
basics blog

basics

Jun 25, 2025

CONTENTS

I	Introduction to Statistics	3
1	Introduction to Statistics	5
II	Probability Theory	7
2	Introduction to probability theory	9
2.1	Definition of stochastic variable	10
2.2	Discrete stochastic variables	10
2.3	Continuous stochastic variables	10
2.4	Multi-dimensional stochastic variables	11
2.5	Transformations of probability functions	14
2.6	Characteristic functions	14
2.7	Convergence in statistics	15
2.8	Independent identically distributed random variables	16
2.9	Heavy-tailed distributions	18
3	Stochastic processes	19
3.1	Wiener process - Brownian motion	20
3.2	White noise	23
3.3	Stochastic calculus	23
III	Statistical Inference	27
4	Introduction to Statistical Inference	29
IV	Introduction to Machine Learning	31
5	Introduction to Machine Learning	33
5.1	Models in Machine Learning	34
5.2	Good Practices in Machine Learning	34
6	Supervised Learning	35
6.1	SL: theory	35
7	Unsupervised Learning	37
8	Reinforcement Learning	39

This material is part of the **basics-books project**.

Contents.

Introduction to statistics

Different approaches to statistics and **descriptive statistics**

Probability theory

Inferential statistics

Inferential and Bayesian statistics

Introduction to Machine Learning: SL, UL, ML

Machine learning (ML) is a branch of artificial intelligence (AI) focused on designing systems that can learn from data to improve their performance on a task. ML frameworks include supervised learning (e.g., regression and classification), unsupervised learning (e.g., clustering, compression, principal component analysis), and reinforcement learning (e.g., planning and control). ML emphasizes practical problem-solving, grounded in statistical methods, numerical optimization, and enabled by advances in computing hardware.

Part I

Introduction to Statistics

INTRODUCTION TO STATISTICS

Part II

Probability Theory

INTRODUCTION TO PROBABILITY THEORY

Probability theory is an axiomatic approach to probability, assigning

Stochastic variables

- Definition of stochastic variable
- Discrete and continuous stochastic variables
 - Probability functions, moments (if they exists, see heavy-tailed distribution), and examples
- Multi-dimensional stochastic variables:
 - joint, conditional, marginal probability
 - Bayes' theorem
 - independence
 - moments: covariance, correlation
- Generators...
- I.i.d. variables: law of large numbers, central limit theorem; convergence of statistics (reference to measure in the definition of a sthochastic variable)
- Sampling
- Extra:
 - heavy tails probability functions

Stochastic processes

- Definition of stochastic process
- Time-continuous/time-discrete
- Ergodicity and stationarity:
 - moments, correlation,...
 - analysis in time and Fourier domains of time-signals
- Applications:
 - example of processes:
 - * white noise

- * Wiener process (Brownian motion): definition, application, relation with
- * discrete-time Markov process (useful in RL, can be interpreted as a discretized continuous process)
- response of LTI to random input

Stochastic fields

2.1 Definition of stochastic variable

2.2 Discrete stochastic variables

2.3 Continuous stochastic variables

2.3.1 Examples

Here some common examples of continuous random variables are introduced. Their functional dependence on the value of the r.v. is quite easy to remember, while the normalization factor could look quite “esoteric”.

Normal distribution, $\mathcal{N}(\mu, \sigma^2)$

pdf is

$$f(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \propto e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Moment	Value
Expected value	μ
Variance	σ^2

Expected value, μ ; variance, σ^2 .

Unitarity

$$\int_{x=-\infty}^{+\infty} e^{-\frac{(x-\mu)^2}{2\sigma^2}} = \sqrt{2\pi\sigma^2}$$

todo integral $\int_{-\infty}^{+\infty} e^{-\alpha x^2} dx$

Expected value

$$\mathbb{E}[X] = \int_{x=-\infty}^{+\infty} x \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Variance

$$\mathbb{E}[(X - \mu)^2] = \dots$$

Chi-square, χ_N^2

$$\chi_N^2 := \sum_{n=1}^N X_n^2$$

pdf is

$$f(x; n) = \dots \propto x^{\frac{n}{2}-1} e^{-\frac{x}{2}}$$

Student- t distribution, t_ν

$$t_\nu = \frac{Z}{\sqrt{\frac{K}{\nu}}},$$

with $Z \sim \mathcal{N}(0, 1)$, and $K \sim \chi_\nu^2$.

pdf is

$$f(x; n) = \dots \propto \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}}$$

2.4 Multi-dimensional stochastic variables

- joint distribution

$$p_{XY}(x, y)$$

- marginal distribution. For continuous variables

$$p_X(x) := \int_y p_{XY}(x, y) dy$$

while for discrete variables

$$p_X(x_i) = \sum_j p_{XY}(x_i, y_j)$$

- conditional distribution, $p_{X|Y}(x|y)$. The following holds

$$p_{XY} = p_{X|Y} p_Y = p_{Y|X} p_X$$

For continuous r.v., integrating over x the relation $p(x, y) = p(x|y)p(y)$

$$\int_x p(x, y) dx = \int_x p(x|y) p(y) dx = p(y) \underbrace{\int_x p(x|y) dx}_{=1} = p(y),$$

as the normalization condition holds for conditional distribution $p(x|y)$.

Property 2.4.1

$$p(i, j) = p(i|j)p(j)$$

$$\sum_i p(i, j) = \sum_i \underbrace{p(i|j)}_{=1} p(j) = p(j)$$

2.4.1 Moments

- expected value

$$\mu_{\mathbf{X}} := \mathbb{E} [\mathbf{X}] = \int_{\mathbf{x}} p(\mathbf{x}) \mathbf{x} d\mathbf{x}$$

- covariance

$$\sigma_{\mathbf{X}}^2 := \mathbb{E} [\Delta \mathbf{X} \Delta \mathbf{X}^T] = \int_{\mathbf{x}} p(\mathbf{x}) \Delta \mathbf{x} \Delta \mathbf{x}^T d\mathbf{x} ,$$

with $\Delta \mathbf{X} := \mathbf{X} - \mu_{\mathbf{X}}$, and $\Delta \mathbf{x} = \mathbf{x} - \mu_{\mathbf{X}}$.

Taking a pair of components X_i, X_j of the random vector \mathbf{X} , their covariance is the ij component of the array σ^2 ,

$$\sigma_{ij}^2 := \mathbb{E} [\Delta X_i \Delta X_j] =: \rho_{ij} \sigma_i \sigma_j ,$$

having introduced **(Pearson) correlation**, ρ_{ij} , between random variable X_i and X_j , and being σ_i the standard deviation of variable X_i , square root of its variance σ_i^2 ,

$$\begin{aligned} \sigma_i^2 &= \mathbb{E} [(X_i - \mu_i)^2] = \\ &= \int_{\mathbf{x}} (x_i - \mu_i)^2 p_{\mathbf{X}}(\mathbf{x}) d\mathbf{x} = \\ &= \int_{x_i} (x_i - \mu_i)^2 p_i(x_i) dx_i \end{aligned}$$

Here the integrals read

$$\begin{aligned} \mu_i &= \int_{\mathbf{x}} x_i p_{\mathbf{X}}(\mathbf{x}) d\mathbf{x} = \\ &= \int_{\mathbf{x}} x_i p(x_1, x_2, \dots, x_i, \dots, x_n) dx_1 dx_2 \dots dx_i \dots dx_n = \\ &= \int_{\mathbf{x}} x_i p(x_i) p(x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n | x_i) dx_1 dx_2 \dots dx_i \dots dx_n = \\ &= \int_{x_i} x_i p(x_i) \underbrace{\int_{x_1} \dots \int_{x_n} p(x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n | x_i) dx_1 \dots dx_{i-1} dx_{i+1} \dots dx_n}_{=1 \forall x_i} dx_i = \\ &= \int_{x_i} x_i p(x_i) dx_i . \end{aligned}$$

Property of correlation. $|\rho_{XY}| \leq 1$. Proof with Cauchy-Schwartz inequality **todo**

Notation

Here, covariance is indicated as σ^2 . This is not a power 2, but just a symbol, at most recalling that covariance matrix is **semi-definite positive**.

Properties of covariance.

- symmetric
- semi-definite positive
- spectrum...

2.4.2 Bayes' theorem

Theorem 2.4.1 (Bayes' theorem)

Where $p_Y(y) \neq 0$,

$$p_{X|Y}(x|y) = \frac{p_{XY}(x,y)}{p_Y(y)}$$

2.4.3 Statistical independence

Definition 2.4.1 (Independent random variables)

Given two random variables X, Y with joint distribution, the random variable X is independent from Y if its conditional probability equals its marginal probability,

$$p_{X|Y} = p_X ,$$

i.e. the probability of X doesn't depend on Y .

Independence implies no correlation

Given two random variables X, Y are independent if $p(x|y) = p(x)$ and thus $p(x, y) = p(x)p(y)$. Covariance of two random variable reads

$$\sigma_{xy}^2 = \mathbb{E}[(X - \mu_X)(Y - \mu_Y)] ,$$

and if they're independent, it immediately follows that their covariance σ_{XY}^2 is zero (and so their correlation ρ_{XY})

$$\sigma_{xy}^2 = \underbrace{\mathbb{E}[X - \mu_X]}_{=0} \underbrace{\mathbb{E}[Y - \mu_Y]}_{=0} = 0 ,$$

as the expected value of the deviation from the expected value is zero, $\mathbb{E}[X - \mathbb{E}[X]] = 0$.

Proof for continuous r.v.

$$\begin{aligned}\sigma_{xy}^2 &= \mathbb{E}[(X - \mu_X)(Y - \mu_Y)] = \\ &= \int_{x,y} (x - \mu_X)(y - \mu_Y)p(x,y) dx dy = \quad (1)\end{aligned}$$

$$\begin{aligned}&= \int_{x,y} (x - \mu_X)(y - \mu_Y)p(x)p(y) dx dy = \\ &= \int_x (x - \mu_X)p(x)dx \int_y (y - \mu_Y)p(y)dy = \quad (2)\end{aligned}$$

having used here the common notation abuse $p_X(x) = p(x)$ and (1) statistical independence, $p(x,y) = p(x)p(y)$, and (2) $\mathbb{E}[X - \mathbb{E}[X]] = 0$.

Proof for discrete r.v.

Repeat the proof for continuous r.v. using summations instead of integrals.

2.5 Transformations of probability functions

2.6 Characteristic functions

Characteristic function of a random variable X is defined as

$$\varphi_X(t) := \mathbb{E}[e^{itX}] .$$

Characteristic function of a continuous random variable with probability density function $f(x)$ thus reads

$$\varphi_X(t) := \mathbb{E}[e^{itX}] = \int_{x \in D_x} f(x)e^{itx} dx ,$$

i.e. its the **Fourier transform** of its pdf.

Example 2.6.1 (Characteristic function of a multi-dimensional variable)

$$Z(\mathbf{Y})$$

$$\varphi_{Z(\mathbf{Y})} := \mathbb{E}[e^{itZ(\mathbf{Y})}] = \int_{\mathbf{y}} e^{itZ(\mathbf{y})} f(\mathbf{y}) d\mathbf{y}$$

Example 2.6.2 (Characteristic function of a linear combination of independent variables)

$$Z(\mathbf{Y}) = a_1 Y_1 + \dots a_n Y_n ,$$

with

$$f(\mathbf{y}) = f(y_1, \dots, y_n) = f_1(y_1) \dots f_n(y_n) .$$

$$\begin{aligned}
\varphi_{Z(\mathbf{Y})} &:= \mathbb{E} [e^{itZ(\mathbf{Y})}] = \\
&= \int_{\mathbf{y}} e^{it(\sum_k a_k y_k)} f(\mathbf{y}) d\mathbf{y} = \\
&= \int_{y_1} e^{it a_1 y_1} f_1(y_1) dy_1 \dots \int_{y_n} e^{it a_n y_n} f_n(y_n) dy_n = \\
&= \varphi_{Y_1}(a_1 t) \dots \varphi_{Y_n}(a_n t) .
\end{aligned}$$

Example 2.6.3 (Taylor expansion of characteristic function)

For “small” values of t , an approximation of the characteristic function is provided by Taylor expansion around $t = 0$,

$$\begin{aligned}
\int e^{iyt} f(y) dy &= \int \left[1 + iyt - \frac{1}{2}(yt)^2 + o(t^2) \right] f(y) dy = \quad (1) \\
&= 1 + i\mu t - \frac{1}{2}t^2 (\sigma^2 + \mu^2) + o(t^2)
\end{aligned}$$

as (1) $\sigma^2 = \mathbb{E}[(y - \mu)^2] = \mathbb{E}[y^2] - \mu^2$

Example 2.6.4 (Characteristic function of a normal distribution $\mathcal{N}(0, 1)$)

$$\begin{aligned}
f(x) &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right) \\
\int_{x=-\infty}^{+\infty} e^{ixt} f(x) dx &= \frac{1}{\sqrt{2\pi}} \int_{x=-\infty}^{+\infty} e^{ixt - \frac{x^2}{2}} dx = \quad (1) \\
&= \frac{1}{\sqrt{2\pi}} \int_{x=-\infty}^{+\infty} e^{-\frac{(x-it)^2}{2}} dx e^{-\frac{t^2}{2}} = \quad (2) \\
&= \frac{1}{\sqrt{2\pi}} \sqrt{2\pi} e^{-\frac{t^2}{2}} = e^{-\frac{t^2}{2}}
\end{aligned}$$

having (1) completed the square $(x - it)^2 = x^2 - i2xt - t^2$, and evaluated the integral **todo** (it's similar to the standard result $\int_{-\infty}^{+\infty} e^{-x^2} dx = \sqrt{2\pi}$, but with complex variable. Link to math material, complex calculus).

2.7 Convergence in statistics

2.7.1 Convergence in distribution - weak convergence

A sequence of X_i of real-valued random variables, cumulative distribution functions F_i , converges in distribution to a random variable X with cumulative distribution F is

$$\lim_{n \rightarrow +\infty} F_n(x) = F(x) ,$$

for $\forall x \in \mathbb{R}$ where $F(x)$ is continuous.

For multi-valued random variables, the condition reads

$$\lim_{n \rightarrow +\infty} P(X_n \in A) = P(X \in A) ,$$

for every $A \subset \mathbb{R}^n$...**todo**

2.7.2 Convergence in probability

$$\lim_{n \rightarrow +\infty} P(|X_n - X| > \varepsilon) = 0$$

Warning: Convergence in probability and convergence in distribution

Convergence in probability \rightarrow convergence in distribution, but not viceversa.

Example taken from wikipedia

2.7.3 Almost sure convergence - strong convergence

$$P\left(\lim_{n \rightarrow +\infty} X_n = X\right) = 1$$

i.e. events for which X_n doesn't converge to X has probability 0,

$$P\left(\omega \in \Omega : \lim_{n \rightarrow +\infty} X_n(\omega) = X(\omega)\right) = 1$$

2.7.4 Sure convergence - pointwise convergence

$$\left\{\omega \in \Omega : \lim_{n \rightarrow +\infty} X_n(\omega) = X(\omega)\right\} = \Omega .$$

The same definition of almost sure convergence, without allowing the existence of sets with zero probability where convergence is not satisfied. Thus, it's likely there is no point in using sure convergence instead of almost sure convergence in probability theory.

2.7.5 Convergence in absolute moments: mean,...

$$\lim_{n \rightarrow +\infty} \mathbb{E}(|X_n - X|^r) = 0$$

2.8 Independent identically distributed random variables

Definition 2.8.1 (Independent identically distributed (iid) random variables)

2.8.1 Law of the large numbers

Weak form

todo

Strong form

todo

2.8.2 Central Limit Theorem

Theorem 2.8.1 (CLT)

Let $\{X_k\}_{k=1:n}$ a sequence of iid random variables with average value $\mathbb{E}[X_k] = \mu$ and **finite**¹ variance $\mathbb{E}[(X_k - \mu)^2] = \sigma^2 < \infty$, then the **sample average**

$$\bar{X}_n := \frac{1}{n} \sum_{k=1}^n X_k ,$$

converges in distribution - or weakly converges - to the normal distribution $\mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$,

$$\bar{X}_n \rightarrow^d \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right) .$$

Proof of CLT

Let $\{X_k\}_{k=1:n}$ the sequence of iid random variables. Thus, $\sum_{k=1}^n X_k$ has expected value $n\mu$ and variance $n\sigma^2$. Let

$$Z_n := \frac{\sum_{k=1}^n X_k - n\mu}{\sqrt{n\sigma^2}} = \sum_{k=1}^n \frac{X_k - \mu}{\sqrt{n\sigma^2}} =: \sum_{k=1}^n \frac{Y_k}{\sqrt{n}} .$$

Expected value and variance of variables Y_k are respectively $\mathbb{E}[Y_k] = 0$ and $\mathbb{E}[Y_k^2] = 1$. The *characteristic function* of Z_n , see [Example 2.6.2](#) for the linear combination of independent variables, reads

$$\begin{aligned} \varphi_{Z_n}(t) &= \varphi_{Y_1}\left(\frac{t}{\sqrt{n}}\right) \dots \varphi_{Y_n}\left(\frac{t}{\sqrt{n}}\right) = \quad (1) \\ &= \left[\varphi_{Y_1}\left(\frac{t}{\sqrt{n}}\right) \right]^n , \end{aligned}$$

as (1) the variables are not only independent but identically distributed: as they have the same pdf, they also have the same characteristic function. Expanding in Taylor series, see [example Example 2.6.3](#) for $\frac{t}{\sqrt{n}} \rightarrow 0$, the approximation of the characteristic function reads (remembering that Y_n have zero expected value and unit variance),

$$\varphi_{Y_1}\left(\frac{t}{\sqrt{n}}\right) \sim 1 - \frac{t^2}{2n} ,$$

while

$$\varphi_{Z_n}(t) = \left[\varphi_{Y_1}\left(\frac{t}{\sqrt{n}}\right) \right]^n \sim \left[1 - \frac{t^2}{2n} \right]^n \sim e^{-\frac{t^2}{2}} ,$$

i.e. it converges to the characteristic function of a normal distribution $\mathcal{N}(0, 1)$, see [Example 2.6.4](#).

Levy's continuity theorem completes the proof. **todo**

¹ Does the CLT hold for *heavy-tailed distributions*?

2.9 Heavy-tailed distributions

- Does the *CLT* hold for heavy-tailed distributions?

2.9.1 References

- Clauset, Aaron and Woodard, Ryan, “Estimating the historical and future probabilities of large terrorist events.” *Annals of Applied Statistics* 7(4), 1838-1865 (2013).

STOCHASTIC PROCESSES

Definition of stochastic process.

Examples.

- White noise, $\xi(t)$, is a zero-mean process with no correlation between its values at different times

$$\mathbb{E} [\xi(t) \xi(s)] = \delta(t - s)$$

- Wiener process (Brownian motion), $W(t)$

$$W(0) = 0$$

$W(t)$ has independent increments

$$W(t) - W(s) \sim N(0, t - s) \text{ for } t > s$$

$W(t)$ are continuous but nowhere differentiable

Informal relation between Wiener process and white noise signal

$$W(t) - W(s) = \int_s^t \xi(\tau) d\tau$$
$$\text{``} \frac{dW(t)}{dt} = \xi(t) \text{''}$$

where the derivative relation doesn't hold in the classical sense, as $W(t)$ is nowhere differentiable

- time-discrete Markov processes

Applications

- LTI
- Stochastic differential equations...

$$dX(t) = \mu(t) dt + \sigma(t) dW(t)$$

Assumptions.

- Stationarity
- Ergodicity

$$k_{xy}(\tau) := \mathbb{E}[x(t)y(t-\tau)] = \lim_{T \rightarrow +\infty} \left\{ \frac{1}{2T} \int_{t=-T}^T x(t) y(t-\tau) dt \right\}$$

3.1 Wiener process - Brownian motion

- Introduction: history and relation with other problems (diffusion?)
- Definition and some theory
- Simulation of Wiener process, demonstration of properties shown in theory section

3.1.1 Definition

Definition 3.1.1 (Wiener process - Brownian motion)

A Wiener process is a random process $W(t)$ with

1. initial condition, almost surely

$$W(0) = 0$$

2. increments with zero-mean *normal distribution*

$$W(t) - W(s) \sim \mathcal{N}(0, |t - s|)$$

3. W has independent increments: $W(t) - W(t + u)$ is independent from W_s , $s < t$
 4. in $W(t)$ is almost surely continuous in t
-

Almost sure convergence in statistics

“Almost surely” here means *almost sure convergence* and it is explained in the section dealing with *convergence in statistics*, and used below to prove some properties of a Wiener process.

3.1.2 Properties

Property 3.1.1 (Covariance of increments)

Covariance of an increment follows the definition of Wiener process and the definition of *normal distribution*,

$$\mathbb{E} [(W(t) - W(s))^2] = \mathbb{E} [\mathcal{N}(0, |t - s|)] = |t - s|. \quad (3.1)$$

Covariance of independent increments - on non-overlapping ranges - is zero, as *independence implies no correlation*, i.e. zero covariance. Thus, if $a \leq b \leq c \leq d$, $W(b) - W(a)$ and $W(d) - W(c)$ are independent by property (3) in *Definition 3.1.1* of Wiener process, and thus their covariance - and correlation - is zero,

$$\mathbb{E} [(W(b) - W(a))(W(d) - W(c))] = 0 \quad (3.2)$$

Covariance of two generic increments reads

$$\mathbb{E} [(W(t_1) - W(s_1))(W(t_2) - W(s_2))] = |[s_1, t_1] \cap [s_2, t_2]| \quad (3.3)$$

as it's proved below.

Proof of the covariance of two generic increments

If $s_1 \leq s_2 \leq t_2 \leq t_1$,

$$\begin{aligned}
 \mathbb{E}[(W(t_1) - W(s_1))(W(t_2) - W(s_2))] &= \\
 &= \mathbb{E}[(W(t_1) - W(t_2) + W(t_2) - W(s_2) + W(s_2) - W(s_1))(W(t_2) - W(s_2))] = \\
 &= \underbrace{\mathbb{E}[(W(t_1) - W(t_2))(W(t_2) - W(s_2))]}_{=0} + \underbrace{\mathbb{E}[(W(t_2) - W(s_2))(W(t_2) - W(s_2))]}_{=|t_2 - s_2|} + \\
 &\quad + \underbrace{\mathbb{E}[(W(s_2) - W(s_1))(W(t_2) - W(s_2))]}_{=0} = \\
 &= 0 + |t_2 - s_2| + 0 = |[s_1, t_1] \cap [s_2, t_2]|.
 \end{aligned}$$

Similarly, if $s_1 \leq s_2 \leq t_1 \leq t_2$,

$$\begin{aligned}
 \mathbb{E}[(W(t_1) - W(s_1))(W(t_2) - W(s_2))] &= \\
 &= \mathbb{E}[(W(t_1) - W(s_2) + W(s_2) - W(s_1))(W(t_2) - W(t_1) + W(t_1) - W(s_2))] = \\
 &= \underbrace{\mathbb{E}[(W(t_1) - W(s_2))(W(t_2) - W(t_1))]}_{=0} + \underbrace{\mathbb{E}[(W(t_1) - W(s_2))(W(t_1) - W(s_2))]}_{=|t_1 - s_2|} + \\
 &\quad + \underbrace{\mathbb{E}[(W(s_2) - W(s_1))(W(t_2) - W(t_1))]}_{=0} + \underbrace{\mathbb{E}[(W(s_2) - W(s_1))(W(t_1) - W(s_2))]}_{=0} = \\
 &= 0 + |t_1 - s_2| + 0 = |[s_1, t_1] \cap [s_2, t_2]|.
 \end{aligned}$$

All the other situations can be proved in the same way.

Property 3.1.2 (Statistics of maximum)

For $a \geq 0$,

$$P(M(t) \geq a) = 2P(W(t) \geq a) = 2 - 2\phi\left(\frac{a}{\sqrt{t}}\right),$$

with

$$M(t) = \max_{0 \leq \tau \leq t} W(\tau)$$

and

$$\phi(x) = \int_{y=-\infty}^x p_{\mathcal{N}(0,1)}(x) dx$$

is the cumulative probability function of a normal distribution $\mathcal{N}(0, 1)$.

Proof.

The second inequality immediately follows from the very definition of Wiener process with initial conditions $W(0)$,

$$P(W(t) - W(0) \geq a) = P(\mathcal{N}(0, t) \geq a) = \quad (1)$$

$$\begin{aligned}
 &= P\left(\mathcal{N}(0, 1) > \frac{a}{\sqrt{t}}\right) = \\
 &= \int_{x=\frac{a}{\sqrt{t}}}^{+\infty} p(y) dy = \quad (2) \\
 &= 1 - \int_{x=-\infty}^{\frac{a}{\sqrt{t}}} p(y) dy = 1 - \phi\left(\frac{a}{\sqrt{t}}\right)
 \end{aligned}$$

having used (1) scaling rule for *transformation of probability functions*... **todo**, and (2) the normalization condition of the probability density $1 = \int_{x=-\infty}^{+\infty} p(x) dx$, and the definition of cumulative probability function.

First inequality. In order to prove the first inequality, it could be useful to introduce the definition of **stepping time**, τ_a , as the random variable defined as

$$\tau_a = \min_s \{s : W(s) = a\} .$$

Using *reflection principle*, it follows

$$P(M(t) \geq a) = \quad (1)$$

$$= P(M(t) \geq a, W(t) \geq a) + P(M(t) \geq a, W(t) < a) = \quad (2)$$

$$= P(W(t) \geq a) + P(M(t) \geq a, W(t) - W(\tau_a) < 0) = \quad (3)$$

$$= P(W(t) \geq a) + P(M(t) \geq a, W'(t - \tau_a) < 0) = \quad (4)$$

$$= P(W(t) \geq a) + P(M(t) \geq a)P(W'(t - \tau_a) < 0) = \quad (5)$$

$$= P(W(t) \geq a) + \frac{1}{2}P(M(t) \geq a) .$$

haing (1) used “marginalization” to write $P(A) = P(A, B) + P(A, \bar{B})$, (2) recognized that if $B : W(t) \geq a$ then $A : M(t) \geq a$ or $B \subseteq A$, and thus $P(A, B) = P(B)$, and that $a = W(\tau_a)$, (3) defined the Wiener process $W'(t - \tau_a) := W(t) - W(\tau_a)$, independent from $W(s)$, $0 \leq s \leq \tau_a$, (4) exploited the independence of the two conditions (**todo be more explicit, proof needed?**), (5) and the symmetry of Wiener process to get $P(W'(t - \tau_a) < 0) = \frac{1}{2}$.

Thus, it follows the requied relation

$$P(M(t) \geq a) = 2P(W(t) \geq a) .$$

Property 3.1.3 ($W(t)$ is almost surely not differentiable)

For all time t , a Wiener process is almost surely not differentiable, i.e. ...**todo**

Proof.

todo check details

Wiener process is differentiable in t if the limit

$$\lim_{h \rightarrow 0} \frac{W(t+h) - W(t)}{h} = \ell$$

exists finite. Definition of limit reads,

$$\forall \varepsilon > 0 \quad \exists U_{0,\delta} \quad \text{s.t.} \quad \left| \frac{W(t+h) - W(t)}{h} - \ell \right| < \varepsilon \quad \forall h \in U_{0,\delta} \setminus \{0\}$$

todo how to go from this definition to the following one?

Let E_{ε,A,t_0} be the event s.t. for a given t_0 , $W(t)$ is differentiable in t_0 , i.e. $\exists A, \varepsilon_0$ const. s.t. $\forall \varepsilon$ s.t. $0 < \varepsilon < \varepsilon_0$, $W(t) - W(t_0) \leq A\varepsilon$ holds for $\forall \varepsilon, 0 < t - t_0 \leq \varepsilon$.

Let $E_{A,t_0} = \cap_{\varepsilon} E_{\varepsilon,A,t_0}$. Then

$$P(E_{\varepsilon,A,t_0}) = P(|W(t) - W(t_0)| \leq A\varepsilon \text{ for } \forall t - t_0 \text{ s.t. } 0 < t - t_0 \leq \varepsilon) = \quad (1)$$

$$= P(M(t - t_0) \leq A\varepsilon) = \quad (2)$$

$$= 1 - P(M(t - t_0) \geq A\varepsilon) = \quad (3)$$

$$= 1 - \left[2 - 2\phi\left(\frac{A\varepsilon}{\sqrt{\Delta t}}\right) \right] =$$

$$= -1 + 2\phi\left(\frac{A\varepsilon}{\sqrt{\Delta t}}\right),$$

having used (1)..., (2)..., (3)...

Now, being $\varepsilon \leq \Delta t$, it follows that $\frac{\varepsilon}{\sqrt{\Delta t}} \leq \sqrt{\Delta t}$. As $\varepsilon \rightarrow 0$, then $\frac{\varepsilon}{\sqrt{\Delta t}} \rightarrow 0$, and $\phi\left(\frac{A\varepsilon}{\sqrt{\Delta t}}\right) \rightarrow \frac{1}{2}$, and $P(E_{\varepsilon,A,t_0}) \rightarrow 0$

3.2 White noise

Definition 3.2.1 (White noise - properties)

A white noise is a random process with

- zero expected value

$$\mathbb{E}[\xi(t)] = 0$$

- Dirac delta correlation

$$\mathbb{E}[\xi(t)\xi(s)] = \delta(t - s)$$

todo link to *math:functional-analysis:distributions*

Definition 3.2.2 (White noise - time derivative of Wiener process $W(t)$ in the sense of distributions)

3.3 Stochastic calculus

3.3.1 Ito's lemma

It allows to find the differential of a time-dependent function of a stochastic process. Let $f(t, x)$ be a twice-differentiable scalar function. Its Taylor series gives

$$\Delta f = \frac{\partial f}{\partial t} \Delta t + \frac{\partial f}{\partial x} \Delta x + \frac{1}{2} \frac{\partial^2 f}{\partial t^2} \Delta t^2 + \frac{\partial^2 f}{\partial t \partial x} \Delta t \Delta x + \frac{1}{2} \frac{\partial^2 f}{\partial x^2} \Delta x^2$$

If the argument x of the function f is chosen to be a random process X_t satisfying *Ito drift-diffusion process*,

$$dX_t = \mu_t dt + \sigma_t dW_t,$$

the differential of function $f(t, X_t)$ results from the limit of Taylor series

$$\begin{aligned} df &= \lim_{dt \rightarrow 0, dW_t \rightarrow 0} \{\Delta f\} = \\ &= \lim_{dt \rightarrow 0, dW_t \rightarrow 0} \left\{ \partial_t f dt + \partial_x f dX_t + \frac{1}{2} [\partial_{tt} f dt^2 + 2\partial_{xt} f dt dX_t + \partial_{xx} f dX_t^2] \right\} = \\ &= \lim_{dt \rightarrow 0, dW_t \rightarrow 0} \left\{ \partial_t f dt + \partial_x f (\mu_t dt + \sigma_t dW_t) + \frac{1}{2} [\partial_{tt} f dt^2 + 2\partial_{xt} f dt (\mu_t dt + \sigma_t dW_t) + \partial_{xx} f (\mu_t dt + \sigma_t dW_t)^2] \right\} = \end{aligned}$$

For $dt \rightarrow 0$, $(dW_t)^2 = O(dt)$; keeping only terms of order lower than or equal to $O(dt)$, the differential becomes,

$$df = (\partial_t f + \mu_t \partial_x f) dt + \sigma_t \partial_x f dW_t + \frac{\sigma_t^2}{2} \partial_{xx} f dW_t^2 .$$

Replacing dW_t^2 with dt **todo** why?, and recalling the SDE of the Ito drift-diffusion process,

$$\begin{aligned} df &= \left(\partial_t f + \mu_t \partial_x f + \frac{\sigma_t^2}{2} \partial_{xx} f \right) dt + \sigma_t \partial_x f dW_t = \\ &= \left(\partial_t f + \frac{\sigma_t^2}{2} \partial_{xx} f \right) dt + \partial_x f (\mu dt + \sigma_t dW_t) = \\ &= \left(\partial_t f + \frac{\sigma_t^2}{2} \partial_{xx} f \right) dt + \partial_x f dX_t . \end{aligned}$$

3.3.2 Ito's calculus

Integration w.r.t. Brownian motion produces a random variable that can be defined as

$$\int_0^t F dW := \lim_{n \rightarrow +\infty} \sum_{[t_{i-1}, t_i] \in \pi_n} F_{t_{i-1}} (W_{t_i} - W_{t_{i-1}}) ,$$

being π_n a partition of interval $[0, t]$, and H a random process **todo** with some characteristics...

Example 3.3.1 (Integral of a Brownian motion w.r.t. itself)

$$Y(t) = \int_{s=0}^t W_s dW_s = \frac{1}{2} W_t^2 - \frac{t}{2} .$$

The expected value for each t of the random process Y_t is zero for all t , $\mathbb{E}[W_t^2] = 0$, as the expected value of W_t^2 is the variance of W_t , and thus t by definition of the Wiener process.

Evaluation of the integral

Let $f(t, x) = x^2$. Let's find the differential df evaluated for $x = W_t$ using *Ito's lemma*, retaining only terms with order up to $O(dt)$. Since $\partial_t f \equiv 0$,

$$df = \partial_x f|_{x=W_t} dW_t + \frac{1}{2} \partial_{xx} f|_{x=W_t} dW_t^2$$

and thus, replacing $dW_t^2 = dt$,

$$dW_t^2 = 2W_t dW_t + dt .$$

or

$$W_t dW_t = d\left(\frac{W_t^2}{2}\right) - \frac{dt}{2}.$$

Thus (**todo** add details if needed. A bit too much freedom in using differentials over stochastic processes here),

$$\begin{aligned} Y(t) &= \int_{s=0}^t W_s dW_s ds = \\ &= \int_{s=0}^t \left(\frac{W_s}{2}\right) ds - \int_{s=0}^t \frac{1}{2} ds = \\ &= \frac{1}{2} (W_t^2 - W_0^2) - \frac{t}{2}. \end{aligned}$$

3.3.3 Ito processes

Ito drift-diffusion process

An Ito drift-diffusion process is a stochastic process satisfying the stochastic differential equation (SDE)

$$dX_t = \mu_t dt + \sigma_t dW_t, \quad (3.4)$$

with W_t a Wiener process. If $\mu_t = \mu$, $\sigma_t = \sigma$ are constant a closed-form solution can be found using [Ito's lemma](#), for $f(t, x) = x$, or by direct (stochastic) integration of the SDE (3.4), as

$$\begin{aligned} \int_{s=0}^t dX_s &= \int_{s=0}^t \mu ds + \int_{s=0}^t \sigma dW_s \\ X_t - X_0 &= \mu t + \sigma (W_t - W_0), \end{aligned}$$

so that $X_t - X_0 \sim \mathcal{N}(\mu t, \sigma^2 t)$.

Scaling of a Wiener process

Term σW_t represents a scaling of a Wiener process $W_t \sim \mathcal{N}(0, t)$ with zero expected value and variance t . Multiplication by factor σ results in a multiplication of the expected value by σ and variance by σ^2 .

Geometric Brownian Motion

A geometric Brownian motion is a stochastic process satisfying the SDE

$$dX_t = \mu X_t dt + \sigma X_t dW_t.$$

Let $f(x) = \ln x$ be evaluated for $x = X_t$. Ito's lemma, with $\partial_t f \equiv 0$, provides the expression of the differential

$$\begin{aligned} df &= \partial_x f|_{X_t} dX_t + \frac{1}{2} \partial_{xx} f|_{X_t} dX_t^2 = \\ &= \partial_x f|_{X_t} (\mu X_t dt + \sigma X_t dW_t) + \frac{1}{2} \partial_{xx} f|_{X_t} (\mu X_t dt + \sigma X_t dW_t)^2 = \\ &= \frac{1}{X_t} (\mu X_t dt + \sigma X_t dW_t) - \frac{1}{2} \frac{1}{X_t^2} \sigma^2 X_t^2 dW_t^2 = \\ d(\ln X_t) &= \left(\mu - \frac{\sigma^2}{2} \right) dt + \sigma dW_t, \end{aligned}$$

whose solution after integration reads

$$\ln X_t = \ln X_0 + \left(\mu - \frac{\sigma^2}{2} \right) t + \sigma W_t ,$$

or

$$X_t = X_0 e^{\left(\mu - \frac{\sigma^2}{2} \right) t + \sigma W_t} .$$

Part III

Statistical Inference

INTRODUCTION TO STATISTICAL INFERENCE

Part IV

Introduction to Machine Learning

INTRODUCTION TO MACHINE LEARNING

Artificial intelligence can be broadly defined as a field dealing with making machines perform tasks that require intelligence, when performed by humans, like: reasoning, perception, representation, language processing, planning, learning

Machine learning is a branch of AI focused on statistical algorithms that can **learn from data** and **generalize to unseen data** and perform tasks, without explicit instructions.¹

Three core paradigms. Algorithms in machine learning can be divided into three paradigms:

- **Supervised Learning, SL:** algorithm learns from labelled data; many applications can be reduced to 2 main tasks: **regression** (or function approximation) and **classification**.
- **Unsupervised Learning, UL:** algorithm learns pattern from un-labelled data; examples of tasks in UL are clustering, dimensionality reduction (and recognition of *main* components in data), compression (retaining only relevant components in data). Some historical algorithms and linear algebra decompositions can be interpreted or generalized as unsupervised learning.
- **Reinforcement Learning, RL:** an algorithm (**agent**) learns a **policy** - i.e. the way to behave - interacting with an **environment**, and maximizing some performance to efficiently perform required tasks. Applications of RL includes **planning** and **control**.

Goals and methodology. ML is mainly an engineering-oriented and an application-focused discipline, relying on statistical inference (**todo be more explicit**). A **ML model** usually takes an input **u**, and produces an output **y**, depending on its own structure and a set of parameters θ and hyper-parameters μ . Learning usually relies on **optimization** of an objective function

$$L(\theta; \mu) ,$$

w.r.t. parameters θ , whose value is learned/adjusted towards an optimal solution θ^* that makes $L(\theta^*; \mu)$ extreme. The choice of hyper-parameters μ instead influences the training process and model behavior. Optimization usually relies on gradient methods, updating the parameters in the direction of the gradient of the objective function w.r.t. the parameters,

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} L(\theta; \mu) .$$

Optimization of model parameters is made fast by the use of **back-propagation** and **automatic differentiation** (AD), which efficiently compute gradients of the cost function with respect to the model's parameters, and technically feasible for large-dimensional models - as the ones used in multi-layered neural networks, in deep learning² - by recent hardware improvement. These algorithms are not only feasible but also particularly well-suited (being a major driver for new designs) to modern processing architectures, such as **GPUs** and **TPUs**, that accelerate the large-scale matrix and tensor computations involved in both the forward and backward passes of training.

todo Show NVIDIA, TSMC revenues

¹ "Without explicit instructions" means that a system has no user-coded behavior, but learns it usually via **optimization**, usually either involving minimization of an error function or maximization of an objective function or energy/information content.

² Deep learning can be roughly defined as that branch of machine learning using multi-layered neural networks, indeed.

todo Add references: Bishop,...

5.1 Models in Machine Learning

Linear models

Kernel methods

Decision trees and ensembles

Neural networks

...probabilistic models, clustering models, dimensionality reduction models,...

Reinforcement learning models: Q-learning (tabular, and DQN), Policy gradient, Actor-Critic, Proximal Policy Optimization,...

5.2 Good Practices in Machine Learning

SUPERVISED LEARNING

Theory.

Examples.

6.1 SL: theory

Supervised learning can be thought as a **function approximation problem**. Given a set of data

$$\{(x_i, y_i)\}_{i=1:N} ,$$

supervised learning can be formulated as the evaluation of a function $\hat{y}(x, \theta)$, or a **model**, that *approximates well* the relation between input x_i and output y_i ,

$$y_i \simeq \hat{y}(x_i; \theta) .$$

Two main tasks of SL can be distinguished on the output of the function: **regression** can be formulated as function approximation with continuous output, while in **classification** the function maps inputs to discrete output/**labels**

Learning process aims at finding values of the parameters θ (and hyper-parameters μ), that minimize a “prediction” error function, e.g. for a scalar output function,

$$E(\theta) = \frac{1}{2} \sum_{i \in D_{Tr}} |\hat{y}(x_i; \theta) - y_i|^2 ,$$

being D_{Tr} the set of indices belonging to the *training set*. Minimization usually relies on gradient methods of the error function w.r.t. the parameters θ ,

$$\begin{aligned} \nabla_{\theta} E(\theta, \mathbf{x}_{Tr}, \mathbf{y}_{Tr}) &= \sum_{i \in D_{Tr}} (\hat{y}(x_i; \theta) - y_i) \nabla_{\theta} \hat{y}(x_i; \theta) \\ \theta &\leftarrow \theta - \alpha \nabla_{\theta} E(\theta, \mathbf{x}_{Tr}, \mathbf{y}_{Tr}) , \end{aligned}$$

with α an hyper-parameter called *learning rate*, governing the “length” of the update step. Other objective functions to be maximised or minimized can be used. Slight variations to objective functions allow for regularization (e.g. parameter weighting)

Dataset. Available data $\{x_i, y_i\}_i$ is divided in different sets:

- training set: for learning/tuning model parameters, minimizing an error function
- validation set: for early stopping, and hyper-parameter tuning (e.g. to avoid
- test set: to evaluate model performance

UNSUPERVISED LEARNING

REINFORCEMENT LEARNING

PROOF INDEX

definition-0

definition-0 (*ch/prob/iid*), 16

definition-2

definition-2 (*ch/prob/rv-multi-dimensional*), 13

ex:char-fun:independent

ex:char-fun:independent
(*ch/prob/characteristic-fun*), 14

ex:char-fun:multidimensional

ex:char-fun:multidimensional
(*ch/prob/characteristic-fun*), 14

ex:char-fun:normal

ex:char-fun:normal (*ch/prob/characteristic-fun*),
15

ex:char-fun:taylor

ex:char-fun:taylor (*ch/prob/characteristic-fun*),
15

example-0

example-0 (*ch/prob/processes-calculus*), 24

property-0

property-0 (*ch/prob/rv-multi-dimensional*), 12

property-1

property-1 (*ch/prob/wiener*), 20

property-2

property-2 (*ch/prob/wiener*), 21

property-3

property-3 (*ch/prob/wiener*), 22

theorem-1

theorem-1 (*ch/prob/rv-multi-dimensional*), 13

thm:clt

thm:clt (*ch/prob/iid*), 17

wiener:def

wiener:def (*ch/prob/wiener*), 20

wn:def:derivative

wn:def:derivative (*ch/prob/white-noise*), 23

wn:def:properties

wn:def:properties (*ch/prob/white-noise*), 23