

P22 Background

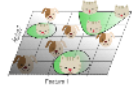
Overfitting \rightarrow ① Add data ② Regularization ③ Dimensionality reduction

DR \rightarrow ① Feature Selection ② Linear DR (PCA, MPS) ③ Nonlinear DR

Curse of dimensionality: Cause overfitting and data sparsity

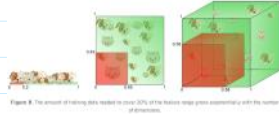
Increasing the dimensionality may lead to better result. But at the same time, it will make number of samples per unit smaller. (Data Sparsity ①)

Obtaining a better classifier by increasing dimensionality is equivalent to using a more complex nonlinear classifier in low-dimensional space

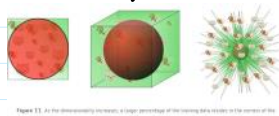


(Data Sparsity ②): As the dimensionality increases, more data is needed for each dimension to cover the proportion of space

This means that increasing the dimensionality also requires exponentially increasing the data



Data Sparsity ③: (Data is distributed differently in high-dimensional space) Assuming that the mean of the eigenvalues is taken as the centroid of the space in two dimensions. Make a circle with the feature range as the diameter, and the samples outside the circle are distributed in the corners of the space.



The characteristics of the data in the corners are widely disparate. These data is more difficult to classify.

$V_{\text{residual}} = 1$ $V_{n\text{-sphere}} = k \cdot 2(0.5)^D$ $D \rightarrow \infty \therefore V_{n\text{-sphere}} \rightarrow 0 \therefore$ Almost all of the samples are distributed in the corners.

dist: Euclidean distance from the sample point to the center of the space. $\lim_{D \rightarrow \infty} \frac{\text{dist}_{\text{max}} - \text{dist}_{\text{min}}}{\text{dist}_{\text{max}}} = 0 \therefore$ All samples in high-dimensional space are far from the center.

\therefore No difference in maximum and minimum distance

\therefore Euclidean, Manhattan, Mahalanobis distance and other methods gradually fail in high-dimensional space.

P23 Sample Mean & Variance

$$\bar{X} = \frac{1}{N} \sum_{i=1}^N X^i = \frac{1}{N} (X^1, X^2, \dots, X^N) \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}_{N \times 1} = \frac{1}{N} X^T \mathbf{1}_N$$

$$S = \frac{1}{N} \sum_{i=1}^N (X^i - \bar{X})(X^i - \bar{X})^T \Rightarrow \frac{1}{N} \sum_{i=1}^N (X^i - \bar{X}) (X^i - \bar{X})^T = X^T - \bar{X} \mathbf{1}_N^T = X^T (I_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T)$$

$$\text{Let } I_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T = H_{\text{center}} \Rightarrow \text{Centering Matrix} \therefore S = \frac{1}{N} X^T H H^T X \quad H = \begin{bmatrix} \frac{N-1}{N} & -\frac{1}{N} & \dots & -\frac{1}{N} \\ -\frac{1}{N} & \frac{N-1}{N} & \dots & -\frac{1}{N} \\ \vdots & \vdots & \ddots & \vdots \\ -\frac{1}{N} & -\frac{1}{N} & \dots & \frac{N-1}{N} \end{bmatrix}$$

$$\therefore H^T = H \quad H^N = H \quad \therefore S = \frac{1}{N} X^T H X$$

P24 PCA - Maximum Projection Variance

PCA main task: Reconstructing the original eigenspace

Two methods: ① Minimum Projection Variance ② Minimum Reconfiguration Cost

①: Project the data onto linearly independent basis and make the variance of the projected point set as small as possible.

$\vec{a} \cdot \vec{b} = |\vec{a}| |\vec{b}| \cos \theta \therefore \|\vec{u}\| = 1 \quad \vec{a}^T \vec{u} = \|\vec{a}\| \cos \theta$, which is the projection of \vec{a} on \vec{u}

First make X zero-centered $\Rightarrow X^i - \bar{X} \Rightarrow E(X - \bar{X}) = 0 \Rightarrow E[(X - \bar{X})^T] = E[(X - \bar{X})^T] E[u] = 0$

$$J = D[(X - \bar{X})^T u] = E[(X - \bar{X})^T u]^2 = \frac{1}{N} \sum_{i=1}^N [(X^i - \bar{X})^T u]^2 \quad \text{s.t. } u^T u = 1$$

$$= u^T \left[\frac{1}{N} \sum_{i=1}^N (X^i - \bar{X})(X^i - \bar{X})^T \right] u = u^T S u$$

$$\therefore \begin{cases} \hat{u} = \arg \max_u u^T S u \\ \text{s.t. } u^T u = 1 \end{cases} \quad \mathcal{L}(u, \lambda) = u^T S u + \lambda(1 - u^T u)$$

$$\frac{\partial \mathcal{L}}{\partial u} = 2Su - 2\lambda u = 0$$

$$\therefore S\hat{u} = \lambda \hat{u} \quad \therefore \text{PCA is to calculate the eigenvectors of } S \text{ and sort them by eigenvalues.}$$

And select the first q as required.

P25 Minimum Reconfiguration Cost

$$x^i = (\underbrace{x^{i^T} u^1}_{\text{The coordinates of } x^i \text{ on } u^1}) u^1 + (\underbrace{x^{i^T} u^2}_{\text{The coordinates of } x^i \text{ on } u^2}) u^2 \quad u^{1^T} u^1 = 1$$

$$\text{Original } p \text{ dimensions: } x^i = \sum_{k=1}^p [(x^i - \bar{x})^T u^k] u^k \quad \text{Left } q \text{ dimensions: } \hat{x}^i = \sum_{k=1}^q [(x^i - \bar{x})^T u^k] u^k$$

$$\textcircled{2} \text{ Reconfiguration cost: } J = \frac{1}{N} \sum_{i=1}^N \|x^i - \hat{x}^i\|^2 = \frac{1}{N} \sum_{i=1}^N \left\| \sum_{k=q+1}^p [(x^i - \bar{x})^T u^k] u^k \right\|^2 = \frac{1}{N} \sum_{i=1}^N \sum_{k=q+1}^p [(x^i - \bar{x})^T u^k]^2$$

$$\begin{cases} u^k = \arg \min_{u^k} \sum_{i=1}^N [(x^i - \bar{x})^T u^k]^2 \\ \text{s.t. } u^{k^T} u^k = 1 \end{cases} \quad \therefore \text{Linearly independent among } u^1, u^2, \dots, u^q \therefore \text{Solve optimization problems one by one.}$$

① find the q principal components, while ② find the $p-q$ principal components to be deleted

P26 Singular Value Decomposition (SVD) and PCA & PCoA

$$\textcircled{1} H_N = I_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T \quad HX \text{ (centralized)} = U \Sigma V^T \Rightarrow \text{SVD} \begin{cases} U^T U = I \\ V^T V = I \\ \Sigma \text{ diagonal matrix} \end{cases} \quad (V \text{ is eigenvector})$$

$$\textcircled{2} S_{\text{PCoA}} = \frac{1}{N} X^T H X = \frac{1}{N} X^T H^T H X = \frac{1}{N} V \Sigma^T U^T U \Sigma V^T = \frac{1}{N} V \Sigma^2 V^T \therefore S \text{ can be calculated from the SVD of the centralized data.}$$

$$\textcircled{3} \text{ at } T_{\dots} = HX X^T H^T = U \Sigma V^T V \Sigma U^T = U \Sigma^2 U^T$$

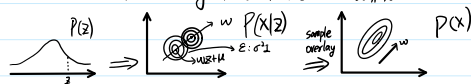
- ① $H_N = I_N - \frac{1}{N} \mathbf{1} \mathbf{1}^T$ H_X (centralized) $= U \Sigma V^T \Rightarrow \text{SVD} \begin{cases} V^T V = V V^T = I \\ \Sigma \text{ diagonal matrix} \end{cases}$ (V is eigenvector)
- ② $S_{PP} = \frac{1}{N} X^T H_X X = \frac{1}{N} X^T H^T H X = \frac{1}{N} V \Sigma U^T U \Sigma V^T = \frac{1}{N} V \Sigma^2 V^T \therefore S$ can be calculated from the SVD of the centralized data.
- ③ Let $T_{NN} = H_X X^T H^T = U \Sigma V^T V \Sigma U^T = U \Sigma^+ U^T$
- S: Calculate the direction (principal component V), and then calculate the coordinates under the new stage b $H_X \cdot V \leftarrow \text{PCA}$
- T: Find the coordinates directly $\Rightarrow H_X \cdot V = U \Sigma V^T V = U \Sigma \leftarrow$ (principle coordinate analysis PCA)
- $\therefore U^T U \Sigma = U \Sigma^+ U^T U \Sigma = U \Sigma \Sigma^+ \therefore U \Sigma$ is coordinate matrix
- eigenvector eigenvalue

When the data dimension is too large (PXP is too large), PCoA (T) is better

P27 Probabilistic Perspective PCA (P-PCA)

$X \in \mathbb{R}^P$ $Z \in \mathbb{R}^q$ $q < P$ Let $Z \sim N(0_q, I_q)$ $X = WZ + \mu + \epsilon$ $X \sim \text{CPA}$

$\epsilon \sim N(0, \sigma^2 I_P)$ P-PCA $\begin{cases} \text{Inference: } P(Z|X) \\ \text{Learning: } w, \mu, \sigma = \underset{w, \mu, \sigma}{\text{argmax}} \log P(Z|X) \end{cases}$ No calculation here



Each Z determines a Gaussian distribution with $wz + \mu$ as expectation and $E(\sigma^2 I)$ as variance.

$E(X|Z) = E(WZ + \mu + \epsilon) \therefore \text{find } X|Z \therefore Z \text{ is constant now} \therefore E(WZ + \mu + \epsilon) = WZ + \mu$

$\text{Var}(X|Z) = \text{Var}(WZ + \mu + \epsilon) = \sigma^2 I \quad X|Z \sim N(WZ + \mu, \sigma^2 I)$

$E(X) = E(WZ + \mu + \epsilon) = \mu \quad \text{Var}(X) = \text{Var}(WZ + \mu + \epsilon) = W I W^T + \sigma^2 I = W W^T + \sigma^2 I$

Review P7: $X = \begin{pmatrix} X_a \\ X_b \end{pmatrix} \quad X \sim N \left(\begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix}, \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix} \right)$

$\begin{cases} X_{b|a} = X_b - \Sigma_{ba} \Sigma_{aa}^{-1} X_a \\ \mu_{b|a} = \mu_b - \Sigma_{ba} \Sigma_{aa}^{-1} \mu_a \\ \Sigma_{b|a} = \Sigma_{bb} - \Sigma_{ba} \Sigma_{aa}^{-1} \Sigma_{ab} \end{cases}$

$\therefore X_{b|a} \sim N(\mu_{b|a}, \Sigma_{b|a})$

$X_b = X_{b|a} + \Sigma_{ba} \Sigma_{aa}^{-1} X_a \quad E(X_b|X_a) = \mu_{b|a} + \Sigma_{ba} \Sigma_{aa}^{-1} X_a = \mu_b + \Sigma_{ba} \Sigma_{aa}^{-1} (X_a - \mu_a)$

$\text{Var}(X_b|X_a) = \Sigma_{b|a} \quad \therefore X_b|X_a \sim N(\mu_b + \Sigma_{ba} \Sigma_{aa}^{-1} (X_a - \mu_a), \Sigma_{bb} - \Sigma_{ba} \Sigma_{aa}^{-1} \Sigma_{ab})$

$\text{cov}(X, Z) = E[(X - \mu)(Z - 0)^T] = E[(WZ + \epsilon)(Z^T)] = E(WZ Z^T) + E(\epsilon Z^T) = W \text{Var}(Z) = W I = W$

$\therefore P(Z|X)$ is found from the above equation and the calculated $E(X) E(Z) \text{Var}(X) \text{Var}(Z) \text{cov}(X, Z)$

Difference with GMM \Rightarrow GMM: $Z \rightarrow \text{discrete: } 1, 2, \dots, K$ $P\text{-PCA: } Z \rightarrow \text{continuous: } Z \sim N$

$X \downarrow \quad X|Z \sim N$ $X \downarrow \quad X|Z \sim N$