

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

May 6, 2020

spojeni velmi dobre	25	86%
dostacujici	3	10%
spatne	1	3%

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

May 6, 2020

spojeni velmi dobre	25	86%
dostacujici	3	10%
spatne	1	3%

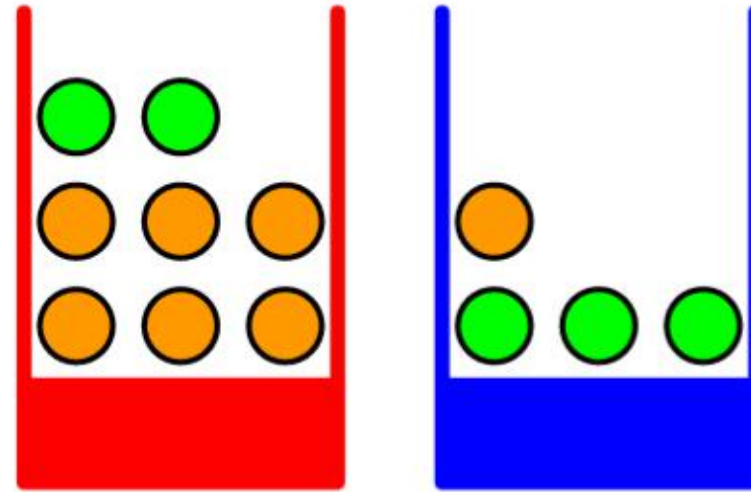
1 / 24

(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

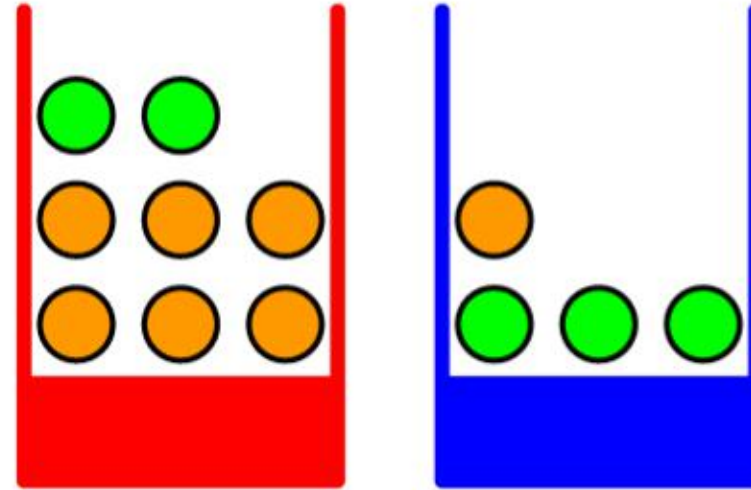


- ▶ Scenario: Pick a box (say red box in 40% cases), then pick a fruit at random
- ▶ (Frequent) questions:
 - ▶ What is the overall probability that the selection procedure will pick an apple?
 - ▶ Given that we have chosen an orange, what is the probability that the box we chose was the blue one?

Example from Chapter 1.2 [1]

Picking fruits. What is the probability that ...?

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

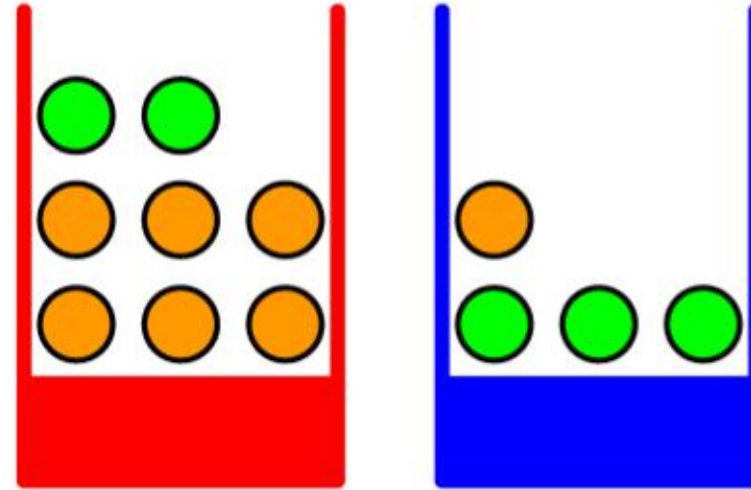


Scenario: Pick a box (say red box in 40% cases), then pick a fruit at random. What is the overall probability that the selection procedure will pick an apple?

- A: $11/20$
- B: $6/8$
- C: $1/2$
- D: Different value.

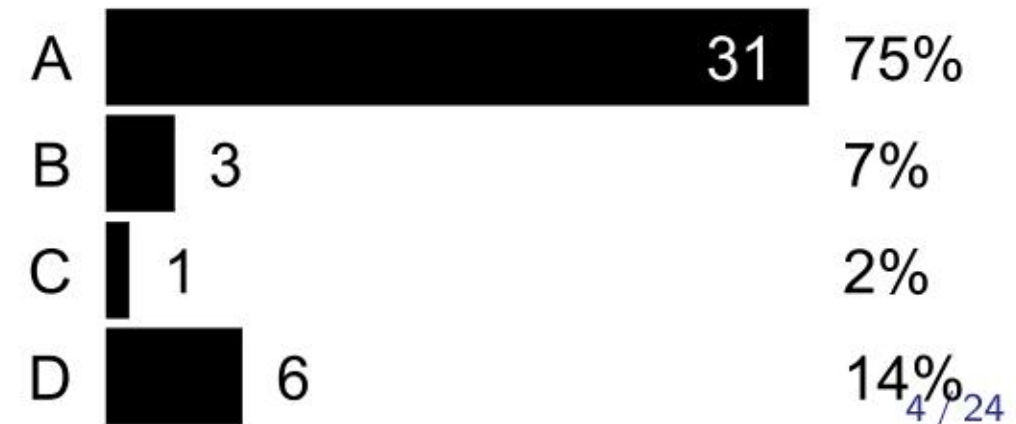
Picking fruits. What is the probability that ...?

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



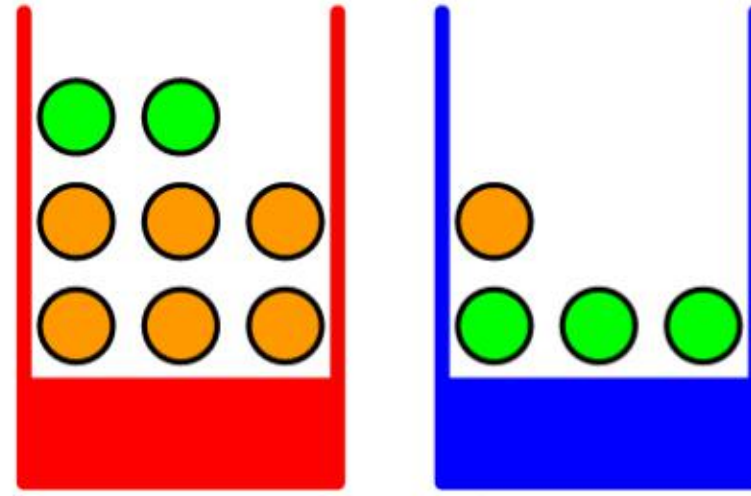
Scenario: Pick a box (say red box in 40% cases), then pick a fruit at random. What is the overall probability that the selection procedure will pick an apple?

- A: $11/20$
- B: $6/8$
- C: $1/2$
- D: Different value.



Picking fruits. What is the probability that ...?

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



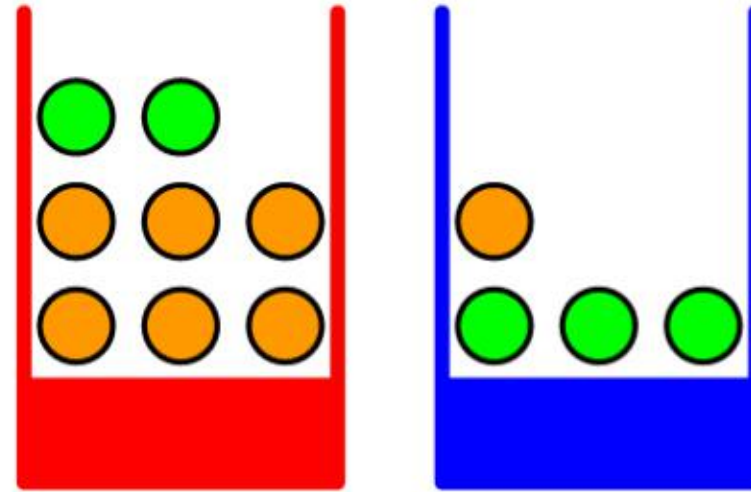
Scenario: Pick a box (say red box in 40% cases), then pick a fruit at random. What is the overall probability that the selection procedure will pick an apple?

- A: $11/20$
- B: $6/8$
- C: $1/2$
- D: Different value.

What is the probability that the selection procedure will pick an apple?

Picking fruits. What is the probability that ... ?

- ▶ 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. What is the overall probability that the selection procedure will pick an apple?

A: $11/20$

B: $6/8$

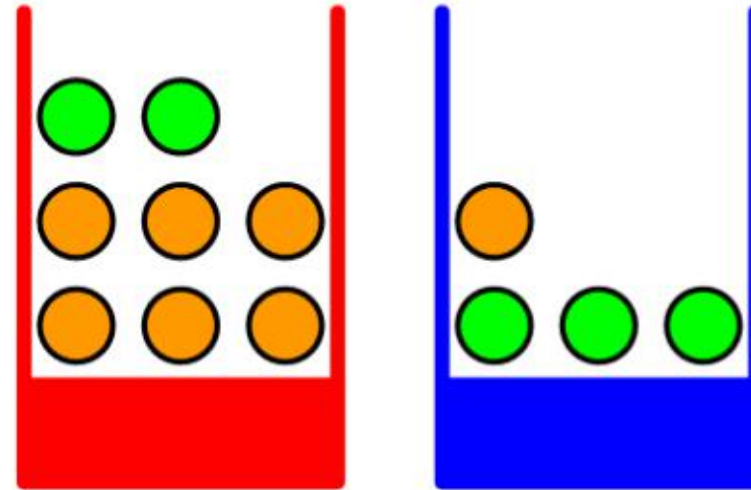
C: $1/2$

D: Different value.

What is the probability that the selection procedure will pick an ~~apple~~ *orange* $9/20$

Picking fruits. What is the probability that ...?

- ▶ 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.

Given that we have chosen an orange, what is the probability that the box we chose was the blue one?

- A: $1/4$
- B: $3/5$
- C: $1/3$
- D: Different value.

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
- ▶ $P(X = x_i, Y = y_j)$ – probability that X takes the value x_i and Y takes y_j – 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 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
- ▶ $P(X = x_i, Y = y_j)$ – probability that X takes the value x_i and Y takes y_j – 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: $P(X) = \sum_Y P(X, Y)$

A handwritten diagram illustrating joint and marginal probabilities. It features a table with two rows and two columns. The columns are labeled 'apple' and 'orange' at the top. The rows are labeled 'X = blue' and 'X = red' on the left. Each cell in the table contains a dot. To the right of the table, there is a vertical line followed by the Greek letter sigma (Σ), indicating a summation operation.

	apple	orange	
X = blue	.	.	Σ
X = red	.	.	

Rules of probability and notation II

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

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

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

- ▶ **Independence** : $P(X, Y) = P(X)P(Y)$

Rules of probability and notation II

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

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

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

- ▶ **Independence** : $P(X, Y) = P(X)P(Y)$

Rules of probability and notation II

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

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

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

- ▶ **Independence** : $P(X, Y) = P(X)P(Y)$

Rules of probability and notation II

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

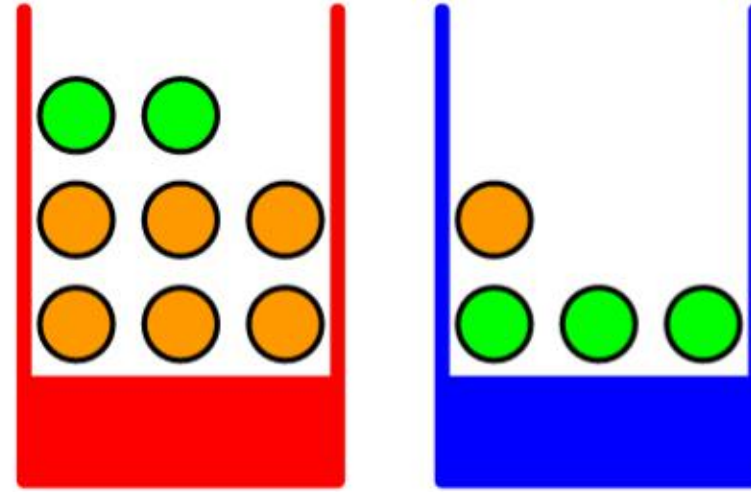
$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

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

- ▶ **Independence** : $P(X, Y) = P(X)P(Y)$ $P(B=r, F=\sigma) = P(B=r)P(F=\sigma)$

Boxes and Fruits: posterior? likelihood? prior? evidence?

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



▶ posterior
after observation

▶ $P(B)$

▶ likelihood
of an observation

▶ $P(F)$

▶ prior
before observation

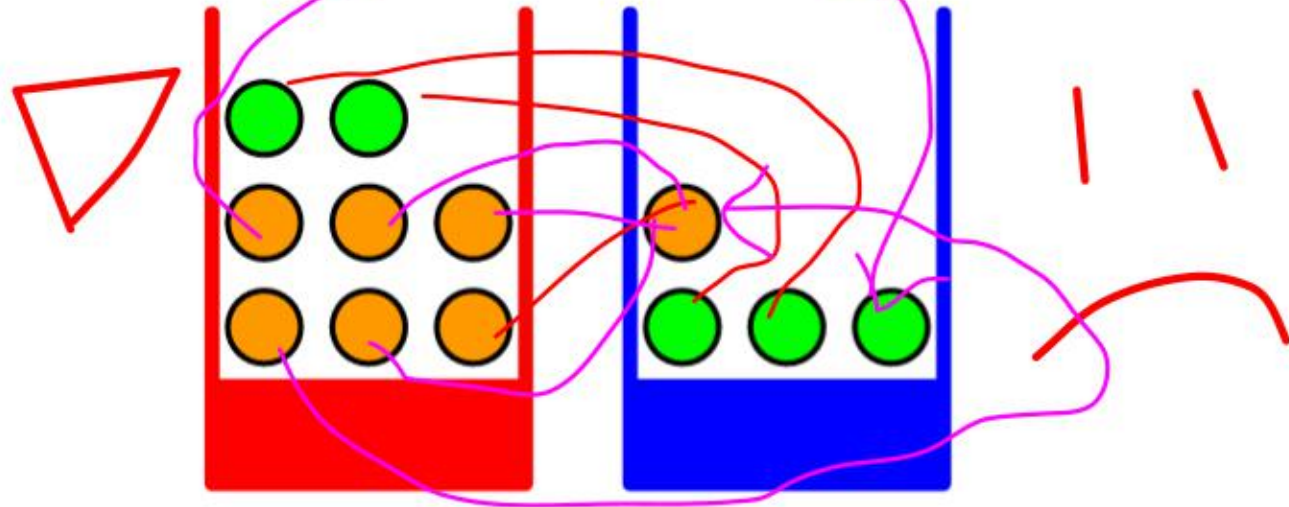
▶ $P(F | B)$

▶ evidence
total observations

▶ $P(B | F)$

Boxes and Fruits: posterior? likelihood? prior? evidence?

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



▶ posterior
after observation

$$P(B) = \frac{1}{6}$$

▶ likelihood
of an observation

$$P(F) = \frac{1}{6}$$

▶ prior
before observation

$$P(F | B)$$

▶ evidence
total observations

$$P(B | F)$$



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

A doctor calls: “Your HIV test is positive, 999/1000 you will die in 10 years. I’m sorry ...”.
Insurance company does not want to insure a 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, has family, no drugs, no risk behavior.

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

AIDS

A doctor calls: “Your HIV test is positive, 999/1000 you will die in 10 years. I’m sorry ...”.
Insurance company does not want to insure a 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, has family, no drugs, no risk behavior.

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

A doctor calls: “Your HIV test is positive, 999/1000 you will die in 10 years. I’m sorry ...”.
Insurance company does not want to insure a 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, has family, no drugs, no risk behavior.

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

A doctor calls: “Your HIV test is positive, 999/1000 you will die in 10 years. I’m sorry ...”.
Insurance company does not want to insure a 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, has family, no drugs, no risk behavior.

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

A doctor calls: “Your HIV test is positive, 999/1000 you will die in 10 years. I’m sorry ...”.
Insurance company does not want to insure a 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, has family, no drugs, no risk behavior.

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

A doctor calls: “Your HIV test is positive, 999/1000 you will die in 10 years. I’m sorry ...”.
Insurance company does not want to insure a 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, has family, no drugs, no risk behavior.

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

A doctor calls: "Your HIV test is positive, 999/1000 you will die in 10 years. I'm sorry ...".
Insurance company does not want to insure a married couple.

▶ Was the doctor right?

▶ Was the insurance company rational?

$$T = \{P, N\}$$
$$H = \{h, s\}$$
$$P(H=h | T=p) = ?$$

What the doctor (and the company) knew:

▶ HIV test falsely positive only in 1 case out of 1000.

▶ Heterosexual male, has family, no drugs, no risk behavior.

$$P(T=p | H=h) = \frac{1}{1000}$$

$$P(T=n | H=s) \rightarrow \emptyset$$

$$P(h|p) = \frac{P(p|h)P(h)}{P(p)}$$

Decision: guilty or not? (people of CA vs Collins, 1968) [4]

- ▶ Robbery, LA 1964, fuzzy evidence of the offenders:
 - ▶ female, around 65 kg
 - ▶ wearing something dark
 - ▶ hairs of light color, between light and dark blond, in a ponytail
- ▶ At the same time, additional evidence close to the crime scene:
 - ▶ loud scream, yelling, looking at the this direction
 - ...
 - ▶ a woman sitting into a yellow car
 - ▶ car starts immediately and passes close to the additional witness
 - ▶ a black man with beard and moustache was driving
- ▶ No more evidence
- ▶ Testimony of both the victim and the witness not unambiguous (didn't recognize suspects)
- ▶ Still, the suspects were sentenced to jail.

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

A doctor calls: "Your HIV test is positive, 999/1000 you will die in 10 years. I'm sorry ...".
Insurance company does not want to insure a married couple.

► Was the doctor right?

► Was the insurance company rational?

$$T = \{P, N\}$$
$$H = \{h, s\}$$
$$P(H=h | T=p) = ?$$

What the doctor (and the company) knew:

► HIV test falsely positive only in 1 case out of 1000.

► Heterosexual male, has family, no drugs, no risk behavior.

$$P(T=p | H=h) = \frac{1}{1000}$$

$$P(T=n | H=s) \rightarrow \emptyset$$

$$P(h|p) = \frac{P(p|h)P(h)}{P(p)}$$

Decision: guilty or not? (people of CA vs Collins, 1968) [4]

- ▶ Robbery, LA 1964, fuzzy evidence of the offenders:
 - ▶ female, around 65 kg
 - ▶ wearing something dark
 - ▶ hairs of light color, between light and dark blond, in a ponytail
- ▶ At the same time, additional evidence close to the crime scene:
 - ▶ loud scream, yelling, looking at the this direction
 - ...
 - ▶ a woman sitting into a yellow car
 - ▶ car starts immediately and passes close to the additional witness
 - ▶ a black man with beard and moustache was driving
- ▶ No more evidence
- ▶ Testimony of both the victim and the witness not unambiguous (didn't recognize suspects)
- ▶ Still, the suspects were sentenced to jail.

Decision: guilty or not? (people of CA vs Collins, 1968) [4]

- ▶ Robbery, LA 1964, fuzzy evidence of the offenders:
 - ▶ female, around 65 kg
 - ▶ wearing something dark
 - ▶ hairs of light color, between light and dark blond, in a ponytail
- ▶ At the same time, additional evidence close to the crime scene:
 - ▶ loud scream, yelling, looking at the this direction
 - ...
 - ▶ a woman sitting into a yellow car
 - ▶ car starts immediately and passes close to the additional witness
 - ▶ a black man with beard and moustache was driving
- ▶ No more evidence
- ▶ Testimony of both the victim and the witness not unambiguous (didn't recognize suspects)
- ▶ Still, the suspects were sentenced to jail.

$$P(\text{yellow}) = \frac{1}{10}$$

$$P(M=\text{beard}) = \frac{1}{4}$$

$$P(\bar{C} = \dots) = \frac{1}{10}$$

$$P(F \text{ ponytail}) = \frac{1}{10}$$

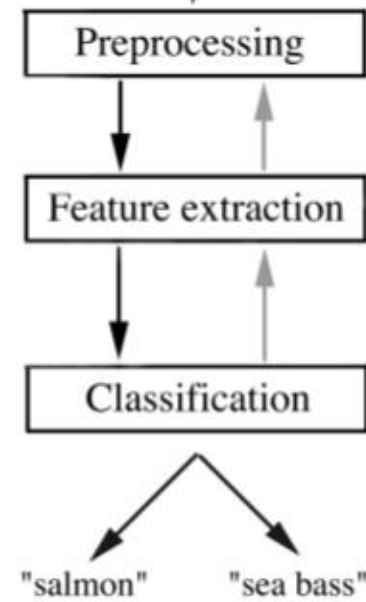
$$P(\text{blond}) = \frac{1}{3}$$

$$P(\text{mix race pair}) = \frac{1}{1000}$$

$$P(\text{randomly sel. pair}) = \frac{1}{125000}$$

$$P(\text{guilty pair}) = \frac{1}{3}$$

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

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

- ▶ 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} :

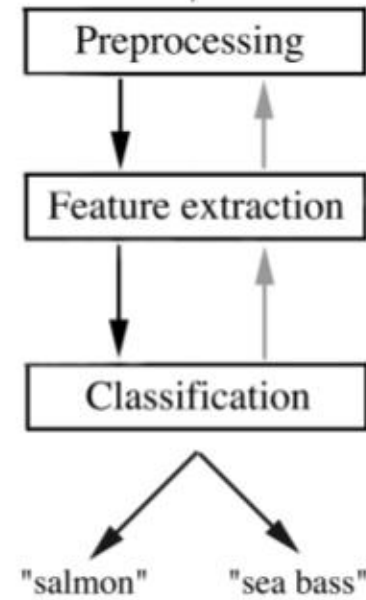
$$\delta^*(\vec{x}) = \arg \max_j P(s_j | \vec{x})$$

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

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

- ▶ Can we do (classify) better?

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

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

- ▶ 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:

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

▶ Can we do (classify) better?

Fish – classification using probability

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

- ▶ 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:

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

- ▶ Can we do (classify) better?

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 $(\vec{x}_1, s_1), (\vec{x}_2, s_2), \dots, (\vec{x}_l, s_l)$
 - ▶ 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 $\underset{\text{arg, max } s_j}{P}(s|\vec{x})$
- ▶ It has to be estimated from already classified examples – training data
- ▶ For discrete \vec{x} , training examples $(\vec{x}_1, s_1), (\vec{x}_2, s_2), \dots, (\vec{x}_l, s_l)$
 - ▶ 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 $(\vec{x}_1, s_1), (\vec{x}_2, s_2), \dots, (\vec{x}_l, s_l)$
 - ▶ 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:

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

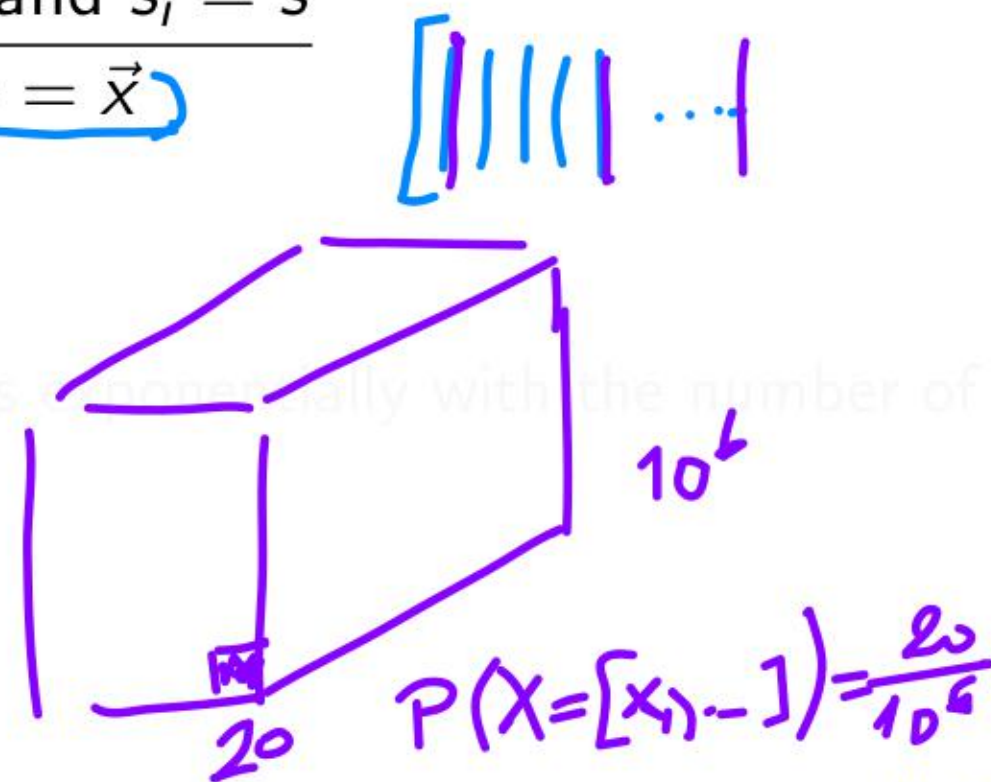
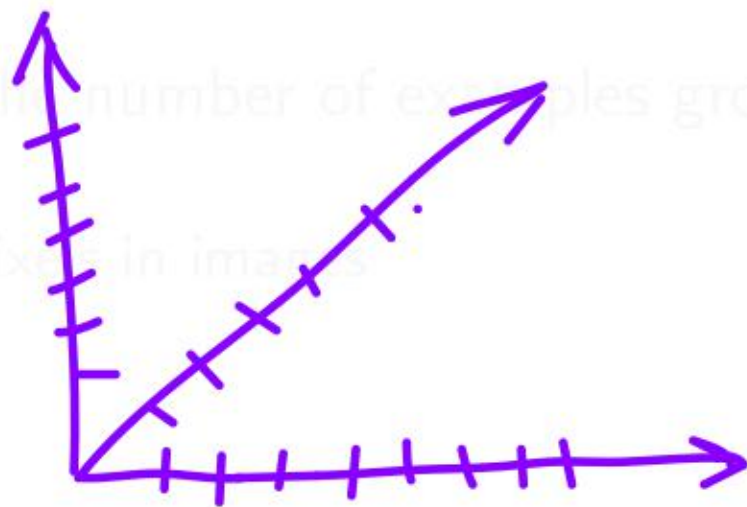
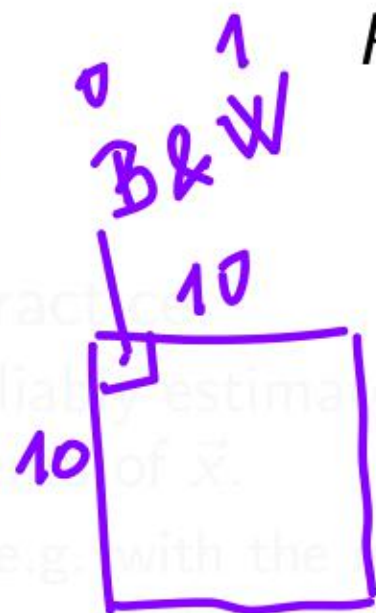
- ▶ Hard in practice:
 - ▶ To reliably estimate $P(s|\vec{x})$, the number of examples grows exponentially with the number of elements of \vec{x} .
 - ▶ e.g. with the number of pixels in images
 - ▶ curse of dimensionality
 - ▶ denominator often 0

Bayes classification in practice

- ▶ Usually we are not given $P(s|\vec{x}) = \frac{P(\vec{x}|s)P(s)}{P(\vec{x})}$
- ▶ It has to be estimated from already classified examples – training data
- ▶ For discrete \vec{x} , training examples $(\vec{x}_1, s_1), (\vec{x}_2, s_2), \dots, (\vec{x}_l, s_l)$
 - ▶ 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:

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

$\dim(\vec{x}) = 100$
 2^{100}



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 $(\vec{x}_1, s_1), (\vec{x}_2, s_2), \dots, (\vec{x}_l, s_l)$
 - ▶ 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

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

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

$2^{100} \rightarrow 100.2$

- ▶ 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?

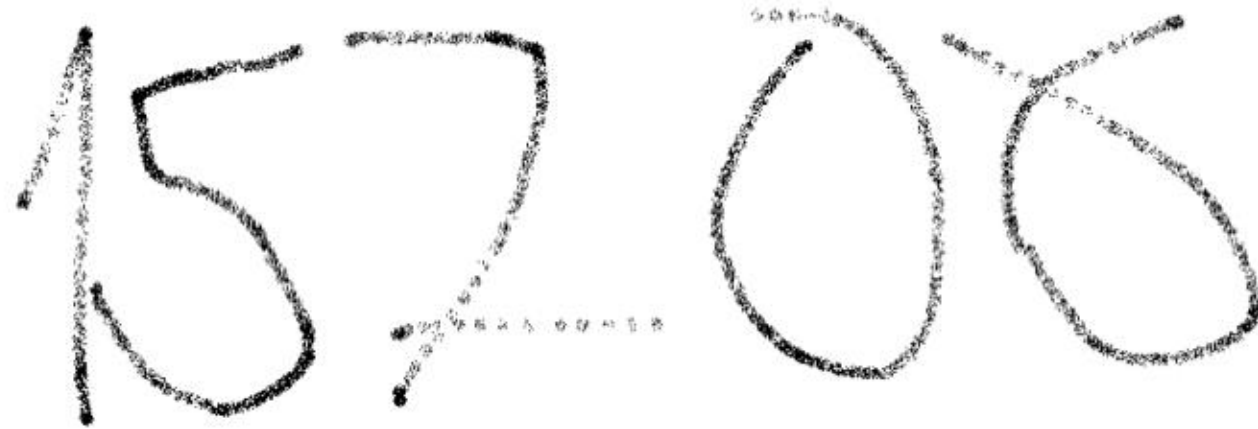
▶ 15700? 15706? 15200? 15206?

▶ What is the optimal decision ?

▶ Both examples fall into the same framework.

Decision making under uncertainty

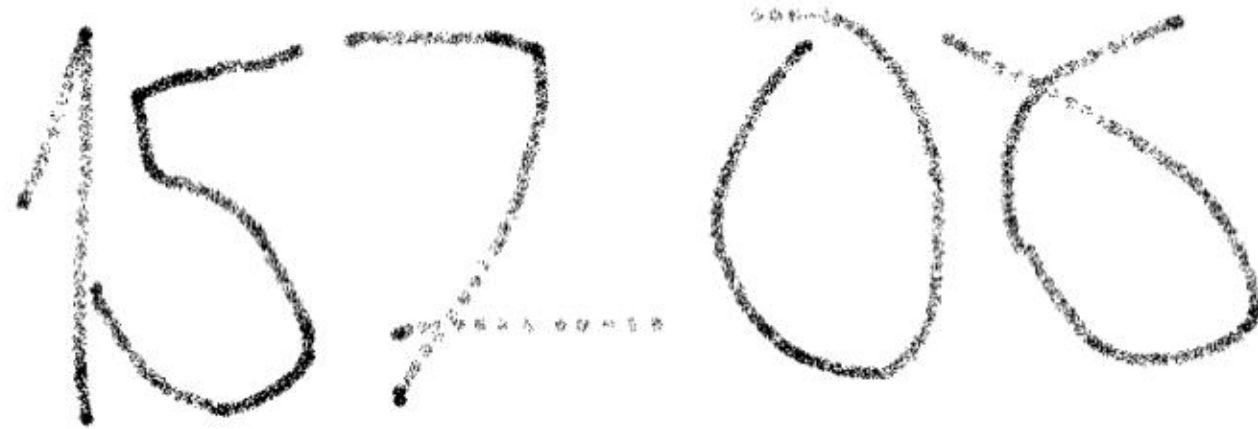
- ▶ An important feature of intelligent systems
 - ▶ make the best possible decision
 - ▶ in uncertain conditions.
- ▶ **Example:** Take a tram OR subway from *A* to *B*?
 - ▶ Tram: timetables imply a quicker route, but adherence uncertain.
 - ▶ Subway: longer route, but adherence almost certain.
- ▶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

A handwritten ZIP code '157006' is shown in a dark, grainy font. The digits are somewhat irregular and connected, with some noise or artifacts around the lines, making it difficult to read precisely. The '1' is a simple vertical stroke, '5' is a loop, '7' has a horizontal top bar and a diagonal stem, and the '006' consists of two loops and a final stroke.

- ▶ 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?

A handwritten ZIP code '157006' is shown. The digits are somewhat blurry and the ink is dark. The '1' is a simple vertical stroke. The '5' has a curved top. The '7' has a horizontal top bar and a diagonal stroke. The '0' is a simple oval. The '0' is a simple oval. The '6' has a top loop and a tail that curves to the right.

- ▶ 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 $l(d, s)$:

$l(s, d)$	$d = nothing$	$d = pizza$	$d = g.T.c.$
$s = good$	0	2	4
$s = average$	5	3	5
$s = bad$	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

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

= 1

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:

$\delta(x)$	$x = \text{mild}$	$x = \text{irritated}$	$x = \text{upset}$	$x = \text{alarming}$
$\delta_1(x) =$	<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:

$\delta(x)$	$x = \text{mild}$	$x = \text{irritated}$	$x = \text{upset}$	$x = \text{alarming}$
$\delta_1(x) =$	<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** .

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

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

$I(s, d)$	$d = \textit{nothing}$	$d = \textit{pizza}$	$d = \textit{g.T.c.}$
$s = \textit{good}$	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

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

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

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

$I(s, d)$	$d = \textit{nothing}$	$d = \textit{pizza}$	$d = \textit{g.T.c.}$	
$s = \textit{good}$	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 $r(\delta) = \sum_x \sum_s I(s, \delta(x)) P(x, s)$

$I(s, d)$	$d = \text{nothing}$	$d = \text{pizza}$	$d = \text{g.T.c.}$	
$s = \text{good}$	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

$\delta(x)$	$x = \text{mild}$	$x = \text{irritated}$	$x = \text{upset}$	$x = \text{alarming}$
$\delta_1(x) =$	<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? $P(x, s) = P(s|x)P(x)$