

Pattern Recognition

Probability Theory

Probability Space

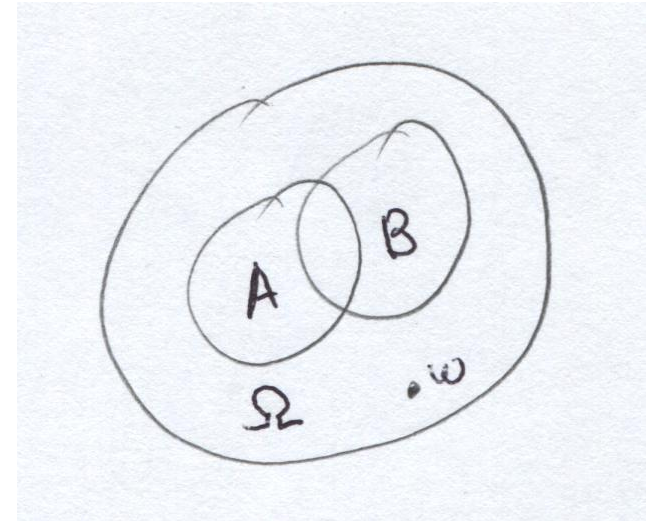
is a three-tuple (Ω, σ, P) with:

- Ω – the set of elementary events
- σ – algebra
- P – probability measure

σ -algebra over Ω is a system of subsets, i.e. $\sigma \subseteq \mathcal{P}(\Omega)$ (\mathcal{P} is the power set) with:

- $\Omega \in \sigma$
- $A \in \sigma \Rightarrow \Omega/A \in \sigma$
- $A_i \in \sigma, i = 1, \dots, n \Rightarrow \bigcup_{i=1}^n A_i \in \sigma$

σ is closed with respect to the complement and countable conjunction
It follows – $\emptyset \in \sigma$, countable disjunction (due to the De Morgan's laws)



Probability Space

Examples:

- $\sigma = \{\emptyset, \Omega\}$ (smallest) and $\sigma = \mathcal{P}(\Omega)$ (largest) σ -algebras over Ω
- the minimal σ -algebra over Ω containing a particular subset $A \subset \Omega$ is $\sigma = \{\emptyset, A, \Omega \setminus A, \Omega\}$
- Ω discrete and finite, $\sigma = 2^\Omega$
- $\Omega = \mathbb{R}$, the Borel-algebra (contains all intervals amongst others)
- etc.

Probability Measure

$P : \sigma \rightarrow [0, 1]$ Is a “measure” (Π) with the normalizing $P(\Omega) = 1$

σ -additivity: let $A_i \in \sigma$ be pairwise disjoint subsets, i.e. $A_i \cap A_{i'} = \emptyset$, then

$$P\left(\bigcup_i A_i\right) = \sum_i P(A_i)$$

Note: there are sets, for which there is no measure.

Examples: the set of irrational numbers, function spaces \mathbb{R}^∞ etc.

Banach–Tarski paradox:



(For us) practically relevant cases

- The set Ω is “good-natured”, i.e. \mathbb{R}^n , discrete finite sets etc.
- $\sigma = \mathcal{P}(\Omega)$, i.e. the algebra is the power set
- We often consider a (composite) “event” $A \subseteq \Omega$ as the union of the elementary ones
- Probability of an event is

$$P(A) = \sum_{\omega \in A} P(\omega)$$

Random variables

Here a special case – **real-valued** random variables.

A random variable ξ for a probability space (Ω, σ, P) is a mapping $\xi : \Omega \rightarrow \mathbb{R}$, satisfying

$$\{\omega : \xi(\omega) \leq r\} \in \sigma \quad \forall r \in \mathbb{R}$$

(always holds for power sets $\sigma = \mathcal{P}(\Omega)$).

Note: elementary events are **not numbers** – they are elements of an abstract set Ω

Random variables in contrast are numbers, i.e. they can be summed up, subtracted, squared etc.

Distributions

Cumulative distribution function of a random variable ξ :

$$F_{\xi}(r) = P(\{\omega : \xi(\omega) \leq r\})$$

Probability distribution of a discrete random variable $\xi : \Omega \rightarrow \mathbb{Z}$:

$$p_{\xi}(r) = P(\{\omega : \xi(\omega) = r\})$$

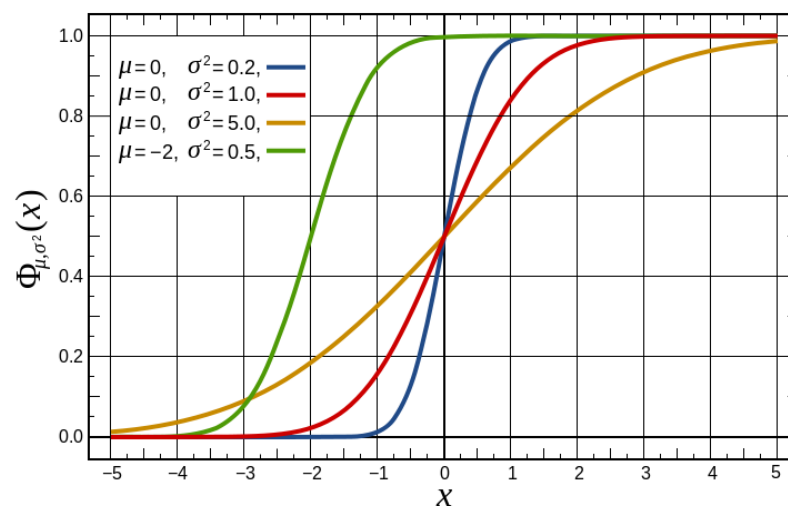
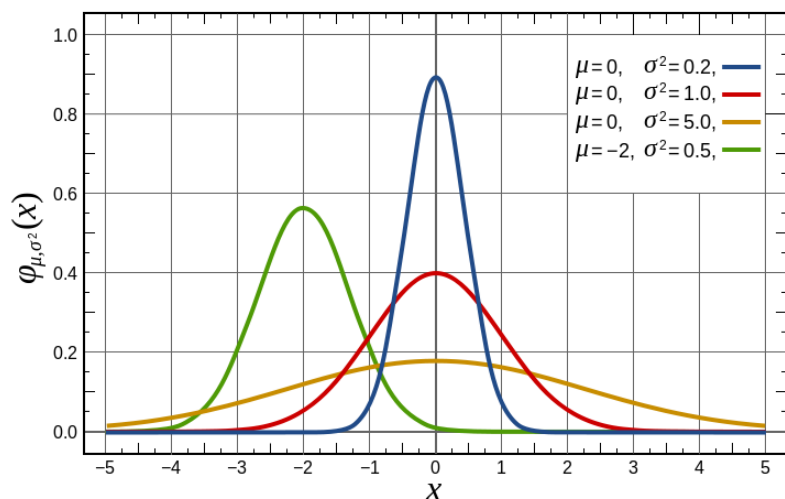
Probability density of a continuous random variable $\xi : \Omega \rightarrow \mathbb{R}$:

$$p_{\xi}(r) = \frac{\partial F_{\xi}(r)}{\partial r}$$

Distributions

Why it is necessary to do it so complicated (through the cumulative distribution function)?

Example – a Gaussian.



Probability of any particular real value is zero \rightarrow a “direct” definition of a “probability distribution” is senseless ☹

It is indeed possible through the cumulative distribution function.

Mean

A mean (average, expectation...) of a random variable ξ is

$$\mathbb{E}_P(\xi) = \sum_{\omega \in \Omega} P(\omega) \cdot \xi(\omega) = \sum_r \sum_{\omega: \xi(\omega)=r} P(\omega) \cdot r = \sum_r p_\xi(r) \cdot r$$

Arithmetic mean is a special case:

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i = \sum_r p_\xi(r) \cdot r$$

with

$$x \equiv r \quad \text{and} \quad p_\xi(r) = \frac{1}{N}$$

(uniform probability distribution)

Mean

The probability of an event $A \subset \Omega$ can be expressed as the mean value of a corresponding “indicator”-variable:

$$P(A) = \sum_{\omega \in A} P(\omega) = \sum_{\omega \in \Omega} P(\omega) \cdot \xi(\omega)$$

with

$$\xi(\omega) = \begin{cases} 1 & \text{if } \omega \in A \\ 0 & \text{otherwise} \end{cases}$$

Often, the set of elementary events can be associated with a random variable (just enumerate all $\omega \in \Omega$).

Then one can speak about a “probability distribution over Ω ” (instead of the probability measure).

Example 1 – numbers of a die

The set of elementary events: $\Omega = \{a, b, c, d, e, f\}$

Probability measure: $P(\{a\}) = 1/6$, $P(\{c, f\}) = 1/3$...

Random variable: $\xi(a) = 1$, $\xi(b) = 2$... $\xi(f) = 6$

Cumulative distribution: $F_\xi(3) = 1/2$, $F_\xi(4.5) = 2/3$...

Probability distribution: $p_\xi(1) = p_\xi(2) = \dots = p_\xi(6) = 1/6$

Mean value: $\mathbb{E}_P(\xi) = 3.5$

Another random variable (squared numbers of a die):

$$\xi'(a) = 1, \xi'(b) = 4 \dots \xi'(f) = 36$$

Mean value:

$$\mathbb{E}_P(\xi') = 15\frac{1}{6}$$

Note: $\mathbb{E}_P(\xi') \neq \mathbb{E}_P^2(\xi)$

Example 2 – two independent dice numbers

The set of elementary events (6x6 faces):

$$\Omega = \{a, b, c, d, e, f\} \times \{a, b, c, d, e, f\}$$

Probability measure: $P(\{ab\}) = 1/36$, $P(\{cd, fa\}) = 1/18$...

Two random variables:

1. The number of the first die:

$$\xi_1(ab) = 1, \xi_1(ac) = 1 \dots \xi_1(ef) = 5 \dots$$

2. The number of the second die

$$\xi_2(ab) = 2, \xi_2(ac) = 3 \dots \xi_2(ef) = 6 \dots$$

Probability distributions:

$$p_{\xi_1}(1) = p_{\xi_1}(2) = \dots = p_{\xi_1}(6) = 1/6$$

$$p_{\xi_2}(1) = p_{\xi_2}(2) = \dots = p_{\xi_2}(6) = 1/6$$

Example 2 – two independent dice numbers

Consider the new random variable

$$\xi = \xi_1 + \xi_2$$

The probability distribution p_ξ is not uniform anymore 😊

$$p_\xi \sim (1, 2, 3, 4, 5, 6, 5, 4, 3, 2, 1)$$

Mean value is $\mathbb{E}_P(\xi) = 7$

6	7	8	9	10	11	12
5	6	7	8	9	10	11
4	5	6	7	8	9	10
3	4	5	6	7	8	9
2	3	4	5	6	7	8
$\xi_2 = 1$	2	3	4	5	6	7
	$\xi_1 = 1$	2	3	4	5	6

In general for mean values:

$$\mathbb{E}_P(\xi_1 + \xi_2) = \sum_{\omega \in \Omega} P(\omega) \cdot (\xi_1(\omega) + \xi_2(\omega)) = \mathbb{E}_P(\xi_1) + \mathbb{E}_P(\xi_2)$$

Random variables of higher dimension

Analogously: Let $\xi : \Omega \rightarrow \mathbb{R}^n$ be a mapping ($n = 2$ for simplicity), with $\xi = (\xi_1, \xi_2)$, $\xi_1 : \Omega \rightarrow \mathbb{R}$ and $\xi_2 : \Omega \rightarrow \mathbb{R}$

Cumulative distribution function:

$$F_{\xi}(r, s) = P(\{\omega : \xi_1(\omega) \leq r\} \cap \{\omega : \xi_2(\omega) \leq s\})$$

Joint probability distribution (discrete):

$$p_{\xi=(\xi_1, \xi_2)}(r, s) = P(\{\omega : \xi_1(\omega) = r\} \cap \{\omega : \xi_2(\omega) = s\})$$

Joint probability density (continuous):

$$p_{\xi=(\xi_1, \xi_2)}(r, s) = \frac{\partial F_{\xi}(r, s)}{\partial r \partial s}$$

Independence

Two events $A \subset \Omega$ and $B \subset \Omega$ are **independent**, if

$$P(A \cap B) = P(A) \cdot P(B)$$

Interesting:

Events A and $\bar{B} = \Omega/B$ are independent, if A and B are independent.

Two random variables are independent, if

$$F_{\xi=(\xi_1, \xi_2)}(r, s) = F_{\xi_1}(r) \cdot F_{\xi_2}(s) \quad \forall r, s$$

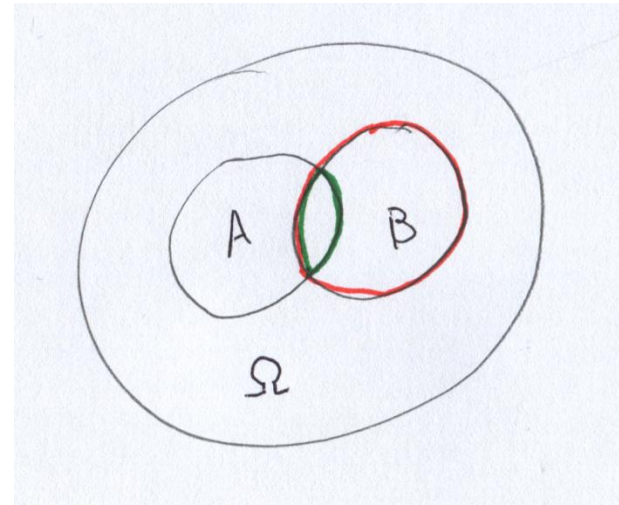
It follows (example for continuous ξ)

$$p_{\xi}(r, s) = \frac{\partial^2 F_{\xi}(r, s)}{\partial r \partial s} = \frac{\partial F_{\xi_1}}{\partial r} \cdot \frac{\partial F_{\xi_2}}{\partial s} = p_{\xi_1}(r) \cdot p_{\xi_2}(s)$$

Conditional Probabilities

Conditional probability:

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$



Independence (“almost” equivalent): A and B are independent, if

$$P(A|B) = P(A) \quad \text{and/or} \quad P(B|A) = P(B)$$

Bayes’ theorem (formula, rule):

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

Further definitions (for random variables)

Shorthand: $p(x, y) \equiv p_{\xi}(x, y)$

Marginal probability distribution:

$$p(x) = \sum_y p(x, y)$$

Conditional probability distribution:

$$p(x|y) = \frac{p(x, y)}{p(y)}$$

Note: $\sum_x p(x|y) = 1$

Independent probability distributions:

$$p(x, y) = p(x) \cdot p(y)$$

Example

Let the probability to be taken ill be

$$P(ill) = 0.02$$

Let the conditional probability to have a temperature in that case is

$$P(temp|ill) = 0.9$$

However, one may have a temperature without any illness, i.e.

$$P(temp|\overline{ill}) = 0.05$$

What is the probability to be taken ill provided that one has a temperature?

Example

Bayes' rule:

$$P(ill|temp) = \frac{P(temp|ill) \cdot P(ill)}{P(temp)} =$$

(marginal probability in the denominator)

$$\frac{P(temp|ill) \cdot P(ill)}{P(temp|ill) \cdot P(ill) + P(temp|\overline{ill}) \cdot P(\overline{ill})}$$

$$= \frac{0.9 \cdot 0.02}{0.9 \cdot 0.02 + 0.05 \cdot 0.98} \approx 0.27$$

– not so high as expected ☹, the reason – very low **prior** probability to be taken ill

Further topics

The model

Let two random variables be given:

- The first one is typically discrete (i.e. $k \in K$) and is called “class”
- The second one is often continuous ($x \in \mathbb{R}^n$) and is called “observation”

Let the joint probability distribution $p(x, k)$ be “given”.

As k is discrete it is often specified by $p(x, k) = p(k) \cdot p(x|k)$

The recognition task: given x , estimate k .

Usual problems (questions):

- How to estimate k from x ?
- The joint probability is not always explicitly specified.
- The set K is sometimes huge (remember the Hopfield-Networks)

Further topics

The learning task:

Often (almost always) the probability distribution is known up to free parameters. How to choose them (learn from examples)?

Next themes:

1. Recognition, Bayesian Decision Theory
2. Probabilistic (generative) learning, Maximum-Likelihood principle
3. Discriminative models, recognition and learning
4. Support Vector Machines