

Classification

Volker Blobel – University of Hamburg
March 2005

Given objects (e.g. particle tracks), which have certain features (e.g. momentum \mathbf{p} , specific energy loss dE/dx) and which belong to one of certain classes (e.g. particle types).

Classification = assignment of object to classes

1. Strategies for classification
2. Minimizing the probability for misclassification
3. Risk minimization

Christopher M. Bishop, Neural Networks for Pattern Recognition, Clarendon Press Oxford

Keys during display: enter = next page; → = next page; ← = previous page; home = first page; end = last page (index, clickable); C-← = back; C-N = goto page; C-L = full screen (or back); C-+ = zoom in; C-- = zoom out; C-0 = fit in window; C-M = zoom to; C-F = find; C-P = print; C-Q = exit.

Strategies for classification

In statistical classification problems an object belongs to one of several classes \mathcal{C}_k .

Each object has to be assigned to **one** of the classes on the basis of certain features of the object.

Example from particle physics analysis: identification of the particle type in a high energy particle reaction; here the classes are the particles types *pion*, *kaon*, *proton*

Classification strategies can be based on

- minimization of the probability of misclassification, or on
- risk minimization; this requirement applies to cases with rare classes, where rejecting the correct class has to be avoided.

Formally **a-priori probabilities** $P(\mathcal{C}_k)$ of an object belonging to each of the classes \mathcal{C}_k are introduced. This probability $P(\mathcal{C}_k)$ corresponds to the fractions of objects in each class in the limit of an infinite number of observations. If a classification is required, but nothing is known about the features of the objects, the best strategy is to assign it to the class having the highest a-priori probability. This minimizes the probability of misclassification, although some classifications must be wrong. In the particle identification example one would classify each particle as a *pion*, the most common particle created in a high energy reaction.

Information in terms of values x of certain features connected to each object is of statistical nature and has to be used to improve the classification. The distribution of values x follows a certain probability density function $p(x)$. The probability of x lying in an interval (a, b) of x is given by

$$P(x \in [a, b]) = \int_a^b p(x) \, dx .$$

If there are several feature variables x_1, x_2, \dots , they are grouped into a vector \mathbf{x} corresponding to a point in multidimensional feature space. The probability of a multidimensional \mathbf{x} lying in a region \mathcal{R} of \mathbf{x} -space is given by

$$P(\mathbf{x} \in \mathcal{R}) = \int_{\mathcal{R}} p(\mathbf{x}) \, d\mathbf{x} .$$

Class-conditional probability

The class-conditional probability density $p(\mathbf{x}|\mathcal{C}_k)$ is the density of \mathbf{x} *given* that it belongs to class \mathcal{C}_k . The unconditional density function $p(\mathbf{x})$ is the density function for \mathbf{x} , irrespective of the class, and is given by the sum

$$p(\mathbf{x}) = \sum_k p(\mathbf{x}|\mathcal{C}_k) \cdot P(\mathcal{C}_k) .$$

If, for a given object, the value \mathbf{x} of the feature vector is known, one has to combine this information with the a-priori probabilities.

This is done by Bayes' theorem

$$\text{a-posteriori probability for class } k : \quad P(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k) \cdot P(\mathcal{C}_k)}{p(\mathbf{x})} .$$

The denominator $p(\mathbf{x})$, the unconditional density function, ensures the correct normalization of the a-posteriori probabilities $P(\mathcal{C}_k|\mathbf{x})$:

$$\text{Normalization:} \quad \sum_k P(\mathcal{C}_k|\mathbf{x}) = 1 .$$

Notation

Lower-case letters are used for probability densities, and upper-case letters for probabilities.

k -th class	\mathcal{C}_k
point in feature space	x, \mathbf{x}
probability density	$p(x), p(\mathbf{x})$
class-conditional probability density	$p(\mathbf{x} \mathcal{C}_k)$
a-priori probability for class k	$P(\mathcal{C}_k)$
a-posteriori probability for class k	$P(\mathcal{C}_k \mathbf{x})$
discriminant function for class \mathcal{C}_k	$y_k(\mathbf{x})$
loss matrix	\mathbf{L} with elements L_{kj}
rejection threshold	P_{th}

Minimizing the probability of misclassification

After the observation of a feature vector \mathbf{x} , the a-posteriori probabilities $P(\mathcal{C}_k|\mathbf{x})$ give the probability of the pattern belonging to class \mathcal{C}_k .

The probability of misclassification is minimized by selecting class \mathcal{C}_k if

$$P(\mathcal{C}_k|\mathbf{x}) > P(\mathcal{C}_j|\mathbf{x}) \quad \text{for all } j \neq k$$

For this comparison the unconditional density $p(\mathbf{x})$ may be removed from the Bayes' formula, and the criterion can be rewritten in the form

$$\boxed{P(\mathbf{x}|\mathcal{C}_k) \cdot P(\mathcal{C}_k) > P(\mathbf{x}|\mathcal{C}_j) \cdot P(\mathcal{C}_j) \quad \text{for all } j \neq k}.$$

This is the fundamental decision formula, which minimizes the *probability of misclassification*.

Each point of feature space is assigned to one of c classes. The feature space is divided up into c decision regions $\mathcal{R}_1, \mathcal{R}_2 \dots \mathcal{R}_c$; a feature vector \mathbf{x} falling in region \mathcal{R}_k is assigned to class \mathcal{C}_k .

The decision regions \mathcal{R}_k for a particular class \mathcal{C}_k may be contiguous or may be divided into disjoint regions.

The boundaries between these regions are called decision boundaries (surfaces).

Decision probabilities

The probability for a correct decision is given by

$$\begin{aligned} P(\text{correct}) &= \sum_k P(\mathbf{x} \in \mathcal{R}_k, \mathcal{C}_k) = \sum_k P(\mathbf{x} \in \mathcal{R}_k | \mathcal{C}_k) \cdot P(\mathcal{C}_k) \\ &= \sum_k \int_{\mathcal{R}_k} p(\mathbf{x} \in \mathcal{R}_k | \mathcal{C}_k) \cdot P(\mathcal{C}_k) \, d\mathbf{x} . \end{aligned}$$

This probability is maximized by choosing the regions $\{\mathcal{R}_k\}$ such that each \mathbf{x} is assigned to the class for which the integrand is a maximum.

Error probabilities in the case of only two classes:

$$\begin{aligned} P(\text{error}) &= P(\mathbf{x} \in \mathcal{R}_2, \mathcal{C}_1) + P(\mathbf{x} \in \mathcal{R}_1, \mathcal{C}_2) = P(\mathbf{x} \in \mathcal{R}_2 | \mathcal{C}_1) \cdot P(\mathcal{C}_1) + P(\mathbf{x} \in \mathcal{R}_1 | \mathcal{C}_2) \cdot P(\mathcal{C}_2) \\ &= \int_{\mathcal{R}_2} p(\mathbf{x} | \mathcal{C}_1) \cdot P(\mathcal{C}_1) \, d\mathbf{x} + \int_{\mathcal{R}_1} p(\mathbf{x} | \mathcal{C}_2) \cdot P(\mathcal{C}_2) \, d\mathbf{x} \end{aligned}$$

$P(\mathbf{x} \in \mathcal{R}_2, \mathcal{C}_1)$ is the joint probability of \mathbf{x} being assigned to class \mathcal{C}_2 , and the true class being \mathcal{C}_1 .

If $p(\mathbf{x} | \mathcal{C}_1) \cdot P(\mathcal{C}_1) > p(\mathbf{x} | \mathcal{C}_2) \cdot P(\mathcal{C}_2)$ for a given \mathbf{x} , the regions \mathcal{R}_1 and \mathcal{R}_2 have to be chosen such that \mathbf{x} is in region \mathcal{R}_1 , since this gives the smaller contribution to the error.

Discriminant functions

Decision on class membership has been based solely on the relative sizes of the probabilities. The classification can be reformulated in terms of a set of *discriminant functions* $y_1(\mathbf{x})$, $y_2(\mathbf{x}) \dots$

A feature vector \mathbf{x} is assigned to class \mathcal{C}_k if

$$y_k(\mathbf{x}) > y_j(\mathbf{x}) \quad \text{for all } j \neq k .$$

The decision rule for minimizing the probability of misclassification can be made by choosing

$$y_k(\mathbf{x}) = P(\mathcal{C}_k|\mathbf{x}) .$$

Since the unconditional density $p(\mathbf{x})$ does not affect the classification decision an equivalent choice is

$$y_k(\mathbf{x}) = p(\mathbf{x}|\mathcal{C}_k) \cdot P(\mathcal{C}_k) .$$

Furthermore, since only the relative magnitude of the discriminant function are important in determining the class, replacing $y(\mathbf{x})$ by $g(y(\mathbf{x}))$, where $g(\cdot)$ is any monotonic function, will not modify the classification.

Taking the logarithms:

$$y_k(\mathbf{x}) = \ln p(\mathbf{x}|\mathcal{C}_k) + \ln P(\mathcal{C}_k) .$$

Decision boundaries are given by the regions where the discriminant functions are equal. If the regions \mathcal{R}_k and \mathcal{R}_j are contiguous then the decision boundary is given by $y_k(\mathbf{x}) = y_j(\mathbf{x})$.

Two-class discriminant functions

In this special case the two discriminant functions $y_1(\mathbf{x})$ and $y_2(\mathbf{x})$ can alternatively be combined to a single discriminant function

$$y(\mathbf{x}) = y_1(\mathbf{x}) - y_2(\mathbf{x}) .$$

The pattern \mathbf{x} is assigned to class \mathcal{C}_1 if $y(\mathbf{x}) > 0$ and to class \mathcal{C}_2 if $y(\mathbf{x}) < 0$. Alternative forms of the single discriminant function are:

$$\begin{aligned} y(\mathbf{x}) &= P(\mathcal{C}_1|\mathbf{x}) - P(\mathcal{C}_2|\mathbf{x}) \\ \text{and} \quad y(\mathbf{x}) &= \ln \frac{p(\mathbf{x}|\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)} + \ln \frac{P(\mathcal{C}_1)}{P(\mathcal{C}_2)} . \end{aligned}$$

Risk minimization

In many applications the minimization of the probability of misclassifying is not the most appropriate criterion, especially if there are rare processes, where rejection of the correct class has serious consequences. This may be taken into account in the following way.

A loss matrix \mathbf{L} is defined, where the elements L_{kj} have the meaning of a **penalty** associated with assigning a pattern to class \mathcal{C}_j , when it in fact belongs to class \mathcal{C}_k . Then the expected (average) loss for those patterns is given by

$$R_k = \sum_j L_{kj} \int_{\mathcal{R}_j} p(\mathbf{x}|\mathcal{C}_k) \, d\mathbf{x}.$$

The overall expected loss, or *risk*, for patterns from all classes is

$$R = \sum_k R_k \cdot P(\mathcal{C}_k) = \sum_j \int_{\mathcal{R}_j} \left\{ \sum_k L_{kj} p(\mathbf{x}|\mathcal{C}_k) \cdot P(\mathcal{C}_k) \right\} d\mathbf{x}.$$

The *risk is minimized* if the integrand is minimized at each point \mathbf{x} , i.e. if the regions \mathcal{R}_j are chosen such that $\mathbf{x} \in \mathcal{R}_j$ when

$$\sum_k L_{kj} p(\mathbf{x}|\mathcal{C}_k) \cdot P(\mathcal{C}_k) < \sum_k L_{ki} p(\mathbf{x}|\mathcal{C}_k) \cdot P(\mathcal{C}_k) \quad \text{for all } i \neq j.$$

Loss matrix

This is a generalization of the decision rule for minimization the probability of misclassification. The general formula reduces to the special formula for the special assignment

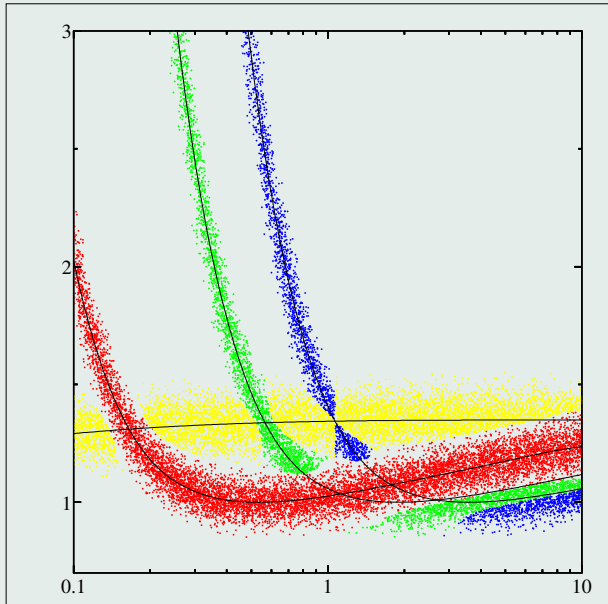
$$L_{kj} = \begin{cases} 0 & \text{for } k = j \quad (\text{correct class}) \\ 1 & \text{for } k \neq j \quad (\text{wrong class}) \end{cases}.$$

In general the elements of the loss matrix have to be estimated based on **experience**. If class \mathcal{C}_j is a **rare**, but important class, which should have a small loss, then for $k \neq j$ the elements L_{kj} should be small; otherwise the acceptance of class \mathcal{C}_j is small because of the small value of the a-priori probability $P(\mathcal{C}_j)$ for rare classes.

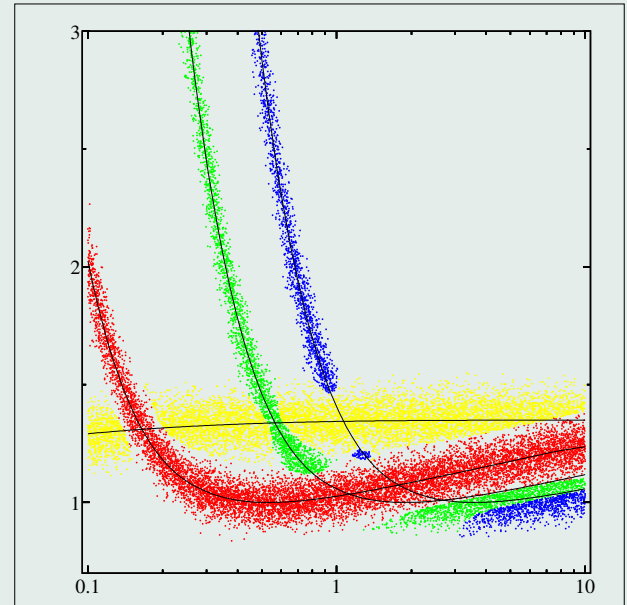
In any case the classification has to be studied carefully in order to check the values of the a-priori probabilities and the loss matrix \mathbf{L} .

Example: particle identification

dE/dx -curves for π (red), K^\pm (green), p (blue) and electron (yellow) with a-priori probabilities 100 : 10 : 5 : 5



Default loss matrix



Loss matrix elements in favour of electron with element 0.2 instead of 1 for p , 0.8 for π and K^\pm .

Rejection thresholds

Most classification errors are expected to occur in those regions of \mathbf{x} -space where the largest of the a-posteriori probabilities is relatively **low**, due to a strong overlap between different classes.

It may be better **not** to make a classification decision in those case, and to reject the pattern instead.

Rejection can be made on the basis of a threshold P_{th} in the range 0 to 1, according to

$$\text{if } \max_k P(\mathcal{C}_k|\mathbf{x}) \quad \begin{cases} \geq P_{\text{th}}, & \text{then classify } \mathbf{x} \\ < P_{\text{th}}, & \text{then reject } \mathbf{x} . \end{cases}$$

Normal distribution and discriminant functions

If the class-conditional density functions are independent normal distributions, then the discriminant function has the form

$$y_k(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \mathbf{V}_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k) - \frac{1}{2} \ln |\mathbf{V}_k| + \ln P(\mathcal{C}_k) ,$$

where constant terms have been dropped. Decision boundaries given by conditions $y_k(\mathbf{x}) = y_j(\mathbf{x})$ are general quadratic functions in the feature space.

The discriminant functions are simplified if the covariance matrices for the various classes are equal ($\mathbf{V}_k \equiv \mathbf{V}$). Then the terms with the determinant $|\mathbf{V}|$ are identical for all classes and can be dropped. The quadratic contribution $\mathbf{x}^T \mathbf{V}^{-1} \mathbf{x}$ is class-independent too and can be dropped. Because $\mathbf{x}^T \mathbf{V}^{-1} \boldsymbol{\mu}_k = \boldsymbol{\mu}_k^T \mathbf{V}^{-1} \mathbf{x}$ the discriminant functions can be simplified to the form

$$\begin{aligned} y_k(\mathbf{x}) &= \mathbf{w}_k^T \mathbf{x} + w_{k0} \quad \text{with} \quad \mathbf{w}_k^T &= \boldsymbol{\mu}_k^T \mathbf{V}^{-1} \\ w_{k0} &= -\frac{1}{2} \boldsymbol{\mu}_k^T \mathbf{V}^{-1} \boldsymbol{\mu}_k + \ln P(\mathcal{C}_k) . \end{aligned}$$

The functions are *linear* in \mathbf{x} ; decision boundaries corresponding to $y_k(\mathbf{x}) = y_j(\mathbf{x})$ are then hyperplanes.

Classification

Strategies for classification	2
Probabilities and classification	3
Class-conditional probability	4
Notation	5
 Minimizing the probability of misclassification	 6
Decision probabilities	7
Discriminant functions	8
Two-class discriminant functions	9
 Risk minimization	 10
Loss matrix	11
Example: particle identification	12
Rejection thresholds	13
Normal distribution and discriminant functions	14