

Machine Learning Homework Sheet 10

Dimensionality Reduction & Matrix Factorization

1 PCA & SVD

Problem 1: 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 $\boldsymbol{\mu}_{ML}$, \mathbf{W}_{ML} , and $\boldsymbol{\Phi}_{ML}$ represent the maximum likelihood solution corresponding to the original untransformed data, show that $\mathbf{A}\boldsymbol{\mu}_{ML}$, $\mathbf{A}\mathbf{W}_{ML}$, and $\mathbf{A}\boldsymbol{\Phi}_{ML}\mathbf{A}^T$ 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. Given that the distribution of \mathbf{x} is known, we therefore 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 $\mathbf{W}\mathbf{W}^T + \boldsymbol{\Phi}$. And 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 $\boldsymbol{\Phi}_x$, by simple *matching patterns* the MLE solutions for \mathbf{y} are $\mathbf{A}\boldsymbol{\mu}_x$, $\mathbf{A}\mathbf{W}_x$ 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 2: 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?

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.

	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

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

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

Problem 3: Load the notebook `homework_10_dim_reduction.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.

Note: We suggest that you use Anaconda for installing Python and Jupyter, as well as for managing packages. We recommend that you use Python 3.

For more information on Jupyter notebooks consult the Jupyter documentation. Instructions for converting the Jupyter notebooks to PDF are provided within the notebook.

2 Matrix Factorization

Problem 4: Load the notebook `homework_10_matrix_factorization.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.

Upload a single PDF file with your solution to Moodle by 20.01.2019, 11:59pm CET. We recommend to typeset your solution (using L^AT_EX 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.

Note: We suggest that you use Anaconda for installing Python and Jupyter, as well as for managing packages. We recommend that you use Python 3.

For more information on Jupyter notebooks consult the Jupyter documentation. Instructions for converting the Jupyter notebooks to PDF are provided within the notebook.

3 Autoencoders

Problem 5: We train a linear autoencoder to D -dimensional data. The autoencoder has a single K -dimensional hidden layer, there are no biases, and all activation functions are identity ($\sigma(x) = x$).

- Why is it usually impossible to get zero reconstruction error in this setting if $K < D$?
- Under which conditions is this possible?

We have $f(\mathbf{x}) = \mathbf{X}\mathbf{W}_1\mathbf{W}_2$ where \mathbf{X} is the data matrix and the dimensions of the weight matrices are $D \times K$ for \mathbf{W}_1 and $K \times D$ for \mathbf{W}_2 .

The final multiplication \mathbf{W}_2 brings points from K -dimensions up into D -dimensions but the points will still all be in a K -dimensional linear subspace. Unless the data happen to lie exactly in a K -dimensional linear subspace, they can't be exactly fitted.