

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 $M + 1$.

Suppose that the result holds for projection spaces of dimensionality M . The $M + 1$ dimensional principal subspace will be defined by the M principal eigenvectors $\mathbf{u}_1, \dots, \mathbf{u}_M$ together with an additional direction vector \mathbf{u} whose value we wish to determine.

Following the argument given in section 12.1.1 for \mathbf{u}_1 , we see that the variance of the data projected onto the subspace spanned by $\{\mathbf{u}_1, \dots, \mathbf{u}_M, \mathbf{u}\}$ is given by $\mathbf{u}^T \mathbf{S} \mathbf{u} + \sum_{i=1}^M \mathbf{u}_i^T \mathbf{S} \mathbf{u}_i$. Since $\{\mathbf{u}_1, \dots, \mathbf{u}_M\}$ are already chosen and fixed, our goal is to maximize $\mathbf{u}^T \mathbf{S} \mathbf{u}$.

We must constrain \mathbf{u} such that it cannot be linearly related to $\mathbf{u}_1, \dots, \mathbf{u}_M$ (otherwise it will lie in the M -dimensional projection space instead of defining an $M + 1$ independent direction). For this, we introduce M constraints that require \mathbf{u} to be orthogonal to $\mathbf{u}_1, \dots, \mathbf{u}_M$. These constraints are enforced using Lagrange multipliers η_1, \dots, η_M .

Finally, we introduce a constraint with Lagrange multiplier λ that ensures that $\mathbf{u}^T \mathbf{u} = 1$.

Putting everything together, we obtain the Lagrangian

$$L(\mathbf{u}, \lambda, \eta_1, \dots, \eta_M) = \mathbf{u}^T \mathbf{S} \mathbf{u} + \lambda(1 - \mathbf{u}^T \mathbf{u}) + \sum_{i=1}^M \eta_i \mathbf{u}_i^T \mathbf{u}$$

The solution to a constrained optimization problem is found at the points $(\mathbf{u}, \lambda, \eta_1, \dots, \eta_M)$, where the gradient of the Lagrangian w.r.t. \mathbf{u} is zero and all other constraints are fulfilled, that is

$$\begin{cases} \nabla_{\mathbf{u}} L(\mathbf{u}, \lambda, \eta_1, \dots, \eta_M) = 2\mathbf{S}\mathbf{u} - 2\lambda\mathbf{u} + \sum_{i=1}^M \eta_i \mathbf{u}_i = \mathbf{0} \\ \mathbf{u}^T \mathbf{u} = 1 \\ \mathbf{u}^T \mathbf{u}_i = 0 \text{ for } i = 1, \dots, M \end{cases}$$

By left-multiplying the first equality with \mathbf{u}_j^T , we can conclude that it must hold $\eta_j = 0$ for $j = 1, \dots, M$. This means we can simplify the problem to

$$\begin{cases} \mathbf{S}\mathbf{u} = \lambda\mathbf{u} \\ \mathbf{u}^T\mathbf{u} = 1 \\ \mathbf{u}^T\mathbf{u}_i = 0 \text{ for } i = 1, \dots, M \end{cases}$$

This set of equalities is fulfilled for any eigenvalue-eigvector pair $(\lambda_i, \mathbf{u}_i)$ of \mathbf{S} , except the first M eigenvalues/eigenvectors, since those are not orthogonal to $\mathbf{u}_1, \dots, \mathbf{u}_M$.

Now we just need to select one of feasible solutions $\{\mathbf{u}_{M+1}, \dots, \mathbf{u}_D\}$ that maximizes the variance:

$$\arg \max_{i \in \{M+1, \dots, D\}} \mathbf{u}_i^T \mathbf{S} \mathbf{u}_i$$

This happens to be the $M + 1$'s eigenvector \mathbf{u}_{M+1} that corresponds to variance λ_{M+1} , since we assumed w.l.o.g. that the eigenvalues are sorted in decreasing order.

To summarize, the variance in the projected space is maximized by choosing \mathbf{u} to be the eigenvector having the largest eigenvalue amongst those not previously selected. Thus the result holds also for projection spaces of dimensionality $M + 1$, which completes the inductive step. Since we have already shown this result explicitly for $M = 1$ it follows that the result must hold for any $M \leq D$.

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

See Bishop chapter 12.1.2.

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.

- a) 70%. All eigenvalues are scaled by the same amount λ^2 , so the fraction doesn't change.
- b) 70%. \mathbf{R} is a rotation/reflection/permutation matrix. The direction of the eigenvectors of the covariance matrix is changed, but the eigenvalues stay the same.
- c) 70%. This is just combination of (a) and (b). All data points are scaled by 5 (i.e. eigenvalues of $\mathbf{X}^T\mathbf{X}$ are all scaled by 25), and some dimensions are reflected around origin, but the fraction of variance explained by the first K components stays the same.
- d) We cannot tell without additional information. since each column (i.e. each dimension) is scaled by a different amount.
- e) 70%. All data points are shifted by $\boldsymbol{\mu}$. But since we center the data as the first step of PCA, shifting has no effect.
- f) 100%. Since $\text{rank}(\mathbf{A}) = 5$, $\text{rank}(\mathbf{Y}_6) \leq 5$ as well. This means that the data lies in a ≤ 5 dimensional subspace, and the first 5 principal components capture all the variance.

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.

First we center the data. The mean is $\bar{\mathbf{x}} = [2, 1, 1]$, thus we have

$$\mathbf{X}_c = \mathbf{X} - \bar{\mathbf{x}} = \begin{bmatrix} 2 & 2 & 1 \\ 0 & 0 & -3 \\ 2 & -2 & 1 \\ -4 & 0 & 1 \end{bmatrix}$$

Then we compute the covariance matrix.

$$\Sigma_{X_c} = \frac{1}{N} \mathbf{X}_c^T \mathbf{X}_c = \begin{bmatrix} 6 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix}$$

Since Σ_{X_c} it is already in a diagonal form we can conclude that $\mathbf{\Lambda} = \Sigma_{X_c}$ and $\mathbf{\Gamma} = \mathbf{I}_3$, and that it holds $\Sigma_{X_c} = \mathbf{\Gamma} \mathbf{\Lambda} \mathbf{\Gamma}^T$. The principal components are the canonical basis vectors.

- 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} ?

The projection matrix is:

$$\mathbf{\Gamma}_{trunc} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{bmatrix}$$

since we pick the first and the third principal vector corresponding to the two largest eigenvalues. Thus, we have

$$\mathbf{Y} = \mathbf{X} \mathbf{\Gamma}_{trunc} = \begin{bmatrix} 2 & 1 \\ 0 & -3 \\ 2 & 1 \\ -4 & 1 \end{bmatrix}$$

We preserve $\frac{6+3}{6+2+3} = \frac{9}{11}$ of the variance.

- 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).

Let $\mathbf{x}_5 = \bar{\mathbf{x}}$, i.e. the new data point equals the mean before including \mathbf{x}_5 to the dataset. Therefore, the new mean including \mathbf{x}_5 is equal to the old mean. We have:

$$\mathbf{X}_c = \mathbf{X} - \bar{\mathbf{x}} = \begin{bmatrix} 2 & 2 & 1 \\ 0 & 0 & -3 \\ 2 & -2 & 1 \\ -4 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}$$

which leads to the same Σ_{X_c} as in (a) up to a difference in the multiplicative constant. In (a) we had $\frac{1}{4} \mathbf{X}_c^T \mathbf{X}_c$ and here we have $\frac{1}{5} \mathbf{X}_c^T \mathbf{X}_c$. While this difference leads to different eigenvalues, the eigenvectors and thus the principal components stay the same.

SVD

Problem 5: Use the SVD shown below. Suppose a new user Leslie assigns rating 3 to Alien and

	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

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

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?

The projection is given by $\mathbf{P} = \mathbf{M} \cdot \mathbf{V}$, thus the representation of Leslie in concept space is given by $[0, 3, 0, 0, 4] \cdot \mathbf{V} = [1.74, 2.84]$. It seems that Leslie has a higher preference for "classic" movies (the score is 2.84) such as "Titanic" and "Casablanca" compared to the "sci-fi" movies (the score is 1.74). Thus, since she already saw "Titanic", "Casablanca" would be a reasonable recommendation.

In general, if $\hat{\mathbf{U}}, \hat{\mathbf{\Sigma}}, \hat{\mathbf{V}}^T$ are the full singular values/vectors of \mathbf{M} (obtained by performing full SVD on \mathbf{M}) and $\mathbf{U}, \mathbf{\Sigma}, \mathbf{V}^T$ are the respective truncated versions (i.e. by taking only the top K singular values/vectors) it holds that the projected data \mathbf{P} can be obtained in two alternative and equivalent ways: $\mathbf{P} = \mathbf{U} \cdot \mathbf{\Sigma}$ or $\mathbf{P} = \mathbf{M} \cdot \mathbf{V}$. We usually prefer the second way since we only need to compute the top k singular vectors.

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 $\mathcal{O}(M)$

$$\begin{aligned}
 \mathbf{w}^* &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\
 &= ((\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^T (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T))^{-1} (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^T \mathbf{y} \\
 &= (\mathbf{V}\mathbf{\Sigma}\mathbf{U}^T \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^{-1} \mathbf{V}\mathbf{\Sigma}\mathbf{U}^T \mathbf{y} \\
 &= (\mathbf{V}\mathbf{\Sigma}^2 \mathbf{V}^T)^{-1} \mathbf{V}\mathbf{\Sigma}\mathbf{U}^T \mathbf{y} \\
 &= (\mathbf{V}^T)^{-1} (\mathbf{\Sigma}^2)^{-1} \mathbf{V}^{-1} \mathbf{V}\mathbf{\Sigma}\mathbf{U}^T \mathbf{y} \\
 &= \mathbf{V}\mathbf{\Sigma}^{-2} \mathbf{V}^T \mathbf{V}\mathbf{\Sigma}\mathbf{U}^T \mathbf{y} \\
 &= \mathbf{V}\mathbf{\Sigma}^{-2} \mathbf{V}^T \mathbf{V}\mathbf{\Sigma}\mathbf{U}^T \mathbf{y} \\
 &= \mathbf{V}\mathbf{\Sigma}^{-1} \mathbf{U}^T \mathbf{y}
 \end{aligned}$$

Multiplication $\mathbf{a} = \mathbf{U}^T \mathbf{y}$ takes $\mathcal{O}(N \cdot D \cdot 1)$

Multiplication $\mathbf{b} = \mathbf{\Sigma}^{-1} \mathbf{a}$ takes $\mathcal{O}(D)$

Multiplication $\mathbf{w} = \mathbf{V}\mathbf{b}$ takes $\mathcal{O}(D^2)$

In total, $\mathcal{O}(ND + D + D^2) = \mathcal{O}(ND)$ if $N > D$.

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.

The solution notebook is uploaded on Moodle.