

CS/ECE/ME532 Activity 12

Estimated time: 5 mins for Q1, 15 mins for Q2 (review), 20 mins for Q3, and 20 mins for Q4.

1. Let the n -by- p rank- r ($n > p > r$) matrix \mathbf{X} have SVD $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ where \mathbf{U} is n -by- r , $\mathbf{\Sigma}$ is r -by- r , and \mathbf{V} is p -by- r .

a) Find the SVD of $\mathbf{Z} = \mathbf{X}^T$ in terms of \mathbf{U} , $\mathbf{\Sigma}$, and \mathbf{V} . $\mathbf{Z} = \mathbf{V}\mathbf{\Sigma}^T\mathbf{U}^T$

b) Find the orthonormal basis for the best rank-1 subspace to approximate the rows of \mathbf{Z} in terms of \mathbf{U} , \mathbf{V} , and $\mathbf{\Sigma}$.

basis = \mathbf{u}_1

first column of \mathbf{u}

2. **Uniqueness of solutions and Tikhonov regularization (ridge regression).**

The least-squares problem is $\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$. Assume \mathbf{X} is n -by- p with $p < n$.

a) Under what conditions is the solution to the least-squares problem not unique?

b) The Tikhonov-regularized least-squares problem is $\text{rank}(\mathbf{X}) < p$

$$\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

Show that this can be written as an ordinary least-squares problem $\min_{\mathbf{w}} \|\hat{\mathbf{y}} - \hat{\mathbf{X}}\mathbf{w}\|_2^2$ and find $\hat{\mathbf{y}}$ and $\hat{\mathbf{X}}$.

$$\hat{\mathbf{y}} = \begin{bmatrix} \mathbf{y} \\ 0 \end{bmatrix} \quad \hat{\mathbf{X}} = \begin{bmatrix} \mathbf{X} \\ \lambda^{1/2} \mathbf{I} \end{bmatrix}$$

c) Use the results from the previous part to determine the conditions for which the Tikhonov-regularized least-squares problem has a unique solution.

when $\lambda > 0$, this ensures the diagonal of $\hat{\mathbf{X}}^T \hat{\mathbf{X}}$ has no 0s, so it is now invertible and has a unique solution. $\hat{\mathbf{X}}$ also becomes a full rank matrix with $\text{rank}(\hat{\mathbf{X}}) = p$.

3. **Pseudoinverse and truncated SVD.** The solution to the ridge regression problem

$$\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

is given by $\mathbf{w}^* = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$. The *pseudoinverse* of \mathbf{X} , denoted \mathbf{X}^\dagger , can be defined by looking at the limit of the ridge regression solution as $\lambda \rightarrow 0$ (from above):

$$\mathbf{X}^\dagger = \lim_{\lambda \downarrow 0} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T.$$

$$\begin{aligned}
 (V \Sigma^2 V^T + V \lambda I V^T)^{-1} X^T &= V (\Sigma^2 + \lambda I)^{-1} V^T V \Sigma^T U^T \\
 (V (\Sigma^2 + \lambda I) V^T)^{-1} X^T &= V (\Sigma^2 + \lambda I)^{-1} \Sigma U^T \\
 V (\Sigma^2 + \lambda I)^{-1} V^T X^T &= \sum_{i=1}^p \frac{\sigma_i}{\sigma_i^2 + \lambda} v_i u_i^T
 \end{aligned}$$

a) Let $X \in \mathbb{R}^{n \times p}$, $p \leq n$, have SVD $X = U \Sigma V^T = \sum_{i=1}^p \sigma_i u_i v_i^T$. Show that

$$(X^T X + \lambda I)^{-1} X^T = \sum_{i=1}^p \frac{\sigma_i}{\sigma_i^2 + \lambda} v_i u_i^T.$$

Hint: Note that $X^T X = V \Sigma^2 V^T$ and $\lambda I = V \lambda I V^T$.

- b) Using the limit definition of the psuedoinverse above, show that when $X^T X$ is invertible, then $X^\dagger = (X^T X)^{-1} X^T$. When $X^T X$ is invertible, we don't need the λ term to keep the denominator > 0 , so $(X^T X + \lambda I)^{-1} X^T$ becomes $(X^T X)^{-1} X^T$.
- c) Argue that when X is square and invertible, then $X^\dagger = X^{-1}$.
- d) Argue that if X is rank $r < p$, then for $\lambda > 0$,

The $\lambda > 0$ ensures denominator > 0 , so that the pseudoimverse is actually invertible.

$$(X^T X + \lambda I)^{-1} X^T = \sum_{i=1}^r \frac{\sigma_i}{\sigma_i^2 + \lambda} v_i u_i^T.$$

e) Now argue that if X is rank $r < p$,

$$X^\dagger = \sum_{i=1}^r \frac{1}{\sigma_i} v_i u_i^T = V \Sigma_r^{-1} U^T$$

where Σ_r^{-1} is a matrix with $1/\sigma_i$ on the diagonal for $i = 1, \dots, r$, and zero elsewhere.

rank r means that the first r singular values are > 0 . By only using those singular values, the denominator will always be > 0 . It gives the same X^\dagger because the other

4. The data file is available with a matrix X of 100 three-dimensional data points. A script is available with code to assist you with visualizing and fitting this data. Use the results of the SVD to find a , a basis for the best (minimum sum of squared distances) one-dimensional subspace for the data.

- a) Run the code to display the data in Figure the first figure. Use the rotate tool to inspect the scatter plot from different angles. Does the data appear to lie very close to a one-dimensional subspace? Does the data appear to be zero mean?
- Yes No, not centered at (0,0,0)
- b) Figure 2 depicts the centered data and the one-dimensional subspace that contains the dominant feature you identified using the SVD. Use the rotate tool to inspect the data and one-dimensional subspace from different angles. Is a one-dimensional subspace a reasonable fit to the data? Comment on the error. The subspace is in the center of the data cloud and follows the trend of the data.
- c) Now comment out (insert %) the line of code that subtracts the mean of the data. Does the dominant feature identified by SVD continue to be a good fit to the data? Comment on the importance of removing the mean before performing PCA. The principal component now points in a completely different direction from the trend of the data and is not even in the data cloud.