

PRML chapter12

1. Principal Component Analysis

(a) Maximum variance formulation

- Our goal is to project data D onto a space having a dimensionality M
- each data point \mathbf{x}_n is projected onto $\mathbf{u}_1^T \mathbf{x}_n$
- mean of the projected data is $\mathbf{u}_1^T \bar{\mathbf{x}}$
- the variance of the projected data is $\frac{1}{N} \sum_{n=1}^N \{\mathbf{u}_1^T \mathbf{x}_n - \mathbf{u}_1^T \bar{\mathbf{x}}\}^2 = \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$
- maximize $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$ with respect to \mathbf{u}

$$\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^T \mathbf{u}_1)$$

- the variance will be a maximum when we set \mathbf{u} equal to the maximum eigen vector
- principal component analysis requires the mean and variance, also needs largest M eigen-vectors

(b) Minimum-error formulation

- Introduce a complete orthonormal set of D -dimensional vectors \mathbf{u}

$$\mathbf{x}_n = \sum_{i=1}^D (\mathbf{x}_n^T \mathbf{u}_i) \mathbf{u}_i$$

- we approximate the data

$$\tilde{\mathbf{x}}_n = \sum_{i=1}^M z_{ni} \mathbf{u}_i + \sum_{i=M+1}^D b_i \mathbf{u}_i$$

- our goal is to minimize

$$J = \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2$$

- first, minimize it with respect to \mathbf{z} , then minimize it with respect to \mathbf{u}

(c) Applications of PCA

- comprehension of the data

$$\tilde{\mathbf{x}}_n = \sum_{i=1}^M (\mathbf{x}_n^T \mathbf{u}_i) \mathbf{u}_i + \sum_{i=M+1}^D (\bar{\mathbf{x}}^T \mathbf{u}_i) \mathbf{u}_i = \bar{\mathbf{x}} + \sum_{i=1}^M (\mathbf{x}_n^T \mathbf{u}_i - \bar{\mathbf{x}}^T \mathbf{u}_i) \mathbf{u}_i$$

- data pre-processing

(d) PCA for high-dimensional data

- if the number of data points is smaller than the dimensionality of the data space, we have to take different approach because the computational cost $O(D^3)$

- take the following algorithms and solve it with $O(N^3)$
- let \mathbf{X} be a matrix whose n th row is $(\mathbf{x}_n - \bar{\mathbf{x}})^T$
- under $\mathbf{v}_i = \mathbf{X}\mathbf{u}_i$,

$$\frac{1}{N} \mathbf{X}\mathbf{X}^T \mathbf{v}_i = \lambda_i \mathbf{v}_i$$

2. Probabilistic PCA

- PCA can also be expressed as the maximum likelihood solution of a probabilistic latent variable model
- First we gave the prior distribution over \mathbf{z} as $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$
- the conditional distribution of observed value \mathbf{x} , $p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2\mathbf{I})$
- $P(x) = \int P(x|z)p(z)dz$ provide us the parameters

(a) Maximum likelihood PCA

- we want to maximize log likelihood function

$$\begin{aligned} \ln p(\mathbf{X}|\boldsymbol{\mu}, \mathbf{W}, \sigma^2) &= \sum_{n=1}^N \ln p(\mathbf{x}_n|\mathbf{W}, \boldsymbol{\mu}, \sigma^2) \\ &= -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\mathbf{C}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \end{aligned}$$

- the calculation is very complex, Tipping and Bishop(1999)

$$\mathbf{W}_{ML} = \mathbf{U}_M (\mathbf{L}_M - \sigma^2\mathbf{I})^{1/2} \mathbf{R}$$

- the number of independent parameters are controlled automatically

(b) EM algorithm for PCA

- it has an advantage when we treat high-dimensional data
- complete log-likelihood function takes the form

$$\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\mu}, \mathbf{W}, \sigma^2) = \sum_{n=1}^N \{\ln p(\mathbf{x}_n|\mathbf{z}_n) + \ln p(\mathbf{z}_n)\}$$

- E Step

$$E[\mathbf{z}_n] = \mathbf{M}^{-1} \mathbf{W}^T (\mathbf{x}_n - \bar{\mathbf{x}})$$

$$E[\mathbf{z}_n \mathbf{z}_n^T] = \sigma^2 \mathbf{M}^{-1} + E[\mathbf{z}_n] E[\mathbf{z}_n]^T$$

- M Step

$$\mathbf{W}_{new} = \left[\sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}}) E[\mathbf{z}_n]^T \right] \left[\sum_{n=1}^N E[\mathbf{z}_n \mathbf{z}_n^T] \right]^{-1}$$

$$\sigma_{new}^2 = \frac{1}{ND} \sum_{n=1}^N \left\{ \|\mathbf{x}_n - \bar{\mathbf{x}}\|^2 - 2E[\mathbf{z}_n]^T \mathbf{W}_{new}^T (\mathbf{x}_n - \bar{\mathbf{x}}) + \text{Tr} (E[\mathbf{z}_n \mathbf{z}_n^T] \mathbf{W}_{new}^T \mathbf{W}_{new}) \right\}$$

(c) Bayesian PCA

- we want to decide M with Bayesian approach

- choose model with Bayesian approach

(d) Factor analysis

- factor analysis has

$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \boldsymbol{\Psi})$$

- We can determine $\boldsymbol{\mu}$, \mathbf{W} , $\boldsymbol{\Psi}$ in the factor analysis model by maximum likelihood
- use EM algorithm

3. Kernel PCA

- we want to obtain non-linear generalization
- the principal vector is defined as

$$\mathbf{S}\mathbf{u}_i = \lambda_i \mathbf{u}_i$$

where

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T$$

- sample covariance matrix is

$$\mathbf{C} = \frac{1}{N} \sum_{n=1}^N \boldsymbol{\phi}(\mathbf{x}_n) \boldsymbol{\phi}(\mathbf{x}_n)^T$$

- eigen vector of matrix C is

$$\mathbf{v}_i = \sum_{n=1}^N a_{in} \boldsymbol{\phi}(\mathbf{x}_n)$$

Kernel function gives us a solution for a by solving the following eigenvalue problem

$$\mathbf{K}\mathbf{a}_i = \lambda_i N \mathbf{a}_i$$

- the projection of \mathbf{x} onto eigenvector i is written as

$$y_i(\mathbf{x}) = \boldsymbol{\phi}(\mathbf{x})^T \mathbf{v}_i = \sum_{n=1}^N a_{in} \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{x}_n) = \sum_{n=1}^N a_{in} k(\mathbf{x}, \mathbf{x}_n)$$

4. Nonlinear latent value model

- consider the models based on non-linear and non-Gaussian distributions

(a) Independent component analysis

- the example of non-linear latent variables model
- In this models, observed variables are related linearly to the latent variables but the latent distribution is non-Gaussian.
- latent variables are independent so

$$p(\mathbf{z}) = \prod_{j=1}^M p(z_j)$$

- no need to consider the noise, because the number of observed variables and latent variables are same
- The success of this approach requires that the latent variables have non-Gaussian distributions.

(b) Auto associative neural networks

- Create the neural network model whose input and output are D dimension
- The model tells us the features of the dataset.
- tells us more information than PCA but calculation amount is large