

Probability for machine learning

Rokhlin D.B.

IMMCS, SFU, 2021

Probability space (Ω, \mathcal{F}, P)

- ▶ Ω : sample space
- ▶ \mathcal{F} : σ -algebra of events (subsets of Ω)
- ▶ P : probability measure

\mathcal{F} satisfies the axioms

- ▶ $\emptyset \in \mathcal{F}$;
- ▶ the condition $A \in \mathcal{F}$ implies that $A^c := \Omega \setminus A \in \mathcal{F}$;
- ▶ the condition $A_i \in \mathcal{F}$, $i \in \mathbb{Z}_+$ implies that $\bigcup_{i=1}^{\infty} A_i \in \mathcal{F}$.

A function $P : \mathcal{F} \mapsto [0, 1]$ is called a probability measure, if

- ▶ $P(\emptyset) = 0$, $P(\Omega) = 1$;
- ▶ for any sequence of pairwise *disjoint* sets $A_i \in \mathcal{F}$, $i = 1, \dots, \infty$ we have the equality

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i).$$

Examples

► $\Omega = \{H, T\}, \mathcal{F} = \{\emptyset, \Omega, \{H\}, \{T\}\} = 2^\Omega,$

$$P(\{H\}) = p \in [0, 1], \quad P(\{T\}) = 1 - p.$$

► $\Omega = [0, 1], \mathcal{F} = \mathcal{B}(0, 1), P((a, b)) = b - a.$

► $\Omega = \mathbb{R}, \mathcal{F} = \mathcal{B}(\mathbb{R}), P(a, b) = \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx.$

Random variables

A function $\xi : \Omega \mapsto \mathbb{R}^d$ is called \mathcal{F} -measurable if

$$\xi^{-1}(B) = \{\omega : \xi(\omega) \in B\} \in \mathcal{F}.$$

Such functions are called random variables.

► $\Omega = \{(a_1, a_2) : a_i \in \{H, T\}\}$, $\mathcal{F}_1 = \sigma(\{H\}, \{T\})$,

$$\xi_1 = I_{\{a_1=H\}}, \quad \xi_2 = I_{\{a_2=H\}},$$

$$I_A(\omega) = \begin{cases} 1, & \omega \in A, \\ 0, & \omega \notin A. \end{cases}$$

ξ_1 is \mathcal{F}_1 -measurable, ξ_2 is not \mathcal{F}_1 -measurable.

- ▶ $(\Omega, \mathcal{F}) = ([0, 1], (\emptyset, [0, 1]))$.
 $\xi : [0, 1] \mapsto \mathbb{R}$ is \mathcal{F} -measurable $\iff \xi = \text{const.}$
- ▶ $(\Omega, \mathcal{F}) = ([0, 1], \mathcal{B}(0, 1))$.
 $\xi : [0, 1] \mapsto \mathbb{R}$ is \mathcal{F} -measurable $\iff \xi$ is Borel.
- ▶ $\Omega = \{(a_1, \dots, a_m) : a_i \in \{H, T\}\}$,
 $\mathcal{F} = \sigma\{\{H, H\}, \{H, T\}, \{T, H\}, \{T, T\}\}$.
 $\xi : \Omega \mapsto \mathbb{R}$ is \mathcal{F} -measurable $\iff \xi$ does not depend on a_3, \dots, a_m .

Integral

Simple functions:

$$\xi(\omega) = \sum_{j=1}^n c_j I_{A_j}(\omega), \quad A_j \in \mathcal{F}, \quad c_j \in \mathbb{R}$$

Clearly, such functions are \mathcal{F} -measurable. Integral of a simple function:

$$\int_{\Omega} \xi(\omega) \mathbf{P}(d\omega) = \int_{\Omega} \xi d\mathbf{P} := \sum_{j=1}^n c_j \mu(A_j).$$

Integral of a non-negative measurable function:

$$\int_{\Omega} \xi d\mathbf{P} = \sup \left\{ \int_{\Omega} \varphi d\mathbf{P}, \text{ where } \varphi \text{ simple and } 0 \leq \varphi \leq \xi \right\}.$$

Infinite values are acceptable.

Furthermore, any \mathcal{F} -measurable function can be represented as a difference of two nonnegative \mathcal{F} -measurable functions: $\xi = \xi^+ - \xi^-$, where

$$\xi^+(\omega) := \max\{\xi, 0\}, \quad \xi^-(\omega) := -\max\{-\xi(\omega), 0\}.$$

If both values $\int_{\Omega} \xi^+ dP$, $\int_{\Omega} \xi^- dP$ are finite, then ξ is called *integrable* and the integral is defined by the formula

$$\int_{\Omega} \xi d\mu = \int_{\Omega} \xi^+ d\mu - \int_{\Omega} \xi^- d\mu.$$

- ▶ expectation: $E\xi = \int_{\Omega} \xi dP$,
- ▶ variance: $\text{Var}(\xi) = E(\xi - E\xi)^2$,
- ▶ standard deviation: $\sqrt{\text{Var}(\xi)}$.

Distribution

An random variable $\xi : \Omega \mapsto \mathbb{R}^d$ induces the probability space $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d), P_\xi)$, where P_ξ is the image measure:

$$P_\xi(B) = P(\xi^{-1}(B)), \quad B \in \mathcal{B}(\mathbb{R}^d),$$

which is called the distribution of ξ .

- ▶ ξ is called discrete if there exist a finite or countable set Z and a probability mass function (p.m.f.) $p_\xi : Z \mapsto [0, 1]$ such that

$$P_\xi(B) = \sum_{z_i \in Z \cap B} p_\xi(z_i), \quad B \in \mathcal{B}(\mathbb{R}^d).$$

- ▶ ξ is called continuous if there exists a Borel function $p_\xi : \mathbb{R}^d \mapsto \mathbb{R}_+$ such that

$$P_\xi(B) = \int_B p_\xi(x) dx, \quad B \in \mathcal{B}(\mathbb{R}^d).$$

Change of variables formula

Let g be a (bounded) Borel function. Then

$$\mathbb{E}g(\xi) = \int_{\mathbb{R}} g(x) \mathbb{P}_{\xi}(dx).$$

- ▶ for discrete ξ : $\mathbb{E}g(\xi) = \sum_{x \in Z} g(x)p_{\xi}(x)$,
- ▶ for continuous ξ : $\mathbb{E}g(\xi) = \int_{\mathbb{R}} g(x)f_{\xi}(x) dx$.

In particular,

- ▶ for discrete ξ : $\mathbb{E}\xi = \sum_{x \in Z} xp_{\xi}(x)$,
- ▶ for continuous ξ : $\mathbb{E}\xi = \int_{\mathbb{R}} xf_{\xi}(x) dx$.

Examples

- Bernoulli distribution: $\xi \sim \text{Ber}(p)$, $Z = \{0, 1\}$, $p_\xi(0) = 1$, $p_\xi(1) = 1 - p$,

$$\mathbb{E}\xi = 0 \cdot p_\xi(0) + 1 \cdot p_\xi(1) = p,$$

$$\mathbb{E}(\xi^2) = 0 \cdot (1 - p) + 1 \cdot p = p,$$

$$\text{Var}(\xi) = \mathbb{E}(\xi^2) - (\mathbb{E}\xi)^2 = p - p^2 = p(1 - p).$$

- Poisson distribution: $\xi \sim \Pi(\lambda)$, $Z = \mathbb{Z}_+$,

$$p_\xi(k) = e^{-\lambda} \frac{\lambda^k}{k!}, \quad k \in \mathbb{Z}_+.$$

$$\mathbb{E}\xi = \sum_{k=0}^{\infty} k p_\xi(k) = e^{-\lambda} \sum_{k=1}^{\infty} \frac{\lambda^k}{(k-1)!} = \lambda \sum_{k=0}^{\infty} p_\xi(k) = \lambda,$$

$$\text{Var}(\xi) = \mathbb{E}\xi^2 - (\mathbb{E}\xi)^2 = \sum_{k=0}^{\infty} k^2 p_\xi(k) - (\mathbb{E}\xi)^2 = \lambda(\lambda + 1) - \lambda^2 = \lambda.$$

- Uniform distribution: $\xi \sim U(a, b)$, $a < b$,

$$f_{\xi}(x) = \begin{cases} \frac{1}{b-a}, & x \in [a, b], \\ 0, & x \notin [a, b]. \end{cases}$$

$$\mathbb{E}\xi = \frac{a+b}{2}, \quad \text{Var}(\xi) = \frac{(b-a)^2}{12}.$$

- Gaussian (normal) distribution: $\xi \sim N(0, 1)$,

$$p_{\xi}(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right), \quad x \in \mathbb{R}.$$

$$\mathbb{E}\xi = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x e^{-x^2/2} dx = 0.$$

$$\begin{aligned} \text{Var}(\xi) = \mathbb{E}\xi^2 &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x^2 e^{-x^2/2} dx = -\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x de^{-x^2/2} \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-x^2/2} dx = 1. \end{aligned}$$

Distribution functions

$F_\xi(x) = P(\xi \leq x)$ is called a cumulative distribution function (cdf) of ξ .

$$F_\xi(x) = \begin{cases} \sum_{y \leq x} p_\xi(y), & \xi \text{ discrete,} \\ \int_{-\infty}^x f_\xi(y) dy, & \xi \text{ continuous.} \end{cases}$$

$$F_\xi(x) \leq F_\xi(y), \quad \text{if } x \leq y,$$

$$\lim_{x \rightarrow -\infty} F_\xi(x) = 0, \quad \lim_{x \rightarrow +\infty} F_\xi(x) = 1,$$

Moreover, the distribution functions are right-continuous:

$$\lim_{y \downarrow x} F_\xi(y) = \lim_{Q \ni y \downarrow x} P(\xi \leq y) = P(\cap_{y \geq x, y \in Q} \{\xi \leq y\}) = P(\xi \leq x) = F_\xi(x).$$

The second equality uses the continuity of the probability measure.

For a continuous random variable $F'_\xi(x) = f_\xi(x)$.

Let $\xi \sim N(0, 1)$, then

$$\Phi(x) := F_{\xi}(x) = P(\xi \leq x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-y^2/2} dy.$$

“Three-sigma” rule. $\Phi(3) \approx 0.9987$. Let $\eta \in N(\mu, \sigma^2)$. Then

$$\xi = \frac{\eta - \mu}{\sigma} \in N(0, 1).$$

$$\begin{aligned} P(|\eta - \mu| \leq 3\sigma) &= P(|\xi| \leq 3) = P(\xi \leq 3) - P(\xi \leq -3) \\ &= P(\xi \leq 3) - (1 - P(\xi \leq 3)) = 2P(\xi \leq 3) - 1 \\ &= 2\Phi(3) - 1 \approx 0.997 \end{aligned}$$

Conditional probability

Let $P(B) > 0$. The *conditional probability* of A given B is the number

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

This formula shows what portion of the probability attributed to B relates to A .

Example

Roll 2 symmetric 6-sided dices. The first die shows 5. Find the conditional probability that the sum of the obtained points will be 10.

$$\Omega = \{(a_1, a_2) : a_i \in \{1, \dots, 6\}\}, \quad \mathcal{F} = 2^\Omega,$$

$$B = \{\omega : a_1 = 5\}, \quad A = \{\omega : a_1 + a_2 = 10\},$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(\{\omega : a_1 = 5, a_2 = 5\})}{P(\{\omega : a_1 = 5\})} = \frac{1/36}{1/6} = 1/6.$$

Multiplication rule

Let $P(A_1 \cap \dots \cap A_m) > 0$. The equality

$$P(A_1 \cap \dots \cap A_m) = P(A_1)P(A_2|A_1) \dots P(A_m|A_1 \cap \dots \cap A_{m-1})$$

is called the multiplication rule.

Proof.

$(m = 3)$

$$\begin{aligned} P(A_1 \cap A_2 \cap A_3) &= P(A_3|A_1 \cap A_2)P(A_1 \cap A_2) \\ &= P(A_3|A_1 \cap A_2)P(A_2|A_1)P(A_1) \end{aligned}$$

Example

There is a 36-card deck. 3 cards are drawn at random. What is the probability that there will get no spades?

Let A_i be the event that the card drawn at step i is not a spade. We are interested in the probability of the event $A_1 \cap A_2 \cap A_3$. By the multiplication rule,

$$P(A_1 \cap A_2 \cap A_3) = P(A_1)P(A_2|A_1)P(A_3|A_1 \cap A_2).$$

Clearly,

$$P(A_1) = \frac{27}{36}, \quad P(A_2|A_1) = \frac{26}{35}, \quad P(A_3|A_1 \cap A_2) = \frac{25}{34}.$$

$$P(A_1 \cap A_2 \cap A_3) \approx 0.41.$$

Example

There is a 36-card deck. 3 cards are drawn at random. What is the probability that there will get no spades?

Let A_i be the event that the card drawn at step i is not a spade. We are interested in the probability of the event $A_1 \cap A_2 \cap A_3$. By the multiplication rule,

$$P(A_1 \cap A_2 \cap A_3) = P(A_1)P(A_2|A_1)P(A_3|A_1 \cap A_2).$$

Clearly,

$$P(A_1) = \frac{27}{36}, \quad P(A_2|A_1) = \frac{26}{35}, \quad P(A_3|A_1 \cap A_2) = \frac{25}{34}.$$

$$P(A_1 \cap A_2 \cap A_3) \approx 0.41.$$

Total probability formula

Consider the partition $\mathcal{D} = \{A_1, \dots, A_m\}$. The equality

$$P(B) = \sum_{i=1}^m P(B|A_i)P(A_i)$$

is called the *total probability formula*.

Proof.

$$\sum_{i=1}^m P(B|A_i)P(A_i) = \sum_{i=1}^m P(B \cap A_i) = P\left(B \cap \bigcup_{i=1}^m A_i\right) = P(B).$$

Example

There are three urns, each of which can contain white and black balls:

- ▶ in the first urn there are 3 white and 2 black ones;
- ▶ in the second urn there are 0 white and 4 black;
- ▶ in the third urn there are 5 white and 2 black ones.

In the first step of the experiment, one of the urns is selected at random. In the second step, one ball is drawn from it at random. What is the probability to get a white ball?

Example (continued)

Let A_i be the event that the urn with the number i is selected, and B be the event that the white ball is drawn. The probabilities that B will occur under the condition A_i are easy to calculate since we know the number of white and black balls in the corresponding urn:

$$P(B|A_1) = 3/5, \quad P(B|A_2) = 0, \quad P(B|A_3) = 5/7.$$

By the total probability formula:

$$P(B) = \sum_{i=1}^3 P(B|A_i)P(A_i) = \frac{3}{5} \cdot \frac{1}{3} + \frac{5}{7} \cdot \frac{1}{3} = \frac{46}{105} \approx 0.44.$$

Example (Monty Hall)

There is a prize behind one of the three closed doors, there is nothing behind the other two.

- ▶ The player selects one of the doors.
- ▶ The host opens one of the two remaining doors, behind which there is nothing.
- ▶ The player can keep or change his original choice.

Which strategy is better?

Example (continued)

If the player retains his original choice, then the probability that he will receive the prize is $1/3$.

Denote by A the event that the prize is behind the door originally chosen by the player. Let B be the event that a prize is won. If the player changes his original choice, then

$$P(B|A) = 0, \quad P(B|A^c) = 1.$$

By the total probability formula,

$$P(B) = P(B|A)P(A) + P(B|A^c)P(A^c) = 0 \cdot \frac{1}{3} + 1 \cdot \frac{2}{3} = \frac{2}{3}.$$

Thus, the second strategy is much better.

The Bayes formula

Consider a partition $\mathcal{D} = \{A_1, \dots, A_m\}$. The equality

$$P(A_i|B) = \frac{P(B|A_i)P(A_i)}{\sum_{j=1}^m P(B|A_j)P(A_j)}$$

is called *the Bayes formula*.

Proof.

$$P(A_i|B) = \frac{P(A_i \cap B)}{P(B)} = \frac{P(B|A_i)P(A_i)}{\sum_{j=1}^m P(B|A_j)P(A_j)} \quad \square$$

Bayes' formula allows to take into account information concerning the occurrence of an event B to re-evaluate the probabilities of other events.

- ▶ $P(A_i)$ — prior probabilities (before an experiment)
- ▶ $P(A_i|B) = P_B(A_i)$ — posterior probabilities (after the experiment)

Example

A patient has a rare disease with probability 0.001. A medical test gives the correct answer in 95% of cases. What is the probability that the patient is ill if the test showed the positive result?

- ▶ H the hypothesis that the patient is ill, $\mathcal{D} = \{H, H^c\}$
- ▶ B the event, which is that the test result is positive.

The Bayes formula gives a surprising result:

$$\begin{aligned} P(H|B) &= \frac{P(B|H)P(H)}{P(B|H)P(H) + P(B|H^c)P(H^c)} \\ &= \frac{0.95 \cdot 0.001}{0.95 \cdot 0.001 + 0.05 \cdot 0.999} \approx 0.0187, \end{aligned}$$

that is, the patient is ill less than in 2% of cases. This means that a much more accurate test is required to determine such a rare disease or test should be repeated.

Example

A patient has a rare disease with probability 0.001. A medical test gives the correct answer in 95% of cases. What is the probability that the patient is ill if the test showed the positive result?

- ▶ H the hypothesis that the patient is ill, $\mathcal{D} = \{H, H^c\}$
- ▶ B the event, which is that the test result is positive.

The Bayes formula gives a surprising result:

$$\begin{aligned} P(H|B) &= \frac{P(B|H)P(H)}{P(B|H)P(H) + P(B|H^c)P(H^c)} \\ &= \frac{0.95 \cdot 0.001}{0.95 \cdot 0.001 + 0.05 \cdot 0.999} \approx 0.0187, \end{aligned}$$

that is, the patient is ill less than in 2% of cases. This means that a much more accurate test is required to determine such a rare disease or test should be repeated.

Independence

Events A , B are independent if

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

If $P(B) > 0$, then the independence is equivalent to the condition

$$P(A|B) = P(A).$$

Similarly, if $P(A) > 0$, then the independence is equivalent to the condition

$$P(B|A) = P(B).$$

Example

One card is taken at random from the deck of 36 cards. Are the events $A = \{\text{queen}\}$, $B = \{\text{spades}\}$ independent? Will the answer change for a 52-card deck?

► 36 cards

$$P(A) = 1/9, \quad P(B) = 1/4,$$

$$P(A \cap B) = P(\text{queen of spades}) = 1/36 = P(A)P(B).$$

► 52 cards

$$P(A) = 1/13, \quad P(B) = 1/4,$$

$$P(A \cap B) = P(\text{queen of spades}) = 1/52 \neq P(A)P(B).$$

Example

One card is taken at random from the deck of 36 cards. Are the events $A = \{\text{queen}\}$, $B = \{\text{spades}\}$ independent? Will the answer change for a 52-card deck?

► 36 cards

$$P(A) = 1/9, \quad P(B) = 1/4,$$

$$P(A \cap B) = P(\text{queen of spades}) = 1/36 = P(A)P(B).$$

► 52 cards

$$P(A) = 1/13, \quad P(B) = 1/4,$$

$$P(A \cap B) = P(\text{queen of spades}) = 1/52 \neq P(A)P(B).$$

Independence

σ -algebras $\mathcal{F}_1, \dots, \mathcal{F}_n$ are called independent, if

$$P(A_1 \cap \dots \cap A_n) = P(A_1) \dots P(A_n).$$

ξ is independent of \mathcal{H} , if \mathcal{F}_ξ and \mathcal{H} are independent

$$P(\{\xi \in B\} \cap A) = P(\xi \in B)P(A), \quad B \in \mathcal{B}(\mathbb{R}), A \in \mathcal{H}.$$

$(\xi_i)_{i=1}^n$ are independent, if σ -algebras \mathcal{F}_{ξ_i} are independent:

$$P(\xi_1 \in B_1, \dots, \xi_n \in B_n) = P(\xi_1 \in B_1) \dots P(\xi_n \in B_n)$$

In other words,

$$P_{\xi_1, \dots, \xi_n}(B_1, \dots, B_n) = P_{\xi_1}(B_1) \dots P_{\xi_n}(B_n).$$

For discrete ξ_i :

$$\begin{aligned} p_{\xi_1, \dots, \xi_n}(x_1, \dots, x_n) &= P_{\xi_1, \dots, \xi_n}(x_1, \dots, x_n) = P_{\xi_1}(x_1) \dots P_{\xi_n}(x_n) \\ &= p_{\xi_1}(x_1) \dots p_{\xi_n}(x_n). \end{aligned}$$

For continuous ξ_i :

$$\begin{aligned} \int_{B_1} \dots \int_{B_n} p_{\xi_1, \dots, \xi_n}(x_1, \dots, x_n) dx_1 \dots dx_n &= P_{\xi_1, \dots, \xi_n}(B_1, \dots, B_n) \\ &= P_{\xi_1}(B_1) \dots P_{\xi_n}(B_n) = \int_{B_1} p_{\xi_1}(x_1) dx_1 \dots \int_{B_n} p_{\xi_n}(x_n) dx_n \\ &= \int_{B_1} \dots \int_{B_n} p_{\xi_1}(x_1) \dots p_{\xi_n}(x_n) dx_1 \dots dx_n. \end{aligned}$$

Thus,

$$p_{\xi_1, \dots, \xi_n}(x_1, \dots, x_n) = p_{\xi_1}(x_1) \dots p_{\xi_n}(x_n).$$

Let ξ, η be independent, then $E(\xi\eta) = E(\xi)E(\eta)$.

In particular, if $\xi = I_A, \eta = I_B$, then

$$P(A \cap B) = EI_{A \cap B} = E(I_A I_B) = EI_A EI_B = P(A)P(B).$$

Furthermore,

$$\begin{aligned} \text{Var}(\xi + \eta) &= \text{Var}(\xi - E\xi + \eta - E\eta) = E(\xi - E\xi + \eta - E\eta)^2 \\ &= E(\xi - E\xi)^2 + 2E[(\xi - E\xi)(\eta - E\eta)] + E(\eta - E\eta)^2 \\ &= \text{Var}(\xi - E\xi) + \text{Var}(\eta - E\eta) = \text{Var}(\xi) + \text{Var}(\eta). \end{aligned}$$

Example: variance of the binomial distribution

$\xi \sim B(n, p)$: $p_\xi(k) = C_n^k p^k (1-p)^{n-k}$. Find $\text{Var}(\xi)$.

Sum of independent Bernoulli random variables

$$\eta = \sum_{i=1}^n \xi_i, \quad \xi_i = I_{\{a_i=H\}}, \quad \xi_i \sim \text{Ber}(p)$$

has the binomial distribution: $\eta \sim B(n, p)$. Hence,

$$\text{Var}(\xi) = \text{Var}(\eta) = \sum_{i=1}^n \text{Var}(\xi_i) = np(1-p).$$

Covariance

$$\begin{aligned}\text{Cov}(\xi, \eta) &:= \mathbb{E}[(\xi - \mathbb{E}\xi)(\eta - \mathbb{E}\eta)] = \int_{\mathbb{R}^2} (x - \mathbb{E}\xi)(y - \mathbb{E}\eta) P_{\xi}(dxdy) \\ &= \mathbb{E}(\xi\eta) - \mathbb{E}\xi \cdot \mathbb{E}\eta.\end{aligned}$$

The covariance is positive (negative) if the deviations of ξ and η from the mean values “more often ” have the same sign (different signs).

If ξ, η are independent, then

$$\text{Cov}(\xi, \eta) := \mathbb{E}[\xi - \mathbb{E}\xi] \cdot \mathbb{E}[\eta - \mathbb{E}\eta] = 0.$$

The converse is false.

Properties of the covariance

$$\text{Cov}(\xi, \xi) = \text{Var}(\xi),$$

$$\text{Cov}(\xi, \eta) = \text{Cov}(\eta, \xi),$$

$$\text{Cov}(\alpha\xi + \beta\eta, \zeta) = \alpha\text{Cov}(\xi, \zeta) + \beta\text{Cov}(\eta, \zeta).$$

$$\text{Var}(\xi + \eta) = \text{Cov}(\xi + \eta, \xi + \eta) = \text{Var}(\xi) + \text{Var}(\eta) + 2\text{Cov}(\xi, \eta)$$

$$\text{Var}\left(\sum_{i=1}^n \xi_i\right) = \sum_{i=1}^n \text{Var}(\xi_i) + \sum_{i \neq j} \text{Cov}(\xi_i, \xi_j)$$

Correlation coefficient

Cauchy-Schwarz inequality

$$|\mathbb{E}(\xi\eta)| \leq \sqrt{\mathbb{E}(\xi^2)}\sqrt{\mathbb{E}(\eta^2)}.$$

Put $\sigma(\xi) = \sqrt{\text{Var}(\xi)}$. Correlation coefficient

$$\rho(\xi, \eta) := \frac{\text{Cov}(\xi, \eta)}{\sigma(\xi)\sigma(\eta)}, \quad \sigma(\xi), \sigma(\eta) > 0.$$

is a dimensionless version of the covariance coefficient. By the Cauchy-Schwarz inequality,

$$|\text{Cov}(\xi, \eta)| \leq \sqrt{\mathbb{E}(\xi - \mathbb{E}\xi)^2} \sqrt{\mathbb{E}(\eta - \mathbb{E}\eta)^2} = \sigma(\xi)\sigma(\eta).$$

Hence, $\rho(\xi, \eta) \in [-1, 1]$.

Density of transformation of a multidimensional random variable

Let $g : \mathbb{R}^d \mapsto \mathbb{R}^d$ be a one-to-one differentiable mapping. Let $\eta = (\eta_1, \dots, \eta_d) = g(\xi_1, \dots, \xi_d) = g(\xi)$. Then

$$\begin{aligned} \mathbb{E}h(\eta) &= \int_{\mathbb{R}^d} h(y) f_{\eta}(y) dy = \mathbb{E}h(g(\xi)) = \int_{\mathbb{R}^d} h(g(x)) f_{\xi}(x) dx \\ &= \int_{\mathbb{R}^d} h(y) f_{\xi}(g^{-1}(y)) |\det [(g^{-1})']| dy, \end{aligned}$$

where $(g^{-1})'$ is the Jacobian of the mapping g^{-1} . Thus,

$$f_{\eta}(y) = f_{\xi}(g^{-1}(y)) |\det [(g^{-1})']|.$$

In particular, if $\eta = A\xi + b$, where A is an invertible $d \times d$ matrix and $b \in \mathbb{R}^d$, then

$$f_{\eta}(y) = f_{\xi}(A^{-1}(y - b)) |\det [A^{-1}]| = \frac{1}{|\det A|} f_{\xi}(A^{-1}(y - b)).$$

Multidimensional normal distribution

Let $\xi_i \sim N(0, 1)$, $i = 1, \dots, d$ be independent. *Multivariate standard normal (Gaussian)* distribution is the joint distribution of ξ_i :

$$\begin{aligned} f_{\xi_1, \dots, \xi_d}(x_1, \dots, x_d) &= \prod_{i=1}^d f_{\xi_i}(x_i) = \frac{1}{(\sqrt{2\pi})^d} \prod_{i=1}^d e^{-x_i^2/2} \\ &= \left(\frac{1}{\sqrt{2\pi}} \right)^d e^{-\langle x, x \rangle / 2} \end{aligned}$$

Multidimensional normal (Gaussian) distribution is the distribution of the vector $\eta = A\xi + \mu$ (the matrix A is assumed to be invertible for now). According to the density transformation formula,

$$\begin{aligned} f_{\eta}(y) &= \left(\frac{1}{\sqrt{2\pi}} \right)^d \frac{1}{|\det A|} e^{-\langle A^{-1}(y-\mu), A^{-1}(y-\mu) \rangle / 2} \\ &= \left(\frac{1}{\sqrt{2\pi}} \right)^d \frac{1}{|\det A|} e^{-\langle (A^{-1})^T A^{-1}(y-\mu), y-\mu \rangle / 2} \end{aligned}$$

$$(A^{-1})^T A^{-1} = (A^T)^{-1} A^{-1} = (AA^T)^{-1}.$$

Matrix AA^T is symmetric and positive definite:

$$\langle AA^T x, x \rangle = \|A^T x\|^2 > 0, \quad x \neq 0.$$

Put $\Sigma = AA^T$. Since

$$\det \Sigma = \det(AA^T) = (\det A)^2,$$

we have

$$f_{\eta}(y) = \left(\frac{1}{\sqrt{2\pi}} \right)^d \frac{1}{\sqrt{\det \Sigma}} e^{-\langle \Sigma^{-1}(y-\mu), y-\mu \rangle / 2}$$

For the distribution of a Gaussian vector η , the notation $\eta \sim N(\mu, \Sigma)$ is used. Consider the vector of expectations μ :

$$E\eta = E(\mu + A\xi) = \mu$$

and the *covariance matrix*

$$\begin{aligned} \text{Cov}(\eta_i, \eta_j) &= E(\langle a_i, \xi \rangle \cdot \langle a_j, \xi \rangle) = E \sum_{kr} a_{ik} a_{jr} \xi_k \xi_r \\ &= \sum_k a_{ik} a_{jk} = (AA^T)_{ij} = \Sigma_{ij}. \end{aligned}$$

These formulas indicate the meaning of the parameters μ , Σ .

If η_i are uncorrelated, then the matrix Σ is diagonal and

$$\begin{aligned} f_{\eta}(y) &= \left(\frac{1}{\sqrt{2\pi}} \right)^d \frac{1}{\sqrt{\prod_{i=1}^d \Sigma_{ii}}} \exp \left(- \sum_{i=1}^d \frac{(y_i - \mu_i)^2}{2\Sigma_{ii}} \right) \\ &= \prod_{i=1}^d \sqrt{\frac{1}{2\pi\Sigma_{ii}}} \exp \left(- \frac{(y_i - \mu_i)^2}{2\Sigma_{ii}} \right) \end{aligned}$$

Thus, $f_{\eta}(y)$ decomposes into the product of marginal densities of one-dimensional normal random variables $\eta_i \sim N(\mu_i, \Sigma_{ii})$. Hence the components of η are independent. Thus, for a Gaussian vector, the components are uncorrelated if and only if they are independent.

Characteristic functions

Characteristic function of a random variable $\xi : \Omega \mapsto \mathbb{R}^d$ is the Fourier transform of its distribution:

$$\varphi_{\xi}(u) = \mathbb{E} e^{i\langle u, \xi \rangle} = \int_{\mathbb{R}^d} e^{i\langle u, x \rangle} P_{\xi}(dx).$$

Here $\langle u, x \rangle$ is the usual scalar product in \mathbb{R}^d .

- ▶ The distribution P_{ξ} is uniquely determined by its characteristic function (uniqueness theorem).

Since the characteristic function is uniquely determined by the distribution, it can also be called the characteristic function of the distribution .

The characteristic function can be used to calculate the mean and variance, since (in the one-dimensional case)

$$\varphi'_{\xi}(u) = i\mathbb{E}(\xi e^{iu\xi}), \quad \varphi''_{\xi}(0) = -\mathbb{E}(\xi^2 e^{iu\xi}),$$

$$\varphi'_{\xi}(0) = i\mathbb{E}\xi, \quad \varphi''_{\xi}(0) = -\mathbb{E}\xi^2.$$

The characteristic function of the Poisson distribution

Let us find the expectation and variance of the Poisson distribution using the characteristic functions. Let $\xi \sim \Pi(\lambda)$. Then

$$\varphi_{\xi}(u) = \mathbb{E}e^{iu\xi} = \sum_{k=0}^{\infty} e^{iuk} \frac{\lambda^k}{k!} e^{-\lambda} = e^{-\lambda} \sum_{k=0}^{\infty} \frac{(\lambda e^{iu})^k}{k!} = \exp(-\lambda + \lambda e^{iu}),$$

$$\varphi'_{\xi}(u) = \lambda i e^{iu} \exp(\lambda(e^{iu} - 1)),$$

$$\varphi''_{\xi}(u) = (-\lambda e^{iu} + (\lambda i e^{iu})^2) \exp(\lambda(e^{iu} - 1)).$$

Hence,

$$\begin{aligned} \mathbb{E}\xi &= -i\varphi'_{\xi}(0) = \lambda, & \mathbb{E}\xi^2 &= -\varphi''_{\xi}(0) = \lambda + \lambda^2, \\ \text{Var}(\xi) &= \mathbb{E}(\xi^2) - (\mathbb{E}\xi)^2 = \lambda. \end{aligned}$$

The characteristic function of the normal distribution

Let $\eta \sim N(0, 1)$. Then

$$\varphi_{\eta}(u) = \mathbb{E}e^{iu\eta} = \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{iux} e^{-x^2/2} dx,$$

$$\begin{aligned}\varphi'_{\eta}(u) &= i\mathbb{E}(\eta e^{iu\eta}) = \frac{i}{\sqrt{2\pi}} \int_{\mathbb{R}} x e^{iux} e^{-x^2/2} dx = -\frac{i}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{iux} d e^{-x^2/2} \\ &= -\frac{u}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{iux} e^{-x^2/2} dx = -u\varphi_{\eta}(u).\end{aligned}$$

Solving the differential equation $\varphi'_{\eta}(u) = -u\varphi_{\eta}(u)$ with the initial condition $\varphi_{\eta}(0) = 1$, we get

$$\varphi_{\eta}(u) = e^{-u^2/2}.$$

To find the characteristic function of an arbitrary normal distribution, consider a random variable $\xi = \mu + \sigma\eta \sim N(\mu, \sigma^2)$:

$$\varphi_{\xi}(u) = e^{iu\mu} \mathbf{E} e^{iu\sigma\eta} = e^{i\mu u - \sigma^2 u^2/2}.$$

$$\varphi'_{\xi}(u) = (i\mu - \sigma^2 u) e^{i\mu u - \sigma^2 u^2/2}, \quad \varphi''_{\xi}(u) = [(i\mu - \sigma^2 u)^2 - \sigma^2] e^{i\mu u - \sigma^2 u^2/2}.$$

Hence,

$$\begin{aligned} \mathbf{E}\xi &= -i\varphi'_{\xi}(0) = \mu, & \mathbf{E}\xi^2 &= -\varphi''_{\xi}(0) = \mu^2 + \sigma^2, \\ \text{Var}(\xi) &= \mathbf{E}\xi^2 - (\mathbf{E}\xi)^2 = \sigma^2. \end{aligned}$$

Characteristic function of the sum of independent random variables

The characteristic function of the joint distribution of a random vector is equal to the product of the characteristic functions of the components:

$$\varphi_{\xi}(u) = \mathbb{E}e^{i\langle u, \xi \rangle} = \mathbb{E} \prod_{k=1}^n e^{iu_k \xi_k} = \prod_{k=1}^n \mathbb{E}e^{iu_k \xi_k} = \prod_{k=1}^n \varphi_{\xi_k}(u_k).$$

From the mentioned uniqueness theorem it follows that the converse statement is also true, i.e. this equality is a criterion of independence of (ξ_1, \dots, ξ_n) .

Characteristic function of a multivariate normal distribution

$$\eta = \mu + A\xi:$$

$$\begin{aligned}\varphi_{\eta}(u) &= e^{i\langle u, \mu \rangle} \mathbb{E} e^{i\langle u, A\xi \rangle} = e^{i\langle u, \mu \rangle} \mathbb{E} e^{i\langle A^T u, \xi \rangle} = e^{i\langle u, \mu \rangle} \prod_{k=1}^d \mathbb{E} e^{i(A^T u)_k \xi_k} \\ &= e^{i\langle u, \mu \rangle} \prod_{i=1}^d e^{-(A^T u)_i^2 / 2} = e^{i\langle u, \mu \rangle} e^{-\|A^T u\|^2 / 2} = e^{i\langle u, \mu \rangle - \langle \Sigma u, u \rangle / 2},\end{aligned}$$

where $\Sigma = AA^T$. Thus, the vector η is *Gaussian*, if and only if

$$\varphi_{\eta}(u) = e^{i\langle u, \mu \rangle - \langle \Sigma u, u \rangle / 2},$$

where $\mu \in \mathbb{R}^d$ and Σ is a symmetric nonnegative definite matrix. This property can be taken as a definition of a multidimensional Gaussian random variable. Moreover, it suffices to assume that Σ is non-negative definite: its invertibility is not required.

Distribution of a sum of one-dimensional independent normal r.v.

$\xi_i \sim N(\mu_i, \sigma_i^2)$ are independent, $\zeta = \xi_1 + \xi_2$.

$$\begin{aligned}\varphi_\zeta(u) &= \varphi_{\xi_1}(u)\varphi_{\xi_2}(u) = e^{iu\mu_1 - \sigma_1^2 u^2/2} e^{iu\mu_2 - \sigma_2^2 u^2/2} \\ &= e^{iu(\mu_1 + \mu_2) - (\sigma_1^2 + \sigma_2^2)u^2/2}.\end{aligned}$$

The uniqueness theorem implies that $\zeta \sim N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$. It is much more difficult to obtain this result by calculating convolution. It was formulated above, but calculations were not performed.

Convergence of random variables

A sequence of random variables ξ_n converges to ξ

- ▶ *in probability*: $\xi_n \xrightarrow{P} \xi$, if

$$\lim_{n \rightarrow \infty} P\{\omega : |\xi_n(\omega) - \xi(\omega)| \geq \varepsilon\} = 0;$$

- ▶ *in distribution*: $\xi_n \xrightarrow{d} \xi$, if $\lim_{n \rightarrow \infty} Eg(\xi_n) = Eg(\xi)$ for any bounded continuous function $g : \mathbb{R}^d \mapsto \mathbb{R}$.
- ▶ convergence in probability \implies convergence in distribution;
- ▶ *Continuity theorem*: if $\varphi_{\xi_n}(u) \rightarrow \varphi(u)$, $n \rightarrow \infty$ for all u , then $\xi_n \xrightarrow{d} \xi$.

Weak law of large numbers

Let $(\xi_i)_{i=1}^{\infty}$ be i.i.d. random variables, $E\xi_i = \mu$, $\text{Var}(\xi_i) = \sigma^2$.

► *Weak law of large numbers:*

$$M_n = \frac{1}{n} \sum_{i=1}^n \xi_i \xrightarrow{P} \mu, \quad n \rightarrow \infty.$$

Proof.

$$EM_n = \mu, \quad \text{Var}(M_n) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(\xi_i) = \frac{\sigma^2}{n}.$$

$$P(|M_n - \mu| \geq \varepsilon) \leq \frac{\text{Var}(M_n)}{\varepsilon^2} = \frac{\sigma^2}{n\varepsilon^2} \rightarrow 0, \quad n \rightarrow \infty. \quad \square$$

In particular, the frequency of occurrence of an event in a series of independent experiments approaches its probability: A_i are independent, $P(A_i) = p$,

$$\nu_n = \frac{1}{n} \sum_{i=1}^n I_{A_i} \xrightarrow{P} p$$

Central limit theorem

Let $(\xi_k)_{k=0}^{\infty}$ be i.i.d. random variables, $E\xi_n = 0$, $\text{Var}(\xi_n) = 1$.

► *Central limit theorem:*

$$\frac{S_n}{\sqrt{n}} := \frac{1}{\sqrt{n}}(\xi_1 + \cdots + \xi_n) \xrightarrow{d} \eta \sim N(0, 1).$$

Proof. Let $\xi_j \sim \xi$. Then

$$\begin{aligned}\varphi_{S_n/\sqrt{n}}(u) &= E \prod_{j=1}^n e^{iu\xi_j/\sqrt{n}} = \prod_{j=1}^n E e^{iu\xi_j/\sqrt{n}} = (E e^{iu\xi/\sqrt{n}})^n \\&= \left(E \left(1 + \frac{iu\xi}{\sqrt{n}} - \frac{1}{2} \frac{u^2 \xi^2}{n} + o\left(\frac{1}{n}\right) \right) \right)^n = \left(1 - \frac{1}{2} \frac{u^2}{n} + o\left(\frac{1}{n}\right) \right)^n \\&= \exp \left(n \ln \left(1 - \frac{u^2}{2n} + o\left(\frac{1}{n}\right) \right) \right) \rightarrow e^{-u^2/2}, \quad n \rightarrow \infty \quad \square\end{aligned}$$

► $\eta_n \xrightarrow{d} \eta$, if and only if

$$F_{\eta_n}(x) \rightarrow F_{\eta}(x)$$

at all points where F_{η} is continuous.

In particular, in the central limit theorem

$$\mathbf{P}\left(\frac{\xi_1 + \cdots + \xi_n}{\sqrt{n}} \leq u\right) \rightarrow \Phi(x) := \mathbf{P}(\zeta \leq x),$$

$$\zeta \in N(0, 1), \quad \mathbf{E}\xi_i = 0, \quad \text{Var}(\xi_i) = 1.$$

Monte-Carlo integration

Let $X_i \sim U(0, 1)$. Then

$$\frac{1}{n} \sum_{i=1}^n g(X_i) \xrightarrow{P} \mathbb{E}g(X_1) = \int_0^1 g(x) dx,$$

$$\sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n g(X_i) - \int_0^1 g(x) dx \right) \xrightarrow{d} N(0, \text{Var}(g(X_1))).$$

If X_i has a p.d.f. f on \mathbb{R} , then

$$\frac{1}{n} \sum_{i=1}^n g(X_i) \xrightarrow{P} \mathbb{E}g(X_1) = \int_{\mathbb{R}} g(x)f(x) dx.$$

Thus, if a computer can generate a realization of a sample with a given distribution, then we get a very simple method for the approximate calculation of integrals.

It is actually enough to have a generator for the uniform distribution: $\zeta \sim U(0, 1)$. Let $F_\xi(x) = P(\xi \leq x)$ be the c.d.f. of ξ . Put

$$F_\xi^{-1}(u) = \min\{x : F_\xi(x) \geq u\}, \quad u \in (0, 1)$$

(the minimum is attained, since F_ξ is right-continuous).

► Inverse-transform method:

$$F_\xi^{-1}(\zeta) \sim P_\xi.$$

Proof for a continuous and strictly increasing F_ξ . The equation

$$F_\xi(x) = u$$

has a unique solution for any $u \in (0, 1)$. In this case F_ξ^{-1} coincides with the ordinary inverse function and

$$P(F_\xi^{-1}(\zeta) \leq x) = P(\zeta \leq F_\xi(x)) = F_\xi(x).$$

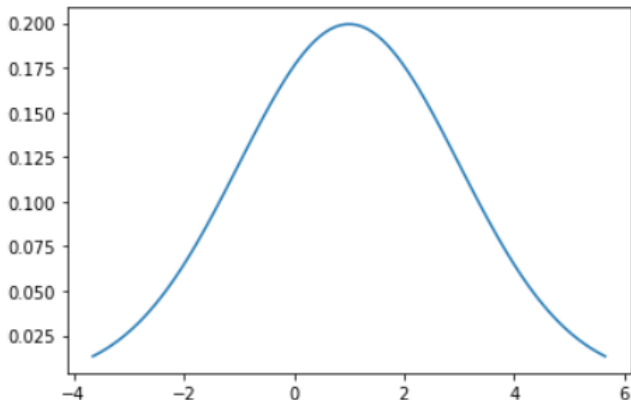
```
In [10]: import numpy as np
import matplotlib.pyplot as plt
from scipy import stats
rv = stats.norm(1,2)
```

```
In [11]: print('mean:', rv.mean())
print('variance:', rv.var())
print('standard deviation:', rv.std())
print('median:', rv.median())
print('0.01 percentile:', rv.ppf(0.01))
print('0.99 percentile:', rv.ppf(0.99))
```

```
mean: 1.0
variance: 4.0
standard deviation: 2.0
median: 1.0
0.01 percentile: -3.6526957480816815
0.99 percentile: 5.6526957480816815
```

```
In [12]: x=np.linspace(rv.ppf(0.01),rv.ppf(0.99),100)
plt.plot(x,rv.pdf(x))
```

```
Out[12]: [<matplotlib.lines.Line2D at 0x1d0295ca128>]
```



```
In [14]: X_samples=rv.rvs(1000)
print(X_samples[:4])
```

```
[-1.26321354 -1.12273014  0.94553942  1.94721138]
```

```
In [6]: print('sample mean:', X_samples.mean())
print('sample variance:', X_samples.var())
print('sample standard deviation:', X_samples.std())
print('sample median:', np.median(X_samples))
print('0.01 percentile:', np.percentile(X_samples,1.))
print('0.99 percentile:', np.percentile(X_samples,99.))
```

```
sample mean: 1.0234708179992147
```

```
sample variance: 4.20289766675017
```

```
sample standard deviation: 2.0500969895958994
```

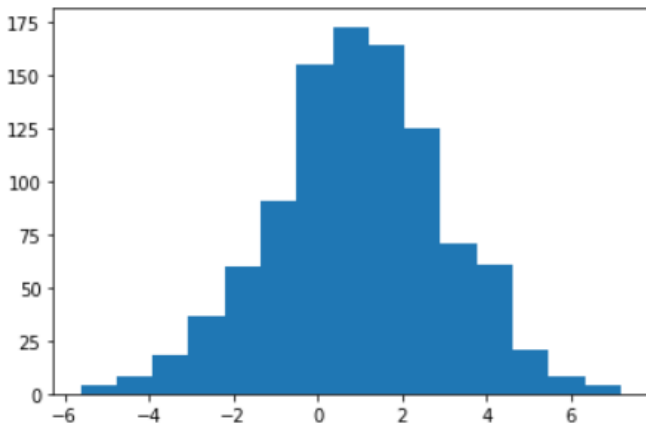
```
sample median: 1.0317126998060728
```

```
0.01 percentile: -3.7197014673708044
```

```
0.99 percentile: 5.67401027625593
```



```
In [15]: plt.hist(X_samples, bins=15)  
plt.show()
```



Approximate calculation of the integral $\int_0^1 x^2 dx = 1/3$

by the Monte-Carlo method:

```
In [17]: X_samples=stats.uniform.rvs(size=1000)
         np.mean(X_samples**2)
```

```
Out[17]: 0.34342367972776844
```

Approximate calculation of the integral $\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x^2 e^{-x^2} dx = 1$

by the Monte-Carlo method:

```
In [18]: X_samples=stats.norm.rvs(size=1000)
         np.mean(X_samples**2)
```

```
Out[18]: 0.9752834638543979
```

Parameter estimation in classical statistics

Data x_1, \dots, x_n are assumed to be the realizations of i.i.d. random variables $X_1, \dots, X_n \sim \mathcal{P}$, where \mathcal{P} is an unknown probability measure from a known set \mathcal{P} . The aim is to identify \mathcal{P} using the data.

Parametric model: $\mathcal{P} = \{\mathcal{P}_\theta : \theta \in \Theta\}$. Usually one considers families of p.m.f or p.d.f.:

$$\{x \mapsto f(x; \theta) : \theta \in \Theta\},$$

which we will denote by the same letter \mathcal{P} . For example, the family of normal p.d.f.:

$$f(x; \theta) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/(2\sigma^2)}, \quad \theta = (\mu, \sigma) \in \mathbb{R} \times (0, \infty).$$

We want to construct a statistics $\hat{\theta}$ (any function of data) such that

$$\theta \approx \hat{\theta} = \hat{\theta}(X_1, \dots, X_n).$$

Maximum likelihood method

Likelihood function

$$L_n(x; \theta) = \prod_{i=1}^n p(x_i; \theta).$$

For a p.m.f. family, $L_n(x; \theta) = P_\theta(X_1 = x_1, \dots, X_n = x_n)$.

Logarithmic likelihood function:

$$l_n(x; \theta) = \ln L_n(x; \theta).$$

MLE (maximum likelihood estimate):

$$\hat{\theta} \in \arg \max_{\theta} L_n(x; \theta) = \arg \max_{\theta} l_n(x; \theta).$$

In connection with this optimization problem, we can assume that the likelihood function is determined up to a positive factor $C(x)$, independent of θ :

$$L_n(x; \theta) \propto C(x) L_n(x; \theta).$$

Example. $X_i \sim \text{Ber}(\theta)$,

$$p_{\theta}(x) = \begin{cases} \theta, & x = 1, \\ 1 - \theta, & x = 0. \end{cases}$$

$$L_n(x; \theta) = \prod_{i=1}^n p_{\theta}(x_i) = \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{1 - x_i} = \theta^S (1 - \theta)^{n - S}, \quad S = \sum_{i=1}^n x_i,$$

$$l_n(x; \theta) = S \ln \theta + (n - S) \ln(1 - \theta),$$

$$\frac{\partial l_n(x; \theta)}{\partial \theta} = \frac{S}{\theta} - \frac{n - S}{1 - \theta} = 0, \quad S = n\theta,$$

$$\hat{\theta} = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X}.$$

Example. $X \sim N(\mu, 1)$, $\theta = \mu$,

$$f_{\theta}(y) = \frac{1}{\sqrt{2\pi}} e^{-(y-\mu)^2/2},$$

$$L_n(x; \mu) \propto \prod_{i=1}^n e^{-(x_i - \mu)^2/2} \propto e^{-n(\bar{x} - \mu)^2/2}, \quad \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i.$$

Indeed,

$$\sum_{i=1}^n (x_i - \mu)^2 = \sum_{i=1}^n x_i^2 - 2n\mu\bar{x} + n\mu^2 = \sum_{i=1}^n x_i^2 - n\bar{x}^2 + n(\bar{x} - \mu)^2.$$

Hence (neglecting terms independent of μ),

$$l_n(x; \mu) = -n(\bar{x} - \mu)^2/2,$$

Thus,

$$\hat{\mu} = \bar{X}.$$

Example. $X_i \sim N(\mu, \sigma^2)$, $\theta = (\mu, \sigma)$,

$$L_n(x; \theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-(x_i - \mu)^2 / (2\sigma^2)} \propto \frac{1}{\sigma^n} \prod_{i=1}^n e^{-(x_i - \mu)^2 / (2\sigma^2)},$$

$$l_n(x; \theta) = \ln L_n(x; \theta) = -n \ln \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

$$(l_n)_\mu = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = \frac{1}{\sigma^2} \left(\sum_{i=1}^n x_i - n\mu \right) = 0,$$

$$(l_n)_\sigma = -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - \mu)^2 = 0.$$

It follows that

$$\hat{\mu} = \bar{X}, \quad \hat{\sigma} = \left(\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \right)^{1/2}.$$

By the invariance principle for MLE for σ^2 coincides with $\hat{\sigma}^2$. The same estimates were obtained by the method of moments.

Example. $X_i \sim U(0, \theta)$,

$$f(y; \theta) = \begin{cases} 1/\theta, & y \in (0, \theta), \\ 0, & \text{otherwise,} \end{cases}$$

$$L_n(x; \theta) = \prod_{i=1}^n f(x_i; \theta) = \begin{cases} \frac{1}{\theta^n}, & x_i < \theta, i = 1, \dots, n, \\ 0, & \text{otherwise.} \end{cases}.$$

The MLE is the smallest θ , which is larger than all x_i :

$$\hat{\theta} = X_{(n)} := \max_{1 \leq i \leq n} X_i.$$

The method of moments gave another estimate:

$$\hat{\theta} = 2\bar{X}.$$

Example. $X_i \sim \Pi(\lambda)$,

$$f(x; \lambda) = e^{-\lambda} \frac{\lambda^x}{x!}, \quad x \in \mathbb{Z}_+,$$

$$L_n(x; \lambda) = \prod_{i=1}^n e^{-\lambda} \frac{\lambda^{x_i}}{x_i!} = e^{-n\lambda} \lambda^{\sum_{i=1}^n x_i} \prod_{i=1}^n \frac{1}{x_i!},$$

$$l_n(x; \lambda) = -n\lambda + \sum_{i=1}^n x_i \ln \lambda + \text{Const},$$

$$\frac{\partial l_n}{\partial \lambda} = -n + \frac{1}{\lambda} \sum_{i=1}^n x_i = 0.$$

MLE:

$$\hat{\lambda} = \frac{1}{n} \sum_{i=1}^n X_i = \overline{X}.$$

Conditional expectation

Consider a σ -algebra $\mathcal{H} \subset \mathcal{F}$. Let $\xi \in L^1(\mathcal{F}) := L^1(\Omega, \mathcal{F}, \mathbf{P}, \mathbb{R})$.

- ▶ A random variable $\eta \in L^1(\mathcal{H})$ is called a *conditional expectation* of ξ given \mathcal{H} , if

$$\mathbf{E}(\xi I_A) = \mathbf{E}(\eta I_A), \quad A \in \mathcal{H}.$$

Notation: $\eta = \mathbf{E}(\xi | \mathcal{H})$.

- ▶ If $\eta' = \eta$ a.s., $\mathbf{P}(\eta = \eta') = 1$, then

$$\mathbf{E}(\eta I_A) = \mathbf{E}(\eta' I_A).$$

Hence, the conditional expectation can be unique only up to a set of measure zero.

Theorem

For any $\xi \in L^1(\mathcal{F})$ the conditional expectation

$$\eta = E(\xi|\mathcal{H}) \in L^1(\mathcal{H})$$

exists and is unique (as an equivalence class).

Evident properties:

- ▶ Let $\xi = C$, then $E(\xi|\mathcal{H}) = \xi$ (by the definition).
- ▶ Let $\xi \in L^1(\mathcal{H})$, then $E(\xi|\mathcal{H}) = \xi$ (by the definition).
- ▶ Let $\mathcal{H} = \{\emptyset, \Omega\}$. Then $E(\xi|\mathcal{H}) = E\xi$.

Example

Consider a partition $(A_i)_{i=1}^m$ with $P(A_i) > 0$, $i = 1, \dots, m$. Assume that \mathcal{H} is generated by this partition. Then

$$E(\xi|\mathcal{H}) = \sum_{i=1}^m \frac{E(\xi I_{A_i})}{P(A_i)} I_{A_i}, \quad \xi \in L^1(\mathcal{F}).$$

Indeed, denote the left-hand side by η . For any atom A_j we have

$$E(\eta I_{A_j}) = E\left(\sum_{i=1}^m \frac{E(\xi I_{A_i})}{P(A_i)} I_{A_i} I_{A_j}\right) = E\left(\frac{E(\xi I_{A_j})}{P(A_j)} I_{A_j}\right) = E(\xi I_{A_j}).$$

But any $A \in \mathcal{H}$ is the union of some A_j .

Example

Coin tossing model.

- ▶ $\Omega = \{(a_1, \dots, a_n) : a_i \in \{H, T\}\},$
- ▶ $\mathcal{F}_t = \sigma(\{\omega : (a_1, \dots, a_t) \text{ fixed}\}),$
- ▶ $P(\omega) = p^{\nu(\omega)} q^{n-\nu(\omega)}, \nu(\omega) - \text{the number of } H \text{ in the sequence } \omega = (a_1, \dots, a_n), p \in (0, 1), q = 1 - p.$
- ▶ $A_H := H_1 = \{\omega : a_1 = H\},$
- ▶ $A_{HH} := H_1 \cap H_2 = \{\omega : a_1 = a_2 = H\},$
- ▶ $A_{HHH} := H_1 \cap H_2 \cap H_3 = \{\omega : a_1 = a_2 = a_3 = H\}.$

$$\begin{aligned} E I_{A_{HHH}} &= P(A_{HHH}) = p^3, \\ E(I_{A_{HHH}} | \mathcal{F}_1) &= \frac{E(I_{A_{HHH}} I_{A_H})}{P(A_H)} I_{A_H} + \frac{E(I_{A_{HHH}} I_{A_T})}{P(A_T)} I_{A_T} \\ &= \frac{P(A_{HHH})}{P(A_H)} I_{A_H} = \frac{p^3}{p} I_{A_H} = p^2 I_{A_H}. \end{aligned}$$

Example (continued)

$$\begin{aligned} \mathbb{E}(I_{A_{HHH}}|\mathcal{F}_2) &= \frac{\mathbb{E}(I_{A_{HHH}}I_{A_{HH}})}{\mathbb{P}(A_{HH})}I_{A_{HH}} + \frac{\mathbb{E}(I_{A_{HHH}}I_{A_{HT}})}{\mathbb{P}(A_{HT})}I_{A_{HT}} \\ &\quad + \frac{\mathbb{E}(I_{A_{HHH}}I_{A_{TH}})}{\mathbb{P}(A_{TH})}I_{A_{TH}} + \frac{\mathbb{E}(I_{A_{HHH}}I_{A_{TT}})}{\mathbb{P}(A_{TT})}I_{A_{TT}} \\ &= \frac{\mathbb{E}(I_{A_{HHH} \cap A_{HH}})}{\mathbb{P}(A_{HH})}I_{A_{HH}} = \frac{\mathbb{P}(A_{HHH})}{\mathbb{P}(A_{HH})}I_{A_{HH}} = \frac{p^3}{p^2}I_{A_{HH}} \\ &= pI_{A_{HH}}, \\ \mathbb{E}(I_{A_{HHH}}|\mathcal{F}_3) &= I_{A_{HHH}}. \end{aligned}$$

Example

On the same probability space consider the stochastic process

$$S_{t+1} = S_t(uI_{\{a_{t+1}=H\}} + dI_{\{a_{t+1}=T\}}),$$

where S_0 is a constant and $u > d$. This is the simplest model of a risky asset (binomial model, Cox-Ross-Rubinstein model). Let us find $E(S_2|\mathcal{F}_1)$:

$$E(S_2|\mathcal{F}_1) = \frac{E(S_2 I_{A_H})}{P(A_H)} I_{A_H} + \frac{E(S_2 I_{A_T})}{P(A_T)} I_{A_T}.$$

Since $S_2 = S_0(u^2 I_{A_{HH}} + ud I_{A_{HT}} + du I_{A_{TH}} + d^2 I_{A_{TT}})$, we get

$$E(S_2 I_{A_H}) = S_0 E(u^2 I_{A_{HH}} + ud I_{A_{HT}}) = S_0(u^2 p^2 + udpq),$$

$$E(S_2 I_{A_T}) = S_0 E(du I_{A_{TH}} + d^2 I_{A_{TT}}) = S_0(udpq + d^2 q^2).$$

Example (continued)

Thus,

$$\begin{aligned} E(S_2|\mathcal{F}_1) &= \frac{S_0(u^2p^2 + udpq)}{p}I_{A_H} + \frac{S_0(udpq + d^2q^2)}{q}I_{A_T} \\ &= S_0u(up + dq)I_{A_H} + S_0d(up + dq)I_{A_T} \\ &= (up + dq)S_1 \end{aligned}$$

Properties of conditional expectation

- (1) Conditional expectation is a linear operator from $L^1(\mathcal{F})$ to $L^1(\mathcal{H})$:

$$\mathbb{E}(a_1\xi_1 + a_2\xi_2|\mathcal{H}) = a_1\mathbb{E}(\xi_1|\mathcal{H}) + a_2\mathbb{E}(\xi_2|\mathcal{H}).$$

- (2) Let $\xi \in L^1(\mathcal{F})$ and let η be \mathcal{H} -measurable and bounded. Then

$$\mathbb{E}(\xi\eta|\mathcal{H}) = \eta\mathbb{E}(\xi|\mathcal{H}).$$

Proof for indicators. Let $\eta = I_B$, $B \in \mathcal{H}$. Then the left and right sides are equal to

$$\mathbb{E}(\xi\eta I_A) = \mathbb{E}(\xi I_B I_A) = \mathbb{E}(\xi I_{A \cap B})$$

$$\mathbb{E}(\eta \mathbb{E}(\xi|\mathcal{H}) I_A) = \mathbb{E}(\mathbb{E}(\xi|\mathcal{H}) I_{A \cap B}) = \mathbb{E}(\xi I_{A \cap B})$$

respectively.

Properties of conditional expectation

(3) If $\xi \in L^1(\mathcal{F})$ is independent from \mathcal{H} , then

$$E(\xi|\mathcal{H}) = E\xi.$$

Proof.

$$E(\xi I_A) = E\xi \cdot EI_A = E(I_A E\xi), \quad A \in \mathcal{H}.$$

(4) Let $\xi \in L^1(\mathcal{F})$. Then

$$EE(\xi|\mathcal{H}) = E\xi$$

(law of total expectation).

Prof. Put $A = \Omega$ in the definition of conditional expectation:

$$E\xi = E(\xi I_\Omega) = E(E(\xi|\mathcal{H})I_\Omega) = EE(\xi|\mathcal{H}).$$

Properties of conditional expectation

(5) Let $\xi \in L^1(\mathcal{F})$ and $\mathcal{H} \subset \mathcal{G} \subset \mathcal{F}$. Then

$$E(E(\xi|\mathcal{G})|\mathcal{H}) = E(\xi|\mathcal{H})$$

Telescopic property, (tower property, law of iterated expectations): generalizes (4).

Proof. We need to prove that

$$E(\xi I_A) = E(E(E(\xi|\mathcal{G})|\mathcal{H})I_A), \quad A \in \mathcal{H}.$$

Use (2) and (4) in right-hand side:

$$\begin{aligned} E(E(E(\xi|\mathcal{G})|\mathcal{H})I_A) &\stackrel{(2)}{=} E(E(E(\xi|\mathcal{G})I_A|\mathcal{H})) \stackrel{(4)}{=} E(E(\xi|\mathcal{G})I_A) \\ &\stackrel{(2)}{=} E(E(\xi I_A|\mathcal{G})) \stackrel{(4)}{=} E(\xi I_A). \end{aligned}$$

Properties of conditional expectation

- (6) Let $\xi \in L^2(\mathcal{F})$. Then $\eta = E(\xi|\mathcal{H}) \in L^2(\mathcal{H})$ is the best L^2 -approximation of ξ on the basis of information contained in \mathcal{H} :

$$E(\xi - \eta)^2 = \min_{\zeta \in L^2(\mathcal{H})} E(\xi - \zeta)^2.$$

In other words, $\xi \mapsto E(\xi|\mathcal{H})$ is the orthogonal projection on $L^2(\mathcal{H})$.

Proof.

$$\begin{aligned} E((\xi - \zeta)^2|\mathcal{H}) &= E(\xi^2|\mathcal{H}) - 2\zeta E(\xi|\mathcal{H}) + \zeta^2 \\ &= E(\xi^2|\mathcal{H}) + (E(\xi|\mathcal{H}) - \zeta)^2 - (E(\xi|\mathcal{H}))^2, \\ E((\xi - \zeta)^2|\mathcal{H}) &= E(E(\xi|\mathcal{H}) - \zeta)^2 + E\xi^2 - E(E(\xi|\mathcal{H}))^2. \end{aligned}$$

The minimal value is attained at $\zeta^* = E(\xi|\mathcal{H})$.

Example

Let us find $E(S_{t+1}|\mathcal{F}_t)$, ES_{t+1} in the binomial model.

In the formula

$$S_{t+1} = S_t(uI_{\{a_{t+1}=H\}} + dI_{\{a_{t+1}=T\}}),$$

the first multiplier is \mathcal{F}_t -measurable, and the second one does not depend on \mathcal{F}_t . Hence,

$$\begin{aligned} E(S_{t+1}|\mathcal{F}_t) &\stackrel{(2)}{=} S_t E(uI_{\{a_{t+1}=H\}} + dI_{\{a_{t+1}=T\}}|\mathcal{F}_t) \\ &\stackrel{(3)}{=} S_t E(uI_{\{a_{t+1}=H\}} + dI_{\{a_{t+1}=T\}}) \\ &= S_t(uEI_{\{a_{t+1}=H\}} + dEI_{\{a_{t+1}=T\}}) = (up + dq)S_t. \\ ES_{t+1} &\stackrel{(4)}{=} EE(S_{t+1}|\mathcal{F}_t) = (up + dq)ES_t = (up + dq)^2 ES_{t-1} \\ &= (up + dq)^{t+1} S_0. \end{aligned}$$

Conditional expectation with respect to a random variable

Let $\xi \in L^1(\mathcal{F})$. Be the definition,

$$E(\xi|\eta) = E(\xi|\mathcal{F}_\eta),$$

where $\mathcal{F}_\eta := \sigma(\eta) = \{\eta^{-1}(B), B \in \mathcal{B}(\mathbb{R})\}$ is the σ -algebra, generated by η .

Example

Let $\eta : \Omega \mapsto \{z_1, \dots, z_m\}$. Then

$$\mathcal{F}_\eta = \sigma(A_1, \dots, A_m), \quad A_i = \eta^{-1}(\{z_i\}).$$

$$E(\xi|\eta) = \sum_{i=1}^m \frac{E(\xi I_{A_i})}{P(A_i)} I_{A_i} = \sum_{i=1}^m \frac{E(\xi I_{\{\eta=z_i\}})}{P(\eta = z_i)} I_{\{\eta=z_i\}}.$$

- ▶ For a random variable η with finite number of values we obtained the formula

$$E(\xi|\eta) = \sum_{i=1}^m c_i I_{\{\eta=z_i\}} = g(\eta).$$

- ▶ It is known that if ζ is $\sigma(\eta)$ -measurable, then $\zeta = g(\eta)$ for some Borel function g (*Doob-Dynkin lemma*).
- ▶ It follows that for any r.v. η there exists a Borel function g such that $E(\xi|\eta) = g(\eta)$. The following notation is used:

$$E(\xi|\eta = x) = g(x).$$

So we obtain the possibility to calculate the conditional expectation with respect to events with zero probabilities.

- ▶ g is not unique. We could define it in an arbitrary way outside the set $\{z_1, \dots, z_m\}$.

Example

Let ξ_1, \dots, ξ_n be i.i.d integrable random variables. Put $\xi = \xi_1 + \dots + \xi_n$ and find $E(\xi_1|\xi)$, $E(\xi|\xi_1)$. From the symmetry,

$$E(\xi_1|\xi) = \dots = E(\xi_n|\xi).$$

From the definition it follows that

$$E(\xi_1|\xi) = \frac{1}{n} \sum_{i=1}^n E(\xi_i|\xi) = \frac{1}{n} E(\xi_1 + \dots + \xi_n|\xi) = \frac{1}{n} E(\xi|\xi) = \frac{\xi}{n}.$$

Furthermore,

$$\begin{aligned} E(\xi|\xi_1) &= E(\xi_2 + \dots + \xi_n|\xi_1) + E(\xi_1|\xi_1) = E(\xi_2 + \dots + \xi_n) + \xi_1 \\ &= (n-1)E\xi_1 + \xi_1. \end{aligned}$$

Example

Let ξ_1, \dots, ξ_n be i.i.d integrable random variables. Put $\xi = \xi_1 + \dots + \xi_n$ and find $E(\xi_1|\xi)$, $E(\xi|\xi_1)$. From the symmetry,

$$E(\xi_1|\xi) = \dots = E(\xi_n|\xi).$$

From the definition it follows that

$$E(\xi_1|\xi) = \frac{1}{n} \sum_{i=1}^n E(\xi_i|\xi) = \frac{1}{n} E(\xi_1 + \dots + \xi_n|\xi) = \frac{1}{n} E(\xi|\xi) = \frac{\xi}{n}.$$

Furthermore,

$$\begin{aligned} E(\xi|\xi_1) &= E(\xi_2 + \dots + \xi_n|\xi_1) + E(\xi_1|\xi_1) = E(\xi_2 + \dots + \xi_n) + \xi_1 \\ &= (n-1)E\xi_1 + \xi_1. \end{aligned}$$

Example

Let ξ_1, \dots, ξ_n be i.i.d integrable random variables. Put $\xi = \xi_1 + \dots + \xi_n$ and find $E(\xi_1|\xi)$, $E(\xi|\xi_1)$. From the symmetry,

$$E(\xi_1|\xi) = \dots = E(\xi_n|\xi).$$

From the definition it follows that

$$E(\xi_1|\xi) = \frac{1}{n} \sum_{i=1}^n E(\xi_i|\xi) = \frac{1}{n} E(\xi_1 + \dots + \xi_n|\xi) = \frac{1}{n} E(\xi|\xi) = \frac{\xi}{n}.$$

Furthermore,

$$\begin{aligned} E(\xi|\xi_1) &= E(\xi_2 + \dots + \xi_n|\xi_1) + E(\xi_1|\xi_1) = E(\xi_2 + \dots + \xi_n) + \xi_1 \\ &= (n-1)E\xi_1 + \xi_1. \end{aligned}$$

The independence (freezing) lemma

Let ξ be independent from η and $g : \mathbb{R}^2 \mapsto \mathbb{R}$ be a (bounded) Borel function. Then

$$\mathbb{E}(g(\xi, \eta)|\eta) = h(\eta), \quad h(y) = \mathbb{E}g(\xi, y).$$

Proof for indicators. $g = I_{B_1 \times B_2}$:

$$\begin{aligned}\mathbb{E}(g(\xi, \eta)|\eta) &= \mathbb{E}(I_{B_1 \times B_2}(\xi, \eta)|\eta) = \mathbb{E}(I_{B_1}(\xi)I_{B_2}(\eta)|\eta) \\ &= I_{B_2}(\eta)\mathbb{E}(I_{B_1}(\xi)|\eta) = I_{B_2}(\eta)\mathbb{E}(I_{B_1}(\xi)) \\ &= \mathbb{E}(I_{B_1}(\xi)I_{B_2}(y))|_{y=\eta} = h(\eta), \\ h(y) &= \mathbb{E}(I_{B_1}(\xi)I_{B_2}(y)) = \mathbb{E}g(\xi, y).\end{aligned}$$

Example

Let ξ and η be independent exponential r.v. with parameters λ and μ respectively. Let us find $P(\xi \leq \eta)$:

$$P(\xi \leq \eta) = \mathbb{E}E(I_{\{\xi \leq \eta\}}|\eta) = \mathbb{E}h(\eta),$$

$$h(y) = \mathbb{E}I_{\{\xi \leq y\}} = P(\xi \leq y) = \int_0^y \lambda e^{-\lambda x} dx = -e^{-\lambda x} \Big|_{x=0}^{x=y} = 1 - e^{-\lambda y},$$

$$P(\xi \leq \eta) = 1 - \mathbb{E}e^{-\lambda \eta} = 1 - \int_0^\infty e^{-\lambda y} \mu e^{-\mu y} dy = 1 - \frac{\mu}{\lambda + \mu} = \frac{\lambda}{\lambda + \mu}.$$

Transition probability

A function $Q : \mathbb{R}^m \times \mathcal{B}(\mathbb{R}^d) \mapsto [0, 1]$ is called the *transition probability* from \mathbb{R}^m to \mathbb{R}^d , if

- ▶ for any $B \in \mathcal{B}(\mathbb{R}^d)$ the function $y \mapsto Q(y, B)$ is $\mathcal{B}(\mathbb{R}^m)$ -measurable;
- ▶ for any $y \in \mathbb{R}^m$ the function $B \mapsto Q(y, B)$ is a probability measure on $\mathcal{B}(\mathbb{R}^d)$.

Conditional distribution

Consider the random variables $\xi : \Omega \mapsto \mathbb{R}^d$, $\eta : \Omega \mapsto \mathbb{R}^m$. A transition probability $P_{\xi|\eta}(B|y) = Q(y, B)$ from \mathbb{R}^m to \mathbb{R}^d is called a *(regular) conditional distribution of ξ with respect to η* , if

$$E(g(\xi)|\eta)(\omega) = \int_{\mathbb{R}^d} g(x) P_{\xi|\eta}(dx|\eta) \quad \text{P-a.s.}$$

for any bounded Borel function on \mathbb{R}^d . Regular conditional distribution always exists.

The definition above also can be written as follows:

$$E(g(\xi, \eta)|\eta) = \int_{\mathbb{R}^d} g(x, \eta) P_{\xi|\eta}(dx|\eta).$$

Analysis of two-dimensional discrete random variables

Consider two dimensional random variable $(\xi, \eta) : \Omega \mapsto \mathbb{R}^2$ with the distribution

$$P_{\xi, \eta}(B) = P((\xi, \eta) \in B).$$

Assume that there exists a finite or countable set $Z \subset \mathbb{R}^2$ such that

$$P((\xi, \eta) \in Z) = P_{\xi, \eta}(Z) = 1.$$

Let $p_{\xi, \eta}(x, y) = P_{\xi, \eta}(x, y)$ be the correspondent p.m.f. By the change of variables formula

$$Eg(\xi, \eta) = \sum_{(x, y) \in Z} g(x, y) p_{\xi, \eta}(x, y).$$

Let us find $E(g(\xi)|\eta)$ and the conditional expectation of ξ with respect to η .

By the general theory,

$$\mathbb{E}(g(\xi)|\eta) = r(\eta).$$

From the definition of conditional expectation it follows that

$$\mathbb{E}(g(\xi)h(\eta)) = \mathbb{E}(r(\eta)h(\eta))$$

for any bounded Borel function h . By the change of variables formula,

$$\mathbb{E}(g(\xi)h(\eta)) = \sum_{x,y} g(x)h(y)p_{\xi,\eta}(x,y)$$

$$\mathbb{E}(r(\eta)h(\eta)) = \sum_y r(y)h(y)p_{\eta}(y),$$

where $p_{\eta}(y)$ is the marginal distribution of η . Since h is arbitrary, we get

$$r(y)p_{\eta}(y) = \sum_x g(x)p_{\xi,\eta}(x,y).$$

Hence,

$$E(g(\xi)|\eta) = r(\eta) = \sum_x g(x) \frac{p_{\xi,\eta}(x, \eta)}{p_\eta(\eta)}.$$

Introduce the *conditional probability mass*

$$p_{\xi|\eta}(x|y) = \frac{p_{\xi,\eta}(x, y)}{p_\eta(y)} \left(= \frac{P(\xi = x, \eta = y)}{P(\eta = y)} \right).$$

Then

$$E(g(\xi)|\eta) = \sum_x g(x) p_{\xi|\eta}(x|\eta).$$

Hence,

$$P_{\xi|\eta}(A|y) = \sum_{x \in A} p_{\xi|\eta}(x|y)$$

is the conditional distribution of ξ with respect to η .

Total probability formula

$$\begin{aligned} p_{\xi}(x) &= P_{\xi}(x) = E I_{\{\xi=x\}} = E E(I_{\{\xi=x\}}|\eta) = E p_{\xi|\eta}(x|\eta) \\ &= \sum_y p_{\xi|\eta}(x|y) p_{\eta}(y). \end{aligned}$$

Total expectation formula

$$Eg(\xi) = E E(g(\xi)|\eta) = E \sum_x g(x) p_{\xi|\eta}(x|\eta) = \sum_y \sum_x g(x) p_{\xi|\eta}(x|y) p_{\eta}(y).$$

The Bayes formula

$$\begin{aligned} p_{\xi,\eta}(x,y) &= p_{\xi|\eta}(x|y) p_{\eta}(y) = p_{\eta|\xi}(y|x) p_{\xi}(x), \\ p_{\eta|\xi}(y|x) &= \frac{p_{\xi|\eta}(x|y) p_{\eta}(y)}{p_{\xi}(x)} = \frac{p_{\xi|\eta}(x|y) p_{\eta}(y)}{\sum_{y'} p_{\xi|\eta}(x|y') p_{\eta}(y')}. \end{aligned}$$

Analysis of two dimensional continuous random variables

Consider some two dimensional continuous random variable $(\xi, \eta) : \Omega \mapsto \mathbb{R}^2$ with the distribution

$$P_{\xi, \eta}(B) = P((\xi, \eta) \in B).$$

Assume that there exists a p.d.f. $f_{\xi, \eta} : \mathbb{R}^2 \mapsto \mathbb{R}_+$ such that

$$P((\xi, \eta) \in B) = \int_B f_{\xi, \eta}(x, y) dx dy.$$

By the change of variables formula

$$Eg(\xi, \eta) = \int_{\mathbb{R}^2} g(x, y) f_{\xi, \eta}(x, y) dx dy.$$

Let us find $E(g(\xi)|\eta)$ and the conditional distribution ξ with respect to η .

$$\mathbb{E}(g(\xi)|\eta) = r(\eta)$$

$$\mathbb{E}(g(\xi)h(\eta)) = \mathbb{E}(r(\eta)h(\eta)),$$

for any bounded Borel function h . By the change of variables formula,

$$\begin{aligned}\mathbb{E}(g(\xi)h(\eta)) &= \int_{\mathbb{R}^2} g(x)h(y)f_{\xi,\eta}(x,y) \, dx dy \\ &= \int_{\mathbb{R}} h(y) \left(\int_{\mathbb{R}} g(x)f_{\xi,\eta}(x,y) \, dx \right) dy\end{aligned}$$

$$\mathbb{E}(r(\eta)h(\eta)) = \int_{\mathbb{R}} r(y)h(y)f_{\eta}(y) \, dy,$$

where $f_{\eta}(y)$ is the marginal density of η . Since h is arbitrary, we get

$$r(y)f_{\eta}(y) = \int_{\mathbb{R}} g(x)f_{\xi,\eta}(x,y) \, dx.$$

Hence,

$$\mathbb{E}(g(\xi)|\eta) = r(\eta) = \int_{\mathbb{R}} g(x) \frac{f_{\xi,\eta}(x,\eta)}{f_{\eta}(\eta)} dx.$$

Let us introduce the *conditional density*

$$f_{\xi|\eta}(x|y) = \frac{f_{\xi,\eta}(x,y)}{f_{\eta}(y)}.$$

Then

$$\mathbb{E}(g(\xi)|\eta) = \int_{\mathbb{R}} g(x) f_{\xi|\eta}(x|\eta) dx.$$

Hence,

$$\mathbb{P}_{\xi|\eta}(A|y) = \int_{x \in A} f_{\xi|\eta}(x|y) dx$$

is the conditional distribution of ξ with respect to η .

Total probability formula

$$f_{\xi}(x) = \int_{\mathbb{R}} f_{\xi,\eta}(x, y) dy = \int_{\mathbb{R}} f_{\xi|\eta}(x|y) f_{\eta}(y) dy$$

Total expectation formula

$$\begin{aligned} \mathbb{E}g(\xi) &= \mathbb{E}\mathbb{E}(g(\xi)|\eta) = \mathbb{E} \int_{\mathbb{R}} g(x) f_{\xi|\eta}(x|\eta) dx \\ &= \int_{\mathbb{R}} \left(\int_{\mathbb{R}} g(x) f_{\xi|\eta}(x|y) dx \right) f_{\eta}(y) dy \\ &= \int_{\mathbb{R}} \mathbb{E}(g(\xi)|\eta = y) f_{\eta}(y) dy \end{aligned}$$

The Bayes formula

$$f_{\xi,\eta}(x,y) = f_{\xi|\eta}(x|y)f_{\eta}(y) = f_{\eta|\xi}(y|x)f_{\xi}(x),$$

$$f_{\eta|\xi}(y|x) = \frac{f_{\xi|\eta}(x|y)f_{\eta}(y)}{f_{\xi}(x)} = \frac{f_{\xi|\eta}(x|y)f_{\eta}(y)}{\int_{\mathbb{R}} f_{\xi|\eta}(x|y')f_{\eta}(y') dy'}.$$

Freezing lemma for conditional distributions

- ▶ Let ξ, η be independent, then the conditional distribution $\zeta = \psi(\xi, \eta)$ with respect to $\eta = y$ coincides with the distribution of the random variable $\psi(\xi, y)$:

$$P_{\zeta|\eta}(dz|y) = P_{\psi(\xi,y)}(dz).$$

Proof.

$$\begin{aligned} E(g(\psi(\xi, \eta))|\eta) &= E g(\psi(\xi, y)) \Big|_{y=\eta} = \int_{\mathbb{R}} g(z) P_{\psi(\xi,y)}(dz) \Big|_{y=\eta} \\ &= \int_{\mathbb{R}} g(z) P_{\zeta|\eta}(dz|y) \Big|_{y=\eta}. \end{aligned}$$

Example: signal detection

One-dimensional signal $\eta \sim N(0, 1)$. Noise $w \sim N(0, \sigma^2)$ is independent from η . Only the sum $\xi = \eta + w$ is observable. Find the conditional distribution η with respect to ξ .

Solution. To find $f_{\xi|\eta}(x|y)$ fix (freeze) η and write down the distribution of $y + w$:

$$f_{\xi|\eta}(x|y) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-y)^2/(2\sigma^2)}.$$

By the Bayes formula,

$$\begin{aligned} f_{\eta|\xi}(y|x) &= c_1(x) f_{\xi|\eta}(x|y) f_{\eta}(y) = c_2(x) \exp\left(-\frac{1}{2} \left(y^2 + \frac{(x-y)^2}{\sigma^2}\right)\right) \\ &= c_3(x) \exp\left(-\frac{1+\sigma^2}{2\sigma^2} \left(y - \frac{x}{1+\sigma^2}\right)^2\right), \end{aligned}$$

Hence for fixed $\xi = x$ the signal $\eta \sim N\left(\frac{x}{1+\sigma^2}, \frac{\sigma^2}{1+\sigma^2}\right)$.

ξ is continuous, η is discrete

Let ξ be continuous with the p.d.f. $f_\xi(x)$, and let η be discrete with the p.m.f. $p_\eta(y)$. *Introduce* the conditional density $f_{\xi|\eta}(x|y)$,

$$f_{\xi|\eta}(x|y) \geq 0, \quad \int_{\mathbb{R}} f_{\xi|\eta}(x|y) dx = 1.$$

In other words, assume that the conditional distribution of ξ with respect to η is determined by the formula

$$\mathbb{E}(g(\xi, \eta)|\eta) = \int_{\mathbb{R}} g(x, \eta) f_{\xi|\eta}(x|\eta) dx.$$

Let us find the joint distribution:

$$\begin{aligned} \mathbb{E}g(\xi, \eta) &= \mathbb{E}\mathbb{E}(g(\xi, \eta)|\eta) = \mathbb{E} \int_{\mathbb{R}} g(x, \eta) f_{\xi|\eta}(x|\eta) dx \\ &= \sum_y \int_{\mathbb{R}} g(x, y) f_{\xi|\eta}(x|y) dx \cdot p_\eta(y). \end{aligned}$$

In particular we get the total expectation formula:

$$\mathbb{E}g(\xi) = \sum_y \int_{\mathbb{R}} g(x) f_{\xi|\eta}(x|y) dx \cdot p_{\eta}(y) = \int_{\mathbb{R}} g(x) \sum_y f_{\xi|\eta}(x|y) p_{\eta}(y) dx$$

and the total probability formula:

$$f_{\xi}(x) = \sum_y f_{\xi|\eta}(x|y) p_{\eta}(y).$$

Let us find $\mathbb{E}(g(\eta)|\xi) = r(\xi)$. By the definition,

$$\mathbb{E}(g(\eta)h(\xi)) = \mathbb{E}(r(\xi)h(\xi))$$

for any bounded Borel function h .

$$\begin{aligned}
E(g(\eta)h(\xi)) &= \sum_y \int_{\mathbb{R}} g(y)h(x)f_{\xi|\eta}(x|y) dx \cdot p_{\eta}(y) \\
&= \int_{\mathbb{R}} \left(\sum_y g(y)f_{\xi|\eta}(x|y)p_{\eta}(y) \right) h(x) dx, \\
E(r(\xi)h(\xi)) &= \int_{\mathbb{R}} r(x)h(x)f_{\xi}(x) dx.
\end{aligned}$$

Hence,

$$\begin{aligned}
r(x) &= \sum_y g(y) \frac{f_{\xi|\eta}(x|y)p_{\eta}(y)}{f_{\xi}(x)}, \\
E(g(\eta)|\xi) = r(\xi) &= \sum_y g(y) \frac{f_{\xi|\eta}(\xi|y)p_{\eta}(y)}{f_{\xi}(x)} = \sum_y g(y)p_{\eta|\xi}(y|\xi),
\end{aligned}$$

where the conditional density $p_{\eta|\xi}(y|x)$ is determined by the Bayes formula:

$$p_{\eta|\xi}(y|x) = \frac{f_{\xi|\eta}(x|y)p_{\eta}(y)}{f_{\xi}(x)} = \frac{f_{\xi|\eta}(x|y)p_{\eta}(y)}{\sum_{y'} f_{\xi|\eta}(x|y')p_{\eta}(y')}.$$

Example: signal detection

A binary signal is described by the random variable η :

$$p_{\eta}(1) = p, \quad p_{\eta}(-1) = q = 1 - p.$$

The signal is affected by the standard Gaussian noise $w \sim N(0, 1)$, which does not depend on η . The received signal has the form

$$\xi = \eta + w.$$

Find the probability that $\eta = 1$ under the condition that the observed value of ξ equals to x .

Solution. $f_{\xi|\eta}(x|y)$ coincides with the density of the random variable $y + w \sim N(y, 1)$:

$$f_{\xi|\eta}(x|y) = \frac{1}{\sqrt{2\pi}} e^{-(x-y)^2/2}.$$

By the Bayes formula

$$p_{\eta|\xi}(y|x) = \frac{f_{\xi|\eta}(x|y)p_{\eta}(y)}{\sum_{y'} f_{\xi|\eta}(x|y')p_{\eta}(y')}.$$

In particular,

$$p_{\eta|\xi}(1|x) = \frac{e^{-(x-1)^2/2}p}{e^{-(x-1)^2/2}p + e^{-(x+1)^2/2}q} = \frac{p}{p + qe^{-2x}}.$$

Note that the function $x \mapsto p_{\eta|\xi}(1|x)$ is increasing and

$$\lim_{x \rightarrow -\infty} p_{\eta|\xi}(1|x) = 0, \quad \lim_{x \rightarrow \infty} p_{\eta|\xi}(1|x) = 1.$$

ξ is discrete, η is continuous

Let ξ be discrete with the p.m.f. $p_\xi(x)$, and η is continuous with the density $f_\eta(y)$. *Introduce* the conditional p.m.f. $p_{\xi|\eta}(x|y)$:

$$p_{\xi|\eta}(x|y) \geq 0, \quad \sum_x p_{\xi|\eta}(x|y) = 1.$$

In other words, assume that the conditional distribution of ξ with respect to η is determined by the formula

$$\mathbb{E}(g(\xi, \eta)|\eta) = \sum_x g(x, \eta) p_{\xi|\eta}(x|\eta).$$

Let us find the joint distribution:

$$\begin{aligned} \mathbb{E}g(\xi, \eta) &= \mathbb{E}\mathbb{E}(g(\xi, \eta)|\eta) = \mathbb{E} \sum_x g(x, \eta) p_{\xi|\eta}(x|\eta) \\ &= \int_{\mathbb{R}} \sum_x g(x, y) p_{\xi|\eta}(x|y) f_\eta(y) dy. \end{aligned}$$

A special case is the total expectation formula:

$$\mathbb{E}g(\xi) = \sum_x g(x) \int_{\mathbb{R}} p_{\xi|\eta}(x|y) f_{\eta}(y) dy$$

which implies the total probability formula

$$p_{\xi}(x) = \int_{\mathbb{R}} p_{\xi|\eta}(x|y) f_{\eta}(y) dy.$$

Let us find $\mathbb{E}(g(\eta)|\xi) = r(\xi)$. By the definition,

$$\mathbb{E}(g(\eta)h(\xi)) = \mathbb{E}(r(\xi)h(\xi))$$

for any bounded Borel function h .

$$\begin{aligned}
\mathbb{E}(g(\eta)h(\xi)) &= \int_{\mathbb{R}} \sum_x g(y)h(x)p_{\xi|\eta}(x|y)f_{\eta}(y) dy \\
&= \sum_x \left(\int_{\mathbb{R}} g(y)p_{\xi|\eta}(x|y)f_{\eta}(y) dy \right) h(x), \\
\mathbb{E}(r(\xi)h(\xi)) &= \sum_x r(x)h(x)p_{\xi}(x).
\end{aligned}$$

Hence,

$$\begin{aligned}
r(x) &= \int_{\mathbb{R}} g(y) \frac{p_{\xi|\eta}(x|y)f_{\eta}(y)}{p_{\xi}(x)} dy, \\
\mathbb{E}(g(\eta)|\xi) = r(\xi) &= \int_{\mathbb{R}} g(y) \frac{p_{\xi|\eta}(\xi|y)f_{\eta}(y)}{p_{\xi}(\xi)} dy = \int_{\mathbb{R}} g(y)f_{\eta|\xi}(y|\xi) dy,
\end{aligned}$$

where the conditional density $f_{\eta|\xi}(y|x)$ is determined by the Bayes formula:

$$f_{\eta|\xi}(y|x) = \frac{p_{\xi|\eta}(x|y)f_{\eta}(y)}{p_{\xi}(x)} = \frac{p_{\xi|\eta}(x|y)f_{\eta}(y)}{\int_{\mathbb{R}} p_{\xi|\eta}(x|y')f_{\eta}(y') dy'}.$$