

$$p(\omega_j)$$

Exercises

$$p(\omega_j|x)$$



1. How to use the **prior** and **likelihood** to calculate the **posterior**? What is the formula?

We can use the Bayes' formula to answer the question:

$$P(\omega_j|x) = \frac{p(x|\omega_j)P(\omega_j)}{p(x)}, \quad (1)$$

where in this case of two categories

$$p(x) = \sum_{j=1}^2 p(x|\omega_j)P(\omega_j). \quad (2)$$

Bayes' formula can be expressed informally in English by saying that

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}. \quad (3)$$



2. What's the difference in the ideas of the **minimum error Bayesian decision** and **minimum risk Bayesian decision**? What's the condition that makes the minimum error Bayesian decision identical to the minimum risk Bayesian decision?

the minimum error Bayesian decision: to minimize the classification error of the Bayesian decision.

the minimum risk Bayesian decision: to minimize the risk of the Bayesian decision.

~~$$R(\omega_1|x) = \lambda_{11}P(\omega_1|x) + \lambda_{12}P(\omega_2|x)$$~~

~~$$R(\omega_2|x) = \lambda_{21}P(\omega_1|x) + \lambda_{22}P(\omega_2|x)$$~~

~~$$\text{if } R(\omega_1|x) < R(\omega_2|x)$$~~

~~"decide ω_1 " is taken~~

~~$$(\lambda_{21} - \lambda_{11})P(\omega_1|x) > (\lambda_{12} - \lambda_{22})P(\omega_2|x)$$~~

~~"decide ω_1 " is taken~~
~~Condition. factor $(\lambda_{21} - \lambda_{11})$, $(\lambda_{12} - \lambda_{22})$ both are positive, and $(\lambda_{12} - \lambda_{22}) > (\lambda_{21} - \lambda_{11})$~~


3. A person takes a lab test of nuclear radiation and the result is positive. The test returns a correct positive result in 99% of the cases in which the nuclear radiation is actually present, and a correct negative result in 95% of the cases in which the nuclear radiation is not present. Furthermore, 3% of the entire population are radioactively contaminated. Is this person contaminated?

• **Solution:**

Given: $P(+|\text{contaminated})=0.99$ $P(-|\text{no contaminated})=0.95$

$P(\text{contaminated})=0.03$ $P(\text{no contaminated})=0.97$

Compute: $P(\text{no contaminated} | +)$, $P(\text{contaminated} | +)$,

$P(\text{contaminated} | +) = P(+|\text{contaminated}) * P(\text{contaminated}) / p(+) = 0.99 * 0.03 / p(+) = 0.0297 / p(+)$

$P(\text{no contaminated} | +) = P(+|\text{no contaminated}) * P(\text{no contaminated}) / P(+)$

$= (1 - 0.95) * (1 - 0.03) / P(+) = 0.0485 / P(+)$

Since $P(\text{no contaminated} | +) > P(\text{contaminated} | +)$, we decide that the patient does not have cancer (Bayes decision rule)



4. Please present the basic ideas of the **maximum likelihood estimation** method and **Bayesian estimation** method. When do these two methods have similar results?

(1) General principle

Suppose that \mathcal{D} contains n samples, $\mathbf{x}_1, \dots, \mathbf{x}_n$. Then, since the samples were drawn independently, we have

$$p(\mathcal{D}|\boldsymbol{\theta}) = \prod_{k=1}^n p(\mathbf{x}_k|\boldsymbol{\theta}). \quad (1)$$

For analytical purposes, it is usually easier to work with the logarithm of the likelihood than with the likelihood itself. If the number of parameters to be set is p , then we let $\boldsymbol{\theta}$ denote the p -component vector $\boldsymbol{\theta} = (\theta_1, \dots, \theta_p)^t$, and $\nabla_{\boldsymbol{\theta}}$ be the gradient operator

$$\nabla_{\boldsymbol{\theta}} \equiv \begin{bmatrix} \frac{\partial}{\partial \theta_1} \\ \vdots \\ \frac{\partial}{\partial \theta_p} \end{bmatrix}. \quad (2)$$

We define $l(\boldsymbol{\theta})$ as the *log-likelihood* function*

$$l(\boldsymbol{\theta}) \equiv \ln p(\mathcal{D}|\boldsymbol{\theta}). \quad (3)$$

We can then write our solution formally as the argument $\boldsymbol{\theta}$ that maximizes the log-likelihood, i.e.,

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} l(\boldsymbol{\theta}), \quad (4)$$

where the dependence on the data set \mathcal{D} is implicit. Thus we have from Eq. 1

$$l(\boldsymbol{\theta}) = \sum_{k=1}^n \ln p(\mathbf{x}_k|\boldsymbol{\theta}) \quad (5)$$

and

$$\nabla_{\boldsymbol{\theta}} l = \sum_{k=1}^n \nabla_{\boldsymbol{\theta}} \ln p(\mathbf{x}_k|\boldsymbol{\theta}). \quad (6)$$

Thus, a set of necessary conditions for the maximum likelihood estimate for $\boldsymbol{\theta}$ can be obtained from the set of p equations

$$\boxed{\nabla_{\boldsymbol{\theta}} l = \mathbf{0}.} \quad (7)$$

If all solutions are found, we are guaranteed that one represents the true maximum, though we might have to check each solution individually (or calculate second derivatives) to identify which is the global optimum.

(2)Goal: compute $P(\omega_i | \mathbf{x}, \mathcal{D})$

Given the sample \mathcal{D} , Bayes formula can be written

$$P(\omega_i | \mathbf{x}, \mathcal{D}) = \frac{p(\mathbf{x} | \omega_i, \mathcal{D}).P(\omega_i | \mathcal{D})}{\sum_{j=1}^c p(\mathbf{x} | \omega_j, \mathcal{D}).P(\omega_j | \mathcal{D})}$$

To demonstrate the preceding equation, use:

$$D = D_1 \cup D_2 \dots \cup D_c \quad x \in D_i \rightarrow x \text{ is } \omega_i$$

D_i has no influence on $p(x | \omega_j, D_j)$ if $i \neq j$

$P(\omega_i) = P(\omega_i | D)$ (Training sample provides this!)

Thus :

$$P(\omega_i | x, D) = \frac{P(x | \omega_i, D_i) \cdot P(\omega_i)}{\sum_{j=1}^c P(x | \omega_j, D_j) \cdot P(\omega_j)}$$

Parameter Distribution

$P(x)$ is unknown, we assume it has a known parametric form $P(x|\theta)$, and value of parameter θ is unknown

Know prior density $P(\theta)$

Training data convert $P(\theta)$ to a posterior density $P(\theta|D)$

Our path:

$$\begin{aligned} p(x | \omega, D) &= p(x) \cong p(x | D) \\ &= \int p(x, \theta | D) d\theta \\ &= \int p(x | \theta) p(\theta | D) d\theta \end{aligned}$$

In virtually every case, maximum likelihood and Bayes solutions are equivalent in the asymptotic limit of infinite training data.



5. Please present the nature of **principal component analysis**.

- 1) PCA is an unsupervised method.
- 2) PCA is a good dimension reduction method. Usually, even if the dimension of the sample is greatly reduced, the main information can be still stored.
- 3) PCA is a de-correlation method. After the PCA transform, the obtained components are statistically uncorrelated.
- 4) PCA can be used as a compression method. PCA can achieve a high compression ratio.
- 5) PCA is the optimal representation method, which allows us to obtain the minimum reconstruction error.
- 6) As the transform axes of PCA are orthogonal, it is also referred to as an orthogonal transform method.

Disadvantage: need to a training phase.

FDA



6. Describe the basic idea and possible advantage of **Fisher discriminant analysis**.

Basic idea:

- 1) Fisher discriminant analysis's task is to find the best such direction w to obtain accurate classification.
- 2) Fisher discriminant analysis is to maximize the between-class distance and minimize the within-class distance
- 3) Fisher discriminant analysis exploits the training sample to produce transform axes.

Advantage:

- 1) The time complex is $O(d^2n)$. It is linear.
- 2) Fisher discriminant analysis is a linear discriminant method, thus it is easy to generalize.
- 3) Fisher discriminant analysis is supervised; it can exploit the training sample to produce the transform axes.



7. What is the **K nearest neighbor classifier**? Is it reasonable?

We can center a cell about x and let it grow until it captures kn samples, where kn is some specified function of n . These samples are the kn nearest-neighbors of x . If the density is high near x , the cell will be relatively small, which leads to good resolution. If the density is low, it is true that the cell will grow large, but it will stop soon after it enters regions of higher density.

Classify x by assigning it the label most frequently represented among the k nearest

samples and use a voting scheme

It is reasonable. The nearest-neighbor rule leads to an error rate greater than the minimum possible value of the Bayes rate; If the number of prototype is large (unlimited), the error rate of the nearest-neighbor classifier is never worse than twice the Bayes rate.

8. Is it possible that a classifier can obtain a higher accuracy for any dataset than any other classifiers?

No. According to the No Free Lunch Theorem, there are no context-independent or problem-independent reasons to favor one learning or classification method over another.

If one algorithm seems to outperform another in a particular situation, it is a consequence of it fit to the particular pattern recognition problem, not the general superiority of the algorithm.

9. Please describe the over-fitting problem.

Over-fitting generally occurs when a model is excessively complex, such as having too many parameters relative to the number of observations. A model which has been over-fit will generally have poor predictive performance, as it can exaggerate minor fluctuations in the data.

- 1) Over-fitting: The classification performs good on training dataset but bad on test dataset.
- 2) Reason: (i)The training sample is inadequate;(ii)There are difference between training samples and test samples;(iii)The classification focus on detail;(iv)Samples exist noise.
- 3) How to handle: (i)reduce dimensionality(ii)suppose all classes share same covariance matrix(iii) look for a better estimate for covariance matrix.

10. Usually a more complex learning algorithm can obtain a higher accuracy in the training stage. So, should a more complex learning algorithm be favored ?

No context-independent or usage-independent reasons to favor one learning or classification method over another to obtain good generalization performance.

When confronting a new pattern recognition problem, we need focus on the aspects — prior information, data distribution, amount of training data and cost or reward functions.

Ugly Duckling Theorem: an analogous theorem, addresses features and patterns. shows that in the absence of assumptions we should not prefer any learning or classification algorithm over another.

11. Under the condition that the number of the training samples approaches to the infinity, the estimation of the mean obtained using Bayesian estimation method is almost identical to that obtained using the maximum likelihood estimation method. Is this statement correct ?

For the reasonable prior distributions that does not preclude the true solution: Maximum likelihood and Bayesian solution are equivalent in the asymptotic limit of infinite training data.

12. Can the minimum squared error procedure be used for binary classification ?

Yes, the minimum squared error procedure can be used for binary classification.

$$\begin{pmatrix} y_{10} & y_{11} & \cdots & y_{1d} \\ y_{20} & y_{21} & \cdots & y_{2d} \\ \cdots & \cdots & \cdots & \cdots \\ y_{n0} & y_{n1} & \cdots & y_{nd} \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \\ \cdots \\ a_d \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \cdots \\ b_n \end{pmatrix}$$

$$Ya = b, \quad Y = \begin{bmatrix} Y_1^T \\ \vdots \\ Y_n^T \end{bmatrix}, \quad Y_i^T = \begin{bmatrix} y_{i0} \\ \vdots \\ y_{id} \end{bmatrix}.$$

A simple way to set b : if Y_i is from the first class, then b_i is set to 1; if Y_i is from the second class, then b_i is set to -1.

Another simple way to set b : if Y_i is from the first class, then b_i is set to $\frac{n}{n_1}$; if Y_i is from the second class, then b_i is set to $-\frac{n}{n_2}$.

13. Can you devise a minimum squared error procedure to perform multiclass classification ?

Yes, a minimum squared error procedure to perform multiclass classification.

■ Generalization for MSE Procedure
consider multicategory case as a set of c two-class problem

$a_i^t y = 1$ for all $y \in Y_i$
 $a_i^t y = 0$ for all $y \notin Y_i$

$$A = [a_1 \ a_2 \ \dots \ a_c] = \begin{bmatrix} a_{11} & a_{21} & \dots & a_{c1} \\ a_{12} & a_{22} & \dots & a_{c2} \\ \dots & \dots & \dots & \dots \\ a_{1d} & a_{2d} & \dots & a_{cd} \end{bmatrix}$$

$$Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \dots \\ Y_c \end{bmatrix} = \begin{bmatrix} y_{111} & y_{112} & \dots & y_{11d} \\ y_{121} & y_{122} & \dots & y_{12d} \\ \dots & \dots & \dots & \dots \\ y_{211} & y_{212} & \dots & y_{21d} \\ y_{221} & y_{222} & \dots & y_{22d} \\ \dots & \dots & \dots & \dots \\ y_{c11} & y_{c12} & \dots & y_{c1d} \\ y_{c21} & y_{c22} & \dots & y_{c2d} \end{bmatrix}$$

$$B = \begin{bmatrix} B_1 \\ B_2 \\ \dots \\ B_c \end{bmatrix} = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 1 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 1 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 1 \\ 0 & 0 & \dots & 1 \\ \dots & \dots & \dots & \dots \end{bmatrix}$$

$YA = B$
 $A = Y^+ B$

8

where $Y_i = \{y_{i11}, y_{i12}, \dots, y_{i1d}\}$ $a_i^t = [a_{i1}, a_{i2}, \dots, a_{id}]^t$, $b_i = [0 \ 0 \ \dots i \ \dots 0]$
 $y_{i21}, y_{i22}, \dots, y_{i2d}$



14. Which kind of applications is the Markov model suitable for ? (了解) p105

In problems that have an inherent temporality (时间性) — that is, consist of a process that unfolds in time — we may have states at time t that are influenced directly by a state at $t - 1$.

Applications: speech recognition, gesture recognition, parts of speech tagging (词性标注) and DNA sequencing

HMM: evaluation, decoding, learning problems

15. For minimum squared error procedure based on $Y a = b$ (Y is the matrix consisting of all the training samples), if we have proper b and criterion function, then this minimum squared error

procedure might be equivalent to Fisher discriminant analysis. Is this presentation correct ?

Yes,

$$\mathbf{b} = \begin{bmatrix} \frac{n}{n_1} \mathbf{1}_1 \\ \frac{n}{n_2} \mathbf{1}_2 \end{bmatrix}.$$

We shall now show that this special choice for \mathbf{b} links the MSE solution to Fisher's linear discriminant.

$$\mathbf{w} = \alpha n \mathbf{S}_W^{-1} (\mathbf{m}_1 - \mathbf{m}_2), \quad (54)$$

which, except for an unimportant scale factor, is identical to the solution for Fisher's linear discriminant. In addition, we obtain the threshold weight w_0 and the following decision rule: Decide ω_1 if $\mathbf{w}^t(\mathbf{x} - \mathbf{m}) > 0$; otherwise decide ω_2 .

16. Suppose that the number of the training samples approaches to the infinity, then the minimum error Bayesian decision will perform better than any other classifier achieving a lower classification error rate. Do you agree on this?

No.

According to the No Free Lunch Theorem, there are no context-independent or problem-independent reasons to favor one learning or classification method over another.

If one algorithm seems to outperform another in a particular situation, it is a consequence of it fit to the particular pattern recognition problem, not the general superiority of the algorithm.

Yes, it is correct. The premise is the number of the training samples approaches to the infinity.

How to prove?

17. What are the upper and lower bound of the classification error rate of the K nearest neighbor classifier ?

(1) The nearest-neighbor rule leads to an error rate greater than the minimum possible value of the Bayes rate

(2) If the number of prototype is large (unlimited), the error rate of the nearest-neighbor classifier is never worse than twice the Bayes rate (it can be demonstrated!)

$$P^*(e | x) = 1 - P(\omega_m | x)$$

$$P^*(e) = \int P^*(e | x) p(x) dx$$

$$P(e) = \lim_{n \rightarrow \infty} P_n(e)$$

$$P^*(e) \leq P(e) \leq P^*(e) \left(2 - \frac{c}{c-1} P^*(e) \right)$$

$$P^* \leq P \leq P^* \left(2 - \frac{c}{c-1} P^* \right).$$

Here we show that

If different sets of n samples are used to classify \mathbf{x} , different vectors \mathbf{x}_n will be obtained for the nearest-neighbor of \mathbf{x} . Since the decision rule depends on this nearest-neighbor, we have a conditional probability of error $P(e | \mathbf{x}, \mathbf{x}_n)$ that depends on both \mathbf{x} and \mathbf{x}_n . By averaging over \mathbf{x}_n , we

$$P(e|x) = \int P(e|x, x') p(x'|x) dx'$$

obtain

As n goes to infinity we expect $p(x|x)$ to approach a delta function centered at x . Assume that at the given x , $p(\cdot)$ is continuous and not equal to zero. the probability that any sample falls within a hypersphere S centered about x is some positive number P_s :

$$P_s = \int_{x' \in S} p(x') dx'$$

Thus x converges to x in probability, and $p(x|x)$ approaches a delta function, as expected.

18. Can you demonstrate that a **statistics-based classifier** usually cannot lead to a classification accuracy of 100% ?

- (1) The training sample is limited, but the real data is infinite and we cannot get all the samples.
- (2) We cannot ensure the correct of the training sample. It may has noise.

19. What is **representation-based classification**? Please present the characteristics of Representation-based classification.

(1) The representation-based methods attempts to exploit the training samples to provide local representation for the test sample. Typical examples include sparse representation.

(2) The method first expresses the test sample as a linear combination of a number of training samples and exploits the representation result to perform classification. In other words, it is regarded that the test sample is related with each training sample. The test sample will be classified into the class that has the most "relationship".

Basic idea $a_1 \sim 0; a_i = 0$

Testing sample: y

Training samples: x_1, \dots, x_n Suppose that $y = a_1 x_1 + \dots + a_n x_n$

$$A = (X^T X)^{-1} X^T y \quad \text{or} \quad A = (X^T X + uI)^{-1} X^T y$$

Suppose that training samples x_p, \dots, x_q are from the s -th class.

The representation error of the s -th class is evaluated by $E_s = \|y - a_p x_p - \dots - a_q x_q\|$

(3) The method considers that testing sample y is from the class that has the smallest representation error.

(4) A good performance can be achieved.

20. A simple **representation-based classification** method is presented as follows:

This method seeks to represent the test sample as a linear combination of all training samples and uses the representation result to classify the test sample:

$$y = b_1 \tilde{x}_1 + \dots + b_M \tilde{x}_M, \quad (1)$$

where \tilde{x}_i ($i = 1, 2, \dots, M$) denote all the training samples and b_i ($i = 1, 2, \dots, M$) are the coefficients.

We rewrite Eq.(1) into $y = \tilde{X} B$,

where $B = [b_1 \dots b_M]^T$, $\tilde{X} = [\tilde{x}_1 \dots \tilde{x}_M]$. If \tilde{X} is not singular,

we can solve B using $B = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T y$; otherwise, we can solve it using

$$B = (\tilde{X}^T \tilde{X} + \mathcal{M})^{-1} \tilde{X}^T y, \quad (3)$$

$$a = (Y^T Y)^{-1} Y^T b$$

where γ is a positive constant and I is the identity matrix. After we obtain B , we refer to $\tilde{X}B$ as the representation result of our method. We can convert the representation result into a two-dimensional image having the same size of the original sample image.

We exploit the sum of the contribution, to representing the test sample, of the training samples from a class, to classify the test sample. For example, if all the training samples from the r th ($r \in C$) class are $\tilde{x}_s, \dots, \tilde{x}_t$, then the sum of the contribution, to representing the test sample, of the

r th class will be

$$g_r = a_s \tilde{x}_s + \dots + a_t \tilde{x}_t. \quad (4)$$

We calculate the deviation of g_r from y using

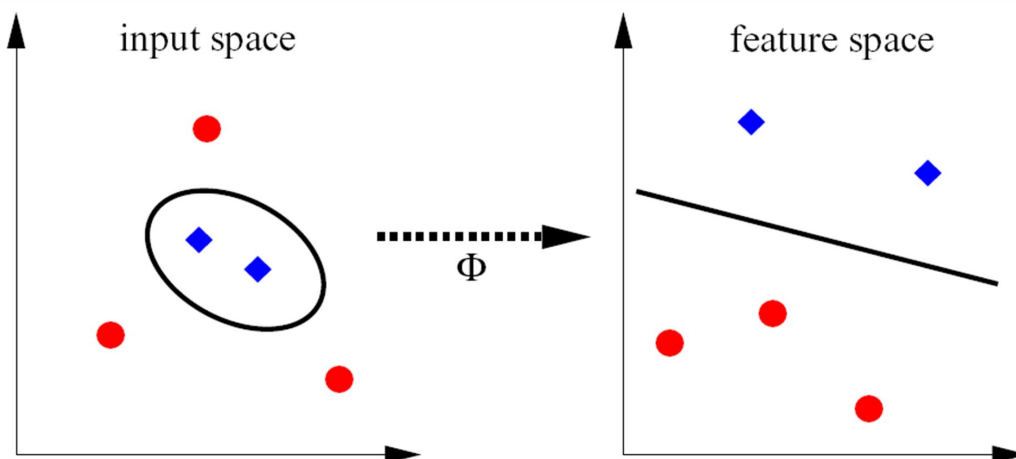
$$D_r = \|y - g_r\|^2, r \in C. \quad (5)$$

We can also convert g_r into a two-dimensional matrix having the same size of the original sample image. If we do so, we refer to the matrix as the two-dimensional image corresponding to the contribution of the r th class. The smaller the deviation D_r , the greater the contribution to representing the test sample of the r th class. In other words, if $D_q = \min D_r (q, r \in C)$, the test sample will be classified into the q th class.

From the above presentation, we know that representation-based classification method is a novel method and totally different from previous classifiers ! It performs very well in image-based classification, such as face recognition and palmprint recognition.

We should understand its nature and advantages.

Please describe the difference between linear and nonlinear discriminant functions? What potential advantage does nonlinear discriminant function have in comparison with linear discriminant function?



The above figure is just auxiliary for the question !

It was agreed that simplicity is generally best. The use of linear method is recommended wherever

possible. It also was agreed that nonlinear methods in some applications can provide better result, particularly with complex and/or other very large data.

A linear discriminant function is defined by a hyperplane's vector w and offset b , the decision boundary is $\{x | wx+b=0\}$.

Linear discriminant functions are generally more robust than nonlinear discriminant functions. Because of having only limited flexibility and less prone to overfitting.

A linear algorithm is applied in some appropriate(kernel) feature space.

The linear classification in feature space corresponds to a (powerful) non-linear decision function in input space.

The non-linear methods often involve a number of parameters.

Linear methods seem ideal when limited data and limited knowledge about the data is available. If there are large amount of data, nonlinear methods are suitable to find potentially more complex structure in a data.

22. What is the naïve Bayes rule ?

The naïve Bayes rule is a simple probabilistic classifier based on applying Bayes theorem with strong independence assumption that all attributions are independent another.

The procedures are shown as follows:

(1) Given unknown tuple $X=\{x_1, x_2, \dots, x_n\}$, where x_i are attributions.

(2) Class $C=\{w_1, w_2, \dots, w_n\}$

(3) Calculate $P(w_1|X), P(w_2|X), \dots, P(w_n|X)$,

(4) If $P(w_k|X)=\text{Max}\{P(w_1|X), P(w_2|X), \dots, P(w_n|X)\}$ then decide X belong to w_k .

$$p(\omega_k|X) \propto \prod_{i=1}^d p(x_i|\omega_k).$$

$$\begin{aligned} p(w_k|x) &= p(x|w_k) \\ &= p(x_1) \end{aligned}$$

23. What is the difference between supervised and unsupervised learning methods? Please show two examples of supervised and unsupervised learning methods.

In supervised learning, a teacher provides a category label or cost for each pattern in a training set, and we seek to reduce the sum of the costs for these patterns.

In unsupervised learning or clustering there is no explicit teacher, and the system forms clusters or "natural groupings" of the input patterns. "Natural" is always defined explicitly or implicitly in the clustering system itself, and given a particular set of patterns or cost function, different clustering algorithms lead to different clusters. Often the user will set the hypothesized number of different clusters ahead of time.

For example, for supervised learning, there are maximum likelihood estimation and Bayesian learning; for unsupervised learning, there are PCA and Kmeans.

24. In some special real-world classification applications, the Bayesian decision theory might perform badly. What are possible reasons ?

The possible reasons will be that we can not get all the knowledge of the probability structure of the problem or only get some fuzzy or general knowledge about the problem.

25. Suppose that we are applying a linear discriminant function to a nonlinear separable problem, what means can we adopt to obtain an optimal solution?

(1) If problem is non-linear separable problem in low-dimension, we can transform it into high-dimension such that it is a linear separable.

Take an example.

The discriminant function $d(x)$ is quadratic polynomial function, $x=(x_1, x_2)^t$,

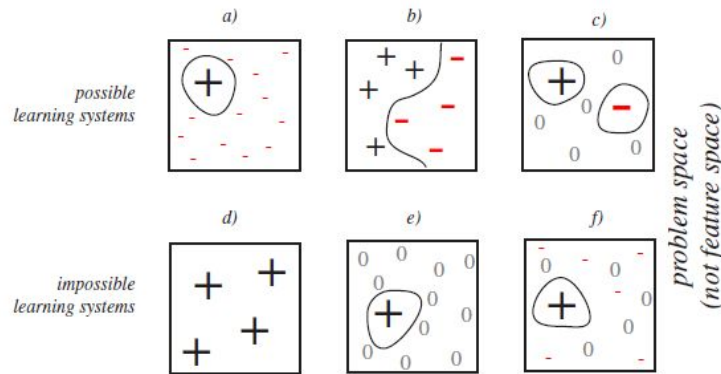
$$d(x) = w_{11}x_1^2 + w_{12}x_1x_2 + w_{22}x_2^2 + w_1x_1 + w_2x_2 + w_3$$

we can linearize it:

$$x^* = (x_1^2, x_1x_2, x_2^2, x_1, x_2, 1)^t \quad w = (w_{11}, w_{12}, w_{22}, w_1, w_2, w_3)^t$$

(2) According to the generalized linear discriminant function, we can transfer any low-dimension nonlinear function into high-dimension linear function. We can adopt kernel function to solve this problem such as KPCA, KMSE, SVM.

26. Please present possible generalization capability in the sample space of a method.



① The no free lunch theorem shows the generalization performance on off-training set data (d, e, f) that can be achieved (a, b, c) and also show the performance that can't be achieved.

② Each square represents all possible classification problem consistent with the training data.

③ "+" indicates that the classification algorithm has generalization higher than average, "-" indicates lower than average.

④ a) shows that it's possible for an algorithm to have high accuracy on a small set of problems so long as it has poor performance on all other problem.

b) shows that it's possible to have excellent performance throughout a large range of problem, but this will be balanced by very poor performance on other problems.

⑤ c) shows that it's possible to have higher-than-average performance on a small set of problems, but this will be balanced by poor performance on another small set of problems, and average performance on remaining.

⑥ It is impossible to have good performance throughout the full range of problems. It is also impossible to have higher-than-average performance on some problem while having average performance on every where else.



27. Apply model $Y\mathbf{a}=\mathbf{b}$ to perform classification.

$$\begin{pmatrix} y_{10} & y_{11} & \dots & y_{1d} \\ y_{20} & y_{21} & \dots & y_{2d} \\ \dots & \dots & \dots & \dots \\ y_{n0} & y_{n1} & \dots & y_{nd} \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \\ \dots \\ a_d \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \dots \\ b_n \end{pmatrix} \Leftrightarrow Y\mathbf{a} = \mathbf{b}$$

If we define the error vector \mathbf{e} by

$$\mathbf{e} = Y\mathbf{a} - \mathbf{b}$$

then one approach is to try to minimize the squared length of the error vector. This is equivalent to minimizing the sum-of-squared-error criterion function

$$J_s(\mathbf{a}) = \|\mathbf{Y}\mathbf{a} - \mathbf{b}\|^2 = \sum_{i=1}^n (\mathbf{a}^t \mathbf{y}_i - b_i)^2.$$

A simple closed-form solution can also be found by forming the gradient

$$\nabla J_s = \sum_{i=1}^n 2(\mathbf{a}^t \mathbf{y}_i - b_i) \mathbf{y}_i = 2\mathbf{Y}^t(\mathbf{Y}\mathbf{a} - \mathbf{b})$$

and setting it equal to zero. This yields the necessary condition

$$\mathbf{Y}^t \mathbf{Y} \mathbf{a} = \mathbf{Y}^t \mathbf{b}$$

and in this way we have converted the problem of solving $\mathbf{Y}\mathbf{a} = \mathbf{b}$ to that of solving

$$\mathbf{Y}^t \mathbf{Y} \mathbf{a} = \mathbf{Y}^t \mathbf{b}.$$

$$\mathbf{a} = (\mathbf{Y}^t \mathbf{Y})^{-1} \mathbf{Y}^t \mathbf{b}$$

$$\mathbf{a} = (\mathbf{Y}^t \mathbf{Y})^{-1} \mathbf{Y}^t \mathbf{b} \text{ is an MSE solution to } \mathbf{Y}\mathbf{a} = \mathbf{b}.$$

28. How to extend the **binary minimum squared error** procedure to the multiclass minimum squared error procedure?

Example of Linear Classifier by Pseudoinverse

- $\omega_1: (1,2)^t$ and $(2,0)^t$
- $\omega_2: (3,1)^t$ and $(2,3)^t$

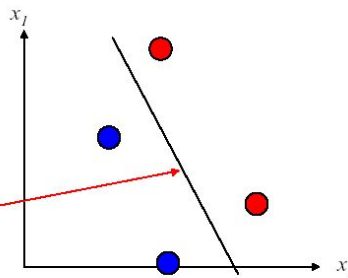
Sample Matrix ($d = 1+2, n = 4$)

$$Y = \begin{bmatrix} 1 & 1 & 2 \\ 1 & 2 & 0 \\ -1 & -3 & -1 \\ -1 & -2 & -3 \end{bmatrix}$$

Pseudo-inverse

$$Y^* = (Y^t Y)^{-1} Y^t = \begin{bmatrix} 5/4 & 13/12 & 3/4 & 7/12 \\ -1/2 & -1/6 & -1/2 & -1/6 \\ 0 & -1/3 & 0 & -1/3 \end{bmatrix}$$

$$\mathbf{a}^t \begin{pmatrix} 1 \\ x_1 \\ x_2 \end{pmatrix} = 0$$



$$\text{Assuming } \mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

$$\text{our solution is } \mathbf{a} = Y^t \mathbf{b} = \begin{bmatrix} 11/3 \\ -4/3 \\ -2/3 \end{bmatrix}$$

- Generalization for MSE Procedure

consider multicategory case as a set of c two-class problems

$$a_i^t y = 1 \text{ for all } y \in Y_i$$

$$a_i^t y = 0 \text{ for all } y \notin Y_i$$

$$A = [a_1 \ a_2 \ \dots \ a_c] = \begin{bmatrix} a_{11} & a_{21} & \dots & a_{c1} \\ a_{12} & a_{22} & \dots & a_{c2} \\ \dots & \dots & \dots & \dots \\ a_{1d} & a_{2d} & \dots & a_{cd} \end{bmatrix}$$

$$Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \dots \\ Y_c \end{bmatrix} = \begin{bmatrix} y_{111} & y_{112} & \dots & y_{11d} \\ y_{121} & y_{122} & \dots & y_{12d} \\ \dots & \dots & \dots & \dots \\ y_{211} & y_{212} & \dots & y_{21d} \\ y_{221} & y_{222} & \dots & y_{22d} \\ \dots & \dots & \dots & \dots \\ y_{c11} & y_{c12} & \dots & y_{c1d} \\ y_{c21} & y_{c22} & \dots & y_{c2d} \end{bmatrix}$$

$$B = \begin{bmatrix} B_1 \\ B_2 \\ \dots \\ B_c \end{bmatrix} = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 1 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 1 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 1 \\ 0 & 0 & \dots & 1 \\ \dots & \dots & \dots & \dots \end{bmatrix}$$

$$YA = B$$

$$A = Y^+ B$$