

Problem 1

Wednesday, 22 January 2020 14:22

Problem 1: Use the SVD shown below. Suppose a new user Leslie assigns rating 3 to Alien and

| | Joe | Jim | John | Jack | Jill | Jenny | Jane | Star Wars | Alien | Casablanca | Titanic |
|--------|-----|-----|------|------|------|-------|------|-----------|-------|------------|---------|
| Matrix | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 |
| Joe | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 |
| Jim | 3 | 3 | 3 | 0 | 0 | 0 | 0 | 4 | 4 | 4 | 0 |
| John | 4 | 4 | 4 | 0 | 0 | 0 | 0 | 5 | 5 | 5 | 0 |
| Jack | 5 | 5 | 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 4 |
| Jill | 0 | 0 | 0 | 4 | 4 | 4 | 4 | 0 | 0 | 0 | 5 |
| Jenny | 0 | 0 | 0 | 5 | 5 | 5 | 5 | 0 | 0 | 0 | 2 |
| Jane | 0 | 0 | 0 | 2 | 2 | 2 | 2 | 0 | 0 | 0 | 2 |

Figure 11.6: Ratings of movies by users

$$M = U \Sigma V^T$$

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 0 & 0 & 2 & 2 \end{bmatrix} = \begin{bmatrix} .14 & 0 \\ .42 & 0 \\ .56 & 0 \\ .70 & 0 \\ 0 & .60 \\ 0 & .75 \\ 0 & .30 \end{bmatrix} \begin{bmatrix} 12.4 & 0 \\ 0 & 9.5 \end{bmatrix} \begin{bmatrix} .58 & .58 & .58 & 0 & 0 \\ 0 & 0 & 0 & .71 & .71 \end{bmatrix}$$

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?

$$P = M \cdot V = U \cdot \Sigma \cdot V$$

$$[0, 3, 0, 0, 4] \cdot V = [7.74, 2.84]$$

Look at V: sci-fi "classic"

\Rightarrow preference for classic

Problem 2: Consider the latent space distribution

$$p(z) = \mathcal{N}(z | 0, I)$$

and a conditional distribution for the observed variable $x \in \mathbb{R}^d$,

$$p(x|z) = \mathcal{N}(x | Wz + \mu, \Phi)$$

where Φ is an arbitrary symmetric, positive-definite noise covariance variable. Now suppose that we make a nonsingular linear transformation of the data variables $y = Ax$ where A is a non-singular $d \times d$ matrix. If μ_{MLE} , W_{MLE} , and Φ_{MLE} represent the maximum likelihood solution corresponding to the original untransformed data, show that $A\mu_{MLE}$, AW_{MLE} , and $A\Phi_{MLE}A^T$ will represent the corresponding maximum likelihood solution for the transformed data set. Finally, show that the form of the model is preserved if A is orthogonal and Φ is proportional to the unit matrix so $\Phi = \sigma^2 I$ (i.e. probabilistic PCA). The transformed Φ 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.

$y = Ax$ with A regular $\hat{=}$ noiseless transformation

$$p(x) = \mathcal{N}(x | \mu, Ww^T + \Phi)$$

Bishop: $p(x') = \mathcal{N}(x' | \mu', A^{-1})$

$$p(y' | x') = \mathcal{N}(y' | A'x' + b, L^{-1})$$

$$\Rightarrow p(y') = \mathcal{N}(y' | A'\mu' + b, L^{-1} + A'A^{-1})$$

$$p(y | x) = \mathcal{N}(y | Ax, 0)$$

$$p(y|x) = \mathcal{N}(y|Ax, 0)$$

$$y^i = y_i, x^i = x_i, \mu^i = \mu_i, \Sigma^{-1} = WW^T + \phi, A^i = A, b = 0, L^{-1} = 0$$

$$p(y) = \mathcal{N}(y|A\mu, A(WW^T + \phi)A^T) =$$

$$= \mathcal{N}(y|A\mu, AWW^TA^T + A\phi A^T)$$

previously: $\mathcal{N}(x|\mu, WW^T + \phi)$

$$\text{MLE: } \mu_{Y,\text{MLE}} = A\mu_{nLE}$$

$$W_{Y,\text{MLE}} = AW_{nLE}$$

$$\boxed{\phi_{Y,\text{MLE}} = A\phi_{nLE} A^T}$$

$$\begin{aligned} & \mathcal{N}(y|A\mu_{nLE}, AW_{nLE}W_{nLE}^TA^T + A\phi_{nLE}A^T) \\ & \mathcal{N}(y|\mu_{Y,\text{MLE}}, W_{Y,\text{MLE}}W_{Y,\text{MLE}}^T + \phi_{Y,\text{MLE}}) \end{aligned}$$

$$p_Y(As) = p_X(s) \quad s \in \mathbb{R}^d$$

$$\begin{aligned} A \text{ orthogonal: } \phi_{Y,\text{MLE}} &= A\phi_{nLE}A^T = A\sigma^2 I A^T = \\ &= \sigma^2 AA^T = \sigma^2 I = \phi_{nLE} \end{aligned}$$

$$p(y|z) = \mathcal{N}(y|Aw(z) + A\mu, \sigma^2 I) \quad \text{(isotropic)}$$

$$\text{isotropic: } p(Rz) = p(z)$$

Problem 3: Let the matrix $X \in \mathbb{R}^{N \times D}$ represent N data points of dimension $D = 10$ (samples stored as rows). We applied PCA to X . By using the $K = 2$ top principal components, we transformed/projected X into $\tilde{X} \in \mathbb{R}^{N \times K}$. We computed that \tilde{X} preserves 70% of the variance of the original data X .

Suppose now we apply PCA on the following matrices:

- a) $Y_1 = XS$ where $S = \lambda I$, with $\lambda \in \mathbb{R}$ and $I \in \mathbb{R}^{D \times D}$ is the identity matrix
- b) $Y_2 = XR$ where $R \in \mathbb{R}^{D \times D}$ and $RR^T = I$
- c) $Y_3 = XP$ where $P = \text{diag}(+5, -5, \dots, +5, -5)$ is a $D \times D$ diagonal matrix
- d) $Y_4 = XQ$ where $Q = \text{diag}(1, 2, 3, \dots, D-1, D)$ is a $D \times D$ diagonal matrix
- e) $Y_5 = X + 1_N\mu^T$ where $\mu \in \mathbb{R}^D$ and 1_N is an N -dimensional column vector of all ones
- f) $Y_6 = XA$ where $A \in \mathbb{R}^{D \times D}$ and $\text{rank}(A) = 5$

and obtain the projected data $\tilde{Y}_1, \dots, \tilde{Y}_6 \in \mathbb{R}^{N \times K}$ using the principal components corresponding to the top $K = 5$ largest eigenvalues of the respective Y_i .

What fraction of variance of each Y_i will be preserved by each respective \tilde{Y}_i ? Justify your answer.

The answer "cannot tell without additional information" is also valid if you provide a justification.

$$\alpha) X^T X = U \Sigma V^T \quad (\text{wlog } X \text{ centered, ignore } \frac{1}{n})$$

$$S^T X^T S = U \Sigma V^T \cdot \lambda^2 I = U \lambda^2 \Sigma V^T$$

$$\text{Variance: } \text{tr}(\Sigma), \quad \text{tr}(\Sigma) = \sum_i \lambda_i = \sum_{i=1}^5 \lambda_i / 0.7$$

$$\Sigma' = \lambda^2 \Sigma = \sum_i \lambda^2 \lambda_i = \lambda^2 \sum_{i=1}^5 \lambda_i / 0.7$$

$$\Sigma' = \sum_i \lambda_i' = \sum_i \lambda \lambda_i = \lambda^2 \sum_{i=1}^5 \lambda_i / 0.7$$

$$\beta) Y_2^T Y_2 = R^T X^T X R = \underbrace{R^T}_{V} \underbrace{U \Sigma V^T}_{\Sigma'} \underbrace{R}_{V^T} = V \Sigma V^T$$

$$VV^T = I_n$$

$$\gamma) P = S \cdot D(1, -1, 1, \dots) = 5 \cdot D_{5,1}$$

$$\begin{pmatrix} 1 & -1 & 1 & \dots \end{pmatrix} \Rightarrow \text{orthogonal}$$

$$\overbrace{\begin{pmatrix} 1 & & \\ & 1 & \\ & & \ddots \end{pmatrix}}^S \Rightarrow \text{orthogonal}$$

$$Y_3^T Y_3 = S \cdot D_{1,1} X^T \times D_{1,1} \cdot S = \\ = 2S \cdot \underbrace{D_{1,1} U}_{V^T} \underbrace{\Sigma V^T D_{1,1}}_{V^T} = 2S \cdot V^T \Sigma V^T$$

d) $Q = \begin{pmatrix} 1 & & 0 \\ 0 & 2 & 0 \\ 0 & 0 & \ddots \end{pmatrix}$ Different scale in each original direction

$$Q X^T X Q = \\ = Q^T U \Sigma V^T Q \\ \cancel{ZV \text{ not orthogonal!}}$$

e) $Y_S = X + \vec{1} \vec{\mu}^T$

PCA: Eigendecomp of $(X - \vec{1} \vec{\mu}^T)^T (X - \vec{1} \vec{\mu}^T)$
Subtract mean \Rightarrow no effect

f) $Y_G = X \vec{A}$
 $\text{rank}(S) \Rightarrow X \cdot \vec{A}$ in S -dim space
 $\Rightarrow \vec{A}^T X^T X \vec{A} = V^T \underbrace{\Sigma^T}_{\begin{pmatrix} \sigma_1 & & \\ \sigma_2 & \ddots & \\ & \ddots & 0 \end{pmatrix}} V^T$

$\Rightarrow 100\% \text{ of variance}$

Problem 4: You are given $N = 4$ data points: $\{x_i\}_{i=1}^4, x_i \in \mathbb{R}^3$, represented with the matrix $X \in \mathbb{R}^{4 \times 3}$.

$$X = \begin{bmatrix} 4 & 3 & 2 \\ 2 & 1 & -2 \\ 4 & -1 & 2 \\ -2 & 1 & 2 \end{bmatrix}$$

Hint: In this task the results of all (final and intermediate) computations happen to be integers.

- Perform principal component analysis (PCA) of the data X , i.e. find the principal components and their associated variances in the transformed coordinate system. Show your work.
- Project the data to two dimensions, i.e. write down the transformed data matrix $Y \in \mathbb{R}^{4 \times 2}$ using the top-2 principal components you computed in (a). What fraction of variance of X is preserved by Y ?
- Let $x_5 \in \mathbb{R}^3$ be a new data point. Specify the vector x_5 such that performing PCA on the data including the new data point $\{x_i\}_{i=1}^5$ leads to exactly the same principal components as in (a).

In-class Exercises

Wednesday, 22 January 2020 14:21

$$x_i \in \mathbb{R}^D \quad i=1..N$$

find $z_i \in \{1..K\}$ s.t. $z_i = z_j \Rightarrow x_i$ similar to x_j

K centroids $\mu_k \in \mathbb{R}^D$

$$\mathcal{J}(Z, \mu) = \sum_{i=1}^N \sum_{k=1}^K z_{ik} \|x_i - \mu_k\|_2^2$$

(Lloyd's algorithm

1. initialize μ_k

2. find $z = \min_z \mathcal{J}(Z, \mu)$

3. find $\mu = \min_\mu \mathcal{J}(Z, \mu)$

4. iterate 2+3 until convergence

Problem 6: Consider a modified version of the K -means objective, where we use L_1 distance instead.

$$\mathcal{J}(X, Z, \mu) = \sum_{i=1}^N \sum_{k=1}^K z_{ik} \|x_i - \mu_k\|_1$$

This variation of the algorithm is called K -medians. Derive the Lloyd's algorithm for this model.

Step 2

$$z_{ik} = \begin{cases} 1 & \text{if } k = \arg \min_j \|x_i - \mu_j\|_1 \\ 0 & \text{else} \end{cases}$$

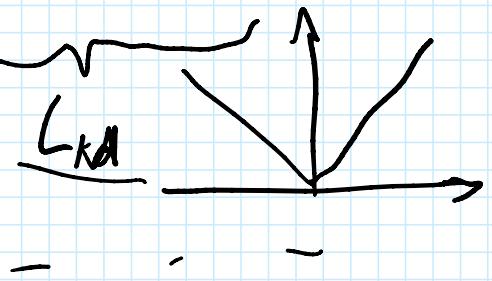
Step 3

$$\mu_k = \arg \min_\mu \sum_{i=1}^N z_{ik} \|x_i - \mu_k\|_1$$

$$\mathcal{J} = \sum_{i=1}^N \sum_{k=1}^K z_{ik} \sum_{d=1}^D |x_{id} - \mu_{kd}| = \sum_{k=1}^K \sum_{d=1}^D \sum_{i=1}^N z_{ik} |x_{id} - \mu_{kd}|$$

$$J = \sum_{i=1}^N \sum_{k=1}^K z_{ik} \sum_{d=1}^D |x_{id} - \mu_{kd}| = \sum_{k=1}^K \sum_{d=1}^D \sum_{i=1}^N z_{ik} |x_{id} - \mu_{kd}|$$

$$\frac{\partial L_{kd}}{\partial \mu_{kd}} = \sum_{i=1}^N z_{ik} \frac{\partial}{\partial \mu_{kd}} |x_{id} - \mu_{kd}| =$$



$$\frac{\partial |x_{id} - \mu_{kd}|}{\partial \mu_{kd}} = \begin{cases} -1 & \text{if } \mu_{kd} > x_{id} \\ 1 & \text{if } \mu_{kd} < x_{id} \\ 0 & \text{if } \mu_{kd} = x_{id} \end{cases}$$

$$= \sum_{i=1}^N z_{ik} [\mu_{kd} > x_{id}] + \sum_{i=1}^N z_{ik} [\mu_{kd} < x_{id}] = 0$$

A = B

$$\frac{\partial L_{kd}}{\partial \mu_{kd}} = 0 \Leftrightarrow \mu_{kd} = \text{median} \{ x_{id} \mid z_{ik} = 1 \}$$