

Machine Learning Exercise Sheet 11

Dimensionality Reduction & Matrix Factorization

Homework

PCA & SVD

Problem 1: 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

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

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 $\hat{\mathbf{U}}, \hat{\Sigma}, \hat{\mathbf{V}}^T$ are the full singular values/vectors of \mathbf{M} (obtained by performing full SVD on \mathbf{M}) and $\mathbf{U}, \Sigma, \mathbf{V}^T$ 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 2: Consider the latent space distribution

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$$

and a conditional distribution for the observed variable $\mathbf{x} \in \mathbb{R}^d$,

$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \boldsymbol{\Phi})$$

where $\boldsymbol{\Phi}$ is an arbitrary symmetric, positive-definite noise covariance variable. Now suppose that we make a nonsingular linear transformation of the data variables $\mathbf{y} = \mathbf{A}\mathbf{x}$ where \mathbf{A} is a non-singular $d \times d$ matrix. If $\boldsymbol{\mu}_{\text{MLE}}$, \mathbf{W}_{MLE} , and $\boldsymbol{\Phi}_{\text{MLE}}$ represent the maximum likelihood solution corresponding to the original untransformed data, show that $\mathbf{A}\boldsymbol{\mu}_{\text{MLE}}$, $\mathbf{A}\mathbf{W}_{\text{MLE}}$, and $\mathbf{A}\boldsymbol{\Phi}_{\text{MLE}}\mathbf{A}^T$ will represent the corresponding maximum likelihood solution for the transformed data set. Finally, show that the form of the model is preserved if \mathbf{A} is orthogonal and $\boldsymbol{\Phi}$ is proportional to the unit matrix so $\boldsymbol{\Phi} = \sigma^2 \mathbf{I}$ (i.e. probabilistic PCA). The transformed $\boldsymbol{\Phi}$ matrix remains proportional to the unit matrix, and hence probabilistic PCA is covariant under a rotation of the axes of data space, as is the case for conventional PCA.

The model for \mathbf{y} is a *noiseless* linear transformation. Since the distribution of \mathbf{x} is known we also know the distribution of \mathbf{y} . Because of the definitions for \mathbf{z} and $\mathbf{x}|\mathbf{z}$ we know that \mathbf{x} is a Gaussian with mean $\boldsymbol{\mu}$ and covariance $\mathbf{W}\mathbf{W}^T + \boldsymbol{\Phi}$. Thus, \mathbf{y} is also Gaussian with mean $\mathbf{A}\boldsymbol{\mu}$ and covariance $\mathbf{A}\mathbf{W}\mathbf{W}^T\mathbf{A}^T + \mathbf{A}\boldsymbol{\Phi}\mathbf{A}^T$. Now, assuming that the maximum likelihood solutions for the conditional model for \mathbf{x} are $\boldsymbol{\mu}_x$, \mathbf{W}_x and $\boldsymbol{\Phi}_x$, we obtain the MLE solutions for \mathbf{y} by simple *pattern matching* as $\mathbf{A}\boldsymbol{\mu}_x$, $\mathbf{A}\mathbf{W}_x$ and $\mathbf{A}\boldsymbol{\Phi}_x\mathbf{A}^T$.

Now, if \mathbf{A} is orthogonal and $\boldsymbol{\Phi}$ a scaled identity matrix, the model characteristics are also preserved since $\mathbf{A}\boldsymbol{\Phi}_x\mathbf{A}^T = \sigma^2 \mathbf{I}\mathbf{A}\mathbf{A}^T = \sigma^2 \mathbf{I}^2 = \sigma^2 \mathbf{I}$.

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.

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

Since $\boldsymbol{\Sigma}_{X_c}$ it is already in a diagonal form we can conclude that $\boldsymbol{\Lambda} = \boldsymbol{\Sigma}_{X_c}$ and $\boldsymbol{\Gamma} = \mathbf{I}_3$, and that it holds $\boldsymbol{\Sigma}_{X_c} = \boldsymbol{\Gamma} \boldsymbol{\Lambda} \boldsymbol{\Gamma}^T$. 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 $\mathbf{\Sigma}_{X_c}$ 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.

Problem 5: Download the notebook `exercise_11_notebook.ipynb` from Piazza. 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 Piazza.

In-class Exercises

Probabilistic PCA

Problem 6: For pPCA, we consider the latent space distribution

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$$

Upload a single PDF file with your homework solution to Moodle by 19.01.2020, 11:59pm CET. We recommend to typeset your solution (using \LaTeX or Word), but handwritten solutions are also accepted. If your handwritten solution is illegible, it won't be graded and you waive your right to dispute that.

and a conditional distribution for the observed variable $\mathbf{x} \in \mathbb{R}^d$,

$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I})$$

- Verify that the covariance of the marginal distribution $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{C})$ is given by $\mathbf{C} = \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}$. What is the interpretation of this result?
- Verify that the model is unidentifiable, i.e. that the matrix \mathbf{W} is only defined up to a rotation \mathbf{R} . What is the interpretation of this result?
- Derive an expression for the posterior of the latent variables $p(\mathbf{z}|\mathbf{x})$!

- a) We know from probability theory that $p(\mathbf{x}) = \int d\mathbf{z} p(\mathbf{x}, \mathbf{z}) = \int d\mathbf{z} p(\mathbf{z})p(\mathbf{x}|\mathbf{z})$. From the expression for marginals of linear Gaussian systems (Murphy equation 4.126 or Bishop equation 2.115) we derive

$$p(\mathbf{x}) = \int d\mathbf{z} \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I}) \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}).$$

A direct way to verify this result is given by observing that the distribution of $\mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\varepsilon}$ needs to be a Gaussian: \mathbf{z} is Gaussian ($p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$), a linear transformation $\mathbf{W}\mathbf{z} + \boldsymbol{\mu}$ of a Gaussian RV is also Gaussian, $\boldsymbol{\varepsilon}$ is Gaussian ($p(\boldsymbol{\varepsilon}) = \mathcal{N}(\boldsymbol{\varepsilon}|\mathbf{0}, \sigma^2 \mathbf{I})$), and the sum of two independent Gaussian RVs $(\mathbf{W}\mathbf{z} + \boldsymbol{\mu}) + \boldsymbol{\varepsilon}$ is also Gaussian.

Computing the mean and variance of the resulting Gaussian yields the same result as before:

$$\mathbb{E}[\mathbf{x}] = \mathbb{E}[\mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\varepsilon}] = \mathbf{W}\mathbb{E}[\mathbf{z}] + \mathbb{E}[\boldsymbol{\mu}] + \mathbb{E}[\boldsymbol{\varepsilon}] = \boldsymbol{\mu}$$

because of the linearity of the Expectation and the definitions given above. For the Covariance we obtain

$$\begin{aligned} \text{Cov}[\mathbf{x}] &= \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] = \mathbb{E}[(\mathbf{W}\mathbf{z} + \boldsymbol{\varepsilon})(\mathbf{W}\mathbf{z} + \boldsymbol{\varepsilon})^T] \\ &= \mathbb{E}[\mathbf{W}\mathbf{z}\mathbf{z}^T \mathbf{W}^T + \boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^T] = \mathbf{W}\mathbb{E}[\mathbf{z}\mathbf{z}^T] \mathbf{W}^T + \mathbb{E}[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^T] \\ &= \mathbf{W}\mathbf{I}\mathbf{W}^T + \sigma^2 \mathbf{I} = \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I} \end{aligned}$$

Interpretation: If the matrix \mathbf{W} is a $D \times M$ matrix this means that we get a low-rank approximation of a full Covariance $\boldsymbol{\Sigma}$. The main idea of pPCA is to “force the explanation” of the connection between the components x_i into the latent \mathbf{z} and the matrix \mathbf{W} by using a diagonal covariance in the Gaussian $p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I})$, instead of learning a general covariance as e.g. in $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ to explain the correlations between the x_i .

- b) Setting $\widetilde{\mathbf{W}} = \mathbf{W}\mathbf{R}$ with \mathbf{R} being an (orthonormal) rotation, we see that $\widetilde{\mathbf{W}}\widetilde{\mathbf{W}}^T = \mathbf{W}\mathbf{R}(\mathbf{W}\mathbf{R})^T = \mathbf{W}\mathbf{R}\mathbf{R}^T \mathbf{W}^T = \mathbf{W}\mathbf{I}\mathbf{W}^T = \mathbf{W}\mathbf{W}^T$. By that, the marginal distribution is unaltered:

$$p(\mathbf{x}|\boldsymbol{\mu}, \widetilde{\mathbf{W}}, \sigma) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \widetilde{\mathbf{W}}\widetilde{\mathbf{W}}^T + \sigma^2 \mathbf{I}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}) = p(\mathbf{x}|\boldsymbol{\mu}, \mathbf{W}, \sigma).$$

Interpretation: \mathbf{R} is a rotation in the latent space (the space of \mathbf{z}). Since the “generative story” of pPCA is “sample \mathbf{z} from $\mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$ and then sample \mathbf{x} from $\mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I})$ ”, we see that rotations in \mathbf{z} do not alter the model, because rotating \mathbf{z} before applying \mathbf{W} corresponds to having sampled a rotated \mathbf{z} from $\mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$, which is equally likely because the prior Gaussian is isotropic (has variance $\text{Cov}[\mathbf{z}] \propto \mathbf{I}$, here even $\text{Cov}[\mathbf{z}] = \mathbf{I}$).

- c) Using Murphy equation 4.125 or Bishop equation 2.116 for the conditional distribution, we can directly read off $p(\mathbf{z}|\mathbf{x})$ as $p(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_{\mathbf{z}|\mathbf{x}}, \boldsymbol{\Sigma}_{\mathbf{z}|\mathbf{x}})$ where

$$\begin{aligned}\boldsymbol{\Sigma}_{\mathbf{z}|\mathbf{x}}^{-1} &= \sigma^{-2}(\sigma^2 \mathbf{I} + \mathbf{W}^T \mathbf{W}) \\ \boldsymbol{\mu}_{\mathbf{z}|\mathbf{x}} &= \boldsymbol{\Sigma}_{\mathbf{z}|\mathbf{x}} [\sigma^{-2} \mathbf{W}^T (\mathbf{x} - \boldsymbol{\mu})]\end{aligned}$$

Using $\mathbf{M} = \sigma^2 \mathbf{I} + \mathbf{W}^T \mathbf{W}$ we can simplify these equations to obtain

$$\begin{aligned}\boldsymbol{\Sigma}_{\mathbf{z}|\mathbf{x}} &= \sigma^2 \mathbf{M}^{-1} \\ \boldsymbol{\mu}_{\mathbf{z}|\mathbf{x}} &= \mathbf{M}^{-1} \mathbf{W}^T (\mathbf{x} - \boldsymbol{\mu})\end{aligned}$$