

Probabilistic classification

Tomáš Svoboda and Matěj Hoffmann
thanks to, Daniel Novák and Filip Železný

Vision for Robots and Autonomous Systems, Center for Machine Perception
Department of Cybernetics
Faculty of Electrical Engineering, Czech Technical University in Prague

April 23, 2019

(Re-)introduction uncertainty/probability

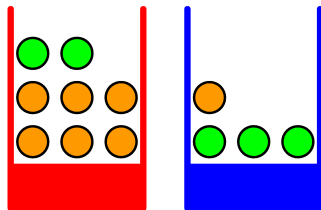
- ▶ Markov Decision Processes - uncertainty about outcome of **actions**
- ▶ Now: uncertainty may be also associated with **states**
 - ▶ Different states may have different **prior probabilities**
 - ▶ The states $s \in S$ may not be directly observable
 - ▶ They need to be inferred from **features** $x \in X$
- ▶ This is addressed by the rules of probability (*such as Bayes theorem*) and leads on to
 - ▶ Bayesian classification
 - ▶ Bayesian decision making

(Re-)introduction uncertainty/probability

- ▶ Markov Decision Processes - uncertainty about outcome of **actions**
- ▶ Now: uncertainty may be also associated with **states**
 - ▶ Different states may have different **prior probabilities**
 - ▶ The states $s \in S$ may not be directly observable
 - ▶ They need to be inferred from **features** $x \in X$
- ▶ This is addressed by the rules of probability (*such as Bayes theorem*) and leads on to
 - ▶ Bayesian classification
 - ▶ Bayesian decision making

Probability example: Picking fruits

- ▶ red box: 2 apples, 6 oranges
- ▶ blue box: 3 apples, 1 orange



- ▶ Scenario: Pick a box (say red box in 40% cases), then pick a fruit at random
- ▶ (Frequent) questions:
 - ▶ What is the overall probability that the selection procedure will pick an apple?
 - ▶ Given that we have chosen an orange, what is the probability that the box we chose was the blue one?

Example from Chapter 1.2 [1]

Rules of probability and notation I

- ▶ random variables X, Y
- ▶ x_i where $i = 1, \dots, M$ – values taken by variable X
- ▶ y_j where $j = 1, \dots, L$ – values taken by variable Y
- ▶ $P(X = x_i, Y = y_i)$ – probability that X takes the value x_i and Y takes y_i – joint probability
- ▶ $P(X = x_i)$ – probability that X takes the value x_i
- ▶ Sum rule of probability :
 - ▶ $P(X = x_i) = \sum_{j=1}^L P(X = x_i, Y = y_j)$
 - ▶ $P(X = x_i)$ is sometimes called marginal probability – obtained by marginalizing / summing out the other variables
 - ▶ general rule, compact notation: $P(X) = \sum_Y P(X, Y)$

Rules of probability and notation II

- ▶ **Conditional probability** : $P(Y = y_j | X = x_i)$
- ▶ **Product rule of probability** :
 - ▶ $P(X = x_i, Y = y_j) = P(Y = y_j | X = x_i)P(X = x_i)$
 - ▶ general rule, compact notation: $P(X, Y) = P(Y|X)P(X)$
- ▶ **Bayes theorem** :
 - ▶ from $P(X, Y) = P(Y, X)$ and product rule

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

$$posterior = \frac{likelihood \times prior}{evidence}$$

Decision example: Insure or not? (from late 1980s) [4]

A doctor calls: “HIV test positive, 999/1000 you die in 10 years, I’m sorry ...”. Insurance company does not want to insure married couple.

- ▶ Was the doctor right?
- ▶ Was the insurance company rational?

What the doctor (and the company) knew:

- ▶ HIV test falsely positive only in 1 case out of 1000.
- ▶ Heterosexual male, have family, no drugs, no risk behavior.

Decision example: Insure or not? (from late 1980s) [4]

A doctor calls: “HIV test positive, 999/1000 you die in 10 years, I’m sorry ...”. Insurance company does not want to insure married couple.

- ▶ Was the doctor right?
- ▶ Was the insurance company rational?

What the doctor (and the company) knew:

- ▶ HIV test falsely positive only in 1 case out of 1000.
- ▶ Heterosexual male, have family, no drugs, no risk behavior.

Decision example: Insure or not? (from late 1980s) [4]

A doctor calls: “HIV test positive, 999/1000 you die in 10 years, I’m sorry ...”. Insurance company does not want to insure married couple.

- ▶ Was the doctor right?
- ▶ Was the insurance company rational?

What the doctor (and the company) knew:

- ▶ HIV test falsely positive only in 1 case out of 1000.
- ▶ Heterosexual male, have family, no drugs, no risk behavior.

Decision example: Insure or not? (from late 1980s) [4]

A doctor calls: “HIV test positive, 999/1000 you die in 10 years, I’m sorry ...”. Insurance company does not want to insure married couple.

- ▶ Was the doctor right?
- ▶ Was the insurance company rational?

What the doctor (and the company) knew:

- ▶ HIV test falsely positive only in 1 case out of 1000.
- ▶ Heterosexual male, have family, no drugs, no risk behavior.

Decision example: Insure or not? (from late 1980s) [4]

A doctor calls: “HIV test positive, 999/1000 you die in 10 years, I’m sorry ...”. Insurance company does not want to insure married couple.

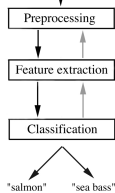
- ▶ Was the doctor right?
- ▶ Was the insurance company rational?

What the doctor (and the company) knew:

- ▶ HIV test falsely positive only in 1 case out of 1000.
- ▶ Heterosexual male, have family, no drugs, no risk behavior.

Decision example: guilty or not

Classification example: What's the fish?



- ▶ Factory for fish processing
- ▶ 2 classes $s_{1,2}$:
 - ▶ salmon
 - ▶ sea bass
- ▶ Features \vec{x} : length, width, lightness etc. from a camera

Fish – classification using probability

$$\textit{posterior} = \frac{\textit{likelihood} \times \textit{prior}}{\textit{evidence}}$$

- ▶ Notation for classification problem
 - ▶ Classes $s_j \in S$ (e.g., salmon, sea bass)
 - ▶ Features $x_i \in X$ or feature vectors (\vec{x}_i) (also called attributes)
- ▶ Optimal classification of \vec{x} :

$$\delta^*(\vec{x}) = \arg \max_j P(s_j | \vec{x})$$

- ▶ We thus choose the most probable class for a given feature vector.
- ▶ Both likelihood and prior are taken into account – recall Bayes rule:

$$P(s_j | \vec{x}) = \frac{P(\vec{x} | s_j) P(s_j)}{P(\vec{x})}$$

Fish – classification using probability

$$\textit{posterior} = \frac{\textit{likelihood} \times \textit{prior}}{\textit{evidence}}$$

- ▶ Notation for classification problem
 - ▶ Classes $s_j \in S$ (e.g., salmon, sea bass)
 - ▶ Features $x_i \in X$ or feature vectors (\vec{x}_i) (also called attributes)
- ▶ Optimal classification of \vec{x} :

$$\delta^*(\vec{x}) = \arg \max_j P(s_j | \vec{x})$$

- ▶ We thus choose the most probable class for a given feature vector.
- ▶ Both likelihood and prior are taken into account – recall Bayes rule:

$$P(s_j | \vec{x}) = \frac{P(\vec{x} | s_j) P(s_j)}{P(\vec{x})}$$

Bayes classification in practice

- ▶ Usually we are not given $P(s|\vec{x})$
 - ▶ It has to be estimated from already classified examples – training data
 - ▶ For discrete \vec{x} , training examples $(\vec{x}_1, s_1), (\vec{x}_2, s_2), \dots, (\vec{x}_I, s_I)$
 - ▶ so-called i.i.d (independent, identically distributed) multiset
 - ▶ every (\vec{x}_i, s_i) is drawn independently from $P(\vec{x}, s)$
 - ▶ Without knowing anything about the distribution, a non-parametric estimate:

$$P(s|\vec{x}) \approx \frac{\# \text{ examples where } \vec{x}_i = \vec{x} \text{ and } s_i = s}{\# \text{ examples where } \vec{x}_i = \vec{x}}$$

- ▶ Hard in practice:
 - ▶ To reliably estimate $P(s|\vec{x})$, the number of examples grows exponentially with the number of elements of \vec{x} .
 - ▶ e.g. with the number of pixels in images
 - ▶ curse of dimensionality
 - ▶ denominator often 0

Bayes classification in practice

- ▶ Usually we are not given $P(s|\vec{x})$
- ▶ It has to be estimated from already classified examples – training data
- ▶ For discrete \vec{x} , training examples $(\vec{x}_1, s_1), (\vec{x}_2, s_2), \dots, (\vec{x}_I, s_I)$
 - ▶ so-called i.i.d (independent, identically distributed) multiset
 - ▶ every (\vec{x}_i, s) is drawn independently from $P(\vec{x}, s)$
- ▶ Without knowing anything about the distribution, a non-parametric estimate:

$$P(s|\vec{x}) \approx \frac{\# \text{ examples where } \vec{x}_i = \vec{x} \text{ and } s_i = s}{\# \text{ examples where } \vec{x}_i = \vec{x}}$$

- ▶ Hard in practice:
 - ▶ To reliably estimate $P(s|\vec{x})$, the number of examples grows exponentially with the number of elements of \vec{x} .
 - ▶ e.g. with the number of pixels in images
 - ▶ curse of dimensionality
 - ▶ denominator often 0

Bayes classification in practice

- ▶ Usually we are not given $P(s|\vec{x})$
- ▶ It has to be estimated from already classified examples – training data
- ▶ For discrete \vec{x} , training examples $(\vec{x}_1, s_1), (\vec{x}_2, s_2), \dots, (\vec{x}_I, s_I)$
 - ▶ so-called i.i.d (independent, identically distributed) multiset
 - ▶ every (\vec{x}_i, s) is drawn independently from $P(\vec{x}, s)$
- ▶ Without knowing anything about the distribution, a non-parametric estimate:

$$P(s|\vec{x}) \approx \frac{\# \text{ examples where } \vec{x}_i = \vec{x} \text{ and } s_i = s}{\# \text{ examples where } \vec{x}_i = \vec{x}}$$

- ▶ Hard in practice:
 - ▶ To reliably estimate $P(s|\vec{x})$, the number of examples grows exponentially with the number of elements of \vec{x} .
 - ▶ e.g. with the number of pixels in images
 - ▶ curse of dimensionality
 - ▶ denominator often 0

Bayes classification in practice

- ▶ Usually we are not given $P(s|\vec{x})$
- ▶ It has to be estimated from already classified examples – training data
- ▶ For discrete \vec{x} , training examples $(\vec{x}_1, s_1), (\vec{x}_2, s_2), \dots, (\vec{x}_I, s_I)$
 - ▶ so-called i.i.d (independent, identically distributed) multiset
 - ▶ every (\vec{x}_i, s) is drawn independently from $P(\vec{x}, s)$
- ▶ Without knowing anything about the distribution, a non-parametric estimate:

$$P(s|\vec{x}) \approx \frac{\# \text{ examples where } \vec{x}_i = \vec{x} \text{ and } s_i = s}{\# \text{ examples where } \vec{x}_i = \vec{x}}$$

- ▶ Hard in practice:
 - ▶ To reliably estimate $P(s|\vec{x})$, the number of examples grows exponentially with the number of elements of \vec{x} .
 - ▶ e.g. with the number of pixels in images
 - ▶ curse of dimensionality
 - ▶ denominator often 0

Naïve Bayes classification

- ▶ For efficient classification we must thus rely on additional assumptions.
- ▶ In the exceptional case of **statistical independence** between \vec{x} components for each class s it holds

$$P(\vec{x}|s) = P(x[1]|s) \cdot P(x[2]|s) \cdot \dots$$

- ▶ Use simple Bayes law and maximize:

$$P(s|\vec{x}) = \frac{P(\vec{x}|s)P(s)}{P(\vec{x})} = \frac{P(s)}{P(\vec{x})} P(x[1]|s) \cdot P(x[2]|s) \cdot \dots =$$

- ▶ No combinatorial curse in estimating $P(s)$ and $P(x[i]|s)$ separately for each i and s .
- ▶ No need to estimate $P(\vec{x})$. (Why?)
- ▶ $P(s)$ may be provided apriori.
- ▶ **naïve** = when used despite statistical dependence

Decision making under uncertainty

- ▶ An important feature of intelligent systems
 - ▶ make the best possible decision
 - ▶ in uncertain conditions.
- ▶ **Example:** Take a tram OR subway from *A* to *B*?
 - ▶ Tram: timetables imply a quicker route, but adherence uncertain.
 - ▶ Subway: longer route, but adherence almost certain.
- ▶ **Example:** where to route a letter with this ZIP?
 - ▶ 15700? 15706? 15200? 15206?
- ▶ What is the optimal decision ?
- ▶ Both examples fall into the same framework.

Decision making under uncertainty

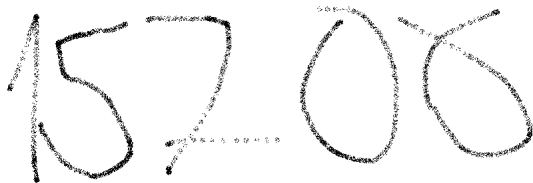
- ▶ An important feature of intelligent systems
 - ▶ make the best possible decision
 - ▶ in uncertain conditions.
- ▶ **Example:** Take a tram OR subway from A to B ?
 - ▶ Tram: timetables imply a quicker route, but adherence uncertain.
 - ▶ Subway: longer route, but adherence almost certain.
- ▶ Example: where to route a letter with this ZIP?

▶ 15700? 15706? 15200? 15206?

- ▶ What is the optimal decision ?
- ▶ Both examples fall into the same framework.

Decision making under uncertainty

- ▶ An important feature of intelligent systems
 - ▶ make the best possible decision
 - ▶ in uncertain conditions.
- ▶ **Example:** Take a tram OR subway from *A* to *B*?
 - ▶ Tram: timetables imply a quicker route, but adherence uncertain.
 - ▶ Subway: longer route, but adherence almost certain.
- ▶ **Example:** where to route a letter with this ZIP?



- ▶ 15700? 15706? 15200? 15206?
- ▶ What is the optimal decision ?
- ▶ Both examples fall into the same framework.

Decision making under uncertainty

- ▶ An important feature of intelligent systems
 - ▶ make the best possible decision
 - ▶ in uncertain conditions.
- ▶ **Example:** Take a tram OR subway from *A* to *B*?
 - ▶ Tram: timetables imply a quicker route, but adherence uncertain.
 - ▶ Subway: longer route, but adherence almost certain.
- ▶ **Example:** where to route a letter with this ZIP?



- ▶ 15700? 15706? 15200? 15206?
- ▶ What is the **optimal decision** ?
- ▶ Both examples fall into the same framework.

Example: What to cook for a dinner [3]

- ▶ Wife coming back from work. Husband: what to cook for dinner?
- ▶ 3 dishes (decisions) in his repertoire:
 - ▶ *nothing* ... don't bother cooking \Rightarrow no work but makes wife upset
 - ▶ *pizza* ... microwave a frozen pizza \Rightarrow not much work but won't impress
 - ▶ *g.T.c.* ... general Tso's chicken \Rightarrow will make her day, but very laborious.
- ▶ Hassle incurred by the individual options depends wife's feeling
- ▶ For each of the 9 possible situation (3 possible decisions \times 3 possible states) the hassle is quantified by a loss function $l(d, s)$:

$l(s, d)$	$d = \textit{nothing}$	$d = \textit{pizza}$	$d = \textit{g.T.c.}$
$s = \textit{good}$	0	2	4
$s = \textit{average}$	5	3	5
$s = \textit{bad}$	10	9	6

Wife's state of mind is an uncertain state.

Example: What to cook for a dinner [3]

- ▶ Wife coming back from work. Husband: what to cook for dinner?
- ▶ 3 dishes (**decisions**) in his repertoire:
 - ▶ *nothing* ... **don't bother cooking** \Rightarrow no work but makes wife upset
 - ▶ *pizza* ... **microwave a frozen pizza** \Rightarrow not much work but won't impress
 - ▶ *g.T.c.* ... **general Tso's chicken** \Rightarrow will make her day, but very laborious.
- ▶ Hassle incurred by the individual options depends wife's feeling
- ▶ For each of the 9 possible situation (3 possible decisions \times 3 possible states) the hassle is quantified by a loss function $l(d, s)$:

$l(s, d)$	$d = \textit{nothing}$	$d = \textit{pizza}$	$d = \textit{g.T.c.}$
$s = \textit{good}$	0	2	4
$s = \textit{average}$	5	3	5
$s = \textit{bad}$	10	9	6

Wife's state of mind is an uncertain state.

Example: What to cook for a dinner [3]

- ▶ Wife coming back from work. Husband: what to cook for dinner?
- ▶ 3 dishes (**decisions**) in his repertoire:
 - ▶ *nothing* ... **don't bother cooking** \Rightarrow no work but makes wife upset
 - ▶ *pizza* ... **microwave a frozen pizza** \Rightarrow not much work but won't impress
 - ▶ *g.T.c.* ... **general Tso's chicken** \Rightarrow will make her day, but very laborious.
- ▶ Hassle incurred by the individual options depends wife's feeling
- ▶ For each of the 9 possible situation (3 possible decisions \times 3 possible states) the hassle is quantified by a **loss function** $l(d, s)$:

$l(s, d)$	$d = \textit{nothing}$	$d = \textit{pizza}$	$d = \textit{g.T.c.}$
$s = \textit{good}$	0	2	4
$s = \textit{average}$	5	3	5
$s = \textit{bad}$	10	9	6

Wife's state of mind is an uncertain state.

Example: What to cook for a dinner [3]

- ▶ Wife coming back from work. Husband: what to cook for dinner?
- ▶ 3 dishes (**decisions**) in his repertoire:
 - ▶ *nothing* ... **don't bother cooking** \Rightarrow no work but makes wife upset
 - ▶ *pizza* ... **microwave a frozen pizza** \Rightarrow not much work but won't impress
 - ▶ *g.T.c.* ... **general Tso's chicken** \Rightarrow will make her day, but very laborious.
- ▶ Hassle incurred by the individual options depends wife's feeling
- ▶ For each of the 9 possible situation (3 possible decisions \times 3 possible states) the hassle is quantified by a **loss function** $l(d, s)$:

$l(s, d)$	$d = \textit{nothing}$	$d = \textit{pizza}$	$d = \textit{g.T.c.}$
$s = \textit{good}$	0	2	4
$s = \textit{average}$	5	3	5
$s = \textit{bad}$	10	9	6

Wife's state of mind is an **uncertain state**.

Example (cont'd), State uncertain, ...

- ▶ Husband's experiment. He tells her he accidentally overtaped their wedding video and observes her reaction
- ▶ Anticipates 4 possible reactions:
 - ▶ *mild* ... all right, we keep our memories.
 - ▶ *irritated* ... how many times do I have to tell you...
 - ▶ *upset* ... Why did I marry this guy?
 - ▶ *alarming* ... silence
- ▶ The reaction is a measurable attribute ("feature") of the mind state.
- ▶ From experience, the husband knows how individual reactions are probable in each state of mind; this is captured by the joint distribution $P(x, s)$.

$P(x, s)$	$x = mild$	$x = irritated$	$x = upset$	$x = alarming$
$s = good$	0.35	0.28	0.07	0.00
$s = average$	0.04	0.10	0.04	0.02
$s = bad$	0.00	0.02	0.05	0.03

Example (cont'd), State uncertain, ...

- ▶ Husband's experiment. He tells her he accidentally overtaped their wedding video and observes her reaction
- ▶ Anticipates 4 possible reactions:
 - ▶ *mild* ... all right, we keep our memories.
 - ▶ *irritated* ... how many times do I have to tell you...
 - ▶ *upset* ... Why did I marry this guy?
 - ▶ *alarming* ... silence
- ▶ The reaction is a measurable attribute ("feature") of the mind state.
- ▶ From experience, the husband knows how individual reactions are probable in each state of mind; this is captured by the joint distribution $P(x, s)$.

$P(x, s)$	$x = mild$	$x = irritated$	$x = upset$	$x = alarming$
$s = good$	0.35	0.28	0.07	0.00
$s = average$	0.04	0.10	0.04	0.02
$s = bad$	0.00	0.02	0.05	0.03

Example (cont'd), State uncertain, ...

- ▶ Husband's experiment. He tells her he accidentally overtaped their wedding video and observes her reaction
- ▶ Anticipates 4 possible reactions:
 - ▶ *mild* ... all right, we keep our memories.
 - ▶ *irritated* ... how many times do I have to tell you....
 - ▶ *upset* ... Why did I marry this guy?
 - ▶ *alarming* ... silence
- ▶ The reaction is a measurable **attribute** (**"feature"**) of the mind state.
- ▶ From experience, the husband knows how individual reactions are probable in each state of mind; this is captured by the joint distribution $P(x, s)$.

$P(x, s)$	$x = mild$	$x = irritated$	$x = upset$	$x = alarming$
$s = good$	0.35	0.28	0.07	0.00
$s = average$	0.04	0.10	0.04	0.02
$s = bad$	0.00	0.02	0.05	0.03

Example (cont'd), State uncertain, ...

- ▶ Husband's experiment. He tells her he accidentally overtaped their wedding video and observes her reaction
- ▶ Anticipates 4 possible reactions:
 - ▶ *mild* ... all right, we keep our memories.
 - ▶ *irritated* ... how many times do I have to tell you....
 - ▶ *upset* ... Why did I marry this guy?
 - ▶ *alarming* ... silence
- ▶ The reaction is a measurable **attribute** (**"feature"**) of the mind state.
- ▶ From experience, the husband knows how individual reactions are probable in each state of mind; this is captured by the **joint distribution $P(x, s)$** .

$P(x, s)$	$x = mild$	$x = irritated$	$x = upset$	$x = alarming$
$s = good$	0.35	0.28	0.07	0.00
$s = average$	0.04	0.10	0.04	0.02
$s = bad$	0.00	0.02	0.05	0.03

Decision strategy

- ▶ **Decision strategy** : a rule selecting a decision for any given value of the measured attribute(s).
- ▶ i.e. function $d = \delta(x)$.
- ▶ Example of husband's possible strategies:

$\delta(x)$	$x = \text{mild}$	$x = \text{irritated}$	$x = \text{upset}$	$x = \text{alarming}$
$\delta_1(x) =$	<i>nothing</i>	<i>nothing</i>	<i>pizza</i>	<i>g.T.c.</i>
$\delta_2(x) =$	<i>nothing</i>	<i>pizza</i>	<i>g.T.c.</i>	<i>g.T.c.</i>
$\delta_3(x) =$	<i>g.T.c.</i>	<i>g.T.c.</i>	<i>g.T.c.</i>	<i>g.T.c.</i>
$\delta_4(x) =$	<i>nothing</i>	<i>nothing</i>	<i>nothing</i>	<i>nothing</i>

- ▶ How many strategies?
- ▶ How to define which strategy is best? How to sort them by quality?
- ▶ Define the risk of a strategy as a mean (expected) loss value .

$$r(\delta) = \sum_x \sum_s l(s, \delta(x)) P(x, s)$$

Decision strategy

- ▶ **Decision strategy** : a rule selecting a decision for any given value of the measured attribute(s).
- ▶ i.e. function $d = \delta(x)$.
- ▶ Example of husband's possible strategies:

$\delta(x)$	$x = mild$	$x = irritated$	$x = upset$	$x = alarming$
$\delta_1(x) =$	<i>nothing</i>	<i>nothing</i>	<i>pizza</i>	<i>g.T.c.</i>
$\delta_2(x) =$	<i>nothing</i>	<i>pizza</i>	<i>g.T.c.</i>	<i>g.T.c.</i>
$\delta_3(x) =$	<i>g.T.c.</i>	<i>g.T.c.</i>	<i>g.T.c.</i>	<i>g.T.c.</i>
$\delta_4(x) =$	<i>nothing</i>	<i>nothing</i>	<i>nothing</i>	<i>nothing</i>

- ▶ How many strategies?
- ▶ How to define which strategy is best? How to sort them by quality?
- ▶ Define the risk of a strategy as a mean (expected) loss value .

$$r(\delta) = \sum_x \sum_s l(s, \delta(x)) P(x, s)$$

Decision strategy

- ▶ **Decision strategy** : a rule selecting a decision for any given value of the measured attribute(s).
- ▶ i.e. function $d = \delta(x)$.
- ▶ Example of husband's possible strategies:

$\delta(x)$	$x = mild$	$x = irritated$	$x = upset$	$x = alarming$
$\delta_1(x) =$	<i>nothing</i>	<i>nothing</i>	<i>pizza</i>	<i>g.T.c.</i>
$\delta_2(x) =$	<i>nothing</i>	<i>pizza</i>	<i>g.T.c.</i>	<i>g.T.c.</i>
$\delta_3(x) =$	<i>g.T.c.</i>	<i>g.T.c.</i>	<i>g.T.c.</i>	<i>g.T.c.</i>
$\delta_4(x) =$	<i>nothing</i>	<i>nothing</i>	<i>nothing</i>	<i>nothing</i>

- ▶ How many strategies?
- ▶ How to define which strategy is best? How to sort them by quality?
- ▶ Define the **risk of a strategy** as a **mean (expected) loss value** .

$$r(\delta) = \sum_x \sum_s l(s, \delta(x)) P(x, s)$$

Calculating $r(\delta) = \sum_x \sum_s l(s, \delta(x))P(x, s)$

$l(s, d)$	$d = \textit{nothing}$	$d = \textit{pizza}$	$d = \textit{g.T.c.}$
$s = \textit{good}$	0	2	4
$s = \textit{average}$	5	3	5
$s = \textit{bad}$	10	9	6

$P(x, s)$	$x = \textit{mild}$	$x = \textit{irritated}$	$x = \textit{upset}$	$x = \textit{alarming}$
$s = \textit{good}$	0.35	0.28	0.07	0.00
$s = \textit{average}$	0.04	0.10	0.04	0.02
$s = \textit{bad}$	0.00	0.02	0.05	0.03

$\delta(x)$	$x = \textit{mild}$	$x = \textit{irritated}$	$x = \textit{upset}$	$x = \textit{alarming}$
$\delta_1(x) =$	<i>nothing</i>	<i>nothing</i>	<i>pizza</i>	<i>g.T.c.</i>
$\delta_2(x) =$	<i>nothing</i>	<i>pizza</i>	<i>g.T.c.</i>	<i>g.T.c.</i>
$\delta_3(x) =$	<i>g.T.c.</i>	<i>g.T.c.</i>	<i>g.T.c.</i>	<i>g.T.c.</i>
\vdots	\vdots	\vdots	\vdots	\vdots

Do we need to evaluate all possible strategies? $P(x, s) = P(s|x)P(x)$

Calculating $r(\delta) = \sum_x \sum_s l(s, \delta(x))P(x, s)$

$l(s, d)$	$d = \textit{nothing}$	$d = \textit{pizza}$	$d = \textit{g.T.c.}$
$s = \textit{good}$	0	2	4
$s = \textit{average}$	5	3	5
$s = \textit{bad}$	10	9	6

$P(x, s)$	$x = \textit{mild}$	$x = \textit{irritated}$	$x = \textit{upset}$	$x = \textit{alarming}$
$s = \textit{good}$	0.35	0.28	0.07	0.00
$s = \textit{average}$	0.04	0.10	0.04	0.02
$s = \textit{bad}$	0.00	0.02	0.05	0.03

$\delta(x)$	$x = \textit{mild}$	$x = \textit{irritated}$	$x = \textit{upset}$	$x = \textit{alarming}$
$\delta_1(x) =$	<i>nothing</i>	<i>nothing</i>	<i>pizza</i>	<i>g.T.c.</i>
$\delta_2(x) =$	<i>nothing</i>	<i>pizza</i>	<i>g.T.c.</i>	<i>g.T.c.</i>
$\delta_3(x) =$	<i>g.T.c.</i>	<i>g.T.c.</i>	<i>g.T.c.</i>	<i>g.T.c.</i>
\vdots	\vdots	\vdots	\vdots	\vdots

Do we need to evaluate all possible strategies? $P(x, s) = P(s|x)P(x)$

Calculating $r(\delta) = \sum_x \sum_s l(s, \delta(x))P(x, s)$

$l(s, d)$	$d = \textit{nothing}$	$d = \textit{pizza}$	$d = \textit{g.T.c.}$
$s = \textit{good}$	0	2	4
$s = \textit{average}$	5	3	5
$s = \textit{bad}$	10	9	6

$P(x, s)$	$x = \textit{mild}$	$x = \textit{irritated}$	$x = \textit{upset}$	$x = \textit{alarming}$
$s = \textit{good}$	0.35	0.28	0.07	0.00
$s = \textit{average}$	0.04	0.10	0.04	0.02
$s = \textit{bad}$	0.00	0.02	0.05	0.03

$\delta(x)$	$x = \textit{mild}$	$x = \textit{irritated}$	$x = \textit{upset}$	$x = \textit{alarming}$
$\delta_1(x) =$	<i>nothing</i>	<i>nothing</i>	<i>pizza</i>	<i>g.T.c.</i>
$\delta_2(x) =$	<i>nothing</i>	<i>pizza</i>	<i>g.T.c.</i>	<i>g.T.c.</i>
$\delta_3(x) =$	<i>g.T.c.</i>	<i>g.T.c.</i>	<i>g.T.c.</i>	<i>g.T.c.</i>
\vdots	\vdots	\vdots	\vdots	\vdots

Do we need to evaluate all possible strategies? $P(x, s) = P(s|x)P(x)$

Calculating $r(\delta) = \sum_x \sum_s l(s, \delta(x))P(x, s)$

$l(s, d)$	$d = \textit{nothing}$	$d = \textit{pizza}$	$d = \textit{g.T.c.}$
$s = \textit{good}$	0	2	4
$s = \textit{average}$	5	3	5
$s = \textit{bad}$	10	9	6

$P(x, s)$	$x = \textit{mild}$	$x = \textit{irritated}$	$x = \textit{upset}$	$x = \textit{alarming}$
$s = \textit{good}$	0.35	0.28	0.07	0.00
$s = \textit{average}$	0.04	0.10	0.04	0.02
$s = \textit{bad}$	0.00	0.02	0.05	0.03

$\delta(x)$	$x = \textit{mild}$	$x = \textit{irritated}$	$x = \textit{upset}$	$x = \textit{alarming}$
$\delta_1(x) =$	<i>nothing</i>	<i>nothing</i>	<i>pizza</i>	<i>g.T.c.</i>
$\delta_2(x) =$	<i>nothing</i>	<i>pizza</i>	<i>g.T.c.</i>	<i>g.T.c.</i>
$\delta_3(x) =$	<i>g.T.c.</i>	<i>g.T.c.</i>	<i>g.T.c.</i>	<i>g.T.c.</i>
\vdots	\vdots	\vdots	\vdots	\vdots

Do we need to evaluate all possible strategies?

$$P(x, s) = P(s|x)P(x)$$

Calculating $r(\delta) = \sum_x \sum_s l(s, \delta(x))P(x, s)$

$l(s, d)$	$d = \textit{nothing}$	$d = \textit{pizza}$	$d = \textit{g.T.c.}$
$s = \textit{good}$	0	2	4
$s = \textit{average}$	5	3	5
$s = \textit{bad}$	10	9	6

$P(x, s)$	$x = \textit{mild}$	$x = \textit{irritated}$	$x = \textit{upset}$	$x = \textit{alarming}$
$s = \textit{good}$	0.35	0.28	0.07	0.00
$s = \textit{average}$	0.04	0.10	0.04	0.02
$s = \textit{bad}$	0.00	0.02	0.05	0.03

$\delta(x)$	$x = \textit{mild}$	$x = \textit{irritated}$	$x = \textit{upset}$	$x = \textit{alarming}$
$\delta_1(x) =$	<i>nothing</i>	<i>nothing</i>	<i>pizza</i>	<i>g.T.c.</i>
$\delta_2(x) =$	<i>nothing</i>	<i>pizza</i>	<i>g.T.c.</i>	<i>g.T.c.</i>
$\delta_3(x) =$	<i>g.T.c.</i>	<i>g.T.c.</i>	<i>g.T.c.</i>	<i>g.T.c.</i>
\vdots	\vdots	\vdots	\vdots	\vdots

Do we need to evaluate all possible strategies? $P(x, s) = P(s|x)P(x)$

Bayes optimal strategy

- ▶ The **Bayes optimal strategy** : one minimizing mean risk.

$$\delta^* = \arg \min_{\delta} r(\delta)$$

- ▶ From $P(x, s) = P(s|x)P(x)$ (Bayes rule), we have

$$\begin{aligned} r(\delta) &= \sum_x \sum_s l(s, \delta(x)) P(x, s) = \sum_s \sum_x l(s, \delta(x)) P(s|x) P(x) \\ &= \sum_x P(x) \underbrace{\sum_s l(s, \delta(x)) P(s|x)}_{\text{Conditional risk}} \end{aligned}$$

- ▶ The optimal strategy is obtained by minimizing the conditional risk *separately* for each x :

$$\delta^*(x) = \arg \min_d \sum_s l(s, d) P(s|x)$$

Optimal strategy: $\delta^*(x) = \arg \min_d \sum_s l(s, d)P(s|x)$

$l(s, d)$	$d = \text{nothing}$	$d = \text{pizza}$	$d = \text{g.T.c.}$
$s = \text{good}$	0	2	4
$s = \text{average}$	5	3	5
$s = \text{bad}$	10	9	6

$P(x, s)$	$x = \text{mild}$	$x = \text{irritated}$	$x = \text{upset}$	$x = \text{alarming}$
$s = \text{good}$	0.35	0.28	0.07	0.00
$s = \text{average}$	0.04	0.10	0.04	0.02
$s = \text{bad}$	0.00	0.02	0.05	0.03

$\delta(x)$	$x = \text{mild}$	$x = \text{irritated}$	$x = \text{upset}$	$x = \text{alarming}$
$\delta^*(x) =$??	??	??	??

Statistical decision making: wrapping up

▶ Given:

- ▶ A set of possible **states** : \mathcal{S}
- ▶ A set of possible **decisions** : \mathcal{D}
- ▶ A **loss function** $l : \mathcal{D} \times \mathcal{S} \rightarrow \mathfrak{R}$
- ▶ The range \mathcal{X} of the **attribute**
- ▶ Distribution $P(x, s)$, $x \in \mathcal{X}, s \in \mathcal{S}$.

▶ Define:

- ▶ **Strategy** : function $\delta : \mathcal{X} \rightarrow \mathcal{D}$
- ▶ **Risk of strategy** $\delta : r(\delta) = \sum_x \sum_s l(s, \delta(x))P(x, s)$

▶ Bayes problem:

- ▶ Goal: find the optimal strategy $\delta^* = \arg \min_{\delta \in \Delta} r(\delta)$
- ▶ Solution: $\delta^*(x) = \arg \min_d \sum_s l(s, d)P(s|x)$

A special case - Bayesian *classification*

- ▶ Bayesian classification is a special case of statistical decision theory:
 - ▶ Attribute vector $\vec{x} = (x_1, x_2, \dots)$: pixels 1, 2, ...
 - ▶ **State set \mathcal{S} = decision set $\mathcal{D} = \{0, 1, \dots, 9\}$.**
 - ▶ **State = actual class, Decision = recognized class**
 - ▶ Loss function:

$$l(s, d) = \begin{cases} 0, & d = s \\ 1, & d \neq s \end{cases}$$

$$\delta^*(\vec{x}) = \arg \min_d \sum_s \underbrace{l(s, d)}_{0 \text{ if } d=s} P(s|\vec{x}) = \arg \min_d \sum_{s \neq d} P(s|\vec{x})$$

Obviously $\sum_s P(s|\vec{x}) = 1$, then:

$$P(d|\vec{x}) + \sum_{s \neq d} P(s|\vec{x}) = 1$$

Inserting into above:

$$\delta^*(\vec{x}) = \arg \min_d [1 - P(d|\vec{x})] = \arg \max_d P(d|\vec{x})$$

A special case - Bayesian *classification*

- ▶ Bayesian classification is a special case of statistical decision theory:
 - ▶ Attribute vector $\vec{x} = (x_1, x_2, \dots)$: pixels 1, 2, ...
 - ▶ **State set \mathcal{S} = decision set $\mathcal{D} = \{0, 1, \dots, 9\}$.**
 - ▶ **State = actual class, Decision = recognized class**
 - ▶ Loss function:

$$l(s, d) = \begin{cases} 0, & d = s \\ 1, & d \neq s \end{cases}$$

$$\delta^*(\vec{x}) = \arg \min_d \sum_s \underbrace{l(s, d)}_{0 \text{ if } d=s} P(s|\vec{x}) = \arg \min_d \sum_{s \neq d} P(s|\vec{x})$$

Obviously $\sum_s P(s|\vec{x}) = 1$, then:

$$P(d|\vec{x}) + \sum_{s \neq d} P(s|\vec{x}) = 1$$

Inserting into above:

$$\delta^*(\vec{x}) = \arg \min_d [1 - P(d|\vec{x})] = \arg \max_d P(d|\vec{x})$$

A special case - Bayesian *classification*

- ▶ Bayesian classification is a special case of statistical decision theory:
 - ▶ Attribute vector $\vec{x} = (x_1, x_2, \dots)$: pixels 1, 2, \dots
 - ▶ **State set \mathcal{S} = decision set $\mathcal{D} = \{0, 1, \dots, 9\}$.**
 - ▶ **State = actual class, Decision = recognized class**
 - ▶ Loss function:

$$l(s, d) = \begin{cases} 0, & d = s \\ 1, & d \neq s \end{cases}$$

$$\delta^*(\vec{x}) = \arg \min_d \sum_s \underbrace{l(s, d)}_{0 \text{ if } d=s} P(s|\vec{x}) = \arg \min_d \sum_{s \neq d} P(s|\vec{x})$$

Obviously $\sum_s P(s|\vec{x}) = 1$, then:

$$P(d|\vec{x}) + \sum_{s \neq d} P(s|\vec{x}) = 1$$

Inserting into above:

$$\delta^*(\vec{x}) = \arg \min_d [1 - P(d|\vec{x})] = \arg \max_d P(d|\vec{x})$$

A special case - Bayesian *classification*

- ▶ Bayesian classification is a special case of statistical decision theory:
 - ▶ Attribute vector $\vec{x} = (x_1, x_2, \dots)$: pixels 1, 2, ...
 - ▶ **State set \mathcal{S} = decision set $\mathcal{D} = \{0, 1, \dots, 9\}$.**
 - ▶ **State = actual class, Decision = recognized class**
 - ▶ Loss function:

$$l(s, d) = \begin{cases} 0, & d = s \\ 1, & d \neq s \end{cases}$$

$$\delta^*(\vec{x}) = \arg \min_d \sum_s \underbrace{l(s, d)}_{0 \text{ if } d=s} P(s|\vec{x}) = \arg \min_d \sum_{s \neq d} P(s|\vec{x})$$

Obviously $\sum_s P(s|\vec{x}) = 1$, then:

$$P(d|\vec{x}) + \sum_{s \neq d} P(s|\vec{x}) = 1$$

Inserting into above:

$$\delta^*(\vec{x}) = \arg \min_d [1 - P(d|\vec{x})] = \arg \max_d P(d|\vec{x})$$

A special case - Bayesian *classification*

- ▶ Bayesian classification is a special case of statistical decision theory:
 - ▶ Attribute vector $\vec{x} = (x_1, x_2, \dots)$: pixels 1, 2, ...
 - ▶ **State set \mathcal{S} = decision set $\mathcal{D} = \{0, 1, \dots, 9\}$.**
 - ▶ **State = actual class, Decision = recognized class**
 - ▶ Loss function:

$$l(s, d) = \begin{cases} 0, & d = s \\ 1, & d \neq s \end{cases}$$

$$\delta^*(\vec{x}) = \arg \min_d \sum_s \underbrace{l(s, d)}_{0 \text{ if } d=s} P(s|\vec{x}) = \arg \min_d \sum_{s \neq d} P(s|\vec{x})$$

Obviously $\sum_s P(s|\vec{x}) = 1$, then:

$$P(d|\vec{x}) + \sum_{s \neq d} P(s|\vec{x}) = 1$$

Inserting into above:

$$\delta^*(\vec{x}) = \arg \min_d [1 - P(d|\vec{x})] = \arg \max_d P(d|\vec{x})$$

References I

Further reading: Chapter 13 and 14 of [6]. Books [1] and [2] are classical textbooks in the field of pattern recognition and machine learning. An interesting insights into how people think and interact with probabilities are presented in [4] (in Czech as [5])

[1] Christopher M. Bishop.

Pattern Recognition and Machine Learning.

Springer Science+Business Media, New York, NY, 2006.

PDF freely downloadable.

[2] Richard O. Duda, Peter E. Hart, and David G. Stork.

Pattern Classification.

John Wiley & Sons, 2nd edition, 2001.

[3] Zdeněk Kotek, Petr Vysoký, and Zdeněk Zdráhal.

Kybernetika.

SNTL, 1990.

References II

[4] Leonard Mlodinow.

The Drunkard's Walk. How Randomness Rules Our Lives.

Vintage Books, 2008.

[5] Leonard Mlodinow.

Život je jen náhoda. Jak náhoda ovlivňuje naše životy.

Slovart, 2009.

[6] Stuart Russell and Peter Norvig.

Artificial Intelligence: A Modern Approach.

Prentice Hall, 3rd edition, 2010.

<http://aima.cs.berkeley.edu/>.