

Pattern Recognition. Bayesian and non-Bayesian Tasks.

Petr Pošík

This lecture is based on the book
Ten Lectures on Statistical and Structural Pattern Recognition
by Michail I. Schlesinger and Václav Hlaváč (Kluwer, 2002).

(V české verzi kniha vyšla pod názvem
Deset přednášek z teorie statistického a strukturálního rozpoznávání
ve vydavatelství ČVUT v roce 1999.)

Pattern Recognition	2
Concepts.....	3
Notes.....	4
PR task examples	5
Two types of PR	6
Bayesian DT	7
Bayesian dec. task.....	8
Characteristics of q^*	9
Two special cases	10
Limitations.....	11
Non-Bayesian DT	13
Non-Bayesian tasks	14
Neyman-Pearson	15
Minimax task.....	16
Wald task	17
Linnik tasks	19
Summary of PR.....	20
Reference	21

Definitions of concepts

An **object** of interest is characterized by the following parameters:

- **observation** $x \in X$ (vector of numbers, graph, picture, sound, ECG, ...), and
- **hidden state** $k \in K$.
- k is often viewed as the object **class**, but it may be something different, e.g. when we seek for the location k of an object based on the picture x taken by a camera.

Joint probability distribution $p_{XK} : X \times K \rightarrow \langle 0, 1 \rangle$

- $p_{XK}(x, k)$ is the joint probability that the object is in the state k and we observe x .
- $p_{XK}(x, k) = p_{X|K}(x|k) \cdot p_K(k)$

Decision strategy (or function or rule) $q : X \rightarrow D$

- D is a set of possible decisions. (Very often $D = K$.)
- q is a function that assigns a decision $d = q(x), d \in D$, to each $x \in X$.

Penalty function (or loss function) $W : K \times D \rightarrow \mathcal{R}$ (real numbers)

- $W(k, d)$ is a penalty for decision d if the object is in state k .

Risk $R : Q \rightarrow \mathcal{R}$

- the mathematical expectation of the penalty which must be paid when using the strategy q .

Notes to decision tasks

In the following, we consider decision tasks where

- the decisions do not influence the state of nature (unlike *game theory* or *control theory*).
- a single decision is made, issues of time are ignored in the model (unlike *control theory*, where decisions are typically taken continuously in real time).
- the costs of obtaining the observations are not modelled (unlike *sequential decision theory*).

The *hidden parameter* k (*state*, *class*) is considered not observable. Common situations are:

- k can be observed, but at a high cost.
- k is a future state (e.g. price of gold) and will be observed later.

Pattern recognition task examples

The description of the concepts is very general—so far we did not specify what the items of the X , K , and D sets actually are, how they are represented.

Application	Observation (measurement)	Decisions
Coin value in a slot machine	$x \in \mathcal{R}^n$	Value
Cancerous tissue detection	Gene-expression profile, $x \in \mathcal{R}^n$	{yes, no}
Medical diagnostics	Results of medical tests, $x \in \mathcal{R}^n$	Diagnosis
Optical character recognition	2D bitmap, intensity image	Words, numbers
License plate recognition	2D bitmap, grey-level image	Characters, numbers
Fingerprint recognition	2D bitmap, grey-level image	Personal identity
Face detection	2D bitmap	{yes, no}
Speech recognition	$x(t)$	Words
Speaker identification	$x(t)$	Personal identity
Speaker verification	$x(t)$	{yes, no}
EEG, ECG analysis	$\mathbf{x}(t)$	Diagnosis
Forfeit detection	Various	{yes, no}

Two types of pattern recognition

1. Statistical pattern recognition

- Objects are represented as points in a vector space.
- The point (vector) \mathbf{x} contains the individual observations (in a numerical form) as its coordinates.

2. Structural pattern recognition

- The object observations contain a structure which is represented and used for recognition.
- A typical example of the representation of a structure is *a grammar*.

Bayesian decision task

Given the sets X, K , and D , and functions $p_{XK} : X \times K \rightarrow \langle 0, 1 \rangle$ and $W : K \times D \rightarrow \mathcal{R}$, find a strategy $q : X \rightarrow D$ which minimizes the **Bayesian risk** of the strategy q

$$R(q) = \sum_{x \in X} \sum_{k \in K} p_{XK}(x, k) \cdot W(k, q(x)).$$

The optimal strategy q , denoted as q^* , is then called **the Bayesian strategy**. The Bayesian risk can be expressed as

$$\begin{aligned} R(q) &= \sum_{x \in X} \sum_{k \in K} p_{K|X}(k|x) \cdot p_X(x) \cdot W(k, q(x)) = \\ &= \sum_{x \in X} p_X(x) \sum_{k \in K} p_{K|X}(k|x) \cdot W(k, q(x)) = \\ &= \sum_{x \in X} p_X(x) \cdot R(q(x), x), \text{ where} \\ R(d, x) &= \sum_{k \in K} p_{K|X}(k|x) \cdot W(k, d) \end{aligned}$$

is the **partial risk**, i.e. the expected penalty for decision d given the observation x . The minimization of the Bayesian risk can be formulated as

$$R(q^*) = \min_{q \in Q} R(q) = \sum_{x \in X} p_X(x) \cdot \min_{d \in D} R(d, x),$$

i.e. the Bayesian strategy can be constructed by choosing the decision d^* that minimizes the partial risk for each observation x .

Bayesian strategy characteristics

Bayesian strategy can be derived for infinite X, K and/or D by replacing summation with integration and probability mass function with probability density function in the formulation of Bayesian decision task.

Bayesian strategy is deterministic.

- q provides the same decision $d = q(x)$ for the same x , despite k may be different.
- What if we used a randomized strategy q of the form $q(d|x)$, i.e. if the decision d would be chosen randomly using the probability distribution $q(d|x)$?
- The risk of the randomized strategy $q(d|x)$ is equal or greater than the risk of the deterministic Bayesian strategy $q(x)$.

Bayesian strategy divides the probability space to $|D|$ convex cones $C(d)$.

- **Probability space?** Any observation x is mapped to a point in a $|K|$ -dimensional linear space (delimited by the positive coordinates) with the coordinates $(p_{X|1}(x|1), p_{X|2}(x|2), \dots, p_{X|k}(x|k))$.
- **Cone?** Let S be a linear space. Any subspace $C \subset S$ is a **cone** if for each $x \in C$ also $\alpha x \in C$ for any real number $\alpha > 0$.
- **Convex cone?** For any 2 points $x_1 \in C$ and $x_2 \in C$, and for any point x lying on the line between x_1 and x_2 , also $x \in C$.
- The individual $C(d)$ are **linearly separable!!!**

Two special cases of the Bayesian decision task

Probability of error when estimating k

- The task is to decide the object state k , i.e. $D = K$.
- The goal is to minimize $Pr(q(x) \neq k)$.
- $Pr(q(x) \neq k) = R(q)$ if

$$W(k, q(x)) = \begin{cases} 0 & \text{if } q(x) = k, \\ 1 & \text{otherwise.} \end{cases}$$

- In this case:

$$\begin{aligned} q(x) &= \arg \min_{d \in D} \sum_{k \in K} p_{XK}(x, k) W(k, d) = \\ &= \arg \max_{d \in D} p_{K|X}(k|x), \end{aligned} \quad (1)$$

i.e. compute *posterior probabilities* of all states k given the observation x , and decide for the most probable state.

- **Maximum posterior (MAP) estimation.**

Bayesian strategy with the dontknow decision

- Using the partial risk $R(d, x) = \sum_{k \in K} p_{K|X}(k|x) \cdot W(k, d)$, for each observation x , we shall provide the decision d minimizing $R(d, x)$.
- But even this optimal $R(d, x)$ may not be sufficiently low, i.e. x does not convey sufficient information for a low-risk decision.
- Let's use $D = K \cup \{\text{dontknow}\}$ and define

$$W(k, d) = \begin{cases} 0 & \text{if } d = k, \\ 1 & \text{if } d \neq k \text{ and } d \neq \text{dontknow} \\ \epsilon & \text{if } d = \text{dontknow.} \end{cases}$$

- In this case:

$$q(x) = \begin{cases} \arg \max_{k \in K} p_{K|X}(k|x) & \text{if } \max_{k \in K} p_{K|X}(k|x) > 1 - \epsilon, \\ \text{dontknow} & \text{if } \max_{k \in K} p_{K|X}(k|x) \leq 1 - \epsilon. \end{cases}$$

Limitations of the Bayesian approach

To use the Bayesian approach we need to know:

1. The penalty function W .
2. The *a priori* probabilities of states $p_K(k)$.
3. The conditional probabilities of observations $p_{X|K}(x|k)$.

Penalty function:

- Important: $W(k, d) \in \mathbb{R}$
- We cannot use the Bayesian formulation for tasks where identifying the penalties with R substantially deforms the task, i.e. *when the penalties cannot be measured in (or easily transformed to) the same units.*
- How do you compare the following penalties:
 - games, fairy tales:
loose your horse vs. loose your sword vs. loose your fiancée
 - system diagnostics, health diagnosis:
false alarm (costs you some money) vs. overlooked danger (may cost you a human life)
 - judicial error:
to convict an innocent (huge harm for 1 innocent person) vs. to free a killer (potential harm to many innocent persons)

Limitations of the Bayesian approach (cont.)

Prior probabilities of states:

- Probabilities $p_K(k)$
 - may be unknown (then we can determine them by further study), or
 - may not exist at all (if the state k is not random).
- E.g. we observe a plane x and we want to decide if it is an enemy aircraft or not.
 - $p_{X|K}(x|k)$ may be quite complex, but known (it at least exists).
 - $p_K(k)$, however, do not exist—the frequency of enemy plane observation does not converge to any number.

Conditional probabilities of observations:

- Again, probabilities $p_{X|K}(x|k)$ may not be known or may not exist.
- E.g. if we want to decide what characters are on paper cards written by several persons, the observation x of the state k is influenced by an unobservable non-random intervention—by the writer z .
 - We can only talk about $p_{X|K,Z}(x|k,z)$, not about $p_{X|K}(x|k)$.
 - If Z was random and if we knew $p_Z(z)$, then we could compute also $p_{X|K}(x|k)$.

Non-Bayesian Decision Theory

13 / 21

Non-Bayesian decision tasks

When?

- Tasks where W , p_K , or $p_{X|K}$ are not known.
- Even if all the events are random and all probabilities are known, it is sometimes helpful to approach the problem as a non-Bayesian task.
- In practical tasks, it can be more intuitive for the customer to express the desired strategy properties as allowed rates of false positives (false alarm) and false negatives (overlooked danger).

There are several special cases of practically useful non-Bayesian formulations for which the solution is known:

- The strategies that solve these non-Bayesian tasks are of the same form as Bayesian strategies—they **divide the probability space to a set of convex cones**.
- These non-Bayesian tasks can be formulated as linear programs and **solved by linear programming methods**.

There are many other non-Bayesian tasks for which the solution is not known yet.

Neyman-Pearson task

Situation:

- Observation $x \in X$, states $k = 1$ (normal), $k = 2$ (dangerous), $K = \{1, 2\}$.
- The probability distribution $p_{X|K}(x|k)$ exists and is known.
- Given the observation x , the task is to decide k , i.e. if the object is in normal or dangerous state.
- The set X is to be divided to 2 subsets X_1 and X_2 , $X = X_1 \cup X_2$.
- In this formulation, $p_K(k)$ and $W(k, d)$ is not needed.

Each strategy q is characterized by 2 numbers:

- Probability of false positive (false alarm):

$$\omega(1) = \sum_{x \in X_2} p_{X|K}(x|1)$$

- Probability of false negative (overlooked danger):

$$\omega(2) = \sum_{x \in X_1} p_{X|K}(x|2)$$

Neyman-Pearson task formulation:

Find a strategy q , i.e. a decomposition of X into X_1 and X_2 , such that

- the probability of overlooked danger (FN) is not larger than a predefined value ϵ , i.e.

$$\omega(2) = \sum_{x \in X_1} p_{X|K}(x|2) \leq \epsilon,$$

- and the probability of false alarm (FP) is minimal, i.e.

$$\text{minimize } \omega(1) = \sum_{x \in X_2} p_{X|K}(x|1),$$

- under the additional conditions

$$X_1 \cap X_2 = \emptyset, X_1 \cup X_2 = X.$$

Solution: The optimal strategy q^* separates X_1 and X_2 according to the *likelihood ratio*:

$$q^*(x) = \begin{cases} 1 & \text{iff } \frac{p_{X|K}(x|1)}{p_{X|K}(x|2)} > \theta, \\ 2 & \text{iff } \frac{p_{X|K}(x|1)}{p_{X|K}(x|2)} < \theta. \end{cases}$$

Minimax task

Situation:

- Observation $x \in X$, states $k \in K$.
- $q : X \rightarrow K$ — given the observation x , the strategy decides the object state k .
- The set X is to be divided to $|K|$ subsets $X_1, \dots, X_{|K|}$, $X = X_1 \cup \dots \cup X_{|K|}$.
- Again, $p_K(k)$ and $W(k, d)$ are not required.

Each strategy is described by $|K|$ numbers

$$\omega(k) = \sum_{x \notin X_k} p_{X|K}(x|k),$$

i.e. by the conditional probabilities of a wrong decision under the condition that the true hidden state is k .

Minimax task formulation:

Find a strategy q^* which minimizes

$$\max_{k \in K} \omega(k)$$

Solution:

- The solution is of the same form as the Bayesian strategies.
- The solution for the $|K| = 2$ case is similar to the Neyman-Pearson task, with the exception that in minimax task the probability of FN cannot be controlled explicitly.

Wald task

Motivation:

- The Neyman-Pearson task is asymmetric: the prob. of FN is controlled explicitly, while the probability of FP is minimized (but can be quite high).
- Can we find a strategy for which *both* the error probabilities would not exceed a predefined ϵ ? No, the demands often cannot be accomplished in the same time.

Wald's relaxation:

- Decompose X into X_1 , X_2 , and X_0 corresponding to $D = K \cup \{\text{dontknow}\}$.
- Strategy of this form is characterized by 4 numbers:

- the conditional prob. of a wrong decision about the state k ,

$$\omega(1) = \sum_{x \in X_2} p_{X|K}(x|1) \quad \text{and} \quad \omega(2) = \sum_{x \in X_1} p_{X|K}(x|2),$$

- the conditional prob. of the dontknow decision when the object state is k ,

$$\chi(1) = \sum_{x \in X_0} p_{X|K}(x|1) \quad \text{and} \quad \chi(2) = \sum_{x \in X_0} p_{X|K}(x|2).$$

- The requirements $\omega(1) \leq \epsilon$ and $\omega(2) \leq \epsilon$ are no longer contradictory for an arbitrarily small $\epsilon > 0$, since the strategy $X_0 = X, X_1 = \emptyset, X_2 = \emptyset$ is plausible.
- Each strategy fulfilling $\omega(1) \leq \epsilon$ and $\omega(2) \leq \epsilon$ is then characterized by how often the strategy refuses to decide, i.e. by the number $\max(\chi(1), \chi(2))$.

Wald task (cont.)

Wald task formulation:

Find a strategy q^* which minimizes

$$\max(\chi(1), \chi(2))$$

subject to conditions $\omega(1) \leq \epsilon$ and $\omega(2) \leq \epsilon$.

Solution: The optimal decision is based on the likelihood ratio and 2 thresholds $\theta_1 > \theta_2$:

$$q^*(x) = \begin{cases} 1 & \text{iff } \frac{p_{X|K}(x|1)}{p_{X|K}(x|2)} > \theta_1, \\ 2 & \text{iff } \frac{p_{X|K}(x|1)}{p_{X|K}(x|2)} < \theta_2, \\ \text{dontknow} & \text{otherwise.} \end{cases}$$

In [SH02], also the generalization for $|K| > 2$ is given.

Linnik tasks

a.k.a. statistical decisions with non-random interventions
a.k.a. evaluations of complex hypotheses.

Previous non-Bayesian tasks did not require

- the a priori probabilities of the states $p_K(k)$, and
- the penalty function $W(k, d)$ to be known.

In Linnik tasks,

- the conditional probabilities $p_{X|K}(x|k)$ do not exist,
- the a priori probabilities $p_K(k)$ may exist (it depends on the fact if the state k is a random variable or not),
- but the conditional probabilities $p_{X|K,Z}(x|k, z)$ do exist, i.e. the random observation x depends not only on the (random or non-random) object state k , but also on a non-random intervention z .

Goal:

- find a strategy that minimizes the probability of incorrect decision in case of the worst intervention z .

See examples in [SH02].

Summary of PR

- The aim of PR is to design decision strategies (classifiers) which—given an observation x of an object with a hidden state k —provide a decision d such that this decision making process is optimal with respect to a certain criterion.
- If the statistical properties of (x, k) are completely known, and if we are able to design a suitable penalty function $W(k, d)$, we should solve the task in the *Bayesian framework* and search for the *Bayesian strategy* which optimizes the *Bayesian risk* of the strategy.
 - The minimization of the probability of an error is a special case, the resulting Bayesian strategy decides for the state with the *maximum a posteriori probability*.
- If the statistical properties are known only partially, or are not known at all, or if a reasonable penalty function cannot be constructed, we face a *non-Bayesian task*.
 - Several practically important special cases of non-Bayesian tasks are well-analyzed and solved (Neyman-Pearson, minimax, Wald, ...).
 - There are plenty of non-Bayesian tasks we can say nothing about.

Reference

[SH02] Michail I. Schlesinger and Václav Hlaváč. *Ten Lectures on Statistical and Structural Pattern Recognition*. Kluwer Academic Publishers, Dodrecht, The Netherlands, 2002.