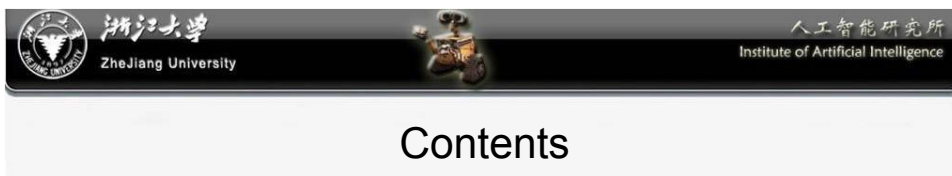




Artificial Intelligence

Course Review

Donghui Wang
AI Institute@ZJU
2014.4



Contents

- About final exam
- Basic concepts
- Important formulas and derivations
- Learning Algorithms

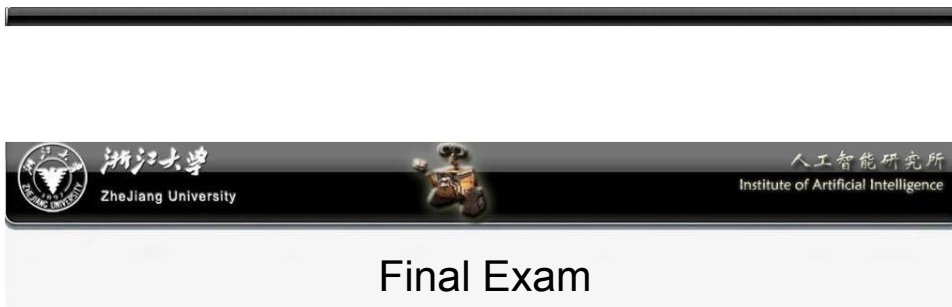
References:

1. AI Course slides, <http://10.15.62.79/cv>
2. Christopher M. Bishop. "Pattern Recognition and Machine Learning", 2006, Springer.
3. Stuart J. Russell and Peter Norvig. "Artificial Intelligence: A Modern Approach (Third Edition)". 2009, Prentice Hall.
4. <http://cs229.stanford.edu/info.html>, by Prof. Andrew Ng



Course Review

ABOUT FINAL EXAM



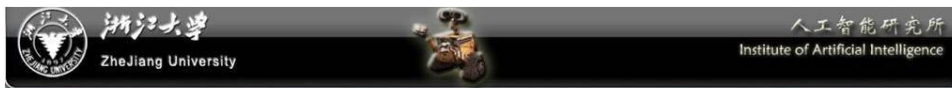
1. Fill in the blank (30 points, 2pt/per)
 - Basic concept, definition and fundamental knowledge
2. Multiple Choice (20 points, 2pt/per)
 - Same as first part, but focus on the difference among important concepts and definitions.
3. Calculus, Analysis and Proof (50 points, 10~15pt/per)
4. Algorithm Design (10 points)



Final Exam

- Information:

Date and Time	22 April, 8:00am~10:00am (2 hours)
Location	Teaching Building #7, Room 302 (Multimedia) Yuquan Campus
Number of Examination	40
Examination form	Closed Book
Final Q&A Time	17 April, After class (15:40-17:30, 19:00~21:00)
Final Q&A Room	Business Administration Building, Room 218



Course Review

Basic concepts



Basic concepts

1. What is AI?

- **Acting humanly**: the **Turing test** approach
- **Thinking humanly**: the cognitive modeling approach
- **Thinking rationally**: the “laws of thought” approach
- **Acting rationally**: the rational agent approach

2. Turing test

- A computer passes the test if a human interrogator, after posing some written questions, cannot tell whether the written responses come from a person or from a computer.

3. Strong AI and Weak AI

- Philosophers use the term **weak AI** for the hypothesis that machines could possibly behave intelligently, and **strong AI** for the hypothesis that such machines would count as having actual minds.



Basic concepts

4. Supervised learning approach (predictive approach)

- Regression, classification (SVM)
- Training phase (learning phase), prediction phase
- Training set (with target vector), test set
- Binary classification, multiclass classification

5. Unsupervised learning approach (descriptive approach)

- Knowledge discovery
- Clustering, density estimation, dimensionality reduction (manifold learning)
- Training set (without target vector), **test set**

6. Reinforcement learning



Basic concepts

7. Error function

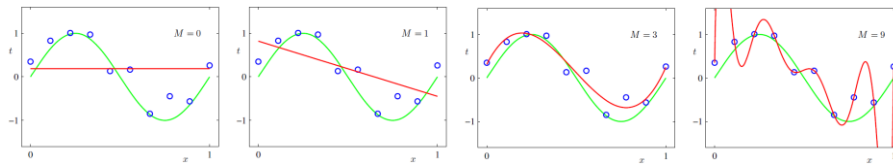
- SSE (the sum-of-squares error)

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2$$

- RMS (the root-mean-square error)

$$E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^*)/N}$$

8. Over-fitting phenomenon



- How to control?
- *Regularization (penalty term), Bayesian approach (prior), CV...*
- *Shrinkage methods: ridge regression (weight decay)*

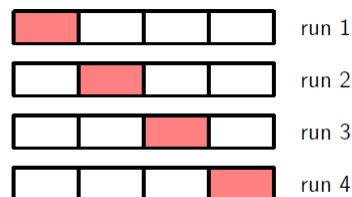


Basic concepts

9. Model comparison (model selection)

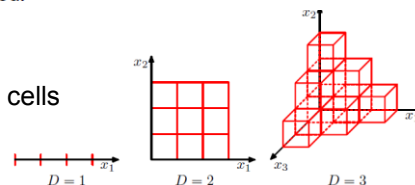
- validation set, Cross-validation (CV)

The technique of S -fold cross-validation, illustrated here for the case of $S = 4$, involves taking the available data and partitioning it into S groups (in the simplest case these are of equal size). Then $S - 1$ of the groups are used to train a set of models that are then evaluated on the remaining group. This procedure is then repeated for all S possible choices for the held-out group, indicated here by the red blocks, and the performance scores from the S runs are then averaged.



10. The curse of dimensionality

- Divide the input space into regular cells





Basic concepts

11. Frequentist statistics vs. Bayesian statistics

- View probabilities in terms of the frequencies of random, repeatable events.
- Probabilities provide a quantification of uncertainty and make rational coherent inference.

12. Likelihood function

$$p(\mathcal{D}|\mathbf{w})$$

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{p(\mathcal{D})}$$

- It expresses how probable the observed data set is for different settings of the parameter vector \mathbf{w}
- The likelihood is not a probability distribution over \mathbf{w} , and its integral with respect to \mathbf{w} does not (necessarily) equal one.
- Maximum Likelihood – ML (widely used frequentist estimator)
 - *Maximizing the likelihood is equivalent to minimizing the error (e.g. SSE).*
 - *i.i.d (independent and identically distributed)*
- Maximum posterior – MAP



Basic concepts

13. Decision regions, decision boundaries (decision surface)

14. Linearly separable

15. Loss function (cost function, - utility function), loss matrix

- Minimize the average loss:

$$L_{kj}=0, \text{ for } k=j. \text{ others } 1.$$

$$\mathbb{E}[L] = \sum_k \sum_j \int_{\mathcal{R}_j} L_{kj} p(\mathbf{x}, \mathcal{C}_k) d\mathbf{x}.$$

16. Reject option

17. Discriminant function

18. Generative models and discriminative models

19. Relative entropy (Kullback-Leibler divergence or KL divergence)

$$\text{KL}(q\|p) = - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z}|\mathbf{X}, \theta)}{q(\mathbf{Z})} \right\}$$



Basic concepts

19. Conjugacy and Conjugate priors

- Posterior distribution has the same functional form as the prior.

20. Exponential family of distribution

21. Mahalanobis distance $\Delta^2 = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$

22. Covariance matrix and precision matrix

- The inverse of the covariance matrix

23. Parametric and non-parametric models

- Probability distributions have specific functional forms governed by a small number of parameters.
- There are few assumptions about the form of the distribution.

24. Jacobian matrix



Basic concepts

25. Linear models

- Linear model $y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_D x_D$ where $\mathbf{x} = (x_1, \dots, x_D)^T$
- Linear basis function model $y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$
- Generalized linear models
 - Activation function and link function $y(\mathbf{x}) = f(\mathbf{w}^T \mathbf{x} + w_0)$
 - Logistic sigmoid function and logit function

26. The Fisher's criterion

- maximize the separation between the projected class means as well as the inverse of the total within-class variance.
- Generalized Rayleigh quotient, Between-class covariance matrix and Within-class covariance matrix

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \quad J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}$$

27. The perceptron criterion: $E_P(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \mathbf{w}^T \boldsymbol{\phi}_n t_n$





Basic concepts

28. The probit regression:

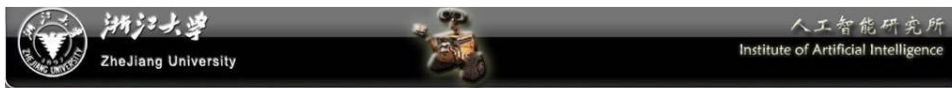
- Use CDF of $N(x|0,1)$ to construct an activation function (Inverse probit function)

29. The Laplace Approximation

- Find a Gaussian approximation $q(z)$ which is centred on a mode of the distribution $p(z)$

30. kernel function $k(x, x') = \phi(x)^T \phi(x')$

- symmetric function
- How to construct kernel functions?
- Kernel trick
- The simplest kernel function: $k(x,y) = x^T y$



Course Review

IMPORTANT FORMULAS AND DERIVATIONS





Probability theory

- Marginal probability

$$p(X) = \sum_Y p(X, Y)$$

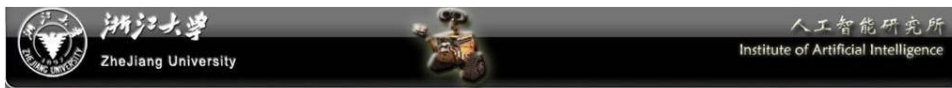
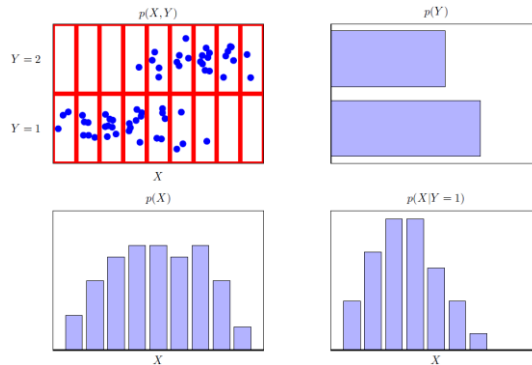
- Joint probability

$$p(X, Y) = p(Y|X)p(X)$$

- Conditional probability

- Bayes' theorem

$$p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$$



Multivariate Gaussian Distribution

- Definition: $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right\}$

$$\Delta^2 = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})$$

- The matrix $\boldsymbol{\Sigma}$ can be taken to be symmetric, without loss of generality. $\boldsymbol{\Sigma}^{-1}$ is symmetric.
- The eigenvector equation for the covariance matrix:

$$\boldsymbol{\Sigma} \mathbf{u}_i = \lambda_i \mathbf{u}_i \quad \text{where } i = 1, \dots, D \quad \mathbf{u}_i^T \mathbf{u}_j = I_{ij} \quad I_{ij} = \begin{cases} 1, & \text{if } i = j \\ 0, & \text{otherwise.} \end{cases} \quad \mathbf{U} \mathbf{U}^T = \mathbf{I}$$

$$\Rightarrow \boldsymbol{\Sigma}^{-1} = \mathbf{U} \boldsymbol{\Lambda}^{-1} \mathbf{U}^T \quad \Rightarrow \Delta^2 = \sum_{i=1}^D \frac{y_i^2}{\lambda_i} \quad y_i = \mathbf{u}_i^T (\mathbf{x} - \boldsymbol{\mu}) \quad \mathbf{y} = \mathbf{U}(\mathbf{x} - \boldsymbol{\mu})$$

$$\Rightarrow |\boldsymbol{\Sigma}|^{1/2} = \prod_{j=1}^D \lambda_j^{1/2} \quad \Rightarrow |\mathbf{J}| = 1$$

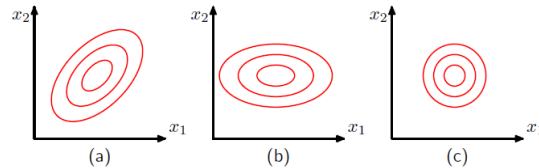


Multivariate Gaussian Distribution

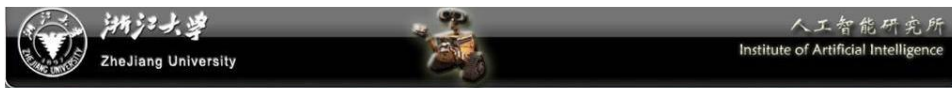
• Definition: $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right\}$

$$\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu} \quad \mathbb{E}[\mathbf{x}\mathbf{x}^T] = \boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\Sigma} \quad \text{cov}[\mathbf{x}] = \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T]$$

→ $\text{cov}[\mathbf{x}] = \boldsymbol{\Sigma}$



- A general symmetric covariance matrix $\boldsymbol{\Sigma}$ will have $D(D+1)/2$ independent parameters, and there are another D independent parameters in $\boldsymbol{\mu}$, giving $D(D+3)/2$ parameters in total.
 - 2D independent parameters $\boldsymbol{\Sigma} = \text{diag}(\sigma_i^2)$
 - isotropic covariance, $D+1$ independent parameters $\boldsymbol{\Sigma} = \sigma^2 \mathbf{I}$



Multivariate Gaussian Distribution

- Bayes' Theorem for Gaussian Variables:

$$\begin{aligned} p(\mathbf{x}) &= \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \\ p(\mathbf{y}|\mathbf{x}) &= \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1}) \end{aligned} \quad \Rightarrow \quad p(\mathbf{z}) = p(\mathbf{y}|\mathbf{x})p(\mathbf{x}) \quad \mathbf{z} = \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix}$$

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{M} & -\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \\ -\mathbf{D}^{-1}\mathbf{C}\mathbf{M} & \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{C}\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \end{pmatrix}$$

$$\mathbf{M} = (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}$$

$$\begin{aligned} \ln p(\mathbf{z}) &= \ln p(\mathbf{x}) + \ln p(\mathbf{y}|\mathbf{x}) \\ &= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Lambda}(\mathbf{x} - \boldsymbol{\mu}) - \frac{1}{2}(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b})^T \mathbf{L}(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b}) + \text{const} \end{aligned}$$

$$\begin{aligned} &-\frac{1}{2}\mathbf{x}^T(\boldsymbol{\Lambda} + \mathbf{A}^T\mathbf{L}\mathbf{A})\mathbf{x} - \frac{1}{2}\mathbf{y}^T\mathbf{L}\mathbf{y} + \frac{1}{2}\mathbf{y}^T\mathbf{L}\mathbf{A}\mathbf{x} + \frac{1}{2}\mathbf{x}^T\mathbf{A}^T\mathbf{L}\mathbf{y} \\ &= -\frac{1}{2}\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix}^T \begin{pmatrix} \boldsymbol{\Lambda} + \mathbf{A}^T\mathbf{L}\mathbf{A} & -\mathbf{A}^T\mathbf{L} \\ -\mathbf{L}\mathbf{A} & \mathbf{L} \end{pmatrix} \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} = -\frac{1}{2}\mathbf{z}^T \mathbf{R} \mathbf{z} \end{aligned} \quad \Rightarrow \quad \mathbf{R} = \begin{pmatrix} \boldsymbol{\Lambda} + \mathbf{A}^T\mathbf{L}\mathbf{A} & -\mathbf{A}^T\mathbf{L} \\ -\mathbf{L}\mathbf{A} & \mathbf{L} \end{pmatrix}$$

$$\text{cov}[\mathbf{z}] = \mathbf{R}^{-1} = \begin{pmatrix} \boldsymbol{\Lambda}^{-1} & \boldsymbol{\Lambda}^{-1}\mathbf{A}^T \\ \mathbf{L}\mathbf{A}^{-1} & \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T \end{pmatrix}$$

$$\mathbf{x}^T\boldsymbol{\Lambda}\boldsymbol{\mu} - \mathbf{x}^T\mathbf{A}^T\mathbf{L}\mathbf{b} + \mathbf{y}^T\mathbf{L}\mathbf{b} = \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix}^T \begin{pmatrix} \boldsymbol{\Lambda}\boldsymbol{\mu} - \mathbf{A}^T\mathbf{L}\mathbf{b} \\ \mathbf{L}\mathbf{b} \end{pmatrix} \quad \Rightarrow \quad \mathbb{E}[\mathbf{z}] = \mathbf{R}^{-1} \begin{pmatrix} \boldsymbol{\Lambda}\boldsymbol{\mu} - \mathbf{A}^T\mathbf{L}\mathbf{b} \\ \mathbf{L}\mathbf{b} \end{pmatrix} = \begin{pmatrix} \boldsymbol{\mu} \\ \mathbf{A}\boldsymbol{\mu} + \mathbf{b} \end{pmatrix}$$



Maximum likelihood for the Gaussian

- Given a data set $\mathbf{X} = (x_1, \dots, x_N)^T$ in which the observations $\{x_n\}$ are assumed to be drawn independently from a multivariate Gaussian distribution:

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right\}$$

$$\Rightarrow \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu})$$

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^T$$

$$\frac{\partial \mathbf{a}^T \mathbf{X}^{-1} \mathbf{b}}{\partial \mathbf{X}} = -\mathbf{X}^{-T} \mathbf{a} \mathbf{b}^T \mathbf{X}^{-T}$$

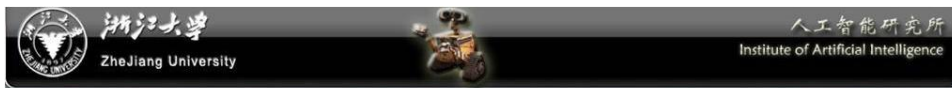
$$\frac{\partial}{\partial \boldsymbol{\mu}} \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu}) \Rightarrow \boldsymbol{\mu}_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

$$\frac{\partial}{\partial \mathbf{A}} \ln |\mathbf{A}| = (\mathbf{A}^{-1})^T$$

$$\frac{\partial}{\partial \boldsymbol{\Sigma}} \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{N}{2} \frac{\partial}{\partial \boldsymbol{\Sigma}} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \frac{\partial}{\partial \boldsymbol{\Sigma}} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu}) \Rightarrow \boldsymbol{\Sigma} = \mathbf{S}$$

$$\Rightarrow \boldsymbol{\Sigma}_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{\text{ML}})(\mathbf{x}_n - \boldsymbol{\mu}_{\text{ML}})^T$$

$$\begin{aligned} \mathbb{E}[\boldsymbol{\mu}_{\text{ML}}] &= \boldsymbol{\mu} \\ \mathbb{E}[\boldsymbol{\Sigma}_{\text{ML}}] &= \frac{N-1}{N} \boldsymbol{\Sigma} \end{aligned}$$



Bayesian Inference for the Gaussian

- Known the variance, to infer the mean:

$$\text{Likelihood: } p(\mathbf{X}|\boldsymbol{\mu}) = \prod_{n=1}^N p(x_n|\boldsymbol{\mu}) = \frac{1}{(2\pi\sigma^2)^{N/2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \boldsymbol{\mu})^2 \right\}$$

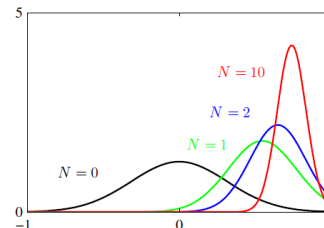
$$\text{Prior: } p(\boldsymbol{\mu}) = \mathcal{N}(\boldsymbol{\mu}|\boldsymbol{\mu}_0, \sigma_0^2)$$

$$\text{Posterior: } p(\boldsymbol{\mu}|\mathbf{X}) \propto p(\mathbf{X}|\boldsymbol{\mu})p(\boldsymbol{\mu})$$

$$p(\boldsymbol{\mu}|\mathbf{X}) = \mathcal{N}(\boldsymbol{\mu}|\boldsymbol{\mu}_N, \sigma_N^2)$$

$$\boldsymbol{\mu}_N = \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \boldsymbol{\mu}_0 + \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} \boldsymbol{\mu}_{\text{ML}}$$

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}$$



Likelihood: $p(\mathbf{X}|\mu) = \prod_{n=1}^N p(x_n|\mu) = \frac{1}{(2\pi\sigma^2)^{N/2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 \right\}$

Prior: $p(\mu) = \mathcal{N}(\mu|\mu_0, \sigma_0^2)$

Posterior: $p(\mu|\mathbf{X}) \propto p(\mathbf{X}|\mu)p(\mu)$ $p(\mu|\mathbf{X}) = \mathcal{N}(\mu|\mu_N, \sigma_N^2)$

$$-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) = -\frac{1}{2}\mathbf{x}^T \boldsymbol{\Sigma}^{-1}\mathbf{x} + \mathbf{x}^T \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} + \text{const}$$

$$\begin{aligned} -\frac{1}{2\sigma_0^2}(\mu - \mu_0)^2 - \frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 \\ = -\frac{\mu^2}{2} \left(\frac{1}{\sigma_0^2} + \frac{N}{\sigma^2} \right) + \mu \left(\frac{\mu_0}{\sigma_0^2} + \frac{1}{\sigma^2} \sum_{n=1}^N x_n \right) + \text{const} \end{aligned}$$

$$\begin{aligned} \frac{1}{\sigma_N^2} &= \frac{N}{\sigma^2} + \frac{1}{\sigma_0^2} & \mu_N &= \left(\frac{N}{\sigma^2} + \frac{1}{\sigma_0^2} \right)^{-1} \left(\frac{\mu_0}{\sigma_0^2} + \frac{1}{\sigma^2} \sum_{n=1}^N x_n \right) \\ \mu_{\text{ML}} &= \frac{1}{N} \sum_{n=1}^N x_n & &= \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \mu_0 + \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} \mu_{\text{ML}}. \end{aligned}$$



Bayesian Inference for the Gaussian

2. Known the mean, to infer the variance: $\lambda \equiv 1/\sigma^2$

Likelihood: $p(\mathbf{X}|\lambda) = \prod_{n=1}^N \mathcal{N}(x_n|\mu, \lambda^{-1}) \propto \lambda^{N/2} \exp \left\{ -\frac{\lambda}{2} \sum_{n=1}^N (x_n - \mu)^2 \right\}$

Prior: $\text{Gam}(\lambda|a_0, b_0)$ *gamma distribution*

Posterior: $p(\lambda|\mathbf{X}) \propto p(\mathbf{X}|\lambda) \text{Gam}(\lambda|a_0, b_0)$ $\text{Gam}(\lambda|a, b) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda)$

$$p(\lambda|\mathbf{X}) \propto \lambda^{a_0-1} \lambda^{N/2} \exp \left\{ -b_0\lambda - \frac{\lambda}{2} \sum_{n=1}^N (x_n - \mu)^2 \right\} \Rightarrow \text{Gam}(\lambda|a_N, b_N)$$

$$a_N = a_0 + \frac{N}{2}$$

$$b_N = b_0 + \frac{1}{2} \sum_{n=1}^N (x_n - \mu)^2 = b_0 + \frac{N}{2} \sigma_{\text{ML}}^2$$



Maximum likelihood and least squares

- Assume: $t = y(\mathbf{x}, \mathbf{w}) + \epsilon$ $y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$
- Thus: $p(t|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1}) \rightarrow \mathbb{E}[t|\mathbf{x}] = \int t p(t|\mathbf{x}) dt = y(\mathbf{x}, \mathbf{w})$
- For data set $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ and target vector $\mathbf{t} = (t_1, \dots, t_N)^T$, the likelihood function:

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1})$$

$$\ln p(\mathbf{t}|\mathbf{w}, \beta) = \sum_{n=1}^N \ln \mathcal{N}(t_n | \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}) = \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\mathbf{w})$$

SSE: sum-of-squares
error function

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n)\}^2$$

Maximum likelihood and least squares

$$\ln p(\mathbf{t}|\mathbf{w}, \beta) = \sum_{n=1}^N \ln \mathcal{N}(t_n | \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}) = \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\mathbf{w})$$

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n)\}^2$$

- Solving \mathbf{w} by ML:

$$\nabla \ln p(\mathbf{t}|\mathbf{w}, \beta) = \sum_{n=1}^N \{t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n)\} \boldsymbol{\phi}(\mathbf{x}_n)^T$$

$$0 = \sum_{n=1}^N t_n \boldsymbol{\phi}(\mathbf{x}_n)^T - \mathbf{w}^T \left(\sum_{n=1}^N \boldsymbol{\phi}(\mathbf{x}_n) \boldsymbol{\phi}(\mathbf{x}_n)^T \right)$$

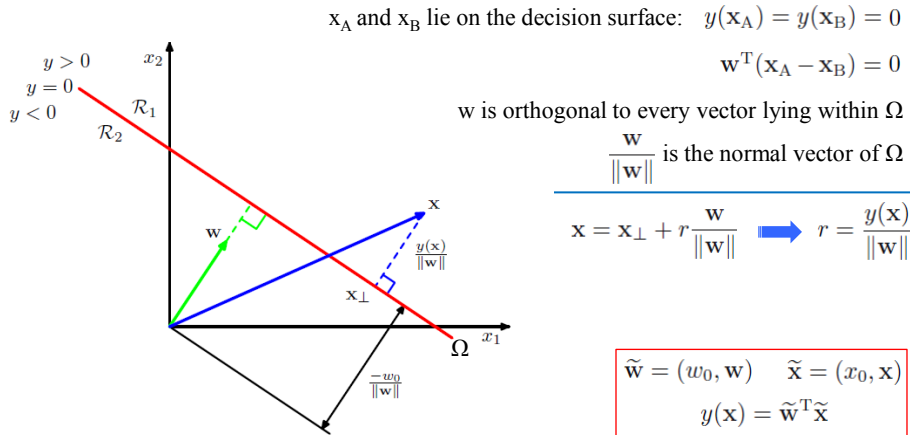
$$\Rightarrow \mathbf{w}_{\text{ML}} = (\boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^T \mathbf{t}$$

$$\boldsymbol{\Phi} = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix} \quad N \times M \text{ design matrix}$$

$$\boldsymbol{\Phi}^\dagger \equiv (\boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^T \quad \text{Moore-Penrose pseudo-inverse}$$

Two classes

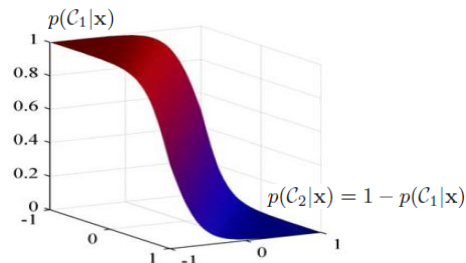
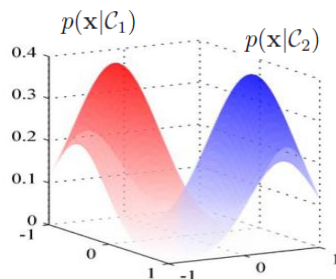
- Linear discriminant function: $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$
 - if $y(\mathbf{x}) \geq 0$, assign \mathbf{x} to class C_1 , else class C_2
 - decision surface Ω : $y(\mathbf{x}) = 0$
 - the normal distance from the origin to the decision surface: $\frac{\mathbf{w}^T \mathbf{x}}{\|\mathbf{w}\|} = -\frac{w_0}{\|\mathbf{w}\|}$



Probabilistic Generative Models: Continuous inputs

- Assume: $p(\mathbf{x}|\mathcal{C}_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \mu_k)^T \Sigma^{-1}(\mathbf{x} - \mu_k) \right\}$
- 2 classes: $p(\mathcal{C}_1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$

$\Rightarrow \mathbf{w} = \Sigma^{-1}(\mu_1 - \mu_2) \quad w_0 = -\frac{1}{2}\mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2}\mu_2^T \Sigma^{-1} \mu_2 + \ln \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)}$



Logistic regression

- Logistic regression model:

- Only M parameters need to be estimated.

logistic sigmoid function

$$p(C_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi) \quad p(C_2|\phi) = 1 - p(C_1|\phi) \quad \sigma(a) = \frac{1}{1 + \exp(-a)}$$

- For a data set $\{\phi_n, t_n\}$, where $t_n \in \{0, 1\}$ and $\phi_n = \phi(\mathbf{x}_n)$, $n=1, \dots, N$, the likelihood function can be written

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n} \quad \text{where } \mathbf{t} = (t_1, \dots, t_N)^T \text{ and } y_n = p(C_1|\phi_n).$$

- Cross-entropy error function:

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$$

$$y_n = \sigma(a_n)$$

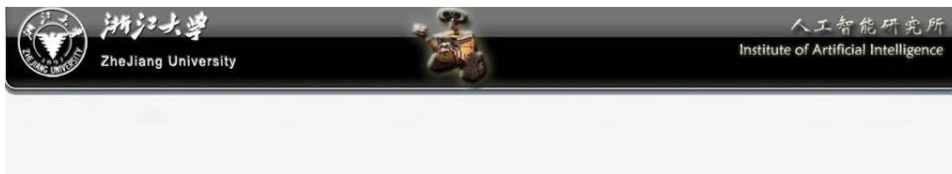
$$a_n = \mathbf{w}^T \phi_n$$

$$\frac{d\sigma}{da} = \sigma(1 - \sigma)$$

$$\left\{ \begin{array}{l} \frac{\partial E}{\partial y_n} = \frac{1 - t_n}{1 - y_n} - \frac{t_n}{y_n} = \frac{y_n(1 - t_n) - t_n(1 - y_n)}{y_n(1 - y_n)} = \frac{y_n - t_n}{y_n(1 - y_n)} \\ \frac{\partial y_n}{\partial a_n} = \frac{\partial \sigma(a_n)}{\partial a_n} = \sigma(a_n)(1 - \sigma(a_n)) = y_n(1 - y_n) \end{array} \right.$$

No closed-form solution

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N \frac{\partial E}{\partial y_n} \frac{\partial y_n}{\partial a_n} \nabla a_n = \sum_{n=1}^N (y_n - t_n) \phi_n = \sum_{n=1}^N (\sigma(\mathbf{w}^T \phi_n) - t_n) \phi_n$$



Course Review

LEARNING ALGORITHMS

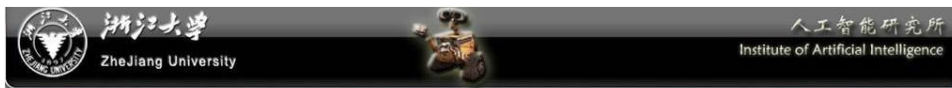
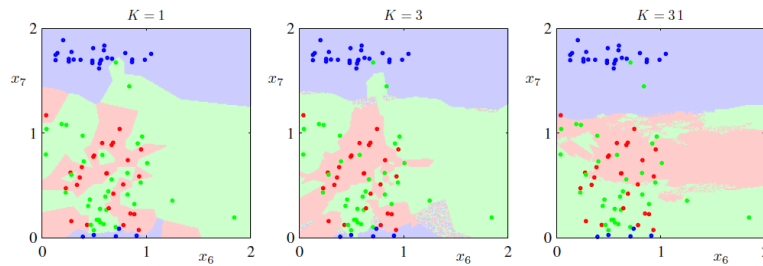


Nearest-neighbour methods

- KNN density estimation
 - K govern the radius of the sphere
- KNN classifier

$$p(\mathbf{x}|C_k) = \frac{K_k}{N_k V} \quad p(\mathbf{x}) = \frac{K}{NV} \quad p(C_k) = \frac{N_k}{N}$$

$$\Rightarrow p(C_k|\mathbf{x}) = \frac{p(\mathbf{x}|C_k)p(C_k)}{p(\mathbf{x})} = \frac{K_k}{K}$$



Sequential learning

- Stochastic gradient descent (sequential gradient descent)

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n \quad n = 1, 2, \dots, N$$

Learning rate

Error function

- least-mean-squares or the LMS algorithm

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 \quad \Rightarrow \quad E_n(\mathbf{w}) = \frac{1}{2} \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2$$

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n \quad \Rightarrow \quad \mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} + \eta (t_n - \mathbf{w}^{(\tau)T} \phi_n) \phi_n$$

Maximum margin classifiers

- Support Vector Machines (SVM) learning algorithm:**

1. Choose a kernel function, e.g. Gaussian kernel function.

2. Use SMO algorithm to solve $\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$

3. Select support vectors with $a_n > 0$

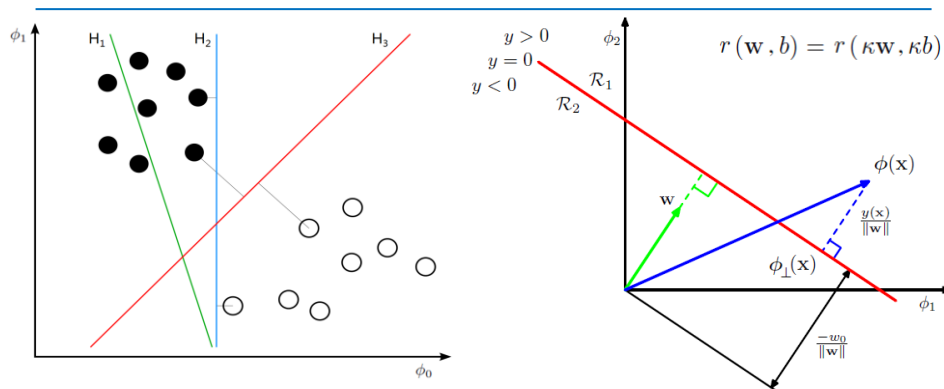
4. Compute the threshold parameter b by using SV set: $b = \frac{1}{N_S} \sum_{n \in S} \left(t_n - \sum_{m \in S} a_m t_m k(\mathbf{x}_n, \mathbf{x}_m) \right)$

5. Use SV set to classify new data point: $y(\mathbf{x}) = \sum_{m \in S} a_m t_m k(\mathbf{x}_n, \mathbf{x}_m) + b$

Maximum margin classifiers

- For the two-class classification problem using linear models: $y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$
 - Assume the training data set is linearly separable in **feature space**;
 - And all data points are correctly classified, so that: $t_n y(\mathbf{x}_n) > 0$ $t_n \in \{-1, 1\}$
- The definition of the **Margin** (rescaling \mathbf{w} and b doesn't change r)

$$r(\mathbf{w}, b) = \min_n \left\{ \frac{t_n y(\mathbf{x}_n)}{\|\mathbf{w}\|} \right\} = \min_n \left\{ \frac{t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|} \right\}$$



Maximum margin classifiers

- **Optimization problem:** find the solution of the maximum margin

$$\arg \max_{\mathbf{w}, b} r(\mathbf{w}, b) \longleftrightarrow \arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n [t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b)] \right\}$$

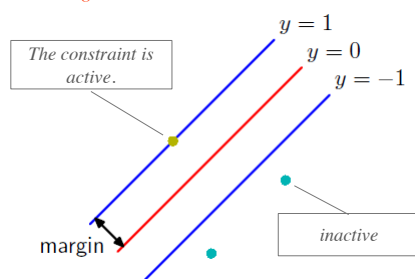
- Because $r(\mathbf{w}, b) = r(\kappa \mathbf{w}, \kappa b)$, so we can set $t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) = 1$ for the point that is closest to the decision surface. Then, all data points will satisfy the constraints: $t_n y(\mathbf{x}_n) \geq 1$

- **Equivalent constrained optimization problem:**

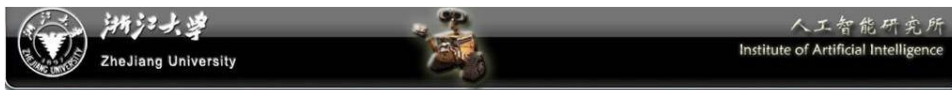
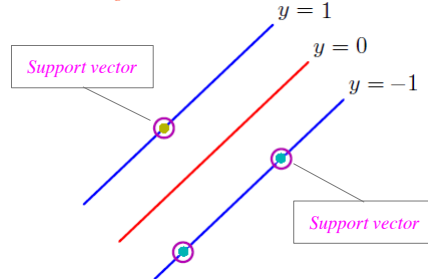
$$\arg \max_{\mathbf{w}, b} r(\mathbf{w}, b) \longleftrightarrow \arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \right\} \longleftrightarrow \arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

subject to $t_n y(\mathbf{x}_n) \geq 1, n = 1, \dots, N.$

The margin has not been maximized.



The margin has been maximized.



K-means clustering

- Distortion measure (responsibilities):

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2 \quad \xrightarrow{\frac{\partial J}{\partial \boldsymbol{\mu}_k} = 2 \sum_{n=1}^N r_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k) = 0} \quad \boldsymbol{\mu}_k = \frac{\sum_n r_{nk} \mathbf{x}_n}{\sum_n r_{nk}}$$

Responsibilities

Data

Prototypes
(expected value)

$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \|\mathbf{x}_n - \boldsymbol{\mu}_j\|^2 \\ 0 & \text{otherwise.} \end{cases}$$

$$r_{n,k} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix}$$

Example: 5 data points
and 3 clusters

K-means algorithm (batch version):

1. Pick number of clusters k
2. Randomly scatter k “cluster centers” in data space
3. Repeat:
 - a. Assign each data point to its closest cluster center
 - b. Move each cluster center to the mean of the points assigned to it



K-means clustering

- Online k-means algorithm (sequential k-means):

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$$

$$\boldsymbol{\mu}_k = \frac{\sum_n r_{nk} \mathbf{x}_n}{\sum_n r_{nk}} \quad \rightarrow \quad \boldsymbol{\mu}_k^{\text{new}} = \boldsymbol{\mu}_k^{\text{old}} + \eta_n (\mathbf{x}_n - \boldsymbol{\mu}_k^{\text{old}})$$

*The nearest
prototype to \mathbf{x}_n*

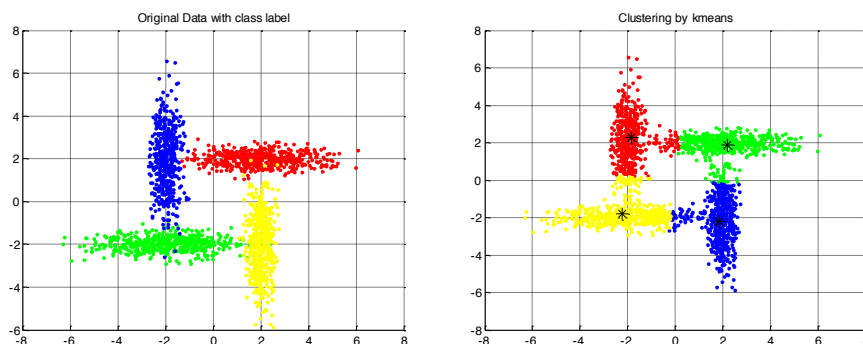
- K-medoids algorithm:
 - Chooses input data points as centers;
 - Works with an arbitrary matrix of distances between data points instead of Euclidean distance.
 - E.g. Manhattan distance or Minkowski distance

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2 \quad \rightarrow \quad \tilde{J} = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \mathcal{V}(\mathbf{x}_n, \boldsymbol{\mu}_k)$$



The limitation of K-means clustering

- The K-means algorithm adopts the hard assignment and doesn't consider the data density and probabilistic distribution.



Expectation-Maximization algorithm for GMM

EM for Gaussian Mixtures

1. Initialize the means μ_k , covariances Σ_k and mixing coefficients π_k , and evaluate the initial value of the log likelihood.

2. **E step.** Evaluate the responsibilities using the current parameter values

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_n | \mu_j, \Sigma_j)}$$

3. **M step.** Re-estimate the parameters using the current responsibilities

$$\mu_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n$$

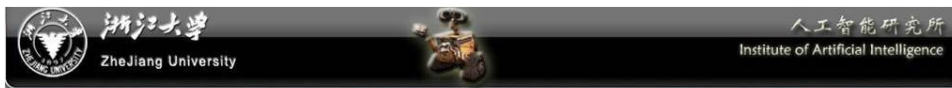
$$\Sigma_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \mu_k^{\text{new}}) (\mathbf{x}_n - \mu_k^{\text{new}})^T$$

$$\pi_k^{\text{new}} = \frac{N_k}{N} \quad \text{where} \quad N_k = \sum_{n=1}^N \gamma(z_{nk}).$$

4. Evaluate the log likelihood

$$\ln p(\mathbf{X} | \mu, \Sigma, \pi) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k) \right\}$$

and check for convergence of either the parameters or the log likelihood. If the convergence criterion is not satisfied return to step 2.



The general EM algorithm

The General EM Algorithm

Given a joint distribution $p(\mathbf{X}, \mathbf{Z} | \theta)$ over observed variables \mathbf{X} and latent variables \mathbf{Z} , governed by parameters θ , the goal is to maximize the likelihood function $p(\mathbf{X} | \theta)$ with respect to θ .

1. Choose an initial setting for the parameters θ^{old} .

2. **E step** Evaluate $p(\mathbf{Z} | \mathbf{X}, \theta^{\text{old}})$.

3. **M step** Evaluate θ^{new} given by $\theta^{\text{new}} = \arg \max_{\theta} Q(\theta, \theta^{\text{old}})$

$$\text{where} \quad Q(\theta, \theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z} | \theta). \quad \Rightarrow \quad Q(\theta, \theta^{\text{old}}) + \ln p(\theta)$$

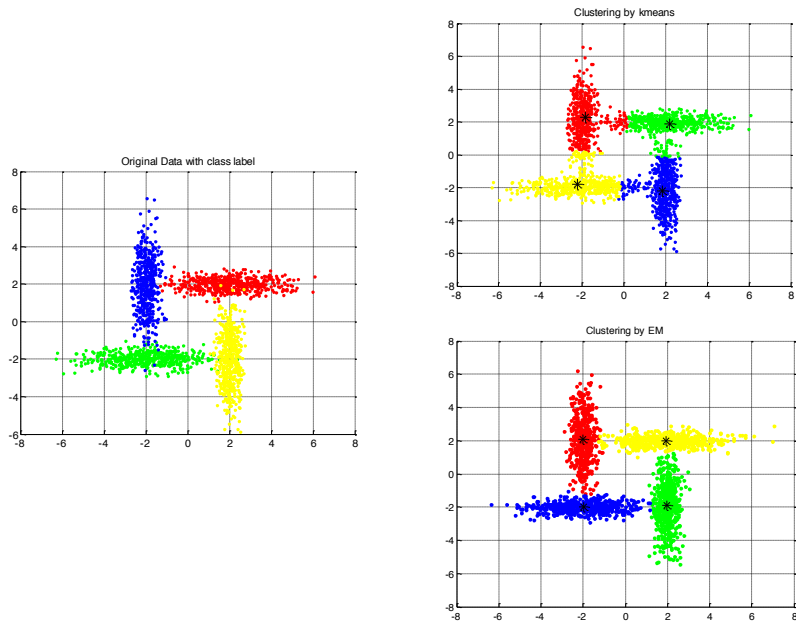
EM algorithm can be used to find MAP solution

4. Check for convergence of either the log likelihood or the parameter values. If the convergence criterion is not satisfied, then let

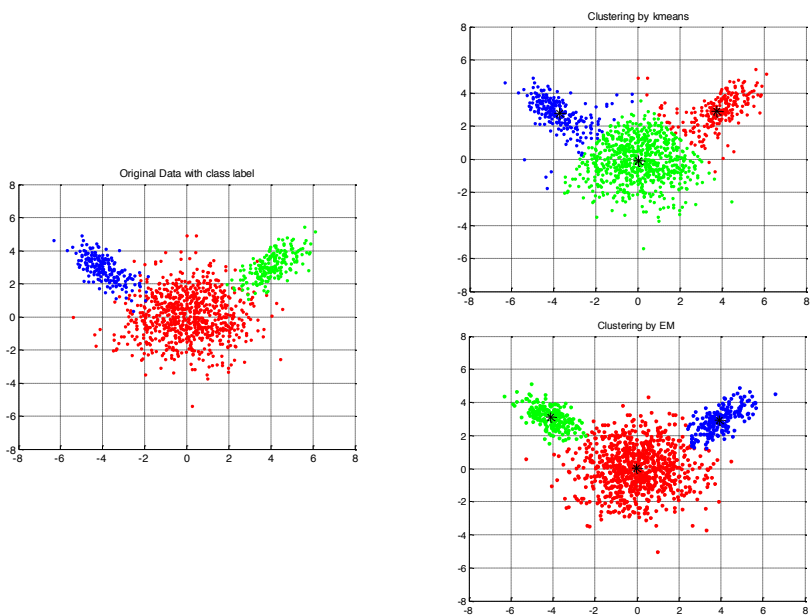
$$\theta^{\text{old}} \leftarrow \theta^{\text{new}}$$

and return to step 2.

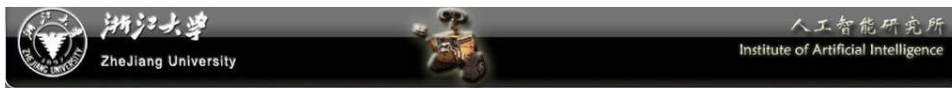
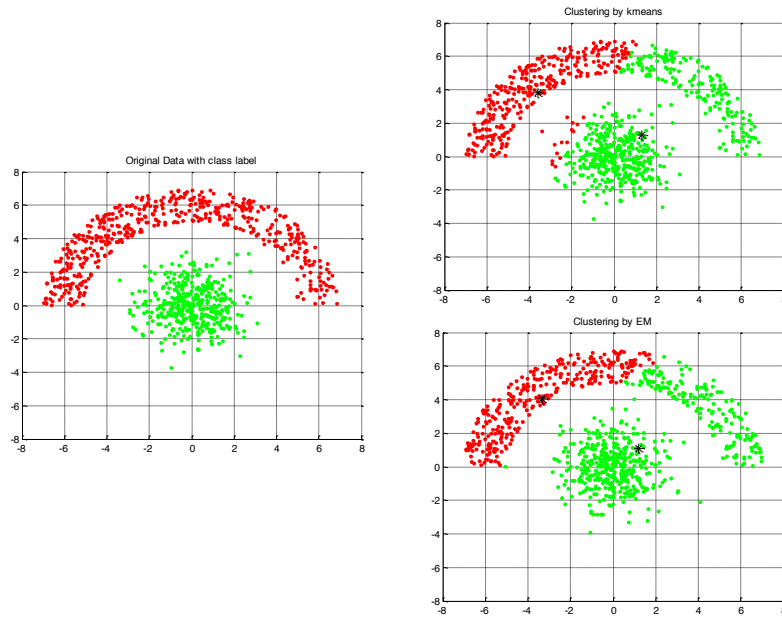
EM for GMM vs. K-means



EM for GMM vs. K-means

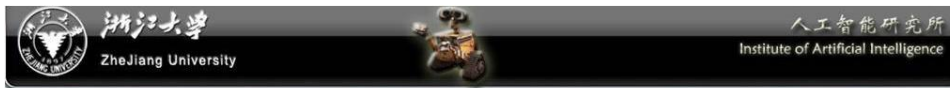


EM for GMM vs. K-means



Applications of PCA

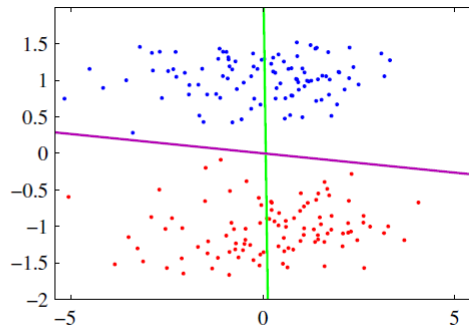
- Dimensionality reduction
 - Avoid the curse of dimensionality
- Lossy data compression
- Feature extraction
- Data visualization
 - How to visualize high-dimensional data?
- ...



Applications of PCA

- PCA vs. Fisher's LDA:

A comparison of principal component analysis with Fisher's linear discriminant for linear dimensionality reduction. Here the data in two dimensions, belonging to two classes shown in red and blue, is to be projected onto a single dimension. PCA chooses the direction of maximum variance, shown by the magenta curve, which leads to strong class overlap, whereas the Fisher linear discriminant takes account of the class labels and leads to a projection onto the green curve giving much better class separation.



<http://www.face-rec.org/algorithms/PCA/jcn.pdf>

<http://www.cs.jhu.edu/~hager/Public/teaching/cs461/pami97-eigenfaces.pdf>



Modeling nonlinear manifolds

- Two nonprobabilistic methods for dimensionality reduction and data visualization:
 - *Isometric feature mapping (ISOMAP): global method*
 - project the data to a lower-dimensional space using MDS, but where the dissimilarities are defined in terms of the *geodesic distances* measured along the manifold.
 - *Locally linear embedding (LLE): local method*
 - Map the high-dimensional data points down to a lower dimensional space while preserving coefficients.

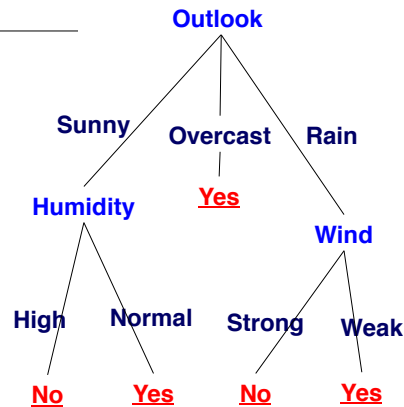
1. Tenenbaum J.B., Silva V. De , Langford J. C., A global geometric framework for nonlinear dimensionality reduction, *Science*, 2000, 290 (5500): 2219-2323
2. Sam Roweis, Lawrence Saul, Nonlinear dimensionality reduction by locally linear embedding, *Science*, 2000, 290(5500):2323-2326



Basic Decision Trees Learning Algorithm

- Data is processed in Batch (i.e., all the data is available).
- Recursively build a decision tree top-down.

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No



Information Gain

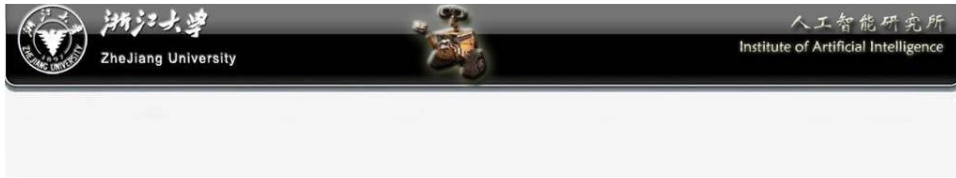
- The information gain of an attribute a is the expected reduction in entropy caused by partitioning on this attribute.

$$Gain(S, a) = Entropy(S) - \sum_{v \in values(s)} \frac{|S_v|}{|S|} Entropy(S_v)$$

Where S_v is the subset of S for which attribute a has value v

and the entropy of partitioning the data is calculated by weighing the entropy of each partition by its size relative to the original set

Partitions of low entropy lead to high gain



Course grade: 40% on homework + 60% on final exam

GOOD LUCK!

