

Applied Statistics Qualifying Exams Coaching

Michael Howes*

Summer 2023

Contents

0	Applied 2022	2
---	--------------	---

*With lots of content credit given to previous applied quals coaches

0 Applied 2022¹

Problem 1: R-squared and PCA

Key ideas:

- Connections between principal components analysis and the singular value decomposition.
 - Orthogonality of principal components and principal component directions.
- (a) We are given that $\mathbf{X} \in \mathbb{R}^{n \times p}$ is standardized so that the columns have mean zero and variance one. We are also given that $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^\top$ is the SVD of \mathbf{X} . Let $d_1 \geq d_2 \geq \dots \geq d_p \geq 0$ be the diagonal entries of \mathbf{D} . The columns of \mathbf{V} are thus the principal component directions and $\frac{1}{n}d_j^2 = \frac{1}{n}\|\mathbf{X}\mathbf{v}_j\|_2^2$ is the variance of \mathbf{X} in the j th principal component direction. The cumulative percent variance explained sequence is thus,

$$\rho_k = 100 \times \frac{\sum_{s=1}^k \frac{1}{n}d_s^2}{\sum_{s=1}^p \frac{1}{n}d_s^2} = 100 \times \frac{\sum_{s=1}^k d_s^2}{\sum_{s=1}^p d_s^2},$$

for $k = 1, \dots, p$.

- (b) We are given a response $\mathbf{y} \in \mathbb{R}^n$ with mean zero and variance one. The fitted values are $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} \in \mathbb{R}^n$. We know that the fitted values $\hat{\mathbf{y}}$ are orthogonal to the residuals $\mathbf{y} - \hat{\mathbf{y}}$. Thus,

$$\|\mathbf{y}\|_2^2 = \|\mathbf{y} - \hat{\mathbf{y}} + \hat{\mathbf{y}}\|_2^2 = \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 + \|\hat{\mathbf{y}}\|_2^2.$$

And so

$$\frac{1}{n} \sum_{i=1}^n y_i^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \frac{1}{n} \sum_{i=1}^n \hat{y}_i^2.$$

Since \mathbf{y} has mean zero and variance one, this implies that

$$1 = \text{MSE}_0 = \text{MSE} + \text{MSS}.$$

And so,

$$R^2 = 1 - \frac{\text{MSE}}{\text{MSE}_0} = 1 - \text{MSE} = \text{MSS}.$$

- (c) We now regress the j th column of \mathbf{X} on the first k principal components of \mathbf{X} . Let $\mathbf{U}_k = [\mathbf{u}_1, \dots, \mathbf{u}_k] \in \mathbb{R}^{n \times k}$ be the matrix containing the first k columns of \mathbf{U} . Since $\mathbf{U}_k^\top \mathbf{U}_k = \mathbf{I}_k$, the fitted values from regressing \mathbf{x}_j on \mathbf{U}_k are

$$\hat{\mathbf{x}}_j = \mathbf{U}_k(\mathbf{U}_k^\top \mathbf{U}_k)^{-1} \mathbf{U}_k^\top \mathbf{x}_j = \mathbf{U}_k \mathbf{U}_k^\top \mathbf{x}_j.$$

¹Michael Howes

The average regression sum of squares is thus,

$$\begin{aligned}\text{MSS}_j &= \frac{1}{n} \hat{\mathbf{x}}_j^\top \hat{\mathbf{x}}_j \\ &= \frac{1}{n} \mathbf{x}_j^\top \mathbf{U}_k \mathbf{U}_k^\top \mathbf{U}_k \mathbf{U}_k^\top \mathbf{x}_j \\ &= \frac{1}{n} \mathbf{x}_j^\top \mathbf{U}_k \mathbf{U}_k^\top \mathbf{x}_j.\end{aligned}$$

Since \mathbf{x}_j is the j th column of \mathbf{X} we $\mathbf{x}_j = \mathbf{X} \mathbf{e}_j$ where $\mathbf{e}_j \in \mathbb{R}^p$ is the j th standard basis vector.

$$\begin{aligned}\mathbf{U}_k^\top \mathbf{x}_j &= \mathbf{U}_k^\top \mathbf{X} \mathbf{e}_j \\ &= \mathbf{U}_k^\top \mathbf{U} \mathbf{D} \mathbf{V}^\top \mathbf{e}_j.\end{aligned}$$

We know that $\mathbf{U}_k^\top \mathbf{U} = [\mathbf{I}_k, \mathbf{0}_{k \times (p-k)}] \in \mathbb{R}^{k \times p}$ where $\mathbf{0}_{k \times (p-k)}$ is a matrix of all zeros of size $k \times (p-k)$. It follows that

$$\begin{aligned}\mathbf{U}_k^\top \mathbf{U} \mathbf{D} \mathbf{V}^\top &= [\mathbf{I}_k, \mathbf{0}_{k \times (p-k)}] \mathbf{D} \mathbf{V}^\top \\ &= \mathbf{D}_k \mathbf{V}_k^\top,\end{aligned}$$

where $\mathbf{D}_k = \text{diag}(d_1, \dots, d_k) \in \mathbb{R}^{k \times k}$ and $\mathbf{V}_k \in \mathbb{R}^{p \times k}$ is equal to the first k rows of \mathbf{V} . Thus,

$$\begin{aligned}\mathbf{U}_k^\top \hat{\mathbf{x}}_j &= \mathbf{D}_k \mathbf{V}_k^\top \mathbf{e}_j \\ &= \sum_{s=1}^k d_s (\mathbf{v}_s^\top \mathbf{e}_j) \mathbf{e}_s \\ &= \sum_{s=1}^k d_s v_{sj} \mathbf{e}_s\end{aligned}$$

where v_{sj} is the entry of \mathbf{V} in row s and column j . We thus have

$$\begin{aligned}\text{MSS}_j &= \frac{1}{n} \|\mathbf{U}_k^\top \hat{\mathbf{x}}_j\|_2^2 \\ &= \frac{1}{n} \left\| \sum_{s=1}^k d_s v_{sj} \mathbf{e}_s \right\|_2^2 \\ &= \frac{1}{n} \sum_{s=1}^k d_s^2 v_{sj}^2.\end{aligned}$$

Since \mathbf{x}_j has mean zero and variance one, we are in the setting of part (b) and hence

$$R_j^2 = \text{MSS}_j = \frac{1}{n} \sum_{s=1}^k d_s^2 v_{sj}^2.$$

(d) Note that

$$\begin{aligned}
\sum_{j=1}^p \text{MSS}_j &= \sum_{j=1}^p \frac{1}{n} \sum_{s=1}^k d_s^2 v_{sj}^2 &= \sum_{s=1}^k \sum_{j=1}^p \frac{1}{n} d_s^2 v_{sj}^2 \\
&= \sum_{s=1}^k \frac{1}{n} d_s^2 \sum_{j=1}^p v_{sj}^2 \\
&= \sum_{s=1}^k \frac{1}{n} d_s^2 \|\mathbf{v}_s\|_2^2 \\
&= \frac{1}{n} \sum_{s=1}^k d_s^2,
\end{aligned}$$

since all rows of \mathbf{V} have norm one. We know that each column of \mathbf{X} has variance one and mean zero. Thus,

$$\begin{aligned}
p &= \sum_{j=1}^p \frac{1}{n} \mathbf{x}_j^\top \mathbf{x}_j \\
&= \frac{1}{n} \sum_{j=1}^p \text{tr}(\mathbf{x}_j^\top \mathbf{x}_j) \\
&= \frac{1}{n} \sum_{j=1}^p \text{tr}(\mathbf{x}_j \mathbf{x}_j^\top) \\
&= \frac{1}{n} \text{tr} \left(\sum_{j=1}^p \mathbf{x}_j \mathbf{x}_j^\top \right) \\
&= \frac{1}{n} \text{tr}(\mathbf{X} \mathbf{X}^\top) \\
&= \frac{1}{n} \text{tr}(\mathbf{U} \mathbf{D} \mathbf{V}^\top \mathbf{V} \mathbf{D} \mathbf{U}^\top) \\
&= \frac{1}{n} \text{tr}(\mathbf{D}^2) \\
&= \frac{1}{n} \sum_{s=1}^p d_s^2.
\end{aligned}$$

Thus,

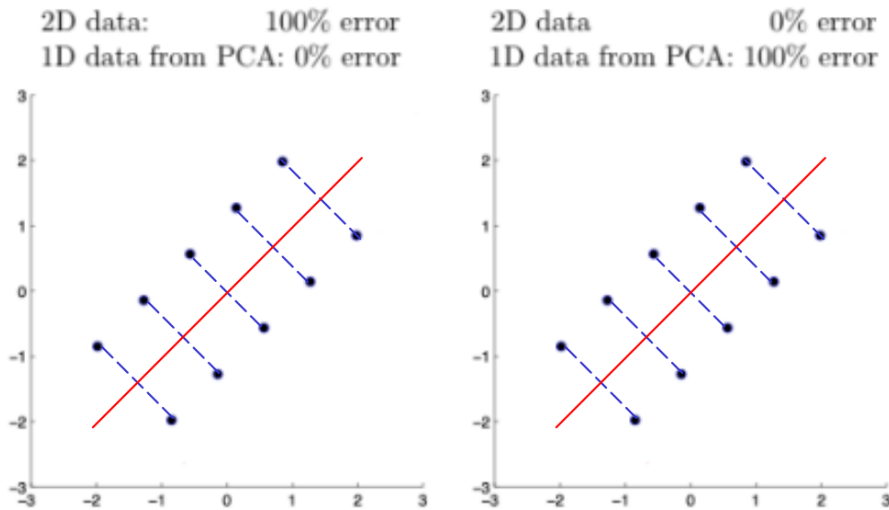
$$\frac{100}{p} \sum_{j=1}^p \text{MSS}_j = \frac{100}{np} \sum_{s=1}^k d_s^2 = 100 \times \frac{\sum_{s=1}^k d_s^2}{\sum_{s=1}^p d_s^2} = \rho_k.$$

So $\frac{100}{p} \sum_{j=1}^p \text{MSS}_j$ is exactly the cumulative percent variance explained sequence.

Problem 2: LOO CV, PCA and 1NN

Key ideas:

- “Eyeballing” principal component directions.
 - Working out nearest neighbor classifiers.
1. We are asked to draw a line corresponding to the first principal component. This does not have to be exact, but we can see that a roughly 45 degree line gives the direction with the most variance. Here the principal component direction is represented as a red solid line. In both plots I have also included dashed lines showing the projection onto the first principal direction.



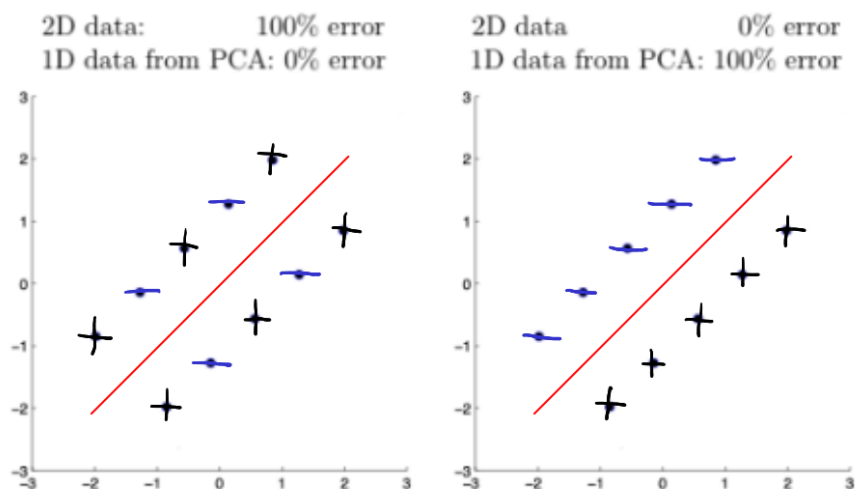
2. Looking at the previous figure we can make two observations.
 - Using the 2D data, the nearest neighbor of a point is one of the adjacent points on the line of points parallel to the PCA direction. For example, for the point at roughly $(-1, -2)$, the nearest neighbor is the point at roughly $(0, -1.25)$.
 - Using the projected 1D data from PCA, the nearest neighbor of a point is the point across the PCA direction. This is because these points get projected onto the same value when we perform PCA. For example, for the point at roughly $(-1, -2)$, the 1D projected nearest neighbor is the point at roughly $(-2, -1)$.

To have 100% error on the 2D data we want an alternating sequence of “+”’s and “-”’s along the two lines of points parallel to the PCA direction. If we use

the same sequence of “+”’s and “-”’s on both lines of points, then we will also have 0% error when using the 1D data. This is because the 1D-nearest neighbor points will have the same labels.

To have 0% error on the 2D data we want adjacent points on the two lines to all have the same label. If we label only line of points with “+”’s and the other with “-”’s, then the 1D data will also have 100% error. This is because the 1D nearest neighbors will have the opposite labels.

In summary, the labelling below has the specified error rates.



References