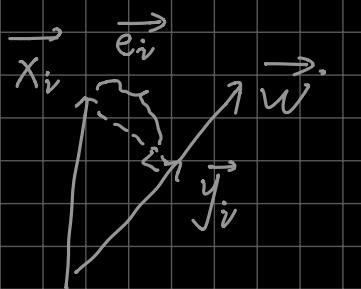


Lecture 8 PCA

Derivation of PCA:

I. Through MSE (Mean Square Error)



$$y = \underbrace{w^T x}_{\text{PCA变换矩阵}} \Rightarrow \text{降维后数据}$$

$$\vec{e}_i = \overrightarrow{x_i} - \underbrace{\vec{y}_i \cdot w}_{\text{重构数据}} w \Rightarrow \text{重构损失.}$$

$$\text{MLE: } J(w) = \frac{1}{m} \cdot \sum_{i=1}^m (x_i - (x_i^T w) \cdot w)^2$$

$$\text{s.t. } \|w\| = 1$$

Lagrangian Multiplier Method:

$$J(w, \lambda) = \frac{1}{m} \cdot \sum_{i=1}^m (x_i - (x_i^T w) \cdot w)^2 - \lambda \cdot (w^T w - 1)$$

$$\frac{\partial J(w, \lambda)}{\partial w} = -2 \cdot \underbrace{\sum_{i=1}^m x_i \cdot x_i^T w}_{\Sigma_x} - 2 \cdot \lambda \cdot w = 0$$

$$\Sigma_x \cdot w = -\lambda \cdot w$$

w 为 Σ_x 的 eigenvector (对上式作特征值解得)

因此 PCA components 即 Σ_x 的 eigen vectors,
即 w.

II. Through Eigen value Decomposition (特征分解)

$$\Sigma_X = \sum_{i=1}^N \cdot X_i \cdot X_i^T \Rightarrow \text{covariance matrix}$$

(协方差矩阵)

$$\Sigma_X = U \cdot \Lambda \cdot U^T \Rightarrow \text{特征分解}$$

U 即 PCA components.

III. Through SVD (Singular Value Decomposition)

$$X = U \cdot \Sigma \cdot V^T \Rightarrow SVD$$

where : U is the eigenvector of $X \cdot X^T$ and $U \cdot U^T = I$

and V is the eigenvector of $X^T \cdot X$ and $V^T \cdot V = I$

$$\Rightarrow \Phi X \cdot X^T = U \cdot \Sigma \cdot V^T \cdot V \cdot \Sigma \cdot U^T = U \cdot \Sigma^2 \cdot U^T$$
$$= \underbrace{U \cdot \Lambda \cdot U^T}_{\text{特征分解}}$$

$$\textcircled{2} X^T \cdot X = V^T \cdot \Sigma \cdot U^T \cdot U \cdot \Sigma \cdot V$$
$$= \underbrace{V^T \cdot \Sigma^2 \cdot V}_{\text{特征分解}}$$

U 即 PCA components.

IV. Hebbian learning. Oja. LMSER

$$J(w) = \frac{1}{m} \sum_{i=1}^m \underbrace{\left(x_i - (x_i^T w) \cdot w \right)^2}_{\text{实际上为: 上式只是另一种写法}}$$

$$\left\| x_i - (x_i^T w) \cdot w \right\|_2^2$$

$$\left\| x_i - (x_i^T w) \cdot w \right\|_2^2 = \text{tr} \left((x_i - (x_i^T w) \cdot w)^T \cdot (x_i - (x_i^T w) \cdot w) \right)$$

$$(AB)^T = B^T A^T$$

$$(ABC)^T = C^T B^T A^T$$

$$= \text{tr} \left(x_i^T - w^T w^T x_i \right) (x_i - (x_i^T w) \cdot w)$$

$$= \text{tr} \left(x_i^T \cdot x_i - w^T w^T x_i \cdot x_i - x_i^T x_i^T w w + w^T w^T x_i \cdot x_i^T w w \right)$$

tr中可移动

$$= \text{tr} (x_i^T x_i) - \text{tr} (\underbrace{w^T w^T x_i}_{w^T x_i \cdot x_i w^T} \cdot x_i) - \cancel{\text{tr} (x_i^T x_i^T w w)} + \cancel{\text{tr} (w^T x_i x_i^T w)}$$

$$w^T x_i \cdot x_i w^T = w^T x_i \cdot x_i^T w$$

$$w^T x_i x_i^T w = w^T x_i \cdot x_i^T w$$

$$= \underbrace{\text{tr} (x_i^T x_i)}_{\text{常数}} - \text{tr} (w^T \cdot x_i \cdot x_i^T \cdot w)$$

$$= C - \left\| w^T x_i \right\|_2^2$$

$$= C - y_i \cdot y_i^T$$

$$\therefore J(w) = \frac{1}{m} \cdot \sum_{i=1}^m \left(x_i - (x_i^T w) \cdot w \right)^2$$

$$= C - \sum_{i=1}^m y_i \cdot y_i^T$$

$$= C - \sum y \text{ covariance of } y,$$

因此 PCA minimize $J(w)$ 相当于 maximize
变换后数据的方差。

Which vector is the PCA solution?

$$\therefore C - \text{tr}(W^T X \cdot X^T W)$$

$$= C - \text{tr}(W^T \cdot \underbrace{\sum_X w}_\Downarrow) \quad \text{上式的} \quad \sum_X w = (-\vec{\lambda})w \\ = C - \text{tr}(W^T \cdot \underbrace{\lambda \cdot w}_\Downarrow) \quad = \lambda w$$

$$= C - \text{tr}(\lambda \cdot \underbrace{w w^T}_\Downarrow)$$

$$= C - \lambda \quad \Downarrow \|w\| = 1$$

$$\therefore \text{minimize } J(w) \Leftrightarrow C - \lambda \Leftrightarrow$$

maximize λ , 因此 PCA 的 solution 是 eigenvalue

λ 的前几个。 PCA component.

$$\triangle \text{ 也就是说: } \frac{\sum_{i=1}^k (x_i - W^T x)^2}{\sum_{i=1}^m \|x_i\|^2} \leq 0.01 \Leftrightarrow \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^m \lambda_i} \geq 0.99$$

的证明方法。

Lecture 9 FA

Why FA:

从朴素生成模型的角度思考PCA.

y 为一个低维的向量， X 为原始的数据，若能找这个隐变量 y 使其线性变换后 ($X = A^T y + u + e$) 能接近原始的分布，那么这个 y 即可代表降维后的数据。

因此 FA 的目的即找此 y 满足上述，方法为 EM。

Derivation:

$$y \sim \mathcal{N}(0, \Sigma_y) \xrightarrow{\text{I}^2}, \quad e \sim \mathcal{N}(0, \Gamma^2 I)$$

y 为标准分布 e 与 y 独立为 noise

$$P(X|\theta) \sim \mathcal{N}(u, A A^T + \Gamma^2 I)$$

~~~~~

$$Y = AX + b, \text{ 若 } X \sim N(u, \Sigma),$$

$$\text{则 } Y \sim N(Au + b, A \Sigma A^T)$$

$\therefore X = Ay + u + e$ , 因此依上式得此。

E step:

$$p(y|x) = \frac{p(y|\theta) \cdot p(x|y, \theta)}{p(x|\theta)} = \frac{\mathcal{N}(0, \Sigma_y) \cdot \mathcal{N}(Ay + u, \Gamma^2 I)}{\mathcal{N}(u, A^T A + \Gamma^2 I)}$$

Conditional Gaussian:

$$X_1 \sim N(u_1, \Sigma_1), \quad X_2 \sim N(u_2, \Sigma_2)$$

$$(X_1 | X_2 = a) \sim N(\bar{u} - \bar{\Sigma}), \text{ where:}$$

条件高斯

$$\left\{ \begin{array}{l} \bar{u} = u_1 + \Sigma_{12} \cdot \Sigma_{22}^{-1} (a - u_2) \end{array} \right.$$

$$\left. \begin{array}{l} \bar{\Sigma} = \Sigma_{11} - \Sigma_{12} \cdot \Sigma_{22}^{-1} \Sigma_{21} \end{array} \right.$$

相关系数:  $P = \frac{\text{cov}(X_1, X_2)}{\sigma_1 \sigma_2}$

$$\Sigma_1 = \Gamma^2 I \quad \underbrace{\Gamma}_{P \cdot \Gamma \cdot \Gamma}$$

$$\Rightarrow \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \quad \begin{matrix} \Sigma_{11} \\ \Sigma_{21} \\ \Sigma_{22} \end{matrix} \quad \begin{matrix} \Gamma \\ \Gamma \cdot \Gamma \\ \Gamma \cdot \Gamma \end{matrix} \quad \Sigma_{12} = \Sigma_{21} = \Sigma_{22} = \Gamma^2$$

$$\Rightarrow E(y|x) = \bar{u} = 0 + \underbrace{A^T(AA^T + \sigma^2 I)^{-1}}_{\text{令其为 } W} \cdot x$$

$$= Wx$$

$$E(yy^T|x) = \sum_{k=1}^{\infty} + E(y|x)^2$$

$$= I - W \cdot A + \underbrace{Wx \cdot (Wx)^T}_{Wx \cdot x^T \cdot W^T}$$

M step:

$$\underset{\theta}{\operatorname{argmax}} \underset{\sim}{P(y, x | \theta)}$$

$$u, \Gamma, A \quad \downarrow$$

$$\sum_{i=1}^N \log P(y_i, x_i | \theta)$$

$$= \sum_{i=1}^N \log \frac{P(y_i, x_i | \theta) \cdot P(y_i | x_i; \theta)}{P(y_i | x_i; \theta)}$$

$$= \sum_{i=1}^N \left[ P(y_i | x_i; \theta) - \log \frac{P(y_i | \theta) \cdot P(x_i | y_i; \theta)}{P(y_i | x_i; \theta)} \right]$$

$$Q(y_i) / Q(z_i)$$

General EM  
↑ 考虑的部 分

隐变量

$$\theta = \arg\max_{\theta} \sum_i \cdot \int_{z^{(i)}} Q_i(z^{(i)}) \cdot \log \frac{P(X^{(i)}, z^{(i)}|\theta)}{Q_i(z^{(i)})}$$

$$= \sum_i E_{z^{(i)} \sim Q_i} [\log \frac{P(X^{(i)}, z^{(i)}|\theta)}{Q_i(z^{(i)})}] \xrightarrow{\text{Gaussian form:}} e^{-\cdot}$$

$$= -\sum_i \left[ \frac{1}{2} \right]$$

Square form, can directly take  $\Delta u \theta$

# Lecture 10

# ICA

- Independence :  $p(y_1, y_2) = p(y_1) \cdot p(y_2)$  ||  $E(h_1(y_1), h_2(y_2)) = E(h_1(y_1)) \cdot E(h_2(y_2))$   
 $\Rightarrow$  no information can be gained from each other.
- Uncorrelatedness :  $E(y_1 y_2) = E(y_1) \cdot E(y_2)$
- Conditional Independence:  $P(X, Y | Z) = P(X | Z) \cdot P(Y | Z)$
- Dependence
  - Common cause
  - Causal relation
  - Common effect (conditional dependence)

Gaussian Distribution properties:

$$E(XY) = E(X) \cdot E(Y)$$

① Uncorrelatedness  $\Leftrightarrow$  Independence

$$E(f(x)g(y)) = E(f(x)) \cdot E(g(y)), V_{fg}$$

② Central Limit Theorem (CLT) means:  $\sum_{i=1}^n X_i \xrightarrow{D} S$

If  $X_1, X_2, \dots, X_n$  are independent random variables, the sum is approximately Gaussian Distribution.

③ Cramer's decomposition theorem:  $\xi = \xi_1 + \xi_2$ .

If  $\xi$  is normally distributed &  $\xi_1, \xi_2$  are independent random variables Gaussian distribution

, then  $\{\cdot\}_1, \{\cdot\}_2$  are both normally distributed.

④ Sub-Gaussian and Super-Gaussian  $\Rightarrow$   
e.g.: uniform distribution Laplace distribution  
(double side exp.) 

Vector is noiseless, then it's more Super-Gaussian;  
Otherwise, if a large amount of noise exists, then it's more  
tend to Gaussian distribution.

⑤ A gaussian variable has the largest entropy among all random variables of equal variance.  $\Rightarrow$

More noise, larger entropy (according to  $\text{H}(\cdot)$ )

$$\text{Kurtosis} = \text{kurt}(y) = E(y^4) - 3E^2(y^2)$$

$\Rightarrow \begin{cases} \text{kurt}(y) > 0 \Leftrightarrow \text{super-Gaussian} \\ \text{kurt}(y) < 0 \Leftrightarrow \text{sub-Gaussian} \end{cases}$

## Measures

of non-Gaussian

[ closer to 0  $\Rightarrow$  more ]

$I(y_1, y_2 \dots y_n) = \sum_{i=1}^n H(y_i) - H(y)$

# Independent Component Analysis (ICA)

Intuition:

$$X = A \underbrace{y}_{\substack{\text{raw data} \\ (\text{观测数据})}} \quad \Leftrightarrow \quad \begin{aligned} & \text{transform matrix} \\ & A^{-1} \cdot X = y \Leftrightarrow y = W^T \cdot X \\ & = W^T \cdot A \cdot S \\ & = \underbrace{z^T \cdot S}_{\substack{\text{more Gaussian than} \\ \text{any of the } s_i}} \end{aligned}$$

original independent components  
(target data)

Thus, we want to find  $w$  that maximizes the non-Gaussianity measured by negentropy

原始的  $s$  中  
1个

$\Rightarrow$  ICA allows  $\leq 1$  Gaussian variables. (complete symmetry, no information on the direction of  $A$ ) : 当  $> 1$  个原始

Fast ICA algorithm:

① Randomly initialize  $w$

② while  $w$  is not converged :

get from the derivative of MLE.

③  $w = E [x \cdot g_1(w^T x)] - E [g_2(w^T x)]$

where  $\begin{cases} g_1(x) = \frac{1}{a_1} \cdot \log \cosh(a_1 x) \\ g_2(x) = x \cdot e^{-\frac{x^2}{2}} \end{cases}$

④  $w = \frac{w}{\|w\|}$  # normalized

end

CDF (Cumulative Distribution Functions)

PDF (Probability Distribution Function)

