

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} .
- 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 $\mathbf{\Sigma}$.

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

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

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 $\hat{\mathbf{X}}$.

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

3. Psuedoinverse 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 *psuedoinverse* 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):

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

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

$$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T = \sum_{i=1}^p \frac{\sigma_i}{\sigma_i^2 + \lambda} \mathbf{v}_i \mathbf{u}_i^T.$$

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

- b) Using the limit definition of the psuedoinverse above, show that when $\mathbf{X}^T \mathbf{X}$ is invertible, then $\mathbf{X}^\dagger = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$.
- c) Argue that when \mathbf{X} is square and invertible, then $\mathbf{X}^\dagger = \mathbf{X}^{-1}$.
- d) Argue that if \mathbf{X} is rank $r < p$, then for $\lambda > 0$,

$$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T = \sum_{i=1}^r \frac{\sigma_i}{\sigma_i^2 + \lambda} \mathbf{v}_i \mathbf{u}_i^T.$$

- e) Now argue that if \mathbf{X} is rank $r < p$,

$$\mathbf{X}^\dagger = \sum_{i=1}^r \frac{1}{\sigma_i} \mathbf{v}_i \mathbf{u}_i^T = \mathbf{V} \Sigma_r^{-1} \mathbf{U}^T$$

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

4. The data file is available with a matrix \mathbf{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 \mathbf{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?
- 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.
- 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.

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

- Find the SVD of $Z = X^T$ in terms of U , Σ , and V .
- Find the orthonormal basis for the best rank-1 subspace to approximate the rows of Z in terms of U , V , and Σ .

$$\begin{array}{c} \underline{X} \\ \left[\begin{array}{c} n \times p \end{array} \right] \\ \text{rank } r \end{array} = \begin{array}{ccc} \underline{U} & \underline{\Sigma} & \underline{V}^T \\ n \times r & r \times r & r \times p \end{array}$$

Best approx for rows

- Left singular vectors
- (columns of left matrix)

Best approx for columns

- Right singular vectors
- (rows of right matrix)

$$\text{a.) } \begin{array}{ccccc} \underline{Z} & = & \underline{X}^T & = & \underline{V} \underline{\Sigma} \underline{U}^T \\ p \times n & & p \times n & & p \times r \quad r \times r \quad r \times n \end{array}$$

b.) This is given by the first left singular vector of the decomposition.

In the case of \underline{Z} , that is the first column of \underline{V} .

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

The least-squares problem is $\min_w \|y - Xw\|_2^2$. Assume X is n -by- p with $p < n$.

- Under what conditions is the solution to the least-squares problem not unique?
- The Tikhonov-regularized least-squares problem is

$$\min_w \|y - Xw\|_2^2 + \lambda \|w\|_2^2$$

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

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

$$a.) \begin{matrix} X \\ n \times p \end{matrix} \quad \begin{bmatrix} X \\ \end{bmatrix}, \quad \min_w \|y - Xw\|_2^2$$

→ The least-squares problem does not have a unique solution when X is not full column rank.

• This can occur when $Xw = y$ has no solution or when $Xw = y$ has infinitely many solutions.

$$b.) \min_w \|Xw - y\|_2^2 + \lambda \|w\|_2^2$$

$$\rightarrow \text{use identity: } \|a\|_2^2 + \|b\|_2^2 = \left\| \begin{bmatrix} a \\ b \end{bmatrix} \right\|_2^2$$

$$\rightarrow \left\| \begin{bmatrix} Xw - y \\ \lambda^{1/2} w \end{bmatrix} \right\|_2^2 \rightarrow \left\| \underbrace{\begin{bmatrix} X \\ \lambda^{1/2} I \end{bmatrix}}_{\tilde{X}} w - \underbrace{\begin{bmatrix} y \\ 0 \end{bmatrix}}_{\tilde{y}} \right\|_2^2$$

$$\rightarrow \|\tilde{X}w - \tilde{y}\|_2^2$$

- The Tikhonov-reg LS equation always has a unique solution because \tilde{X} is guaranteed to be full rank due to the addition of λI .

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

$$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T = \sum_{i=1}^p \frac{\sigma_i}{\sigma_i^2 + \lambda} \mathbf{v}_i \mathbf{u}_i^T.$$

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

- b) Using the limit definition of the pseudoinverse above, show that when $\mathbf{X}^T \mathbf{X}$ is invertible, then $\mathbf{X}^\dagger = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$.
- c) Argue that when \mathbf{X} is square and invertible, then $\mathbf{X}^\dagger = \mathbf{X}^{-1}$.
- d) Argue that if \mathbf{X} is rank $r < p$, then for $\lambda > 0$,

$$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T = \sum_{i=1}^r \frac{\sigma_i}{\sigma_i^2 + \lambda} \mathbf{v}_i \mathbf{u}_i^T.$$

- e) Now argue that if \mathbf{X} is rank $r < p$,

$$\mathbf{X}^\dagger = \sum_{i=1}^r \frac{1}{\sigma_i} \mathbf{v}_i \mathbf{u}_i^T = \mathbf{V} \Sigma_r^{-1} \mathbf{U}^T$$

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

$$\begin{aligned} \text{a.) } & (\underline{\mathbf{X}}^T \underline{\mathbf{X}} + \lambda \underline{\mathbf{I}})^{-1} \underline{\mathbf{X}}^T \\ &= (\underline{\mathbf{V}} \underline{\Sigma} \cancel{\underline{\mathbf{U}}^T} \underline{\mathbf{U}} \underline{\Sigma} \underline{\mathbf{V}}^T + \lambda \underline{\mathbf{I}})^{-1} \underline{\mathbf{V}} \underline{\Sigma} \underline{\mathbf{U}}^T \\ &= (\underline{\mathbf{V}} \underline{\Sigma}^2 \underline{\mathbf{V}}^T + \underline{\mathbf{V}} \lambda \underline{\mathbf{I}} \underline{\mathbf{V}}^T)^{-1} \underline{\mathbf{V}} \underline{\Sigma} \underline{\mathbf{U}}^T \\ &= (\underline{\mathbf{V}} (\underline{\Sigma}^2 + \lambda \underline{\mathbf{I}}) \underline{\mathbf{V}}^T)^{-1} \underline{\mathbf{V}} \underline{\Sigma} \underline{\mathbf{U}}^T \\ &= \underline{\mathbf{V}} (\underline{\Sigma}^2 + \lambda \underline{\mathbf{I}})^{-1} \cancel{\underline{\mathbf{V}}^T} \underline{\mathbf{V}} \underline{\Sigma} \underline{\mathbf{U}}^T \\ &= \underline{\mathbf{V}} (\underline{\Sigma}^2 + \lambda \underline{\mathbf{I}})^{-1} \underline{\Sigma} \underline{\mathbf{U}}^T \\ &= \sum_{i=1}^p \frac{\sigma_i}{\sigma_i^2 + \lambda} \underline{\mathbf{v}}_i \underline{\mathbf{u}}_i^T \end{aligned}$$

$$\begin{aligned} & \textcircled{1} \begin{bmatrix} 1/\sigma_1^2 + \lambda & 0 \\ 0 & \ddots & 1/\sigma_p^2 + \lambda \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 \\ 0 & \ddots & \sigma_p \end{bmatrix} \\ &= \begin{bmatrix} \sigma_1/\sigma_1^2 + \lambda & 0 \\ 0 & \ddots & \sigma_p/\sigma_p^2 + \lambda \end{bmatrix} \end{aligned}$$

This is just a scaling matrix that scales each $\underline{\mathbf{v}}_i$ and $\underline{\mathbf{u}}_i^T$ by $\sigma_i / (\sigma_i^2 + \lambda)$

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

if $\underline{\underline{X}}^T \underline{\underline{X}}$ is invertible, then as $\lambda \downarrow 0$,

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

c.) When $\underline{\underline{X}}$ is square and invertible,

$$\underline{\underline{X}}^\dagger = (\underline{\underline{X}}^T \underline{\underline{X}})^{-1} \underline{\underline{X}}^T = \underline{\underline{X}}^{-1} \underbrace{\underline{\underline{X}}^{-T} \underline{\underline{X}}^T}_{\text{this becomes } \underline{\underline{I}}} = \underline{\underline{X}}^{-1}.$$

this becomes $\underline{\underline{I}}$ because $\underline{\underline{X}}^{-1} \underline{\underline{X}} = \underline{\underline{I}}$, and this property holds for the transpose as well.

$$\underline{\underline{X}} \underline{\underline{X}}^{-1} = \underline{\underline{I}} \Leftrightarrow (\underline{\underline{X}} \underline{\underline{X}}^{-1})^T = \underline{\underline{I}}^T$$

$$\Rightarrow \underline{\underline{X}}^{-T} \underline{\underline{X}}^T = \underline{\underline{I}}.$$

d.) In this case, return to the term:

$$(\underline{X}^T \underline{X} + \lambda \underline{I})^{-1} \underline{X}^T = \underline{V} \underbrace{(\underline{\Sigma}^2 + \lambda \underline{I})^{-1} \underline{\Sigma}}_{\textcircled{1}} \underline{U}^T$$

$$\textcircled{1} \begin{bmatrix} 1/\delta_1^2 + \lambda & & 0 \\ & \ddots & \\ 0 & & 1/\delta_r^2 + \lambda \\ & & 0 & \ddots \end{bmatrix} \begin{bmatrix} \delta_1 & & 0 \\ & \ddots & \\ 0 & & \delta_r \\ & & 0 & \ddots \end{bmatrix} = \begin{bmatrix} \delta_1/\delta_1^2 + \lambda & & 0 \\ & \ddots & \\ 0 & & \delta_r/\delta_r^2 + \lambda \\ & & 0 & \ddots \end{bmatrix}$$

Because now there are only $r < p$ nonzero singular values.

$$(\underline{X}^T \underline{X} + \lambda \underline{I})^{-1} \underline{X}^T = \sum_{i=1}^r \frac{\delta_i}{\delta_i^2 + \lambda} \underline{V}_i \underline{U}_i^T$$

e.) \underline{X}^+ is the limit as $\lambda \downarrow 0$ for the expression above.

$$\underline{X}^+ = \sum_{i=1}^r \frac{\delta_i}{\delta_i^2 + 0} \underline{V}_i \underline{U}_i^T = \sum_{i=1}^r \frac{1}{\delta_i} \underline{V}_i \underline{U}_i^T = \underline{V} \underline{\Sigma}^{-1} \underline{U}^T$$

Back to that step, but with $\lambda = 0 \dots$

$$(\underline{X}^T \underline{X})^{-1} \underline{X}^T = \underline{V} (\underline{\Sigma}^2)^{-1} \underline{\Sigma} \underline{U}^T = \underline{V} \underline{\Sigma}^{-2} \underline{\Sigma} \underline{U}^T = \underline{V} \underline{\Sigma}^{-1} \underline{U}^T$$

4a.) Yes, the data appears to lie very close to a line (1-d subspace)

No, the data appears to be skewed off the mean in the $-ve$ x_1 direction.

4b.) The line produced captures a lot of the variance in the data.

There is still a decent amount of variance in the x_3 direction.

4c.) No, the line now points in a completely different direction.

Removing the mean before performing PCA is crucial to ensure the correctly-oriented principal components are identified.

$$\underline{X} \underline{w} = \underline{y} \quad \text{rank}[\underline{X}] \leq \text{rank}[\underline{X} : \underline{y}] \rightarrow \text{No solution}$$

$$\text{rank}[\underline{X} : \underline{y}] = \text{rank}[\underline{X}] \rightarrow \geq 1 \text{ exact solution}$$

$$\hookrightarrow \begin{cases} \text{rank}(\underline{X}) < P & \rightarrow \infty \text{ many solutions} \\ \text{rank}(\underline{X}) = P & \rightarrow \text{Unique solution.} \end{cases}$$

(rank deficient)
(full rank)

$$\min_{\underline{w}} \|\underline{X}\underline{w} - \underline{y}\|_2^2$$

System has an exact, unique solution:

$$\underline{w}_{\min} = \underline{X}^{-1} \underline{y}, \text{ unique}$$

System has no solution:

if \underline{X} is full column rank $\rightarrow \underline{X}^T \underline{X}$ is invertible $\rightarrow \underline{w}_{\min}$ unique

\underline{X} is not full column rank:

if \underline{X} is ~~not~~ full column rank $\rightarrow \underline{X}^T \underline{X}$ is not invertible \rightarrow