

Machine Learning Exercise Sheet 10

Dimensionality Reduction & Matrix Factorization, Part 1

In-class Exercises

Problem 1: In this exercise, we use proof by induction to show that the linear projection onto an M -dimensional subspace that maximizes the variance of the projected data is defined by the M eigenvectors of the data covariance matrix \mathbf{S} , given by

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T \quad \bar{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

corresponding to the M largest eigenvalues. In Section 12.1 in Bishop this result was proven for the case of $M = 1$. Now suppose the result holds for some general value of M and show that it consequently holds for dimensionality $M + 1$.

Problem 2: Proof that minimizing the error is equivalent to maximizing the variance.

① Mean and Variance $u_i^T \mathbf{x}_n$

$$E(u_i^T \mathbf{X}) = \frac{1}{N} \sum_{n=1}^N u_i^T \mathbf{x}_n = u_i^T \left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \right) = u_i^T \bar{\mathbf{x}}$$

$$\text{Var}(u_i^T \mathbf{X}) = \frac{1}{N} \sum_{n=1}^N (u_i^T \mathbf{x}_n - u_i^T \bar{\mathbf{x}})(u_i^T \mathbf{x}_n - u_i^T \bar{\mathbf{x}})^T = u_i^T \underbrace{\left[\frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T \right]}_{\mathbf{S}} u_i = u_i^T \mathbf{S} u_i$$

① construct Lagrangian $\max_{u_{M+1}} \text{Var}(u_{M+1}^T \mathbf{X})$
orthogonality & normalization constraints

$$L(u_{M+1}, \lambda_{M+1}, \eta_{1:N}) = \underbrace{u_{M+1}^T \mathbf{S} u_{M+1}}_{\text{scalar}} + \underbrace{\lambda_{M+1} (1 - u_{M+1}^T u_{M+1})}_{\text{variables}} + \sum_{i=1}^M \eta_i u_{M+1}^T u_i$$

= 0 iff $\|u_{M+1}\|=1$ = 0 iff $u_{M+1}^\perp u_i$

$$\frac{\partial L}{\partial u_{M+1}} = 2 \mathbf{S} u_{M+1} - 2 \lambda_{M+1} u_{M+1} + \sum_{i=1}^M \eta_i u_i = 0$$

$$2 u_{M+1}^T \mathbf{S} u_{M+1} - 2 \lambda_{M+1} u_{M+1}^T u_{M+1} + \sum_{i=1}^M \eta_i u_i^T u_{M+1} = 0$$

$$2 u_{M+1}^T \mathbf{S} u_{M+1} = 2 \lambda_{M+1}$$

$$u_{M+1}^T \mathbf{S} u_{M+1} = \lambda_{M+1}$$

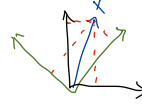
$$\textcircled{2} \max u_{M+1}^T \mathbf{S} u_{M+1} = \max \lambda_{M+1}$$

\Rightarrow to maximize choose λ_{M+1} to be the longest "remaining" eigenvalue

$\Rightarrow u_{M+1}$ must be an eigenvector & λ_{M+1} its corresponding eigenvalue

2. We have a complete orthonormal set of D -dim basis vectors $\{u_i\}$ where $i=1, \dots, D$
 $\vec{u}_i^T \cdot \vec{u}_i = 1$

$$\begin{aligned}\vec{x}_n &= \sum_{i=1}^D x_{ni} \vec{u}_i \\ &= \sum_{i=1}^D (\vec{u}_i^T \vec{x}_n) \vec{u}_i\end{aligned}$$



We approximate $\vec{x}_n \approx \tilde{x}_n = \sum_{i=1}^M z_{ni} \vec{u}_i + \sum_{i=M+1}^D b_i \vec{u}_i$

\uparrow \uparrow
 a variable for each sample n one variable for the complete data

objective: $\min_{\tilde{x}} \sum_{n=1}^N \|\vec{x} - \tilde{x}\|^2$
 $\text{rank}(\tilde{x}) = M$

$$\vec{x} - \tilde{x} = \sum_{i=1}^M (\vec{u}_i^T \vec{x} - z_{ni}) \vec{u}_i + \sum_{i=M+1}^D (\vec{u}_i^T \vec{x} - b_i) \vec{u}_i$$

$$\frac{\partial \|\vec{x} - \tilde{x}\|^2}{\partial z_{n1}} = \dots = 0$$

$\Rightarrow \sum_{n=1}^N \|\vec{x}_n - \tilde{x}\|^2$ is minimized if the M "first" components are 0

$\Rightarrow \vec{x} - \tilde{x} = \sum_{i=M+1}^D (\vec{u}_i^T \vec{x} - b_i) \vec{u}_i$

$$J = \sum_{n=1}^N \|\vec{x}_n - \tilde{x}\|^2 = \sum_{n=1}^N \sum_{i=M+1}^D [(\vec{u}_i^T \vec{x}_n - b_i) \vec{u}_i]^T [(\vec{u}_i^T \vec{x}_n - b_i) \vec{u}_i]$$

$$\begin{aligned}\frac{\partial J}{\partial b_i} &= \sum_{n=1}^N \frac{\partial}{\partial b_i} [(\vec{u}_i^T \vec{x}_n - b_i) \vec{u}_i]^T (-\vec{u}_i) = 0 \\ &= \sum_{n=1}^N (b_i \vec{u}_i^T \vec{u}_i - (\vec{u}_i^T \vec{x}_n) \vec{u}_i^T \vec{u}_i) \\ &= N b_i - \sum_{n=1}^N \vec{u}_i^T \vec{x}_n = 0 \\ b_i &= \frac{1}{N} \sum_{n=1}^N \vec{u}_i^T \vec{x}_n = \vec{u}_i^T \bar{x}\end{aligned}$$

$$J = \sum_{n=1}^N \sum_{i=M+1}^D (\vec{u}_i^T \vec{x}_n - \vec{u}_i^T \bar{x})^2 = \vec{u}_i^T S \vec{u}_i$$

Lagrangian:

$$\mathcal{L}(u_{M+1}, \lambda_{M+1}, \eta_{M+2:D}) = u_{M+1}^T S u_{M+1} + \lambda_{M+1} (1 - \vec{u}_{M+1}^T \vec{u}_{M+1}) + \sum_{i=M+2}^D \eta_i \vec{u}_{M+1}^T \vec{u}_i$$

... take smallest eigenvalue
 Problem: \nexists start with $M=D-1$ (equiv to $M=1$ in Problem 1)

Homework

PCA

Problem 3: Let the matrix $\mathbf{X} \in \mathbb{R}^{N \times D}$ represent N data points of dimension $D = 10$ (samples stored as rows). We applied PCA to \mathbf{X} . By using the $K = 5$ top principal components, we transformed/projected \mathbf{X} into $\tilde{\mathbf{X}} \in \mathbb{R}^{N \times K}$. We computed that $\tilde{\mathbf{X}}$ preserves 70% of the variance of the original data \mathbf{X} .

Suppose now we apply PCA on the following matrices:

- a) $\mathbf{Y}_1 = \mathbf{X}\mathbf{S}$ where $\mathbf{S} = \lambda\mathbf{I}$, with $\lambda \in \mathbb{R}$ and $\mathbf{I} \in \mathbb{R}^{D \times D}$ is the identity matrix
- b) $\mathbf{Y}_2 = \mathbf{X}\mathbf{R}$ where $\mathbf{R} \in \mathbb{R}^{D \times D}$ and $\mathbf{R}\mathbf{R}^T = \mathbf{I}$
- c) $\mathbf{Y}_3 = \mathbf{X}\mathbf{P}$ where $\mathbf{P} = \text{diag}(+5, -5, \dots, +5, -5)$ is a $D \times D$ diagonal matrix
- d) $\mathbf{Y}_4 = \mathbf{X}\mathbf{Q}$ where $\mathbf{Q} = \text{diag}(1, 2, 3, \dots, D-1, D)$ is a $D \times D$ diagonal matrix
- e) $\mathbf{Y}_5 = \mathbf{X} + \mathbf{1}_N \boldsymbol{\mu}^T$ where $\boldsymbol{\mu} \in \mathbb{R}^D$ and $\mathbf{1}_N$ is an N -dimensional column vector of all ones
- f) $\mathbf{Y}_6 = \mathbf{X}\mathbf{A}$ where $\mathbf{A} \in \mathbb{R}^{D \times D}$ and $\text{rank}(\mathbf{A}) = 5$

and obtain the projected data $\tilde{\mathbf{Y}}_1, \dots, \tilde{\mathbf{Y}}_6 \in \mathbb{R}^{N \times K}$ using the principal components corresponding to the top $K = 5$ largest eigenvalues of the respective \mathbf{Y}_i .

What fraction of variance of each \mathbf{Y}_i will be preserved by each respective $\tilde{\mathbf{Y}}_i$? *Justify your answer.*

The answer “cannot tell without additional information” is also valid if you provide a justification.

Problem 4: You are given $N = 4$ data points: $\{\mathbf{x}_i\}_{i=1}^4, \mathbf{x}_i \in \mathbb{R}^3$, represented with the matrix $\mathbf{X} \in \mathbb{R}^{4 \times 3}$.

$$\mathbf{X} = \begin{bmatrix} 4 & 3 & 2 \\ 2 & 1 & -2 \\ 4 & -1 & 2 \\ -2 & 1 & 2 \end{bmatrix}$$

Hint: In this task the results of all (final and intermediate) computations happen to be integers.

- a) Perform principal component analysis (PCA) of the data \mathbf{X} , i.e. find the principal components and their associated variances in the transformed coordinate system. Show your work.
- b) Project the data to two dimensions, i.e. write down the transformed data matrix $\mathbf{Y} \in \mathbb{R}^{4 \times 2}$ using the top-2 principal components you computed in (a). What fraction of variance of \mathbf{X} is preserved by \mathbf{Y} ?
- c) Let $\mathbf{x}_5 \in \mathbb{R}^3$ be a new data point. Specify the vector \mathbf{x}_5 such that performing PCA on the data including the new data point $\{\mathbf{x}_i\}_{i=1}^5$ leads to exactly the same principal components as in (a).

SVD

Problem 5: Use the SVD shown below. Suppose a new user Leslie assigns rating 3 to Alien and rating 4 to Titanic, giving us a representation of Leslie in the 'original space' of $[0, 3, 0, 0, 4]$. Find the representation of Leslie in concept space. What does that representation predict about how well Leslie would like the other movies appearing in our example data?

	Matrix	Alien	Star Wars	Casablanca	Titanic
Joe	1	1	1	0	0
Jim	3	3	3	0	0
John	4	4	4	0	0
Jack	5	5	5	0	0
Jill	0	0	0	4	4
Jenny	0	0	0	5	5
Jane	0	0	0	2	2

Figure 11.6: Ratings of movies by users

$$\begin{array}{c}
 \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 0 & 0 & 2 & 2 \end{bmatrix} = \begin{bmatrix} .14 & 0 \\ .42 & 0 \\ .56 & 0 \\ .70 & 0 \\ 0 & .60 \\ 0 & .75 \\ 0 & .30 \end{bmatrix} \begin{bmatrix} 12.4 & 0 \\ 0 & 9.5 \end{bmatrix} \begin{bmatrix} .58 & .58 & .58 & 0 & 0 \\ 0 & 0 & 0 & .71 & .71 \end{bmatrix} \\
 M \qquad \qquad U \qquad \qquad \Sigma \qquad \qquad V^T
 \end{array}$$

Problem 6: You want to perform linear regression on a data set with features $\mathbf{X} \in \mathbb{R}^{N \times D}$ and targets $\mathbf{y} \in \mathbb{R}^N$. Assume that you have already computed the SVD of the feature matrix $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$. Additionally, assume that \mathbf{X} has full rank and $N > D$.

Show how we can compute the optimal linear regression weights \mathbf{w}^* in $\mathcal{O}(ND)$ operations by using the result of the SVD.

Hint: Matrix operations have the following asymptotic complexity

- Matrix multiplication \mathbf{AB} for arbitrary $\mathbf{A} \in \mathbb{R}^{P \times Q}$ and $\mathbf{B} \in \mathbb{R}^{Q \times R}$ takes $\mathcal{O}(PQR)$
- Matrix multiplication \mathbf{AD} for an arbitrary $\mathbf{A} \in \mathbb{R}^{P \times Q}$ and a diagonal $\mathbf{D} \in \mathbb{R}^{Q \times Q}$ takes $\mathcal{O}(PQ)$
- Matrix inversion \mathbf{C}^{-1} for an arbitrary matrix $\mathbf{C} \in \mathbb{R}^{M \times M}$ takes $\mathcal{O}(M^3)$
- Matrix inversion \mathbf{D}^{-1} for a diagonal matrix $\mathbf{D} \in \mathbb{R}^{M \times M}$ takes $\mathcal{O}(M)$

Coding

Problem 7: Download the notebook `exercise_10_notebook.ipynb` from Moodle. Fill in the missing code and run the notebook. Convert the evaluated notebook to pdf and add it to the printout of your homework.