Tyler Olivieri   HW6  #1      Consider a classification problem with $N$ different classes.

let prior probability of class $n \in N$ be $\pi_n$

denote $f_n(x) = Pr(X=x \mid Y=n)$   observations in class $n$ are drawn from $N(\varphi_n, \Sigma_n)$ where $\Sigma_n = \Sigma$ $\forall n$.

a) Use bayes theorem to find $Pr(Y=n \mid X=x)$

$$Pr\ Pr(X,Y) = Pr(X,Y)$$

$$Pr(X \mid y=n) p(y=n) = Pr(y \mid X=x) Pr(X=x)$$

$$Pr(y \mid X=x) = \frac{P(X \mid y=n) p(y=n)}{Pr(X=x)}$$

$$\Rightarrow P(y=n \mid X=x) = \frac{P(X=x \mid y=n) p(y=n)}{P(X=x)}$$

$$= \frac{f_n(x) \pi_n}{P(X=x)}$$

$$= \frac{f_n(x) \pi_n}{P \sum_{i=1}^{k} P(X=x, y=i)}$$

$$= \frac{f_n(x) \pi_n}{\sum_{i=1}^{k} P(X=x \mid y=i) P(y=i)}$$

$$P(y=n \mid X=x) = \frac{f_n(x) \pi_n}{\sum^{k} P(X=x \mid y=i) \pi_i}$$

$$Pr\left(y=n\,|\,x=x\right) = \frac{f_n(x)\,\Pi_n}{\sum\limits_{j=1}^{F} f_j(x)\,\Pi_j}$$
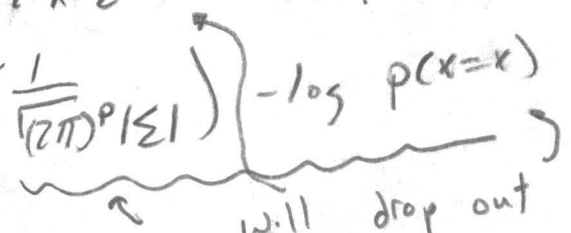
b) Derive the linear descriminant function, $f_n(x)$ and write the classification rule for the predicted class, $\hat{y}$ for an LDA in terms of $f_n(x)$.

$$f_n(x) = \frac{1}{\sqrt{(2\pi)^p |\Sigma|}}\exp\left(-\tfrac{1}{2}(x-\varphi_n)\Sigma^{-1}(x-\varphi_n)\right)$$

$$Pr\left(y=n\,|\,x=x\right) = \frac{\frac{1}{\sqrt{(2\pi)^p|\Sigma|}}\exp\left(-\tfrac{1}{2}(x-\mu_n)\Sigma^{-1}(x-\mu_n)\right)\Pi_n}{\sum\limits_{i=1}^{\phantom{F}} \frac{p(x=x)}{(2\pi)^p|\Sigma|}\exp\left(-\tfrac{1}{2}(x-\mu_n)\Sigma^{-1}(x-\mu_n)\right)}$$

$$\log Pr(y=n\,|\,x=x) = \frac{\log \Pi_n + \log\left[\frac{1}{\sqrt{(2\pi)^p|\Sigma|}}\exp\left(-\tfrac{1}{2}(x-\mu_n)^T\Sigma^{-1}(x-\mu_n)\right)\right]}{}$$

$$-\log p(x=x)$$

$$= \log \Pi_n + \log\left(\frac{1}{\sqrt{(2\pi)^p|\Sigma|}}\right) + \log\left(\exp\left(-\tfrac{1}{2}x^T\Sigma^{-1}x + x_n^T\Sigma^{-1}\mu_n\right.\right.$$
$$\left.\left. -\tfrac{1}{2}\mu_n^T\Sigma^{-1}\varphi_n\right)\right)$$

$$-\log p(x=x)$$

$$f_n(x) = \log Pr(y=n\,|\,x=x) = \log \Pi_n -\tfrac{1}{2}x^T\Sigma^{-1}x + x^T\Sigma^{-1}\mu_n -\tfrac{1}{2}\mu_n^T\Sigma^{-1}\mu_n$$
$$+ \log\left(\frac{1}{\sqrt{(2\pi)^p|\Sigma|}}\right)\underbrace{-\log p(x=x)}_{\text{will drop out}}$$

the classification rule should be the LDA in the case of

$$\hat{y} = \underset{y \in \{1,\dots,n\}}{\text{argmax}} \; \delta_n(x)$$

taking the max wrt $y$ # removes two constants in $\delta_n(x)$.

$$\hat{y} = \underset{y \in \{1,\dots,n\}}{\text{argmax}} \left( \log \pi_n + x^T \Sigma^{-1} \mu_n - \tfrac{1}{2} \mu_n^T \Sigma^{-1} \mu_n \right)$$

where we estimate $\pi_n, \mu_n$ with $\hat{\pi}_n, \hat{\mu}_n$

c) Derive the decision boundary for the LDA in the case of two classes, $a$ and $b$.

ÿ The decision boundary for two classes $a, b$ is the set of points where

$$\delta_a(x) = \delta_b(x)$$

$$\log \pi_a + x^T \Sigma^{-1} \mu_a - \tfrac{1}{2}\mu_a^T \Sigma^{-1} \mu_a = \log \pi_b + x^T \Sigma^{-1} \mu_b - \tfrac{1}{2}\mu_b^T \Sigma^{-1} \mu_b$$

This is an equation of a line in $x$.

(substitue estimates for $\pi_a, \pi_b, \mu_a, \mu_b$)

$$\hat{\pi}_a, \hat{\pi}_b, \hat{\mu}_a, \hat{\mu}_b$$

d) Consider two classes, $a=1$ and $b=2$. You are

given $\hat{\pi}_a = .6$, $\hat{\pi}_b = .4$  $\hat{\mu}_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$, $\hat{\mu}_2 = \begin{bmatrix} -2 \\ 1 \end{bmatrix}$ and

$$\hat{\Sigma} = \begin{bmatrix} 1.5 & .1 \\ .1 & 1 \end{bmatrix}$$

Find the decision boundary and classification rule of the

corresponding LDA. How would the observation $x = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ be classified?

The decision boundary are the points $x$ s.t

$$\log(.6) + x^T \begin{bmatrix} 1.5 & .1 \\ .1 & 1 \end{bmatrix}^{-1} \begin{bmatrix} 1 \\ 0 \end{bmatrix} - \frac{1}{2} \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} 1.5 & .1 \\ .1 & 1 \end{bmatrix}^{-1} \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$= \log(.4) + x^T \begin{bmatrix} 1.5 & .1 \\ .1 & 1 \end{bmatrix}^{-1} \begin{bmatrix} -2 \\ 1 \end{bmatrix} - \frac{1}{2} \begin{bmatrix} -2 & 1 \end{bmatrix} \begin{bmatrix} 1.5 & .1 \\ .1 & 1 \end{bmatrix}^{-1} \begin{bmatrix} -2 \\ 1 \end{bmatrix}$$

$$-.51 + x^T \begin{bmatrix} .71 \\ -.07 \end{bmatrix} - .36 = -.92 + x^T \begin{bmatrix} -2.14 \\ 1.21 \end{bmatrix} - 2.75$$

$$-.87 + x^T \begin{bmatrix} .71 \\ -.07 \end{bmatrix} = -3.67 + x^T \begin{bmatrix} -2.14 \\ 1.21 \end{bmatrix}$$

classify $x$ as class $a$ if

$$-.87 + x^T \begin{bmatrix} .71 \\ -.07 \end{bmatrix} > -3.67 + x^T \begin{bmatrix} -2.14 \\ 1.21 \end{bmatrix}$$

and $b$ otherwise.

$$-.87 + \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} .71 \\ -.07 \end{bmatrix} = -.23$$

$$-3.67 + \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} -2.14 \\ 1.21 \end{bmatrix} = -4.6$$

since $-.23 > -4.6$ $\Rightarrow$ classify $x$ as $a$.

e) This relates to logistic regression under the binary case in the following:

with two classes the decision boundary for gaussian bayes is when $\log p(y=1|x) = \log p(y=0|x)$

so choose class 1 when $\log p(y=1|x) - \log p(y=0|x) > 0$

$\not\equiv \log \frac{p(y=1|x)}{p(y=0|x)} > 0$   whic

the quantity $\log \frac{p(y=1|x)}{p(y=0|x)}$ is the log odds

which is the basis of logistic regression

c) The naive Bayes classifier is similar to LDA, except it assumes each predictor is conditionally independent of every other predictor given class $n$. Derive the classification rule for $\hat{y}$ under this classifier. How does this relate to logistic regression in the binary case (i.e for two classes?)

$$p(X|y) = \prod_{i=1}^{p} p(x_i|y) \qquad \text{under naive bayes conditionally independent assumption.}$$

$$p(y=n|x=x) = \frac{f_n(x)\,\Pi_n}{p(x=x)} = \frac{\prod_{i=1}^{p} p(x_i|y=n)\,\Pi_n}{p(x=x)}$$

$$\log p(y=n|X=x) = \log\left(\prod_{i=1}^{p} p(x_i|y=n)\right) + \log\Pi_n - \log p(x=x)$$

$$= \sum_{i=1}^{p} \left(\log p(x_i|y=n)\right) + \log\Pi_n + \log p(x=x)$$

$$= \sum_{i=1}^{p} \log \frac{1}{\sqrt{(2\pi)^p}\,\sigma_n^2}\exp\left(-\frac{1}{2\sigma_n^2}\left(x_{i}-\mu_n\right)^2\right) + \underbrace{\log\Pi_n + \log p(x=x)}_{\text{constant wrt } y.}$$

$$\hat{y} = \underset{y\in\{1,2,\ldots,n\}}{\arg\max}\ \log p(y=n|x=x) \qquad \text{is classification rule}$$

Substitute above expression except for $\log p(x=x)$ as it does not effect maximization wrt $y$. Substitue estimates for necessary parameters.

Tyler Olivieri   HW6   #2   k-means clustering

a) Show that setting the objective function the sum of the squared
Euclidean distances of points from the center of their clusters

$$obj = \sum_{k=1}^{K} \sum_{x \in C_k} \sum_{i=1}^{P} (C_{ki} - x_i)^2$$

results in an update rule where the optimal centroid is the
mean of the points in the cluster.

min obj wrt $C_k$

$$\frac{d\,obj}{dc_k} = \sum_{k=1}^{K} \sum_{x \in C_k} \sum_{i=1}^{P} 2(C_{ki} - x_i) = 0$$

$$2 \sum_{k=1}^{K} \sum_{x \in C_k} \sum_{i=1}^{P} (C_{ki} - x_i) = 0$$

$$\sum_{k=1}^{K} \sum_{x \in C_k} \sum_{i=1}^{P} (C_{ki} - x_i) = 0$$

$$\sum_{k=1}^{K} \sum_{x \in C_k} \sum_{i=1}^{P} C_{ki} - \sum_{k} \sum_{x \in C_k} \sum_{i=1}^{P} x_i = 0$$

$$\sum_{k} \sum_{x \in C_k} \sum_{i=1}^{P} C_{ki} = \sum_{k} \sum_{x \in C_k} \sum_{i=1}^{P} x_i$$

let $n_k$ be
# points assigned
to cluster k.
$x \in C_k$

$$n_k \sum_{x \in C_k} \sum_{i=1}^{P} C_{ki} = \sum_{k} \sum_{x \in C_k} \sum_{i=1}^{P} x_i$$

$$\frac{1}{p}\sum_{i=1}^{p} c_{ki}\Bigg) = \frac{\sum_{x \in C_k}\left(\sum_{i=1}^{p} x_i\right)}{n_k} \xleftarrow{} \text{total } \#\text{'s of points} \atop \text{assigned to cluster } k. = \frac{\sum_{i=1}^{p}\sum_{x \in C_k} x_i}{n_k}$$

So set $c_{ki}$ to be the mean of all data points assigned to cluster $k$.

$$c_{ki} = \frac{\sum_{x \in C_k} x_i}{n_k}$$

ignoring $p$ - feature notation

for each feature, set cluster center to be mean of data points assigned to the cluster centr.

b) Show that setting the objective function to the sum of the manhattan distances of points from the center of their clusters,

$$obj = \sum_{k=1}^{k} \sum_{x \in C_k} \sum_{i=1}^{P} |C_{ki} - X_i|$$

results in an update rule where the optimal centroid is the median of the cluster.

$$\frac{d\,obj}{dC_K} = \sum_{x \in C_k} \sum_{i=1}^{P} \frac{C_{ki} - X_i}{|C_{ki} - X_i|} = 0$$

when $\sum_{x \in C_k} C_{ki} - X_i > 0$ then $\frac{C_{ki} - X_i}{|C_{ki} - X_i|} = \frac{C_{ki} - X_i}{C_{ki} - X_i} = 1$

when $C_{ki} - X_i < 0$ then $\frac{C_{ki} - X_i}{|C_{ki} - X_i|} = \frac{C_{ki} - X_i}{X_i - C_{ki}} = \frac{-(C_{ki} - X_i)}{-(C_{ki} - X_i)}$

$$= \frac{1}{-1} = -1$$

$$\Rightarrow \frac{C_{ki} - X_i}{|C_{ki} - X_i|} = sign(C_{ki} - X_i)$$ where $sign(x) = 1$ when $x > 0$

$$= -1 \text{ when } x < 0$$

$$\sum_{\substack{x \in \\ i=1}}^{P} \left( \sum_{x \in C_k} sign(C_{ki} - X_i) \right) = 0$$

This can only equal zero when the number of the positive elements equals the number of negative elements. So for $C_{ki}$

$$C_{ki} = median(X_i \in C_{ki})$$ then

So $C_k = median(X_n)$

Consider the dataset

| X | y |
|---|---|
| 0 | 1 |
| 1 | 1 |
| 2 | 1 |
| 2 | 3 |
| 3 | 2 |
| 3 | 3 |
| 4 | 5 |

a) Normalize the data and derive the two principal components in sorted order.

$$\mu_x = \frac{0+1+2+2+3+3+4}{7} = 2.14$$

$$\mu_y = \frac{1+1+1+3+2+3+5}{7} = 2.29$$

$$\sigma_x^2 = \frac{1}{n-1}\sum_{i=1}^{n}(x_i - \mu_x)^2 = \frac{1}{6}\sum_{i=1}^{7}(x_i - 2.14)^2 = 1.55$$

$$\sigma_x = \sqrt{\sigma_x^2} = \sqrt{1.55} = 1.24$$

$$\sigma_y^2 = \frac{1}{n-1}\sum_{i=1}^{n}(y_i - \mu_y)^2 = \frac{1}{6}\sum_{i=1}^{7}(y_i - 2.29)^2 = 1.92$$

$$\sigma_y = \sqrt{\sigma_y^2} = \sqrt{1.92} = 1.38$$

Normalized data $\quad X_i := \dfrac{X_i - \mu_x}{\sigma_x} \qquad y_i := \dfrac{y_i - \mu_y}{\sigma_y}$

0 mean
std-dev = 1

X-norm $\approx (0,1)$
Y-norm " "

| $X_{norm}$ | y norm |
|------------|--------|
| -1.72 | -.93 |
| -.92 | -.93 |
| -.11 | -.93 |
| -.11 | .52 |
| .69 | -.21 |
| .69 | .52 |

let $X =$ $\begin{bmatrix} -1.72 & -.93 \\ -.92 & -.93 \\ -.11 & -.93 \\ -.11 & .52 \\ .69 & .21 \\ .69 & .52 \\ 1.49 & 1.96 \end{bmatrix}$ be the data matrix.

$A =$

find principal components

$$T_1 = AW \quad \text{where the columns of } W \text{ are}$$
$$= U\Sigma W^T W \quad \text{the eigenvectors of } A^T A.$$

(by SVD)

since $W$ is chosen to be orthonormal $W^T W = I$

$$= U\Sigma$$

SVD of $A$: (Done in python)

$$A = \begin{bmatrix} -.52 & -.54 \\ -.36 & .06 \\ -.2 & .56 \\ .08 & -.42 \\ .18 & .33 \\ .24 & -.11 \\ .68 & -.32 \end{bmatrix} \qquad \overset{\Sigma}{\underset{||}{ }} \begin{bmatrix} 3.6 & 0 \\ 0 & 1.04 \end{bmatrix} \qquad \overset{W^T}{\underset{||}{ }} \begin{bmatrix} .71 & .71 \\ .71 & -.71 \end{bmatrix}$$

Now we have $T = U\Sigma = \begin{bmatrix} -.52 & -.54 \\ -.36 & .06 \\ -.2 & .56 \\ -.08 & -.42 \\ .18 & .33 \\ .24 & .11 \end{bmatrix} \begin{bmatrix} 3.6 & 0 \\ 0 & 1.04 \end{bmatrix}$

$$T = U\Sigma = \begin{bmatrix} -1.88 & -.56 \\ -1.31 & .06 \\ -.74 & .58 \\ -.52 & -.44 \\ .29 & \\ .64 & .34 \\ .86 & .12 \\ -2.44 & -.33 \end{bmatrix}$$

The new transformed dataset using the first principal component is

$$\tilde{X} = \begin{bmatrix} -1.88 \\ -1.31 \\ -.74 \\ .29 \\ .64 \\ -.86 \\ 2.44 \end{bmatrix}$$

Which is the first column of T.

b) repeat the previous analysis but do not normalize the data.

Is pca scale-invariant

SVD of A_unnormalized: (removing n-r columns of U, and 0's rows and columns of Σ)

when you multiply UΣ, they are irrelevant. because python libs give SVD this way except numpy, which I did not use

$$A_{unnormalized} = \begin{bmatrix} .08 & .44 \\ .15 & -.04 \\ .22 & -.51 \\ .37 & .37 \\ .37 & -.55 \\ .45 & -.11 \\ .67 & .30 \end{bmatrix} \quad \begin{bmatrix} 9.52 & 0 \\ 0 & 1.54 \end{bmatrix} \begin{bmatrix} .68 & .73 \\ -.73 & .68 \end{bmatrix}$$

(U under first matrix)  (Σ under second)  ($V^T$ under third)

$$T_{unnormalized} = U\Sigma = \begin{bmatrix} .73 & .68 \\ 1.41 & -.06 \\ 2.09 & -.79 \\ 3.60 & .56 \\ 3.50 & -.85 \\ 4.24 & -.17 \\ 6.39 & .46 \end{bmatrix}$$

$$\tilde{X}_{unnormalized} = \begin{bmatrix} .73 \\ 1.41 \\ 2.09 \\ 3.60 \\ 3.50 \\ 4.24 \\ 6.39 \end{bmatrix}$$

⇒ PCA is not scale-invariant   $\tilde{X} \neq \tilde{X}_{unnormalized}$

$$T = U\Sigma = \begin{bmatrix} -1.88 & -.56 \\ -1.31 & .06 \\ -.74 & .58 \\ .29 & -.44 \\ .64 & .34 \\ .86 & .12 \\ -2.44 & -.33 \end{bmatrix}$$

The new transformed dataset using the first principal component is

$$\check{X} = \begin{bmatrix} -1.88 \\ -1.31 \\ -.74 \\ .29 \\ .64 \\ -.86 \\ 2.44 \end{bmatrix}$$

Which is the first column of T.