. Thus, the expression becomes

$$\phi^T \mathbf{U} \mathbf{S} \mathbf{U}^T \phi = \psi^T \mathbf{S} \psi.$$

The quantity $\psi^T \mathbf{S} \psi$ is maximized when ψ aligns with the eigenvector corresponding to the largest singular value σ_1 . Therefore, the maximum is achieved when $\phi = \mathbf{U}(:, 1)$.

$\hat{\phi} = 0.8\mathbf{U}(:, 1) + 0.6\mathbf{U}(:, 2)$ satisfies $\|\hat{\phi}\|_2 = 1$, the contribution from σ_2 will reduce the value of $\phi^T \mathbf{X}^T \mathbf{X} \phi$ since $\sigma_1 \geq \sigma_2$. Hence, $\mathbf{U}(:, 1)$ is the optimal choice.

—

Part (b)

Show that the solution is not unique if ϕ is an optimal solution, then $-\phi$ is also optimal.

Proof:

Consider the objective function

$$\|\mathbf{X}\phi\|_2^2 = \phi^T \mathbf{X}^T \mathbf{X} \phi$$

Take $-\phi$

$$\|\mathbf{X}(-\phi)\|_2^2 = (-\phi)^T \mathbf{X}^T \mathbf{X} (-\phi) = \phi^T \mathbf{X}^T \mathbf{X} \phi.$$

The value of the objective function remains unchanged, which implies that both ϕ and $-\phi$ are optimal solutions.

—

Part (c)

Show that the first r columns of \mathbf{U} , where $[\mathbf{U}, \mathbf{S}] = \text{svd}(\mathbf{X}^T \mathbf{X})$, maximize $\|\mathbf{X}\mathbf{W}\|_F^2$, s.t. $\mathbf{W}^T \mathbf{W} = \mathbf{I}_r$.

Proof:

The objective function can be written as

$$\|\mathbf{X}\mathbf{W}\|_F^2 = \text{Tr}(\mathbf{W}^T \mathbf{X}^T \mathbf{X} \mathbf{W}).$$

By SVD $\mathbf{X}^T \mathbf{X} = \mathbf{U} \mathbf{S} \mathbf{U}^T$,

$$\text{Tr}(\mathbf{W}^T \mathbf{U} \mathbf{S} \mathbf{U}^T \mathbf{W}).$$

Let $\mathbf{V} = \mathbf{U}^T \mathbf{W}$. Since \mathbf{U} is orthogonal, $\mathbf{V}^T \mathbf{V} = \mathbf{I}_r$. The objective becomes

$$\text{Tr}(\mathbf{V}^T \mathbf{S} \mathbf{V}).$$

The trace $\text{Tr}(\mathbf{V}^T \mathbf{S} \mathbf{V})$ is maximized when \mathbf{V} aligns with the first r columns of \mathbf{U} , corresponding to the largest r singular values $\sigma_1, \dots, \sigma_r$. Therefore, $\mathbf{W} = \mathbf{U}(:, 1:r)$ is the optimal solution.

—

Part (d)

Assuming the singular values in \mathbf{S} are different, how many possible different \mathbf{W} 's will maximize the objective?

Answer:

Given that the singular values in \mathbf{S} are all distinct, the optimization problem of maximizing $\|\mathbf{X}\mathbf{W}\|_F^2$ under the constraint $\mathbf{W}^T\mathbf{W} = \mathbf{I}_r$ has solutions that can be expressed as:

$$\mathbf{W} = \mathbf{U}_r\mathbf{Q},$$

where \mathbf{U}_r consists of the first r columns of \mathbf{U} , and \mathbf{Q} is any orthogonal $r \times r$ matrix such that $\mathbf{Q}^T\mathbf{Q} = \mathbf{I}_r$. The set of all such orthogonal matrices \mathbf{Q} forms the orthogonal group $O(r)$.

Problem 2

Given matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$ (assume each column is centered already), where n denotes sample size while p feature size. To conduct PCA, we need find eigenvectors to the largest eigenvalues of $\mathbf{X}^T \mathbf{X}$, where usually the complexity is $\mathcal{O}(p^3)$. Apparently when $n \ll p$, this is not economic when p is large. Please consider conducting PCA based on $\mathbf{X} \mathbf{X}^T$ and obtain the eigenvectors for $\mathbf{X}^T \mathbf{X}$ accordingly and use experiment to demonstrate the acceleration.

Answer

To efficiently perform PCA when $n \ll p$, the eigenvectors of $\mathbf{X}^T \mathbf{X}$ can be obtained from the eigenvectors of $\mathbf{X} \mathbf{X}^T$. This is computationally cheaper when p is large and n is relatively small.

By decomposition of $\mathbf{X} \mathbf{X}^T$,

$$\mathbf{X} \mathbf{X}^T = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T,$$

The eigenvectors of $\mathbf{X}^T \mathbf{X}$ can be obtained from \mathbf{U} .

$$\mathbf{V} = \mathbf{X}^T \mathbf{U} \mathbf{\Lambda}^{-\frac{1}{2}}$$

Python Code

```
import numpy as np
import time

# Generate a random matrix X with n << p
n, p = 100, 1000
np.random.seed(0)
X = np.random.randn(n, p)

# Method 1: Compute eigenvectors of X^T X directly
start_time = time.time()
XtX = np.dot(X.T, X)
_, V1 = np.linalg.eigh(XtX)
time_direct = time.time() - start_time

# Method 2: Compute eigenvectors of X X^T and transform
start_time = time.time()
XXt = np.dot(X, X.T)
D, U = np.linalg.eigh(XXt)
V2 = np.dot(X.T, U) * (1 / np.sqrt(D)) # Normalize eigenvectors
time_indirect = time.time() - start_time

# Display the computational times
print("Time using direct method (X^T X):", time_direct, "seconds")
print("Time using indirect method (X X^T):", time_indirect, "seconds")
```

Experimental Results

From the experiment, we expect that the indirect method using $\mathbf{X}\mathbf{X}^T$ will be significantly faster than the direct method when $n \ll p$. The computational complexity of finding eigenvectors of an $n \times n$ matrix $\mathbf{X}\mathbf{X}^T$ is $\mathcal{O}(n^3)$, which is much more efficient when $n \ll p$.

Conclusion

This approach allows us to efficiently perform PCA when the number of features p is much larger than the number of samples n . The experiment demonstrates a reduction in computational time using the indirect method.

Problem 3

Let $\theta^* \in \mathbb{R}^d$ be the ground truth linear model parameter and $\mathbf{X} \in \mathbb{R}^{N \times d}$ be the observing matrix and each column of \mathbf{X} is independent. Assume the linear model is $\mathbf{y} = \mathbf{X}\theta^* + \epsilon$ where ϵ follows *Gaussian*(0, $\sigma^2\mathbf{I}$). Assume $\hat{\theta} = \arg \min_{\theta} \|\mathbf{X}\theta - \mathbf{y}\|^2$.

- Please show that $\mathbf{X}^T\mathbf{X}$ is invertible.
- Show that $MSE(\theta^*, \hat{\theta}) := E_{\epsilon}\{\|\theta^* - \hat{\theta}\|^2\} = \sigma^2 \text{trace}((\mathbf{X}^T\mathbf{X})^{-1})$
- Show that as N increases, MSE decreases. (hint: make use of ‘Woodbury matrix identity’)

Solution

Let $\theta^* \in \mathbb{R}^d$ be the ground truth linear model parameter, and $\mathbf{X} \in \mathbb{R}^{N \times d}$ be the observing matrix, where each column of \mathbf{X} is independent. We are given the linear model

$$\mathbf{y} = \mathbf{X}\theta^* + \epsilon$$

where $\epsilon \sim \mathcal{N}(0, \sigma^2\mathbf{I})$. Let

$$\hat{\theta} = \arg \min_{\theta} \|\mathbf{X}\theta - \mathbf{y}\|^2.$$

Part (a)

We need to show that $\mathbf{X}^T\mathbf{X}$ is invertible.

Proof: 1. Since each column of \mathbf{X} is independent, the columns of \mathbf{X} form a linearly independent set. 2. Therefore, $\mathbf{X}^T\mathbf{X}$, which is a $d \times d$ Gram matrix, has full rank. 3. A matrix with full rank is invertible. Hence, $\mathbf{X}^T\mathbf{X}$ is invertible.

—

Part (b)

We want to show that

$$MSE(\theta^*, \hat{\theta}) := \mathbb{E}_{\epsilon} \left\{ \|\theta^* - \hat{\theta}\|^2 \right\} = \sigma^2 \text{trace}((\mathbf{X}^T\mathbf{X})^{-1}).$$

Proof: 1. The estimator $\hat{\theta}$ is given by the least squares solution:

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

Substituting $\mathbf{y} = \mathbf{X}\theta^* + \epsilon$, we get

$$\hat{\theta} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T(\mathbf{X}\theta^* + \epsilon).$$

2. Simplifying, we have

$$\hat{\theta} = \theta^* + (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\epsilon.$$

3. The mean squared error (MSE) is given by

$$\text{MSE}(\theta^*, \hat{\theta}) = \mathbb{E}_{\epsilon} \left\{ \|\theta^* - \hat{\theta}\|^2 \right\} = \mathbb{E}_{\epsilon} \left\{ \|(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon\|^2 \right\}.$$

4. Since $\epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$, we have

$$\mathbb{E}_{\epsilon} \{ \epsilon \epsilon^T \} = \sigma^2 \mathbf{I}.$$

5. Thus,

$$\mathbb{E}_{\epsilon} \{ (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon \epsilon^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}.$$

6. Taking the trace, we get

$$\text{MSE}(\theta^*, \hat{\theta}) = \sigma^2 \text{trace}((\mathbf{X}^T \mathbf{X})^{-1}).$$

—

Part (c)

We want to show that as N increases, the MSE decreases using the Woodbury matrix identity.

Proof: 1. The Woodbury matrix identity states that for matrices \mathbf{A} , \mathbf{U} , \mathbf{C} , and \mathbf{V} of appropriate dimensions:

$$(\mathbf{A} + \mathbf{UCV})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{U} (\mathbf{C}^{-1} + \mathbf{VA}^{-1} \mathbf{U})^{-1} \mathbf{VA}^{-1}.$$

2. As N increases, we add more rows to \mathbf{X} , effectively increasing the information content in $\mathbf{X}^T \mathbf{X}$. 3. This results in $\mathbf{X}^T \mathbf{X}$ becoming better conditioned, meaning that $(\mathbf{X}^T \mathbf{X})^{-1}$ decreases in magnitude. 4. Consequently, $\text{trace}((\mathbf{X}^T \mathbf{X})^{-1})$ decreases, leading to a reduction in the MSE. 5. Hence, as N increases, $\text{MSE}(\theta^*, \hat{\theta})$ decreases, indicating improved estimation accuracy.