

# Probabilistic Artificial Intelligence

## Gaussian Processes (GP)

Kernalizd Bayesian Linear Regression

1. Gaussian Process Regression
2. Model Selection



Proofs

*Proof.* Property 2.3The kernel matrix is positive-semidefinite:  
Let  $\phi : \mathcal{X} \mapsto \mathbb{R}^d$  and  $\Phi = \begin{bmatrix} \phi(\mathbf{x}_1) & \dots & \phi(\mathbf{x}_n) \end{bmatrix}^\top \in \mathbb{R}^{d \times n}$ .  
Thus:  $\mathcal{K} = \Phi^\top \Phi \in \mathbb{R}^{n \times n}$ .  
 $\mathbf{v}^\top \mathcal{K} \mathbf{v} = \mathbf{v}^\top \Phi^\top \Phi \mathbf{v} = (\Phi \mathbf{v})^\top \Phi \mathbf{v} = \|\Phi \mathbf{v}\|_2^2 \geq 0$

□

Examples

**Example 2.1 Calculating the Kernel by hand:**

Let :  $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$   $\phi(\mathbf{x}) \mapsto \{x_1^2, x_2^2, \sqrt{2}x_1, x_2\}$   
 $\phi : \mathbb{R}^{d=2} \mapsto \mathbb{R}^{\bar{D}=3}$

We can now have a decision boundary in this 3-D feature space  $\mathcal{Y}$  of  $\phi$  as:

$$\beta_0 + \beta_1 x_1^2 + \beta_2 x_2^2 + \beta_3 \sqrt{2} x_1 x_2 = 0$$
$$\begin{aligned} &\langle \phi(\mathbf{x}^{(i)}), \phi(\mathbf{x}^{(j)}) \rangle \\ &= \langle \{x_{i1}^2, x_{i2}^2, \sqrt{2}x_{i1}, x_{i2}\}, \{x_{j1}^2, x_{j2}^2, \sqrt{2}x_{j1}, x_{j2}\} \rangle \\ &= x_{i1}^2 x_{j1}^2 + x_{i2}^2 x_{j2}^2 + 2x_{i1} x_{i2} x_{j1} x_{j2} \end{aligned}$$

**Operation Count:**

- 2 · 3 operations to map  $\mathbf{x}_i$  and  $\mathbf{x}_j$  into the 3D space  $\mathcal{Y}$ .
- Calculating an inner product of  $\langle \phi(\mathbf{x}^{(i)}), \phi(\mathbf{x}^{(j)}) \rangle$  with 3 additional operations.

**Example 2.2**  
**Calculating the Kernel using the Kernel Trick:**

$$\begin{aligned} \langle \phi(\mathbf{x}^{(i)}), \phi(\mathbf{x}^{(j)}) \rangle &= \underbrace{\langle \mathbf{x}_i, \mathbf{x}_j \rangle}^2 = \langle \{x_{i1}, x_{i2}\}, \{x_{i1}, x_{i2}\} \rangle^2 \\ &:= \mathbf{k}(\mathbf{x}_i, \mathbf{x}_j) \\ &= (x_{i1} x_{i2} + x_{j1} x_{j2})^2 \\ &= x_{i1}^2 x_{j1}^2 + x_{i2}^2 x_{j2}^2 + 2x_{i1} x_{i2} x_{j1} x_{j2} \end{aligned}$$

**Operation Count:**

- 2 multiplicaitons of  $\mathbf{x}_{i1} x_{j1}$  and  $\mathbf{x}_{i2} x_{j2}$ .
- 1 operation for taking the square of a scalar.

**Conclusion** The Kernel trick needed only 3 in comparison to 9 operations.

**Example 2.3 Stationary Kernels:**

$$\mathbf{k}(\mathbf{x}, \mathbf{y}) = \exp \left( \frac{(\mathbf{x} - \mathbf{y})^\top M (\mathbf{x} - \mathbf{y})}{h^2} \right)$$

is a stationary but not an isotropic kernel.

Math Appendix

Logic  
Set Theory

<b>Definition 4.1 Set</b> $A = \{1, 3, 2\}$ : is a well-defined group of distinct items that are considered as an object in its own right. The arrangement/order of the objects does not matter but each member of the set must be unique.
<b>Definition 4.2 Empty Set</b> $\{\}$ / $\emptyset$ : is the unique set having no elements/cardinality <sup>[def. 4.4]</sup> zero.
<b>Definition 4.3 Multiset/Bag</b> : Is a set-like object in which multiplicity matters, that is we can have multiple elements of the same type. I.e. $\{1, 1, 2, 3\} \neq \{1, 2, 3\}$
<b>Definition 4.4 Cardinality</b> $ S $ : Is the number of elements that are contained in a set.
<b>Definition 4.5 The Power Set</b> $\mathcal{P}(S)/2^S$ : The power set of any set $S$ is the set of all subsets of $S$ , including the empty set and $S$ itself. The cardinality of the power set is $2^S$ is equal to $2^{ S }$ .
<b>Definition 4.6 Closure</b> : A set is <i>closed</i> under an operation $\Omega$ if performance of that operations onto members of the set always produces a member of that set.

1. Number Sets  $\mathbb{R}$   
1.1. The Real Numbers  
1.1.1. Intervals

<b>Definition 4.7 Closed Interval</b> $[a, b]$ : The closed interval of $a$ and $b$ is the set of all real numbers that are within $a$ and $b$ , including $a$ and $b$ : $[a, b] = \{x \in \mathbb{R} \mid a \leq x \leq b\}$ (4.1)
<b>Definition 4.8 Open Interval</b> $(a, b)$ : The open interval of $a$ and $b$ is the set of all real numbers that are within $a$ and $b$ : $(a, b) = \{x \in \mathbb{R} \mid a < x < b\}$ (4.2)

1.2. The Rational Numbers  $\mathbb{Q}$

<b>Example 4.1 Power Set/Cardinality of</b> $S = \{x, y, z\}$ : The subsets of $S$ are: $\{\emptyset\}, \{x\}, \{y\}, \{z\}, \{x, y\}, \{x, z\}, \{y, z\}, \{x, y, z\}$ and hence the power set of $S$ is $\mathcal{P}(S) = \{\{\emptyset\}, \{x\}, \{y\}, \{z\}, \{x, y\}, \{x, z\}, \{y, z\}, \{x, y, z\}\}$ with a cardinality of $ S  = 2^3 = 8$ .
---

Sequences&Series

<b>Definition 5.1 Index Set</b> : Is a set <sup>[def. 4.1]</sup> $A$ , whose members are labels to another set $S$ . In other words its members index member of another set. An index set is build by enumerating the members of $S$ using a function $f$ s.t. $f : A \mapsto S \qquad A \in \mathbb{N}$ (5.1)
<b>Definition 5.2 Sequence</b> $(a_n)_{n \in A}$ : is an by an index set $A$ <i>enumerated</i> multiset <sup>[def. 4.3]</sup> (repetitions are allowed) of objects in which <i>order does matter</i> .
<b>Definition 5.3 Series</b> : is an infinite ordered set of terms combined together by addition.

1. Types of Sequences  
1.1. Arithmetic Sequence

<b>Definition 5.4 Arithmetic Sequence</b> : Is a sequence where the <i>difference</i> between two consecutive terms constant i.e. $(2, 4, 6, 8, 10, 12, \dots)$ . $t_n = t_0 + nd \qquad d$ :difference between two terms (5.2)
--

1.2. Geometric Sequence

<b>Definition 5.5 Geometric Sequence</b> : Is a sequence where the <i>ratio</i> between two consecutive terms constant i.e. $(2, 4, 8, 16, 32, \dots)$ . $t_n = t_0 \cdot r^n \qquad r$ :ratio between two terms (5.3)
---

# Calculus and Analysis

## 1. Building Blocks of Analysis

### 1.1. Polynomials

**Definition 6.1 Polynomial:** A function  $\mathcal{P}_n : \mathbb{R} \rightarrow \mathbb{R}$  is called *Polynomial*, if it can be represented in the form:

$$\mathcal{P}_n(x) = a_0 + a_1x + a_2x^2 + \cdots + a_{n-1}x^{n-1} + a_nx^n \quad (6.1)$$

**Corollary 6.1 Degree n-of a Polynomial  $\deg(\mathcal{P}_n)$ :** the *degree* of the polynomial is the highest exponent of the variable  $x$ , among all non-zero coefficients  $a_i \neq 0$ .

**Definition 6.2 Monomial:** Is a polynomial with only one term.

**Definition 6.3 Quadratic Formula:**  $ax^2 + bx + c = 0$  or in reduced form:  
 $x^2 + px + q = 0$  with  $p = b/a$  and  $q = c/a$

**Definition 6.4 Discriminant:**  $\delta = b^2 - 4ac$

**Definition 6.5 Solution to** <sup>[def. 6.3]</sup>:  
$$x_{\pm} = \frac{-b \pm \sqrt{\delta}}{2a} \quad \text{or} \quad x_{\pm} = \frac{1}{2} \left( -p \pm \sqrt{p^2 - 4q} \right)$$

**Theorem 6.1**  
**Fist Fundamental Theorem of Calculus:** Let  $f$  be a continuous real-valued function defined on a closed interval  $[a, b]$ . Let  $F$  be the function defined  $\forall x \in [a, b]$  by:

$$F(X) = \int_a^x f(t) dt \quad (6.2)$$

Then it follows:  
$$F'(x) = f(x) \quad \forall x \in (a, b) \quad (6.3)$$

**Theorem 6.2**  
**Second Fundamental Theorem of Calculus:** Let  $f$  be a real-valued function on a closed interval  $[a, b]$  and  $F$  an antiderivative of  $f$  in  $[a, b]$ :  $F'(x) = f(x)$ , then it follows if  $f$  is Riemann integrable on  $[a, b]$ :

$$\int_a^b f(t) dt = F(b) - F(a) \iff \int_a^x \frac{\partial}{\partial x} F(t) dt = F(x) \quad (6.4)$$

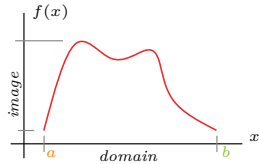
**Definition 6.6 Domain of a function  $\text{dom}(\cdot)$ :**  
**Given** a function  $f : \mathcal{X} \rightarrow \mathcal{Y}$ , the set of all possible input values  $\mathcal{X}$  is called the domain of  $f - \text{dom}(f)$ .

**Definition 6.7**  
**Codomain/target set of a function  $\text{codom}(\cdot)$ :**  
**Given** a function  $f : \mathcal{X} \rightarrow \mathcal{Y}$ , the codaomain of that function is the set  $\mathcal{Y}$  into which all of the output of the function is constrained to fall.

**Definition 6.8 Image (Range) of a function:**  $f[\cdot]$

**Given** a function  $f : \mathcal{X} \rightarrow \mathcal{Y}$ , the image of that function is the set to which the function can actually map:  
$$\{y \in \mathcal{Y} | y = f(x), \quad \forall x \in \mathcal{X}\} := f[\mathcal{X}] \quad (6.5)$$

Evaluating the function  $f$  at each element of a given subset  $A$  of its domain  $\text{dom}(f)$  produces a set called the *image* of  $A$  under (or through)  $f$ . The image is thus a subset of a function's codomain.



**Definition 6.9 Inverse Image/Preimage  $f^{-1}(\cdot)$ :**  
Let  $f : X \mapsto Y$  be a function, and  $A$  a subset set of its codomain  $Y$ .  
Then the preimage of  $A$  under  $f$  is the set of all elements of the domain  $X$ , that map to elements in  $A$  under  $f$ :  
$$f^{-1}(A) = \{x \subseteq X : f(x) \subseteq A\} \quad (6.6)$$

**Example 6.1 :**  
**Given**  $f : \mathbb{R} \rightarrow \mathbb{R}$   
defined by  $f : x \mapsto x^2 \iff f(x) = x^2$   
 $\text{dom}(f) = \mathbb{R}$ ,  $\text{codom}(f) = \mathbb{R}$  but its image is  $f[\mathbb{R}] = \mathbb{R}_+$ .

**Image (Range) of a subset**

The image of a subset  $A \subseteq \mathcal{X}$  under  $f$  is the subset  $f[A] \subseteq \mathcal{Y}$  defined by:  
$$f[A] = \{y \in \mathcal{Y} | y = f(x), \quad \forall x \in A\} \quad (6.7)$$

**Note: Range**

The term range is ambiguous as it may refer to the image or the codomain, depending on the definition. However, modern usage almost always uses range to mean image.

**Definition 6.10 (strictly) Increasing Functions:**  
A function  $f$  is called **monotonically increasing/ increasing/non-decreasing** if:  
$$x \leq y \iff f(x) \leq f(y) \quad \forall x, y \in \text{dom}(f) \quad (6.8)$$
  
And **strictly increasing** if:  
$$x < y \iff f(x) < f(y) \quad \forall x, y \in \text{dom}(f) \quad (6.9)$$

**Definition 6.11 (strictly) Decreasing Functions:**  
A function  $f$  is called monotonically decreasing/decreasing or non-increasing if:  
$$x \geq y \iff f(x) \geq f(y) \quad \forall x, y \in \text{dom}(f) \quad (6.10)$$
  
And **strictly decreasing** if:  
$$x > y \iff f(x) > f(y) \quad \forall x, y \in \text{dom}(f) \quad (6.11)$$

**Definition 6.12 Monotonic Function:** A function  $f$  is called monotonic iff either  $f$  is **increasing** or **decreasing**.

**Definition 6.13 Linear Function:**  
A function  $L : \mathbb{R}^n \mapsto \mathbb{R}^m$  is linear if and only if:  
$$L(\mathbf{x} + \mathbf{y}) = L(\mathbf{x}) + L(\mathbf{y})$$
$$L(\alpha \mathbf{x}) = \alpha L(\mathbf{x}) \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, \quad \alpha \in \mathbb{R}$$

**Corollary 6.2 Linearity of Differentiation:** The derivative of **any** linear combination of functions equals the same linear combination of the derivatives of the functions:  
$$\frac{d}{dx} (af(x) + bg(x)) = a \frac{d}{dx} f(x) + b \frac{d}{dx} g(x) \quad a, b \in \mathbb{R} \quad (6.12)$$

**Definition 6.14 Quadratic Function:**  
A function  $f : \mathbb{R}^n \mapsto \mathbb{R}^m$  is quadratic if it can be written in the form:  
$$f(x) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} + c \quad (6.13)$$

## 2. Continuity and Smoothness

**Definition 6.15 Continuous Function:**

**Definition 6.16 Smoothness of a Function  $\mathcal{C}^k$ :** **Given** a function  $f : \mathcal{X} \rightarrow \mathcal{Y}$ , the function is said to be of class  $k$  if it is differentiable up to order  $k$  **and** continuous, on its entire domain:  
$$f \in \mathcal{C}^k(\mathcal{X}) \iff \exists f', f'', \dots, f^{(k)} \text{ continuous} \quad (6.14)$$

**Note**

- The class  $\mathcal{C}^0$  consists of all continuous functions.
- P.w. continuous  $\neq$  continuous.
- A function of that is  $k$  times differentiable must at least be of class  $\mathcal{C}^{k-1}$ .
- $\mathcal{C}^m(\mathcal{X}) \subset \mathcal{C}^{m-1}, \dots \mathcal{C}^1 \subset \mathcal{C}^0$
- Continuity is implied by the differentiability of all **derivatives** of up to order  $k - 1$ .

**Corollary 6.3 Smooth Function  $\mathcal{C}^\infty$ :** Is a function  $f : \mathcal{X} \rightarrow \mathcal{Y}$  that has derivatives infinitely many times differentiable.  
$$f \in \mathcal{C}^\infty(\mathcal{X}) \iff f', f'', \dots, f^{(\infty)} \quad (6.15)$$

**Corollary 6.4 Continuously Differentiable Function  $\mathcal{C}^1$ :** Is the class of functions that consists of all differentiable functions whose derivative is continuous.  
Hence a function  $f : \mathcal{X} \rightarrow \mathcal{Y}$  of the class must satisfy:  
$$f \in \mathcal{C}^1(\mathcal{X}) \iff f' \text{ continuous} \quad (6.16)$$

Often functions are not differentiable but we still want to state something about the rate of change of a function  $\Rightarrow$  hence we need a weaker notion of differentiability.

**Definition 6.17 Lipschitz Continuity:** A Lipschitz continuous function is a function  $f$  whose rate of change is bound by a Lipschitz Contant  $L$ :  
$$|f(x) - f(y)| \leq L \|\mathbf{x} - \mathbf{y}\|_2^2 \quad \forall \mathbf{x}, \mathbf{y}, \quad L > 0 \quad (6.17)$$

**Note**

This property is useful as it allows us to conclude that a small perturbation in the input (i.e. of an algorithm) will result in small changes of the output  $\Rightarrow$  tells us something about robustness.

**Definition 6.18 Lipschitz Continuous Gradient:**  
A *continuously differentiable* function  $f : \mathbb{R}^d \mapsto \mathbb{R}$  has  $L$ -Lipschitz continuous gradient if it satisfies:  
$$\|\nabla f(x) - \nabla f(y)\| \leq L \|\mathbf{x} - \mathbf{y}\| \quad \forall \mathbf{x}, \mathbf{y} \in \text{dom}(f), \quad L > 0 \quad (6.18)$$

if  $f \in \mathcal{C}^2$ , this is equivalent to:  
$$\nabla^2 f(\mathbf{x}) \leq L \mathbf{I} \quad \forall \mathbf{x} \in \text{dom}(f), \quad L > 0 \quad (6.19)$$

**Lemma 6.1 Descent Lemma:** If a function  $f : \mathbb{R}^d \mapsto \mathbb{R}$  has *Lipschitz continuous gradient* eq. (6.18) over its domain, then it holds that:  
$$|f(x) - f(y) - \nabla f(y)^T(\mathbf{x} - \mathbf{y})| \leq \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2 \quad (6.20)$$

**Note**

If  $f$  is twice differentiable then the largest eigenvalue of the Hessian <sup>[def. 7.5]</sup> of  $f$  is uniformly upper bounded by  $L$

*Proof.* lemma 6.1 for  $\mathcal{C}^1$  functions:  
Let  $g(t) \equiv f(\mathbf{y} + t(\mathbf{x} - \mathbf{y}))$  from the FToC (theorem 6.2) we know that:

$$\begin{aligned} & \int_0^1 g'(t) dt = g(1) - g(0) = f(\mathbf{x}) - f(\mathbf{y}) \\ & \text{It then follows from the reverse:} \\ & |f(\mathbf{x}) - f(\mathbf{y}) - \nabla f(\mathbf{y})^T(\mathbf{x} - \mathbf{y})| \\ & \stackrel{\text{Chain. R}}{\stackrel{\text{FToC}}{=}} \left| \int_0^1 \nabla f(\mathbf{y} + t(\mathbf{x} - \mathbf{y}))^T(\mathbf{x} - \mathbf{y}) dt - \nabla f(\mathbf{y})^T(\mathbf{x} - \mathbf{y}) \right| \\ & = \left| \int_0^1 (\nabla f(\mathbf{y} + t(\mathbf{x} - \mathbf{y})) - \nabla f(\mathbf{y}))^T(\mathbf{x} - \mathbf{y}) dt \right| \\ & = \left| \int_0^1 (\nabla f(\mathbf{y} + t(\mathbf{x} - \mathbf{y})) - \nabla f(\mathbf{y}))^T(\mathbf{x} - \mathbf{y}) dt \right| \\ & \stackrel{\text{C.S.}}{\leq} \left| \int_0^1 \|\nabla f(\mathbf{y} + t(\mathbf{x} - \mathbf{y})) - \nabla f(\mathbf{y})\| \cdot \|\mathbf{x} - \mathbf{y}\| dt \right| \\ & \stackrel{\text{eq. (6.18)}}{=} \left| \int_0^1 L \|\mathbf{y} + t(\mathbf{x} - \mathbf{y}) - \mathbf{y}\| \cdot \|\mathbf{x} - \mathbf{y}\| dt \right| \\ & = \left| L \|\mathbf{x} - \mathbf{y}\|^2 \int_0^1 t dt \right| = \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2_2 \end{aligned}$$

*Proof.* lemma 6.1 for  $\mathcal{C}^2$  functions:

$$f(\mathbf{y}) \stackrel{\text{Taylor}}{=} f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{1}{2}(\mathbf{y} - \mathbf{x})^T \nabla^2 f(z)(\mathbf{y} - \mathbf{x})$$

Now we plug in  $\nabla^2 f(\mathbf{x})$  and recover eq. (6.21):

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{1}{2}(\mathbf{y} - \mathbf{x})^T L(\mathbf{y} - \mathbf{x})$$

□

**Definition 6.19 L-Smoothness:** A  $L$ -smooth function is a function  $f : \mathbb{R}^d \mapsto \mathbb{R}$  that satisfies:

$$f(\mathbf{x}) \leq f(\mathbf{y}) + \nabla f(\mathbf{y})^T(\mathbf{x} - \mathbf{y}) + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2$$

with  $\forall \mathbf{x}, \mathbf{y} \in \text{dom}(f), \quad L > 0 \quad (6.21)$

If  $f$  is a twice differentiable this is equivalent to:  
$$\nabla^2 f(\mathbf{x}) \leq L \mathbf{I} \quad L > 0 \quad (6.22)$$

**Theorem 6.3**  
**L-Smoothness of convex functions:** A *convex* and  $L$ -Smooth function <sup>[def. 6.19]</sup> has a Lipschitz continuous gradient (eq. (6.18)) thus it holds that:

$$f(\mathbf{x}) \leq f(\mathbf{y}) + \nabla f(\mathbf{y})^T(\mathbf{x} - \mathbf{y}) \leq \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2 \quad (6.23)$$

*Proof.* theorem 6.3:

With the definition of convexity for a differentiable function (eq. (6.26)) it follows

$$\begin{aligned} f(\mathbf{x}) - f(\mathbf{y}) + \nabla f(\mathbf{y})^T(\mathbf{x} - \mathbf{y}) & \geq 0 \\ \Rightarrow |f(\mathbf{x}) - f(\mathbf{y}) + \nabla f(\mathbf{y})^T(\mathbf{x} - \mathbf{y})| \\ & \stackrel{\text{if eq. (6.26)}}{=} f(\mathbf{x}) - f(\mathbf{y}) + \nabla f(\mathbf{y})^T(\mathbf{x} - \mathbf{y}) \end{aligned}$$

with lemma 6.1 and <sup>[def. 6.19]</sup> it follows theorem 6.3 □

**Corollary 6.5 :**  $L$ -smoothnes is a weaker condition than  $L$ -Lipschitz continuous gradients

## 3. Convexity

Read stuff about uniqueness and so on again in NPDE/or NUM CSE and add proofs

**Definition 6.20 Convex Functions:**  
A function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is convex if it satisfies:  
$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}) \quad \forall \mathbf{x}, \mathbf{y} \in \text{dom}(f) \quad \forall \lambda \in [0, 1] \quad (6.24)$$

include figure from tika/convexity

**Definition 6.21 Concave Functions:**  
A function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is convex if it satisfies:  
$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \geq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}) \quad \forall \mathbf{x}, \mathbf{y} \in \text{dom}(f) \quad \forall \lambda \in [0, 1] \quad (6.25)$$

**Corollary 6.6 Convexity  $\rightarrow$  global minimima:** Convexity implies that all local minima (if they exist) are global minima.

**Definition 6.22 Stricly Convex Functions:**  
A function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is **strictly** convex if it satisfies:  
$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) < \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}) \quad \forall \mathbf{x}, \mathbf{y} \in \text{dom}(f) \quad \forall \lambda \in [0, 1]$$

add plot

If  $f$  is a differentiable function this is equivalent to:  
$$f(\mathbf{x}) \geq f(\mathbf{y}) + \nabla f(\mathbf{y})^T(\mathbf{x} - \mathbf{y}) \quad \forall \mathbf{x}, \mathbf{y} \in \text{dom}(f) \quad (6.26)$$

If  $f$  is a twice differentiable function this is equivalent to:  
$$\nabla^2 f(\mathbf{x}) \geq 0 \quad \forall \mathbf{x}, \mathbf{y} \in \text{dom}(f) \quad (6.27)$$

**Intuition**

- Convexity implies that a function  $f$  is bound by/below a linear interpolation from  $x$  to  $y$  and strong convexity that  $f$  is strictly bound/below.
- eq. (6.26) implies that  $f(\mathbf{x})$  is above the tangent  $f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x})$  for all  $\mathbf{x}, \mathbf{y} \in \text{dom}(f)$
- ?? implies that  $f(\mathbf{x})$  is flat or curved upwards

**Corollary 6.7 Strict Convexity  $\rightarrow$  Uniqueness:**  
 Strict convexity implies a unique minimizer  $\iff$  at most one global minimum.

**Corollary 6.8 :** A twice differentiable function of one variable  $f : \mathbb{R} \rightarrow \mathbb{R}$  is **convex** on an interval  $\mathcal{X} = [a, b]$  if and only if its second derivative is non-negative on that interval  $\mathcal{X}$ :  

$$f''(x) \geq 0 \quad \forall x \in \mathcal{X} \quad (6.28)$$

**Definition 6.23  $\mu$ -Strong Convexity:**  
 Let  $\mathcal{X}$  be a Banach space over  $\mathbb{K} = \mathbb{R}, \mathbb{C}$ . A function  $f : \mathcal{X} \rightarrow \mathbb{R}$  is called strongly convex iff the following equation holds:  

$$f(tx + (1-t)y) \leq tf(x) + (1-t)f(y) - \frac{t(1-t)}{2} \mu \|x - y\| \quad \forall x, y \in \mathcal{X}, \quad t \in [0, 1], \quad \mu > 0$$

If  $f \in \mathcal{C}^1 \iff f$  is differentiable, this is equivalent to:  

$$f(y) \geq f(x) + \nabla f(x)^\top (y - x) + \frac{\mu}{2} \|y - x\|_2^2 \quad (6.29)$$

If  $f \in \mathcal{C}^2 \iff f$  is twice differentiable, this is equivalent to:  

$$\nabla^2 f(x) \geq \mu \mathbf{I} \quad \forall x, y \in \mathcal{X} \quad \mu > 0 \quad (6.30)$$

**Corollary 6.9 Strong Convexity implies Strict Convexity:**  
<https://math.stackexchange.com/question/2090991/proof-for-strongly-convex-function-is-strictly-convex>

**Property 6.1:**  

$$f(y) \leq f(y) + \nabla f(y)^\top (x - y) + \frac{1}{2\mu} \|\nabla f(x) - \nabla f(y)\|_2^2 \quad (6.31)$$

**Intuition**  
 Strong convexity implies that a function  $f$  is lower bounded by its second order (quadratic) approximation, rather then only its first order (linear) approximation.

**Size of  $\mu$**   
 The parameter  $\mu$  specifies how strongly the bounding quadratic function/approximation is.

*Proof.* eq. (6.30) analogously to **Proof** eq. (6.22)  $\square$

**Note**  
 If  $f$  is twice differentiable then the smallest eigenvalue of the Hessian <sup>([def. 7.51](#))</sup> of  $f$  is uniformly lower bounded by  $\mu$ .  
**Hence** strong convexity can be considered as the analogous to smoothness

**Example 6.2 Quadratic Function:** A quadratic function eq. (6.13) is convex if:  

$$\nabla_{\mathbf{x}}^2 \text{eq. (6.13)} = \mathbf{A} \geq 0 \quad (6.32)$$

**Corollary 6.10 :**  
 Strong convexity  $\implies$  Strict convexity  $\implies$  Convexity

### 3.1. Properties that preserve convexity

**Property 6.2 Non-negative weighted Sums:** Let  $f$  be a convex function then  $g(x)$  is convex as well:  

$$g(x) = \sum_{i=1}^n \alpha_i f_i(x) \quad \forall \alpha_j > 0$$

**Property 6.3 Composition of Affine Mappings:** Let  $f$  be a convex function then  $g(x)$  is convex as well:  

$$g(x) = f(\mathbf{Ax} + \mathbf{b})$$

**Property 6.4 Pointwise Maxima:** Let  $f$  be a convex function then  $g(x)$  is convex as well:  

$$g(x) = \max_i \{f_i(x)\}$$

### Functions

**Even Functions:** have rotational symmetry with respect to the origin.  
 $\implies$  **Geometrically:** its graph remains unchanged after reflection about the y-axis.

$$f(-x) = f(x) \quad (6.33)$$

**Odd Functions:** are symmetric w.r.t. to the y-axis.  
 $\implies$  **Geometrically:** its graph remains unchanged after rotation of 180 degrees about the origin.

$$f(-x) = -f(x) \quad (6.34)$$

**Theorem 6.4 Rules:**  
**Let  $f$  be even and  $f$  odd respectively.**  
 $g =: f \cdot f$  is even  $g =: f \cdot f$  is even  
 $g =: f \cdot f$  is odd the same holds for division

**Examples**  
**Even:**  $\cos x, |x|, \mathbf{c}, x^2, x^4, \dots \exp(-x^2/2)$ .  
**Odd:**  $\sin x, \tan x, x, x^3, x^5, \dots$

**x-Shift:**  $f(x - \mathbf{c}) \Rightarrow$  shift to the right  
 $f(x + \mathbf{c}) \Rightarrow$  shift to the left (6.35)  
**y-Shift:**  $f(x) \pm \mathbf{c} \Rightarrow$  shift up/down (6.36)

*Proof.* eq. (6.35)  $f(x_n - \mathbf{c})$  we take the x-value at  $x_n$  but take the y-value at  $x_o := x_n - \mathbf{c} \implies$  we shift the function to  $x_n$ .  $\square$

**Euler's formula**  

$$e^{\pm ix} = \cos x \pm i \sin x \quad (6.37)$$

**Euler's Identity**  

$$e^{\pm i} = -1 \quad (6.38)$$

**Note**  

$$e^n = 1 \Leftrightarrow n = i2\pi k, \quad k \in \mathbb{N} \quad (6.39)$$

**Corollary 6.11 Every norm is a convex function:** By using definition <sup>([def. 6.20](#))</sup> and the triangular inequality it follows (with the exception of the L0-norm):  

$$\|\lambda x + (1 - \lambda)y\| \leq \lambda \|x\| + (1 - \lambda) \|y\|$$

### 3.2. Taylor Expansion

**Definition 6.24 Taylor Expansion:**  

$$T_n(x) = \sum_{i=0}^n \frac{1}{n!} f^{(i)}(x_0) \cdot (x - x_0)^{(i)} \quad (6.40)$$

$$= f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2} f''(x_0)(x - x_0)^2 + \mathcal{O}(x^3) \quad (6.41)$$

**Definition 6.25 Incremental Taylor:**  
**Goal:** evaluate  $T_n(x)$  (eq. (6.41)) at the point  $x_0 + \Delta x$  in order to propagate the function  $f(x)$  by  $h = \Delta x$ :

$$T_n(x_0 \pm h) = \sum_{i=0}^n \frac{h^i}{n!} f^{(i)}(x_0) i^{-1} \quad (6.42)$$

$$= f(x_0) \pm h f'(x_0) + \frac{h^2}{2} f''(x_0) \pm f'''(x_0)(h)^3 + \mathcal{O}(h^4)$$

**Note**  
 If we chose  $\Delta x$  small enough it is sufficient to look only at the first two terms.

**Definition 6.26 Multidimensional Taylor:** Suppose  $X \in \mathbb{R}^n$  is open,  $\mathbf{x} \in X, f : X \mapsto \mathbb{R}$  and  $f \in \mathcal{C}^2$  then it holds that  

$$f(\mathbf{x}) \approx f(\mathbf{x}_0) + \nabla_{\mathbf{x}} f(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_0)^\top H(\mathbf{x} - \mathbf{x}_0) \quad (6.43)$$

**Definition 6.27 Argmax:** The argmax of a function defined on a set  $D$  is given by:  

$$\arg \max_{x \in D} f(x) = \{x | f(x) \geq f(y), \forall y \in D\} \quad (6.44)$$

**Definition 6.28 Argmin:** The argmin of a function defined on a set  $D$  is given by:  

$$\arg \min_{x \in D} f(x) = \{x | f(x) \leq f(y), \forall y \in D\} \quad (6.45)$$

**Corollary 6.12 Relationship**  $\arg \min \leftrightarrow \arg \max$ :  

$$\arg \min_{x \in D} f(x) = \arg \max_{x \in D} -f(x) \quad (6.46)$$

**Property 6.5 Argmax Identities:**  
**1. Shifting:**  
 $\forall \lambda \text{ const} \quad \arg \max f(x) = \arg \max f(x) + \lambda \quad (6.47)$   
**2. Positive Scaling:**  
 $\forall \lambda > 0 \text{ const} \quad \arg \max f(x) = \arg \max \lambda f(x) \quad (6.48)$   
**3. Negative Scaling:**  
 $\forall \lambda < 0 \text{ const} \quad \arg \max f(x) = \arg \min \lambda f(x) \quad (6.49)$   
**4. Positive Functions:**  
 $\forall \arg \max f(x) > 0, \forall x \in \text{dom}(f)$   

$$\arg \max f(x) = \arg \min \frac{1}{f(x)} \quad (6.50)$$
  
**5. Stricly Monotonic Functions:** for all strictly monotonic increasing functions <sup>([def. 6.10](#))</sup>  $g$  it holds that:  

$$\arg \max g(f(x)) = \arg \max f(x) \quad (6.51)$$

**Definition 6.29 Max:** The maximum of a function  $f$  defined on the set  $D$  is given by:  

$$\max_{x \in D} f(x) = f(x^*) \quad \text{with} \quad \forall x^* \in \arg \max_{x \in D} f(x) \quad (6.52)$$

**Definition 6.30 Min:** The minimum of a function  $f$  defined on the set  $D$  is given by:  

$$\min_{x \in D} f(x) = f(x^*) \quad \text{with} \quad \forall x^* \in \arg \min_{x \in D} f(x) \quad (6.53)$$

**Corollary 6.13 Relationship**  $\min \leftrightarrow \max$ :  

$$\min_{x \in D} f(x) = - \max_{x \in D} -f(x) \quad (6.54)$$

**Property 6.6 Max Identities:**  
**1. Shifting:**  
 $\forall \lambda \text{ const} \quad \max \{f(x) + \lambda\} = \lambda + \max f(x) \quad (6.55)$   
**2. Positive Scaling:**  
 $\forall \lambda > 0 \text{ const} \quad \max \lambda f(x) = \lambda \max f(x) \quad (6.56)$   
**3. Negative Scaling:**  
 $\forall \lambda < 0 \text{ const} \quad \max \lambda f(x) = \lambda \min f(x) \quad (6.57)$   
**4. Positive Functions:**  
 $\forall \arg \max f(x) > 0, \forall x \in \text{dom}(f) \quad \max \frac{1}{f(x)} = \frac{1}{\min f(x)} \quad (6.58)$   
**5. Stricly Monotonic Functions:** for all strictly monotonic increasing functions <sup>([def. 6.10](#))</sup>  $g$  it holds that:  

$$\max g(f(x)) = g(\max f(x)) \quad (6.59)$$

**Definition 6.31 Supremum:** The supremum of a function defined on a set  $D$  is given by:  

$$\sup_{x \in D} f(x) = \{y | y \geq f(x), \forall x \in D\} = \min_{y | y \geq f(x), \forall x \in D} y \quad (6.60)$$
 and is the smallest value  $y$  that is equal or greater  $f(x)$  for any  $x \iff$  smallest upper bound.

**Definition 6.32 Infimum:** The infimum of a function defined on a set  $D$  is given by:  

$$\inf_{x \in D} f(x) = \{y | y \leq f(x), \forall x \in D\} = \max_{y | y \leq f(x), \forall x \in D} y \quad (6.61)$$
 and is the biggest value  $y$  that is equal or smaller  $f(x)$  for any  $x \iff$  largest lower bound.

**Corollary 6.14 Relationship**  $\sup \leftrightarrow \inf$ :  

$$\in_{x \in D} f(x) = - \sup_{x \in D} -f(x) \quad (6.62)$$

**Note**  
 The supremum/infimum is necessary to handle unbound function that seem to converge and for which the max/min does not exist as the argmax/argmin may be empty.  
 E.g. consider  $-e^x/e^x$  for which the max/min converges toward 0 but will never reached s.t. we can always choose a bigger  $x \Rightarrow$  there exists no argmax/argmin  $\Rightarrow$  need to bound the functions from above/below  $\iff$  infimum/supremum.

**Definition 6.33 Time-invariant system (TIS):** A function  $f$  is called time-invariant, if shifting the input in time leads to the same output shifted in time by the same amount.  

$$y(t) = f(x(t), t) \xrightarrow{\text{time-invariance}} y(t - \tau) = f(x(t - \tau), t) \quad \forall \tau \quad (6.63)$$

**Definition 6.34 Inverse Function**  $g = f^{-1}$ :  
 A function  $g$  is the inverse function of the function  $f : A \subset \mathbb{R} \rightarrow B \subset \mathbb{R}$  if  

$$f(g(x)) = x \quad \forall x \in \text{dom}(g) \quad (6.64)$$
 and  

$$g(f(u)) = u \quad \forall u \in \text{dom}(f) \quad (6.65)$$

**Property 6.7 Reflective Property of Inverse Functions:**  $f$  contains  $(a, b)$  if and only if  $f^{-1}$  contains  $(b, a)$ .  
 The line  $y = x$  is a symmetry line for  $f$  and  $f^{-1}$ .

**Theorem 6.5 The Existence of an Inverse Function:**  
 A function has an inverse function if and only if it is one-to-one.

**Corollary 6.15 Inverse functions and strict monotonicity:** If a function  $f$  is **strictly monotonic** <sup>([def. 6.12](#))</sup> on its entire domain, then it is one-to-one and therefore has an inverse function.

## 4. Special Functions

### 4.1. The Gamma Function

**Definition 6.35 The gamma function  $\Gamma(\alpha)$ :** Is extension of the factorial function (??) to the real and complex numbers (with a positive real part):  

$$\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} dx \quad \Re(z) > 0 \quad (6.66)$$

$$\Gamma(n) \stackrel{n \in \mathbb{N}}{\iff} \Gamma(n) = (n - 1)!$$



Differential Calculus

**Definition 7.1 Critical/Stationary Point:** Given a function  $f : \mathbb{R}^n \mapsto \mathbb{R}$ , that is differentiable at a point  $\mathbf{x}_0$  then it is called a **critical point** if the functions derivative vanishes at that point:

$$f'(\mathbf{x}_0) = 0 \iff \nabla_{\mathbf{x}} f(\mathbf{x}_0) = 0$$

**Definition 7.2 Second Derivative**  $\frac{\partial^2}{\partial x_i \partial x_j}$ :

**Corollary 7.1 Second Derivative Test**  $f : \mathbb{R} \mapsto \mathbb{R}$ :  
Suppose  $f : \mathbb{R} \mapsto \mathbb{R}$  is twice differentiable at a stationary point  $x$  [def. 7.1] then it follows that:

- $f''(x) > 0 \iff \begin{matrix} f'(x + \epsilon) > 0 & \text{slope points uphill} \\ f'(x - \epsilon) < 0 & \text{slope points downhill} \\ f(x) \text{ is a local minimum} \end{matrix}$
- $f''(x) < 0 \iff \begin{matrix} f'(x + \epsilon) > 0 & \text{slope points downhill} \\ f'(x - \epsilon) < 0 & \text{slope points uphill} \\ f(x) \text{ is a local maximum} \end{matrix}$

$\epsilon > 0$  sufficiently small enough

**Definition 7.3 Gradient:** Given  $f : n \mapsto \mathbb{R}$  its gradient is defined as:

$$\text{grad}_{\mathbf{x}}(f) = \nabla_{\mathbf{x}} f := \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} & \cdots & \frac{\partial f}{\partial x_n} \end{bmatrix} \quad (7.1)$$

**Definition 7.4 Jacobi Matrix:** Given a vector valued function  $f : \mathbb{R}^n \mapsto \mathbb{R}^m$  its derivative/Jacobian is defined as:

$$\mathbf{J}(f(\mathbf{x})) = \mathbf{J}_f(\mathbf{x}) = \mathbf{D}f = \frac{\partial \mathbf{f}}{\partial \mathbf{x}}(\mathbf{x}) = \frac{\partial (f_1, \dots, f_m)}{\partial (x_1, \dots, x_n)}(\mathbf{x}) = \quad (7.2)$$

$$= \begin{bmatrix} \frac{\partial f_1}{\partial x_1}(\mathbf{x}) & \frac{\partial f_1}{\partial x_2}(\mathbf{x}) & \cdots & \frac{\partial f_1}{\partial x_n}(\mathbf{x}) \\ \frac{\partial f_2}{\partial x_1}(\mathbf{x}) & \frac{\partial f_2}{\partial x_2}(\mathbf{x}) & \cdots & \frac{\partial f_2}{\partial x_n}(\mathbf{x}) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1}(\mathbf{x}) & \frac{\partial f_m}{\partial x_2}(\mathbf{x}) & \cdots & \frac{\partial f_m}{\partial x_n}(\mathbf{x}) \end{bmatrix}$$

**Theorem 7.1 Symmetry of second derivatives/Schwartz's Theorem:**  
Given a continuous and twice differentiable function  $f : \mathbb{R}^n \mapsto \mathbb{R}$  then its second order partial derivatives commute:

$$\frac{\partial}{\partial x_i} \frac{\partial f}{\partial x_j} = \frac{\partial}{\partial x_j} \frac{\partial f}{\partial x_i}$$

**Definition 7.5 Hessian Matrix:**  
Given a function  $f : \mathbb{R} \mapsto \mathbb{R}^n$  its Hessian  $\in \mathbb{R}^{n \times n}$  is defined as:

$$\mathbf{H}(f)(\mathbf{x}) = \mathbf{H}_f(\mathbf{x}) = \mathbf{J}(\nabla f(\mathbf{x}))^T \quad (7.3)$$

$$= \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2}(\mathbf{x}) & \frac{\partial^2 f}{\partial x_1 \partial x_2}(\mathbf{x}) & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n}(\mathbf{x}) \\ \frac{\partial^2 f}{\partial x_2 \partial x_1}(\mathbf{x}) & \frac{\partial^2 f}{\partial x_2^2}(\mathbf{x}) & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n}(\mathbf{x}) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1}(\mathbf{x}) & \frac{\partial^2 f}{\partial x_n \partial x_2}(\mathbf{x}) & \cdots & \frac{\partial^2 f}{\partial x_n^2}(\mathbf{x}) \end{bmatrix}$$

and it corresponds to the Jacobian of the Gradient.  
Due to the differentiability and theorem 7.1 it follows that the Hessian is (if it exists):

- Symmetric
- Real

**Corollary 7.2 Eigenvector basis of the Hessian:** Due to the fact that the Hessian is real and symmetric we can decompose it into a set of real eigenvalues and an orthogonal basis of eigenvectors  $\{(\lambda_1, \mathbf{v}_1), \dots, (\lambda_n, \mathbf{v}_n)\}$ .  
Not let  $\mathbf{d}$  be a directional unit vector then the second derivative in that direction is given by:

$$\mathbf{d}^T \mathbf{H} \mathbf{d} \iff \mathbf{d}^T \sum_{i=1}^n \lambda_i \mathbf{v}_i \iff \text{if } \mathbf{d} = \mathbf{v}_j \quad \mathbf{d}^T \lambda_j \mathbf{v}_j$$

- The eigenvectors that have smaller angle with  $\mathbf{d}$  have bigger weight/eigenvalues
- The minimum/maximum eigenvalue determines the minimum/maximum second derivative

**Corollary 7.3 Second Derivative Test**  $f : \mathbb{R}^n \mapsto \mathbb{R}$ :  
Suppose  $f : \mathbb{R}^n \mapsto \mathbb{R}$  is twice differentiable at a stationary point  $\mathbf{x}$  [def. 7.1] then it follows that:

- If  $\mathbf{H}$  is **p.d**  $\iff \forall \lambda_i > 0 \in \mathbf{H} \rightarrow f(\mathbf{x})$  is a local min.
- If  $\mathbf{H}$  is **n.d**  $\iff \forall \lambda_i < 0 \in \mathbf{H} \rightarrow f(\mathbf{x})$  is a local max.
- If  $\exists \lambda_i > 0 \in \mathbf{H}$  and  $\exists \lambda_i < 0 \in \mathbf{H}$  then  $\mathbf{x}$  is a local maximum in one cross section of  $f$  but a local minimum in another
- If  $\exists \lambda_i = 0 \in \mathbf{H}$  and all other eigenvalues have the same sign the test is inclusive as it is inconclusive in the cross section corresponding to the zero eigenvalue.

Note

If  $\mathbf{H}$  is positive definite for a minima  $\mathbf{x}^*$  of a *quadratic* function  $f$  then this point must be a global minimum of that function.

Integral Calculus

Theorem 8.1 Important Integral Properties:

**Addition**  $\int\limits_a^b f(x) \, dx = \int\limits_a^c f(x) \, dx + \int\limits_c^b f(x) \, dx$  (8.1)

**Reflection**  $\int\limits_a^b f(x) \, dx = - \int\limits_b^a f(x) \, dx$  (8.2)

**Translation**  $\int\limits_a^b f(x) \, dx \stackrel{u:=x\pm c}{=} \int\limits_{a\pm c}^{b\pm c} f(x \mp c) \, dx$  (8.3)

**f Odd**  $\int\limits_{-a}^a f(x) \, dx = 0$  (8.4)

**f Even**  $\int\limits_{-a}^a f(x) \, dx = 2 \int\limits_0^a f(x) \, dx$  (8.5)

Proof. eqs. (8.4) and (8.5)

$$\begin{aligned} I &:= \int\limits_{-a}^a f(x) \, dx = \int\limits_{-a}^0 f(x) \, dx + \int\limits_0^a f(x) \, dx \\ &\stackrel{t=-x}{dt=-dx} = - \int\limits_a^0 f(-x) \, dx + \int\limits_0^a f(x) \, dx \\ &= \int\limits_0^a f(-x) + f(x) \, dx = \begin{cases} 0 & \text{if } f \text{ odd} \\ 2I & \text{if } f \text{ even} \end{cases} \end{aligned}$$

□



Linear Algebra

Given a matrix  $A \in \mathbb{K}^{m,n}$

**Rank:**  $\text{rank}(A) = \dim(\mathfrak{R}(A))$   
of a matrix is the dimension of the vector space generated (or spanned) by its columns/rows.

**Span/Linear Hull:**  $\text{span}(v_1, v_2, \dots, v_n) = \{ \lambda_1 v_1, \lambda_2 v_2, \dots, \lambda_n v_n \} = \{ v \mid v = \sum_{i=1}^n \lambda_i v_i, \lambda_i \in \mathbb{R} \}$

Is the set of vectors tha can be expressed as a linear combination of the vectors  $v_1, \dots, v_n$ .

**Note** these vectors may be linearly independent.

**Generatring Set:** Is the set of vectors which span the  $\mathbb{R}^n$  that is:  $\text{span}(v_1, \dots, v_m) = \mathbb{R}^n$ .  
e.g.  $(4, 0)^T, (0, 5)^T$  span the  $\mathbb{R}^n$ .

**Basis  $\mathfrak{B}$ :** A lin. indep. generating set of the  $\mathbb{R}^n$  is called basis of the  $\mathbb{R}^n$ .

The unit vectors  $e_1, \dots, e_n$  build a standard basis of the  $\mathbb{R}^n$

**Vector Space**

**Image/Range:**  $\mathfrak{R}(A) := \{Ax \mid x \in \mathbb{K}^n\} \subset \mathbb{K}^n$

**Null-Space/Kernel:**  $\mathfrak{N} := \{z \in \mathbb{K}^n \mid Az = 0\}$

**Dimension theorem:**

**Theorem 9.1 Rank-Nullity theorem:** For any  $A \in \mathbb{Q}^{m \times n}$   
 $n = \dim(\mathfrak{N}[A]) + \dim(\mathfrak{R}[A])$

From orthogonality it follows  $x \in \mathfrak{R}(A), y \in \mathfrak{N}(A) \Rightarrow x^\top y = 0$ .

1. Transformations

1.1. Affine Transformations

**Definition 9.1 Affine Transfromation/Map:**  
Let  $x \in \mathbb{R}^n, A \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m$  then:  
 $Y = Ax + b$

is called an affine transformation of  $x$ .

2. Eigenvalues and Vectors

**Formula 9.1 Eigenvalues of a 2x2 matrix:** Given a 2x2-matrix  $A$  its eigenvalues can be calculated by:

$$\{\lambda_1, \lambda_2\} \in \frac{\text{tr}(A) \pm \sqrt{\text{tr}(A)^2 - 4 \det(A)}}{2}$$

with  $\text{tr}(A) = a + d \quad \det(A) = ad - bc$

3. Special Kind of Vectors

**Definition 9.2 Orthogonal Vectors:** Let  $\mathcal{V}$  be an inner-product space<sup>[def. 9.14]</sup>. A set of vectors  $\{u_1, \dots, u_n, \dots\} \in \mathcal{V}$  is called *orthogonal* iff:  
 $\langle u_i, u_j \rangle = 0 \quad \forall i \neq j$

**Definition 9.3 Orthonormal Vectors:** Let  $\mathcal{V}$  be an inner-product space<sup>[def. 9.14]</sup>. A set of vectors  $\{u_1, \dots, u_n, \dots\} \in \mathcal{V}$  is called *orthonormal* iff:  
 $\langle u_i, u_j \rangle = \delta_{ij} = \begin{cases} 1 & i = j \\ 0 & i \neq j \end{cases} \quad \forall i, j$

4. Special Kind of Matrices

**Definition 9.4 Orthogonal Matrix:** A real valued square matrix  $Q \in \mathbb{R}^{n \times n}$  is said to be orthogonal if its row vectors (and respectively its column vectors) build an orthonormal basis:  
 $\langle q_{:,i}, q_{:,j} \rangle = \delta_{ij} \quad \text{and} \quad \langle q_{i,:}, q_{j,:} \rangle = \delta_{ij}$

This is exactly true if the inverse of  $Q$  equals its transpose:  
 $Q^{-1} = Q^\top \iff QQ^\top = Q^\top Q = I$

**Definition 9.5 Unitary/Hermitian Matrices:**  
 $A = A^H$

4.1. Properties of Matrices

4.1.1. Eigendecomposition

**Definition 9.6 Eigendecomposition**  $A = Q\Lambda Q^{-1}$ :

4.1.2. Square Root of p.s.d. Matrices

**Definition 9.7 Square Root:**

4.1.3. Cholesky Decomposition  
5. Spaces and Measures

**Definition 9.8 Bilinear Form/Functional:**  
Is a mapping  $a : \mathcal{V} \times \mathcal{V} \mapsto F$  on a field of scalars  $F \subseteq \mathbb{K}, K = \mathbb{R}$  or  $\mathbb{C}$  that satisfies:  
 $a(\alpha u + \beta v, w) = \alpha a(u, w) + \beta a(v, w)$   
 $a(u, \alpha v + \beta w) = \alpha a(u, v) + \beta a(u, w)$   
 $\forall u, v, w \in \mathcal{V}, \quad \forall \alpha, \beta \in \mathbb{K}$

**Thus:**  $a$  is linear w.r.t. each argument.

**Definition 9.9 Symmetric bilinear form:** A bilinear form  $a$  on  $\mathcal{V}$  is symmetric if and only if:  
 $a(u, v) = a(v, u) \quad \forall u, v \in \mathcal{V}$

**Definition 9.10 Positive (semi) definite bilinear form:**  
A symmetric bilinear form  $a$  on a vector space  $\mathcal{V}$  over a field  $F$  is **positive definite** if and only if:  
 $a(u, u) > 0 \quad \forall u \in \mathcal{V} \setminus \{0\}$   
And **positive semidefinite**  $\iff \geq$

**Corollary 9.1 Matrix induced Bilinear Form:**  
For finite dimensional inner product spaces  $\mathcal{X} \in \mathbb{K}^n$  any *sym-metric* matrix  $A \in \mathbb{R}^{n \times n}$  induces a **bilinear form**:  
 $a(x, x') = x^\top A x' = (A x')^\top x$ ,

**Definition 9.11 Positive (semi) definite Matrix  $>$ :**  
A matrix  $A \in \mathbb{R}^{n \times n}$  is **positive definite** if and only if:  
 $x^\top A x > 0 \iff A > \quad \forall x \in \mathbb{R}^n \setminus \{0\}$   
And **positive semidefinite**  $\iff \geq$

**Corollary 9.2**  
**Eigenvalues of positive (semi) definite matrix:**  
A positive definite matrix is a *symmetric matrix* where every eigenvalue is *strictly* positive and positive semi definite if every eigenvalue is *positive*.  
 $\forall \lambda_i \in \text{eigenv}(A) > 0$   
And **positive semidefinite**  $\iff \geq$

*Proof.* corollary 9.2 (for real matrices):  
Let  $v$  be an eigenvector of  $A$  then it follows:  
 $0 < v^\top A v = v^\top \lambda v = \|v\| \lambda$

**Corollary 9.3 Positive Definiteness and Determinant:**  
The determinant of a positive definite matrix is always positive. **Thus** a positive definite matrix is always *nonsingular*

**Definition 9.12 Negative (semi) definite Matrix  $<$ :**  
A matrix  $A \in \mathbb{R}^{n \times n}$  is **negative definite** if and only if:  
 $x^\top A x < 0 \iff A < 0 \quad \forall x \in \mathbb{R}^n \setminus \{0\}$   
And **negative semidefinite**  $\iff \leq$

**Theorem 9.2 Sylvester's criterion:** Let  $A$  be *symmet-ric/Hermitian* matrix and denote by  $A^{(k)}$  the  $k \times k$  upper left sub-matrix of  $A$ .  
Then it holds that:  

- $A > 0 \iff \det(A^{(k)}) > 0 \quad k = 1, \dots, n$
- $A < 0 \iff (-1)^k \det(A^{(k)}) > 0 \quad k = 1, \dots, n$

- $A$  is indefinite if the first  $\det(A^{(k)})$  that breaks both of the previous patterns is on the wrong side.
- Sylvester's criterion is inconclusive ( $A$  can be anything of the previous three) if the first  $\det(A^{(k)})$  that breaks both patterns is 0.

6. Inner Products

**Definition 9.13 Inner Product:** Let  $\mathcal{V}$  be a vector space over a field  $F \in \mathbb{K}$  of scalars. An inner product on  $\mathcal{V}$  is a map:  
 $\langle \cdot, \cdot \rangle : \mathcal{V} \times \mathcal{V} \mapsto F \subseteq \mathbb{K} \quad K = \mathbb{R} \text{ or } \mathbb{C} \quad (9.18)$   
that satisfies:  
 $\forall x, y, z \in \mathcal{V}, \quad \alpha, \beta \in F$   

- (Conjugate) Stmmetry:**  $\langle x, y \rangle = \overline{\langle y, x \rangle}$ .
- Linearity** in the first argument:  
 $\langle \alpha x + \beta y, z \rangle = \alpha \langle x, z \rangle + \beta \langle y, z \rangle$
- Positive-definiteness:**  
 $\langle x, x \rangle \geq 0 : x = 0 \iff \langle x, x \rangle = 0$

**Definition 9.14 Inner Product Space  $(\mathcal{V}, \langle \cdot, \cdot \rangle_{\mathcal{V}})$ :** Let  $F \in \mathbb{K}$  be a field of scalars.  
An inner product space  $\mathcal{V}$  is a vetor space over a field  $F$  together with an an **inner product**  $\langle \cdot, \cdot \rangle_{\mathcal{V}}$ .

**Corollary 9.4 Inner product $\rightarrow$ S.p.d. Bilinear Form:**  
Let  $\mathcal{V}$  be a vector space over a field  $F \in \mathbb{K}$  of scalar.  
An **inner product** on  $\mathcal{V}$  is a positive definite symmetric bilinear form on  $\mathcal{V}$ .

**Example: scalar prodct**  
Let  $a(u, v) = u^\top I v$  then the standard scalar product can be defined in terms of a bilinear form vice versa the standard scalar product induces a bilinear form.

**Note**  
Inner products must be positive definite by definition  
 $\langle x, x \rangle \geq 0$ , whereas bilinear forms must not.

**Definition 9.15 Norm  $\|\cdot\|_{\mathcal{V}}$ :**  
A norm measures the **size** of its argument.  
**Formally** let  $\mathcal{V}$  be a vector space over a field  $F$ , a norm on  $\mathcal{V}$  is a map:  
 $\|\cdot\|_{\mathcal{V}} : \mathcal{V} \mapsto \mathbb{R}_+ \quad (9.19)$   
that satisfies:  $\forall x, y \in \mathcal{V}, \quad \alpha \in F \subseteq \mathbb{K} \quad K = \mathbb{R} \text{ or } \mathbb{C}$   

- Definitness:**  $\|x\|_{\mathcal{V}} = 0 \iff x = 0$ .
- Homogenity:**  $\|\alpha x\|_{\mathcal{V}} = |\alpha| \|x\|_{\mathcal{V}}$
- Triangular Inequality:**  $\|x + y\|_{\mathcal{V}} \leq \|x\|_{\mathcal{V}} + \|y\|_{\mathcal{V}}$

**Meaning: Triangular Inequality**  
States that for any triangle, the sum of the lengths of any two sides must be greater than or equal to the length of the remaining side.

**Corollary 9.5 Reverse Triangular Inequality:**  
 $-\|x - y\|_{\mathcal{V}} \leq \|x\|_{\mathcal{V}} - \|y\|_{\mathcal{V}} \leq \|x - y\|_{\mathcal{V}}$   
resp.  $|\|x\|_{\mathcal{V}} - \|y\|_{\mathcal{V}}| \leq \|x - y\|_{\mathcal{V}}$

**Semi-norm**  
**Add**

**Corollary 9.6 Normed vector space:** Is a vector space  $\mathcal{V}$  over a field  $F$ , on which a norm  $\|\cdot\|_{\mathcal{V}}$  can be defined.

**Corollary 9.7 Inner product induced norm  $\langle \cdot, \cdot \rangle_{\mathcal{V}} \rightarrow \|\cdot\|_{\mathcal{V}}$ :** Every inner product  $\langle \cdot, \cdot \rangle_{\mathcal{V}}$  induces a norm of the form:  
 $\|x\|_{\mathcal{V}} = \sqrt{\langle x, x \rangle} \quad x \in \mathcal{V}$

**Thus** We can define function spaces by their associated norm  $(\mathcal{V}, \|\cdot\|_{\mathcal{V}})$  and inner product spaces lead to normed vector spaces and vice versa.

**Corollary 9.8 Energy Norm:** A *s.p.d.* bilinear form  $a : \mathcal{V} \times \mathcal{V} \mapsto F$  induces an **energy norm**:  
 $\|x\|_a := (a(x, x))^{\frac{1}{2}} = \sqrt{a(x, x)} \quad x \in \mathcal{V}$

**Definition 9.16 Distance Function/Measure:** Is measuring the **distance** between two things.  
**Formally:** on a set  $S$  is a mapping:  
 $d(\cdot, \cdot) : S \times S \mapsto \mathbb{R}_+$   
that satisfies:  
 $\forall x, y, z \in S$   

- ?**:  $d(x, x) = 0$
- Symmetry:**  $d(x, y) = d(y, x)$
- Triangular Identi**y:  $d(x, z) \leq d(x, y) + d(y, z)$

**Definition 9.17 Metric:** Is a distance measure that additionally satisfies:  
 $\forall x, y \in S$   
**identity of indiscernibles :**  $d(x, y) = 0 \iff x = y$

**Corollary 9.9 Metric $\rightarrow$ Norm:** Every norm  $\|\cdot\|_{\mathcal{V}}$  on a vector space  $\mathcal{V}$  over a field  $F$  induces a metric by:  
 $d(x, y) = \|x - y\|_{\mathcal{V}} \quad \forall x, y \in \mathcal{V}$

metric induced by norms additionally satisfy:  $\forall x, y \in \mathcal{V}, \quad \alpha \in F \subseteq \mathbb{K} \quad K = \mathbb{R} \text{ or } \mathbb{C}$   

- Homogenity/Scaling:**  $d(\alpha x, \alpha y)_{\mathcal{V}} = |\alpha| d(x, y)_{\mathcal{V}}$
- Translational Invariance:**  $d(x + \alpha, y + \alpha) = d(x, y)$

Conversely not every metric induces a norm **but** if a metric  $d$  on a vector space  $\mathcal{V}$  satisfies the properties then it induces a norm of the form:  
 $\|x\|_{\mathcal{V}} := d(x, 0)_{\mathcal{V}}$

**Note**  
Similarity measure is a much weaker notion than a metric as triangular inequality does not have to hold.  
**Hence:** If  $a$  is similar to  $b$  and  $b$  is similar to  $c$  it does not imply that  $a$  is similar to  $c$ .

**Note**  
(bilinear form  $\xrightarrow{\text{induces}}$ )  
inner product  $\xrightarrow{\text{induces}}$  norm  $\xrightarrow{\text{induces}}$  metric.

7. Vector Algebra

**7.1. Planes**  
<https://math.stackexchange.com/questions/1485509/show-that-two-planes-are-parallel-and-find-the-distance-between-them>

8. Derivatives

$\frac{\partial}{\partial \mathbf{x}} (\mathbf{b}^\top \mathbf{x}) = \frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^\top \mathbf{b}) = \mathbf{b}$   
 $\frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^\top \mathbf{x}) = 2\mathbf{x}$   
 $\frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^\top \mathbf{A} \mathbf{x}) = (\mathbf{A}^\top + \mathbf{A})\mathbf{x} \quad \frac{\partial}{\partial \mathbf{x}} (\mathbf{b}^\top \mathbf{A} \mathbf{x}) = \mathbf{A}^\top \mathbf{b}$   
 $\frac{\partial}{\partial \mathbf{X}} (\mathbf{c}^\top \mathbf{X} \mathbf{b}) = \mathbf{c} \mathbf{b}^\top \quad \frac{\partial}{\partial \mathbf{x}} (\|\mathbf{x} - \mathbf{b}\|_2) = \frac{\mathbf{x} - \mathbf{b}}{\|\mathbf{x} - \mathbf{b}\|_2}$   
 $\frac{\partial}{\partial \mathbf{x}} (\|\mathbf{x}\|_2^2) = \frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^\top \mathbf{x}) = 2\mathbf{x} \quad \frac{\partial}{\partial \mathbf{X}} (\|\mathbf{X}\|_F^2) = 2\mathbf{X}$   
 $\frac{\partial}{\partial \mathbf{x}} \|\mathbf{x}\|_1 = \frac{\mathbf{x}}{|\mathbf{x}|}$   
 $\frac{\partial}{\partial \mathbf{x}} (\|\mathbf{A} \mathbf{x} - \mathbf{b}\|_2^2) = 2(\mathbf{A}^\top \mathbf{A} \mathbf{x} - \mathbf{A}^\top \mathbf{b}) \quad \frac{\partial}{\partial \mathbf{X}} (|\mathbf{X}|) = |\mathbf{X}| \cdot \mathbf{X}^{-1}$   
 $\frac{\partial}{\partial \mathbf{x}} (\mathbf{Y}^{-1}) = -\mathbf{Y}^{-1} \frac{\partial \mathbf{Y}}{\partial \mathbf{x}} \mathbf{Y}^{-1}$

Geometry

**Corollary 10.1 Affine Transformation in 1D:** Given: numbers  $x \in \hat{\Omega}$  with  $\hat{\Omega} = [a, b]$   
The **affine transformation** of  $\phi : \hat{\Omega} \rightarrow \Omega$  with  $y \in \Omega = [c, d]$  is defined by:

$$y = \phi(x) = \frac{d - c}{b - a} (x - a) + c \tag{10.1}$$

*Proof.* **corollary 10.1** By <sup>[def. 9.1]</sup> we want a function  $f : [a, b] \rightarrow [c, d]$  that satisfies:

$f(a) = c$                       **and**                       $f(b) = d$

additionally  $f(x)$  has to be a linear function (<sup>[def. 6.13]</sup>), that is the output scales the same way as the input scales.

Thus it follows:  
$$\frac{d - c}{b - a} = \frac{f(x) - f(a)}{x - a} \iff f(x) = \frac{d - c}{b - a} (x - a) + c$$

Trigonometry

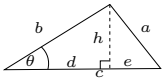
**Law 10.1 Law of Cosine:** relates the side of a triangle to the cosine of its angles.

$$a^2 = b^2 + c^2 - 2bc \cos \theta_{b,c} \tag{10.2}$$

More general for vectors it holds:

$$\|\mathbf{x} - \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 - 2\|\mathbf{x}\|\|\mathbf{y}\|\cos \theta_{\mathbf{x},\mathbf{y}} \tag{10.3}$$

*Proof.* eq. (10.2):  
**We know:**  $\sin \theta = \frac{h}{b} \Rightarrow \underline{h} = b \sin \theta$                       and                       $\cos \theta = \frac{d}{b} \Rightarrow d = b \cos \theta$   
**Thus**  $\underline{e} = c - d = c - b \cos \theta \Rightarrow a^2 = \underline{e}^2 + \underline{h}^2 \Rightarrow a$                        $\square$



*Proof.* eq. (10.3):  
$$\begin{aligned} \|\mathbf{x} - \mathbf{y}\|^2 &= (\mathbf{x} - \mathbf{y}) \cdot (\mathbf{x} - \mathbf{y}) \\ &= \mathbf{x} \cdot \mathbf{x} - 2\mathbf{x} \cdot \mathbf{y} + \mathbf{y} \cdot \mathbf{y} \\ &= \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 - 2(\|\mathbf{x}\|\|\mathbf{y}\|\cos \theta) \end{aligned}$$

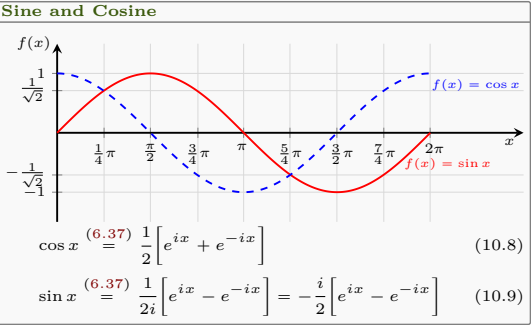
**Law 10.2 Pythagorean theorem:** special case of ?? for right triangle:

$$a^2 = b^2 + c^2 \tag{10.4}$$

**Formula 10.1 Euler's Formula:**  
$$e^{\pm ix} = \cos x \pm i \sin x \tag{10.5}$$

**Formula 10.2 Euler's Identity:**  
$$e^{\pm i} = -1 \tag{10.6}$$

**Note**  
$$e^n = 1 \Leftrightarrow n = i2\pi k, \quad k \in \mathbb{N} \tag{10.7}$$



**Sinh and Cosh**

$$\cosh x \stackrel{(6.37)}{=} \frac{1}{2} \left[ e^x + e^{-x} \right] = \cos(ix) \tag{10.10}$$
$$\sinh x \stackrel{(6.37)}{=} \frac{1}{2} \left[ e^x - e^{-x} \right] = -i \sin(ix) \tag{10.11}$$

**Note**

$$e^x = \cosh x + \sinh x \quad e^{-x} = \cosh x - \sinh x \tag{10.12}$$

**Note**

- $\cosh x$  is strictly positive.
- $\sinh x = 0$  has a unique root at  $x = 0$ .

**Theorem 10.1 Addition Theorems:**

$$\sin(\alpha \pm \beta) = \sin \alpha \cos \beta \pm \cos \alpha \sin \beta \tag{10.13}$$
$$\cos(\alpha \pm \beta) = \cos \alpha \cos \beta \mp \sin \alpha \sin \beta \tag{10.14}$$

**Werner Formulas**

$$\sin \alpha \cos \beta = \frac{1}{2} \left[ \sin(\alpha + \beta) + \sin(\alpha - \beta) \right] \tag{10.15}$$
$$\sin \alpha \sin \beta = \frac{1}{2} \left[ \cos(\alpha - \beta) - \cos(\alpha + \beta) \right] \tag{10.16}$$
$$\cos \alpha \cos \beta = \frac{1}{2} \left[ \cos(\alpha + \beta) + \cos(\alpha - \beta) \right] \tag{10.17}$$

**Note**

Using theorem 10.1 if follows:

$$\cos(\alpha \pm \pi) = -\cos \alpha \quad \text{and} \quad \sin(\alpha \pm \pi) = -\sin \alpha \tag{10.18}$$

Topology

Numerics

1. Machine Arithmetic's

1.1. Machine Numbers

**Definition 12.1 Institute of Electrical and Electronics Engineers (IEEE):** Is a engineering associations that defines a standard on how computers should treat machine numbers in order to have certain guarantees.

**Definition 12.2 Machine/Floating Point Numbers  $\mathbb{F}$ :** Computers are only capable to represent a *finite, discrete* set of the real numbers  $\mathbb{F} \subset \mathbb{R}$

**1.1.1. Floating Point Arithmetic's**  $x\tilde{\Omega}y = \mathfrak{fl}(x\Omega y)$

**Corollary 12.1 Closure:** Machine numbers  $\mathbb{F}$  are not *closed*<sup>[def. 4.6]</sup> under basic arithmetic operations:  
 $\mathbb{F} \Omega \mathbb{F} \mapsto \nmid \mathbb{F} \quad \Omega = \{+, -, *, /\}$  (12.1)

**Note**  
Corollary 12.1 provides a problem as the computer can only represent floating point number  $\mathbb{F}$ .

**Definition 12.3 Floating Point Operation  $\tilde{\Omega}$ :**  
Is a basic arithmetic operation that obtains a number  $x \in \mathbb{F}$  by applying a function rd:  
 $\mathbb{F} \tilde{\Omega} \mathbb{F} \mapsto \mathbb{F} \quad \tilde{\Omega} := \text{rd} \circ \Omega$   
 $\Omega = \{+, -, *, /\}$  (12.2)

**Definition 12.4 Rounding Function rd:**  
Given a real number  $x \in \mathbb{R}$  the rounding function replaces it by the nearest machine number  $\tilde{x} \in \mathbb{F}$ . If this is ambiguous (there are two possibilities), then it takes the larger one:  
 $\text{rd} : \begin{cases} \mathbb{R} \mapsto \mathbb{F} \\ x \mapsto \max_{\tilde{x} \in \mathbb{F}} \min |x - \tilde{x}| \end{cases}$  (12.3)

**Consequence**  
Basic arithmetic rules such as associativity do no longer hold for operations such as addition and subtraction.  
**Axiom 12.1 Axiom of Round off Analysis:**  
Let  $x, y \in \mathbb{F}$  be (normalized) floats and assume that  $x\tilde{\Omega}y \in \mathbb{F}$  (i.e. no over/underflow). Then it holds that:  
 $x\tilde{\Omega}y = (x\Omega y)(1 + \delta) \quad \Omega = \{+, -, *, /\}$   
 $\tilde{f}(x) = f(x)(1 + \delta) \quad f \in \{\exp, \sin, \cos, \log, \dots\}$  (12.4)  
with  $|\delta| < \text{EPS}$

**Explanation 12.1** (axiom 12.1). *gives us a guarantee that for any two floating point numbers  $x, y \in \mathbb{F}$ , any operation involving them will give a floating point result which is within a factor of  $1 + \delta$  of the true result  $x\Omega y$ .*

**Remark**

**Definition 12.5 Overflow:** Result is bigger then the biggest representable floating point number.

**Definition 12.6 Underflow:** Result is smaller then the smaller representable floating point number i.e. to close to zero.

1.2. Roundoff Errors Log-Sum-Exp Trick

The sum exponential trick is at trick that helps to calculate the log-sum-exponential in a robust way by avoiding over/underflow. The log-sum-exponential<sup>[def. 12.7]</sup> is an expression that arises frequently in machine learning i.e. for the cross entropy loss or for calculating the evidence of a posterior prediction.  
The root of the problem is that we need to calculate the exponential  $\exp(x)$ , this comes with two different problems:  
• If  $x$  is large (i.e. 89 for single precision floats) then  $\exp(x)$  will lead to overflow  
• If  $x$  is very negative  $\exp(x)$  will lead to underflow/0. This is not necessarily a problem but if  $\exp(x)$  occurs in the denominator or the logarithm for example this is catastrophic.

**Definition 12.7 Log sum Exponential:**  
 $\text{LogSumExp}(x_1, \dots, x_n) := \log \left( \sum_{i=1}^n e^{x_i} \right)$  (12.5)

**Formula 12.1 Log-Sum-Exp Trick:**  
 $\log \left( \sum_{i=1}^n e^{x_i} \right) = a + \log \sum_{i=1}^n e^{x_i - a} \quad a := \max_{i \in \{1, \dots, n\}} x_i$  (12.6)

**Explanation 12.2** (formula 12.1). *The value  $a$  can be any real value but for robustness one usually chooses the max s.t.*  
• The leading digits are preserved by pulling out the maximum  $a$   
• Inside the log only zero or negative numbers are exponentiated, so there can be no overflow.  
• If there is underflow inside the log we know that at least the leading digits have been returned by the max.

*Proof.*  
$$\begin{aligned} \text{LSE} &= \log \left( \sum_{i=1}^n e^{x_i} \right) = \log \left( \sum_{i=1}^n e^{x_i - a} e^a \right) \\ &= \log \left( e^a \sum_{i=1}^n e^{x_i - a} \right) = \log \left( \sum_{i=1}^n e^{x_i - a} \right) + \log(e^a) \\ &= \log \left( \sum_{i=1}^n e^{x_i - a} \right) + a \end{aligned}$$

**Definition 12.8 Partition  $\Pi$ :**  
Given an interval  $[0, T]$  a sequence of values  $0 < t_0 < \dots < t_n < T$  is called a partition  $\Pi(t_0, \dots, t_n)$  of this interval.

1.3. Convergence for iterative methods

**Definition 12.9 Linear/Exponential Convergence:** A sequence  $\{x^{(k)}\}_k \in \mathbb{R}^n$  converges linearly to  $x^*$  if in the asymptotic limit  $k \rightarrow \infty$  it satisfies:  
 $\left\| x^{k+1} - x^* \right\| \leq \rho \left\| x^{(k)} - x^* \right\| \quad \rho \in (0, 1), \forall k \in \mathbb{N}_0$  (12.7)

**Exponential Convergence**  
Linear convergence is sometimes called exponential convergence. This is due to the fact that:  
1. We often have expressions of the form:  
 $\left\| x^{k+1} - x^* \right\| \leq \underbrace{(1 - \alpha)}_{:= \rho} \left\| x^{(k)} - x^* \right\|$   
2. and that  $(1 - \alpha) = \exp(-\alpha)$  from which follows that:  
eq. (12.8)  $\iff \left\| x^{k+1} - x^* \right\| \leq e^{-\alpha} \left\| x^{(k)} - x^* \right\|$

**Definition 12.10 Rate of Convergence:** Is a way to measure the rate of convergence of a sequence  $\{x^{(k)}\}_k \in \mathbb{R}^n$  to a value to  $x^*$ . Let  $\rho \in [0, 1]$  be the *rate of convergence* and define:

$$\lim_{k \mapsto \infty} \frac{\left\| x^{k+1} - x^* \right\|}{\left\| x^{(k)} - x^* \right\|} = \rho \quad (12.8)$$

- $\rho = 1 \iff$  **Sublinear Rate** i.e. slower than linear
- $\rho \in (0, 1) \iff$  **Linear Rate**
- $\rho = 0 \iff$  **Superlinear Rate** i.e. faster then linear

**Definition 12.11 Convergence of order  $p$ :** In order to distinguish *superlinear convergence* we define the order of convergence.  
A sequence  $\{x^{(k)}\}_k \in \mathbb{R}^n$  converges superlinear with order  $p \in \{2, \dots\}$  to  $x^*$  if it satisfies:  
$$\lim_{k \mapsto \infty} \frac{\left\| x^{k+1} - x^* \right\|}{\left\| x^{(k)} - x^* \right\|^p} = C \quad C < 1 \quad (12.9)$$

Does this even exist/check if this is true

**Definition 12.12 Exponential Convergence:** A sequence  $\{x^{(k)}\}_k \in \mathbb{R}^n$  converges exponentially with rate  $\rho$  to  $x^*$  if in the asymptotic limit  $k \rightarrow \infty$  it satisfies:  
 $\left\| x^{k+1} - x^* \right\| \leq \rho^k \left\| x^{(k)} - x^* \right\| \quad \rho < 1 \quad (12.10)$

1.4. Convention for discretization methods 2. Numerical Quadrature

**Definition 12.13 Order of a Quadrature Rule:** The order of a quadrature rule  $\mathcal{Q}_n : \mathcal{C}^0([a, b]) \rightarrow \mathbb{R}$  is defined as:  
 $\text{order}(\mathcal{Q}_n) := \max \left\{ n \in \mathbb{N}_0 : \mathcal{Q}_n(p) = \int_a^b p(t) dt \quad \forall p \in \mathcal{P}_n \right\} + 1$  (12.11)

Thus it is the maximal degree+1 of polynomials (of degree maximal degree)  $\mathcal{P}$  maximal degree for which the quadrature rule yields exact results.

**Note**  
Is a quality measure for quadrature rules.  
**2.1. Composite Quadrature**

**Definition 12.14 Composite Quadrature:**  
Given a mesh  $\mathcal{M} = \{a = x_0 < x_1 < \dots < x_m = b\}$  apply a Q.R.  $\mathcal{Q}_n$  to each of the mesh cells  $I_j := [x_{j-1}, x_j] \quad \forall j = 1, \dots, m \triangleq \text{p.w.}$  Quadrature:  
$$\int_a^b f(t) dt = \sum_{j=1}^m \int_{x_{j-1}}^{x_j} f(t) dt = \sum_{j=1}^m \mathcal{Q}_n(f|_{I_j}) \quad (12.12)$$

**Lemma 12.1 Error of Composite quadrature Rules:**  
Given a function  $f \in \mathcal{C}^k([a, b])$  with integration domain:  
$$\sum_{i=1}^m h_i = |b - a| \quad \text{for } \mathcal{M} = \{x_j\}_{j=1}^m$$
  
Let:  $h_{\mathcal{M}} = \max_j |x_j, x_{j-1}|$  be the **mesh-width**  
Assume an equal number of quadrature nodes for each interval  $I_j = [x_{j-1}, x_j]$  of the mesh  $\mathcal{M}$  i.e.  $n_j = n$ .  
Then the error of a quadrature rule  $\mathcal{Q}_n(f)$  of order  $q$  is given by:  
$$\epsilon_n(f) = \mathcal{O} \left( n^{-\min\{k, q\}} \right) = \mathcal{O} \left( h_{\mathcal{M}}^{\min\{k, q\}} \right) \quad \text{for } n \rightarrow \infty$$
  
corollary 6.3  $\mathcal{O} \left( n^{-q} \right) = \mathcal{O} \left( h_{\mathcal{M}}^q \right) \quad \text{with } h_{\mathcal{M}} = \frac{1}{n}$  (12.13)

**Definition 12.15 Complexity  $W$ :** Is the number of function evaluations  $\triangleq$  number of quadrature points.  
 $W(\mathcal{Q}(f)_n) = \#f\text{-eval} \triangleq n$  (12.14)

**Lemma 12.2 Error-Complexity  $W(\epsilon_n(f))$ :** Relates the complexity to the quadrature error.  
**Assuming** and quadrature error of the form :  
$$\epsilon_n(f) = \mathcal{O}(n^{-q}) \iff \epsilon_n(f) = cn^{-q} \quad c \in \mathbb{R}_+$$
  
the error complexity is **algebraic** (??) and is given by:  
$$W(\epsilon_n(f)) = \mathcal{O}(\epsilon_n^{1/q}) = \mathcal{O} \left( \sqrt[q]{\epsilon_n} \right) \quad (12.15)$$

*Proof.* lemma 12.2: **Assume:** we want to reduce the error by a factor of  $\epsilon_n$  by increasing the number of quadrature points  $n_{\text{new}} = a \cdot n_{\text{old}}$ .  
**Question:** what is the additional effort ( $\#f\text{-eval}$ ) needed in order to achieve this reduction in error?  
$$\frac{c \cdot n_n^q}{c \cdot n_o^q} = \frac{1}{\epsilon_n} \implies n_n = n_o \cdot \sqrt[q]{\epsilon_n} = \mathcal{O} \left( \sqrt[q]{\epsilon_n} \right) \quad (12.16)$$

Optimization

**Definition 13.1 Fist Order Method:** A first-order method is an algorithm that chooses the  $k$ -th iterate in  
 $w_0 + \text{span}\{\nabla f(w_0), \dots, \nabla f(w_{k-1})\} \quad \forall k = 1, 2, \dots$  (13.1)

**Note**  
Gradient descent is a first order method

1. Lagrangian Optimization Theory

**Add:** derivation of lagrange function

**Definition 13.2 (Primal) Constraint Optimization:**  
Given an optimization problem with domain  $\Omega \subseteq \mathbb{R}^d$ :  
$$\begin{aligned} &\min_{w \in \Omega} f(w) \\ \text{s.t.} \quad &g_i(w) \leq 0 \quad 1 \leq i \leq k \\ &h_j(w) = 0 \quad 1 \leq j \leq m \end{aligned}$$

**Definition 13.3 Lagrange Function:**  
$$\mathcal{L}(\alpha, \beta, w) := f(w) + \alpha g(w) + \beta h(w) \quad (13.2)$$

**Extremal Conditions**  
$$\begin{aligned} \nabla \mathcal{L}(x) &\stackrel{!}{=} 0 && \text{Extremal point } x^* \\ \frac{\partial}{\partial \beta} \mathcal{L}(x) &= h(x) \stackrel{!}{=} 0 && \text{Constraint satisfaction} \end{aligned}$$
  
For the inequality constraints  $g(x) \leq 0$  we distinguish two situations:  
Case I :  $g(x^*) < 0$  switch const. off  
Case II :  $g(x^*) \geq 0$  optimize using active eq. constr.  
$$\frac{\partial}{\partial \alpha} \mathcal{L}(x) = g(x) \stackrel{!}{=} 0 \quad \text{Constraint satisfaction}$$

**Definition 13.4 Lagrangian Dual Problem:** Is given by:  
Find  $\max_{\alpha, \beta} \theta(\alpha, \beta) = \inf_{w \in \Omega} \mathcal{L}(w, \alpha, \beta)$   
s.t.  $\alpha_i \geq 0 \quad 1 \leq i \leq k$

**Solution Strategy**  
1. Find the extremal point  $w^*$  of  $\mathcal{L}(w, \alpha, \beta)$ :  
$$\frac{\partial \mathcal{L}}{\partial w} \Big|_{w=w^*} \stackrel{!}{=} 0 \quad (13.3)$$
  
2. Insert  $w^*$  into  $\mathcal{L}$  and find the extremal point  $\beta^*$  of the resulting dual Lagrangian  $\theta(\alpha, \beta)$  for the active constraints:  
$$\frac{\partial \theta}{\partial \beta} \Big|_{\beta=\beta^*} \stackrel{!}{=} 0 \quad (13.4)$$
  
3. Calculate the solution  $w^*(\beta^*)$  of the constraint minimization problem.

**Value of the Problem**  
**Value of the problem:** the value  $\theta(\alpha^*, \beta^*)$  is called the value of problem  $(\alpha^*, \beta^*)$ .

**Theorem 13.1 Upper Bound Dual Cost:** Let  $w \in \Omega$  be a feasible solution of the primal problem<sup>[def. 13.2]</sup> and  $(\alpha, \beta)$  a feasible solution of the respective dual problem<sup>[def. 13.4]</sup>.  
Then it holds that:  
$$f(w) \geq \theta(\alpha, \beta) \quad (13.5)$$

*Proof.*  
$$\begin{aligned} \theta(\alpha, \beta) &= \inf_{u \in \Omega} \mathcal{L}(u, \alpha, \beta) \leq \mathcal{L}(w, \alpha, \beta) \\ &= f(w) + \sum_{i=1}^k \underbrace{\alpha_i}_{\geq 0} g_i(w) + \sum_{j=1}^m \underbrace{\beta_j}_{\leq 0} \underbrace{h_j(w)}_{=0} \\ &\leq f(w) \end{aligned}$$

**Corollary 13.1 Duality Gap Corollary:** The value of the dual problem is upper bounded by the value of the primal problem:  
$$\sup \{\theta(\alpha, \beta) : \alpha \geq 0\} \leq \inf \{f(w) : g(w) \leq 0, h(w) = 0\} \quad (13.6)$$

**Theorem 13.2 Optimality:** The triple  $(w^*, \alpha^*, \beta^*)$  is a saddle point of the Lagrangian function for the primal problem, if and only if its components are optimal solutions of the primal and dual problems and if there is no duality gap, that is, the primal and dual problems having the same value:  
$$f(w^*) = \theta(\alpha^*, \beta^*) \quad (13.7)$$

**Definition 13.5 Convex Optimization:** Given: a **convex function**  $f$  and a **convex set**  $S$  solve:

$$\begin{aligned} \min_{\mathbf{x} \in S} f(\mathbf{x}) \end{aligned} \tag{13.8}$$

Often  $S$  is specified using linear inequalities:

e.g.  $S = \left\{ \mathbf{x} \in \mathbb{R}^d : \mathbf{A}\mathbf{x} \leq \mathbf{b} \right\}$

**Theorem 13.3 Strong Duality:** Given an convex optimization problem:

$$\begin{aligned} \min_{\mathbf{w} \in \Omega} f(\mathbf{w}) \\ \text{s.t.} \quad & g_i(\mathbf{w}) \leq 0 & 1 \leq i \leq k \\ & h_j(\mathbf{w}) = 0 & 1 \leq j \leq m \end{aligned}$$

where  $g_i, h_i$  can be written as affine functions:  $y(\mathbf{w}) = \mathbf{A}\mathbf{w} - \mathbf{b}$ .

Then it holds that the **duality gap** is zero and we obtain an optimal solution.

**Theorem 13.4 Kuhn-Tucker Conditions:** Given an optimization problem with convex domain  $\Omega \subseteq \mathbb{R}^d$ ,

$$\begin{aligned} \min_{\mathbf{w} \in \Omega} f(\mathbf{w}) \\ \text{s.t.} \quad & g_i(\mathbf{w}) \leq 0 & 1 \leq i \leq k \\ & h_j(\mathbf{w}) = 0 & 1 \leq j \leq m \end{aligned}$$

with  $f \in C^1$  convex and  $g_i, h_i$  affine.

**Necessary and sufficient conditions** for a normal point  $\mathbf{w}^*$  to be an optimum are the existence of  $\alpha^*, \beta^*$  s.t.:

$$\frac{\partial \mathcal{L}(\mathbf{w}, \alpha, \beta)}{\partial \mathbf{w}} \stackrel{!}{=} 0 \qquad \frac{\partial \mathcal{L}(\mathbf{w}^*, \alpha, \beta)}{\partial \beta} \stackrel{!}{=} 0 \tag{13.9}$$

under the conditions that:

- $\forall i_1, \dots, k \quad \alpha_i^* g_i(\mathbf{w}^*) = 0$ , s.t.:
  - Inactive Constraint:  $g_i(\mathbf{w}^*) < 0 \rightarrow \alpha_i = 0$ .
  - Active Constraint:  $g_i(\mathbf{w}^*) \nless 0 \rightarrow \alpha_i \geq 0 \quad \text{s.t.} \quad \alpha_i^* g_i(\mathbf{w}^*) = 0$

**Consequence**

We may become very sparce problems, if a lot of constraints are not active  $\iff \alpha_i = 0$ .

Only a few points, for which  $\alpha_i > 0$  may affact the decision surface.

Stochastics

<b>Definition 13.6 Stochastics:</b> Is a collective term for the areas of <i>probability theory</i> and <i>statistics</i> .
<b>Definition 13.7 Statistics:</b> Is concerned with the analysis of data/experiments in order to draw conclusion of the underlying governing models that describe these experiments.
<b>Definition 13.8 Probability:</b> Is concerned with the quantification of the uncertainty of random experiments by use of statistical models. Hence it is the opposite of statistics.
<b>Definition 13.9 Probability:</b> Probability is the measure of the likelihood that an event will occur in a Random Experiment. Probability is quantified as a number between 0 and 1, where, loosely speaking, 0 indicates impossibility and 1 indicates certainty.
<small>Improve these definitions, maybe ask on quora/hilo</small>
<b>Note: Stochastics vs. Stochastic</b> Stochastics is a noun and is a collective term for the areas of probability theory and statistics, while stochastic is an <i>adjective</i> , describing that a certain phenomena is governed by uncertainty i.e. a process.
<b>Probability Theory</b>
<b>Definition 14.1 Probability Space</b> $W = \{\Omega, \mathcal{F}, \mathbb{P}\}$ : Is the unique triple $\{\Omega, \mathcal{F}, \mathbb{P}\}$ , where $\Omega$ is its sample space, $\mathcal{F}$ is its $\sigma$ -algebra of events, and $\mathbb{P}$ its probability measure.
<b>Definition 14.2 Sample Space <math>\Omega</math>:</b> Is the set of all possible outcomes (elementary events corollary 14.5) of an experiment see example 14.1
<b>Definition 14.3 Event</b> $A$ : An “event” is a subset of the sample space $\Omega$ and is a property which can be observed to hold or not to hold <i>after</i> the experiment is done (example 14.2). Mathematically speaking not every subset of $\Omega$ is an event and has an associated probability. Only those subsets of $\Omega$ that are part of the corresponding $\sigma$ -algebra $\mathcal{F}$ are events and have their assigned probability.
<b>Corollary 14.1 :</b> If the outcome $\omega$ of an experiment is in the subset $A$ , then the event $A$ is said to “have occurred”.
<b>Corollary 14.2 Complement Set</b> $A^C$ : is the contrary event of $A$ .
<b>Corollary 14.3 The Union Set</b> $A \cup B$ : Let $A, B$ be to evenest. The event “ $A$ or $B$ ” is interpreted as the union of both.
<b>Corollary 14.4 The Intersection Set</b> $A \cap B$ : Let $A, B$ be to evenest. The event “ $A$ and $B$ ” is interpreted as the intersection of both.
<b>Corollary 14.5 The Elementary Event</b> $\omega$ : Is a “singleton”, i.e. a subset $\{\omega\}$ containing a single outcome $\omega$ of $\Omega$ .
<b>Corollary 14.6 The Sure Event</b> $\Omega$ : Is equal to the sample space as it contains all possible elementary events.
<b>Corollary 14.7 The Impossible Event</b> $\emptyset$ : The impossible event i.e. nothing is happening is denoted by the empty set.
<b>Definition 14.4 The Family of All Events <math>\mathcal{A}/2^\Omega</math>:</b> The set of all subset of the sample space $\Omega$ called family of all events is given by the power set of the sample space $\mathcal{A} = 2^\Omega$ (for finite sample spaces).

<b>Definition 14.5 Probability</b> $\mathbb{P}(A)$ : Is a number associated with every $A$ , that measures the likelihood of the event to be realized “a priori”. The bigger the number the more likely the event will happen. 1. $0 \leq \mathbb{P}(A) \leq 1$ 2. $\mathbb{P}(\Omega) = 1$ 3. If $A \cap B = \emptyset$ then $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B)$
<b>Note</b> We can think of the probability of an event $A$ as the limit of the “frequency” of repeated experiments: $\mathbb{P}(A) = \lim_{n \rightarrow \infty} \frac{\delta(A)}{n} \quad \text{where} \quad \delta(A) = \begin{cases} 1 & \text{if } \omega \in A \\ 0 & \text{if } \omega \notin A \end{cases}$
<b>0.1. Sigma Algebras</b>
<b>Definition 14.6 Sigma Algebra <math>\sigma</math>:</b> A set $\mathcal{F}$ of subsets of $\Omega$ is called a $\sigma$ -algebra on $\Omega$ if the following properties apply <ul style="list-style-type: none"><li><math>\Omega \in \mathcal{F}</math> and <math>\emptyset \in \mathcal{F}</math></li><li>If <math>A \in \mathcal{F}</math> then <math>\Omega \setminus A = A^C \in \mathcal{F}</math>: The complementary subset of <math>A</math> is also in <math>\Omega</math>.</li><li>For all <math>A_i \in \mathcal{F} : \bigcup_{i=1} A_i \in \mathcal{F}</math></li></ul> See example 14.3.
<b>Corollary 14.8 <math>\mathcal{F}_{\min}</math>:</b> $\mathcal{F} = \{\emptyset, \Omega\}$ is the simplest $\sigma$ -algebra, telling us only if an event happened $\omega \in \Omega$ happened or not but not which one.
<b>Corollary 14.9 <math>\mathcal{F}_{\max}</math>:</b> $\mathcal{F} = 2^\Omega$ consists of all subsets of $\Omega$ and thus corresponds to full information i.e. we know if and which event happened.
<b>Definition 14.7 Measurable Space</b> $(\Omega, \mathcal{F})$ : Is the pair of a set and sigma algebra i.e. a sample space and sigma algebra $\{\Omega, \mathcal{F}\}$ .
<b>Corollary 14.10 <math>\mathcal{F}</math>-measurable Event:</b> The elements $A_i \in \mathcal{F}$ are called <i>measurable sets</i> or <i><math>\mathcal{F}</math>-measurable</i> .
<b>Interpretation</b> The $\sigma$ -algebra represents all of possible events of the experiment that we can detect. Thus we call the sets in $\mathcal{F}$ measurable sets/events. The sigma algebra is the mathematical construct that tells us how much information we obtain once we conduct some experiment.
<b>Definition 14.8 Sigma Algebra generated by a subset of <math>\Omega</math></b> $\sigma(\mathcal{C})$ : Let $\mathcal{C}$ be a class of subsets of $\Omega$ . The $\sigma$ -algebra generated by $\mathcal{C}$ , denoted by $\sigma(\mathcal{C})$ , is the <i>smallest</i> sigma algebra $\mathcal{F}$ that included all elements of $\mathcal{C}$ see example 14.4.
<b>Definition 14.9 Borel <math>\sigma</math>-algebra</b> $\mathcal{B}(\mathbb{R})$ : The Borel $\sigma$ -algebra $\mathcal{B}(\mathbb{R})$ is the smallest $\sigma$ -algebra containing all open intervals in $\mathbb{R}$ . The sets in contained in $\mathcal{B}(\mathbb{R})$ are called Borel sets. The extension to the multi-dimensional case, $\mathcal{B}(\mathbb{R}^n)$ , is straightforward. For all real numbers $a, b \in \mathbb{R}$ , $\mathcal{B}(\mathbb{R})$ contains various sets see example 14.5.
<b>Why do we need Borel Sets</b> So far we only looked at atomic events $\omega$ , with the help of sigma algebras we are now able to measure continuous events s.a. $[0, 1]$ .
<b>Corollary 14.11 :</b> The Borel $\sigma$ -algebra of $\mathbb{R}$ is generated by intervals of the form $(-\infty, a]$ , where $a \in \mathbb{Q}$ ( $\mathbb{Q}$ =rationals). See proof section 13.
<b>Definition 14.10 (<math>\mathbb{P}</math>)-trivial Sigma Algebra:</b> is a $\sigma$ -algebra $\mathcal{F}$ for which each event has a probability of zero or one: $\mathbb{P}(A) \in \{0, 1\} \quad \forall A \in \mathcal{F} \quad (14.1)$

<b>Interpretation</b> A trivial sigma algebra means that all events are almost surely constant and that there exist no non-trivial information. An example of a trivial sigma algebra is $\mathcal{F}_{\min} = \{\Omega, \emptyset\}$ .
<b>0.2. Measures</b>
<b>Definition 14.11 Measure</b> $\mu$ : A measure defined on a measurable space $\{\Omega, \mathcal{F}\}$ is a function/map: $\mu : \mathcal{F} \mapsto [0, \infty]$ (14.2) for which holds: <ul style="list-style-type: none"><li><math>\mu(\emptyset) = 0</math></li><li>countable additivity [def. 14.12]</li></ul>
<b>Definition 14.12 Countable/<math>\sigma</math>-Additive Function:</b> Given a function $\mu$ defined on a $\sigma$ -algebra $\mathcal{F}$ . The function $\mu$ is said to be countable additive if for every countable sequence of pairwise disjoint elements $(F_i)_{i \geq 1}$ of $\mathcal{F}$ it holds that: $\mu\left(\bigcup_{i=1}^{\infty} F_i\right) = \sum_{i=1}^{\infty} \mu(F_i) \quad \text{for all} \quad F_j \cap F_k = \emptyset \quad \forall j \neq k$ (14.3)
<b>Corollary 14.12 Additive Function:</b> A function that satisfies countable additivity, is also additive, meaning that for every $F, G \in \mathcal{F}$ it holds: $\mu(F \cup G) = \mu(F) + \mu(G) \iff F \cap G = \emptyset \quad (14.4)$
<b>Intuition</b> If we take two event that cannot occur simultaneously, then the probability that at least one vent occurs is just the sum of the measure (probabilities) of the original events.
<b>Definition 14.13 Equivalent Measures</b> $\mu \sim \nu$ : Let $\mu$ and $\nu$ be two measures defined on a measurable space [def. 14.7] $(\Omega, \mathcal{F})$ . The two measures are said to be equivalent if it holds that: $\mu(A) > 0 \iff \nu(A) > 0 \quad \forall A \subseteq \mathcal{F} \quad (14.5)$ this is equivalent to $\mu$ and $\nu$ having equivalent null sets: $\mathcal{N}_\mu = \mathcal{N}_\nu \quad \begin{matrix} \mathcal{N}_\mu = \{A \in \mathcal{A}   \mu(A) = 0\} \\ \mathcal{N}_\nu = \{A \in \mathcal{A}   \nu(A) = 0\} \end{matrix} \quad (14.6)$ see example 14.6
<b>Definition 14.14 Measure Space</b> $\{\mathcal{F}, \Omega, \mu\}$ : The triplet of sample space, sigma algebra and a measure is called a measure space.
<b>Definition 14.15 Lebesgue Measure on <math>\mathcal{B}</math></b> $\lambda$ : Is the measure defined on the measurable space $\{\mathbb{R}, \mathcal{B}(\mathbb{R})\}$ which assigns the measure of each interval to be its length: $\lambda([a, b]) = b - a \quad (14.7)$
<b>Corollary 14.13 Lebesgue Measure of Atomitics:</b> <ul style="list-style-type: none"><li>The Lebesgue measure of a set containing only one point must be zero: <math display="block">\lambda(\{a\}) = 0 \quad (14.8)</math></li><li>The Lebesgue measure of a set containing countably many points <math>A = \{a_1, a_2, \dots, a_n\}</math> must be zero: <math display="block">\lambda(A) + \sum_{i=1}^n \lambda(\{a_i\}) = 0 \quad (14.9)</math></li><li>The Lebesgue measure of a set containing uncountably many points <math>A = \{a_1, a_2, \dots\}</math> can be either zero, positive and finite or infinite.</li></ul>
<b>0.3. Probability/Kolomogorov's Axioms</b> 1931
One problem we are still having is the range of $\mu$ , by standardizing the measure we obtain a well defined measure of events.
<b>Axiom 14.1 Non-negativity:</b> The probability of an event is a non-negative real number: If $A \in \mathcal{F}$ then $\mathbb{P}(A) \geq 0 \quad (14.10)$

<b>Axiom 14.2 Unitaarity:</b> The probability that at least one of the elementary events in the entire sample space $\Omega$ will occur is equal to one: $\text{The certain event} \quad \mathbb{P}(\Omega) = 1 \quad (14.11)$
<b>Axiom 14.3 <math>\sigma</math>-additivity:</b> If $A_1, A_2, A_3, \dots \in \mathcal{F}$ are mutually disjoint, then: $\mathbb{P}\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} \mu(A_i) \quad (14.12)$
<b>Corollary 14.14 :</b> As a consequence of this it follows: $\mathbb{P}(\emptyset) = 0 \quad (14.13)$
<b>Corollary 14.15 Complementary Probability:</b> $\mathbb{P}(A^C) = 1 - \mathbb{P}(A) \quad \text{with} \quad A^C = \Omega - A \quad (14.14)$
<b>Definition 14.16 Probability Measure</b> $\mathbb{P}$ : a probability measure is function $\mathbb{P} : \mathcal{F} \mapsto [0, 1]$ defined on a $\sigma$ -algebra $\mathcal{F}$ of a sample space $\Omega$ that satisfies the probability axioms.
<b>1. Conditional Probability</b>
<b>Definition 14.17 Conditional Probability:</b> Let $A, B$ be events, with $\mathbb{P}(B) \neq 0$ . Then the conditional probability of the event $A$ given $B$ is defined as: $\mathbb{P}(A B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} \quad \mathbb{P}(B) \neq 0 \quad (14.15)$
<b>2. Independent Events</b>
<b>Theorem 14.1 Independent Events:</b> Let $A, B$ be two events. $A$ and $B$ are said to be independent iff: $\mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(B) \quad \begin{matrix} \mathbb{P}(A B) = \mathbb{P}(A), & \mathbb{P}(B) > 0 \\ \mathbb{P}(B A) = \mathbb{P}(B), & \mathbb{P}(A) > 0 \end{matrix} \quad (14.16)$
<b>Note</b> The requirement of no impossible events follows from [def. 14.17]
<b>Corollary 14.16 Pairwise Independent Evenest:</b> A finite set of events $\{A_i\}_{i=1}^n \in \mathcal{A}$ is <i>pairwise independent</i> if every pair of events is independent: $\mathbb{P}(A_i \cap A_j) = \mathbb{P}(A_i) \mathbb{P}(A_j) \quad i \neq j, \quad \forall i, j \in \mathcal{A} \quad (14.17)$
<b>Corollary 14.17 Mutal Independent Evenest:</b> A finite set of events $\{A_i\}_{i=1}^n \in \mathcal{A}$ is <i>mutal independent</i> if every event $A_j$ is independent of any intersection of the other events: $\mathbb{P}\left(\bigcap_{i=i}^k B_i\right) = \prod_{i=1}^k \mathbb{P}(B_i) \quad \begin{matrix} \forall \{B_i\}_{i=1}^k \subseteq \{A_i\}_{i=1}^n \\ k \leq n, \quad \{A_i\}_{i=1}^n \in \mathcal{A} \end{matrix} \quad (14.18)$
<b>3. Product Rule</b>
<b>Law 14.1 Product Rule:</b> Let $A, B$ be two events then the probability of both events occurring simultaneously is given by: $\mathbb{P}(A \cap B) = \mathbb{P}(B A)\mathbb{P}(A) = \mathbb{P}(A B)\mathbb{P}(B) \quad (14.19)$
<b>4. Law of Total Probability</b>
<b>Definition 14.18 Complete Event Field:</b> A complete event field $\{A_i : i \in I \subseteq \mathbb{N}\}$ is a countable or finite partition of $\Omega$ that is the partitions $\{A_i : i \in I \subseteq \mathbb{N}\}$ are a <i>disjoint union</i> the sample space: $\bigcup_{i \in I} A_i = \Omega \quad A_i \cap A_j = \emptyset \quad i \neq j, \forall i, j \in I \quad (14.20)$
<b>Theorem 14.2 Law of Total Probability/Partition Equation:</b> Let $\{A_i : i \in I\}$ be a complete event field [def. 14.18] then it holds for $B \in \mathcal{B}$ : $\mathbb{P}(B) = \sum_{i \in I} \mathbb{P}(B A_i)\mathbb{P}(A_i) \quad (14.21)$



## 5. Bayes Theorem

**Law 14.2 Bayes Rule:** Let  $A, B$  be two events s.t.  $\mathbb{P}(B) > 0$  then it holds:

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(B|A)\mathbb{P}(A)}{\mathbb{P}(B)} \quad \mathbb{P}(B) > 0 \quad (14.22)$$

follows directly from eq. (14.19).

**Theorem 14.3 Bayes Theorem:** Let  $\{A_i : i \in I\}$  be a complete event field<sup>[def. 14.18]</sup> and  $B \in \mathcal{B}$  a random event s.t.  $\mathbb{P}(B) > 0$ , then it holds:

$$\mathbb{P}(A_j|B) = \frac{\mathbb{P}(B|A_j)\mathbb{P}(A_j)}{\sum_{i \in I} \mathbb{P}(B|A_i)\mathbb{P}(A_i)} \quad (14.23)$$

proof section 13

## Distributions on $\mathbb{R}$

### 6.1. Distribution Function

**Definition 14.19 Distribution Function of  $\mathbb{P}$**   $F$ : The *distribution function*  $F$  induced by a probability  $\mathbb{P}$  on  $(\mathbb{R}, \mathcal{B})$  is the function:

$$F(x) = \mathbb{P}((-\infty, x]) \quad (14.24)$$

**Theorem 14.4** : A function  $F$  is the distribution function of a (unique) probability on  $(\mathbb{R}, \mathcal{B})$  iff:

- $F$  is non-decreasing
- $F$  is right continuous
- $\lim_{x \rightarrow -\infty} F(x) = 0$  and  $\lim_{x \rightarrow +\infty} F(x) = 1$

**Corollary 14.18** : A probability  $\mathbb{P}$  is uniquely determined by a distribution function  $F$ . That is if there exist another probability  $\mathbb{Q}$  s.t.

$$G(x) = \mathbb{Q}((-\infty, x])$$

and if  $F = G$  then it follows  $\mathbb{P} = \mathbb{Q}$ .

### 6.2. Random Variables

A random variable  $X$  is a quantity that is not a variable in the classical sense but a variable with respect to the outcome of an experiment. Thus it is actually not a variable but a function/map.

Its value is determined in two steps:

- ① The outcome of an experiment is a random quantity  $\omega \in \Omega$
- ② The outcome  $\omega$  determines (possibly various) quantities of interests  $\iff$  random variables

Thus a random variable  $X$ , defined on a probability space  $\{\Omega, \mathcal{F}, \mathbb{P}\}$  is a mapping from  $\Omega$  into another space  $\mathcal{E}$ , usually  $\mathcal{E} = \mathbb{R}$  or  $\mathcal{E} = \mathbb{R}^n$ :

$$X : \Omega \mapsto \mathcal{E} \quad \omega \mapsto X(\omega)$$

Let now  $E \in \mathcal{E}$  be a quantity of interest, in order to quantify its probability we need to map it back to the original sample space  $\Omega$ :

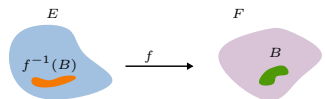
Probability for an event in  $\Omega$

$$\mathbb{P}_X(E) = \mathbb{P}(\{\omega : X(\omega) \in E\}) = \mathbb{P}(X \in E) = \mathbb{P}(X^{-1}(E))$$

Probability for an event in  $E$

**Definition 14.20  $\mathcal{E}$ -measurable function:** Let  $(E, \mathcal{E})$  and  $(F, \mathcal{F})$  be two measurable spaces. A function  $f : E \mapsto F$  is called measurable (relative to  $\mathcal{E}$  and  $\mathcal{F}$ ) if

$$\forall B \in \mathcal{F} : f^{-1}(B) = \{\omega \in E : f(\omega) \in B\} \in \mathcal{E} \quad (14.25)$$



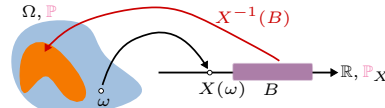
### Interpretation

The pre-image<sup>[def. 6.9]</sup> of  $B$  under  $f$  i.e.  $f^{-1}(B)$  maps all values of the target space  $F$  back to the sample space  $\mathcal{E}$  (for all possible  $B \in \mathcal{F}$ ).

**Definition 14.21 Random Variable:** A real-valued random variable (vector)  $X$ , defined on a probability space  $\{\Omega, \mathcal{E}, \mathbb{P}\}$  is an  $\mathcal{E}$ -measurable function mapping, if it maps its sample space  $\Omega$  into a target space  $(F, \mathcal{F})$ :

$$X : \Omega \mapsto \mathcal{F} \quad (\mathcal{F}^n) \quad (14.26)$$

Since  $X$  is  $\mathcal{E}$ -measurable it holds that  $X^{-1} : \mathcal{F} \mapsto \mathcal{E}$



**Corollary 14.19** : Usually  $F = \mathbb{R}$ , which usually amounts to using the Borel  $\sigma$ -algebra  $\mathcal{B}$  of  $\mathbb{R}$ .

**Corollary 14.20 Random Variables of Borel Sets:** Given that we work with Borel  $\sigma$ -algebras then the definition of a random variable is equivalent to (due to corollary 14.11):

$$X^{-1}(B) = X^{-1}((-\infty, a]) = \{\omega \in \Omega : X(\omega) \leq a\} \in \mathcal{E} \quad \forall a \in \mathbb{R} \quad (14.27)$$

### Definition 14.22

**Realization of a Random Variable**  $x = X(\omega)$ : Is the value of a random variable that is actually observed after an experiment has been conducted. In order to avoid confusion lower case letters are used to indicate actual observations/realization of a random variable.

**Corollary 14.21 Indicator Functions**  $I_A(\omega)$ : An important class of measurable functions that can be used as r.v. are indicator functions:

$$I_A(\omega) = \begin{cases} 1 & \text{if } \omega \in A \\ 0 & \text{if } \omega \notin A \end{cases} \quad (14.28)$$

We know that a probability measure  $\mathbb{P}$  on  $\mathbb{R}$  is characterized by the quantities  $\mathbb{P}((-\infty, a])$ . Thus the quantities.

**Corollary 14.22** : Let  $(F, \mathcal{F}) = (\mathbb{R}, \mathcal{B})$  and let  $(E, \mathcal{E})$  and arbitrary measurable space. Let  $X$  be a real value function on  $E$ .

Then it holds that  $X$  is measurable if and only if  $\{X \leq a\} = \{\omega : X(\omega) \leq a\} = X^{-1}((-\infty, a]) \in \mathcal{E}$ , each  $a \in \mathbb{R}$  or  $\{X < a\} \in \mathcal{E}$ .

**Explanation 14.1** (corollary 14.22). A random variable is a function that is measurable if and only if its distribution function is defined.

### 6.3. The Law of Random Variables

**Definition 14.23 Law/Distribution of  $X$ :** Let  $X$  be a r.v. on  $\{\Omega, \mathcal{F}, \mathbb{P}\}$ , with values in  $(E, \mathcal{E})$ , then the *distribution/law* of  $X$  is defined as:

$$\mathbb{P} : \mathcal{B} \mapsto [0, 1] \quad (14.29)$$

$$\mathbb{P}^X(B) = \mathbb{P}\{X \in B\} = \mathbb{P}(\omega : X(\omega) \in B) \quad \forall B \in \mathcal{E}$$

### Note

- Sometimes  $\mathbb{P}^X$  is also called the *image* of  $\mathbb{P}$  by  $X$
- The law can also be written as:

$$\mathbb{P}^X(B) = \mathbb{P}(X^{-1}(B)) = (\mathbb{P} \circ X^{-1})(B)$$

**Theorem 14.5** : The law/distribution of  $X$  is a probability measure  $\mathbb{P}$  on  $(E, \mathcal{E})$ .

### Definition 14.24

**(Cumulative) Distribution Function**  $F_X$ : Given a real-valued r.v. then its *cumulative distribution function* is defined as:

$$F_X(x) = \mathbb{P}^X((-\infty, x]) = \mathbb{P}(X \leq x) \quad (14.30)$$

**Corollary 14.23** : The distribution of  $\mathbb{P}^X$  of a real valued r.v. is entirely characterized by its cumulative distribution function  $F_X$ <sup>[def. 14.31]</sup>.

### Property 14.1:

$$\mathbb{P}(X > x) = 1 - F_X(x) \quad (14.31)$$

### Property 14.2:

$$\mathbb{P}(a < X \leq b) = F_X(b) - F_X(a) \quad (14.32)$$

### 6.4. Probability Density Function

**Definition 14.25 Continuous Random Variable:** Is a r.v. for which a probability density function  $f_X$  exists.

**Definition 14.26 Probability Density Function:** Let  $X$  be a r.v. with associated cdf  $F_X$ . If  $F_X$  is continuously integrable for all  $x \in \mathbb{R}$  then  $X$  has a *probability density*  $f_X$  defined by:

$$F_X(x) = \int_{-\infty}^x f_X(y) dy \quad (14.33)$$

or alternatively:

$$f_X(x) = \lim_{\epsilon \rightarrow 0} \frac{\mathbb{P}(x \leq X \leq x + \epsilon)}{\epsilon} \quad (14.34)$$

**Corollary 14.24**  $\mathbb{P}(X = b) = 0$ ,  $\forall b \in \mathbb{R}$ :

$$\mathbb{P}(X = b) = \lim_{a \rightarrow b} \mathbb{P}(a < X \leq b) = \lim_{a \rightarrow b} \int_a^b f(x) dx = 0 \quad (14.35)$$

**Corollary 14.25 corollary 14.24:** From corollary 14.24 it follows that the exact borders are not necessary:

$$\mathbb{P}(a < X < b) = \mathbb{P}(a \leq X < b) = \mathbb{P}(a < X \leq b) = \mathbb{P}(a \leq X \leq b)$$

### Corollary 14.26 :

$$\int_{-\infty}^{\infty} f(x) dx = 1 \quad (14.36)$$

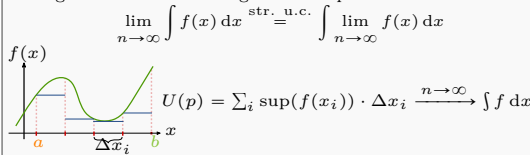
### Notes

- Often the cumulative distribution function is referred to as "cdf" or simply *distribution function*.
- Often the probability density function is referred to as "pdf" or simply *density*.

### 6.5. Lebesgue Integration

#### Problems of Riemann Integration

- Difficult to extend to higher dimensions – general domains of definitions  $f : \Omega \mapsto \mathbb{R}$
- Depends on continuity
- Integration of limit processes require strong uniform convergence in order to integrate limit processes

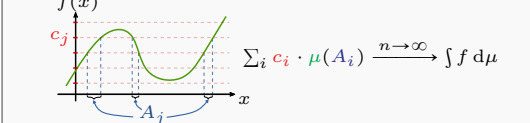


### Idea

Partition domain by function values of equal size i.e. values that lie within the same sets/have the same value  $A_j$  build up the partitions w.r.t. to the variable  $x$ .

**Problem:** we do not know how big those sets/partitions on the  $x$ -axis will be.

**Solution:** we can use the measure  $\mu$  of our measure space  $\{\Omega, \mathcal{A}, \mu\}$  in order to obtain the size of our sets  $A_j \implies$  we do not have to care anymore about discontinuities, as we can measure the size of our sets using our measure.



### Definition 14.27 Lebesgue Integral:

$$\lim_{n \rightarrow \infty} \sum_{i=1}^n c_i \mu(A_i) = \int_{\Omega} f d \mu \quad f(x) \approx c_i \quad \forall x \in A_i \quad (14.37)$$

### Definition 14.28

**Simple Functions (Random Variables):** A r.v.  $X$  is called simple if it takes on only a finite number of values and hence can be written in the form:

$$X = \sum_{i=1}^n a_i \mathbb{I}_{A_i} \quad a_i \in \mathbb{R} \quad \mathcal{A} \ni A_i = \begin{cases} 1 & \text{if } \{X = a_i\} \\ 0 & \text{else} \end{cases} \quad (14.38)$$

## 7. Independent Random Variables

We have seen that two events  $A$  and  $B$  are independent if knowledge that  $B$  has occurred does not change the probability that  $A$  will occur theorem 14.1.

For two random variables  $X, Y$  we want to know if knowledge of  $Y$  leaves the probability of  $X$ , to take on certain values unchanged.

### Definition 14.29 Independent Random Variables:

Two real valued random variables  $X$  and  $Y$  are said to be independent iff:

$$\mathbb{P}(X \leq x | Y \leq y) = \mathbb{P}(X \leq x) \quad \forall x, y \in \mathbb{R} \quad (14.39)$$

which amounts to:

$$F_{X,Y}(x, y) = \mathbb{P}(\{X \leq x\} \cap \{Y \leq y\}) = \mathbb{P}(X \leq x, Y \leq y) = F_X(x) F_Y(y) \quad \forall x, y \in \mathbb{R} \quad (14.40)$$

or alternatively iff:

$$\mathbb{P}(X \in A, Y \in B) = \mathbb{P}(X \in A) \mathbb{P}(Y \in B) \quad \forall A, B \in \mathcal{B} \quad (14.41)$$

### Note

If the joint distribution  $F_{X,Y}(x, y)$  can be factorized into two functions of  $x$  and  $y$  then  $X$  and  $Y$  are independent.

### Definition 14.30

**Independent Identically Distributed:**

## 8. Change Of Variables Formula

### Formula 14.1

**(Scalar Discret) Change of Variables:** Let  $X$  be a discret rv  $X \in \mathcal{X}$  with pmf  $p_X$  and define  $Y \in \mathcal{Y}$  as  $Y = g(x)$  s.t.  $\mathcal{Y} = \{y | y = g(x), \forall x \in \mathcal{X}\}$ . Where  $g$  is an arbitrary strictly monotonic<sup>[def. 6.12]</sup> function.

**Let:**  $\mathcal{X}_y = x_i$  be the set of all  $x_i \in \mathcal{X}$  s.t.  $y = g(x_i)$ .

Then the pmf of  $Y$  is given by:

$$p_Y(y) = \sum_{x_i \in \mathcal{X}_y} p_X(x_i) = \sum_{x \in \mathcal{Y} : g(x) = y} p_X(x) \quad (14.42)$$

see proof section 13

### Formula 14.2

**(Scalar Continuous) Change of Variables:** Let  $X \sim f_X$  be a continuous r.v. and let  $g$  be an arbitrary strictly monotonic<sup>[def. 6.12]</sup> function.

Define a new r.v.  $Y$  as

$$\mathcal{Y} = \{y | y = g(x), \forall x \in \mathcal{X}\} \quad (14.43)$$

then the pdf of  $Y$  is given by:

$$f_Y(y) = f_X(x) \left| \frac{dx}{dy} \right| = f_X(x) \left| \frac{d}{dy} (g^{-1}(y)) \right| \quad (14.44)$$

$$= f_X(x) \frac{1}{\left| \frac{dy}{dx} \right|} = \frac{f_X(g^{-1}(y))}{\left| \frac{dg}{dx} (g^{-1}(y)) \right|} \quad (14.45)$$

**Formula 14.3**  
**(Continuous) Change of Variables:**  
Let  $X = \{X_1, \dots, X_n\} \sim f_X$  be a continuous random vector and let  $g$  be an arbitrary strictly monotonic<sup>[def. 6.12]</sup> function  $g: \mathbb{R}^n \rightarrow \mathbb{R}^m$

Define a new r.v.  $Y$  as  
$$\mathcal{Y} = \{y | y = g(x), \forall x \in \mathcal{X}\} \quad (14.46)$$

and let  $h(x) := g(x)^{-1}$  then the pdf of  $Y$  is given by:  
$$\begin{aligned} f_Y(y) &= f_X(x_1, \dots, x_n) \cdot |J| \\ &= f_X(h_1(y), \dots, h_n(y)) \cdot |J| \\ &= f_X(y) |\det D_y h(x)| \Big|_{x=y}^{-1} \\ &= f_X(g^{-1}(y)) \left| \det \left( \frac{\partial g}{\partial x} \right) \right|^{-1} \end{aligned} \quad (14.47)$$

where  $J = \det Dh$  is the Jacobian<sup>[def. 7.4]</sup>.  
See also proof section 13 and example 14.8

**Note**  
A monotonic function is required in order to satisfy inevitability.

Probability Distributions on  $\mathbb{R}^n$

10. Joint Distribution

**Definition 14.31**  
**Joint (Cumulative) Distribution Function**  $F_X$ :  
Let  $X = (X_1, \dots, X_n)$  be a random vector in  $\mathbb{R}^n$ , then its cumulative distribution function is defined as:  
$$\begin{aligned} F_X(x) &= \mathbb{P}_X((-\infty, x]) = \mathbb{P}(X \leq x) \\ &= \mathbb{P}(X_1 \leq x_1, \dots, X_n \leq x_n) \end{aligned} \quad (14.48)$$

**Definition 14.32 Joint Probability Distribution:**  
Let  $X = (X_1, \dots, X_n)$  be a random vector in  $\mathbb{R}^n$  with associated cdf  $F_X$ . If  $F_X$  is continuously integrable for all  $x \in \mathbb{R}$  then  $X$  has a probability density  $f_X$  defined by:  
$$F_X(x) = \int_{-\infty}^{x_n} \dots \int_{-\infty}^{x_1} f_X(y_1, \dots, y_n) dy_1 \dots dy_n \quad (14.49)$$
  
or alternatively:  
$$f_X(x) = \lim_{\epsilon \rightarrow 0} \frac{\mathbb{P}(x_1 \leq X_1 \leq x_1 + \epsilon, \dots, x_n \leq X_n \leq x_n + \epsilon)}{\epsilon} \quad (14.50)$$

10.1. Marginal Distribution

**Definition 14.33 Marginal Distribution:**

11. The Expectation

**Definition 14.34 Expectation:**  
$$\mathbb{E}[X] = \int_{\Omega} X(\omega) \mathbb{P}(d\omega) = \int_{\Omega} X d\mathbb{P} \quad (14.51)$$

**Corollary 14.27 Expectation of simple r.v.:**  
If  $X$  is a simple<sup>[def. 14.28]</sup> r.v. its expectation is given by:  
$$\mathbb{E}[X] = \sum_{i=1}^n a_i \mathbb{P}(A_i) \quad (14.52)$$

12. Moment Generating Function (MGF)

**Definition 14.35 Moment of Random Variable:** The  $i$ -th moment of a random variable  $X$  is defined as (if it exists):  
$$m_i := \mathbb{E}[X^i] \quad (14.53)$$

**Definition 14.36**  $\psi_X$   
**Moment Generating Function (MGF):**  
$$\psi_X(t) = \mathbb{E}[e^{tX}] \quad t \in \mathbb{R} \quad (14.54)$$

**Corollary 14.28 Sum of MGF:** The moment generating function of a sum of  $n$  independent variables  $(X_j)_{1 \leq j \leq n}$  is the product of the moment generating functions of the components:  
$$\psi_{S_n}(t) = \psi_{X_1}(t) \dots \psi_{X_n}(t) \quad S_n := X_1 + \dots + X_n \quad (14.55)$$

**Corollary 14.29 :** The  $i$ -th moment of a random variable is the  $i$ -th derivative of its associated moment generating function evaluated zero:  
$$\mathbb{E}[X^i] = \psi_X^{(i)}(0) \quad (14.56)$$

13. The Characteristic Function

Transforming probability distributions using the Fourier transform is a handy tool in probability in order to obtain properties or solve problems in another space before transforming them back.

**Definition 14.37**  $\hat{\mu}$   
**Fourier Transformed Probability Measure:**  
$$\hat{\mu} = \int e^{i\langle u, x \rangle} \mu(dx) \quad (14.57)$$

**Corollary 14.30 :** As  $e^{i\langle u, x \rangle}$  can be rewritten using formulaeqs. (10.5) and (10.6) it follows:  
$$\hat{\mu} = \int \cos(\langle u, x \rangle) \mu(dx) + i \int \sin(\langle u, x \rangle) \mu(dx) \quad (14.58)$$
  
where  $x \mapsto \cos(\langle x, u \rangle)$  and  $x \mapsto \sin(\langle x, u \rangle)$  are both bounded and Borel i.e. Lebesgue integrable.

**Definition 14.38 Characteristic Function**  $\varphi_X$ : Let  $X$  be an  $\mathbb{R}^n$ -valued random variable. Its characteristic function  $\varphi_X$  is defined on  $\mathbb{R}^n$  as:  
$$\begin{aligned} \varphi_X(u) &= \int e^{i\langle u, x \rangle} \mathbb{P}^X(dx) = \widehat{\mathbb{P}^X}(u) \\ &= \mathbb{E}[e^{i\langle u, x \rangle}] \end{aligned} \quad (14.59) \quad (14.60)$$

**Corollary 14.31 :** The characteristic function  $\varphi_X$  of a distribution always exists as it is equal to the Fourier transform of the probability measure, which always exists.

**Note**  
This is an advantage over the moment generating function.

**Theorem 14.6 :** Let  $\mu$  be a probability measure on  $\mathbb{R}^n$ . Then  $\hat{\mu}$  is a bounded continuous function with  $\hat{\mu}(0) = 1$ .  
[add proof](#)

**Theorem 14.7 Uniqueness Theorem:** The Fourier Transform  $\hat{\mu}$  of a probability measure  $\mu$  on  $\mathbb{R}^n$  characterizes  $\mu$ . That is, if two probability measures on  $\mathbb{R}^n$  admit the same Fourier transform, they are equal.  
[add proof](#)

**Corollary 14.32 :** Let  $X = (X_1, \dots, X_n)$  be an  $\mathbb{R}^n$ -valued random variable. Then the real valued r.v.'s  $(X_j)_{1 \leq j \leq n}$  are independent if and only if:  
$$\varphi_X(u_1, \dots, u_n) = \prod_{j=1}^n \varphi_{X_j}(u_j) \quad (14.61)$$

Proofs

**Proof.** corollary 14.11: Let  $\mathcal{C}$  denote all open intervals. Since every open set in  $\mathbb{R}$  is the countable union of open intervals<sup>[def. 4.8]</sup>, it holds that  $\sigma(\mathcal{C})$  is the Borel  $\sigma$ -algebra of  $\mathbb{R}$ .  
Let  $\mathcal{D}$  denote all intervals of the form  $(-\infty, a]$ ,  $a \in \mathbb{Q}$ .  
Let  $a, b \in \mathcal{C}$ , and let  
•  $(a_n)_{n \geq 1}$  be a sequence of rationals decreasing to  $a$  and  
•  $(b_n)_{n \geq 1}$  be a sequence of rationals increasing strictly to  $b$   
 $(a, b) = \cup_{n=1}^{\infty} (a_n, b_n] = \cup_{n=1}^{\infty} ((-\infty, b_n] \cap (-\infty, a_n]^C)$

Thus  $\mathcal{C} \subset \sigma(\mathcal{D})$ , whence  $\sigma(\mathcal{C}) \subset \sigma(\mathcal{D})$  but as each element of  $\mathcal{D}$  is a closed subset,  $\sigma(\mathcal{D})$  must also be contained in the Borel sets  $\mathcal{B}$  with

$$\mathcal{B} = \sigma(\mathcal{C}) \subset \sigma(\mathcal{D}) \subset \mathcal{B}$$

□

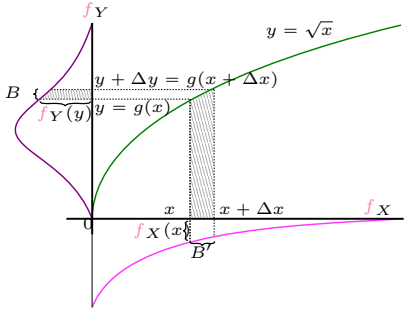
**Proof.** theorem 14.3 Plug eq. (14.21) into the denominator and eq. (14.19) into the nominator and then use<sup>[def. 14.17]</sup>:  
$$\frac{\mathbb{P}(B|A_j)\mathbb{P}(A_j)}{\sum_{i \in I} \mathbb{P}(B|A_i)\mathbb{P}(A_i)} = \frac{\mathbb{P}(B \cap A_j)}{\mathbb{P}(B)} = \mathbb{P}(A_j|B)$$

□

**Proof.** formula 14.1:  
$$Y = g(X) \iff \mathbb{P}(Y = y) = \mathbb{P}(x \in \mathcal{X}_y) = \mathbb{P}_Y(y)$$

□

**Proof.** formula 14.2 (non-formal): The probability contained in a differential area must be invariant under a change of variables that is:  
$$|f_Y(y) dy| = |f_X(x) dx|$$



□

**Proof.** formula 14.2 from CDF:  
$$\mathbb{P}(Y \leq y) = \mathbb{P}(g(X) \leq y) = \begin{cases} \mathbb{P}(X \leq g^{-1}(y)) & \text{if } g \text{ is increas.} \\ \mathbb{P}(X \geq g^{-1}(y)) & \text{if } g \text{ is decreas.} \end{cases}$$

If  $g$  is monotonically increasing:  
$$F_Y(y) = F_X(g^{-1}(y))$$
  
$$f_Y(y) = \frac{d}{dy} F_X(g^{-1}(y)) = f_X(x) \cdot \frac{d}{dy} g^{-1}(y)$$
  
If  $g$  is monotonically decreasing:  
$$F_Y(y) = 1 - F_X(g^{-1}(y))$$
  
$$f_Y(y) = \frac{d}{dy} F_X(g^{-1}(y)) = -f_X(x) \cdot \frac{d}{dy} g^{-1}(y)$$

□

**Proof.** formula 14.2: Let  $B = [x, x + \Delta x]$  and  $B' = [y, y + \Delta y] = [g(x), g(x + \Delta x)]$  we know that the probability of equal events is equal:  
$$y = g(x) \implies \mathbb{P}(y) = \mathbb{P}(g(x)) \text{ (for disc. rv.)}$$
  
Now lets consider the probability for the continuous r.v.s:  
$$\mathbb{P}(X \in B) = \int_x^{x+\Delta x} f_X(t) dt \xrightarrow{\Delta x \rightarrow 0} |\Delta x \cdot f_X(x)|$$
  
For  $y$  we use Taylor (???)  
$$g(x + \Delta x) \stackrel{\text{eq. (6.41)}}{=} g(x) + \frac{dg}{dx} \Delta y \quad \text{for } \Delta x \rightarrow 0$$
  
$$= y + \Delta y \quad \text{with } \Delta y := \frac{dg}{dx} \cdot \Delta x \quad (14.62)$$

**Thus for  $\mathbb{P}(Y \in B')$  it follows:**  
$$\begin{aligned} \mathbb{P}(X \in B') &= \int_y^{y+\Delta y} f_Y(t) dt \xrightarrow{\Delta y \rightarrow 0} |\Delta y \cdot f_Y(y)| \\ &= \left| \frac{dg}{dx}(x) \Delta x \cdot f_Y(y) \right| \end{aligned}$$

Now we simply need to related the surface of the two pdfs:  
$$B = [x, x + \Delta x] \stackrel{\text{same surfaces}}{\propto} [y, y + \Delta y] = B'$$
  
$$\mathbb{P}(Y \in B) = \mathbb{P}(X \in B')$$
  
$$\xrightarrow{\Delta y \rightarrow 0} |f_Y(y) \cdot \Delta y| = \left| f_Y(y) \cdot \frac{dg}{dx}(x) \Delta x \right| = |f_X(x) \cdot \Delta x|$$
  
$$f_Y(y) \cdot \frac{dg}{dx}(x) |\Delta x| = f_X(x) \cdot |\Delta x|$$
  
$$\implies f_Y(y) = \frac{f_X(x)}{\left| \frac{dg}{dx}(x) \right|} = \frac{f_X(g^{-1}(y))}{\left| \frac{dg}{dx} g^{-1}(y) \right|}$$

□

Examples

**Example 14.1 :**  
• Toss of a coin (with head and tail):  $\Omega = \{H, T\}$ .  
• Two tosses of a coin:  $\Omega = \{HH, HT, TH, TT\}$   
• A cubic die:  $\Omega = \{\omega_1, \omega_2, \omega_3, \omega_4, \omega_5, \omega_6\}$   
• The positive integers:  $\Omega = \{1, 2, 3, \dots\}$   
• The reals:  $\Omega = \{\omega | \omega \in \mathbb{R}\}$

**Example 14.2 :**  
• Head in coin toss  $A = \{H\}$   
• Odd number in die roll:  $A = \{\omega_1, \omega_3, \omega_5, \}$   
• The integers smaller five:  $A = \{1, 2, 3, 4\}$

**Example 14.3 :** If the sample space is a die toss  $\Omega = \{\omega_1, \omega_2, \omega_3, \omega_4, \omega_5, \omega_6\}$ , the sample space may be that we are only told whether an even or odd number has been rolled:  
$$\mathcal{F} = \{\emptyset, \{\omega_1, \omega_3, \omega_5\}, \{\omega_2, \omega_4, \omega_6\}\}$$

**Example 14.4 :** If we are only interested in the subset-set  $\mathcal{A} \in \Omega$  of our experiment, then we can look at the corresponding generating  $\sigma$ -algebra  $\sigma(\mathcal{A}) = \{\emptyset, \mathcal{A}, \mathcal{A}^C, \Omega\}$ .

**Example 14.5 :**  
• open half-lines:  $(-\infty, a)$  and  $(a, \infty)$ ,  
• union of open half-lines:  $(a, b) = (-\infty, a) \cup (b, \infty)$ ,  
• closed interval:  $[a, b] = (-\infty, \cup a) \cup (b, \infty)$ ,  
• closed half-lines:  
 $(-\infty, a] = \bigcup_{n=1}^{\infty} [a - n, a]$  and  $[a, \infty) = \bigcup_{n=1}^{\infty} [a, a + n]$ ,  
• half-open and half-closed  $(a, b] = (-\infty, b] \cup (a, \infty)$ ,  
• every set containing only one real number:  
 $\{a\} = \bigcap_{n=1}^{\infty} (a - \frac{1}{n}, a + \frac{1}{n})$ ,  
• every set containing finitely many real numbers:  
 $\{a_1, \dots, a_n\} = \bigcup_{k=1}^n \{a_k\}$ .

**Example 14.6 Equivalent (Probability) Measures:**  
$$\Omega = \{1, 2, 3\} \quad \mathbb{P}(\{1, 2, 3\}) = \{2/3, 1/6, 1/6\}$$
  
$$\tilde{\mathbb{P}}(\{1, 2, 3\}) = \{1/3, 1/3, 1/3\}$$



**Example 14.7 :**  
add example fat book p.1286  
add example prob th book 4

**Example 14.8 formula 14.2:** Let  $X, Y \stackrel{\text{ind.}}{\sim} \mathcal{N}(0, 1)$ .  
**Question:** proof that:  
$$U = X + Y \qquad V = X - 1$$
are indepent and normally distributed:  
$$h(u, v) = \begin{cases} h_1(u, v) = \frac{u+v}{2} \\ h_2(u, v) = \frac{u-v}{2} \end{cases} \quad J = \det \begin{bmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} \end{bmatrix} = -\frac{1}{2}$$
$$\begin{aligned} f_{U,V} &= f_{X,Y}(\underline{x}, \underline{y}) \cdot \frac{1}{2} \\ &\stackrel{\text{indp.}}{=} f_X(x) \cdot f_Y(y) \\ &= \frac{1}{2} \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{y^2}{2}} \\ &= \frac{1}{2} \cdot \frac{1}{\sqrt{2\pi}} e^{-\left\{\left(\frac{u+v}{2}\right)^2 + \left(\frac{u-v}{2}\right)^2 / 2\right\}} \\ &= \frac{1}{\sqrt{2\pi}\sqrt{2}} e^{-\frac{u^2}{4}} \cdot \frac{1}{\sqrt{2\pi}\sqrt{2}} e^{-\frac{v^2}{4}} \end{aligned}$$
Thus  $U, V$  are independent r.v. distributed as  $\mathcal{N}(0, 2)$ .

## Combinatorics

### 0.1. Permutations

**Definition 15.1 Permutation  $n!$ :** Given a set<sup>[def. 4.1]</sup>  $\mathcal{S}$  of  $n$  distinct objects, into how many distinct sequences/orders can we arrange/permutate those distinct objects  
$$P(\mathcal{S}) = n! \iff P(\mathcal{S}) = n \cdot (n-1) \cdot (n-2) \cdot \dots \cdot 1 \quad (15.1)$$

If there exists multiple  $n_j$  objects of the same kind within  $\mathcal{S}$  with  $j \in 1, \dots, n-1$  then we need to divide by those permutations:

$$P(\mathcal{S}) = \frac{n!}{n_1! \cdot \dots \cdot n_k} \quad \text{s.t.} \quad \sum_{i=1}^k n_i \leq n \quad (15.2)$$

#### Note

This is because the sequence/order does not change if we exchange objects of the same kind (e.g. red ball by red ball).

## Statistics

The probability that a discret random variable  $x$  is equal to some value  $\bar{x} \in \mathcal{X}$  is:

$$p_X(\bar{x}) = \mathbb{P}(x = \bar{x})$$

addappt

**Definition 16.1 Almost Surely (a.s.):** Let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a probability space. An event  $\omega \in \mathcal{F}$  happens almost surely iff  
$$\mathbb{P}(\omega) = 1 \iff \omega \text{ happens a.s.} \quad (16.1)$$

**Definition 16.2 Probability Mass Function (PMF):**

**Definition 16.3 Discrete Random Variable (DVR):** The set of possible values  $x$  of  $\mathcal{X}$  is countable of finite.  
 $\mathcal{X} = \{0, 1, 2, 3, 4, \dots, 8\} \quad \mathcal{X} = \mathbb{N} \quad (16.2)$

**Definition 16.4 Probability Density Function (PDF):** Is real function  $f: \mathbb{R}^n \rightarrow [0, \infty)$  that satisfies:

**Non-negativity:**  $f(x) \geq 0, \quad \forall x \in \mathbb{R}^n \quad (16.3)$

**Normalization:**  $\int_{-\infty}^{\infty} f(x) dx \stackrel{!}{=} 1 \quad (16.4)$

**Must be integrable**  $(16.5)$

### Note: why do we need probability density functions

A continuous random variable  $X$  can realise an infinite count of real number values within its support  $B$  (as there are an infinitude of points in a line segment).  
**Thus** we have an infinitude of values whose sum of probabilities must equal one.

Thus these probabilities must each be zero otherwise we would obtain a probability of  $\infty$ . As we can not work with zero probabilities we use the next best thing, infinitesimal probabilities (defined as a limit).

We say they are almost surely equal to zero:  
$$\mathbb{P}(X = x) = 0 \quad \text{a.s.}$$

To have a sensible measure of the magnitude of these infinitesimal quantities, we use the concept of probability density, which yields a probability mass when integrated over an interval.

**Definition 16.5 Continuous Random Variable (CRV):** A real random variable (rrv)  $X$  is said to be (absolutely) continuous if there exists a pdf (<sup>[def. 16.4]</sup>)  $f_X$  s.t. for any subset  $B \subset \mathbb{R}$  it holds:

$$\mathbb{P}(X \in B) = \int_B f_X(x) dx \quad (16.6)$$

**Property 16.1 Zero Probability:** If  $X$  is a continuous rrv (<sup>[def. 16.5]</sup>), then:

$$\mathbb{P}(X = a) = 0 \quad \forall a \in \mathbb{R} \quad (16.7)$$

**Property 16.2 Open vs. Closed Intervals:** For any real numbers  $a$  and  $b$ , with  $a < b$  it holds:

$$\begin{aligned} \mathbb{P}(a \leq X \leq b) &= \mathbb{P}(a \leq X < b) = \mathbb{P}(a < X \leq b) \\ &= \mathbb{P}(a < X < b) \end{aligned} \quad (16.8)$$

$\iff$  including or not the bounds of an interval does not modify the probability of a continuous rrv.

#### Note

Changing the value of a function at finitely many points has no effect on the value of a definite integral.

**Corollary 16.1 :** In particular for any real numbers  $a$  and  $b$  with  $a < b$ , letting  $B = [a, b]$  we obtain:

$$\mathbb{P}(a \leq X \leq b) = \int_a^b f_X(x) dx$$

**Proof.** Property 16.1:

$$\begin{aligned} \mathbb{P}(X = a) &= \lim_{\Delta x \rightarrow 0} \mathbb{P}(X \in [a, a + \Delta x]) \\ &= \lim_{\Delta x \rightarrow 0} \int_a^{a+\Delta x} f_X(x) dx = 0 \end{aligned}$$

$\square$

**Proof.** Property 16.2:

$$\begin{aligned} \mathbb{P}(a \leq X \leq b) &= \mathbb{P}(a \leq X < b) = \mathbb{P}(a < X \leq b) \\ &= \mathbb{P}(a < X < b) = \int_a^b f_X(x) dx \end{aligned}$$

$\square$

**Definition 16.6 Support of a probability density function:** The support of the density of a pdf  $f_X(\cdot)$  is the set of values of the random variable  $X$  s.t. its pdf is non-zero:  
$$\text{supp}(\cdot) f_X := \{x \in \mathcal{X} | f(x) > 0\} \quad (16.9)$$

**Note:** this is not a rigorous definition.

**Theorem 16.1 RVs are defined by a PDFs:** A probability density function  $f_X$  completely determines the distribution of a continuous real-valued random variable  $X$ .

**Corollary 16.2 Identically Distributed:** From theorem 16.1 it follows that to RV  $X$  and  $Y$  that have exactly the same pdf follow the same distribution.  
We say  $X$  and  $Y$  are **identically distributed**.

### 0.1. Cumulative Distribution Fuction

**Definition 16.7 Cumulative distribution function (CDF):** Let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a probability space.

The (cumulative) distribution function of a real-valued random variable  $X$  is the function given by:

$$F_X(x) = \mathbb{P}(X \leq x) \quad \forall x \in \mathbb{R}$$

**Property 16.3:**

**Monotonically Increasing**  $x \leq y \iff F_X(x) \leq F_X(y) \quad \forall x, y \in \mathbb{R}$  (16.10)

**Upper Limit**  $\lim_{x \rightarrow \infty} F_X(x) = 1$  (16.11)

**Lower Limit**  $\lim_{x \rightarrow -\infty} F_X(x) = 0$  (16.12)

**Definition 16.8 CDF of a discret rv  $X$ :** Let  $X$  be discret rv with pdf  $p_X$ , then the CDF of  $X$  is given by:

$$F_X(x) = \mathbb{P}(X \leq x) = \sum_{t=-\infty}^x p_X(t)$$

**Definition 16.9 CDF of a continuous rv  $X$ :** Let  $X$  be continuous rv with pdf  $f_X$ , then the CDF of  $X$  is given by:

$$F_X(x) = \int_{-\infty}^x f_X(t) dt \iff \frac{\partial F_X(x)}{\partial x} = f_X(x)$$

**Lemma 16.1 Probability Interval:** Let  $X$  be a continuous rrv with pdf  $f_X$  and cumulative distribution function  $F_X$ , then it holds that:

$$\mathbb{P}(a \leq X \leq b) = F_X(b) - F_X(a) \quad (16.13)$$

**Proof.** <sup>[def. 16.9]</sup>:

$$F_X(x) = \mathbb{P}(X \leq x) = \mathbb{P}(X \in (-\infty, x)) = \int_{-\infty}^x f_X(t) dt$$

$\square$

**Proof.** lemma 16.1:

$$\mathbb{P}(a \leq X \leq b) = \mathbb{P}(X \leq b) - \mathbb{P}(X \leq a)$$

or by the fundamental theorem of calculus (theorem 6.2):

$$\mathbb{P}(a \leq X \leq b) = \int_a^b f_X(t) dt = \int_a^b \frac{\partial F_X(t)}{\partial t} dt = [F_X(t)]_a^b$$

$\square$

**Theorem 16.2 A continuous rv is fully characterized by its CDF:** A cumulative distribution function completely determines the distribution of a continuous real-valued random variable.

### 1. Key figures

#### 1.1. The Expectation

**Definition 16.10 Expectation (disc. case):**

$$\mu_X := \mathbb{E}_x[x] := \sum_{x \in \mathcal{X}} \bar{x} p_x(\bar{x}) \quad (16.14)$$

**Definition 16.11 Expectation (cont. case):**

$$\mathbb{E}_x[x] := \int_{x \in \mathcal{X}} \bar{x} f_x(\bar{x}) d\bar{x} \quad (16.15)$$

**Law 16.1 Expectation of independent variables:**

$$\mathbb{E}[XY] = \mathbb{E}[X] \mathbb{E}[Y] \quad (16.16)$$

**Property 16.4 Translation and scaling:** If  $X \in \mathbb{R}^n$  and  $Y \in \mathbb{R}^n$  are random vectors, and  $a, b, c \in \mathbb{R}^n$  are constants then it holds:

$$\mathbb{E}[a + bX + cY] = a + b\mathbb{E}[X] + c\mathbb{E}[Y] \quad (16.17)$$

**Thus  $\mathbb{E}$  is a linear operator** (<sup>[def. 6.13]</sup>).

**Note: Expectation of the expectation**

The expectation of a r.v.  $X$  is a constant hence with Property 16.6 it follows:

$$\mathbb{E}[\mathbb{E}[X]] = \mathbb{E}[X] \quad (16.18)$$

**Property 16.5 Matrix×Expectation:** If  $X \in \mathbb{R}^n$  is a random vector and  $A \in \mathbb{R}^{m \times n}, B \in \mathbb{R}^{n \times m}$  are constant matrices then it holds:

$$\mathbb{E}[AXB] = A\mathbb{E}[(XB)] = A\mathbb{E}[X]B \quad (16.19)$$

**Proof.** eq. (16.27):

$$\begin{aligned} \mathbb{E}[XY] &= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{X,Y}(x, y) xy \\ &\stackrel{??}{=} \sum_{x \in \mathcal{X}} p_X(x) x \sum_{y \in \mathcal{Y}} p_Y(y) y = \mathbb{E}[X] \mathbb{E}[Y] \end{aligned}$$

$\square$

**Law 16.2 of the Unconscious Statistician:** Let  $X$  be a random variable  $X \in \mathcal{X}$  and define  $Y \in \mathcal{Y}$  as  $Y = g(x)$  s.t.  $\mathcal{Y} = \{y | y = g(x), \forall x \in \mathcal{X}\}$ , then  $Y$  is a random variable with expectation:

$$\mathbb{E}_Y[y] = \sum_{y \in \mathcal{Y}} y p_Y(y) = \sum_{x \in \mathcal{X}} g(x) p_X(x) \quad \text{or integral for CRV} \quad (16.20)$$

#### Consequence

Hence if we  $p_X$  we do not have to first calculate  $p_Y$  in order to calculate  $\mathbb{E}_Y[y]$ .

**Theorem 16.3 Jensen's Inequality:** If  $X$  is a random variable and  $f$  is a convex function, then it holds that:

$$f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)] \quad (16.21)$$

on the contrary if  $f$  is a concave function it follows:

$$f(\mathbb{E}[X]) \geq \mathbb{E}[f(X)] \quad (16.22)$$

**Definition 16.12 Autocorrelation/Crosscorrelation  $\gamma(t_1, t_2)$ :** Describes the covariance (<sup>[def. 16.16]</sup>) between the two values of a stochastic process  $(X_t)_{t \in T}$  at different time points  $t_1$  and  $t_2$ .

$$\gamma(t_1, t_2) = \text{Cov}[X_{t_1}, X_{t_2}] = \mathbb{E}[(X_{t_1} - \mu_{t_1})(X_{t_2} - \mu_{t_2})] \quad (16.23)$$

For zero time differences  $t_1 = t_2$  the autocorrelation functions equals the variance:

$$\gamma(t, t) = \text{Cov}[X_t, X_t] \stackrel{\text{eq. (16.41)}}{=} \mathbb{V}[X_t] \quad (16.24)$$

#### Notes

- **Hence** the autocorrelation describes the correlation of a function or signal with itself at a previous time point.
- **Given** a random time dependent variable  $x(t)$  the autocorrelation function  $\gamma(t, t - \tau)$  describes how *similar* the time translated function  $x(t - \tau)$  and the original function  $x(t)$  are.
- If there exists some relation between the values of the time series that is non-random then the autocorrelation is non-zero.
- The autocorrelation is maximized/most similar for no translation  $\tau = 0$  at all.

## 2. Key Figures

### 2.1. The Expectation

**Definition 16.13 Expectation (disc. case):**

$$\mu_X := \mathbb{E}_x[x] := \sum_{\bar{x} \in \mathcal{X}} \bar{x} p_x(\bar{x}) \quad (16.25)$$

**Definition 16.14 Expectation (cont. case):**

$$\mathbb{E}_x[x] := \int_{\bar{x} \in \mathcal{X}} \bar{x} f_x(\bar{x}) d\bar{x} \quad (16.26)$$

**Law 16.3 Expectation of independent variables:**

$$\mathbb{E}[XY] = \mathbb{E}[X] \mathbb{E}[Y] \quad (16.27)$$

**Property 16.6 Translation and scaling:** If  $X \in \mathbb{R}^n$  and  $Y \in \mathbb{R}^n$  are random vectors, and  $a, b, c \in \mathbb{R}^n$  are constants then it holds:

$$\mathbb{E}[a + bX + cY] = a + b\mathbb{E}[X] + c\mathbb{E}[Y] \quad (16.28)$$

Thus  $\mathbb{E}$  is a **linear operator**<sup>[def. 6.13]</sup>.

**Property 16.7 Affine Transformation of the Expectation:**  
If  $X \in \mathbb{R}^n$  is a random vector,  $A \in \mathbb{R}^{m \times n}$  a constant matrix and  $b \in \mathbb{R}^m$  then it holds:

$$\mathbb{E}[AX + b] = A\mu + b \quad (16.29)$$

**Note: Expectation of the expectation**

The expectation of a r.v.  $X$  is a constant hence with Property 16.6 it follows:

$$\mathbb{E}[\mathbb{E}[X]] = \mathbb{E}[X] \quad (16.30)$$

**Property 16.8 Matrix×Expectation:** If  $X \in \mathbb{R}^n$  is a random vector and  $A \in \mathbb{R}^{m \times n}$ ,  $B \in \mathbb{R}^{n \times m}$  are constant matrices then it holds:

$$\mathbb{E}[AXB] = A\mathbb{E}[(XB)] = A\mathbb{E}[X]B \quad (16.31)$$

*Proof.* eq. (16.27):

$$\begin{aligned} \mathbb{E}[XY] &= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{X,Y}(x, y) xy \\ &\stackrel{??}{=} \sum_{x \in \mathcal{X}} p_X(x) x \sum_{y \in \mathcal{Y}} p_Y(y) y = \mathbb{E}[X] \mathbb{E}[Y] \end{aligned}$$

**Law 16.4 of the Unconscious Statistician:** Let  $X$  be a random variable  $X \in \mathcal{X}$  and define  $Y \in \mathcal{Y}$  as  $Y = g(x)$  s.t.  $\mathcal{Y} = \{y|y = g(x), \forall x \in \mathcal{X}\}$ , then  $Y$  is a random variable with expectation:

$$\mathbb{E}_Y[y] = \sum_{y \in \mathcal{Y}} y p_Y(y) = \sum_{x \in \mathcal{X}} g(x) p_X(x) \quad \text{or integral for CRV} \quad (16.32)$$

**Consequence**

Hence if we  $p_X$  we do not have to first calculate  $p_Y$  in order to calculate  $\mathbb{E}_Y[y]$ .

**Theorem 16.4 Jensen's Inequality:** If  $X$  is a random variable and  $f$  is a convex function, then it holds that:

$$f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)] \quad (16.33)$$

on the contrary if  $f$  is a concave function it follows:

$$f(\mathbb{E}[X]) \geq \mathbb{E}[f(X)] \quad (16.34)$$

### 2.2. The Variance

**Definition 16.15 Variance  $\mathbb{V}[X]$ :** The variance of a random variable  $X$  is the expected value of the squared deviation from the expectation of  $X$  ( $\mu = \mathbb{E}[X]$ ).  
It is a measure of how much the actual values of a random variable  $X$  fluctuate around its expected value  $\mathbb{E}[X]$  and is defined by:

$$\mathbb{V}[X] := \mathbb{E}[(X - \mathbb{E}[X])^2] \stackrel{\text{see section 3}}{=} \mathbb{E}[X^2] - \mathbb{E}[X]^2 \quad (16.35)$$

### 2.2.1. Properties

**Property 16.9 Variance of a Constant:** If  $a \in \mathbb{R}$  is a constant then it follows that its expected value is deterministic  $\Rightarrow$  we have no uncertainty  $\Rightarrow$  no variance:

$$\mathbb{V}[a] = 0 \quad \text{with} \quad a \in \mathbb{R} \quad (16.36)$$

see shift and scaling for proof section 3

**Property 16.10 Shifting and Scaling:**

$$\mathbb{V}[a + bX] = a^2 \sigma^2 \quad \text{with} \quad a \in \mathbb{R} \quad (16.37)$$

see section 3

**Property 16.11 Affine Transformation of the Variance:**  
If  $X \in \mathbb{R}^n$  is a random vector,  $A \in \mathbb{R}^{m \times n}$  a constant matrix and  $b \in \mathbb{R}^m$  then it holds:

$$\mathbb{V}[AX + b] = A\mathbb{V}[X]A^T \quad (16.38)$$

see section 3.

**Definition 16.16 Covariance:** The Covariance is a measure of how much two or more random variables vary **linearly** with each other.

$$\begin{aligned} \text{Cov}[X, Y] &= \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] \\ &= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] \end{aligned} \quad (16.39)$$

see section 3

**Definition 16.17 Covariance Matrix:** The variance of a  $k$ -dimensional random vector  $X = (X_1 \dots X_k)$  is given by a p.s.d. eq. (9.11) matrix called Covariance Matrix.  
The Covariance is a measure of how much two or more random variables vary **linearly** with each other and the Variance on the diagonal is again a measure of how much a variable varies:

$$\begin{aligned} \mathbb{V}[X] := \Sigma(X) &:= \text{Cov}[X, X] := \\ &= \mathbb{E}[(X - \mathbb{E}[X])(X - \mathbb{E}[X])^T] \\ &= \mathbb{E}[XX^T] - \mathbb{E}[X]\mathbb{E}[X]^T \in [-\infty, \infty] \end{aligned} \quad (16.40)$$

$$\begin{aligned} &= \begin{bmatrix} \mathbb{V}[X_1] & \dots & \text{Cov}[X_1, X_k] \\ \vdots & \ddots & \vdots \\ \text{Cov}[X_k, X_1] & \dots & \mathbb{V}[X_k] \end{bmatrix} \\ &= \begin{bmatrix} \mathbb{E}[(X_1 - \mu_1)(X_1 - \mu_1)] & \dots & \mathbb{E}[(X_1 - \mu_1)(X_k - \mu_k)] \\ \vdots & \ddots & \vdots \\ \mathbb{E}[(X_k - \mu_k)(X_1 - \mu_1)] & \dots & \mathbb{E}[(X_k - \mu_k)(X_k - \mu_k)] \end{bmatrix} \end{aligned}$$

**Note: Covariance and Variance**

The variance is a special case of the covariance in which two variables are identical:

$$\text{Cov}[X, X] = \mathbb{V}[X] \equiv \sigma^2(X) \equiv \sigma_X^2 \quad (16.41)$$

[add http://www.visiondummy.com/2014/04/geometric-interpretation-covariance-matrix/](http://www.visiondummy.com/2014/04/geometric-interpretation-covariance-matrix/)

**Property 16.12 Translation and Scaling:**

$$\text{Cov}(a + bX, c + dY) = bd\text{Cov}(X, Y) \quad (16.42)$$

**Property 16.13 Affine Transformation of the Covariance:**  
If  $X \in \mathbb{R}^n$  is a random vector,  $A \in \mathbb{R}^{m \times n}$  a constant matrix and  $b \in \mathbb{R}^m$  then it holds:

$$\text{Cov}[AX + b] = A\mathbb{V}[X]A^T = A\Sigma(X)A^T \quad (16.43)$$

**Definition 16.18 Correlation Coefficient:** Is the standardized version of the covariance:

$$\begin{aligned} \text{Corr}[X] &:= \frac{\text{Cov}[X]}{\sigma_{X_1} \dots \sigma_{X_k}} \in [-1, 1] \quad (16.44) \\ &= \begin{cases} +1 & \text{if } Y = aX + b \text{ with } a > 0, b \in \mathbb{R} \\ -1 & \text{if } Y = aX + b \text{ with } a < 0, b \in \mathbb{R} \end{cases} \end{aligned}$$

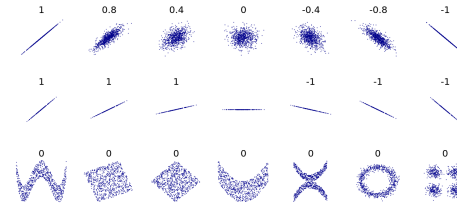


Figure 1: Several sets of  $(x, y)$  points, with their correlation coefficient

**Law 16.5 Translation and Scaling:**

$$\text{Corr}(a + bX, c + dY) = \text{sign}(b)\text{sign}(d)\text{Cov}(X, Y) \quad (16.45)$$

**Note**

- The correlation/covariance reflects the noisiness and direction of a linear relationship (top row fig. 1), **but** not the slope of that relationship (middle row fig. 1) nor many aspects of nonlinear relationships (bottom row)
- The set in the center of fig. 1 has a slope of 0 but in that case the correlation coefficient is undefined because the variance of  $Y$  is zero.
- Zero covariance/correlation  $\text{Cov}(X, Y) = \text{Corr}(X, Y) = 0$  implies that there does not exist a **linear** relationship between the random variables  $X$  and  $Y$ .

**Difference Covariance&Correlation**

- Variance is affected by scaling and covariance not ?? and law 16.5.
- Correlation is dimensionless, whereas the unit of the covariance is obtained by the product of the units of the two RV variables.

**Law 16.6 Covariance of independent RVs:** The covariance/correlation of two independent variable's (??) is zero:

$$\begin{aligned} \text{Cov}[X, Y] &= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] \\ &\stackrel{\text{eq. (16.27)}}{=} \mathbb{E}[X]\mathbb{E}[Y] - \mathbb{E}[X]\mathbb{E}[Y] = 0 \end{aligned}$$

**Zero covariance/correlation  $\Rightarrow$  independence**

$$\text{Cov}(X, Y) = \text{Corr}(X, Y) = 0 \Rightarrow p_{X,Y}(x, y) = p_X(x)p_Y(y)$$

**For example:** let  $X \sim \mathcal{U}([-1, 1])$  and let  $Y = X^2$ .

- Clearly  $X$  and  $Y$  are **dependent**
- But** the covariance/correlation between  $X$  and  $Y$  is non-zero:

$$\begin{aligned} \text{Cov}(X, Y) &= \text{Cov}(X, X^2) = \mathbb{E}[X \cdot X^2] - \mathbb{E}[X]\mathbb{E}[X^2] \\ &= \mathbb{E}[X^3] - \mathbb{E}[X]\mathbb{E}[X^2] \stackrel{\text{eq. (16.67)}}{=} 0 - 0 \cdot \mathbb{E}[X^2] \\ &\stackrel{\text{eq. (16.56)}}{=} 0 \end{aligned}$$

$\Rightarrow$  the relationship between  $Y$  and  $X$  must be non-linear.

**Definition 16.19 Quantile:** Are specific values  $q_\alpha$  in the range<sup>[def. 6.8]</sup> of a random variable  $X$  that are defined as the value for which the cumulative probability is less than  $\alpha \in (0, 1)$ :

$$q_\alpha : \mathbb{P}(X \leq x) = F_X(q_\alpha) = \alpha \xrightarrow{F \text{ invert.}} q_\alpha = F_X^{-1}(\alpha) \quad (16.46)$$

[add figure](#)

## 3. Proofs

*Proof.* eq. (16.35)

$$\begin{aligned} \mathbb{V}[X] &= \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2 - 2X\mathbb{E}[X] + \mathbb{E}[X]^2] \\ &\stackrel{\text{Property 16.6}}{=} \mathbb{E}[X^2] - 2\mathbb{E}[X]\mathbb{E}[X] + \mathbb{E}[X]^2 = \mathbb{E}[X^2] - \mu^2 \end{aligned}$$

*Proof.* Property 16.10

$$\begin{aligned} \mathbb{V}[a + bX] &= \mathbb{E}[(a + bX - \mathbb{E}[a + bX])^2] \\ &= \mathbb{E}[(a + bX - a - b\mathbb{E}[X])^2] \\ &= \mathbb{E}[(bX - b\mathbb{E}[X])^2] \\ &= \mathbb{E}[b^2(X - \mathbb{E}[X])^2] \\ &= b^2\mathbb{E}[(X - \mathbb{E}[X])^2] = b^2\sigma^2 \end{aligned}$$

*Proof.* Property 16.11

$$\begin{aligned} \mathbb{V}(AX + b) &= \mathbb{E}[(AX - \mathbb{E}[AX])^2] + 0 = \\ &= \mathbb{E}[(AX - \mathbb{E}[AX])(AX - \mathbb{E}[AX])^T] \\ &= \mathbb{E}[A(X - \mathbb{E}[X])(X - \mathbb{E}[X])^T A^T] \\ &= \mathbb{E}[A(X - \mathbb{E}[X])(X - \mathbb{E}[X])^T A^T] \\ &= A\mathbb{E}[(X - \mathbb{E}[X])(X - \mathbb{E}[X])^T] A^T = A\mathbb{V}[X]A^T \end{aligned}$$

*Proof.* eq. (16.39)

$$\begin{aligned} \text{Cov}[X, Y] &= \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] \\ &= \mathbb{E}[XY - X\mathbb{E}[Y] - \mathbb{E}[X]Y + \mathbb{E}[X]\mathbb{E}[Y]] \\ &= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] - \mathbb{E}[X]\mathbb{E}[Y] + \mathbb{E}[X]\mathbb{E}[Y] \\ &= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] \end{aligned}$$

Discrete Distributions

**Definition 16.20 Multivariate Distribution:** the variate refers to the number of input variables i.e. a m-variate distribution has m-input variables whereas a uni-variate distribution has only one.

Dimensional vs. Multivariate

The dimension refers to the number of dimensions we need to embed the function. If the variables of a function are independent than the dimension is the same as the number of inputs but the number of input variables can also be less.

4.1. Bernoulli Distribution Bern(p)

**Definition 16.21 Bernoulli Trial:** Is a random experiment with exactly two possible outcomes, success (1) and failure (0), in which the probability of success/failure is constant in every trial i.e. independent trials.

**Definition 16.22 Bernoullidistribution**  $X \sim \text{Bern}(\mathbf{p})$ :  $X$  is a binary variable i.e. can only attain the values 0 (failure) or 1 (success) with a parameter  $\mathbf{p}$  that signifies the success probability:

$$\mathbf{p}(x; \mathbf{p}) = \begin{cases} \mathbf{p} & \text{for } x = 1 \\ 1 - \mathbf{p} & \text{for } x = 0 \end{cases} \iff \begin{cases} \mathbb{P}(X = 1) = \mathbf{p} \\ \mathbb{P}(X = 0) = 1 - \mathbf{p} \end{cases}$$
$$= \mathbf{p}^x \cdot (1 - \mathbf{p})^{1-x} \quad \text{for } x \in \{0, 1\}$$

$$\mathbb{E}[X] = \mathbf{p} \quad (16.47) \quad \mathbb{V}[X] = \mathbf{p}(1 - \mathbf{p}) \quad (16.48)$$

4.2. Binomial Distribution B(n, p)

**Definition 16.23 Binomial Distribution:** Models the probability of exactly  $X$  success given a fixed number  $n$ -Bernoulli experiments<sup>[def. 16.21]</sup>, where the probability of success of a single experiment is given by  $\mathbf{p}$ :

$$\mathbf{p}(x) = \binom{n}{x} \mathbf{p}^x (1 - \mathbf{p})^{n-x} \quad \begin{array}{l} n : \text{nb. of repetitions} \\ x : \text{nb. of successes} \\ \mathbf{p} : \text{probability of success} \end{array}$$
$$\mathbb{E}[X] = n\mathbf{p} \quad (16.49) \quad \mathbb{V}[X] = n\mathbf{p}(1 - \mathbf{p}) \quad (16.50)$$

Note: Binomial Coefficient

The Binomial Coefficient corresponds to the permutation of two classes and not the variations as it seems from the formula.

Lets consider a box of  $n$  balls consisting of black and white balls. If we want to know the probability of drawing first  $x$  white and then  $n - x$  black balls we can simply calculate:

$$\underbrace{(\mathbf{p} \cdots \mathbf{p})}_{x\text{-times}} \cdot \underbrace{(q \cdots q)}_{n-x\text{-times}} = \mathbf{p}^x q^{n-x}$$

But there exists obviously further realization  $X = x$ , that correspond to permutations of the  $n$ -drawn balls.

There exist two classes of  $n_1 = x$ -white and  $n_2 = (n - x)$  black balls s.t.

$$P(n; n_1, n_2) = \frac{n!}{x!(n - x)!} = \binom{n}{x}$$

4.3. Geometric Distribution Geom(p)

**Definition 16.24 Geometric Distribution**Geom(p): Models the probability of the number  $X$  of Bernoulli trials<sup>[def. 16.21]</sup> until the first success

$$\mathbf{p}(x) = \mathbf{p}(1 - \mathbf{p})^{x-1} \quad \begin{array}{l} x : \text{nb. of repetitions until first success} \\ \mathbf{p} : \text{success probability of single Bernoulli experiment} \end{array}$$

$$F(x) = \sum_{i=1}^x \mathbf{p}(1 - \mathbf{p})^{i-1} \stackrel{??}{=} 1 - (1 - \mathbf{p})^x$$

$$\mathbb{E}[X] = \frac{1}{\mathbf{p}} \quad (16.51) \quad \mathbb{V}[X] = \frac{1 - \mathbf{p}}{\mathbf{p}^2} \quad (16.52)$$

Notes

- $\mathbb{E}[X]$  is the mean waiting time until the first success
- the number of trials  $x$  in order to have at least one success with a probability of  $\mathbf{p}(x)$ :

$$x \geq \frac{\mathbf{p}(x)}{1 - \mathbf{p}}$$

- $\log(1 - \mathbf{p}) \approx -\mathbf{p}$  for small  $\mathbf{p}$

4.4. Poisson Distribution Pois(λ)

**Definition 16.25 Poisson Distribution:** Is an extension of the binomial distribution, where the realization  $x$  of the random variable  $X$  may attain values in  $\mathbb{Z}_{\geq 0}$ .

It expresses the probability of a given number of events  $X$  occurring in a fixed interval if those events occur independently of the time since the last event.

$$\mathbf{p}(x) = e^{-\lambda} \frac{\lambda^x}{x!} \quad \begin{array}{l} \lambda > 0 \\ x \in \mathbb{Z}_{\geq 0} \end{array} \quad (16.53)$$

**Event Rate  $\lambda$ :** describes the average number of events in a single interval.

$$\mathbb{E}[X] = \lambda \quad (16.54) \quad \mathbb{V}[X] = \lambda \quad (16.55)$$

Continuous Distributions

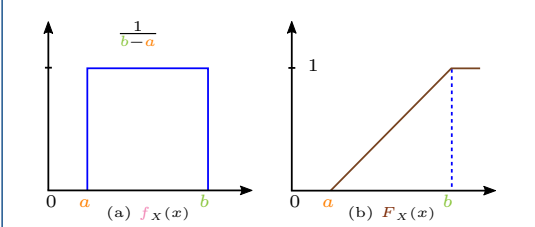
5.1. Uniform Distribution U(a, b)

**Definition 16.26 Uniform Distribution**  $\mathcal{U}(a, b)$ : Is probability distribution, where all intervals of the same length on the distribution's support<sup>(def. 16.6)</sup>  $\text{supp}(\mathcal{U}[a, b]) = [a, b]$  are equally probable/likely.

$$\mathbf{f}(x) = \frac{1}{b - a} \mathbb{1}_{x \in [a; b]} = \begin{cases} \frac{1}{b - a} = \text{const} & \text{if } a \leq x \leq b \\ 0 & \text{else} \end{cases} \quad (16.56)$$

$$F(x) = \begin{cases} 0 & x < a \\ \frac{x - a}{b - a} & \text{if } a \leq x \leq b \\ 1 & x > b \end{cases} \quad (16.57)$$

$$\mathbb{E}[X] = \frac{a + b}{2} \quad \mathbb{V}(X) = \frac{(b - a)^2}{12} \quad (16.58)$$



5.2. Exponential Distribution exp(λ)

**Definition 16.27 Exponential Distribution**  $X \sim \exp(\lambda)$ : Is the continuous analogue to the geometric distribution<sup>[def. 16.24]</sup>.

It describes the probability  $\mathbf{f}(x; \lambda)$  that a continuous Poisson process (i.e., a process in which events occur continuously and independently at a constant average rate) will succeed/change state after a time interval  $x$ .

$$\mathbf{f}(x; \lambda) = \begin{cases} \lambda e^{-\lambda x} & \text{if } x \geq 0 \\ 0 & x < 0 \end{cases} \quad (16.59)$$

$$F(x; \lambda) = \begin{cases} 1 - e^{-\lambda x} & \text{if } x \geq 0 \\ 0 & x < 0 \end{cases} \quad (16.60)$$

$$\mathbb{E}[X] = \frac{1}{\lambda} \quad \mathbb{V}(X) = \frac{1}{\lambda^2} \quad (16.61)$$

5.3. Laplace Distribution

**Definition 16.28 Laplace Distribution:**

$$\text{Laplace Distribution} \quad \mathbf{f}(x; \mu, \sigma) = \frac{1}{2\sigma} \exp\left(-\frac{|x - \mu|}{\sigma}\right) \quad (16.62)$$

5.4. The Normal Distribution N(μ, σ)

**Definition 16.29 Normal Distribution**  $X \sim \mathcal{N}(\mu, \sigma^2)$ : Is a symmetric distribution where the population parameters  $\mu, \sigma^2$  are equal to the expectation and variance of the distribution:

$$\mathbb{E}[X] = \mu \quad \mathbb{V}(X) = \sigma^2 \quad (16.63)$$

$$\mathbf{f}(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2}\left(\frac{x - \mu}{\sigma}\right)^2\right\} \quad (16.64)$$

$$F(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^x \exp\left\{-\frac{1}{2}\left(\frac{u - \mu}{\sigma}\right)^2\right\} du \quad (16.65)$$
$$x \in \mathbb{R} \quad \text{or} \quad -\infty < x < \infty$$

$$\varphi_X(u) = \exp\left\{iu\mu - \frac{u^2\sigma^2}{2}\right\} \quad (16.66)$$

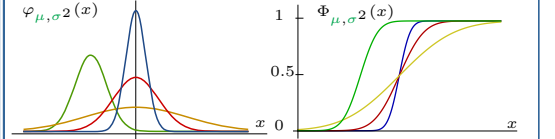


Figure 3:  $\mu = 0 \quad \sigma^2 = 0.2$   $\mu = 0 \quad \sigma^2 = 1.0$   $\mu = 0 \quad \sigma^2 = 5.0$   $\mu = -2 \quad \sigma^2 = 0.5$

**Property 16.14:**  $\mathbb{P}_X(\mu - \sigma \leq x \leq \mu + \sigma) = 0.66$

**Property 16.15:**  $\mathbb{P}_X(\mu - 2\sigma \leq x \leq \mu + 2\sigma) = 0.95$

5.5. The Standard Normal distribution N(0, 1)

**Historic Problem:** the cumulative distribution eq. (16.65) does not have an analytical solution and numerical integration was not always computationally so easy. So how should people calculate the probability of  $x$  falling into certain ranges  $\mathbb{P}(x \in [a, b])$ ?

**Solution:** use a standardized form/set of parameters (by convention)  $\mathcal{N}_{0,1}$  and tabulate many different values for its cumulative distribution  $\phi(x)$  s.t. we can transform all families of Normal Distributions into the standardized version  $\mathcal{N}(\mu, \sigma^2) \xrightarrow{z} \mathcal{N}(0, 1)$  and look up the value in its table.

Definition 16.30

**Standard Normal Distribution**  $X \sim \mathcal{N}(0, 1)$ :

$$\mathbb{E}[X] = 0 \quad \mathbb{V}(X) = 1 \quad (16.67)$$

$$\mathbf{f}(x; 0, 1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} \quad (16.68)$$

$$F(x; 0, 1) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{1}{2}u^2} du \quad (16.69)$$
$$x \in \mathbb{R} \quad \text{or} \quad -\infty < x < \infty$$

$$\psi_X(u) = e^{-\frac{u^2}{2}} \quad \varphi_X(u) = e^{-\frac{u^2}{2}} \quad (16.70)$$

Corollary 16.3

**Standard Normal Distribution Notation:** As the standard normal distribution is so commonly used people often use the letter  $Z$  in order to denote its the *standard* normal distribution and its  $\alpha$ -quantile<sup>[def. 16.19]</sup> is then denoted by:

$$z_\alpha = \Phi^{-1}(\alpha) \quad \alpha \in (0, 1) \quad (16.71)$$

5.5.1. Calculating Probabilities

**Property 16.16 Symmetry:** Let  $z > 0$

$$\mathbb{P}(Z \leq z) = \Phi(z) \quad (16.72)$$

$$\mathbb{P}(Z \leq -z) = \Phi(-z) = 1 - \Phi(z) \quad (16.73)$$

$$\mathbb{P}(-a \leq Z \leq b) = \Phi(b) - \Phi(-a) = \Phi(b) - (1 - \Phi(a))$$
$$\stackrel{a=b=z}{=} 2\Phi(z) - 1 \quad (16.74)$$

### 5.5.2. Linear Transformations of Normal Dist.

**Proposition 16.1 Linear Transformation:** Let  $X$  be a normally distributed random variable  $X \sim \mathcal{N}(\mu, \sigma^2)$ , then the linear transformed r.v.  $Y = a + bX$  is distributed as:

$$Y \sim \mathcal{N}\left(a + b\mu, b^2\sigma^2\right) \iff f_Y(y) = \frac{1}{|b|} f_X\left(\frac{y-a}{b}\right) \quad (16.75)$$

section 1

**Proposition 16.2 Standardization:** Let  $X$  be a normally distributed random variable  $X \sim \mathcal{N}(\mu, \sigma^2)$ , then there exists a linear transformation  $Z = a + bX$  s.t.  $Z$  is a standard normally distributed random variable:

$$X \sim \mathcal{N}(\mu, \sigma^2) \xrightarrow{Z = \frac{X-\mu}{\sigma}} Z \sim \mathcal{N}(0, 1) \quad (16.76)$$

section 1

#### Note

If we know how many standard deviations our distribution is away from our target value then we can characterize it fully by the standard normal distribution.

**Proposition 16.3 Standardization of the CDF:** Let  $F_X(X)$  be the cumulative distribution function of a normally distributed random variable  $X \sim \mathcal{N}(\mu, \sigma^2)$ , then the cumulative distribution function  $\Phi_Z(z)$  of the standardized random normal variable  $Z \sim \mathcal{N}(0, 1)$  is related to  $F_X(X)$  by:

$$F_X(x) = \Phi\left(\frac{x-\mu}{\sigma}\right) \quad (16.77)$$

section 1

## 6. The Multivariate Normal distribution

### Definition 16.31

**Multivariate Normal distribution**  $X \sim \mathcal{N}(\mu, \Sigma)$ :

The  $k$ -multivariate Normal distribution of:

$\mathbf{X} = (x_1 \dots x_k)^\top$  a  $k$ -dimensional random vector with:

$\mu = (\mathbb{E}[x_1] \dots \mathbb{E}[x_k])^\top$  a  $k$ -dim mean vector

and  $k \times k$  **p.s.d.** covariance matrix:

$\Sigma := \mathbb{E}[(\mathbf{X} - \mu)(\mathbf{X} - \mu)^\top] = [\text{Cov}[x_i, x_j], 1 \leq i, j \leq k]$

is given by:

$$\begin{aligned} f_{\mathbf{X}}(x_1, \dots, x_k) &= \frac{1}{\sqrt{(2\pi)^k \det(\Sigma)}} \exp\left(-\frac{1}{2}(\mathbf{X} - \mu)^\top \Sigma^{-1}(\mathbf{X} - \mu)\right) \\ &\quad \text{Normalisation} \end{aligned} \quad (16.78)$$

### Definition 16.32 Jointly Gaussian Random Variables:

Two random variables  $U, V$  both scalars or vectors, are said to be **jointly Gaussian** if the joint vector random variable  $\mathbf{X} = [U \quad V]^\top$  is again a GRV.

**Corollary 16.4 Jointly GRV of GRVs:** If  $\mathbf{x}$  and  $\mathbf{y}$  are both independent GRVs  $\mathbf{x} \sim \mathcal{N}(\mu_x, \Sigma_x)$ ,  $\mathbf{y} \sim \mathcal{N}(\mu_y, \Sigma_y)$ , then they are jointly Gaussian <sup>(def. 16.32)</sup>.

$$\begin{aligned} p(\mathbf{x}, \mathbf{y}) &= p(\mathbf{x})p(\mathbf{y}) \\ &= \exp\left(-\frac{1}{2}\left\{(\mathbf{x} - \mu_x)^\top \Sigma_x^{-1}(\mathbf{x} - \mu_x) + (\mathbf{y} - \mu_y)^\top \Sigma_y^{-1}(\mathbf{y} - \mu_y)\right\}\right) \\ &= \exp\left(-\frac{1}{2}\left[\begin{pmatrix} \mathbf{x} - \mu_x & \mathbf{y} - \mu_y \end{pmatrix}^\top \begin{pmatrix} \Sigma_x^{-1} & 0 \\ 0 & \Sigma_y^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{x} - \mu_x \\ \mathbf{y} - \mu_y \end{pmatrix}\right]\right) \end{aligned} \quad (16.79)$$

**Property 16.17 Scalar Affine Transformation of GRVs:** Let  $\mathbf{y} \in \mathbb{R}^n$  be GRV,  $\mathbf{a} \in \mathbb{R}_+, b \in \mathbb{R}$  and let  $\mathbf{x}$  be defined by the **affine transformation** <sup>(def. 9.1)</sup>:

$$\begin{aligned} \mathbf{x} &= \mathbf{a}\mathbf{y} + b\mathbf{a} \quad \mathbf{a} \in \mathbb{R}_+, b \in \mathbb{R}^d \\ \text{Then } \mathbf{x} \text{ is a GRV with:} \\ \mathbf{x} &\sim \mathcal{N}(\mu_x, \sigma_x^2) = \mathcal{N}(a\mu + b, a^2\sigma^2) \end{aligned} \quad (16.80)$$

**Property 16.18 Affine Transformation of GRVs:** Let  $\mathbf{y} \in \mathbb{R}^n$  be GRV,  $\mathbf{A} \in \mathbb{R}^{d \times n}$ ,  $b \in \mathbb{R}^d$  and let  $\mathbf{x}$  be defined by the **affine transformation** <sup>(def. 9.1)</sup>:

$$\mathbf{x} = \mathbf{A}\mathbf{y} + b \quad \mathbf{A} \in \mathbb{R}^{d \times n}, b \in \mathbb{R}^d$$

Then  $\mathbf{x}$  is a GRV (see Section 1).

**Property 16.19 Linear Combination of jointly GRVs:** Let  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{y} \in \mathbb{R}^m$  two jointly GRVs, and let  $\mathbf{z}$  be defined as:

$$\mathbf{z} = \mathbf{A}_x \mathbf{x} + \mathbf{A}_y \mathbf{y} \quad \mathbf{A}_x \in \mathbb{R}^{d \times n}, \mathbf{A}_y \in \mathbb{R}^{d \times m}$$

Then  $\mathbf{z}$  is GRV (see Section 1).

#### Note

- Joint vs. multivariate:** a joint normal distribution can be a multivariate normal distribution or a product of univariate normal distributions **but**
- Multivariate refers to the number of variables that are placed as inputs to a function.

#### Diagonal Covariance Matrix

For i.i.d. data the covariance matrix becomes diagonal:

$$\Sigma = \begin{bmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_k^2 \end{bmatrix} \quad \text{and} \quad \mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_k \end{bmatrix} \quad (16.81)$$

eq. (16.78) decomposed s.t.  $x_1, \dots, x_k$  become **mutal independent** (??):

$$p(\mathbf{X}) = \prod_{i=1}^k \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{(x_i - \mu_i)^2}{2\sigma_i^2}\right) \quad (16.82)$$

### 6.1. Gamma Distribution

$\Gamma(x, \alpha, \beta)$

**Definition 16.33 Gamma Distribution**  $X \sim \Gamma(x, \alpha, \beta)$ :

Is a widely used distribution that is related to the exponential distribution, Erlang distribution, and chi-squared distribution as well as Normal distribution:

$$f(x; \alpha, \beta) = \begin{cases} \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases} \quad (16.83)$$

$$\Gamma(\alpha) \stackrel{\text{eq. (6.66)}}{=} \int_0^\infty t^{\alpha-1} e^{-t} dt \quad (16.84)$$

with  $\alpha, \beta \in \mathbb{R}_{>0}$

## 7. Student's t-distribution

### Definition 16.34 Student' t-distribution:

add

### 7.1. Delta Distribution

**Definition 16.35 The delta function  $\delta(\mathbf{x})$ :**

The delta/dirac function  $\delta(\mathbf{x})$  is defined by:

$$\int_{\mathbb{R}} \delta(\mathbf{x}) f(\mathbf{x}) d\mathbf{x} = f(0)$$

for any integrable function  $f$  on  $\mathbb{R}$ .

Or alternatively by:

$$\delta(x - x_0) = \lim_{\sigma \rightarrow 0} \mathcal{N}(x|x_0, \sigma) \quad (16.85)$$

$$\approx \infty \mathbb{1}_{\{x=x_0\}} \quad (16.86)$$

**Property 16.20 Properties of  $\delta$ :**

- Normalization:** The delta function integrates to 1:

$$\int_{\mathbb{R}} \delta(x) dx = \int_{\mathbb{R}} \delta(x) \cdot c_1(x) dx = c_1(0) = 1$$

where  $c_1(x) = 1$  is the constant function of value 1.

- Shifting:**

$$\int_{\mathbb{R}} \delta(x - x_0) f(x) dx = f(x_0) \quad (16.87)$$

- Symmetry:**

$$\int_{\mathbb{R}} \delta(-x) f(x) dx = f(0)$$

- Scaling:**

$$\int_{\mathbb{R}} \delta(\alpha x) f(x) dx = \frac{1}{|\alpha|} f(0)$$

#### Note

- In mathematical terms  $\delta$  is not a function but a **generalized function**.
- We may regard  $\delta(x - x_0)$  as a density with all its probability mass centered at the single point  $x_0$ .
- Using a box/indicator function s.t. its surface is one and its width goes to zero, instead of a normal distribution eq. (16.85) would be a non-differentiable/discrete form of the dirac measure.

#### Proofs

*Proof.* proposition 16.1: Let  $X$  be normally distributed with  $X \sim \mathcal{N}(\mu, \sigma^2)$ :

$$\begin{aligned} F_Y(y) &\stackrel{y>0}{=} \mathbb{P}_Y(Y \leq y) = \mathbb{P}(a + bX \leq y) = \mathbb{P}_X\left(X \leq \frac{y-a}{b}\right) \\ &= F_X\left(\frac{y-a}{b}\right) \\ F_Y(y) &\stackrel{y \leq 0}{=} \mathbb{P}_Y(Y \leq y) = \mathbb{P}(a + bX \leq y) = \mathbb{P}_X\left(X \geq \frac{y-a}{b}\right) \\ &= 1 - F_X\left(\frac{y-a}{b}\right) \end{aligned}$$

Differentiating both expressions w.r.t.  $y$  leads to:

$$f_Y(y) = \frac{dF_Y(y)}{dy} = \begin{cases} \frac{1}{b} \frac{dF_X\left(\frac{y-a}{b}\right)}{dy} \\ \frac{1}{-b} \frac{dF_X\left(\frac{y-a}{b}\right)}{dy} \end{cases} = \frac{1}{|b|} f_X(x) \left(\frac{y-a}{b}\right)$$

eq. (16.75)).

in order to prove that  $Y \sim \mathcal{N}(a + b\mu, b^2\sigma^2)$  we simply plug  $f_X$  in the previous expression:

$$\begin{aligned} f_Y(y) &= \frac{1}{\sqrt{2\pi}\sigma|b|} \exp\left\{-\frac{1}{2}\left(\frac{y-a}{b} - \mu\right)^2\right\} \\ &= \frac{1}{\sqrt{2\pi}\sigma|b|} \exp\left\{-\frac{1}{2}\left(\frac{y - (a + b\mu)}{\sigma|b|}\right)^2\right\} \end{aligned} \quad \square$$

*Proof.* proposition 16.2: Let  $X$  be normally distributed with  $X \sim \mathcal{N}(\mu, \sigma^2)$ :

$$\begin{aligned} Z &:= \frac{X - \mu}{\sigma} = \frac{1}{std} X - \frac{\mu}{\sigma} = aX + b \quad \text{with } a = \frac{1}{\sigma}, b = -\frac{\mu}{\sigma} \\ \text{eq. (16.75)} \quad \mathcal{N}(a\mu + b, a^2\sigma^2) &\sim \mathcal{N}\left(\frac{\mu}{\sigma} - \frac{\mu}{\sigma}, \frac{\sigma^2}{\sigma^2}\right) \sim \mathcal{N}(0, 1) \end{aligned} \quad \square$$

*Proof.* proposition 16.3: Let  $X$  be normally distributed with  $X \sim \mathcal{N}(\mu, \sigma^2)$ :

$$\begin{aligned} F_X(x) &= \mathbb{P}(X \leq x) \stackrel{-\mu}{\div \sigma} \mathbb{P}\left(\frac{X - \mu}{\sigma} \leq \frac{x - \mu}{\sigma}\right) \mathbb{P}\left(Z \leq \frac{x - \mu}{\sigma}\right) \\ &= \Phi\left(\frac{x - \mu}{\sigma}\right) \end{aligned} \quad \square$$

*Proof.* Property 16.18 scalar case

Let  $y \sim p(y) = \mathcal{N}(\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y-\mu)^2}{2\sigma^2}\right)$  and define  $\mathbf{x} = \mathbf{a}\mathbf{y} + b$   $\mathbf{a} \in \mathbb{R}_+, b \in \mathbb{R}$

Using the Change of variables formula it follows:

$$\begin{aligned} p_x(\bar{x}) &\stackrel{??}{=} \frac{p_y(\bar{y})}{\left|\frac{d\mathbf{x}}{d\mathbf{y}}\right|} \stackrel{\bar{y} = \frac{\bar{x}-b}{a}}{=} \frac{1}{a} \frac{1}{\sqrt{2\pi\mu^2}} \exp\left(-\frac{1}{2\sigma^2}\left(\frac{\bar{x}-b}{a} - \mu\right)^2\right) \\ &= \frac{1}{\sqrt{2\pi a^2 \mu^2}} \exp\left(-\frac{1}{2\sigma^2 a^2}(\bar{x} - \underbrace{b - a\mu}_{\mu_x})^2\right) \end{aligned}$$

Hence

$$x \sim \mathcal{N}(\mu_x, \sigma_x^2) = \mathcal{N}(a\mu + b, a^2\sigma^2) \quad \square$$

#### Note

We can also verify that we have calculated the right mean and variance by:

$$\begin{aligned} \mathbb{E}[x] &= \mathbb{E}[a\mathbf{y} + b] = a\mathbb{E}[y] + b = a\mu + b \\ \mathbb{V}[x] &= \mathbb{V}[a\mathbf{y} + b] = a^2\mathbb{V}[y] = a^2\sigma^2 \end{aligned}$$

*Proof.* Property 16.19

From Property 16.18 it follows immediately that  $\mathbf{z}$  is GRV  $\mathbf{z} \sim \mathcal{N}(\mu_z, \Sigma_z)$  with:

$$\mathbf{z} = \mathbf{A}\boldsymbol{\xi} \quad \text{with} \quad \mathbf{A} = \begin{bmatrix} \mathbf{A}_x & \mathbf{A}_y \end{bmatrix} \text{ and } \boldsymbol{\xi} = \begin{pmatrix} \mathbf{x} & \mathbf{y} \end{pmatrix}$$

Knowing that  $\mathbf{z}$  is a GRV it is sufficient to calculate  $\mu_z$  and  $\Sigma_z$  in order to characterize its distribution:

$$\begin{aligned} \mathbb{E}[\mathbf{z}] &= \mathbb{E}[\mathbf{A}_x \mathbf{x} + \mathbf{A}_y \mathbf{y}] = \mathbf{A}_x \mu_x + \mathbf{A}_y \mu_y \\ \mathbb{V}[\mathbf{z}] &= \mathbb{V}[\mathbf{A}\boldsymbol{\xi}] \stackrel{??}{=} \mathbf{A} \mathbb{V}[\boldsymbol{\xi}] \mathbf{A}^\top \\ &= \begin{bmatrix} \mathbf{A}_x & \mathbf{A}_y \end{bmatrix} \begin{bmatrix} \mathbb{V}[\mathbf{x}] & \text{Cov}[\mathbf{x}, \mathbf{y}] \\ \text{Cov}[\mathbf{y}, \mathbf{x}] & \mathbb{V}[\mathbf{y}] \end{bmatrix} \begin{bmatrix} \mathbf{A}_x & \mathbf{A}_y \end{bmatrix}^\top \\ &= \begin{bmatrix} \mathbf{A}_x & \mathbf{A}_y \end{bmatrix} \begin{bmatrix} \mathbb{V}[\mathbf{x}] & \text{Cov}[\mathbf{x}, \mathbf{y}] \\ \text{Cov}[\mathbf{y}, \mathbf{x}] & \mathbb{V}[\mathbf{y}] \end{bmatrix} \begin{bmatrix} \mathbf{A}_x^\top \\ \mathbf{A}_y^\top \end{bmatrix} \\ &= \mathbf{A}_x \mathbb{V}[\mathbf{x}] \mathbf{A}_x^\top + \mathbf{A}_y \mathbb{V}[\mathbf{y}] \mathbf{A}_y^\top \\ &\quad + \underbrace{\mathbf{A}_y \text{Cov}[\mathbf{y}, \mathbf{x}] \mathbf{A}_x^\top}_{=0 \text{ by independence}} + \underbrace{\mathbf{A}_x \text{Cov}[\mathbf{x}, \mathbf{y}] \mathbf{A}_y^\top}_{=0 \text{ by independence}} \\ &= \mathbf{A}_x \Sigma_x \mathbf{A}_x^\top + \mathbf{A}_y \Sigma_y \mathbf{A}_y^\top \end{aligned} \quad \square$$

#### Note

Can also be proved by using the normal definition of <sup>(def. 16.15)</sup> and tedious computations.



## 8. Sampling Random Numbers

Most math libraries have uniform **random number generator (RNG)** i.e. functions to generate uniformly distributed random numbers  $U \sim \mathcal{U}[a, b]$  (eq. (16.56)). Furthermore repeated calls to these RNG are independent, that is:

$$\begin{aligned} \mathbb{P}_{U_1, U_2}(u_1, u_2) &\stackrel{??}{=} \mathbb{P}_{U_1}(u_1) \cdot \mathbb{P}_{U_2}(u_2) \\ &= \begin{cases} 1 & \text{if } u_1, u_2 \in [a, b] \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

**Question:** using samples  $\{u_1, \dots, u_n\}$  of these CRVs with uniform distribution, how can we create random numbers with arbitrary discrete or continuous PDFs?

## 9. Inverse-transform Technique

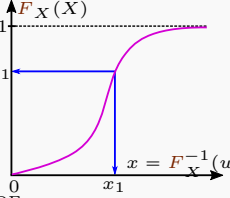
### Idea

Can make use of section 1 and the fact that CDF are increasing functions ([def. 6.10]). **Advantage:**

- Simple to implement
- All discrete distributions can be generated via inverse-transform technique

**Drawback:**

- Not all continuous distributions can be integrated/have closed form solution for their CDF. E.g. Normal-, Gamma-, Beta-distribution.



### 9.1. Continuous Case

**Definition 16.36 One Continuous Variable:** Given: a desired continuous pdf  $f_X$  and uniformly distributed rn  $\{u_1, u_2, \dots\}$ :

1. Integrate the desired pdf  $f_X$  in order to obtain the desired cdf  $F_X$ :

$$F_X(x) = \int_{-\infty}^x f_X(t) dt \quad (16.88)$$

2. Set  $F_X(X) \stackrel{!}{=} U$  on the range of  $X$  with  $U \sim \mathcal{U}[0, 1]$ .
3. Invert this equation/find the inverse  $F_X^{-1}(U)$  i.e. solve:

$$U = F_X(X) = F_X\left(\underbrace{F_X^{-1}(U)}_X\right) \quad (16.89)$$

4. Plug in the uniformly distributed rn:

$$x_i = F_X^{-1}(u_i) \quad \text{s.t.} \quad x_i \sim f_X \quad (16.90)$$

**Definition 16.37 Multiple Continuous Variable:**

**Given:** a pdf of multiple rvs  $f_{X,Y}$ :

1. Use the product rule (??) in order to decompose  $f_{X,Y}$ :

$$f_{X,Y} = f_{X,Y}(x, y) = f_{X|Y}(x|y)f_Y(y) \quad (16.91)$$

2. Use [def. 16.38] to first get a rv for  $y$  of  $Y \sim f_Y(y)$ .
3. Then with this fixed  $y$  use [def. 16.38] again to get a value for  $x$  of  $X \sim f_{X|Y}(x|y)$ .

*Proof.* [def. 16.38]:

**Claim:** if  $U$  is a uniform rv on  $[0, 1]$  then  $F_X^{-1}(U)$  has  $F_X$  as its CDF.

**Assume** that  $F_X$  is strictly increasing ([def. 6.10]). Then for any  $u \in [0, 1]$  there must exist a **unique**  $x$  s.t.  $F_X(x) = u$ .

Thus  $F_X$  must be invertible and we may write  $x = \underline{F_X^{-1}(u)}$ .

**Now** let  $a$  arbitrary:

$$F_X(a) = \mathbb{P}(x \leq a) = \mathbb{P}(F_X^{-1}(U) \leq a)$$

Since  $F_X$  is strictly increasing:

$$\begin{aligned} \mathbb{P}(F_X^{-1}(U) \leq a) &= \mathbb{P}(U \leq F_X(a)) \\ &\stackrel{\text{eq. (16.56)}}{=} \int_0^{F_X(a)} 1 dt = F_X(a) \end{aligned}$$

### Note

Strictly speaking we may not assume that a CDF is **strictly** increasing but we as all CDFs are weakly increasing ([def. 6.10]) we may always define an auxiliary function by its infimum:

$$\hat{F}_X^{-1} := \inf \{x | F_X(x) \geq 0\} \quad u \in [0, 1] \quad (16.92)$$

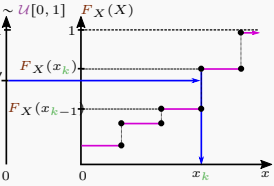
### 9.2. Discret Case

#### Idea

**Given:** a desired  $U \sim \mathcal{U}[0, 1]$  and uniformly distributed rn  $\{u_1, u_2, \dots\}$ .  
**Goal:** given a uniformly distributed rn  $u$  determine  $k$  s.t.:

$$\begin{aligned} k-1 &\leq U \leq k \\ \iff &F_X(x_{k-1}) < u \leq F_X(x_k) \end{aligned} \quad (16.93)$$

and return  $x_k$ .



**Definition 16.38 One Discret Variable:**

1. Compute the CDF of  $p_X$  ([def. 16.8])

$$F_X(x) = \sum_{t=-\infty}^x p_X(t) \quad (16.94)$$

2. Given the uniformly distributed rn  $\{u_i\}_{i=1}^n$  find  $k^i$  ( $\hat{=}$  inversion) s.t.:

$$F_X(x_{k(i)-1}) < u_i \leq F_X(x_{k(i)}) \quad \forall u_i \quad (16.95)$$

*Proof.* ??: First of all notice that we can always solve for an unique  $x_{k_i}$ .

**Ask:** why are Discret CRV always strictly increasing/unique?

**Given** a fixed  $x_{k_i}$  determine the values of  $u$  for which:

$$F_X(x_{k(i)-1}) < u \leq F_X(x_{k(i)}) \quad (16.96)$$

**Now** observe that:

$$\begin{aligned} u &\leq F_X(x_k) = F_X(x_{k-1}) + p_X(x_k) \\ \Rightarrow F_X(x_{k-1}) &< u \leq F_X(x_{k-1}) + p_X(x_k) \end{aligned}$$

The probability of  $U$  being in  $(F_X(x_{k-1}), F_X(x_k)]$  is:

$$\begin{aligned} \mathbb{P}(U \in [F_X(x_{k-1}), F_X(x_k)]) &= \int_{F_X(x_{k-1})}^{F_X(x_k)} p_U(t) dt \\ &= \int_{F_X(x_{k-1})}^{F_X(x_k)} 1 dt = F_X(x_k) - F_X(x_{k-1}) = p_X(x_k) \end{aligned}$$

Hence the random variable  $x_{k_i} \in \mathcal{X}$  has the pdf  $p_X$ .  $\square$

**Definition 16.39**

**Multiple Continuous Variables (Option 1):**

**Given:** a pdf of multiple rvs  $p_{X,Y}$ :

1. Use the product rule (??) in order to decompose  $p_{X,Y}$ :

$$p_{X,Y} = p_{X,Y}(x, y) = p_{X|Y}(x|y)p_Y(y) \quad (16.97)$$

2. Use ?? to first get a rv for  $y$  of  $Y \sim p_Y(y)$ .
3. Then with this fixed  $y$  use ?? again to get a value for  $x$  of  $X \sim p_{X|Y}(x|y)$ .

**Definition 16.40**

**Multiple Continuous Variables (Option 2):**

**Note:** this only works if  $\mathcal{X}$  and  $\mathcal{Y}$  are finite.

**Given:** a pdf of multiple rvs  $p_{X,Y}$  let  $N_x = |\mathcal{X}|$  and  $N_y = |\mathcal{Y}|$  the number of elements in  $\mathcal{X}$  and  $\mathcal{Y}$ .

**Define**  $p_Z(1) = p_{X,Y}(1, 1)$ ,  $p_Z(2) = p_{X,Y}(1, 2)$ ,  $\dots$ ,  $p_Z(N_x \cdot N_y) = p_{X,Y}(N_x, N_y)$

Then simply apply ?? to the auxillary pdf  $p_Z$

1. Use the product rule (??) in order to decompose  $f_{X,Y}$ :

$$f_{X,Y} = f_{X,Y}(x, y) = f_{X|Y}(x|y)f_Y(y) \quad (16.98)$$

2. Use [def. 16.38] to first get a rv for  $y$  of  $Y \sim f_Y(y)$ .
3. Then with this fixed  $y$  use [def. 16.38] again to get a value for  $x$  of  $X \sim f_{X|Y}(x|y)$ .

See examples see comment in code text

10. Descriptive Statistics

10.1. Population Parameters

**Definition 16.41 Population/Statistical Parameter:** Are parameters defining families of probability distributions and thus characteristics of population following such distributions i.e. the normal distribution has two parameters  $\{\mu, \sigma^2\}$

**Definition 16.42 Population Mean:** Given a population  $\{x_i\}_{i=1}^N$  of size  $N$  its variance is defined as:

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i \tag{16.99}$$

**Definition 16.43 Population Variance:** Given a population  $\{x_i\}_{i=1}^N$  of size  $N$  its variance is defined as:  $\{x_i\}_{i=1}^N$

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2 \tag{16.100}$$

**Note**  
The population variance and mean are equally to the mean derived from the true distribution of the population.

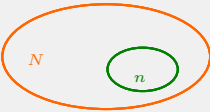
10.2. Sample Estimates

**Definition 16.44 (Sample) Statistic:** A statistic is a measurable function  $f$  that assigns a **single** value  $F$  to a sample of random variables or population:

$$f : \mathbb{R}^n \mapsto \mathbb{R} \qquad F = f(X_1, \dots, X_n)$$

E.g.  $F$  could be the mean, variance,...

**Note**  
The function itself is independent of the sample's distribution; that is, the function can be stated before realization of the data.



**Definition 16.45 (Point) Estimator**  $\hat{\theta} = \hat{\theta}(\mathbf{X})$ :  
**Given:** n-samples  $\mathbf{x}_1, \dots, \mathbf{x}_n \sim \mathbf{X}$  an estimator  
$$\hat{\theta} = h(\mathbf{x}_1, \dots, \mathbf{x}_n) \tag{16.101}$$

is a statistic/random variable used to estimate a true (population) parameter  $\theta$ <sup>[def. 16.41]</sup>.

**Note**  
The other kind of estimators are interval estimators which do not calculate a statistic **but** an interval of plausible values of an unknown population parameter  $\theta$ .  
The most prevalent forms of interval estimation are:

- Confidence intervals (frequentist method).
- Credible intervals (Bayesian method).

**Definition 16.46 Degrees of freedom of a Statistic:** Is the number of values in the final calculation of a statistic that are free to vary.

10.2.1. Empirical Mean

**Definition 16.47 Sample/Empirical Mean**  $\bar{x}$ :  
The sample mean is an estimate/statistic of the population mean<sup>[def. 16.42]</sup> and can be calculated from an observation/sample of the total population  $\{x_i\}_{i=1}^n \subset \{x_i\}_{i=1}^N$ :

$$\bar{x} = \hat{\mu}_X = \frac{1}{n} \sum_{i=1}^n x_i \tag{16.102}$$

**Corollary 16.5 Expectation:** The sample mean estimator is unbiased (see section 14):

$$\mathbb{E}[\hat{\mu}_X] = \mu \tag{16.103}$$

**Corollary 16.6 Variance:** For the variance of the sample mean estimator it holds (see section 14):

$$\mathbb{V}[\hat{\mu}_X] = \frac{1}{n} \sigma_X^2 \tag{16.104}$$

10.2.2. Empirical Variance

**Definition 16.48 Biased Sample Variance:** The sample mean is an estimate/statistic of the population variance<sup>[def. 16.43]</sup> and can be calculated from an observation/sample of the total population  $\{x_i\}_{i=1}^n \subset \{x_i\}_{i=1}^N$ :

$$s_n^2 = \hat{\sigma}_X^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 \tag{16.105}$$

**Definition 16.49 (Unbiased) Sample Variance:**

$$s^2 = \hat{\sigma}_X^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \mu)^2 \tag{16.106}$$

see section 14

**Definition 16.50 Bessel's Correction:** The factor  $\frac{n}{n-1}$  (16.107)

as multiplying the uncorrected population varianceeq. (16.105) by this term yields an unbiased estimated of the variance (not the standard deviation). The reason for this is that are

**Attention:** Usually only unbiased variance is used and also sometimes denoted by  $s_n^2$

*Proof.*  
*finish this*

11. Statistical Tests

**Definition 16.51 Null Hypothesis:** A Null Hypothesis  $H_0$  is usually a commonly accepted fact/view/base hypothesis that researchers try to nullify or disprove.

$$H_0 : \theta = \theta_0 \tag{16.108}$$

**Definition 16.52 Alternative Hypothesis:** The Alternative Hypothesis  $H_A/H_1$  is the opposite of the Null Hypotheses/contradicts it and is what we try to test against the Null Hypothesis.

$$H_A : \theta \begin{cases} > \theta_0 & \text{(one-sided)} \\ < \theta_0 & \text{(one-sided)} \\ \neq \theta_0 & \text{(two-sided)} \end{cases} \tag{16.109}$$

**Definition 16.53 Testing Parameters:**  
**Given:** a parameter  $\theta$  that we want to test.  
Let  $\Theta$  be the set of all possible values that  $\theta$  can achieve. We now split  $\Theta$  in two disjunct sets  $\Theta_0$  and  $\Theta_1$ .  
$$\Theta = \Theta_0 \cup \Theta_1 \qquad \Theta_0 \cap \Theta_1 = \emptyset$$

Null Hypothesis  $H_0 : \theta \in \Theta_0$  (16.110)  
Alternative Hypothesis  $H_A : \theta \in \Theta_1$  (16.111)

11.1. Type I&II Errors

**Definition 16.54 Type I Error:** Is the rejection of a Null Hypothesis, even-tough its true (also known as a "false positive").

**Definition 16.55 Type II Error:** Is the acceptance of a Null Hypothesis, even-tough its false (also known as a "false negative").

Decision	$H_0$ true	$H_0$ false	
Accept	TN	Type II (FN)	
Reject	Type I (FP)	TP	

**Definition 16.56 Critical Value c:** Value from which on the Null-hypothesis  $H_0$  gets rejected.

**Definition 16.57 Statistical significance**  $\alpha$ : A study's defined significance level, denoted  $\alpha$ , is the **probability** of the study rejecting the null hypothesis, given that the null hypothesis were true (Type I Error).

**Definition 16.58 Critical Region**  $K_\alpha$ : Is the set of all values that causes us to reject the Null Hypothesis in favor for the Alternative Hypothesis  $H_A$ . The Critical region is usually chosen s.t. we incur a Type I Error with probability less than  $\alpha$ .

$$K_\alpha \in \Theta : \mathbb{P}(\text{Type I Error}) \leq \alpha \tag{16.112}$$

or 
$$\begin{aligned} &\mathbb{P}(c_2 \leq X \leq c_1) \leq \alpha && \text{two-sided} \\ &\mathbb{P}(c_2 \leq X) \leq \frac{\alpha}{2} \text{ and } \mathbb{P}(X \leq c_1) \leq \frac{\alpha}{2} \\ &\mathbb{P}(c_2 \leq X) \leq \alpha && \text{one-sided} \\ &\mathbb{P}(X \leq c_1) \leq \alpha && \text{one-sided} \end{aligned}$$

**Definition 16.59 Acceptance Region:** Is the region where we accept the null hypothesis  $H_0$ .

**Note**  
see example 16.3.

11.2. Normally Distributed Data

Let us consider a sample of  $\{x_i\}_{i=1}^n$  i.i.d. observations, that follow a normal distribution  $x_i \sim \mathcal{N}(\mu, \sigma^2)$ .

11.2.1. Z-Test  $\sigma$  known  
11.2.2. t-Test  $\sigma$  unknown

12. Inferential Statistics

**Goal of Inference**

- ① What is a good guess of the parameters of my model?
- ② How do I quantify my uncertainty in the guess?

13. Examples

**Example 16.1 ??:** Let  $x$  be uniformly distributed on  $[0, 1]$  (def. 16.26) with pmf  $p_X(x)$  then it follows:  
 $\frac{dy}{dx} = \frac{1}{p_Y(y)} \Rightarrow dx = dy p_Y(y) \Rightarrow x = \int_{-\infty}^y p_Y(t) dt = F_Y(x)$

**Example 16.2 ??:** Let

add <https://www.youtube.com/watch?v=WUUbTVIRagg>

**Example 16.3 Binomialtest:**  
**Given:** a manufacturer claims that a maximum of 10% of its delivered components are substandard goods.  
In a sample of size  $n = 20$  we find  $x = 5$  goods that do not fulfill the standard and are skeptical that the what the manufacture claims is true, so we want to test:  
 $H_0 : p = p_0 = 0.1$  vs.  $H_A : p > 0.1$

We model the number of defective goods using the binomial distribution (def. 16.23)

$$X \sim \mathcal{B}(n, p), n = 20 \quad \mathbb{P}(X \geq x) = \sum_{k=x}^n \binom{n}{k} p^k (1-p)^{n-k}$$

from this we find:  
 $\mathbb{P}_{p_0}(X \geq 4) = 1 - \mathbb{P}_{p_0}(X \leq 3) = 0.13$   
 $\mathbb{P}_{p_0}(X \geq 4) = 1 - \mathbb{P}_{p_0}(X \leq 3) = 0.04 \leq \alpha$   
thus the probability that equal 5 or more then 5 parts out of the 20 are rejects is less then 4%.  
 $\Rightarrow$  throw away null hypothesis for the 5% niveau in favor to the alternative.  
 $\Rightarrow$  the 5% significance niveau is given by  $K = \{5, 6, \dots, 20\}$

Note

If  $x < n/2$  it is faster to calculate  $\mathbb{P}(X \geq x) = 1 - \mathbb{P}(X \leq x-1)$

14. Proofs

Proof. corollary 16.5:

$$\mathbb{E}[\hat{\mu}_X] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \frac{1}{n} \mathbb{E}\left[\sum_{i=1}^n x_i\right] = \frac{1}{n} \mathbb{E}\left[\underbrace{\mu + \dots + \mu}_{1, \dots, n}\right]$$

□

Proof. corollary 16.6:

$$\mathbb{V}[\hat{\mu}_X] = \mathbb{V}\left[\frac{1}{n} \sum_{i=1}^n x_i\right] \stackrel{\text{Property 16.10}}{=} \frac{1}{n^2} \mathbb{V}\left[\sum_{i=1}^n x_i\right]$$
$$\frac{1}{n^2} n \mathbb{V}[X] = \frac{1}{n} \sigma^2$$

□

Proof. definition 16.49:

$$\begin{aligned} \mathbb{E}[\hat{\sigma}_X^2] &= \mathbb{E}\left[\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2\right] \\ &= \frac{1}{n-1} \mathbb{E}\left[\sum_{i=1}^n (x_i^2 - 2x_i \bar{x} + \bar{x}^2)\right] \\ &= \frac{1}{n-1} \mathbb{E}\left[\sum_{i=1}^n x_i^2 - 2\bar{x} \sum_{i=1}^n x_i + \sum_{i=1}^n \bar{x}^2\right] \\ &= \frac{1}{n-1} \mathbb{E}\left[\sum_{i=1}^n x_i^2 - 2n\bar{x} \cdot n\bar{x} + n\bar{x}^2\right] \\ &= \frac{1}{n-1} \mathbb{E}\left[\sum_{i=1}^n x_i^2 - n\bar{x}^2\right] \\ &= \frac{1}{n-1} \left[\sum_{i=1}^n \mathbb{E}[x_i^2] - n\mathbb{E}[\bar{x}^2]\right] \\ &= \frac{1}{n-1} \left[\sum_{i=1}^n (\sigma^2 + \mu^2) - n\mathbb{E}[\bar{x}^2]\right] \\ &= \frac{1}{n-1} \left[\sum_{i=1}^n (\sigma^2 + \mu^2) - n\left(\frac{1}{n}\sigma^2 + \mu^2\right)\right] \\ &= \frac{1}{n-1} \left[(n\sigma^2 + n\mu^2) - (\sigma^2 + n\mu^2)\right] \\ &= \frac{1}{n-1} \left[n\sigma^2 - \sigma^2\right] = \frac{1}{n-1} \left[(n-1)\sigma^2\right] = \sigma^2 \end{aligned}$$

□



Stochastic Calculus

Stochastic Processes

<b>Definition 17.1 Random/Stochastic Process</b> $\{X_t, t \in \mathcal{T} \subseteq \mathbb{R}_+\}$ : is a collection of random variables on the same probability space $(\Omega, \mathcal{A}, \mathbb{P})$ . The index set $\mathcal{T}$ is usually representing time and can be either an interval $[t_1, t_2]$ or a discrete set $\{t_1, t_2, \dots\}$ . Therefore, the random process $X$ can be written as a function: $X: \mathbb{R} \times \Omega \mapsto \mathbb{R} \iff (t, \omega) \mapsto X(t, \omega) \quad (17.1)$
<b>Definition 17.2 Sample path/Trajectory/Realization:</b> Is the <i>stochastic/noise signal</i> $r(\cdot, \omega)$ on the index set $\mathcal{T}$ , that we obtain be sampling $\omega$ from $\Omega$ .
<b>Notation</b> Even though the r.v. $X$ is a function of two variables, most books omit the argument of the sample space $X(t, \omega) := X(t)$
<b>Definition 17.3 Filtration</b> A collection $\{\mathcal{F}_t\}_{t \geq 0}$ of sub $\sigma$ -algebras $\{\mathcal{F}_t\}_{t \geq 0} \subseteq \mathcal{F}$ is called filtration if is <i>increasing</i> : $\mathcal{F}_s \subseteq \mathcal{F}_t \quad \forall s \leq t \quad (17.2)$
<b>Definition 17.4 Adapted Process:</b> A stochastic process $\{X_t: 0 \leq t \leq \infty\}$ is called adapted to a filtration $\mathbb{F}$ if, $X_t$ is $\mathcal{F}_{t-}$ -measurable, i.e. observable at time $t$ .
<b>Definition 17.5 Predictable Process:</b> A stochastic process $\{X_t: 0 \leq t \leq \infty\}$ is called predictable w.r.t. a filtration $\mathbb{F}$ if, $X_t$ is $\{\mathcal{F}_{t-}\}$ -measurable, i.e. the value of $X_t$ is known at time $t - 1$ .
<b>Note</b> The price of a stock will usually be adapted since date $k$ prices are known at date $k$ . On the other hand the interest rate of a bank account is usually already known at the beginning $k - 1$ , s.t. the interest rate $r_t$ ought to be $\mathcal{F}_{k-1}$ measurable, i.e. the process $r = (r_k)_{k=1, \dots, T}$ should be predictable.
<b>Definition 17.6 Filtered Probability Space</b> $(\Omega, \mathcal{F}, \{\mathcal{F}_t\}_{t \geq 0}, \mathbb{P})$ : A probability space $(\Omega, \mathcal{F}, \mathbb{P})$ together with a filtration $\{\mathcal{F}_t\}_{t \geq 0}$ is called a <i>filtered probability space</i> .
<b>Corollary 17.1</b> : The amount of information of an adapted random process is increasing see example 17.1.
<b>Definition 17.7 Martingales:</b> A stochastic process $X(t)$ is a martingale on a <i>filtered probability space</i> $(\Omega, \mathcal{F}, \{\mathcal{F}_t\}_{t \geq 0}, \mathbb{P})$ if the following conditions hold: <ol style="list-style-type: none"><li>Given <math>s \leq t</math> the best prediction of <math>X(t)</math>, with a filtration <math>\{\mathcal{F}_s\}</math> is the current expected value: <math display="block">\forall s \leq t \quad \mathbb{E}[X(t) \mathcal{F}_s] = X(s) \quad \text{a.s.} \quad (17.3)</math></li><li>The expectation is finite: <math display="block">\mathbb{E}[ X(t) ] &lt; \infty \quad \forall t \geq 0 \quad X(t) \text{ is } \{\mathcal{F}_t\}_{t \geq 0} \text{ adapted} \quad (17.4)</math></li></ol>
<b>Interpretation</b> <ul style="list-style-type: none"><li>For any <math>\mathcal{F}_s</math>-adapted process the best prediction of <math>X(t)</math> is the currently known value <math>X(s)</math> i.e. if <math>\mathcal{F}_s = \mathcal{F}_{t-1}</math> then the best prediction is <math>X(t - 1)</math></li><li>A martingale models fair games of limited information.</li></ul>
<b>Definition 17.8 Auto Covariance</b> $\gamma(t_2 - t_1)$ : Describes the covariance <sup>[def. 16.16]</sup> between two values of a stochastic process $(X_t)_{t \in \mathcal{T}}$ at different time points $t_1$ and $t_2$ . $\gamma(t_1, t_2) = \text{Cov}[X_{t_1}, X_{t_2}] = \mathbb{E}[(X_{t_1} - \mu_{t_1})(X_{t_2} - \mu_{t_2})] \quad (17.5)$ For zero time differences $t_1 = t_2$ the autocorrelation functions equals the variance: $\gamma(t, t) = \text{Cov}[X_t, X_t] \stackrel{\text{eq. (16.41)}}{=} \mathbb{V}[X_t] \quad (17.6)$

<b>Notes</b> <ul style="list-style-type: none"><li>Hence the autocorrelation describes the correlation of a function or signal with itself at a previous time point.</li><li>Given a random time dependent variable <math>\mathbf{x}(t)</math> the autocorrelation function <math>\gamma(t, t - \tau)</math> describes how <i>similar</i> the time translated function <math>\mathbf{x}(t - \tau)</math> and the original function <math>\mathbf{x}(t)</math> are.</li><li>If there exists some relation between the values of the time series that is non-random, then the autocorrelation is non-zero.</li><li>The auto covariance is maximized/most similar for no translation <math>\tau = 0</math> at all.</li></ul>
<b>Definition 17.9 Auto Correlation</b> $\rho(t_2 - t_1)$ : Is the scaled version of the auto-covariance <sup>[def. 17.8]</sup> : $\rho(t_2 - t_1) = \frac{\text{Corr}[X_{t_1}, X_{t_2}]}{\sigma_{X_{t_1}} \sigma_{X_{t_2}}} = \frac{\mathbb{E}[(X_{t_1} - \mu_{t_1})(X_{t_2} - \mu_{t_2})]}{\sigma_{X_{t_1}} \sigma_{X_{t_2}}} \quad (17.7)$
<b>1. Different kinds of Processes</b> <b>1.1. Markov Process</b> <b>Definition 17.10 Markov Process:</b> A continuous-time stochastic process $X(t)$ , $t \in T$ , is called a Markov process if for any finite parameter set $\{t_i: t_i < t_{i+1}\} \in T$ it holds: $\mathbb{P}(X(t_{n+1}) \in B   X(t_1), \dots, X(t_n)) = \mathbb{P}(X(t_{n+1}) \in B   X(t_n))$ it thus follows for the <i>transition probability</i> – the probability of $X(t)$ lying in the set $B$ at time $t$ , given the value $x$ of the process at time $s$ : $\mathbb{P}(s, x, t, B) = P(X(t) \in B   X(s) = x) \quad 0 \leq s < t \quad (17.8)$
<b>Interpretation</b> In order to predict the future only the current/last value counts.
<b>Corollary 17.2 Transition Density:</b> The transition probability of a continuous distribution $\mathbf{p}$ can be calculated via: $\mathbb{P}(s, x, t, B) = \int_B \mathbf{p}(s, x, t, y) dy \quad (17.9)$
<b>1.2. Gaussian Process</b> <b>Definition 17.11 Gaussian Process:</b> Is a stochastic process $X(t)$ where the random variables follow a Gaussian distribution: $X(t) \sim \mathcal{N}(\mu(t), \sigma^2(t)) \quad \forall t \in T \quad (17.10)$
<b>1.3. Diffusions</b> <b>Definition 17.12 Diffusion:</b> Is a Markov Process <sup>[def. 17.10]</sup> for which it holds that: $\mu(t, X(t)) = \lim_{t \rightarrow 0} \frac{1}{\Delta t} \mathbb{E}[X(t + \Delta t) - X(t)   X(t)] \quad (17.11)$ $\sigma^2(t, X(t)) = \lim_{t \rightarrow 0} \frac{1}{\Delta t} \mathbb{E}[(X(t + \Delta t) - X(t))^2   X(t)] \quad (17.12)$ See ??/eq. (17.12) for simple proof of eq. (17.11)/??. <ul style="list-style-type: none"><li><math>\mu(t, X(t))</math> is called <b>drift</b></li><li><math>\sigma^2(t, X(t))</math> is called <b>diffusion coefficient</b></li></ul>
<b>Interpretation</b> There exist not discontinuities for the trajectories.

<b>1.4. Brownian Motion/Wiener Process</b> <b>Definition 17.13 d-dim standard Brownian Motion/Wiener Process:</b> Is an $\mathbb{R}^d$ valued <i>stochastic process</i> <sup>[def. 17.1]</sup> $(W_t)_{t \in \mathcal{T}}$ starting at $\mathbf{x}_0 \in \mathbb{R}^d$ that satisfies: <ol style="list-style-type: none"><li><b>Normal Independent Increments:</b> the increments are <i>normally distributed independent random variables</i>: <math display="block">W(t_i) - W(t_{i-1}) \sim \mathcal{N}(0, (t_i - t_{i-1}) \mathbb{1}_{d \times d}) \quad \forall i \in \{1, \dots, T\} \quad (17.13)</math></li><li><b>Stationary increments:</b> <math>W(t + \Delta t) - W(t)</math> is independent of <math>t \in \mathcal{T}</math></li><li><b>Continuity:</b> for a.e. <math>\omega \in \Omega</math>, the function <math>t \mapsto W_t(\omega)</math> is continuous <math display="block">\lim_{t \rightarrow 0} \frac{\mathbb{P}( W(t + \Delta t) - W(t)  \geq \delta)}{\Delta t} = 0 \quad \forall \delta &gt; 0 \quad (17.14)</math></li><li><b>Start</b> <math display="block">W(0) := W_0 = 0 \quad \text{a.s.} \quad (17.15)</math></li></ol>
<b>Notation</b> <ul style="list-style-type: none"><li>In many source the Brownian motion is a synonym for the standard Brownian Motion and it is the same as the Wiener process.</li><li>However in some sources the Wiener process is the standard Brownian Motion, while the Brownian motion denotes a general form <math>\alpha W(t) + \beta</math>.</li></ul>
<b>Corollary 17.3</b> $W_t \sim \mathcal{N}(0, \sigma)$ : The random variable $W_t$ follows the $\mathcal{N}(0, \sigma)$ law $\mathbb{E}[W(t)] = \mu = 0 \quad (17.16)$ $\mathbb{V}[W(t)] = \mathbb{E}[W^2(t)] = \sigma^2 = t \quad (17.17)$ See section 5
<b>1.4.1. Properties of the Wiener Process</b> <b>Property 17.1 Non-Differentiable Trajectories:</b> The sample paths of a Brownian motion are not differentiable: $\frac{dW(t)}{dt} = \lim_{t \rightarrow 0} \mathbb{E} \left[ \left( \frac{W(t + \Delta t) - W(t)}{\Delta t} \right)^2 \right]$ $= \lim_{t \rightarrow 0} \frac{\mathbb{E}[W(t + \Delta t) - W(t)]}{\Delta t} = \lim_{t \rightarrow 0} \frac{\sigma^2}{\Delta t} = \infty$ $\xrightarrow{\text{result}} \text{cannot use normal calculus anymore}$ $\xrightarrow{\text{solution}} \text{Ito Calculus see section 18.}$
<b>Property 17.2 Auto covariance Function:</b> The auto-covariance <sup>[def. 17.8]</sup> for a Wiener process $\mathbb{E}[(W(t) - \mu(t))(W(t') - \mu(t'))] = \min(t, t') \quad (17.18)$
<b>Property 17.3:</b> A standard Brownian motion is a <b>Quadratic Variation</b>
<b>Definition 17.14 Total Variation:</b> The total variation of a function $f: [a, b] \subset \mathbb{R} \mapsto \mathbb{R}$ is defined as: $LV_{[a,b]}(f) = \sup_{\Pi \in \mathcal{S}} \sum_{i=0}^{n_{\Pi}-1}  f(x_{i+1}) - f(x_i)  \quad (17.19)$ $\mathcal{S} = \{\Pi\{x_0, \dots, x_{n_{\Pi}}\} : \Pi \text{ is a partition } \stackrel{[\text{def. 12.8}]}{\text{of}} [a, b]\}$ it is a measure of the (one dimensional) length of a function w.r.t. to the y-axis, when moving alone the function. Hence it is a measure of the variation of a function w.r.t. to the y-axis.
<b>Definition 17.15 Total Quadratic Variation/“sum of squares”:</b> The total quadratic variation of a function $f: [a, b] \subset \mathbb{R} \mapsto \mathbb{R}$ is defined as: $QV_{[a,b]}(f) = \sup_{\Pi \in \mathcal{S}} \sum_{i=0}^{n_{\Pi}-1}  f(x_{i+1}) - f(x_i) ^2 \quad (17.20)$ $\mathcal{S} = \{\Pi\{x_0, \dots, x_{n_{\Pi}}\} : \Pi \text{ is a partition } \stackrel{[\text{def. 12.8}]}{\text{of}} [a, b]\}$

<b>Corollary 17.4 Bounded (quadratic) Variation:</b> The (quadratic) variation <sup>[def. 17.14]</sup> of a function is bounded if it is finite: $\exists M \in \mathbb{R}_+ : LV_{[a,b]}(f) \leq M \quad (QV_{[a,b]}(f) \leq M) \quad \forall \Pi \in \mathcal{S} \quad (17.21)$
<b>Theorem 17.1 Variation of Wiener Process:</b> Almost surely the total variation of a Brownian motion over a interval $[0, T]$ is infinite: $\mathbb{P}(\omega : LV(W(\omega)) < \infty) = 0 \quad (17.22)$
<b>Theorem 17.2 Quadratic Variation of standard Brownian Motion:</b> The quadratic variation of a standard Brownian motion over $[0, T]$ is finite: $\lim_{N \rightarrow \infty} \sum_{k=1}^N \left[ W\left(k \frac{T}{N}\right) - W\left((k-1) \frac{T}{N}\right) \right]^2 = T$ with probability 1 See ??
<b>Corollary 17.5</b> : theorem 17.2 can also be written as: $(dW(t))^2 = dt \quad (17.24)$
<b>1.4.2. Lévy’s Characterization of BM</b> <b>Theorem 17.3 d-dim standard BM/Wiener Process by Paul Lévy:</b> An $\mathbb{R}^d$ valued <i>adapted stochastic process</i> <sup>[def.x, 17.1, 17.3]</sup> $(W_t)_{t \in \mathcal{T}}$ with the filtration $\{\mathcal{F}_t\}_{t \in \mathbb{R}_+}$ , that satisfies: <ol style="list-style-type: none"><li><b>Start</b> <math display="block">W(0) := W_0 = 0 \quad \text{a.s.} \quad (17.25)</math></li><li><b>Continuous Martingale:</b> <math>W_t</math> is an a.s. <i>continuous martingale</i><sup>[def. 17.7]</sup> w.r.t. the filtration <math>(\mathcal{F}_t)_{t \in \mathcal{T}}</math> under <math>\mathbb{P}</math>.</li><li><b>Quadratic Variation:</b> <math display="block">W_t^2 - t \text{ is also an martingale} \iff QV(W_t) = t \quad (17.26)</math></li></ol>
is a standard Brownian motion <sup>[def. 17.18]</sup> . Proof see section 5
<b>Further Stochastic Processes</b> <b>1.4.3. White Noise</b> <div>understand script and add</div> <b>Definition 17.16 Discrete-time white noise:</b> Is a random signal $\{\epsilon_t\}_{t \in T_{\text{discret}}}$ having equal intensity at different frequencies and is defined by: <ul style="list-style-type: none"><li>Having zero tendencies/expectation (otherwise the signal would not be random): <math display="block">\mathbb{E}[\epsilon[k]] = 0 \quad \forall k \in T_{\text{discret}} \quad (17.27)</math></li><li>Zero autocorrelation<sup>[def. 17.9]</sup> <math>\gamma</math> i.e. the signals of different times are in no-way correlated: <math display="block">\gamma(\epsilon[k], \epsilon[k + n]) = \mathbb{E}[\epsilon[k]\epsilon[k + n]^T] = \mathbb{V}[\epsilon[k]] \delta_{\text{discret}}[n] \quad \forall k, n \in T_{\text{discret}} \quad (17.28)</math></li></ul> <b>With</b> $\delta_{\text{discret}}[n] := \begin{cases} 1 & \text{if } n = 0 \\ 0 & \text{else} \end{cases}$ See proofs
<b>Definition 17.17 Continuous-time white noise:</b> Is a random signal $(\epsilon_t)_{t \in T_{\text{continuous}}}$ having equal intensity at different frequencies and is defined by: <ul style="list-style-type: none"><li>Having zero tendencies/expectation (otherwise the signal would not be random): <math display="block">\mathbb{E}[\epsilon(t)] = 0 \quad \forall t \in T_{\text{continuous}} \quad (17.29)</math></li><li>Zero autocorrelation<sup>[def. 17.9]</sup> <math>\gamma</math> i.e. the signals of different times are in no-way correlated: <math display="block">\gamma(\epsilon(t), \epsilon(t + \tau)) = \mathbb{E}[\epsilon(t)\epsilon(t + \tau)^T] \stackrel{\text{eq. (16.86)}}{=} \mathbb{V}[\epsilon(t)] \delta(t - \tau) = \begin{cases} \mathbb{V}[\epsilon(t)] &amp; \text{if } \tau = 0 \\ 0 &amp; \text{else} \end{cases} \quad \forall t, \tau \in T_{\text{continuous}} \quad (17.31)</math></li></ul>

#### 1.4.4. Generalized Brownian Motion

**Definition 17.18 Brownian Motion:**  
Let  $\{W_t\}_{t \in \mathbb{R}_+}$  be a standard Brownian motion<sup>[def. 17.13]</sup>, and define:

$$X_t = \mu t + \sigma W_t \quad t \in \mathbb{R}_+ \quad \begin{array}{l} \mu \in \mathbb{R} : \text{drift parameter} \\ \sigma \in \mathbb{R}_+ : \text{scale parameter} \end{array} \quad (17.32)$$

then  $\{X_t\}_{t \in \mathbb{R}_+}$  is normally distributed with mean  $\mu t$  and variance  $t\sigma^2$ .  $X_t \sim \mathcal{N}(\mu t, \sigma^2 t)$ .

**Theorem 17.4 Normally Distributed Increments:**  
If  $W(T)$  is a Brownian motion, then  $W(t) - W(0)$  is a normal random variable with mean  $\mu t$  and variance  $\sigma^2 t$ , where  $\mu, \sigma \in \mathbb{R}$ . From this it follows that  $W(t)$  is distributed as:

$$f_{W(t)}(x) \sim \mathcal{N}(\mu t, \sigma^2 t) = \frac{1}{\sqrt{2\pi\sigma^2 t}} \exp\left\{-\frac{(x - \mu t)^2}{2\sigma^2 t}\right\} \quad (17.33)$$

**Corollary 17.6 :** More generally we may define the process:

$$t \mapsto f(t) + \sigma W_t \quad (17.34)$$

which corresponds to a noisy version of  $f$ .

**Corollary 17.7 Brownian Motion as a Solution of an SDE:** A stochastic process  $X_t$  follows a BM with drift  $\mu$  and scale  $\sigma$  if it satisfies the following SDE:

$$\begin{aligned} dX(t) &= \mu dt + \sigma dW(t) & (17.35) \\ X(0) &= 0 & (17.36) \end{aligned}$$

#### 1.4.5. Geometric Brownian Motion (GBM)

For many processes  $X(t)$  it holds that:

- there exists an (exponential) growth
- that the values may not be negative  $X(t) \in \mathbb{R}_+$

**Definition 17.19 Geometric Brownian Motion:**  
Let  $\{W_t\}_{t \in \mathbb{R}_+}$  be a standard Brownian motion<sup>[def. 17.13]</sup> the exponential transform:

$$X(t) = \exp(W(t)) = \exp(\mu t + \sigma W(t)) \quad t \in \mathbb{R}_+ \quad (17.37)$$

is called geometric Brownian motion

**Corollary 17.8 Log-normal Returns:** For a geometric BM we obtain log-normal returns:

$$\ln\left(\frac{S_t}{S_0}\right) = \mu t + \sigma W(t) \iff \mu t + \sigma W(t) \sim \mathcal{N}(\mu t, \sigma^2 t) \quad (17.38)$$

meaning that the mean and the variance of the process (stock) *log-returns* grow over time linearly.

**Corollary 17.9 Geometric BM as a Solution of an SDE:**  
A stochastic process  $X_t$  follows a geometric BM with drift  $\mu$  and scale  $\sigma$  if it satisfies the following SDE:

$$\begin{aligned} dX(t) &= X(t) (\mu dt + \sigma dW(t)) & (17.39) \\ &= \mu X(t) dt + \sigma X(t) dW(t) & (17.40) \\ X(0) &= 0 \end{aligned}$$

#### 1.4.6. Locally Brownian Motion

**Definition 17.20 Locally Brownian Motion:**  
Let  $\{W_t\}_{t \in \mathbb{R}_+}$  be a standard Brownian motion<sup>[def. 17.13]</sup> a local Brownian motion is a stochastic process  $X(t)$  that satisfies the SDE:

$$dX(t) = \mu(X(t), t) dt + \sigma(X(t), t) dW(t) \quad (17.41)$$

##### Note

A local Brownian motion is an generalization of a geometric Brownian motion.

#### 1.4.7. Ornstein-Uhlenbeck Process

**Definition 17.21 Ornstein-Uhlenbeck Process:**  
Let  $\{W_t\}_{t \in \mathbb{R}_+}$  be a standard Brownian motion<sup>[def. 17.13]</sup> a Ornstein-Uhlenbeck Process or exponentially correlated noise is a stochastic process  $X(t)$  that satisfies the SDE:

$$dX(t) = -aX(t) dt + b\sigma dW(t) \quad a > 0 \quad (17.42)$$

#### 1.5. Poisson Processes

**Definition 17.22 Rare/Extreme Events:** Are events that lead to discontinuous in stochastic processes.

##### Problem

A Brownian motion is not sufficient as model in order to describe extreme events s.a. crashes in financial market time series. Need a model that can describe such discontinuities/jumps.

**Definition 17.23 Poisson Process:** A Poisson Process with rate  $\lambda \in \mathbb{R}_{\geq 0}$  is a collection of random variables  $X(t)$ ,  $t \in [0, \infty)$  defined on a probability space  $(\Omega, \mathcal{F}, \{\mathcal{F}_t\}_{t \geq 0}, \mathbb{P})$ , having a discrete state space  $N = \{0, 1, 2, \dots\}$  and satisfies:

1.  $X_0 = 0$
2. The increments follow a Poisson distribution<sup>[def. 16.25]</sup>:

$$\mathbb{P}((X_t - X_s) = k) = \frac{\lambda(t-s)}{k!} e^{-\lambda(t-s)} \quad 0 \leq s < t < \infty \quad \forall k \in \mathbb{N}$$

3. No correlation of (non-overlapping) increments:  
 $\forall t_0 < t_1 < \dots < t_n$  : the increments are independent  
 $X_{t_1} - X_{t_0}, X_{t_2} - X_{t_1}, \dots, X_{t_n} - X_{t_{n-1}}$  (17.43)

##### Interpretation

A Poisson Process is a *continuous-time* process with *discrete*, *positive* realizations in  $\mathbb{N}_{\geq 0}$

**Corollary 17.10 Probability of events:** Using Taylor in order to expand the Poisson distribution one obtains:

$$\mathbb{P}(X_{(t+\Delta t)} - X_t \neq 0) = \lambda \Delta t + o(\Delta t^2) \quad t \text{ small i.e. } t \rightarrow 0 \quad (17.44)$$

1. Thus the probability of an event happening during  $\Delta t$  is proportional to time period and the rate  $\lambda$
2. The probability of two or more events to happen *during*  $\Delta t$  is of order  $o(\Delta t^2)$  and thus extremely small (as  $\Delta t$  is small).

**Definition 17.24 Differential of a Poisson Process:** The differential of a Poisson Process is defined as:

$$dX_t = \lim_{\Delta t \rightarrow dt} (X_{(t+\Delta t)} - X_t) \quad (17.45)$$

**Property 17.4 Probability of Events for differential:**  
With the definition of the differential and using the previous results from the Taylor expansion it follows:

$$\mathbb{P}(dX_t = 0) = 1 - \lambda \quad (17.46)$$

$$\mathbb{P}(|dX_t| = 1) = \lambda \quad (17.47)$$

##### Proofs

**Proof.** eq. (17.11):  
Let by  $\delta$  denote the displacement of a particle at each step, and assume that the particles start at the center i.e.  $x(0) = 0$ , then we have:

$$\begin{aligned} \mathbb{E}[x(n)] &= \mathbb{E}\left[\frac{1}{N} \sum_{i=1}^N x_i(n)\right] = \frac{1}{N} \sum_{i=1}^N \mathbb{E}[x_i(n-1) \pm \delta] \\ &= \frac{1}{N} \sum_{i=1}^N \mathbb{E}[x_i(n-1)] \\ &\stackrel{\text{induction}}{=} \mathbb{E}[x_{n-1}] = \dots \mathbb{E}[x(0)] = 0 \end{aligned}$$

Thus in expectation the particles goes nowhere.  $\square$

**Proof.** eq. (17.12):

Let by  $\delta$  denote the displacement of a particle at each step, and assume that the particles start at the center i.e.  $x(0) = 0$ , then we have:

$$\begin{aligned} \mathbb{E}[x(n)^2] &= \mathbb{E}\left[\frac{1}{N} \sum_{i=1}^N x_i(n)^2\right] = \frac{1}{N} \sum_{i=1}^N \mathbb{E}[x_i(n-1) \pm \delta]^2 \\ &= \frac{1}{N} \sum_{i=1}^N \mathbb{E}[x_i(n-1)^2 \pm 2\delta x_i(n-1) + \delta^2] \\ &\stackrel{\text{ind.}}{=} \mathbb{E}[x_{n-1}^2] + \delta^2 = \mathbb{E}[x_{n-2}^2] + 2\delta^2 = \dots \\ &= \mathbb{E}[x(0)^2] + n\delta^2 = n\delta^2 \end{aligned}$$

as  $n = \frac{\text{time}}{\text{step-size}} = \frac{t}{\Delta x}$  it follows:

$$\sigma^2 = \mathbb{E}[x^2(n)] - \mathbb{E}[x(n)]^2 = \mathbb{E}[x^2(n)] = \frac{\delta^2}{\Delta x} t \quad (17.48)$$

Thus in expectation the particles goes nowhere.  $\square$

**Proof.** eq. (17.30):

$$\begin{aligned} \gamma(\epsilon[k], \epsilon[k+n]) &= \text{Cov}[\epsilon[k], \epsilon[k+1]] \\ &= \mathbb{E}[(\epsilon[k] - \mathbb{E}[\epsilon[k]]) (\epsilon[k+n] - \mathbb{E}[\epsilon[k+n]])^T] \\ &\stackrel{\text{eq. (17.27)}}{=} \mathbb{E}[(\epsilon[k]) (\epsilon[k+n])] \end{aligned} \quad \square$$

**Proof.** corollary 17.3:

Since  $B_t - B_s$  is the increment over the interval  $[s, t]$ , it is the same in distribution as the increment over the interval  $[s-s, t-s] = [0, t-s]$

$$\begin{array}{ccc} \text{Thus} & B_t - B_s \sim B_{t-s} - B_0 \\ \text{but as } B_0 \text{ is a.s. zero by definition eq. (17.15) it follows:} & B_t - B_s \sim B_{t-s} & B_{t-s} \sim \mathcal{N}(0, t-s) \end{array} \quad \square$$

**Proof.** corollary 17.3:

$$\begin{aligned} W(t) &= W(t) - \underbrace{W(0)}_{=0} \sim \mathcal{N}(0, t) \\ \Rightarrow \quad \mathbb{E}[X] &= 0 \quad \mathbb{V}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2 = t \end{aligned} \quad \square$$

**Proof.** theorem 17.2:

$$\begin{aligned} \sum_{k=0}^{N-1} [W(t_k) - W(t_{k-1})]^2 & \quad t_k = k \frac{T}{N} \\ &= \sum_{k=0}^{N-1} X_k^2 & X_k \sim \mathcal{N}\left(0, \frac{T}{N}\right) \\ &= \sum_{k=0}^{N-1} Y_k = n \left(\frac{1}{n} \sum_{k=0}^{N-1} Y_k\right) & \mathbb{E}[Y_k] = \frac{T}{N} \\ &\stackrel{\text{S.L.L.N}}{=} n \frac{T}{n} = T \end{aligned}$$

**Proof.** theorem 17.3 ③:  $W_t^2 - t$  is a martingale?  
Using the binomial formula we can write and adding  $W_s - W_s$ :

$$W_t^2 = (W_t - W_s)^2 + 2W_s(W_t - W_s) + W_s^2$$

using the expectation:

$$\begin{aligned} \mathbb{E}[W_t^2 | \mathcal{F}_s] &= \mathbb{E}[(W_t - W_s)^2 | \mathcal{F}_s] + \mathbb{E}[2W_s(W_t - W_s) | \mathcal{F}_s] \\ &\quad + \mathbb{E}[W_s^2 | \mathcal{F}_s] \\ &\stackrel{\text{eq. (17.49)}}{=} \mathbb{E}[(W_t - W_s)^2] + 2W_s \mathbb{E}[(W_t - W_s)] + W_s^2 \\ &\stackrel{\text{eq. (17.17)}}{=} \mathbb{V}[W_t - W_s] + 0 + W_s^2 \\ &\quad = t - s + W_s^2 \end{aligned}$$

from this it follows that:

$$\mathbb{E}[W_t^2 - t | \mathcal{F}_s] = W_s^2 - s \quad \square$$

understand why  $\mathbb{E}[(W_t - W_s)^2 | \mathcal{F}_s] = \mathbb{E}[(W_t - W_s)^2]$

#### Examples

##### Example 17.1 :

Suppose we have a sample space of four elements:  $\Omega = \{\omega_1, \omega_2, \omega_3, \omega_4\}$ . At time zero, we do not have any information about which  $\omega$  has been chosen. At time  $T/2$  we know whether we have  $\{\omega_1, \omega_2\}$  or  $\{\omega_3, \omega_4\}$ . At time  $T$ , we have full information.

$$\mathcal{F} = \begin{cases} \{\emptyset, \Omega\} & t \in [0, T/2) \\ \{\emptyset, \{\omega_1, \omega_2\}, \{\omega_3, \omega_4\}, \Omega\} & t \in [T/2, T) \\ \mathcal{F}_{\max} = 2^\Omega & t = T \end{cases} \quad (17.50)$$

Thus,  $\mathcal{F}_0$  represents initial information whereas  $\mathcal{F}_\infty$  represents full information (all we will ever know). Hence, a stochastic process is said to be defined on a filtered probability space  $(\Omega, \mathcal{F}, \{\mathcal{F}_t\}_{t \geq 0}, \mathbb{P})$ .

## Ito Calculus

**Proof.** theorem 17.3 ②:

1. first we need to show eq. (17.3):  $\mathbb{E}[W_t | \mathcal{F}_s] = W_s$   
Due to the fact that  $W_t$  is  $\mathcal{F}_t$  measurable i.e.  $W_t \in \mathcal{F}_t$  we know that:

$$\begin{aligned} \mathbb{E}[W_t | \mathcal{F}_t] &= W_t & (17.49) \\ \mathbb{E}[W_t | \mathcal{F}_s] &= \mathbb{E}[W_t - W_s + W_s | \mathcal{F}_s] \\ &= \mathbb{E}[W_t - W_s | \mathcal{F}_s] + \mathbb{E}[W_s | \mathcal{F}_s] \\ &\stackrel{\text{eq. (17.49)}}{=} \mathbb{E}[W_t - W_s] + W_s \\ &\stackrel{W_t - W_s \sim \mathcal{N}(0, t-s)}{=} W_s \end{aligned}$$

2. second we need to show eq. (17.4):  $\mathbb{E}[|X(t)|] < \infty$   
 $\mathbb{E}[|W(t)|]^2 \stackrel{\text{eq. (16.33)}}{\leq} \mathbb{E}[|W(t)|^2] = \mathbb{E}[W^2(t)] = t < \infty$