



Math Appendix

Logic  
Set Theory

<b>Definition 2.1 Collection/Multiset:</b> Is a set-like object in which multiplicity matters (order does not). I.e. $\{1, 1, 2, 3\} \neq \{1, 2, 3\}$
<b>Definition 2.2 Cardinality</b> $ S $ : Is the number of elements that are contained in a set.
<b>Definition 2.3 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>Example 2.1 Power Set/Cardinality of <math>S = \{x, y, z\}</math>:</b> 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 3.1 Index Set:</b> Is a set?? $A$ , 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. $f : A \mapsto S \qquad A \in \mathbb{N} \qquad (3.1)$
<b>Definition 3.2 Sequence</b> $(a_n)_{n \in A}$ : is an by an index set $A$ <i>enumerated</i> collection <sup>[def. 2.1]</sup> of objects in which repetitions are allowed and <i>order does matter</i> .
<b>Definition 3.3 Series:</b> is an infinite ordered set of terms combined together by addition.
<b>1. Types of Sequences</b>
<b>1.1. Arithmetic Sequence</b>
<b>Definition 3.4 Arithmetic Sequence:</b> Is a sequence where the <i>difference</i> between two consecutive terms constant i.e. $(2, 4, 6, 8, 10, 12, \dots)$ . $t_n = t_0 + nd \quad d : \text{difference between two terms} \qquad (3.2)$
<b>1.2. Geometric Sequence</b>
<b>Definition 3.5 Geometric Sequence:</b> Is a sequence where the <i>ratio</i> between two consecutive terms constant i.e. $(2, 4, 8, 16, 32, \dots)$ . $t_n = t_0 \cdot r^n \qquad r : \text{ratio between two terms} \qquad (3.3)$

# Calculus and Analysis

**Definition 4.1 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 4.2 Discriminant:**  $\delta = b^2 - 4ac$

**Definition 4.3 Solution to** <sup>[def. 4.1]</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 4.1**  
**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:

$$F(X) = \int_a^x f(t) dt \quad (4.1)$$

Then it follows:

$$F'(x) = f(x) \quad \forall x \in (a, b) \quad (4.2)$$

**Theorem 4.2**  
**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  $[a, b]$ :

$$\int_a^b f(t) dt = F(b) - F(a) \iff \int_a^x \frac{\partial}{\partial x} F(t) dt = F(x) \quad (4.3)$$

**Definition 4.4 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 4.5**  
**Codomain/target set of a function**  $\text{codom}(\cdot)$ :  
Given a function  $f : \mathcal{X} \rightarrow \mathcal{Y}$ , the codomain of that function is the set  $\mathcal{Y}$  into which all of the output of the function is constrained to fall.

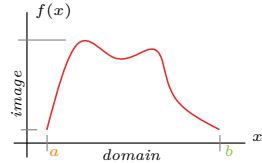
**Definition 4.6 Image (Range) of a function:**  $f[\cdot]$

Given a function  $f : \mathcal{X} \rightarrow \mathcal{Y}$ , the image of that function is the set to which the function can actually map:

$$\{y \in \mathcal{Y} | y = f(x), \quad \forall x \in \mathcal{X}\} := f[\mathcal{X}] \quad (4.4)$$

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 4.7 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  $f$ :

$$f^{-1}(A) = \{x \subseteq X : f(x) \subseteq A\} \quad (4.5)$$

**Example 4.1 :**

**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:

$$f[A] = \{y \in \mathcal{Y} | y = f(x), \quad \forall x \in A\} \quad (4.6)$$

**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 4.8 (strictly) Increasing Functions:**  
A function  $f$  is called **monotonically increasing/ increasing/non-decreasing** if:

$$x \leq y \iff f(x) \leq f(y) \quad \forall x, y \in \text{dom}(f) \quad (4.7)$$

And **strictly increasing** if:

$$x < y \iff f(x) < f(y) \quad \forall x, y \in \text{dom}(f) \quad (4.8)$$

**Definition 4.9 (strictly) Decreasing Functions:**  
A function  $f$  is called monotonically decreasing/decreasing or non-increasing if:

$$x \geq y \iff f(x) \geq f(y) \quad \forall x, y \in \text{dom}(f) \quad (4.9)$$

And **strictly decreasing** if:

$$x > y \iff f(x) > f(y) \quad \forall x, y \in \text{dom}(f) \quad (4.10)$$

**Definition 4.10 Monotonic Function:** A function  $f$  is called monotonic iff either  $f$  is **increasing** or **decreasing**.

**Definition 4.11 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 4.1 Linearity of Differentiation:** The derivative of **any** linear combination of functions equals the same linear combination of the derivatives of the functions:

$$\frac{d}{dx} (a f(x) + b g(x)) = a \frac{d}{dx} f(x) + b \frac{d}{dx} g(x) \quad a, b \in \mathbb{R} \quad (4.11)$$

**Definition 4.12 Quadratic Function:**

A function  $f : \mathbb{R}^n \mapsto \mathbb{R}^m$  is quadratic if it can be written in the form:

$$f(x) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} + c \quad (4.12)$$

## 1. Continuity and Smoothness

**Definition 4.13 Continuous Function:**

**Definition 4.14 Smoothness of a Function  $C^k$ :** 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:  
 $f \in C^k(\mathcal{X}) \iff \exists f', f'', \dots, f^{(k)}$  continuous (4.13)

**Note**

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

**Corollary 4.2 Smooth Function  $C^\infty$ :** Is a function  $f : \mathcal{X} \rightarrow \mathcal{Y}$  that has derivatives infinitely many times differentiable.  
 $f \in C^\infty(\mathcal{X}) \iff f', f'', \dots, f^{(\infty)}$  (4.14)

**Corollary 4.3 Continuously Differentiable Function  $C^1$ :** 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:

$$f \in C^1(\mathcal{X}) \iff f' \text{ continuous} \quad (4.15)$$

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 4.15 Lipschitz Continuity:** A Lipschitz continuous function is a function  $f$  whose rate of change is bound by a Lipschitz Contant  $L$ :

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq L \|\mathbf{x} - \mathbf{y}\|_2^2 \quad \forall \mathbf{x}, \mathbf{y}, \quad L > 0 \quad (4.16)$$

**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 4.16 Lipschitz Continuous Gradient:**

A continuously differentiable function  $f : \mathbb{R}^d \mapsto \mathbb{R}$  has  $L$ -Lipschitz continuous gradient if it satisfies:

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\| \quad \forall \mathbf{x}, \mathbf{y} \in \text{dom}(f), \quad L > 0 \quad (4.17)$$

if  $f \in C^2$ , this is equivalent to:

$$\nabla^2 f(\mathbf{x}) \leq L \mathbf{I} \quad \forall \mathbf{x} \in \text{dom}(f), \quad L > 0 \quad (4.18)$$

**Lemma 4.1 Descent Lemma:** If a function  $f : \mathbb{R}^d \mapsto \mathbb{R}$  has Lipschitz continuous gradient eq. (4.17) over its domain, then it holds that:

$$|f(\mathbf{x}) - f(\mathbf{y}) - \nabla f(\mathbf{y})^T (\mathbf{x} - \mathbf{y})| \leq \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2 \quad (4.19)$$

**Note**

If  $f$  is twice differentiable then the largest eigenvalue of the Hessian (<sup>[def. 5.5]</sup>) of  $f$  is uniformly upper bounded by  $L$

*Proof.* lemma 4.1 for  $C^1$  functions:

Let  $g(t) \equiv f(\mathbf{y} + t(\mathbf{x} - \mathbf{y}))$  from the FToc (theorem 4.2) we know that:

$$\int_0^1 g'(t) dt = g(1) - g(0) = f(\mathbf{x}) - f(\mathbf{y})$$

It then follows from the reverse:

$$\begin{aligned} & |f(\mathbf{x}) - f(\mathbf{y}) - \nabla f(\mathbf{y})^T (\mathbf{x} - \mathbf{y})| \\ & \stackrel{\text{Chain. R.}}{\stackrel{\text{FToc}}{=}} \left| \int_0^1 \nabla f(\mathbf{y} + t(\mathbf{x} - \mathbf{y}))^T (\mathbf{x} - \mathbf{y}) dt - \nabla f(\mathbf{y})^T (\mathbf{x} - \mathbf{y}) \right| \\ & = \left| \int_0^1 (\nabla f(\mathbf{y} + t(\mathbf{x} - \mathbf{y})) - \nabla f(\mathbf{y}))^T (\mathbf{x} - \mathbf{y}) dt \right| \\ & = \left| \int_0^1 (\nabla f(\mathbf{y} + t(\mathbf{x} - \mathbf{y})) - \nabla f(\mathbf{y}))^T (\mathbf{x} - \mathbf{y}) dt \right| \\ & \stackrel{\text{C.S.}}{\leq} \left| \int_0^1 \|\nabla f(\mathbf{y} + t(\mathbf{x} - \mathbf{y})) - \nabla f(\mathbf{y})\| \cdot \|\mathbf{x} - \mathbf{y}\| dt \right| \\ & \stackrel{\text{eq. (4.17)}}{=} \left| \int_0^1 L \|\mathbf{y} + t(\mathbf{x} - \mathbf{y}) - \mathbf{y}\| \cdot \|\mathbf{x} - \mathbf{y}\| dt \right| \\ & = \left| L \|\mathbf{x} - \mathbf{y}\|^2 \int_0^1 t dt \right| = \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2 \end{aligned}$$

*Proof.* lemma 4.1 for  $C^2$  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. (4.20):

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

**Definition 4.17 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

$$\forall \mathbf{x}, \mathbf{y} \in \text{dom}(f), \quad L > 0 \quad (4.20)$$

If  $f$  is a twice differentiable this is equivalent to:

$$\nabla^2 f(\mathbf{x}) \leq L \mathbf{I} \quad L > 0 \quad (4.21)$$

**Theorem 4.3**

**L-Smoothness of convex functions:** A convex and  $L$ -Smooth function (<sup>[def. 4.17]</sup>) has a Lipschitz continuous gradient (eq. (4.17)) thus it holds that:

$$f(\mathbf{x}) \leq f(\mathbf{y}) + \nabla f(\mathbf{y})^T (\mathbf{x} - \mathbf{y}) \leq \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2 \quad (4.22)$$

*Proof.* theorem 4.3:

With the definition of convexity for a differentiable function (eq. (4.25)) it follows

$$\begin{aligned} f(\mathbf{x}) - f(\mathbf{y}) + \nabla f(\mathbf{y})^T (\mathbf{x} - \mathbf{y}) & \geq 0 \\ \Rightarrow |f(\mathbf{x}) - f(\mathbf{y}) + \nabla f(\mathbf{y})^T (\mathbf{x} - \mathbf{y})| & \\ \text{if eq. (4.25)} & f(\mathbf{x}) - f(\mathbf{y}) + \nabla f(\mathbf{y})^T (\mathbf{x} - \mathbf{y}) \end{aligned}$$

with lemma 4.1 and <sup>[def. 4.17]</sup> it follows theorem 4.3  $\square$

**Corollary 4.4 :**  $L$ -smoothness is a weaker condition than  $L$ -Lipschitz continuous gradients

## 2. Convexity

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

**Definition 4.18 Convex Functions:**

A function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is convex if it satisfies:

$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}) \quad \forall \mathbf{x}, \mathbf{y} \in \text{dom}(f) \quad \forall \lambda \in [0, 1] \quad (4.23)$$

*Include figure from the/convexity*

**Definition 4.19 Concave Functions:**

A function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is convex if it satisfies:

$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \geq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}) \quad \forall \mathbf{x}, \mathbf{y} \in \text{dom}(f) \quad \forall \lambda \in [0, 1] \quad (4.24)$$

**Corollary 4.5 Convexity  $\rightarrow$  global minimima:** Convexity implies that all local minima (if they exist) are global minima.

**Definition 4.20 Strictly 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:

$$f(\mathbf{x}) \geq f(\mathbf{y}) + \nabla f(\mathbf{y})^T (\mathbf{x} - \mathbf{y}) \quad \forall \mathbf{x}, \mathbf{y} \in \text{dom}(f) \quad (4.25)$$

If  $f$  is a twice differentiable function this is equivalent to:

$$\nabla^2 f(\mathbf{x}) \geq 0 \quad \forall \mathbf{x}, \mathbf{y} \in \text{dom}(f) \quad (4.26)$$

**Intuition**

- Convexity implies that a function  $f$  is bound by/below a linear interpolation from  $\mathbf{x}$  to  $\mathbf{y}$  and strong convexity that  $f$  is strictly bound/below.
- eq. (4.25) 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 4.6 Strict Convexity  $\rightarrow$  Uniqueness:** Strict convexity implies a unique minimizer  $\iff$  at most one global minimum.

**Corollary 4.7 :** 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  $\mathcal{X}$ :

$$f''(x) \geq 0 \quad \forall x \in \mathcal{X} \quad (4.27)$$

**Definition 4.21  $\mu$ -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:

$$f(tx + (1-t)y) \leq tf(x) + (1-t)f(y) - \frac{t(1-t)}{2} \mu \|x - y\|$$

$$\forall x, y \in \mathcal{X}, \quad t \in [0, 1], \quad \mu > 0$$

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

$$f(y) \geq f(x) + \nabla f(x)^\top (y - x) + \frac{\mu}{2} \|y - x\|_2^2 \quad (4.28)$$

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

$$\nabla^2 f(x) \geq \mu I \quad \forall x, y \in \mathcal{X} \quad \mu > 0 \quad (4.29)$$

**Corollary 4.8 Strong Convexity implies Strict Convexity:**  
<https://math.stackexchange.com/question/2090991/proof-for-strongly-convex-function-is-strictly-convex>

**Property 4.1 :**  
 $f(y) \leq f(y) + \nabla f(y)^\top (x - y) + \frac{1}{2\mu} \|\nabla f(x) - \nabla f(y)\|_2^2 \quad (4.30)$

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

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

*Proof.* eq. (4.29) analogously to **Proof** eq. (4.21) □

**Note**  
 If  $f$  is twice differentiable then the smallest eigenvalue of the Hessian <sup>[def. 5.5]</sup> of  $f$  is uniformly lower bounded by  $\mu$ .  
**Hence** strong convexity can be considered as the analogous to smoothness

**Example 4.2 Quadratic Function:** A quadratic function eq. (4.12) is convex if:

$$\nabla_x^2 \text{eq. (4.12)} = A \geq 0 \quad (4.31)$$

**Corollary 4.9 :**  
 Strong convexity  $\implies$  Strict convexity  $\implies$  Convexity

### 2.1. Properties that preserve convexity

**Property 4.2 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 4.3 Composition of Affine Mappings:** Let  $f$  be a convex function then  $g(x)$  is convex as well:

$$g(x) = f(Ax + b)$$

**Property 4.4 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.

$$f(-x) = f(x) \quad (4.32)$$

**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.

$$f(-x) = -f(x) \quad (4.33)$$

**Theorem 4.4 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:**  $\cos x, |x|, c, x^2, x^4, \dots \exp(-x^2/2)$ .  
**Odd:**  $\sin x, \tan x, x, x^3, x^5, \dots$

**x-Shift:**  $f(x - c) \implies$  shift to the right  
 $f(x + c) \implies$  shift to the left (4.34)  
**y-Shift:**  $f(x) \pm c \implies$  shift up/down (4.35)

*Proof.* eq. (4.34)  $f(x_n - c)$  we take the x-value at  $x_n$  but take the y-value at  $x_o := x_n - c$   
 $\implies$  we shift the function to  $x_n$ . □

**Euler's formula**  

$$e^{\pm ix} = \cos x \pm i \sin x \quad (4.36)$$

**Euler's Identity**  

$$e^{\pm i} = -1 \quad (4.37)$$

**Note**  

$$e^n = 1 \iff n = i2\pi k, \quad k \in \mathbb{N} \quad (4.38)$$

**Corollary 4.10 Every norm is a convex function:** By using definition <sup>[def. 4.18]</sup> and the triangular inequality it follows (with the exception of the L0-norm):

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

### 2.2. Taylor Expansion

**Definition 4.22 Taylor Expansion:**  

$$T_n(x) = \sum_{i=0}^n \frac{1}{n!} f^{(i)}(x_0) \cdot (x - x_0)^{(i)} \quad (4.39)$$

$$= f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2} f''(x_0)(x - x_0)^2 + \mathcal{O}(x^3) \quad (4.40)$$

**Definition 4.23 Incremental Taylor:**  
**Goal:** evaluate  $T_n(x)$  (eq. (4.40)) at the point  $x_0 + \Delta x$  in order to propagate the function  $f(x)$  by  $h = \Delta x$ :

$$T_n(x_0 \pm h) = \sum_{i=0}^n \frac{h^i}{n!} f^{(i)}(x_0) i^{-1} \quad (4.41)$$

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

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

**Definition 4.24 Multidimensional Taylor:** Suppose  $X \in \mathbb{R}^n$  is open,  $x \in X$ ,  $f : X \mapsto \mathbb{R}$  and  $f \in \mathcal{C}^2$  then it holds that

$$f(x) \approx f(x_0) + \nabla_x f(x_0)(x - x_0) + \frac{1}{2} (x - x_0)^\top H(x - x_0) \quad (4.42)$$

**Definition 4.25 Argmax:** The argmax of a function defined on a set  $D$  is given by:

$$\arg \max_{x \in D} f(x) = \{x | f(x) \geq f(y), \forall y \in D\} \quad (4.43)$$

**Definition 4.26 Argmin:** The argmin of a function defined on a set  $D$  is given by:

$$\arg \min_{x \in D} f(x) = \{x | f(x) \leq f(y), \forall y \in D\} \quad (4.44)$$

**Corollary 4.11 Relationship**  $\arg \min \leftrightarrow \arg \max$ :

$$\arg \min_{x \in D} f(x) = \arg \max_{x \in D} -f(x) \quad (4.45)$$

**Property 4.5 Argmax Identities:**

- Shifting:**  $\forall \lambda \text{ const} \quad \arg \max f(x) = \arg \max f(x) + \lambda \quad (4.46)$
- Positive Scaling:**  $\forall \lambda > 0 \text{ const} \quad \arg \max f(x) = \arg \max \lambda f(x) \quad (4.47)$
- Negative Scaling:**  $\forall \lambda < 0 \text{ const} \quad \arg \max f(x) = \arg \min \lambda f(x) \quad (4.48)$
- Positive Functions:**  $\forall \arg \max f(x) > 0, \forall x \in \text{dom}(f)$   

$$\arg \max f(x) = \arg \min \frac{1}{f(x)} \quad (4.49)$$
- Stricly Monotonic Functions:** for all strictly monotonic increasing functions <sup>[def. 4.8]</sup>  $g$  it holds that:  

$$\arg \max g(f(x)) = \arg \max f(x) \quad (4.50)$$

**Definition 4.27 Max:** The maximum of a function  $f$  defined on the set  $D$  is given by:

$$\max_{x \in D} f(x) = f(x^*) \quad \text{with} \quad \forall x^* \in \arg \max f(x) \quad (4.51)$$

**Definition 4.28 Min:** The minimum of a function  $f$  defined on the set  $D$  is given by:

$$\min_{x \in D} f(x) = f(x^*) \quad \text{with} \quad \forall x^* \in \arg \min f(x) \quad (4.52)$$

**Corollary 4.12 Relationship**  $\min \leftrightarrow \max$ :

$$\min_{x \in D} f(x) = - \max_{x \in D} -f(x) \quad (4.53)$$

**Property 4.6 Max Identities:**

- Shifting:**  $\forall \lambda \text{ const} \quad \max \{f(x) + \lambda\} = \lambda + \max f(x) \quad (4.54)$
- Positive Scaling:**  $\forall \lambda > 0 \text{ const} \quad \max \lambda f(x) = \lambda \max f(x) \quad (4.55)$
- Negative Scaling:**  $\forall \lambda < 0 \text{ const} \quad \max \lambda f(x) = \lambda \min f(x) \quad (4.56)$
- Positive Functions:**  $\forall \arg \max f(x) > 0, \forall x \in \text{dom}(f) \quad \max \frac{1}{f(x)} = \frac{1}{\min f(x)} \quad (4.57)$
- Stricly Monotonic Functions:** for all strictly monotonic increasing functions <sup>[def. 4.8]</sup>  $g$  it holds that:  

$$\max g(f(x)) = g(\max f(x)) \quad (4.58)$$

**Definition 4.29 Supremum:** The supremum of a function defined on a set  $D$  is given by:

$$\sup_{x \in D} f(x) = \{y | y \geq f(x), \forall x \in D\} = \min_{y | y \geq f(x), \forall x \in D} y \quad (4.59)$$

and is the smallest value  $y$  that is equal or greater  $f(x)$  for any  $x \iff$  smallest upper bound.

**Definition 4.30 Infimum:** The infimum of a function defined on a set  $D$  is given by:

$$\inf_{x \in D} f(x) = \{y | y \leq f(x), \forall x \in D\} = \max_{y | y \leq f(x), \forall x \in D} y \quad (4.60)$$

and is the biggest value  $y$  that is equal or smaller  $f(x)$  for any  $x \iff$  largest lower bound.

**Corollary 4.13 Relationship**  $\sup \leftrightarrow \inf$ :

$$\inf_{x \in D} f(x) = - \sup_{x \in D} -f(x) \quad (4.61)$$

**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  $x \implies$  there exists no argmax/argmin  $\implies$  need to bound the functions from above/below  $\iff$  infimum/supremum.

**Definition 4.31 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.

$$y(t) = f(x(t), t) \xrightarrow{\text{time-invariance}} y(t - \tau) = f(x(t - \tau), t) \quad \forall \tau \quad (4.62)$$

**Definition 4.32 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

$$f(g(x)) = x \quad \forall x \in \text{dom}(g) \quad (4.63)$$

and

$$g(f(u)) = u \quad \forall u \in \text{dom}(f) \quad (4.64)$$

**Property 4.7 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 4.5 The Existence of an Inverse Function:**  
 A function has an inverse function if and only if it is one-to-one.

**Corollary 4.14 Inverse functions and strict monotonicity:** If a function  $f$  is **strictly monotonic** <sup>[def. 4.10]</sup> on its entire domain, then it is one-to-one and therefore has an inverse function.

## 3. Special Functions

### 3.1. The Gamma Function

**Definition 4.33 The gamma function**  $\Gamma(\alpha)$ : Is extension of the factorial function (??) to the real and complex numbers (with a positive real part):

$$\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} dx \quad \Re(z) > 0 \quad (4.65)$$

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

Differential Calculus

**Definition 5.1 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 5.2 Second Derivative**  $\frac{\partial^2}{\partial x_i \partial x_j}$ :

**Corollary 5.1 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. 5.1] 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 5.3 Gradient:** Given  $f : n \mapsto \mathbb{R}$  its gradient is defined as:

$$\text{grad}_{\mathbf{x}}(f) = \nabla_{\mathbf{x}} f := \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} & \cdots & \frac{\partial f}{\partial x_n} \end{bmatrix} \quad (5.1)$$

**Definition 5.4 Jacobi Matrix:** Given a vector valued function  $f : \mathbb{R}^n \mapsto \mathbb{R}^m$  its derivative/Jacobian is defined as:

$$\mathbf{J}(f(\mathbf{x})) = \mathbf{J}_f(\mathbf{x}) = \mathbf{D}f = \frac{\partial \mathbf{f}}{\partial \mathbf{x}}(\mathbf{x}) = \frac{\partial (f_1, \dots, f_m)}{\partial (x_1, \dots, x_n)}(\mathbf{x}) = \quad (5.2)$$

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

**Theorem 5.1 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 5.5 Hessian Matrix:**  
Given a function  $f : \mathbb{R} \mapsto \mathbb{R}^n$  its Hessian  $\in \mathbb{R}^{n \times n}$  is defined as:

$$\mathbf{H}(f)(\mathbf{x}) = \mathbf{H}_f(\mathbf{x}) = \mathbf{J}(\nabla f(\mathbf{x}))^T \quad (5.3)$$

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

and it corresponds to the Jacobian of the Gradient.  
Due to the differentiability and theorem 5.1 it follows that the Hessian is (if it exists):

- Symmetric
- Real

**Corollary 5.2 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 5.3 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. 5.1] 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 6.1 Important Integral Properties:

**Addition**  $\int\limits_a^b f(x) \, dx = \int\limits_a^c f(x) \, dx + \int\limits_c^b f(x) \, dx$  (6.1)

**Reflection**  $\int\limits_a^b f(x) \, dx = - \int\limits_b^a f(x) \, dx$  (6.2)

**Translation**  $\int\limits_a^b f(x) \, dx \stackrel{u:=x\pm c}{=} \int\limits_{a\pm c}^{b\pm c} f(x \mp c) \, dx$  (6.3)

**f Odd**  $\int\limits_{-a}^a f(x) \, dx = 0$  (6.4)

**f Even**  $\int\limits_{-a}^a f(x) \, dx = 2 \int\limits_0^a f(x) \, dx$  (6.5)

Proof. eqs. (6.4) and (6.5)

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

□

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 7.1 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. Eigenvalues and Vectors

**Formula 7.1 Eigenvalues of a 2x2 matrix:** Given a 2x2-matrix  $A$  its eigenvalues can be calculated by:

$$\{\lambda_1, \lambda_2\} \in \frac{\text{tr}(A) \pm \sqrt{\text{tr}(A)^2 - 4 \det(A)}}{2} \quad (7.1)$$

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

2. Special Kind of Matrices

**Definition 7.1 Hermitian Matrices:**

$$A = A^H \quad (7.2)$$

3. Spaces and Measures

**Definition 7.2 Bilinear Form/Functional:**  
Is a mapping  $a : \mathcal{Y} \times \mathcal{Y} \mapsto F$  on a field of scalars  $F \subseteq \mathbb{K}, K = \mathbb{R}$  or  $\mathbb{C}$  that satisfies:

$$\begin{aligned} a(\alpha u + \beta v, w) &= \alpha a(u, w) + \beta a(v, w) \\ a(u, \alpha v + \beta w) &= \alpha a(u, v) + \beta a(u, w) \end{aligned}$$

$\forall u, v, w \in \mathcal{Y}, \quad \forall \alpha, \beta \in \mathbb{K}$

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

**Definition 7.3 Symmetric bilinear form:** A bilinear form  $a$  on  $\mathcal{Y}$  is symmetric if and only if:

$$a(u, v) = a(v, u) \quad \forall u, v \in \mathcal{Y}$$

**Definition 7.4 Positive (semi) definite bilinear form:** A symmetric bilinear form  $a$  on a vector space  $\mathcal{Y}$  over a field  $F$  is **positive defintie** if and only if:

$$a(u, u) > 0 \quad \forall u \in \mathcal{Y} \setminus \{0\} \quad (7.3)$$

And **positive semidefinte**  $\iff \geq$  (7.4)

**Corollary 7.1 Matrix induced Bilinear Form:**  
For finite dimensional inner product spaces  $\mathcal{X}' \in \mathbb{K}^n$  any matrix  $A \in \mathbb{R}^{n \times n}$  induces a **bilinear form**:

$$a(x, x') = x^\top A x' = (Ax')^\top x,$$

**Definition 7.5 Positive (semi) definite Matrix  $\succ$ :**  
A matrix  $A \in \mathbb{R}^{n \times n}$  is **positive defintie** if and only if:

$$x^\top A x > 0 \iff A \succ \quad \forall x \in \mathbb{R}^n \setminus \{0\} \quad (7.5)$$

And **positive semidefinte**  $\iff \geq$  (7.6)

**Corollary 7.2 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*.

$$\forall \lambda_i \in \text{eigenv}(A) > 0 \quad (7.7)$$

And **positive semidefinte**  $\iff \geq$  (7.8)

*Proof.* corollary 7.2 (for real matrices):  
Let  $v$  be an eigenvector of  $A$  then it follows:

$$0 < v^\top A v = v^\top \lambda v = \lambda \|v\|^2$$

□

**Corollary 7.3 Positive Definiteness and Determinant:**  
The determinant of a positive definite matrix is always positive. **Thus** a positive definite matrix is always *nonsingular*

**Definition 7.6 Negative (semi) definite Matrix  $\prec$ :**  
A matrix  $A \in \mathbb{R}^{n \times n}$  is **negative defintie** if and only if:

$$x^\top A x < 0 \iff A < 0 \quad \forall x \in \mathbb{R}^n \setminus \{0\} \quad (7.9)$$

And **negative semidefinte**  $\iff \leq$  (7.10)

**Theorem 7.2 Sylvester's criterion:** Let  $A$  be *symmetric/Hermitian* matrix and denote by  $A^{(k)}$  the  $k \times k$  upper left sub-matrix of  $A$ .  
Then it holds that:

- $A > 0 \iff \det(A^{(k)}) > 0 \quad k = 1, \dots, n$  (7.11)
- $A < 0 \iff (-1)^k \det(A^{(k)}) > 0 \quad k = 1, \dots, n$  (7.12)

- $A$  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.

4. Inner Products

**Definition 7.7 Inner Product:** Let  $\mathcal{Y}$  be a vector space over a field  $F \in \mathbb{K}$  of scalars. An inner product on  $\mathcal{Y}$  is a map:

$$\langle \cdot, \cdot \rangle : \mathcal{Y} \times \mathcal{Y} \mapsto F \subseteq \mathbb{K} \quad K = \mathbb{R} \text{ or } \mathbb{C} \quad (7.13)$$

that satisfies:

$$\forall x, y, z \in \mathcal{Y}, \quad \alpha, \beta \in F$$

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

**Corollary 7.4 Inner product $\mapsto$ S.p.d. Bilinear Form:**  
Let  $\mathcal{Y}$  be a vector space over a field  $F \in \mathbb{K}$  of scalar.  
An **inner product** on  $\mathcal{Y}$  is a positive definite symmetric bilinear form on  $\mathcal{Y}$ .  
**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 7.9 Norm  $\|\cdot\|_{\mathcal{Y}}$ :** A norm measures the **size** of its argument.  
**Formally** let  $\mathcal{Y}$  be a vector space over a field  $F$ , a norm on  $\mathcal{Y}$  is a map:

$$\|\cdot\|_{\mathcal{Y}} : \mathcal{Y} \mapsto \mathbb{R}_+ \quad (7.14)$$

that satisfies:

$$\forall x, y \in \mathcal{Y}, \quad \alpha \in F \subseteq \mathbb{K} \quad K = \mathbb{R} \text{ or } \mathbb{C}$$

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

**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 7.5 Reverse Triangular Inequality:**

$$-\|x - y\|_{\mathcal{Y}} \leq \|x\|_{\mathcal{Y}} - \|y\|_{\mathcal{Y}} \leq \|x - y\|_{\mathcal{Y}}$$

resp.  $|\|x\|_{\mathcal{Y}} - \|y\|_{\mathcal{Y}}| \leq \|x - y\|_{\mathcal{Y}}$

**Semi-norm**

Gold

**Corollary 7.6 Normed vector space:** Is a vector space  $\mathcal{Y}$  over a field  $F$ , on which a norm  $\|\cdot\|_{\mathcal{Y}}$  can be defined.

**Corollary 7.7 Inner product induced norm  $\langle \cdot, \cdot \rangle_{\mathcal{Y}} \rightarrow \|\cdot\|_{\mathcal{Y}}$ :** Every inner product  $\langle \cdot, \cdot \rangle_{\mathcal{Y}}$  induces a norm of the form:

$$\|x\|_{\mathcal{Y}} = \sqrt{\langle x, x \rangle} \quad x \in \mathcal{Y}$$

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

**Corollary 7.8 Energy Norm:** A *s.p.d.* bilinear form  $a : \mathcal{Y} \times \mathcal{Y} \mapsto F$  induces an energy norm:

$$\|x\|_a := (a(x, x))^{\frac{1}{2}} = \sqrt{a(x, x)} \quad x \in \mathcal{Y}$$

**Definition 7.10 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:**  $d(x, z) \leq d(x, y) + d(y, z)$

**Definition 7.11 Metric:** Is a distance measure that additionally satisfies:

$$\forall x, y \in S$$

**identity of indiscernibles :**  $d(x, y) = 0 \iff x = y$

**Corollary 7.9 Metric $\rightarrow$ Norm:** Every norm  $\|\cdot\|_{\mathcal{Y}}$  on a vector space  $\mathcal{Y}$  over a field  $F$  induces a metric by:

$$d(x, y) = \|x - y\|_{\mathcal{Y}} \quad \forall x, y \in \mathcal{Y}$$

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

- Homogenity/Scaling:**  $d(\alpha x, \alpha y)_{\mathcal{Y}} = |\alpha| d(x, y)_{\mathcal{Y}}$
- 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  $\mathcal{Y}$  satisfies the properties then it induces a norm of the form:

$$\|x\|_{\mathcal{Y}} := d(x, 0)_{\mathcal{Y}}$$

**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  $\xrightarrow{\text{induces}}$ )  
induces  $\xrightarrow{\text{induces}}$  norm  $\xrightarrow{\text{induces}}$  metric.

5. Vector Algebra

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

**6. Derivatives**

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

Geometry

Definition 8.1 Affine Transformation/Map:

Corollary 8.1 Affine Transformation in 1D: Given: numbers  $x \in \Omega$  with  $\Omega = [a, b]$   
The affine transformation of  $\phi : \Omega \rightarrow \Omega$  with  $y \in \Omega = [c, d]$  is defined by:
$$y = \phi(x) = \frac{d - c}{b - a} (x - a) + c \tag{8.1}$$

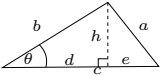
Proof. corollary 8.1 By <sup>[def. 8.1]</sup> we want a function  $f : [a, b] \rightarrow [c, d]$  that satisfies:  
 $f(a) = c$  and  $f(b) = d$   
additionally  $f(x)$  has to be a linear function (<sup>[def. 4.11]</sup>), that is the output scales the same way as the input scales.  
Thus it follows:  
 $\frac{d - c}{b - a} = \frac{f(x) - f(a)}{x - a} \iff f(x) = \frac{d - c}{b - a} (x - a) + c$

Trigonometry

Law 8.1 Law of Cosine: relates the side of a triangle to the cosine of its angles.
$$a^2 = b^2 + c^2 - 2bc \cos \theta_{b,c} \tag{8.2}$$

More general for vectors it holds:  
 $\|x - y\|^2 = \|x\|^2 + \|y\|^2 - 2\|x\|\|y\| \cos \theta_{x,y} \tag{8.3}$

Proof. eq. (8.2):  
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  $\underline{e} = c - d = c - b \cos \theta \Rightarrow a^2 = \underline{e}^2 + \underline{h}^2 \Rightarrow a$



Proof. eq. (8.3):  
 $\|x - y\|^2 = (x - y) \cdot (x - y)$   
 $= x \cdot x - 2x \cdot y + y \cdot y$   
 $= \|x\|^2 + \|y\|^2 - 2(\|x\|\|y\| \cos \theta)$

Law 8.2 Pythagorean theorem: special case of ?? for right triangle:
$$a^2 = b^2 + c^2 \tag{8.4}$$

Euler's formula

$$e^{\pm i x} = \cos x \pm i \sin x \tag{8.5}$$

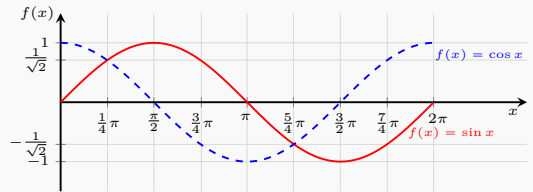
Euler's Identity

$$e^{\pm i} = -1 \tag{8.6}$$

Note

$$e^n = 1 \Leftrightarrow n = i 2 \pi k, \quad k \in \mathbb{N} \tag{8.7}$$

Sine and Cosine


$$\cos x \stackrel{(8.5)}{=} \frac{1}{2} \left[ e^{ix} + e^{-ix} \right] \tag{8.8}$$
$$\sin x \stackrel{(8.5)}{=} \frac{1}{2i} \left[ e^{ix} - e^{-ix} \right] = -\frac{i}{2} \left[ e^{ix} - e^{-ix} \right] \tag{8.9}$$

Sinh and Cosh

$$\cosh x \stackrel{(8.5)}{=} \frac{1}{2} \left[ e^x + e^{-x} \right] = \cos(ix) \tag{8.10}$$
$$\sinh x \stackrel{(8.5)}{=} \frac{1}{2} \left[ e^x - e^{-x} \right] = -i \sin(ix) \tag{8.11}$$

Note

$$e^x = \cosh x + \sinh x \quad e^{-x} = \cosh x - \sinh x \tag{8.12}$$

Note

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

Theorem 8.1 Addition Theorems:

$$\sin(\alpha \pm \beta) = \sin \alpha \cos \beta \pm \cos \alpha \sin \beta \tag{8.13}$$
$$\cos(\alpha \pm \beta) = \cos \alpha \cos \beta \mp \sin \alpha \sin \beta \tag{8.14}$$

Werner Formulas

$$\sin \alpha \cos \beta = \frac{1}{2} \left[ \sin(\alpha + \beta) + \sin(\alpha - \beta) \right] \tag{8.15}$$
$$\sin \alpha \sin \beta = \frac{1}{2} \left[ \cos(\alpha - \beta) - \cos(\alpha + \beta) \right] \tag{8.16}$$
$$\cos \alpha \cos \beta = \frac{1}{2} \left[ \cos(\alpha + \beta) + \cos(\alpha - \beta) \right] \tag{8.17}$$

Note

Using theorem 8.1 if follows:  
$$\cos(\alpha \pm \pi) = -\cos \alpha \quad \text{and} \quad \sin(\alpha \pm \pi) = -\sin \alpha \tag{8.18}$$

Topology



Numerics

**Definition 10.1 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.

0.1. Convention for iterative methods

**Definition 10.2 Linear/Exponential Convergence:** A sequence  $\{\mathbf{x}^{(k)}\}_k \in \mathbb{R}^n$  converges linearly to  $\mathbf{x}^*$  if in the asymptotic limit  $k \rightarrow \infty$  it satisfies:  
$$\left\| \mathbf{x}^{k+1} - \mathbf{x}^* \right\| \leq \rho \left\| \mathbf{x}^{(k)} - \mathbf{x}^* \right\| \quad \rho \in (0, 1), \forall k \in \mathbb{N}_0$$
 (10.1)

Exponential Convergence

Linear convergence is sometimes called exponential convergence. This is due to the fact that:

1. We often have expressions of the form:  
$$\left\| \mathbf{x}^{k+1} - \mathbf{x}^* \right\| \leq \underbrace{(1 - \alpha)}_{:= \rho} \left\| \mathbf{x}^{(k)} - \mathbf{x}^* \right\|$$
  
2. and that  $(1 - \alpha) = \exp(-\alpha)$  from which follows that:  
eq. (10.2)  $\iff \left\| \mathbf{x}^{k+1} - \mathbf{x}^* \right\| \leq e^{-\alpha} \left\| \mathbf{x}^{(k)} - \mathbf{x}^* \right\|$

**Definition 10.3 Rate of Convergence:** Is a way to measure the rate of convergence of a sequence  $\{\mathbf{x}^{(k)}\}_k \in \mathbb{R}^n$  to a value to  $\mathbf{x}^*$ . Let  $\rho \in [0, 1]$  be the rate of convergence and define:

$$\lim_{k \rightarrow \infty} \frac{\left\| \mathbf{x}^{k+1} - \mathbf{x}^* \right\|}{\left\| \mathbf{x}^{(k)} - \mathbf{x}^* \right\|} = \rho \quad (10.2)$$

- $\rho = 1 \iff$  **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 10.4 Convergence of order p:** In order to distinguish *superlinear convergence* we define the order of convergence.

A sequence  $\{\mathbf{x}^{(k)}\}_k \in \mathbb{R}^n$  converges superlinear with order  $p \in \{2, \dots\}$  to  $\mathbf{x}^*$  if it satisfies:

$$\lim_{k \rightarrow \infty} \frac{\left\| \mathbf{x}^{k+1} - \mathbf{x}^* \right\|}{\left\| \mathbf{x}^{(k)} - \mathbf{x}^* \right\|^p} = C \quad C < 1 \quad (10.3)$$

Does this even exist/check if this is true

**Definition 10.5 Exponential Convergence:** A sequence  $\{\mathbf{x}^{(k)}\}_k \in \mathbb{R}^n$  converges exponentially with rate  $\rho$  to  $\mathbf{x}^*$  if in the asymptotic limit  $k \rightarrow \infty$  it satisfies:  
$$\left\| \mathbf{x}^{k+1} - \mathbf{x}^* \right\| \leq \rho^k \left\| \mathbf{x}^{(k)} - \mathbf{x}^* \right\| \quad \rho < 1 \quad (10.4)$$

0.2. Convention for discrization methods

1. Numerical Quadrature

**Definition 10.6 Order of a Quadrature Rule:** The order of a quadrature rule  $\mathcal{Q}_n : C^0([a, b]) \rightarrow \mathbb{R}$  is defined as:  
$$\text{order}(\mathcal{Q}_n) := \max \left\{ n \in \mathbb{N}_0 : \mathcal{Q}_n(p) = \int_a^b p(t) dt \quad \forall p \in \mathcal{P}_n \right\} + 1$$
 (10.5)

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

Note

Is a quality measure for quadrature rules.

1.1. Composite Quadrature

**Definition 10.7 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  $I_j := [x_{j-1}, x_j] \quad \forall j = 1, \dots, m \triangleq$  p.w. Quadrature:

$$\int_a^b f(t) dt = \sum_{j=1}^m \int_{x_{j-1}}^{x_j} f(t) dt = \sum_{j=1}^m \mathcal{Q}_n(f|_{I_j}) \quad (10.6)$$

**Lemma 10.1 Error of Composite quadrature Rules:**  
Given a function  $f \in C^k([a, b])$  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 4.2  $\mathcal{O}\left(n^{-q}\right) = \mathcal{O}\left(h_{\mathcal{M}}^q\right) \quad \text{with } h_{\mathcal{M}} = \frac{1}{n}$  (10.7)

**Definition 10.8 Complexity W:** Is the number of function evaluations  $\triangleq$  number of quadrature points.  
$$W(\mathcal{Q}(f)_n) = \#f\text{-eval} \triangleq n$$
 (10.8)

**Lemma 10.2 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:  
$$W(\epsilon_n(f)) = \mathcal{O}(n^{\frac{1}{q}}) = \mathcal{O}\left(\sqrt[q]{\epsilon_n}\right) \quad (10.9)$$

*Proof.* lemma 10.2: **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?

$$\frac{c \cdot n_n^q}{c \cdot n_o^q} = \frac{1}{\epsilon_n} \implies n_n = n_o \cdot \sqrt[q]{\epsilon_n} = \mathcal{O}\left(\sqrt[q]{\epsilon_n}\right) \quad (10.10)$$

Optimization

**Definition 11.1 Fist Order Method:** A first-order method is an algorithm that chooses the  $k$ -th iterate in  
$$\mathbf{x}_0 + \text{span}\{\nabla f(\mathbf{x}_0), \dots, \nabla f(\mathbf{x}_{k-1})\} \quad \forall k = 1, 2, \dots \quad (11.1)$$

Note

Gradient descent is a first order method

1. Lagrangian Optimization Theory

Add: derivation of lagrange function

**Definition 11.2 (Primal) Constraint Optimization:**  
Given an optimization problem with domain  $\Omega \subseteq \mathbb{R}^d$ :

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

**Definition 11.3 Lagrange Function:**  
$$\mathcal{L}(\alpha, \beta, \mathbf{w}) := f(\mathbf{w}) + \alpha g(\mathbf{w}) + \beta h(\mathbf{w}) \quad (11.2)$$

Extremal Conditions

$$\begin{aligned} \nabla \mathcal{L}(\mathbf{x}) &\stackrel{!}{=} 0 && \text{Extremal point } \mathbf{x}^* \\ \frac{\partial}{\partial \beta} \mathcal{L}(\mathbf{x}) &= h(\mathbf{x}) \stackrel{!}{=} 0 && \text{Constraint satisfaction} \end{aligned}$$

For the inequality constraints  $g(\mathbf{x}) \leq 0$  we distinguish two situations:

Case I :  $g(\mathbf{x}^*) < 0$  switch const. off  
Case II :  $g(\mathbf{x}^*) \geq 0$  optimize using active eq. constr.

$$\frac{\partial}{\partial \alpha} \mathcal{L}(\mathbf{x}) = g(\mathbf{x}) \stackrel{!}{=} 0 \quad \text{Constraint satisfaction}$$

**Definition 11.4 Lagrangian Dual Problem:** Is given by:  
Find  $\max_{\alpha, \beta} \theta(\alpha, \beta) = \inf_{\mathbf{w} \in \Omega} \mathcal{L}(\mathbf{w}, \alpha, \beta)$   
**s.t.**  $\alpha_i \geq 0 \quad 1 \leq i \leq k$

Solution Strategy

1. Find the extremal point  $\mathbf{w}^*$  of  $\mathcal{L}(\mathbf{w}, \alpha, \beta)$ :  
$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} \Big|_{\mathbf{w}=\mathbf{w}^*} \stackrel{!}{=} 0 \quad (11.3)$$
2. Insert  $\mathbf{w}^*$  into  $\mathcal{L}$  and find the extremal point  $\beta^*$  of the resulting dual Lagrangian  $\theta(\alpha, \beta)$  for the active constraints:  
$$\frac{\partial \theta}{\partial \beta} \Big|_{\beta=\beta^*} \stackrel{!}{=} 0 \quad (11.4)$$
3. Calculate the solution  $\mathbf{w}^*(\beta^*)$  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 11.1 Upper Bound Dual Cost:** Let  $\mathbf{w} \in \Omega$  be a feasible solution of the primal problem <sup>[def. 11.2]</sup> and  $(\alpha, \beta)$  a **feasible solution** of the respective dual problem <sup>[def. 11.4]</sup>. Then it holds that:

$$f(\mathbf{w}) \geq \theta(\alpha, \beta) \quad (11.5)$$

*Proof.*

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

**Theorem 11.4 Kuhn-Tucker Conditions:** Given an optimization problem with convex domain  $\Omega \subseteq \mathbb{R}^d$ ,

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

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.:

$$\frac{\partial \mathcal{L}(\mathbf{w}, \alpha, \beta)}{\partial \mathbf{w}} \stackrel{!}{=} 0 \quad \frac{\partial \mathcal{L}(\mathbf{w}^*, \alpha, \beta)}{\partial \beta} \stackrel{!}{=} 0 \quad (11.9)$$

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:  $g_i(\mathbf{w}^*) = 0 \rightarrow \alpha_i \geq 0$  s.t.  $\alpha_i^* g_i(\mathbf{w}^*) = 0$

Consequence

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

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

Stochastics

<b>Definition 11.6 Stochastics:</b> Is a collective term for the areas of <i>probability theory</i> and <i>statistics</i> .
<b>Definition 11.7 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 11.8 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 11.9 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.
<i>Improve these definitions, maybe ask on quora/hk</i>
<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>1. Foundations</b>
<b>Definition 12.1 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 12.2 Sample Space</b> $\Omega$ : Is the set of all possible outcomes (elementary events corollary 12.5) of an experiment see example 12.1
<b>Definition 12.3 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 12.2). 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 12.1</b> : If the outcome $\omega$ of an experiment is in the subset $A$ , then the event $A$ is said to “have occurred”.
<b>Corollary 12.2 Complement Set</b> $A^C$ : is the contrary event of $A$ .
<b>Corollary 12.3 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 12.4 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 12.5 The Elementary Event</b> $\omega$ : Is a “singleton”, i.e. a subset $\{\omega\}$ containing a single outcome $\omega$ of $\Omega$ .
<b>Corollary 12.6 The Sure Event</b> $\Omega$ : Is equal to the sample space as it contains all possible elementary events.
<b>Corollary 12.7 The Impossible Event</b> $\emptyset$ : The impossible event i.e. nothing is happening is denoted by the empty set.
<b>Definition 12.4 The Family of All Events</b> $\mathcal{A}/2^\Omega$ : 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 12.5 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>1.1. Sigma Algebras</b>
<b>Definition 12.6 Sigma Algebra</b> $\sigma$ : A set $\mathcal{F}$ of subsets of $\Omega$ is called a $\sigma$ -algebra on $\Omega$ if the following properties apply • $\Omega \in \mathcal{F}$ and $\emptyset \in \mathcal{F}$ • If $A \in \mathcal{F}$ then $\Omega \setminus A = A^C \in \mathcal{F}$ : The complementary subset of $A$ is also in $\Omega$ . • For all $A_i \in \mathcal{F} : \bigcup_{i=1} A_i \in \mathcal{F}$ See example 12.3.
<b>Corollary 12.8</b> $\mathcal{F}_{\min}$ : $\mathcal{F} = \{\emptyset, \Omega\}$ is the simplest $\sigma$ -algebra, telling us only if an event happened $\omega \in \Omega$ happened or not but not which one.
<b>Corollary 12.9</b> $\mathcal{F}_{\max}$ : $\mathcal{F} = 2^\Omega$ consists of all subsets of $\Omega$ and thus corresponds to full information i.e. we know if and which event happened.
<b>Definition 12.7 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 12.10</b> $\mathcal{F}$ -measurable Event: The elements $A_i \in \mathcal{F}$ are called <i>measurable sets</i> or $\mathcal{F}$ -measurable.
<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 12.8 Sigma Algebra generated by a subset of</b> $\Omega$ $\sigma(C)$ : Let $C$ be a class of subsets of $\Omega$ . The $\sigma$ -algebra generated by $C$ , denoted by $\sigma(C)$ , is the <i>smallest</i> sigma algebra $\mathcal{F}$ that included all elements of $C$ see example 12.4.
<b>Definition 12.9 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 12.5.
<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 12.11</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 at the end of the section.
<b>Definition 12.10 (<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: $\mathbb{P}(A) \in \{0, 1\} \quad \forall A \in \mathcal{F} \quad (12.1)$

<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>1.2. Measures</b>
<b>Definition 12.11 Measure</b> $\mu$ : A measure defined on a measurable space $\{\Omega, \mathcal{F}\}$ is a function/map: $\mu : \mathcal{F} \mapsto [0, \infty] \quad (12.2)$ for which holds: • $\mu(\emptyset) = 0$ • countable additivity [def. 12.12]
<b>Definition 12.12 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: $\mu\left(\bigcup_{i=1}^{\infty} F_i\right) = \sum_{i=1}^{\infty} \mu(F_i) \quad \text{for all} \quad F_j \cap F_k = \emptyset \quad \forall j \neq k \quad (12.3)$
<b>Corollary 12.12 Additive Function:</b> A function that satisfies countable additivity, is also additive, meaning that for every $F, G \in \mathcal{F}$ it holds: $\mu(F \cup G) = \mu(F) + \mu(G) \quad \iff \quad F \cap G = \emptyset \quad (12.4)$
<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 12.13 Equivalent Measures</b> $\mu \sim \nu$ : Let $\mu$ and $\nu$ be two measures defined on a measurable space[def. 12.7] $(\Omega, \mathcal{F})$ . The two measures are said to be equivalent if it holds that: $\mu(A) > 0 \iff \nu(A) > 0 \quad \forall A \subseteq \mathcal{F} \quad (12.5)$ this is equivalent to $\mu$ and $\nu$ having equivalent null sets: $\mathcal{N}_\mu = \mathcal{N}_\nu \quad \begin{matrix} \mathcal{N}_\mu = \{A \in \mathcal{A}   \mu(A) = 0\} \\ \mathcal{N}_\nu = \{A \in \mathcal{A}   \nu(A) = 0\} \end{matrix} \quad (12.6)$ see example 12.6
<b>Definition 12.14 Measure Space</b> $\{\mathcal{F}, \Omega, \mu\}$ : The triplet of sample space, sigma algebra and a measure is called a measure space.
<b>Definition 12.15 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: $\lambda([a, b]) = b - a \quad (12.7)$
<b>Corollary 12.13 Lebesgue Measure of Atomitics:</b> • The Lebesgue measure of a set containing only one point must be zero: $\lambda(\{a\}) = 0 \quad (12.8)$ • The Lebesgue measure of a set containing countably many points $A = \{a_1, a_2, \dots, a_n\}$ must be zero: $\lambda(A) + \sum_{i=1}^n \lambda(\{a_i\}) = 0 \quad (12.9)$ • The Lebesgue measure of a set containing uncountably many points $A = \{a_1, a_2, \dots\}$ can be either zero, positive and finite or infinite.
<b>1.3. 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 12.1 Non-negativity:</b> The probability of an event is a non-negative real number: If $A \in \mathcal{F}$ then $\mathbb{P}(A) \geq 0 \quad (12.10)$

<b>Axiom 12.2 Unitaity:</b> The probability that at least one of the elementary events in the entire sample space $\Omega$ will occur is equal to one: $\text{The certain event} \quad \mathbb{P}(\Omega) = 1 \quad (12.11)$
<b>Axiom 12.3 <math>\sigma</math>-additivity:</b> If $A_1, A_2, A_3, \dots \in \mathcal{F}$ are mutually disjoint, then: $\mathbb{P}\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} \mu(A_i) \quad (12.12)$
<b>Corollary 12.14</b> : As a consequence of this it follows: $\mathbb{P}(\emptyset) = 0 \quad (12.13)$
<b>Corollary 12.15 Complementary Probability:</b> $\mathbb{P}(A^C) = 1 - \mathbb{P}(A) \quad \text{with} \quad A^C = \Omega - A \quad (12.14)$

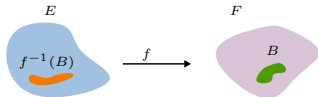
<b>Definition 12.16 Probability Measure</b> $\mathbb{P}$ : a probability measure is function $\mathbb{P} : \mathcal{F} \mapsto [0, 1]$ defined on a $\sigma$ -algebra $\mathcal{F}$ of a sample space $\Omega$ that satisfies the probability axioms.
---

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 $\iff$ <i>random variables</i> 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 $\Omega$ : $\text{Probability for an event in } \Omega$ $\mathbb{P}_X(E) = \mathbb{P}(\{\omega : X(\omega) \in E\}) = \mathbb{P}(X \in E) = \mathbb{P}\left(\overbrace{X^{-1}(E)}\right)$ Probability for an event in $E$
---

<b>Definition 12.17 <math>\mathcal{E}</math>-measurable function:</b> 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
---

$$\forall B \in \mathcal{F} : \quad f^{-1}(B) = \{\omega \in \mathcal{E} : f(\omega) \in B\} \in \mathcal{E} \quad (12.15)$$

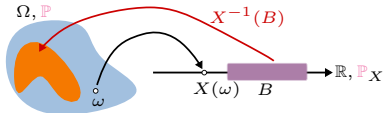


<b>Interpretation</b> The pre-image[def. 4.7] of $B$ under $f$ i.e. $f^{-1}(B)$ maps all values of the target space $F$ back to the sample space $\mathcal{E}$ (for all possible $B \in \mathcal{F}$ ).
--

**Definition 12.18 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  $(F, \mathcal{F})$ :

$$X : \Omega \mapsto \mathcal{F} \quad (\mathcal{F}^n) \quad (12.16)$$

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



**Corollary 12.16 :** Usually  $F = \mathbb{R}$ , which usually amounts to using the Borel  $\sigma$ -algebra  $\mathcal{B}$  of  $\mathbb{R}$ .

**Corollary 12.17 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 12.11):

$$\begin{aligned} X^{-1}(B) &= X^{-1}((-\infty, a]) \\ &= \{\omega \in \Omega : X(\omega) \leq a\} \in \mathcal{E} \quad \forall a \in \mathbb{R} \end{aligned} \quad (12.17)$$

**Definition 12.19 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 12.18 Indicator Functions**  $I_A(\omega)$ : An important class of measurable functions that can be used as r.v. are indicator functions:

$$I_A(\omega) = \begin{cases} 1 & \text{if } \omega \in A \\ 0 & \text{if } \omega \notin A \end{cases} \quad (12.18)$$

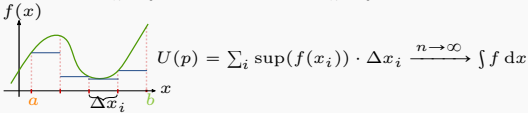
### 3. Lebesgue Integration

#### Problems of Riemann Integration

Verify and deepen knowledge about that

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

$$\lim_{n \rightarrow \infty} \int f(x) dx \stackrel{\text{str. u.c.}}{=} \int \lim_{n \rightarrow \infty} f(x) dx$$

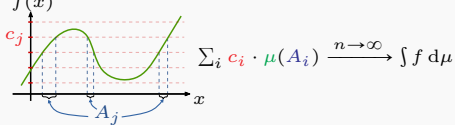


#### Idea

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  $A_j \Rightarrow$  we do not have to care anymore about discontinuities, as we can measure the size of our sets using our measure.



**Definition 12.20 Lebesgue Integral:**

$$\lim_{n \rightarrow \infty} \sum_{i=1}^n c_i \mu(A_i) = \int_{\Omega} f d\mu \quad \begin{matrix} f(x) \approx c_i \\ \forall x \in A_i \end{matrix} \quad (12.19)$$

**Definition 12.21**

**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:

$$X = \sum_{i=1}^n a_i \mathbb{1}_{A_i} \quad a_i \in \mathbb{R} \quad \mathcal{A} \ni A_i = \begin{cases} 1 & \text{if } \{X = a_i\} \\ 0 & \text{else} \end{cases} \quad (12.20)$$

**Definition 12.22 Expectation:**

$$\mathbb{E}[X] = \int_{\Omega} X(\omega) \mathbb{P}(d\omega) = \int_{\Omega} X d\mathbb{P} \quad (12.21)$$

**Corollary 12.19 Expectation of simple r.v.:**

If  $X$  is a simple r.v. its expectation is given by:

$$\mathbb{E}[X] = \sum_{i=1}^n a_i \mathbb{P}(A_i) \quad (12.22)$$

#### Proofs

*Proof.* corollary 12.11: Let  $\mathcal{C}$  denote all open intervals. Since every open set in  $\mathbb{R}$  is the countable union of open intervals, it holds that  $\sigma(\mathcal{C}) =$  the Borel  $\sigma$ -algebra of  $\mathbb{R}$ . Now let  $\mathcal{D}$  denote all intervals of the form  $(-\infty, a]$ , where  $a \in \mathbb{Q}$ . Let  $a, b \in \mathcal{C}$ , and let

See book

#### Examples

**Example 12.1 :**

- 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 12.2 :**

- 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 12.3 :** 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 12.4 :** 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 12.5 :**

- open half-lines:  $(-\infty, a)$  and  $(a, \infty)$ ,
- union of open half-lines:  $(a, b) = (-\infty, a) \cup (b, \infty)$ ,
- closed interval:  $[a, b] = (-\infty, a) \cup (b, \infty)$ ,
- 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 12.6 Equivalent (Probability) Measures:**

$$\begin{aligned} \Omega &= \{1, 2, 3\} & \mathbb{P}(\{1, 2, 3\}) &= \{2/3, 1/6, 1/6\} \\ & & \mathbb{P}(\{1, 2\}) &= \{1/3, 1/3, 1/3\} \end{aligned}$$

**Example 12.7 :**

add example fat book p.1286  
add example prob th book 4

### Combinatorics

#### 0.1. Permutations

**Definition 13.1 Permutation**  $n!$ : Given a set  $\mathcal{S}$  of  $n$  distinct objects, into how many distinct sequences/orders can we arrange/permutate those distinct objects

$$P(\mathcal{S}) = n! \iff P(\mathcal{S}) = n \cdot (n-1) \cdot (n-2) \cdot \dots \cdot 1 \quad (13.1)$$

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:

$$P(\mathcal{S}) = \frac{n!}{n_1! \cdot \dots \cdot n_k!} \quad \text{s.t.} \quad \sum_{i=1}^k n_i \leq n \quad (13.2)$$

#### 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 discrete random variable  $x$  is equal to some value  $\bar{x} \in \mathcal{X}$  is:

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

addispet

**Definition 14.1 Almost Surely (a.s.):** Let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a probability space. An event  $\omega \in \mathcal{F}$  happens almost surely iff  $\mathbb{P}(\omega) = 1 \iff \omega$  happens a.s. (14.1)

**Definition 14.2 Probability Mass Function (PMF):**

**Definition 14.3 Discrete Random Variable (DVR):** The set of possible values  $\bar{x}$  of  $\mathcal{X}$  is countable or finite.  
 $\mathcal{X} = \{0, 1, 2, 3, 4, \dots, 8\} \quad \mathcal{X} = \mathbb{N} \quad (14.2)$

**Definition 14.4 Probability Density Function (PDF):** Is real function  $f : \mathbb{R}^n \rightarrow [0, \infty)$  that satisfies:

$$\text{Non-negativity:} \quad f(x) \geq 0, \quad \forall x \in \mathbb{R}^n \quad (14.3)$$

$$\text{Normalization:} \quad \int_{-\infty}^{\infty} f(x) dx \stackrel{!}{=} 1 \quad (14.4)$$

$$\text{Must be integrable} \quad (14.5)$$

**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 14.5 Continuous Random Variable (CRV):** A real random variable (rrv)  $X$  is said to be (absolutely) continuous if there exists a pdf (def. 14.4)  $f_X$  s.t. for any subset  $B \subset \mathbb{R}$  it holds:

$$\mathbb{P}(X \in B) = \int_B f_X(x) dx \quad (14.6)$$

**Property 14.1 Zero Probability:** If  $X$  is a continuous rrv (def. 14.5), then:

$$\mathbb{P}(X = a) = 0 \quad \forall a \in \mathbb{R} \quad (14.7)$$

**Property 14.2 Open vs. Closed Intervals:** For any real numbers  $a$  and  $b$ , with  $a < b$  it holds:

$$\begin{aligned} \mathbb{P}(a \leq X \leq b) &= \mathbb{P}(a \leq X < b) = \mathbb{P}(a < X \leq b) \\ &= \mathbb{P}(a < X < b) \end{aligned} \quad (14.8)$$

$\iff$  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 14.1 :** 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 14.1:

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

$\square$

*Proof.* Property 14.2:

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

$\square$

**Definition 14.6 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:  
 $\text{supp}(\cdot) f_X := \{x \in \mathcal{X} | f(x) > 0\} \quad (14.9)$

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

**Theorem 14.1 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 14.2 Identically Distributed:** From theorem 14.1 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 Function

**Definition 14.7 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 14.3 :**

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

$$(14.10)$$

$$\text{Upper Limit} \quad \lim_{x \rightarrow \infty} F_X(x) = 1 \quad (14.11)$$

$$\text{Lower Limit} \quad \lim_{x \rightarrow -\infty} F_X(x) = 0 \quad (14.12)$$

**Definition 14.8 CDF of a discrete rv  $X$ :** Let  $X$  be discrete 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 14.9 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 14.1 Probability Interval:** Let  $X$  be a continuous rrv with pdf  $f_X$  and cumulative distribution function  $F_X$ , then it holds that:

$$\mathbb{P}(a \leq X \leq b) = F_X(b) - F_X(a) \quad (14.13)$$

*Proof.* [def. 14.9]:

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

$\square$

*Proof.* lemma 14.1:  
 $\mathbb{P}(a \leq X \leq b) = \mathbb{P}(X \leq b) - \mathbb{P}(X \leq a)$   
 or by the fundamental theorem of calculus (theorem 4.2):  
 $\mathbb{P}(a \leq X \leq b) = \int_a^b f_X(t) dt = \int_a^b \frac{\partial \mathbb{F}_X(t)}{\partial t} dt = [\mathbb{F}_X(t)]|_a^b$

**Theorem 14.2 A continuous rv is fully characterized by its CDF:** A cumulative distribution function completely determines the distribution of a continuous real-valued random variable.

**Theorem 14.3 (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 (def. 4.101) 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:  

$$p_Y(y) = \sum_{x_i \in \mathcal{X}_y} p_X(x_i) = \sum_{x \in \mathcal{Y}: g(x)=y} p_X(x) \quad (14.14)$$

*Proof.* theorem 14.3:  
 $Y = g(X) \iff \mathbb{P}(Y = y) = \mathbb{P}(x \in \mathcal{X}_y) = p_Y(y)$

**Theorem 14.4 (Scalar Continuous) Change of Variables:** Let  $X$  be a continuous rv  $X \in \mathcal{X}$  with pdf  $f_X$  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 (def. 4.101) function.  
 Then the pdf of  $Y$  is given by:  

$$f_Y(y) = f_X(x) \left| \frac{dx}{dy} \right| = f_X(x) \left| \frac{d}{dy} (g^{-1}(y)) \right| \quad (14.15)$$

$$= f_X(x) \frac{1}{\left| \frac{dy}{dx} \right|} = \frac{f_X(g^{-1}(y))}{\left| \frac{dg}{dx}(g^{-1}(y)) \right|} \quad (14.16)$$

**Theorem 14.5 (Continuous) Change of Variables:** Let  $X$  be a continuous rv  $X \in \mathcal{X}$  with pdf  $f_X$  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 (def. 4.101) function.  
 Then the pdf of  $Y$  is given by:  

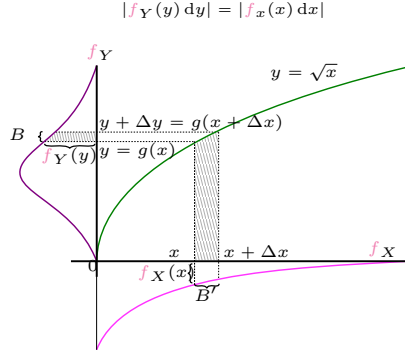
$$f_Y(\mathbf{y}) = f_X(\mathbf{x}) \left| \det \left( \frac{\partial g}{\partial \mathbf{x}} \right) \right|^{-1} \quad (14.17)$$

$$= f_X(g^{-1}(\mathbf{y})) \left| \det \left( \frac{\partial g}{\partial \mathbf{x}} \right) \right|^{-1}$$

**Where**  $\frac{\partial g}{\partial \mathbf{x}}$  is the Jacobian (def. 5.4).

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

*Proof.* theorem 14.4 (non-formal): The probability contained in a differential area must be invariant under a change of variables that is:



*Proof.* theorem 14.4 from CDF:  
 $\mathbb{P}(Y \leq y) = \mathbb{P}(g(X) \leq y) = \begin{cases} \mathbb{P}(X \leq g^{-1}(y)) & \text{if } g \text{ is increases.} \\ \mathbb{P}(X \geq g^{-1}(y)) & \text{if } g \text{ is decreases.} \end{cases}$

If  $g$  is monotonically increasing:  

$$f_Y(y) = f_X(g^{-1}(y))$$

$$f_Y(y) = \frac{d}{dy} \mathbb{F}_X(g^{-1}(y)) = f_X(x) \cdot \frac{d}{dy} g^{-1}(y)$$

If  $g$  is monotonically decreasing:  

$$f_Y(y) = 1 - \mathbb{F}_X(g^{-1}(y))$$

$$f_Y(y) = \frac{d}{dy} \mathbb{F}_X(g^{-1}(y)) = -f_X(x) \cdot \frac{d}{dy} g^{-1}(y)$$

*Proof.* theorem 14.4: 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 (??)

$$g(x + \Delta x) \stackrel{\text{eq. (4.40)}}{=} g(x) + \frac{dg}{dx} \Delta y \quad \text{for } \Delta x \rightarrow 0$$

$$= y + \Delta y \quad \text{with } \Delta y := \frac{dg}{dx} \cdot \Delta x \quad (14.18)$$

**Thus** for  $\mathbb{P}(Y \in B')$  it follows:  

$$\mathbb{P}(X \in B') = \int_y^{y+\Delta y} f_Y(t) dt \xrightarrow{\Delta y \rightarrow 0} |\Delta y \cdot f_Y(y)|$$

$$= \left| \frac{dg}{dx}(x) \Delta x \cdot f_Y(y) \right|$$

Now we simply need to related the surface of the two pdfs:  
 $B = [x, x + \Delta x] \xrightarrow{\text{same surfaces}} [y, y + \Delta y] = B'$   
 $\mathbb{P}(Y \in B) = \mathbb{P}(X \in B')$   

$$\xrightarrow{\Delta y \rightarrow 0} |f_Y(y) \cdot \Delta y| = \left| f_Y(y) \cdot \frac{dg}{dx}(x) \Delta x \right| = |f_X(x) \cdot \Delta x|$$

$$f_Y(y) \left| \frac{dg}{dx}(x) \right| |\Delta x| = f_X(x) \cdot |\Delta x|$$

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

## Rules of Probability

**Definition 14.10 Marginalization/Sum Rule:**  
**Given:**  $p_{x,y}(\bar{x}, \bar{y}) \quad p_x(\bar{x}) := \sum_{\bar{y} \in \mathcal{Y}} p_{x,y}(\bar{x}, \bar{y}) \quad (14.19)$

**Definition 14.11 Conditioning:**  
**Given:**  $p_{xy} \quad p_{xy}(x|y = \bar{y}) := \frac{p_{xy}(x, y = \bar{y})}{p_y(y = \bar{y})}$   
**if**  $p_y(\bar{y}) \neq 0 \quad (14.20)$

**Definition 14.12 Product Rule:** follows directly from eq. (14.20)  

$$p(x, y) = p(y|x)p_x(x) = p(x|y)p(y) \quad (14.21)$$

**Theorem 14.6 Total Probability Theorem:** **Given:**  $p_{x,y}(\bar{x}, \bar{y})$  with eq. (14.19) and eq. (14.21) it follows:  

$$p_x(\bar{x}) \stackrel{\text{eq. (14.19)}}{=} \sum_{\bar{y} \in \mathcal{Y}} p_{x,y}(\bar{x}, \bar{y})$$

$$\stackrel{\text{eq. (14.21)}}{=} \sum_{y \in \mathcal{Y}} p_{x|y}(\bar{x}|\bar{y}) p_y(\bar{y}) \quad (14.22)$$

**Definition 14.13 Independence:** Two random variables  $x$  and  $y$  are said to be **independent** if:  

$$p(x|y) = p(x) \stackrel{\text{eq. (14.20)}}{\iff} p(x, y) = p(x)p(y) \quad (14.23)$$

**Corollary 14.3** eq. (14.23):  

$$p(x|y) = p(x) \xrightarrow{\text{implies}} p(y|x) = p(y) \quad (14.24)$$

**old mutual independence**

## 0.2. Conditional PDF

**Let**  $x, y, z$  be R.V. (which themselves may be collections of random variables)

**Definition 14.14 Marginalization:**  

$$p_{x|z}(\bar{x}|\bar{z}) = \sum_{\bar{y} \in \mathcal{Y}} p_{xy|z}(\bar{x}, \bar{y}|\bar{z}) \quad (14.25)$$

**Definition 14.15 Conditioning:**  

$$p_{x|y|z}(\bar{x}|\bar{y}, \bar{z}) = \frac{p_{xy|z}(\bar{x}, \bar{y}|\bar{z})}{p_{y|z}(\bar{y}|\bar{z})} \quad (14.26)$$

**Definition 14.16 Product Rule:** follows directly from eq. (14.26)  

$$p_{xyz}(\bar{x}, \bar{y}|\bar{z}) = p_{x|y|z}(\bar{x}|\bar{y}, \bar{z}) p_{y|z}(\bar{y}|\bar{z}) \quad (14.27)$$

**Note**  
 $z$  basically parameterizes the pdf.

**Definition 14.17 Conditional Independence:** Two random variables  $x$  and  $y$  are said to be conditionally independent on  $z$  if  

$$p(x|y, z) = p(x|z) \stackrel{\text{eq. (14.26)}}{=} p(x, y|z) = p(x|z)p(y|z) \quad (14.28)$$
 Hence, knowledge of  $z$  makes  $x$  and  $y$  independent.

**Note**  
 Conditional independence does not imply  $p(x, y) = p(x)p(y)$

**Rule 14.1 (Bayes' Rule).** **Given:** the prior  $p(X)$  and the likelihood  $p(Y|X)$ , we can find the posterior by:  

$$p(X|Y) = \frac{p(Y, X)}{p(Y)} = \frac{p(X)p(Y|X)}{p(Y)}$$

$$\stackrel{\text{normalization}}{=} \frac{p(X)p(Y|X)}{\sum_{X=x} p(X=x)p(Y|X=x)}$$

$$\text{Posterior} = \frac{\text{Prior} \cdot \text{Likelihood}}{\text{Normalization}}$$

*Proof.* Equation (14.25)  

$$p_{x|z}(\bar{x}|\bar{z}) \stackrel{\text{eq. (14.20)}}{=} \sum_{\bar{y} \in \mathcal{Y}} p_{xy|z}(\bar{x}, \bar{y}|\bar{z}) \stackrel{\text{eq. (14.19)}}{=} \sum_{\bar{y} \in \mathcal{Y}} \frac{p_{xy|z}(\bar{x}, \bar{y}|\bar{z})}{p_{y|z}(\bar{y}|\bar{z})} p_{y|z}(\bar{y}|\bar{z})$$

$$\stackrel{\text{eq. (14.21)}}{=} \sum_{\bar{y} \in \mathcal{Y}} \frac{p_{xy|z}(\bar{x}, \bar{y}|\bar{z}) p_{y|z}(\bar{y}|\bar{z})}{p_{y|z}(\bar{y}|\bar{z})} \stackrel{\text{cancel}}{=} \sum_{\bar{y} \in \mathcal{Y}} p_{xy|z}(\bar{x}, \bar{y}|\bar{z})$$

*Proof.* Equation (14.26)  

$$p_{x|y|z}(\bar{x}|\bar{y}, \bar{z}) \stackrel{\text{eq. (14.20)}}{=} \frac{p_{xyz}(\bar{x}, \bar{y}, \bar{z})}{p_{yz}(\bar{y}, \bar{z})}$$

$$\stackrel{\text{eq. (14.21)}}{=} \frac{p_{x|y|z}(\bar{x}, \bar{y}|\bar{z}) p_{y|z}(\bar{y}|\bar{z})}{p_{y|z}(\bar{y}|\bar{z})} \stackrel{\text{cancel}}{=} p_{x|y|z}(\bar{x}, \bar{y}|\bar{z})$$

*Proof.* Equation (14.24)  

$$p(y|x) \stackrel{\text{eq. (14.20)}}{=} \frac{p(x, y)}{p(x)} \stackrel{p(x, y) = p(x)p(y)}{=} \frac{p(x)p(y)}{p(x)} = p(y)$$

*Proof.* Equation (14.28)  

$$p_{x|y|z}(\bar{x}|\bar{y}, \bar{z}) = \frac{p(x, y, z)}{p(y, z)} = \frac{p(x, y|z)p(z)}{p(y, z)} \stackrel{\text{cancel}}{=} \frac{p(x, y|z)}{p(y|z)}$$

$$\Rightarrow p(x, y|z) = p(x|z)p(y|z)$$

## 1. Key figures

### 1.1. The Expectation

**Definition 14.18 Expectation (disc. case):**  

$$\mu_X := \mathbb{E}[x] := \sum_{\bar{x} \in \mathcal{X}} \bar{x} p_x(\bar{x}) \quad (14.29)$$

**Definition 14.19 Expectation (cont. case):**  

$$\mathbb{E}_x[x] := \int_{\bar{x} \in \mathcal{X}} \bar{x} f_x(\bar{x}) d\bar{x} \quad (14.30)$$

**Law 14.1 Expectation of independent variables:**  

$$\mathbb{E}[XY] = \mathbb{E}[X] \mathbb{E}[Y] \quad (14.31)$$

**Property 14.4 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:  

$$\mathbb{E}[a + bX + cY] = a + b\mathbb{E}[X] + c\mathbb{E}[Y] \quad (14.32)$$

**Thus**  $\mathbb{E}$  is a **linear** operator (def. 4.111).

**Note: Expectation of the expectation**  
 The expectation of a r.v.  $X$  is a constant hence with Property 14.9 it follows:  

$$\mathbb{E}[\mathbb{E}[X]] = \mathbb{E}[X] \quad (14.33)$$

**Property 14.5 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:  

$$\mathbb{E}[AXB] = A\mathbb{E}[(XB)] = A\mathbb{E}[X]B \quad (14.34)$$

*Proof.* eq. (14.51):  

$$\mathbb{E}[XY] = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{X,Y}(x, y) xy$$

$$\stackrel{[\text{def. 14.13}]}{=} \sum_{x \in \mathcal{X}} p_X(x) x \sum_{y \in \mathcal{Y}} p_Y(y) y = \mathbb{E}[X] \mathbb{E}[Y]$$



**Law 14.2 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:  

$$\mathbb{E}_Y[y] = \sum_{y \in \mathcal{Y}} y p_Y(y) = \sum_{x \in \mathcal{X}} g(x) p_X(x) \quad \text{or integral for CRV} \quad (14.35)$$

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

**Theorem 14.7 Jensen's Inequality:** If  $X$  is a random variable and  $f$  is a convex function, then it holds that:  

$$f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)] \quad (14.36)$$
on the contrary if  $f$  is a concave function it follows:  

$$f(\mathbb{E}[X]) \geq \mathbb{E}[f(X)] \quad (14.37)$$

## 1.2. The Variance

**Definition 14.20 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 executed value  $\mathbb{E}[X]$  and is defined by:

$$\mathbb{V}[X] := \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2] - \mathbb{E}[X]^2 \quad (14.38)$$

*Proof.* eq. (14.58)

$$\begin{aligned} \mathbb{V}[X] &= \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2 - 2X\mathbb{E}[X] + \mathbb{E}[X]^2] \\ &\stackrel{\text{Property 14.9}}{=} \mathbb{E}[X^2] - 2\mathbb{E}[X]\mathbb{E}[X] + \mathbb{E}[X]^2 = \mathbb{E}[X^2] - \mu^2 \end{aligned}$$

**Property 14.6 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:  

$$\mathbb{V}[a] = 0 \quad \text{with} \quad a \in \mathbb{R} \quad (14.39)$$

**Property 14.7 Affine Transformation:** If  $\mathbf{X} \in \mathbb{R}^n$  is a random vector,  $\mathbf{A} \in \mathbb{R}^{m \times n}$  a constant matrix and  $b \in \mathbb{R}^m$  then it holds:  

$$\mathbb{V}[\mathbf{A}\mathbf{X} + b] = \mathbf{A}\mathbb{V}[\mathbf{X}]\mathbf{A}^T \quad (14.40)$$

*Proof.* Property 14.13

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

**Definition 14.21 Covariance:** The Covariance is a measure of how much two or more random variables vary linearly with each other.

$$\begin{aligned} \text{Cov}[X, Y] &= \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] \\ &= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] \end{aligned} \quad (14.41)$$

*Proof.* eq. (14.62)

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

**Definition 14.22 Covariance Matrix:** The variance of a  $k$ -dimensional random vector  $\mathbf{X} = (X_1 \dots X_k)$  is given by the 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:

$$\begin{aligned} \mathbb{V}[\mathbf{X}] &:= \Sigma(\mathbf{X}) := \text{Cov}[\mathbf{X}, \mathbf{X}] := \\ &= \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])^T] \\ &= \mathbb{E}[\mathbf{X}\mathbf{X}^T] - \mathbb{E}[\mathbf{X}]\mathbb{E}[\mathbf{X}]^T \in [-\infty, \infty] \end{aligned} \quad (14.42)$$

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

**Note: Covariance and Variance**

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

$$\text{Cov}[X, X] = \mathbb{V}[X] \equiv \sigma^2(X) \equiv \sigma_X^2 \quad (14.43)$$

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

**Property 14.8 Translation and Scaling:**

$$\text{Cov}(a + bX, c + dY) = bd\text{Cov}(X, Y) \quad (14.44)$$

**Definition 14.23 Correlation Coefficient:** Is the standardized version of the covariance:

$$\begin{aligned} \text{Corr}[X] &:= \frac{\text{Cov}[X]}{\sigma_{X_1} \dots \sigma_{X_k}} \in [-1, 1] \\ &= \begin{cases} +1 & \text{if } Y = aX + b \text{ with } a > 0, b \in \mathbb{R} \\ -1 & \text{if } Y = aX + b \text{ with } a < 0, b \in \mathbb{R} \end{cases} \end{aligned} \quad (14.45)$$

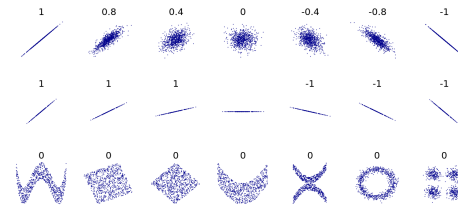


Figure 1: Several sets of  $(x, y)$  points, with their correlation coefficient

**Law 14.3 Translation and Scaling:**

$$\text{Corr}(a + bX, c + dY) = \text{sign}(b)\text{sign}(d)\text{Cov}(X, Y) \quad (14.46)$$

**Note**

- The correlation/covariance reflects the noisiness and direction of a linear relationship (top row fig. 2), **but** not the slope of that relationship (middle row fig. 2) nor many aspects of nonlinear relationships (bottom row)
- The set in the center of fig. 2 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 14.7.
- Correlation is dimensionless, whereas the unit of the covariance is obtained by the product of the units of the two RV variables.

**Law 14.4 Covariance of independent RVs:** The covariance/correlation of two independent variable's (def. 14.13) is zero:

$$\begin{aligned} \text{Cov}[X, Y] &= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] \\ &\stackrel{\text{eq. (14.51)}}{=} \mathbb{E}[X]\mathbb{E}[Y] - \mathbb{E}[X]\mathbb{E}[Y] = 0 \end{aligned}$$

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

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

**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:  

$$\begin{aligned} \text{Cov}(X, Y) &= \text{Cov}(X, X^2) = \mathbb{E}[X \cdot X^2] - \mathbb{E}[X]\mathbb{E}[X^2] \\ &= \mathbb{E}[X^3] - \mathbb{E}[X]\mathbb{E}[X^2] \stackrel{\text{eq. (14.88)}}{=} 0 - 0 \cdot \mathbb{E}[X^2] \\ &\stackrel{\text{eq. (14.78)}}{=} 0 \end{aligned}$$
 $\Rightarrow$  the relationship between  $Y$  and  $X$  must be non-linear.

**Definition 14.24 Autocorrelation/Crosscorrelation  $\gamma(t_1, t_2)$ :** Describes the covariance (def. 14.28) between the two values of a stochastic process  $(\mathbf{X}_t)_{t \in T}$  at different time points  $t_1$  and  $t_2$ .

$$\gamma(t_1, t_2) = \text{Cov}[\mathbf{X}_{t_1}, \mathbf{X}_{t_2}] = \mathbb{E}[(\mathbf{X}_{t_1} - \mu_{t_1})(\mathbf{X}_{t_2} - \mu_{t_2})] \quad (14.47)$$

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

$$\gamma(t, t) = \text{Cov}[\mathbf{X}_t, \mathbf{X}_t] \stackrel{\text{eq. (14.64)}}{=} \mathbb{V}[\mathbf{X}_t] \quad (14.48)$$

**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  $\mathbf{x}(t)$  the autocorrelation function  $\gamma(t, t - \tau)$  describes how *similar* the time translated function  $\mathbf{x}(t - \tau)$  and the original function  $\mathbf{x}(t)$  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 14.25 Expectation (disc. case):**

$$\mu_X := \mathbb{E}_X[x] := \sum_{\bar{x} \in \mathcal{X}} \bar{x} p_X(\bar{x}) \quad (14.49)$$

**Definition 14.26 Expectation (cont. case):**

$$\mathbb{E}_X[x] := \int_{\bar{x} \in \mathcal{X}} \bar{x} f_X(\bar{x}) d\bar{x} \quad (14.50)$$

**Law 14.5 Expectation of independent variables:**

$$\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y] \quad (14.51)$$

**Property 14.9 Translation and scaling:** If  $\mathbf{X} \in \mathbb{R}^n$  and  $\mathbf{Y} \in \mathbb{R}^n$  are random vectors, and  $a, b, c \in \mathbb{R}^n$  are constants then it holds:

$$\mathbb{E}[a + b\mathbf{X} + c\mathbf{Y}] = a + b\mathbb{E}[\mathbf{X}] + c\mathbb{E}[\mathbf{Y}] \quad (14.52)$$

**Thus  $\mathbb{E}$  is a linear operator** (def. 4.11).

**Note: Expectation of the expectation**

The expectation of a r.v.  $X$  is a constant hence with Property 14.9 it follows:

$$\mathbb{E}[\mathbb{E}[X]] = \mathbb{E}[X] \quad (14.53)$$

**Property 14.10 Matrix×Expectation:** If  $\mathbf{X} \in \mathbb{R}^n$  is a random vector and  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{B} \in \mathbb{R}^{n \times m}$  are constant matrices then it holds:

$$\mathbb{E}[\mathbf{A}\mathbf{X}\mathbf{B}] = \mathbf{A}\mathbb{E}[(\mathbf{X}\mathbf{B})] = \mathbf{A}\mathbb{E}[\mathbf{X}]\mathbf{B} \quad (14.54)$$

*Proof.* eq. (14.51):

$$\begin{aligned} \mathbb{E}[XY] &= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{X,Y}(x, y) xy \\ &\stackrel{[\text{def. 14.13}]}{=} \sum_{x \in \mathcal{X}} p_X(x) x \sum_{y \in \mathcal{Y}} p_Y(y) y = \mathbb{E}[X]\mathbb{E}[Y] \end{aligned}$$

**Law 14.6 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:

$$\mathbb{E}_Y[y] = \sum_{y \in \mathcal{Y}} y p_Y(y) = \sum_{x \in \mathcal{X}} g(x) p_X(x) \quad \text{or integral for CRV} \quad (14.55)$$

**Consequence**

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

**Theorem 14.8 Jensen's Inequality:** If  $X$  is a random variable and  $f$  is a convex function, then it holds that:

$$f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)] \quad (14.56)$$

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

$$f(\mathbb{E}[X]) \geq \mathbb{E}[f(X)] \quad (14.57)$$

### 2.2. The Variance

**Definition 14.27 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 executed value  $\mathbb{E}[X]$  and is defined by:

$$\mathbb{V}[X] := \mathbb{E}[(X - \mathbb{E}[X])^2] \stackrel{\text{see section 3}}{=} \mathbb{E}[X^2] - \mathbb{E}[X]^2 \quad (14.58)$$

#### 2.2.1. Properties

**Property 14.11 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:

$$\mathbb{V}[a] = 0 \quad \text{with} \quad a \in \mathbb{R} \quad (14.59)$$

see shift and scaling for proof section 3

**Property 14.12 Shifting and Scaling:**  
 $\mathbb{V}[a + bX] = a^2 \sigma^2$  with  $a \in \mathbb{R}$  (14.60)  
 see section 3

**Property 14.13 Affine Transformation:** 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:  
 $\mathbb{V}[AX + b] = A\mathbb{V}[X]A^T$  (14.61)

*Proof.* Property 14.13  

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

**Definition 14.28 Covariance:** The Covariance is a measure of how much two or more random variables vary **linearly** with each other.  

$$\begin{aligned} \text{Cov}[X, Y] &= \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] \\ &= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] \end{aligned} \quad (14.62)$$

*Proof.* eq. (14.62)  

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

**Definition 14.29 Covariance Matrix:** The variance of a  $k$ -dimensional random vector  $X = (X_1 \dots X_k)$  is given by the 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:

$$\begin{aligned} \mathbb{V}[X] &:= \Sigma(X) := \text{Cov}[X, X] := \\ &= \mathbb{E}[(X - \mathbb{E}[X])(X - \mathbb{E}[X])^T] \\ &= \mathbb{E}[XX^T] - \mathbb{E}[X]\mathbb{E}[X]^T \in [-\infty, \infty] \end{aligned} \quad (14.63)$$

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

**Note: Covariance and Variance**  
 The variance is a special case of the covariance in which two variables are identical:  
 $\text{Cov}[X, X] = \mathbb{V}[X] \equiv \sigma^2(X) \equiv \sigma_X^2$  (14.64)

add <http://www.visiondummy.com/2014/04/geometric-interpretation-covariance-matrix/>

**Property 14.14 Translation and Scaling:**  
 $\text{Cov}(a + bX, c + dY) = bd\text{Cov}(X, Y)$  (14.65)

**Definition 14.30 Correlation Coefficient:** Is the standardized version of the covariance:  

$$\begin{aligned} \text{Corr}[X] &:= \frac{\text{Cov}[X]}{\sigma_{X_1} \dots \sigma_{X_k}} \in [-1, 1] \\ &= \begin{cases} +1 & \text{if } Y = aX + b \text{ with } a > 0, b \in \mathbb{R} \\ -1 & \text{if } Y = aX + b \text{ with } a < 0, b \in \mathbb{R} \end{cases} \end{aligned} \quad (14.66)$$

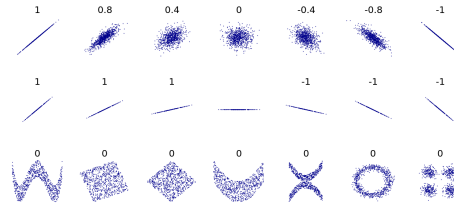


Figure 2: Several sets of  $(x, y)$  points, with their correlation coefficient

**Law 14.7 Translation and Scaling:**  
 $\text{Corr}(a + bX, c + dY) = \text{sign}(b)\text{sign}(d)\text{Cov}(X, Y)$  (14.67)

- Note**
- The correlation/covariance reflects the noisiness and direction of a linear relationship (top row fig. 2), **but** not the slope of that relationship (middle row fig. 2) nor many aspects of nonlinear relationships (bottom row)
  - The set in the center of fig. 2 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 14.7.
  - Correlation is dimensionless, whereas the unit of the covariance is obtained by the product of the units of the two RV variables.

**Law 14.8 Covariance of independent RVs:** The covariance/correlation of two independent variable's (def. 14.13) is zero:  

$$\begin{aligned} \text{Cov}[X, Y] &= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] \\ &\stackrel{\text{eq. (14.51)}}{=} \mathbb{E}[X]\mathbb{E}[Y] - \mathbb{E}[X]\mathbb{E}[Y] = 0 \end{aligned}$$

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

$\text{Cov}(X, Y) = \text{Corr}(X, Y) = 0 \nRightarrow \mathbb{P}_{X,Y}(x, y) = \mathbb{P}_X(x)\mathbb{P}_Y(y)$   
**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:  

$$\begin{aligned} \text{Cov}(X, Y) &= \text{Cov}(X, X^2) = \mathbb{E}[X \cdot X^2] - \mathbb{E}[X]\mathbb{E}[X^2] \\ &= \mathbb{E}[X^3] - \mathbb{E}[X]\mathbb{E}[X^2] \stackrel{\text{eq. (14.88)}}{=} 0 - 0 \cdot \mathbb{E}[X^2] \\ &\stackrel{\text{eq. (14.78)}}{=} 0 \end{aligned}$$
 $\Rightarrow$  the relationship between Y and X must be non-linear.

**Definition 14.31 Quantile:** Are specific values  $q_\alpha$  in the range<sup>[def. 4.6]</sup> of a random variable X that are defined as the value for which the cumulative probability is less than  $\alpha \in (0, 1)$ :  

$$q_\alpha : \mathbb{P}(X \leq x) = \mathbb{F}_X(q_\alpha) = \alpha \xrightarrow{\mathbb{F} \text{ invert.}} q_\alpha = \mathbb{F}_X^{-1}(\alpha) \quad (14.68)$$

add figure

### 3. Proofs

*Proof.* eq. (14.58)  

$$\begin{aligned} \mathbb{V}[X] &= \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2 - 2X\mathbb{E}[X] + \mathbb{E}[X]^2] \\ &\stackrel{\text{Property 14.9}}{=} \mathbb{E}[X^2] - 2\mathbb{E}[X]\mathbb{E}[X] + \mathbb{E}[X]^2 = \mathbb{E}[X^2] - \mu^2 \end{aligned}$$
 □

*Proof.* Property 14.12  

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

## Discrete Distributions

### 4.1. Bernoulli Distribution

Bern(p)

**Definition 14.32 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 14.33 Bernoulli distribution**  $X \sim \text{Bern}(p)$ :  $X$  is a binary variable i.e. can only attain the values 0 (failure) or 1 (success) with a parameter  $p$  that signifies the success probability:

$$p(x; p) = \begin{cases} p & \text{for } x = 1 \\ 1 - p & \text{for } x = 0 \end{cases} \iff \begin{cases} \mathbb{P}(X = 1) = p \\ \mathbb{P}(X = 0) = 1 - p \end{cases}$$

$$= p^x \cdot (1 - p)^{1-x} \quad \text{for } x \in \{0, 1\}$$

$$\mathbb{E}[X] = p \quad (14.69) \quad \mathbb{V}[X] = p(1 - p) \quad (14.70)$$

### 4.2. Binomial Distribution

B(n, p)

**Definition 14.34 Binomial Distribution:** Models the probability of exactly  $X$  success given a fixed number  $n$ -Bernoulli experiments<sup>(def. 14.32)</sup>, where the probability of success of a single experiment is given by  $p$ :

$$p(x) = \binom{n}{x} p^x (1 - p)^{n-x} \quad \begin{array}{l} n: \text{nb. of repetitions} \\ x: \text{nb. of successes} \\ p: \text{probability of success} \end{array}$$

$$\mathbb{E}[X] = np \quad (14.71) \quad \mathbb{V}[X] = np(1 - p) \quad (14.72)$$

#### 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:

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

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 14.35 Geometric Distribution**  $\text{Geom}(p)$ : Models the probability of the number  $X$  of Bernoulli trials<sup>(def. 14.32)</sup> until the first success

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

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

$$\mathbb{E}[X] = \frac{1}{p} \quad (14.73) \quad \mathbb{V}[X] = \frac{1 - p}{p^2} \quad (14.74)$$

#### 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  $p(x)$ :

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

- $\log(1 - p) \approx -p$  for small  $p$

### 4.4. Poisson Distribution

Pois( $\lambda$ )

**Definition 14.36 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.

$$p(x) = e^{-\lambda} \frac{\lambda^x}{x!} \quad \begin{array}{l} \lambda > 0 \\ x \in \mathbb{Z}_{\geq 0} \end{array} \quad (14.75)$$

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

$$\mathbb{E}[X] = \lambda \quad (14.76) \quad \mathbb{V}[X] = \lambda \quad (14.77)$$

## Continuous Distributions

### 5.1. Uniform Distribution

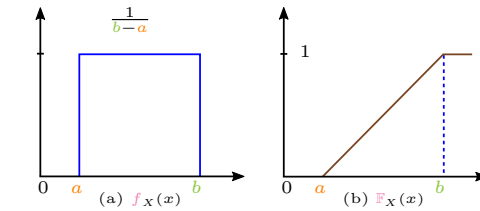
$\mathcal{U}(a, b)$

**Definition 14.37 Uniform Distribution**  $\mathcal{U}(a, b)$ : Is probability distribution, where all intervals of the same length on the distribution's support<sup>(def. 14.6)</sup>  $\text{supp}(\mathcal{U}[a, b]) = [a, b]$  are equally probable/likely.

$$f(x) = \frac{1}{b - a} \mathbf{1}_{x \in [a; b]} = \begin{cases} \frac{1}{b - a} = \text{const} & \text{if } a \leq x \leq b \\ 0 & \text{else} \end{cases} \quad (14.78)$$

$$\mathbb{F}(x) = \begin{cases} 0 & x < a \\ \frac{x - a}{b - a} & \text{if } a \leq x \leq b \\ 1 & x > b \end{cases} \quad (14.79)$$

$$\mathbb{E}[X] = \frac{a + b}{2} \quad \mathbb{V}(X) = \frac{(b - a)^2}{12} \quad (14.80)$$



### 5.2. Exponential Distribution

$\exp(\lambda)$

**Definition 14.38 Exponential Distribution**  $X \sim \exp(\lambda)$ : Is the continuous analogue to the geometric distribution<sup>(def. 14.35)</sup>.

It describes the probability  $f(x; \lambda)$  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  $x$ .

$$f(x; \lambda) = \begin{cases} \lambda e^{-\lambda x} & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases} \quad (14.81)$$

$$\mathbb{F}(x; \lambda) = \begin{cases} 1 - e^{-\lambda x} & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases} \quad (14.82)$$

$$\mathbb{E}[X] = \frac{1}{\lambda} \quad \mathbb{V}(X) = \frac{1}{\lambda^2} \quad (14.83)$$

### 5.3. Laplace Distribution

**Definition 14.39 Laplace Distribution:**

$$\text{Laplace Distribution} \quad f(x; \mu, \sigma) = \frac{1}{2\sigma} \exp\left(-\frac{|x - \mu|}{\sigma}\right) \quad (14.84)$$

### 5.4. The Normal Distribution

$\mathcal{N}(\mu, \sigma)$

**Definition 14.40 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:

$$\mathbb{E}[X] = \mu \quad \mathbb{V}(X) = \sigma^2 \quad (14.85)$$

$$f(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2}\left(\frac{x - \mu}{\sigma}\right)^2\right\} \quad (14.86)$$

$$\mathbb{F}(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^x \exp\left\{-\frac{1}{2}\left(\frac{u - \mu}{\sigma}\right)^2\right\} du \quad (14.87)$$

$$x \in \mathbb{R} \quad \text{or} \quad -\infty < x < \infty$$

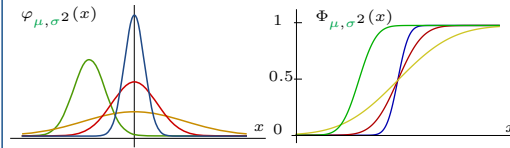


Figure 4:

$$\begin{array}{llll} \mu = 0 & \mu = 0 & \mu = 0 & \mu = -2 \\ \sigma^2 = 0.2 & \sigma^2 = 1.0 & \sigma^2 = 5.0 & \sigma^2 = 0.5 \end{array}$$

$$\text{Property 14.15 : } \mathbb{P}_X(\mu - \sigma \leq x \leq \mu + \sigma) = 0.66$$

$$\text{Property 14.16 : } \mathbb{P}_X(\mu - 2\sigma \leq x \leq \mu + 2\sigma) = 0.95$$

### 5.5. The Standard Normal distribution

$\mathcal{N}(0, 1)$

**Historic Problem:** the cumulative distribution eq. (14.87) 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 14.41**

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

$$\mathbb{E}[X] = 0 \quad \mathbb{V}(X) = 1 \quad (14.88)$$

$$f(x; 0, 1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} \quad (14.89)$$

$$\mathbb{F}(x; 0, 1) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{1}{2}u^2} du \quad (14.90)$$

$$x \in \mathbb{R} \quad \text{or} \quad -\infty < x < \infty$$

#### Corollary 14.4

**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. 14.31)</sup> is then denoted by:

$$z_\alpha = \Phi^{-1}(\alpha) \quad \alpha \in (0, 1) \quad (14.91)$$

#### 5.5.1. Calculating Probabilities

**Property 14.17 Symmetry:** Let  $z > 0$

$$\mathbb{P}(Z \leq z) = \Phi(z) \quad (14.92)$$

$$\mathbb{P}(Z \leq -z) = \Phi(-z) = 1 - \Phi(z) \quad (14.93)$$

$$\mathbb{P}(-a \leq Z \leq b) = \Phi(b) - \Phi(-a) = \Phi(b) - (1 - \Phi(a))$$

$$\stackrel{a=b=z}{=} 2\Phi(z) - 1 \quad (14.94)$$

#### 5.5.2. Linear Transformations of Normal Dist.

**Proposition 14.1 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:

$$Y \sim \mathcal{N}\left(a + b\mu, b^2\sigma^2\right) \iff f_Y(y) = \frac{1}{|b|} f_X\left(\frac{y - a}{b}\right) \quad (14.95)$$

section 2

**Proposition 14.2 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:

$$X \sim \mathcal{N}(\mu, \sigma^2) \xrightarrow{\frac{X - \mu}{\sigma}} Z \sim \mathcal{N}(0, 1) \quad (14.96)$$

section 2

#### 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 14.3 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 normal variable  $Z \sim \mathcal{N}(0, 1)$  is related to  $F_X(X)$  by:

$$F_X(x) = \Phi\left(\frac{x - \mu}{\sigma}\right) \quad (14.97)$$

section 2

## 6. The Multivariate Normal distribution

### Definition 14.42

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

The  $k$ -multivariate Normal distribution of:

$X = (x_1 \dots x_k)^\top$  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:

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

is given by:

$$f_X(x_1, \dots, x_k) = \frac{1}{\sqrt{(2\pi)^k \det(\Sigma)}} \exp\left(-\frac{1}{2}(X - \mu)^\top \Sigma^{-1}(X - \mu)\right) \quad (14.98)$$

Normalisation

### Definition 14.43 Jointly Gaussian Random Variables:

Two random variables  $x, y$  both scalars or vectors, are said to be **jointly Gaussian** if the joint vector random variable  $\begin{bmatrix} x & y \end{bmatrix}^\top$  is again a GRV.

**Corollary 14.5 Jointly GRV of GRVs:** If  $x$  and  $y$  are both independent GRVs  $x \sim \mathcal{N}(\mu_x, \Sigma_x)$ ,  $y \sim \mathcal{N}(\mu_y, \Sigma_y)$ , then they are jointly Gaussian (def. 14.43).

$$p(x, y) = p(x)p(y) \quad (14.99)$$

$$\propto \exp\left(-\frac{1}{2}\left\{\begin{bmatrix} x - \mu_x \\ y - \mu_y \end{bmatrix}^\top \begin{bmatrix} \Sigma_x^{-1} & 0 \\ 0 & \Sigma_y^{-1} \end{bmatrix} \begin{bmatrix} x - \mu_x \\ y - \mu_y \end{bmatrix}\right\}\right)$$

$$= \exp\left(-\frac{1}{2}\begin{bmatrix} x - \mu_x & y - \mu_y \end{bmatrix}^\top \begin{bmatrix} 0 & \Sigma_x^{-1} \\ \Sigma_y^{-1} & 0 \end{bmatrix} \begin{bmatrix} x - \mu_x \\ y - \mu_y \end{bmatrix}\right)$$

**Property 14.18 Scalar Affine Transformation of GRVs:** Let  $y \in \mathbb{R}^n$  be GRV,  $a \in \mathbb{R}_+, b \in \mathbb{R}$  and let  $x$  be defined by the **affine transformation** (def. 8.1):

$$x = ay + b \quad a \in \mathbb{R}_+, b \in \mathbb{R}^d$$

Then  $x$  is a GRV with:

$$x \sim \mathcal{N}(\mu_x, \sigma_x^2) = \mathcal{N}(a\mu + b, a^2\sigma^2) \quad (14.100)$$

**Property 14.19 Affine Transformation of GRVs:** Let  $y \in \mathbb{R}^n$  be GRV,  $A \in \mathbb{R}^{d \times n}, b \in \mathbb{R}^d$  and let  $x$  be defined by the **affine transformation** (def. 8.1):

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

Then  $x$  is a GRV (see Section 2).

**Property 14.20 Linear Combination of jointly GRVs:** Let  $x \in \mathbb{R}^n, y \in \mathbb{R}^m$  two jointly GRVs, and let  $z$  be defined as:

$$z = Ax + Ay \quad A_x \in \mathbb{R}^{d \times n}, A_y \in \mathbb{R}^{d \times m}$$

Then  $z$  is GRV (see Section 2).

### 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:

$$\Sigma = \begin{bmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_k^2 \end{bmatrix} \quad \text{and} \quad \mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_k \end{bmatrix} \quad (14.101)$$

eq. (14.98) decomposed s.t.  $x_1, \dots, x_k$  become **mutal independent** (??):

$$p(X) = \prod_{i=1}^k \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{(x_i - \mu_i)^2}{2\sigma_i^2}\right) \quad (14.102)$$

## 6.1. Gamma Distribution

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

**Definition 14.44 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:

$$f(x; \alpha, \beta) = \begin{cases} \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases} \quad (14.103)$$

$$\Gamma(\alpha) \stackrel{\text{eq. (4.65)}}{=} \int_0^\infty t^{\alpha-1} e^{-t} dt \quad (14.104)$$

with

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

## 6.2. Delta Distribution

**Definition 14.45 The delta function  $\delta(x)$ :**

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

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

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

Or alternatively by:

$$\delta(x - x_0) = \lim_{\sigma \rightarrow 0} \mathcal{N}(x|x_0, \sigma) \quad (14.105)$$

$$\approx \infty \mathbb{1}_{\{x=x_0\}} \quad (14.106)$$

**Property 14.21 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:**

$$\int_{\mathbb{R}} \delta(x - x_0) f(x) dx = f(x_0) \quad (14.107)$$

- Symmetry:**

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

- Scaling:**

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

### 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. (14.105) would be a non-differentiable/discrete form of the dirac measure.

## Proofs

**Proof.** proposition 14.1: Let  $X$  be normally distributed with  $X \sim \mathcal{N}(\mu, \sigma^2)$ :

$$f_Y(y) \stackrel{y \geq 0}{=} \mathbb{P}_Y(Y \leq y) = \mathbb{P}(a + bX \leq y) = \mathbb{P}_X\left(X \leq \frac{y-a}{b}\right)$$

$$= \mathbb{F}_X\left(\frac{y-a}{b}\right)$$

$$f_Y(y) \stackrel{y \leq 0}{=} \mathbb{P}_Y(Y \leq y) = \mathbb{P}(a + bX \leq y) = \mathbb{P}_X\left(X \geq \frac{y-a}{b}\right)$$

$$= 1 - \mathbb{F}_X\left(\frac{y-a}{b}\right)$$

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

$$f_Y(y) = \frac{d\mathbb{F}_Y(y)}{dy} = \begin{cases} \frac{1}{b} \frac{d\mathbb{F}_X\left(\frac{y-a}{b}\right)}{dy} \\ \frac{1}{-b} \frac{d\mathbb{F}_X\left(\frac{y-a}{b}\right)}{dy} \end{cases} = \frac{1}{|b|} f_X(x) \left(\frac{y-a}{b}\right)$$

eq. (14.95)).

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:

$$f_Y(y) = \frac{1}{\sqrt{2\pi\sigma}|b|} \exp\left\{-\frac{1}{2}\left(\frac{\frac{y-a}{b} - \mu}{\sigma}\right)^2\right\}$$

$$= \frac{1}{\sqrt{2\pi\sigma}|b|} \exp\left\{-\frac{1}{2}\left(\frac{y - (a + b\mu)}{\sigma|b|}\right)^2\right\}$$

□

**Proof.** proposition 14.2: Let  $X$  be normally distributed with  $X \sim \mathcal{N}(\mu, \sigma^2)$ :

$$Z := \frac{X - \mu}{\sigma} = \frac{1}{std} X - \frac{\mu}{\sigma} = aX + b \quad \text{with } a = \frac{1}{\sigma}, b = -\frac{\mu}{\sigma}$$

eq. (14.95)  $\mathcal{N}(a\mu + b, a^2\sigma^2) \sim \mathcal{N}\left(\frac{\mu}{\sigma} - \frac{\mu}{\sigma}, \frac{\sigma^2}{\sigma^2}\right) \sim \mathcal{N}(0, 1)$

□

**Proof.** proposition 14.3: Let  $X$  be normally distributed with  $X \sim \mathcal{N}(\mu, \sigma^2)$ :

$$F_X(x) = \mathbb{P}(X \leq x) \stackrel{-\mu}{=} \mathbb{P}\left(\frac{X - \mu}{\sigma} \leq \frac{x - \mu}{\sigma}\right) \mathbb{P}\left(Z \leq \frac{x - \mu}{\sigma}\right)$$

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

□

**Proof.** Property 14.19 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  $x = ay + b$   $a \in \mathbb{R}_+, b \in \mathbb{R}$

**Using** the Change of variables formula it follows:

$$p_x(\bar{x}) \stackrel{??}{=} \frac{p_y(\bar{y})}{\left|\frac{dx}{dy}\right|} \stackrel{\bar{y} = \frac{\bar{x}-b}{a}}{=} \frac{1}{a} \frac{1}{\sqrt{2\pi\mu^2}} \exp\left(-\frac{1}{2\sigma^2}\left(\frac{\bar{x}-b}{a} - \mu\right)^2\right)$$

$$= \frac{1}{\sqrt{2\pi a^2 \mu^2}} \exp\left(-\frac{1}{2\sigma^2 a^2}(\bar{x} - b - a\mu)^2\right)$$

$\mu_x$

Hence

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

□

### Note

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

$$\mathbb{E}[x] = \mathbb{E}[ay + b] = a\mathbb{E}[y] + b = a\mu + b$$

$$\mathbb{V}[x] = \mathbb{V}[ay + b] = a^2\mathbb{V}[y] = a^2\sigma^2$$

**Proof.** Property 14.20

From Property 14.19 it follows immediately that  $z$  is GRV  $z \sim \mathcal{N}(\mu_z, \Sigma_z)$  with:

$$z = A\xi \quad \text{with} \quad A = [A_x \quad A_y] \quad \text{and} \quad \xi = \begin{pmatrix} x & y \end{pmatrix}$$

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

$$\mathbb{E}[z] = \mathbb{E}[A_x x + A_y y] = A_x \mu_x + A_y \mu_y$$

$$\mathbb{V}[z] = \mathbb{V}[A\xi] \stackrel{\text{Property 14.13}}{=} A\mathbb{V}[\xi]A^\top$$

$$= [A_x \quad A_y] \begin{bmatrix} \mathbb{V}[x] & \text{Cov}[x, y] \\ \text{Cov}[y, x] & \mathbb{V}[y] \end{bmatrix} [A_x \quad A_y]^\top$$

$$= [A_x \quad A_y] \begin{bmatrix} \mathbb{V}[x] & \text{Cov}[x, y] \\ \text{Cov}[y, x] & \mathbb{V}[y] \end{bmatrix} \begin{bmatrix} A_x^\top \\ A_y^\top \end{bmatrix}$$

$$= A_x \mathbb{V}[x] A_x^\top + A_y \mathbb{V}[y] A_y^\top$$

$$+ \underbrace{A_y \text{Cov}[y, x] A_x^\top}_{=0 \text{ by independence}} + \underbrace{A_x \text{Cov}[x, y] A_y^\top}_{=0 \text{ by independence}}$$

$$= A_x \Sigma_x A_x^\top + A_y \Sigma_y A_y^\top$$

□

### Note

Can also be proved by using the normal definition of (def. 14.27) and tedious computations.



## 7. 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. (14.78)). Furthermore repeated calls to these RNG are independent, that is:

$$\begin{aligned} \mathbb{P}_{U_1, U_2}(u_1, u_2) &\stackrel{\text{eq. (14.23)}}{=} \mathbb{P}_{U_1}(u_1) \cdot \mathbb{P}_{U_2}(u_2) \\ &= \begin{cases} 1 & \text{if } u_1, u_2 \in [a, b] \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

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

## 8. Inverse-transform Technique

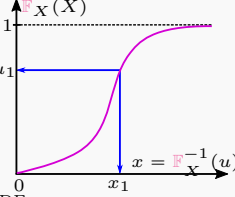
### Idea

Can make use of section 1 and the fact that CDF are increasing functions (<sup>[def. 4.8]</sup>). **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.



### 8.1. Continuous Case

**Definition 14.46 One Continuous Variable:** **Given:** a desired continuous pdf  $f_X$  and uniformly distributed rn  $\{u_1, u_2, \dots\}$ :

- Integrate the desired pdf  $f_X$  in order to obtain the desired cdf  $\mathbb{F}_X$ :

$$\mathbb{F}_X(x) = \int_{-\infty}^x f_X(t) dt \quad (14.108)$$

- Set  $\mathbb{F}_X(X) \stackrel{!}{=} U$  on the range of  $X$  with  $U \sim \mathcal{U}[0, 1]$ .

- Invert this equation/find the inverse  $\mathbb{F}_X^{-1}(U)$  i.e. solve:

$$U = \mathbb{F}_X(X) = \mathbb{F}_X\left(\underbrace{\mathbb{F}_X^{-1}(U)}_X\right) \quad (14.109)$$

- Plug in the uniformly distributed rn:

$$x_i = \mathbb{F}_X^{-1}(u_i) \quad \text{s.t.} \quad x_i \sim f_X \quad (14.110)$$

### Definition 14.47 Multiple Continuous Variable:

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

- Use the product rule (eq. (14.21)) in order to decompose  $f_{X,Y}$ :

$$f_{X,Y} = f_{X,Y}(x, y) = f_{X|Y}(x|y) f_Y(y) \quad (14.111)$$

- Use <sup>[def. 14.48]</sup> to first get a rv for  $y$  of  $Y \sim f_Y(y)$ .

- Then with this fixed  $y$  use <sup>(def. 14.45)</sup> again to get a value for  $x$  of  $X \sim f_{X|Y}(x|y)$ .

*Proof.* <sup>[def. 14.48]</sup>:

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

**Assume** that  $\mathbb{F}_X$  is strictly increasing (<sup>[def. 4.8]</sup>). Then for any  $u \in [0, 1]$  there must exist a **unique**  $x$  s.t.  $\mathbb{F}_X(x) = u$ .

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

Now let  $a$  arbitrary:

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

Since  $\mathbb{F}_X$  is strictly increasing:

$$\begin{aligned} \mathbb{P}(\mathbb{F}_X^{-1}(U) \leq a) &= \mathbb{P}(U \leq \mathbb{F}_X(a)) \\ &\stackrel{\text{eq. (14.78)}}{=} \int_0^{\mathbb{F}_X(a)} 1 dt = \mathbb{F}_X(a) \end{aligned}$$

### Note

Strictly speaking we may not assume that a CDF is **strictly** increasing but we as all CDFs are weakly increasing (<sup>[def. 4.8]</sup>) we may always define an auxiliary function by its infimum:

$$\hat{\mathbb{F}}_X^{-1} := \inf \{x | \mathbb{F}_X(X) \geq 0\} \quad u \in [0, 1] \quad (14.112)$$

### 8.2. Discret Case

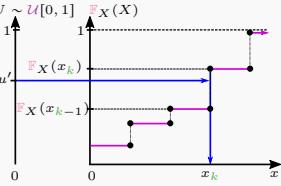
#### Idea

**Given:** a desired  $U \sim \mathcal{U}[0, 1]$   $\mathbb{F}_X(X)$  discret pmf  $\mathbb{P}_X$  s.t.  $\mathbb{F}_X(X = x_i) = p_X(x_i)$  and uniformly distributed rn  $\{u_1, u_2, \dots\}$ .

**Goal:** given a uniformly distributed rn  $u$  determine  $k$  s.t.:

$$\begin{aligned} \sum_{i=1}^{k-1} < U \leq \sum_{i=1}^k &\iff \mathbb{F}_X(x_{k-1}) < u \leq \mathbb{F}_X(x_k) \end{aligned} \quad (14.113)$$

and return  $x_k$ .



### Definition 14.48 One Discret Variable:

- Compute the CDF of  $\mathbb{P}_X$  (<sup>[def. 14.8]</sup>)

$$\mathbb{F}_X(x) = \sum_{t=-\infty}^x \mathbb{P}_X(t) \quad (14.114)$$

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

$$\mathbb{F}_X(x_{k(i)-1}) < u_i \leq \mathbb{F}_X(x_{k(i)}) \quad \forall u_i \quad (14.115)$$

*Proof.*  $??$ : First of all notice that we can always solve for an unique  $x_k$ .

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

**Given** a fixed  $x_k$  determine the values of  $u$  for which:

$$\mathbb{F}_X(x_{k-1}) < u \leq \mathbb{F}_X(x_k) \quad (14.116)$$

Now observe that:

$$\begin{aligned} u &\leq \mathbb{F}_X(x_k) = \mathbb{F}_X(x_{k-1}) + \mathbb{P}_X(x_k) \\ \Rightarrow \mathbb{F}_X(x_{k-1}) &< u \leq \mathbb{F}_X(x_{k-1}) + \mathbb{P}_X(x_k) \end{aligned}$$

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

$$\begin{aligned} \mathbb{P}(U \in [\mathbb{F}_X(x_{k-1}), \mathbb{F}_X(x_k)]) &= \int_{\mathbb{F}_X(x_{k-1})}^{\mathbb{F}_X(x_k)} \mathbb{P}_U(t) dt \\ &= \int_{\mathbb{F}_X(x_{k-1})}^{\mathbb{F}_X(x_k)} 1 dt = \int_{\mathbb{F}_X(x_{k-1})}^{\mathbb{F}_X(x_{k-1}) + \mathbb{P}_X(x_k)} 1 dt = \mathbb{P}_X(x_k) \end{aligned}$$

Hence the random variable  $x_k \in \mathcal{X}$  has the pdf  $\mathbb{P}_X$ .  $\square$

### Definition 14.49

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

**Given:** a pdf of multiple rvs  $\mathbb{P}_{X,Y}$ :

- Use the product rule (eq. (14.21)) in order to decompose  $\mathbb{P}_{X,Y}$ :

$$\mathbb{P}_{X,Y} = \mathbb{P}_{X,Y}(x, y) = \mathbb{P}_{X|Y}(x|y) \mathbb{P}_Y(y) \quad (14.117)$$

- Use  $??$  to first get a rv for  $y$  of  $Y \sim \mathbb{P}_Y(y)$ .

- Then with this fixed  $y$  use  $??$  again to get a value for  $x$  of  $X \sim \mathbb{P}_{X|Y}(x|y)$ .

### Definition 14.50

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

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

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

#### Define

$$\begin{aligned} \mathbb{P}_Z(1) &= \mathbb{P}_{X,Y}(1, 1), \mathbb{P}_Z(2) = \mathbb{P}_{X,Y}(1, 2), \dots \\ \dots, \mathbb{P}_Z(N_x \cdot N_y) &= \mathbb{P}_{X,Y}(N_x, N_y) \end{aligned}$$

Then simply apply  $??$  to the auxiliary pdf  $\mathbb{P}_Z$

- Use the product rule (eq. (14.21)) in order to decompose  $f_{X,Y}$ :

$$f_{X,Y} = f_{X,Y}(x, y) = f_{X|Y}(x|y) f_Y(y) \quad (14.118)$$

- Use <sup>[def. 14.48]</sup> to first get a rv for  $y$  of  $Y \sim f_Y(y)$ .

- Then with this fixed  $y$  use <sup>[def. 14.48]</sup> again to get a value for  $x$  of  $X \sim f_{X|Y}(x|y)$ .

*nice examples see comment in code text*

9. Descriptive Statistics

9.1. Population Parameters

**Definition 14.51 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 14.52 Population Mean:** Given a population  $\{x_i\}_{i=1}^N$  of size  $N$  its variance is defined as:

μ = 1/N ∑ x\_i (14.119)

**Definition 14.53 Population Variance:** Given a population  $\{x_i\}_{i=1}^N$  of size  $N$  its variance is defined as:  $\{x_i\}_{i=1}^N$

σ² = 1/N ∑ (x\_i - μ)² (14.120)

Note

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

9.2. Sample Estimates

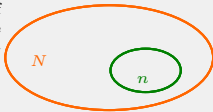
**Definition 14.54 (Sample) Statistic:** A statistic is a measurable function  $f$  that assigns a **single** value  $F$  to a sample of random variables or population:

f : ℝⁿ → ℝ F = f(X₁, ..., Xₙ)

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 14.55 (Point) Estimator**  $\hat{\theta} = \hat{\theta}(\mathbf{X})$ :  
**Given:** n-samples  $\mathbf{x}_1, \dots, \mathbf{x}_n \sim \mathbf{X}$  an estimator  
 $\hat{\theta} = h(\mathbf{x}_1, \dots, \mathbf{x}_n)$  (14.121)

is a statistic/random variable used to estimate a true (population) parameter  $\theta$ <sup>[def. 14.51]</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 14.56 Degrees of freedom of a Statistic:** Is the number of values in the final calculation of a statistic that are free to vary.

9.2.1. Empirical Mean

**Definition 14.57 Sample/Empirical Mean**  $\bar{x}$ :  
The sample mean is an estimate/statistic of the population mean<sup>[def. 14.52]</sup> and can be calculated from an observation/sample of the total population  $\{x_i\}_{i=1}^n \subset \{x_i\}_{i=1}^N$ :

̄x = μ̂\_X = 1/n ∑ x\_i (14.122)

**Corollary 14.6 Expectation:** The sample mean estimator is unbiased (see section 13):

E[μ̂\_X] = μ (14.123)

**Corollary 14.7 Variance:** For the variance of the sample mean estimator it holds (see section 13):

V[μ̂\_X] = 1/n σ\_X² (14.124)

9.2.2. Empirical Variance

**Definition 14.58 Biased Sample Variance:** The sample mean is an estimate/statistic of the population variance<sup>[def. 14.53]</sup> and can be calculated from an observation/sample of the total population  $\{x_i\}_{i=1}^n \subset \{x_i\}_{i=1}^N$ :

s²\_n = σ̂²\_X = 1/n ∑ (x\_i - μ)² (14.125)

**Definition 14.59 (Unbiased) Sample Variance:**

s² = σ̂²\_X = 1/(n-1) ∑ (x\_i - μ)² (14.126)

see section 13

**Definition 14.60 Bessel's Correction:** The factor  $\frac{n}{n-1}$  (14.127)

as multiplying the uncorrected population varianceeq. (14.125) 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.



10. Statistical Tests

**Definition 14.61 Null Hypothesis:** A Null Hypothesis  $H_0$  is usually a commonly accepted fact/view/base hypothesis that researchers try to nullify or disprove.

H₀ : θ = θ₀ (14.128)

**Definition 14.62 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 : θ { > θ₀ (one-sided), < θ₀ (one-sided), ≠ θ₀ (two-sided) } (14.129)

**Definition 14.63 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  $\Theta_1$ .

Θ = Θ₀ ∪ Θ₁ Θ₀ ∩ Θ₁ = ∅

Null Hypothesis H₀ : θ ∈ Θ₀ (14.130)

Alternative Hypothesis H\_A : θ ∈ Θ₁ (14.131)

10.1. Type I&II Errors

**Definition 14.64 Type I Error:** Is the rejection of a Null Hypothesis, even-tough its true (also known as a "false positive").

**Definition 14.65 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 14.66 Critical Value c:** Value from which on the Null-hypothesis  $H_0$  gets rejected.

**Definition 14.67 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 14.68 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\_α ∈ Θ : P(Type I Error) ≤ α (14.132)

or P(c₂ ≤ X ≤ c₁) ≤ α two-sided  
P(c₂ ≤ X) ≤ α/2 and P(X ≤ c₁) ≤ α/2  
P(c₂ ≤ X) ≤ α one-sided  
P(X ≤ c₁) ≤ α one-sided

**Definition 14.69 Acceptance Region:** Is the region where we accept the null hypothesis  $H_0$ .

Note

see example 14.3.

10.2. Normally Distributed Data

Let us consider a sample of  $\{x_i\}_{i=1}^n$  i.i.d. observations, that follow a normal distribution  $x_i \sim \mathcal{N}(\mu, \sigma^2)$ .

10.2.1. Z-Test σ known

10.2.2. t-Test σ unknown

11. 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?

12. Examples

**Example 14.1 Theorem 14.4:** Let  $x$  be uniformly distributed on  $[0, 1]$  (def. 14.37) 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 14.2 Theorem 14.4:** Let

add <https://www.youtube.com/watch?v=WUUhTVIRagg>

**Example 14.3 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:

$$H_0 : p = p_0 = 0.1 \quad \text{vs.} \quad H_A : p > 0.1$$

We model the number of number of defective goods using the binomial distribution<sup>[def. 14.34]</sup>

$$X \sim \mathcal{B}(n, p), n = 20 \quad \mathbb{P}(X \geq x) = \sum_{k=x}^n \binom{n}{k} p^k (1-p)^{n-k}$$

from this we find:

$$\begin{aligned} \mathbb{P}_{p_0}(X \geq 4) &= 1 - \mathbb{P}_{p_0}(X \leq 3) = 0.13 \\ \mathbb{P}_{p_0}(X \geq 4) &= 1 - \mathbb{P}_{p_0}(X \leq 3) = 0.04 \leq \alpha \end{aligned}$$

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  $\mathbb{P}(X \geq x) = 1 - \mathbb{P}(X \leq x-1)$

13. Proofs

*Proof.* corollary 14.6:

$$\mathbb{E}[\hat{\mu}_X] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \frac{1}{n} \mathbb{E}\left[\sum_{i=1}^n x_i\right] = \frac{1}{n} \mathbb{E}\left[\underbrace{\mu + \dots + \mu}_{1, \dots, n}\right]$$

□

*Proof.* corollary 14.7:

$$\begin{aligned} \mathbb{V}[\hat{\mu}_X] &= \mathbb{V}\left[\frac{1}{n} \sum_{i=1}^n x_i\right] \stackrel{\text{Property 14.12}}{=} \frac{1}{n^2} \mathbb{V}\left[\sum_{i=1}^n x_i\right] \\ \frac{1}{n^2} n \mathbb{V}[X] &= \frac{1}{n} \sigma^2 \end{aligned}$$

□

*Proof.* definition 14.59:

$$\begin{aligned} \mathbb{E}[\hat{\sigma}_X^2] &= \mathbb{E}\left[\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2\right] \\ &= \frac{1}{n-1} \mathbb{E}\left[\sum_{i=1}^n (x_i^2 - 2x_i \bar{x} + \bar{x}^2)\right] \\ &= \frac{1}{n-1} \mathbb{E}\left[\sum_{i=1}^n x_i^2 - 2\bar{x} \sum_{i=1}^n x_i + \sum_{i=1}^n \bar{x}^2\right] \\ &= \frac{1}{n-1} \mathbb{E}\left[\sum_{i=1}^n x_i^2 - 2n\bar{x} \cdot n\bar{x} + n\bar{x}^2\right] \\ &= \frac{1}{n-1} \mathbb{E}\left[\sum_{i=1}^n x_i^2 - n\bar{x}^2\right] \\ &= \frac{1}{n-1} \left[\sum_{i=1}^n \mathbb{E}[x_i^2] - n\mathbb{E}[\bar{x}^2]\right] \\ &= \frac{1}{n-1} \left[\sum_{i=1}^n (\sigma^2 + \mu^2) - n\mathbb{E}[\bar{x}^2]\right] \\ &= \frac{1}{n-1} \left[\sum_{i=1}^n (\sigma^2 + \mu^2) - n\left(\frac{1}{n}\sigma^2 + \mu^2\right)\right] \\ &= \frac{1}{n-1} \left[(n\sigma^2 + n\mu^2) - (\sigma^2 + n\mu^2)\right] \\ &= \frac{1}{n-1} \left[n\sigma^2 - \sigma^2\right] = \frac{1}{n-1} \left[(n-1)\sigma^2\right] = \sigma^2 \end{aligned}$$

□

Stochastic Calculus

Stochastic Processes

**Definition 15.1 Random/Stochastic Process**  $\{X_t, t \in \mathcal{T} \subseteq \mathbb{R}_+\}$ : 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:

$$X: \mathbb{R} \times \Omega \mapsto \mathbb{R} \iff (t, \omega) \mapsto X(t, \omega) \quad (15.1)$$

**Definition 15.2 Sample path/Trajectory/Realization:** Is the *stochastic/noise signal*  $r(\cdot, \omega)$  on the index set  $\mathcal{T}$ , that we obtain be sampling  $\omega$  from  $\Omega$ .

**Notation**  
Even though the r.v.  $X$  is a function of two variables, most books omit the argument of the sample space  $X(t, \omega) := X(t)$

**Definition 15.3 Filtration**  $\mathbb{F} = \{\mathcal{F}_t\}_{t \geq 0}$ : 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 *increasing*:

$$\mathcal{F}_s \subseteq \mathcal{F}_t \quad \forall s \leq t \quad (15.2)$$

**Definition 15.4 Adapted Process:** A stochastic process  $\{X_t : 0 \leq t \leq \infty\}$  is called adapted to a filtration  $\mathbb{F}$  if,  $X_t$  is  $\{\mathcal{F}_{t-1}\}$ -measurable, i.e. the value of  $X_t$  is known at time  $t-1$ .

**Definition 15.5 Predictable Process:** A stochastic process  $\{X_t : 0 \leq t \leq \infty\}$  is called predictable w.r.t. a filtration  $\mathbb{F}$  if,  $X_t$  is  $\{\mathcal{F}_{t-1}\}$ -measurable, i.e. the value of  $X_t$  is known at time  $t-1$ .

**Note**  
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.

**Definition 15.6 Filtered Probability Space**  $(\Omega, \mathcal{F}, \{\mathcal{F}_t\}_{t \geq 0}, \mathbb{P})$ : A probability space  $(\Omega, \mathcal{F}, \mathbb{P})$  together with a filtration  $\{\mathcal{F}_t\}_{t \geq 0}$  is called a *filtered probability space*.

**Corollary 15.1** : The amount of information of an adapted random process is increasing see example 15.1.

**Definition 15.7 Martingales:** A stochastic process  $X(t)$  is a martingale on a *filtered probability space*  $(\Omega, \mathcal{F}, \{\mathcal{F}_t\}_{t \geq 0}, \mathbb{P})$  if the following conditions hold:

① Given  $s \leq t$  the best prediction of  $X(t)$ , with a filtration  $\{\mathcal{F}_s\}$  is the current expected value:  
$$\forall s \leq t \quad \mathbb{E}[X(t)|\mathcal{F}_s] = X(s) \quad \text{a.s.} \quad (15.3)$$

② The expectation is finite:  
$$\mathbb{E}[|X(t)|] < \infty \quad \forall t \geq 0 \quad X(t) \text{ is } \{\mathcal{F}_t\}_{t \geq 0} \text{ adapted} \quad (15.4)$$

**Interpretation**

- For any  $\mathcal{F}_s$ -adapted process the best prediction of  $X(t)$  is the currently known value  $X(s)$  i.e. if  $\mathcal{F}_s = \mathcal{F}_{t-1}$  then the best prediction is  $X(t-1)$
- A martingale models fair games of limited information.

**Definition 15.8 Auto Covariance**  $\gamma(t_2 - t_1)$ : Describes the covariance<sup>[def. 14.28]</sup> between two values of a stochastic process  $(X_t)_{t \in \mathcal{T}}$  at different time points  $t_1$  and  $t_2$ .

$$\gamma(t_1, t_2) = \text{Cov}[X_{t_1}, X_{t_2}] = \mathbb{E}[(X_{t_1} - \mu_{t_1})(X_{t_2} - \mu_{t_2})] \quad (15.5)$$

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

$$\gamma(t, t) = \text{Cov}[X_t, X_t] \stackrel{\text{eq. (14.64)}}{=} \mathbb{V}[X_t] \quad (15.6)$$

**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  $\mathbf{x}(t)$  the autocorrelation function  $\gamma(t, t - \tau)$  describes how *similar* the time translated function  $\mathbf{x}(t - \tau)$  and the original function  $\mathbf{x}(t)$  are.
- If there exists some relation between the values of the time series that is non-random, then the autocorrelation is non-zero.
- The auto covariance is maximized/most similar for no translation  $\tau = 0$  at all.

**Definition 15.9 Auto Correlation**  $\rho(t_2 - t_1)$ : Is the scaled version of the auto-covariance<sup>[def. 15.8]</sup>:

$$\begin{aligned} \rho(t_2 - t_1) &= \frac{\text{Cov}[X_{t_1}, X_{t_2}]}{\sigma_{X_{t_1}} \sigma_{X_{t_2}}} = \frac{\mathbb{E}[(X_{t_1} - \mu_{t_1})(X_{t_2} - \mu_{t_2})]}{\sigma_{X_{t_1}} \sigma_{X_{t_2}}} \end{aligned} \quad (15.7)$$

1. Different kinds of Processes

1.1. Markov Process

**Definition 15.10 Markov Process:** 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 *transition probability* – the probability of  $X(t)$  lying in the set  $B$  at time  $t$ , given the value  $x$  of the process at time  $s$ :

$$\mathbb{P}(s, x, t, B) = P(X(t) \in B | X(s) = x) \quad 0 \leq s < t \quad (15.8)$$

**Interpretation**  
In order to predict the future only the current/last value counts.

**Corollary 15.2 Transition Density:** The transition probability of a continuous distribution  $\mathbf{p}$  can be calculated via:

$$\mathbb{P}(s, x, t, B) = \int_B \mathbf{p}(s, x, t, y) dy \quad (15.9)$$

1.2. Gaussian Process

**Definition 15.11 Gaussian Process:** Is a stochastic process  $X(t)$  where the random variables follow a Gaussian distribution:

$$X(t) \sim \mathcal{N}(\mu(t), \sigma^2(t)) \quad \forall t \in T \quad (15.10)$$

1.3. Diffusions

**Definition 15.12 Diffusion:** Is a Markov Process<sup>[def. 15.10]</sup> for which it holds that:

$$\begin{aligned} \mu(t, X(t)) &= \lim_{t \rightarrow 0} \frac{1}{\Delta t} \mathbb{E}[X(t + \Delta t) - X(t) | X(t)] \quad (15.11) \\ \sigma^2(t, X(t)) &= \lim_{t \rightarrow 0} \frac{1}{\Delta t} \mathbb{E}[(X(t + \Delta t) - X(t))^2 | X(t)] \end{aligned} \quad (15.12)$$

See ??/eq. (15.12) for simple proof of eq. (15.11)/??.

- $\mu(t, X(t))$  is called **drift**
- $\sigma^2(t, X(t))$  is called **diffusion coefficient**

**Interpretation**  
There exist not discontinuities for the trajectories.

1.4. Brownian Motion/Wiener Process

**Definition 15.13 d-dim standard Brownian Motion/Wiener Process:** Is an  $\mathbb{R}^d$  valued *stochastic process*<sup>[def. 15.1]</sup>  $(W_t)_{t \in \mathcal{T}}$  starting at  $\mathbf{x}_0 \in \mathbb{R}^d$  that satisfies:

① **Normal Independent Increments:** the increments are *normally distributed independent random variables*:  
$$W(t_i) - W(t_{i-1}) \sim \mathcal{N}(0, (t_i - t_{i-1}) \mathbb{1}_{d \times d}) \quad \forall i \in \{1, \dots, T\} \quad (15.13)$$

② **Stationary increments:**  $W(t + \Delta t) - W(t)$  is independent of  $t \in \mathcal{T}$

③ **Continuity:** for a.e.  $\omega \in \Omega$ , the function  $t \mapsto W_t(\omega)$  is continuous  
$$\lim_{t \rightarrow 0} \frac{\mathbb{P}(|W(t + \Delta t) - W(t)| \geq \delta)}{\Delta t} = 0 \quad \forall \delta > 0 \quad (15.14)$$

④ **Start**  
$$W(0) := W_0 = 0 \quad \text{a.s.} \quad (15.15)$$

**Notation**

- In many source the Brownian motion is a synonym for the standard Brownian Motion and it is the same as the Wiener process.
- However in some sources the Wiener process is the standard Brownian Motion, while the Brownian motion denotes a general form  $\alpha W(t) + \beta$ .

**Corollary 15.3**  $W_t \sim \mathcal{N}(0, \sigma)$ : The random variable  $W_t$  follows the  $\mathcal{N}(0, \sigma)$  law

$$\begin{aligned} \mathbb{E}[W(t)] &= \mu = 0 \quad (15.16) \\ \mathbb{V}[W(t)] &= \mathbb{E}[W^2(t)] = \sigma^2 = t \quad (15.17) \end{aligned}$$

See section 5

1.4.1. Properties of the Wiener Process

**Property 15.1 Non-Differentiable Trajectories:** The sample paths of a Brownian motion are not differentiable:

$$\begin{aligned} \frac{dW(t)}{dt} &= \lim_{t \rightarrow 0} \mathbb{E} \left[ \left( \frac{W(t + \Delta t) - W(t)}{\Delta t} \right)^2 \right] \\ &= \lim_{t \rightarrow 0} \frac{\mathbb{E}[W(t + \Delta t) - W(t)]}{\Delta t} = \lim_{t \rightarrow 0} \frac{\sigma^2}{\Delta t} = \infty \end{aligned}$$

$\xrightarrow{\text{result}}$  cannot use normal calculus anymore  
 $\xrightarrow{\text{solution}}$  Ito Calculus see section 16.

**Property 15.2 Auto covariance Function:** The auto-covariance<sup>[def. 15.8]</sup> for a Wiener process

$$\mathbb{E}[(W(t) - \mu(t))(W(t') - \mu(t'))] = \min(t, t') \quad (15.18)$$

**Property 15.3** : A standard Brownian motion is a

Quadratic Variation

**Definition 15.14 Total Variation:** The total variation of a function  $f: [a, b] \subset \mathbb{R} \mapsto \mathbb{R}$  is defined as:

$$LV_{[a, b]}(f) = \sup_{\Pi \in \mathcal{S}} \sum_{i=0}^{n_{\Pi}-1} |f(x_{i+1}) - f(x_i)| \quad (15.19)$$

$$S = \left\{ \Pi \{x_0, \dots, x_{n_{\Pi}}\} : \Pi \text{ is a partition } \stackrel{[\text{def. 10.1}]}{\text{of}} [a, b] \right\}$$

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.

**Definition 15.15 Total Quadratic Variation/“sum of squares”:** The total quadratic variation of a function  $f: [a, b] \subset \mathbb{R} \mapsto \mathbb{R}$  is defined as:

$$QV_{[a, b]}(f) = \sup_{\Pi \in \mathcal{S}} \sum_{i=0}^{n_{\Pi}-1} |f(x_{i+1}) - f(x_i)|^2 \quad (15.20)$$

$$S = \left\{ \Pi \{x_0, \dots, x_{n_{\Pi}}\} : \Pi \text{ is a partition } \stackrel{[\text{def. 10.1}]}{\text{of}} [a, b] \right\}$$

**Corollary 15.4 Bounded (quadratic) Variation:** The (quadratic) variation<sup>[def. 15.14]</sup> of a function is bounded if it is finite:

$$\exists M \in \mathbb{R}_+ : LV_{[a, b]}(f) \leq M \iff (QV_{[a, b]}(f) \leq M) \quad \forall \Pi \in \mathcal{S} \quad (15.21)$$

**Theorem 15.1 Variation of Wiener Process:** Almost surely the total variation of a Brownian motion over a interval  $[0, T]$  is infinite:

$$\mathbb{P}(\omega : LV(W(\omega)) < \infty) = 0 \quad (15.22)$$

**Theorem 15.2 Quadratic Variation of standard Brownian Motion:** 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

$$(15.23)$$

See ??

**Corollary 15.5** : theorem 15.2 can also be written as:

$$(dW(t))^2 = dt \quad (15.24)$$

1.4.2. Lévy’s Characterization of BM

**Theorem 15.3 d-dim standard BM/Wiener Process by Paul Lévy:** An  $\mathbb{R}^d$  valued *adapted stochastic process*<sup>[def.x, 15.1, 15.3]</sup>  $(W_t)_{t \in \mathcal{T}}$  with the filtration  $\{\mathcal{F}_t\}_{t \in \mathbb{R}_+}$ , that satisfies:

① **Start**  
$$W(0) := W_0 = 0 \quad \text{a.s.} \quad (15.25)$$

② **Continuous Martingale:**  $W_t$  is an a.s. *continuous martingale*<sup>[def. 15.7]</sup> w.r.t. the filtration  $(\mathcal{F}_t)_{t \in \mathcal{T}}$  under  $\mathbb{P}$ .

③ **Quadratic Variation:**  
 $W_t^2 - t$  is also an martingale  $\iff QV(W_t) = t$  (15.26)

is a standard Brownian motion<sup>[def. 15.18]</sup>. Proof see section 5

Further Stochastic Processes

1.4.3. White Noise

understand script and add

**Definition 15.16 Discrete-time white noise:** Is a random signal  $\{\epsilon_t\}_{t \in T_{\text{discret}}}$  having equal intensity at different frequencies and is defined by:

- Having zero tendencies/expectation (otherwise the signal would not be random):  
$$\mathbb{E}[\epsilon[k]] = 0 \quad \forall k \in T_{\text{discret}} \quad (15.27)$$
- Zero autocorrelation<sup>[def. 15.9]</sup>  $\gamma$  i.e. the signals of different times are in no-way correlated:  
$$\gamma(\epsilon[k], \epsilon[k + n]) = \mathbb{E}[\epsilon[k]\epsilon[k + n]^T] = \mathbb{V}[\epsilon[k]] \delta_{\text{discret}}[n] \quad \forall k, n \in T_{\text{discret}} \quad (15.28)$$

**With**  
$$\delta_{\text{discret}}[n] := \begin{cases} 1 & \text{if } n = 0 \\ 0 & \text{else} \end{cases}$$

See proofs

**Definition 15.17 Continuous-time white noise:** Is a random signal  $(\epsilon_t)_{t \in T_{\text{continuous}}}$  having equal intensity at different frequencies and is defined by:

- Having zero tendencies/expectation (otherwise the signal would not be random):  
$$\mathbb{E}[\epsilon(t)] = 0 \quad \forall t \in T_{\text{continuous}} \quad (15.29)$$
- Zero autocorrelation<sup>[def. 15.9]</sup>  $\gamma$  i.e. the signals of different times are in no-way correlated:  
$$\gamma(\epsilon(t), \epsilon(t + \tau)) = \mathbb{E}[\epsilon(t)\epsilon(t + \tau)^T] \quad (15.30)$$

$$\stackrel{\text{eq. (14.106)}}{=} \mathbb{V}[\epsilon(t)] \delta(t - \tau) = \begin{cases} \mathbb{V}[\epsilon(t)] & \text{if } \tau = 0 \\ 0 & \text{else} \end{cases}$$

$$\forall t, \tau \in T_{\text{continuous}} \quad (15.31)$$

#### 1.4.4. Generalized Brownian Motion

**Definition 15.18 Brownian Motion:**  
Let  $\{W_t\}_{t \in \mathbb{R}_+}$  be a standard Brownian motion<sup>[def. 15.13]</sup>, and define:

$$X_t = \mu t + \sigma W_t \quad t \in \mathbb{R}_+ \quad \begin{array}{l} \mu \in \mathbb{R} : \text{drift parameter} \\ \sigma \in \mathbb{R}_+ : \text{scale parameter} \end{array} \quad (15.32)$$

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 15.4 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:

$$f_{W(t)}(x) \sim \mathcal{N}(\mu t, \sigma^2 t) = \frac{1}{\sqrt{2\pi\sigma^2 t}} \exp\left\{-\frac{(x - \mu t)^2}{2\sigma^2 t}\right\} \quad (15.33)$$

**Corollary 15.6 :** More generally we may define the process:

$$t \mapsto f(t) + \sigma W_t \quad (15.34)$$

which corresponds to a noisy version of  $f$ .

**Corollary 15.7 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:

$$\begin{aligned} dX(t) &= \mu dt + \sigma dW(t) & (15.35) \\ X(0) &= 0 & (15.36) \end{aligned}$$

#### 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 15.19 Geometric Brownian Motion:**  
Let  $\{W_t\}_{t \in \mathbb{R}_+}$  be a standard Brownian motion<sup>[def. 15.13]</sup> the exponential transform:

$$X(t) = \exp(W(t)) = \exp(\mu t + \sigma W(t)) \quad t \in \mathbb{R}_+ \quad (15.37)$$

is called geometric Brownian motion

**Corollary 15.8 Log-normal Returns:** For a geometric BM we obtain log-normal returns:

$$\ln\left(\frac{S_t}{S_0}\right) = \mu t + \sigma W(t) \iff \mu t + \sigma W(t) \sim \mathcal{N}(\mu t, \sigma^2 t) \quad (15.38)$$

meaning that the mean and the variance of the process (stock) log-returns grow over time linearly.

**Corollary 15.9 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:

$$\begin{aligned} dX(t) &= X(t) (\mu dt + \sigma dW(t)) \\ &= \mu X(t) dt + \sigma X(t) dW(t) & (15.39) \\ X(0) &= 0 & (15.40) \end{aligned}$$

#### 1.4.6. Locally Brownian Motion

**Definition 15.20 Locally Brownian Motion:**  
Let  $\{W_t\}_{t \in \mathbb{R}_+}$  be a standard Brownian motion<sup>[def. 15.13]</sup> a local Brownian motion is a stochastic process  $X(t)$  that satisfies the SDE:

$$dX(t) = \mu(X(t), t) dt + \sigma(X(t), t) dW(t) \quad (15.41)$$

##### Note

A local Brownian motion is an generalization of a geometric Brownian motion.

#### 1.4.7. Ornstein-Uhlenbeck Process

**Definition 15.21 Ornstein-Uhlenbeck Process:**  
Let  $\{W_t\}_{t \in \mathbb{R}_+}$  be a standard Brownian motion<sup>[def. 15.13]</sup> a Ornstein-Uhlenbeck Process or exponentially correlated noise is a stochastic process  $X(t)$  that satisfies the SDE:

$$dX(t) = -aX(t) dt + b\sigma dW(t) \quad a > 0 \quad (15.42)$$

#### 1.5. Poisson Processes

**Definition 15.22 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 15.23 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. 14.36]</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  
 $X_{t_1} - X_{t_0}, X_{t_2} - X_{t_1}, \dots, X_{t_n} - X_{t_{n-1}} \quad (15.43)$

##### Interpretation

A Poisson Process is a continuous-time process with discrete, positive realizations in  $\mathbb{N}_{\geq 0}$

**Corollary 15.10 Probability of events:** Using Taylor in order to expand the Poisson distribution one obtains:

$$\mathbb{P}(X_{(t+\Delta t)} - X_t \neq 0) = \lambda \Delta t + o(\Delta t^2) \quad t \text{ small i.e. } t \rightarrow 0 \quad (15.44)$$

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 15.24 Differential of a Poisson Process:** The differential of a Poisson Process is defined as:

$$dX_t = \lim_{\Delta t \rightarrow dt} (X_{(t+\Delta t)} - X_t) \quad (15.45)$$

**Property 15.4 Probability of Events for differential:**  
With the definition of the differential and using the previous results from the Taylor expansion it follows:

$$\mathbb{P}(dX_t = 0) = 1 - \lambda \quad (15.46)$$

$$\mathbb{P}(|dX_t| = 1) = \lambda \quad (15.47)$$

##### Proofs

**Proof.** eq. (15.11):  
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. (15.12):

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:

$$\sigma^2 = \mathbb{E}[x^2(n)] - \mathbb{E}[x(n)]^2 = \mathbb{E}[x^2(n)] = \frac{\delta^2}{\Delta x} t \quad (15.48)$$

Thus in expectation the particles goes nowhere.  $\square$

**Proof.** eq. (15.30):

$$\begin{aligned} \gamma(\epsilon[k], \epsilon[k+n]) &= \text{Cov}[\epsilon[k], \epsilon[k+1]] \\ &= \mathbb{E}[(\epsilon[k] - \mathbb{E}[\epsilon[k]]) (\epsilon[k+n] - \mathbb{E}[\epsilon[k+n]])^T] \\ &\stackrel{\text{eq. (15.27)}}{=} \mathbb{E}[(\epsilon[k]) (\epsilon[k+n])] \quad \square \end{aligned}$$

**Proof.** corollary 15.3:

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  $[s-s, t-s] = [0, t-s]$

$$\begin{array}{ll} \text{Thus} & B_t - B_s \sim B_{t-s} - B_0 \\ \text{but as } B_0 \text{ is a.s. zero by definition eq. (15.15) it follows:} & \\ & B_t - B_s \sim B_{t-s} \quad B_{t-s} \sim \mathcal{N}(0, t-s) \quad \square \end{array}$$

**Proof.** corollary 15.3:

$$\begin{aligned} W(t) &= W(t) - \underbrace{W(0)}_{=0} \sim \mathcal{N}(0, t) \\ \Rightarrow \quad \mathbb{E}[X] &= 0 \quad \mathbb{V}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2 = t \quad \square \end{aligned}$$

**Proof.** theorem 15.2:

$$\begin{aligned} \sum_{k=0}^{N-1} [W(t_k) - W(t_{k-1})]^2 & \quad t_k = k \frac{T}{N} \\ &= \sum_{k=0}^{N-1} X_k^2 \quad X_k \sim \mathcal{N}\left(0, \frac{T}{N}\right) \\ &= \sum_{k=0}^{N-1} Y_k = n \left( \frac{1}{n} \sum_{k=0}^{N-1} Y_k \right) \quad \mathbb{E}[Y_k] = \frac{T}{N} \\ &\stackrel{\text{S.L.L.N}}{=} n \frac{T}{n} = T \quad \square \end{aligned}$$

**Proof.** theorem 15.3 ②:

1. first we need to show eq. (15.3):  $\mathbb{E}[W_t | \mathcal{F}_s] = W_s$   
Due to the fact that  $W_t$  is  $\mathcal{F}_t$  measurable i.e.  $W_t \in \mathcal{F}_t$  we know that:

$$\begin{aligned} \mathbb{E}[W_t | \mathcal{F}_t] &= W_t & (15.49) \\ \mathbb{E}[W_t | \mathcal{F}_s] &= \mathbb{E}[W_t - W_s + W_s | \mathcal{F}] \\ &= \mathbb{E}[W_t - W_s | \mathcal{F}_s] + \mathbb{E}[W_s | \mathcal{F}_s] \\ &\stackrel{\text{eq. (15.49)}}{=} \mathbb{E}[W_t - W_s] + W_s \\ &\stackrel{W_t - W_s \sim \mathcal{N}(0, t-s)}{=} W_s \end{aligned}$$

2. second we need to show eq. (15.4):  $\mathbb{E}[|X(t)|] < \infty$   
 $\mathbb{E}[|W(t)|]^2 \stackrel{\text{eq. (14.56)}}{\leq} \mathbb{E}[|W(t)|^2] = \mathbb{E}[W^2(t)] = t < \infty \quad \square$

**Proof.** theorem 15.3 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:

$$\begin{aligned} \mathbb{E}[W_t^2 | \mathcal{F}_s] &= \mathbb{E}[(W_t - W_s)^2 | \mathcal{F}_s] + \mathbb{E}[2W_s(W_t - W_s) | \mathcal{F}_s] \\ &\quad + \mathbb{E}[W_s^2 | \mathcal{F}_s] \\ &\stackrel{\text{eq. (15.49)}}{=} \mathbb{E}[(W_t - W_s)^2] + 2W_s \mathbb{E}[(W_t - W_s)] + W_s^2 \\ &\stackrel{\text{eq. (15.17)}}{=} \mathbb{V}[W_t - W_s] + 0 + W_s^2 \\ &\quad t - s + W_s^2 \end{aligned}$$

from this it follows that:

$$\mathbb{E}[W_t^2 - t | \mathcal{F}_s] = W_s^2 - s \quad \square$$

understand why  $\mathbb{E}[(W_t - W_s)^2 | \mathcal{F}] = \mathbb{E}[(W_t - W_s)^2]$

#### Examples

##### Example 15.1 :

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.

$$\mathcal{F} = \begin{cases} \{\emptyset, \Omega\} & t \in [0, T/2) \\ \{\emptyset, \{\omega_1, \omega_2\}, \{\omega_3, \omega_4\}, \Omega\} & t \in [T/2, T) \\ \mathcal{F}_{\max} = 2^\Omega & t = T \end{cases} \quad (15.50)$$

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