

Probability Notes

Amro Al Baali

Compiled: November 4, 2020

Contents

Notations	III
I Fundamentals	1
1 Single Variable Probability	2
1.1 Definitions and intuition	2
1.2 Basic theorems	5
1.3 Conditional probabilities	6
1.4 Random variables and distributions	8
1.4.1 Properties of random variables	9
1.4.2 Properties of PDFs	10
1.4.3 Conditional distributions	11
1.5 Transformed random variables	13
1.5.1 Method of distribution	13
1.5.2 Method of transformation	14
1.5.3 Discrete transformation	15
1.6 Statistics of random variables	16
1.6.1 Moments	18
1.6.2 Characteristic functions	18
2 Multivariate Probability	20
2.1 Multivariate random variables and distributions	20
2.1.1 Joint cumulative distribution function (JCDF)	20
2.1.2 Joint probability mass function (JPMF)	21
2.1.3 Joint probability density function (JPDF)	22
2.1.4 Independence	23
2.2 Transformed random variables	24
2.2.1 Method of distribution	25
2.2.2 Method of transformation	27
2.2.3 Linear transformation	30
2.2.4 Nonlinear transformation (covariance propagation)	31
2.3 Statistics of random variables	33

2.3.1	Order statistics	33
2.3.2	Mean	34
2.3.3	Covariance and correlation	35
2.4	Moments and characteristic functions	37
2.5	Conditional distributions	39
2.5.1	Properties of conditional probability	40
2.5.2	Conditional covariance	41
2.6	Conditional expectation and variance	41
2.6.1	Properties of conditional variance	43
2.7	Passing measurements through a function	44
3	Distributions	46
3.1	Bernoulli distribution	46
3.2	Geometric random variables	46
3.3	Binomial random variable	47
3.4	Poisson distribution	48
3.5	Uniform random variables	48
3.6	Exponential random variable	48
3.7	Gaussian random variables	48
3.7.1	Properties of Gaussian random variables	49
3.8	Chi-squared distribution	49
3.8.1	Relation to the Mahalanobis distance	50
4	Graphs in probability	53
4.1	Graph definitions and terminology	53
4.2	Graphs in Probability	55
II	Statistical Inference	57
5	Parameter Estimation	58
5.1	Motivation	58
5.2	Problem statement	59
5.3	Definitions and terminology	59
5.4	Method of moments	61
5.5	Maximum likelihood (ML) estimator	62
5.6	Bayesian parameter estimation	64
5.7	Maximum a-posteriori (MAP) estimator	66
5.7.1	MAP estimator of a Markov normally distributed random variable	66
5.7.2	Expressing the nonlinear least squares problem in matrix form	68
5.7.3	Covariance on MAP estimate	69
	Appendices	70

A	Linear algebra	71
A.1	Schur complement	71
A.1.1	Application to solving linear equations	71
A.1.2	Application in probability	72
B	Numerically sampling from a normal distribution	74
B.1	MATLAB's <code>randn</code> function	74
B.2	Sampling (another derivation)	76
B.3	Covariance ellipses	77
	References	77

Preface

Most of the notes in Part I are based off ECSE 509: “*Probability and Random Signals 2*” course taught at McGill University in Fall 2019 by Prof. I. Psaromiligkos [1].

Notations

Notation

S
 \mathcal{P}
 \mathcal{F}
 $P : \mathcal{P} \rightarrow [0, 1]$
 $f : \mathbb{R}^n \rightarrow \mathbb{R}$
 $F : \mathbb{R}^n \rightarrow \mathbb{R}$
 $p : S \subset \mathbb{R}^n \rightarrow \mathbb{R}$
 $\underline{(\cdot)} : S \rightarrow \mathbb{R}^n$
 \mathcal{R}_x
 \underline{x}
 \mathbf{x}
 $\mathbb{E} [\cdot]$
 $\text{Var} [\cdot]$
 j

Definition

Sample space.
 Power set.
 Events algebra.
 Probability function.
 Probability density function (PDF).
 Cumulative distribution function (CDF).
 Cumulative distribution function (CDF).
 Usually used for discrete random variables.
 A random variable (RV). S is referred to as the *domain*.
 Set of possible points that random variable \underline{x} can take.
 Single random variable.
 Multivariate random variable.
 Expectation operator.
 Variance.
 $\sqrt{-1}$.

Part I

Fundamentals

Chapter 1

Single Variable Probability

1.1 Definitions and intuition

Probability is a measure of:

- Our certainty or belief that a statement is true. My belief could be different than yours;
- How frequently an event will occur.

Definition 1.1.1 (*Random experiment*). A *random experiment* is an experiment (or a physical process) whose outcome is not certain but all of its possible outcomes are known and predictable in advance.

A random experiment is modeled as a probability space (check Definition 1.1.7).

Definition 1.1.2 (*Sample space S*). The *sample space* S is the set of *all possible outcomes*.

Definition 1.1.3 (*Event*). An *event* A is a *subset* of S . That is, $A \subseteq S$. It's NOT an element of S ! It is rather an element of the power set \mathcal{P}_S (Definition 1.1.4).

Thus, an event A is a *set* of outcomes.

Definition 1.1.4 (*Power set of S*). The *power set* of the sample space S is the set containing all subsets of a set S . It's denoted by \mathcal{P}_S , or simply \mathcal{P} .

-
- \mathcal{P}_S is the set of all events.
 - $\emptyset \subseteq S$: Impossible event.
 - $S \subseteq S$: Certain event.

Example 1.1.1 (Flipping a coin). Flipping a coin *once* has the sample space

$$S = \{H, T\}, \quad (1.1)$$

H and T are all the possible *outcomes* (NOT events) of the sample space S . The power set is thus given by

$$\mathcal{P}_S = \{\{\emptyset\}, \{H\}, \{T\}, \{H, T\}\}. \quad (1.2)$$

Note that the last element of \mathcal{P}_S , $\{H, T\}$ is the sample space and it is thus a *certain* event. \triangle

Example 1.1.2 (Number of occurrences). Relative frequency interpretation of probability (which is very intuitive). Say a random experiment gives an outcome a in an event A (i.e., $a \in A \subseteq S$) from the sample space. Repeat the experiment N times. Then, define the number of occurrences of event A after repeating the experiment N times by

$$n(A, N) = \# \text{ of times that } A \text{ happens}. \quad (1.3)$$

Note that “# of times that A happens” implies that the outcome a_k of experiment k belongs to the event A . That is, $a_k \in A$. Then,

$$\underbrace{\lim_{N \rightarrow \infty} \frac{n(A, N)}{N}}_{\text{Probability of event } A} = \text{constant}. \quad (1.4)$$

\triangle

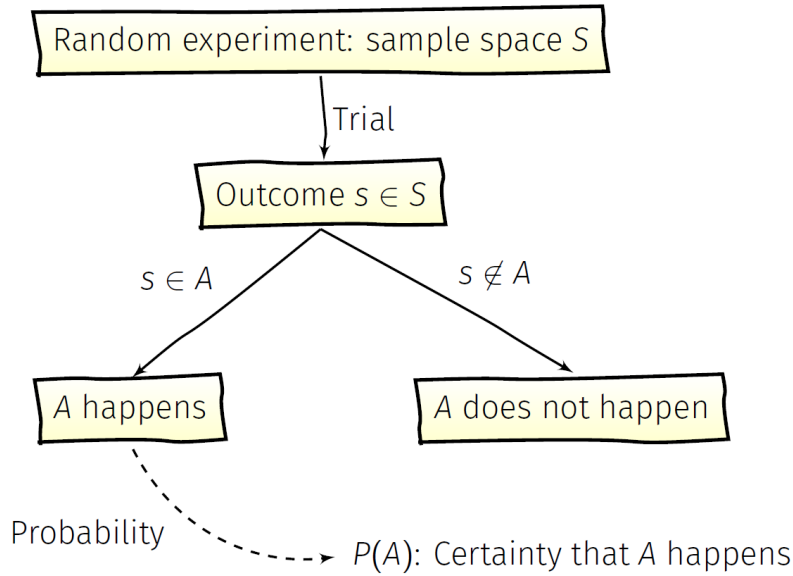


Figure 1.1: Intuition of probability. Figure obtained from [1].

Definition 1.1.5 (Probability). Probability is a set function $P : \mathcal{P}_S \mapsto [0, 1]$ that assigns to an event

$A \in \mathcal{P}_S$ (recall $A \subseteq S$) a number $P(A)$ that satisfies the following three axioms:

1. $P(A) \geq 0$,
2. $P(S) = 1$, and
3. For any sequence of mutually exclusive events A_1, A_2, \dots (i.e., for which $A_i \cap A_j = \emptyset$ for $i \neq j$), we have

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i). \quad (1.5)$$

The number $P(A)$ is called the probability of A .

Remark 1.1.1. It is NOT always possible to define a function $P : \mathcal{P} \mapsto [0, 1]$ that satisfies the three axioms. Workaround: focus on the events in a σ -algebra \mathcal{F} .

Here are some other remarks about the probability function.

- The probability function $P : \mathcal{P}_S \mapsto [0, 1]$ maps the power set \mathcal{P}_S (or σ -algebra) to $[0, 1]$. That is, the argument of P should be an *element* of \mathcal{P}_S . And since \mathcal{P}_S is a *set of sets* of S , then the argument is the event which is denoted by $\{a_1, a_2, \dots, a_n\} = A \in \mathcal{P}_S$, where $a_1, \dots, a_n \in S$. Thus the proper way to write “the probability of event A ” is $P(\{A\})$ but the notation $P(A)$ will often be used for ease of reading and writing.
- Note that $P(A) = 0$ implies that event A is *improbable*; not necessarily impossible! It’s still possible to get A . This might be a little unintuitive, but it can be thought of that the chance of A happening is so small that it seems that event A is highly unlikely to occur.
- Similarly, $P(A) = 1$ does NOT imply that A will certainly occur! It just implies that it is very probably.

Definition 1.1.6 (σ -algebra). A set \mathcal{F} of subsets of S , that is, $\mathcal{F} \subseteq \mathcal{P}$ is σ -algebra if and only if

1. $S \in \mathcal{F}$,
2. $E \in \mathcal{F} \implies E^C \in \mathcal{F}$, and
3. If the sets A_1, A_2, \dots belong to \mathcal{F} , then so does $\bigcup_{i=1}^{\infty} A_i$.

If S is countable, then $\mathcal{F} = \mathcal{P}_S$ is used.

A random experiment is modeled as a probability space.

Definition 1.1.7 (*Probability space*). A probability space is a triplet

$$(S, \mathcal{F}, P), \quad (1.6)$$

where

1. S : is the sample set,
2. \mathcal{F} is the events algebra, and
3. $P : \mathcal{F} \mapsto [0, 1]$ is the probability function.

1.2 Basic theorems

Theorem 1.2.1 (*Basic theorems*). Below is a list of basic theorems of probability.

1. For any event $A \in \mathcal{F}$,

$$P(A^C) = 1 - P(A). \quad (1.7)$$

2. For any event $A \in \mathcal{F}$,

$$0 \leq P(A) \leq 1. \quad (1.8)$$

3. $P(\emptyset) = 0$.

4. For any events $A, B \in \mathcal{F}$,

$$P(A - B) = P(A) - P(A \cap B). \quad (1.9)$$

Special case: if $B \subseteq A$, then

$$P(A - B) = P(A) - P(B). \quad (1.10)$$

5. For any $A, B \in \mathcal{F}$,

$$P(A) = P(AB) + P(AB^C). \quad (1.11)$$

6. For any $A, B \in \mathcal{F}$,

$$P(A \cap B) = P(A) + P(B) - P(A \cup B), \quad (1.12)$$

$$P(A \cup B) = P(A) + P(B) - P(A \cap B). \quad (1.13)$$

1.3 Conditional probabilities

Events can affect other events. Therefore, the question “what’s the probability of event A happening?” could be quite different from “what’s the probability of event A given that event B happened?”¹. The notation

$$P(A|B) \quad (1.14)$$

reads “probability of event A happening given that event B happened”.

Given

1. a random experiment $(S, \mathcal{F}, P())$, and
2. $A, B \in \mathcal{F}$ with $P(B) \neq 0$.

Definition 1.3.1 (Conditional probability). The conditional probability of A given B , $P(A|B)$, defined as

$$P(A|B) = \frac{P(A \cap B)}{P(B)}. \quad (1.15)$$

Theorem 1.3.1 (Conditional probability). Let $B \in \mathcal{F}$ with $P(B) \neq 0$. The mapping

$$P(\cdot|B) : \mathcal{F} \rightarrow \mathbb{R}, \quad (1.16)$$

satisfies the axioms of probability

1. $P(A|B) \geq 0, \quad \forall A \in \mathcal{F}$
2. $P(S|B) = 1$
3. If $A_1, A_2, \dots \in \mathcal{F}$ is a sequence of mutually exclusive events, then

$$P\left(\bigcup_{i=1}^{\infty} A_i | B\right) = \sum_{i=1}^{\infty} P(A_i|B). \quad (1.17)$$

Therefore,

- For a given event B with $P(B) \neq 0$, the mapping

$$P(\cdot|B) : \mathcal{F} \rightarrow \mathbb{R}, \quad (1.18)$$

defines a valid probability function.

- All the basic theorems of probability listed in Theorem 1.2.1 apply to $P(\cdot|B)$.

¹In some cases, the two questions are the same. This brings the notion of independence.

Here are some properties of conditional probability

1. For $A, B \in \mathcal{F}$ with $P(A), P(B) \neq 0$:

$$P(AB) = P(A|B)P(B) \quad (1.19)$$

$$= P(B|A)P(A). \quad (1.20)$$

2. **Total probability.** Let B_1, B_2, \dots, B_n be a partition of S , with $P(B_i) \neq 0, \forall i$:

$$P(A) = P(A|B_1)P(B_1) + P(A|B_2)P(B_2) + \dots + P(A|B_n)P(B_n) \quad (1.21)$$

$$= \sum_{i=1}^n P(A|B_i)P(B_i). \quad (1.22)$$

3. **Bayes rule.** Let B_1, B_2, \dots, B_n be a partition of S , with $P(B_i) \neq 0, \forall i$:

$$P(B_i|A) = \frac{P(A|B_i)P(B_i)}{\sum_{j=1}^n P(A|B_j)P(B_j)}. \quad (1.23)$$

Special case for $n = 1$:

$$P(B|A) = \frac{P(A|B)P(B)}{P(A)}. \quad (1.24)$$

Definition 1.3.2 (Independence). The events $A, B \in \mathcal{F}$ are called *independent* if

$$\underbrace{P(A|B) = P(A)}_{\text{requires } P(B) \neq 0} \text{ or } \underbrace{P(B|A) = P(B)}_{\text{requires } P(A) \neq 0} \text{ or } P(AB) = P(A)P(B) \quad (1.25)$$

Definition 1.3.3 (Independence (multiple events)). The events $A, B, C \in \mathcal{F}$ are called *independent* if

- They are independent in pairs, and
- $P(ABC) = P(A)P(B)P(C)$.

Definition 1.3.4 (Conditional independence). Events $A, B \in \mathcal{F}$ are called *conditionally independent* given $C \in \mathcal{F}$, with $P(C) \neq 0$, if

$$P(AB|C) = P(A|C)P(B|C). \quad (1.26)$$

Remark: Conditional independence is quite different from independence. They do not imply each other. Here are some examples.

Example 1.3.1 (Independent but conditionally dependent). Two fair coins are flipped. Define the fol-

lowing events

- A - Your first coin flip is heads,
- B - Your second coin flip is heads, and
- C - Your first two flips were the same.

Then A and B are independent. However, A and B are conditionally dependent given C , since if you know C then your first coin flip will inform the other one.

Example 1.3.2 (Dependent but conditionally independent). Consider two brothers John and Joseph, both having a genetic disease. These two events are dependent as they are brothers. However, given the condition that Joseph is an adopted son of the family makes the events independent.

1.4 Random variables and distributions

In many random experiments, it is convenient to assign numerical labels to outcomes.

Example 1.4.1. Experiment: pick a car out of a black Escort (Eb), a red Escort (Er), and a red Mazda (Mr). The sample space is then

$$S = \{Eb, Er, Mr\}, \quad (1.27)$$

and the events algebra is $\mathcal{F} = \mathcal{P}$.

Say we are interested in the make of the car and not the color. Then we can assign labels as can be seen in Figure 1.2.

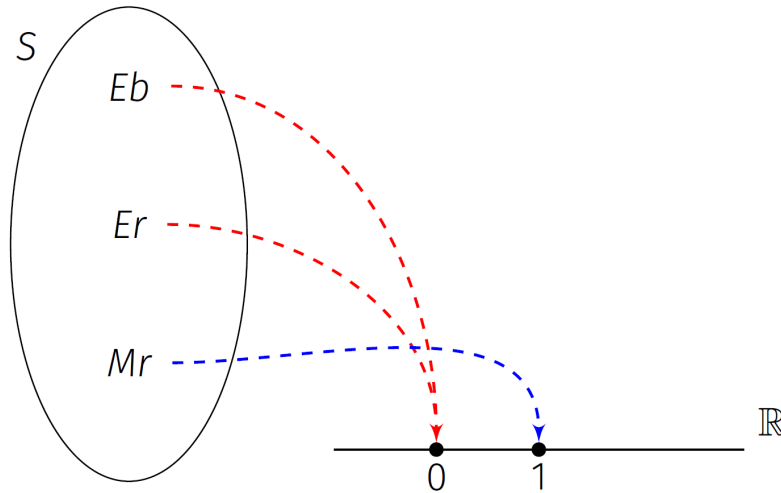


Figure 1.2: Assigning numerical labels to outcomes; each color is an outcome. From [1]

Definition 1.4.1 (Random variable (RV)). A random variable (RV) (denoted by an underline) $\underline{x} : S \rightarrow \mathbb{R}$ is a mapping from S to \mathbb{R} with the following properties.

1. $A(x) = \{s \in S : \underline{x}(s) \leq x\} \in \mathcal{F}, \forall x \in \mathbb{R}$
2. $P(\{s \in S : \underline{x}(s) = \infty\}) = 0$ and $P(\{s \in S : \underline{x}(s) = -\infty\}) = 0$

S is referred to as the *domain* of the RV $\underline{x}(\cdot)$ (since a RV is a mapping, then it has a domain and a range). [The domain can be thought of as the sample space^a.]

^aNot quite sure of this statement.

Definition 1.4.2 (*Cumulative distribution function (CDF)*). The function $F : \mathbb{R} \rightarrow [0, 1]$ defined by

$$F(x) = P(s \in S : \underline{x}(s) \leq x) \quad (1.28)$$

$$= P(\underline{x} \leq x), \quad \forall x \in \mathbb{R} \quad (1.29)$$

is called the *cumulative distribution function (CDF)* of \underline{x} . Note that (1.29) is a shorthand for (1.28).

Note: sometimes the CDF will be denoted with a subscript of the RV (e.g., $F_x(x)$). This is for clarity as to make sure that the CDF is over the RV \underline{x} . This come more in handy when there are multiple variables in hands (such as in integration).

1.4.1 Properties of random variables

Here are some properties of random variables.

1. $\lim_{x \rightarrow \infty} F_x(x) = 1$ and $\lim_{x \rightarrow -\infty} F_x(x) = 0$
2. $F_x(x)$ is non-decreasing
3. $F_x(x)$ is continuous from the right. That is

$$\lim_{x \rightarrow x_0^+} F_x(x) = F_x(x_0) \quad (1.30)$$

4. $P(\{x_1 < \underline{x} \leq x_2\}) = F_x(x_2) - F_x(x_1)$
5. $P(\{\underline{x} = x_1\}) = F_x(x_1) - F_x(x_1^-)$, where $F_x(x_1^-) := \lim_{x \rightarrow x_1^-} F_x(x)$
6. $P(\{x_1 \leq \underline{x} \leq x_2\}) = F_x(x_2) - F_x(x_1^-)$

Definition 1.4.3 (*Types of RVs*). A RV \underline{x} is called

- *Discrete* if $F_x(x)$ is constant except of countable number of discontinuities (piecewise constant). Figure 1.3 shows an example of a CDF of a discrete RV.
- *Continuous* if $F_x(x)$ is
 1. continuous, and
 2. differentiable (with the exception of a countable number of point).
- *Mixed* if it is neither continuous nor discrete.

Definition 1.4.4 (*Probability mass function (PMF)*). Let \underline{x} be a discrete RV. Then,

- $P(\underline{x} = x) = 0$ when $x \neq x_i, i = 1, 2, \dots$,
- $\mathcal{R}_x = \{x_i : i = 1, 2, \dots\}$ is the set of possible points.

The *probability mass function (PMF)* of a RV \underline{x} , $p(x)$ or $p_x(x)$, $x \in \mathbb{R}$, is defined as

$$p(x) = P(\underline{x} = x) = \begin{cases} 0 & \text{if } x \notin \mathcal{R}_x, \\ P(\underline{x} = x_i) & \text{if } x \in \mathcal{R}_x, \text{ i.e., } x = x_i, i = 1, 2, \dots \end{cases} \quad (1.31)$$

Also,

- $\sum_{i=1}^{\infty} p(x_i) = 1$ and $p(x) \geq 0, \forall x$ (defining properties),
- $P(\underline{x} \in \mathcal{A}) = \sum_{x_i \in \mathcal{A}} p(x_i)$.

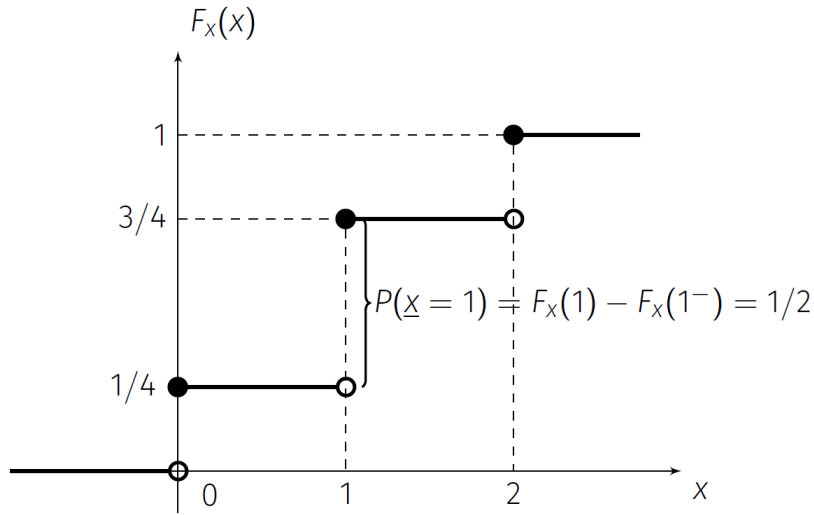


Figure 1.3: Example of a CDF of a discrete RV. From [1].

Definition 1.4.5 (*Probability density function (PDF)*). The *probability density function (PDF)*, $f(x)$ or $f_x(x)$, of a RV \underline{x} is defined as

$$f_x(x) = \frac{dF_x(x)}{dx}. \quad (1.32)$$

1.4.2 Properties of PDFs

If \underline{x} is continuous with PDF $f(x)$

1. $f_x(\underline{x}) \geq 0, \forall x$ (defining property)
2. $\int_{-\infty}^x f_x(\lambda) d\lambda = F_x(x)$
3. $\int_{-\infty}^{\infty} f_x(x) d\lambda = 1$ (defining property)

4.

$$P(\{a \leq \underline{x} \leq b\}) = \int_a^b f_x(\lambda) d\lambda \quad (1.33)$$

$$P(\{a < \underline{x} \leq b\}) = P(\{a \leq \underline{x} < b\}) = P(\{a < \underline{x} < b\}). \quad (1.34)$$

Interpretation of PDF: For a continuous \underline{x} : $P(\underline{x} = x) = 0$ for all $x \in \mathbb{R}$ which implies that the PDF $f(x)$ is *not* a probability.

For small $\epsilon > 0$:

$$P(|\underline{x} - x| < \epsilon) = P(x - \epsilon < \underline{x} < x + \epsilon) \quad (1.35)$$

$$= \int_{x-\epsilon}^{x+\epsilon} f(t) dt \quad (1.36)$$

$$\approx 2\epsilon f(x) \quad (1.37)$$

assuming $f(x)$ is continuous at x . Equivalently,

$$f(x) \approx \frac{1}{2\epsilon} P(x - \epsilon < \underline{x} < x + \epsilon). \quad (1.38)$$

$f(x)$ is proportional to the probability that \underline{x} lies in a small neighbourhood of x .

1.4.3 Conditional distributions

Definition 1.4.6 (*Conditional distributions*). Let $M \in \mathcal{F}$ with $P(M) \neq 0$. Then

1. *Conditional CDF* of \underline{x} given event M is

$$F_x(x|M) = P(\underline{x} \leq x|M) \quad (1.39)$$

$$= \frac{P(\underline{x} \leq x, M)}{P(M)}. \quad (1.40)$$

2. *Conditional PDF* of \underline{x} given event M is

$$f_x(x|M) = \frac{dF_x(x|M)}{dx}. \quad (1.41)$$

3. *Conditional PMF* of \underline{x} (\underline{x} discrete given M) is

$$p(x|M) = \begin{cases} 0, & x \notin \{x_i : i = 1, 2, \dots\} \\ P(\underline{x} = x_i|M), & x \in \{x_i : i = 1, 2, \dots\}. \end{cases} \quad (1.42)$$

Note that conditional CDFs, PDFs, and PMFs behave like unconditional ones.

Properties of conditional distributions

Let B_1, B_2, \dots, B_n be a partition of S , $P(B_i) \neq 0$. Let $A \in \mathcal{F}$, $P(A) \neq 0$.

- **Total probability:**

$$F_x(x) = \sum_{i=1}^n F_x(x|B_i)P(B_i), \quad (1.43)$$

$$f_x(x) = \sum_{i=1}^n f_x(x|B_i)P(B_i), \quad (1.44)$$

$$p_x(x) = \sum_{i=1}^n p_x(x|B_i)P(B_i). \quad (1.45)$$

- **Bayes rule:**

$$P(A|\underline{x} \leq x) = \frac{P(\underline{x} \leq x, A)}{P(\underline{x} \leq x)} \quad (1.46)$$

$$= \frac{P(\underline{x} \leq x|A)P(A)}{P(\underline{x} \leq x)} \quad (1.47)$$

$$= \frac{F_x(\underline{x} \leq x|A)P(A)}{F_x(\underline{x} \leq x)}. \quad (1.48)$$

When \underline{x} is discrete with $\mathcal{R}_x = \{x_1, x_2, \dots\}$, $p(x_i) \neq 0$:

- **Conditional probability** using PMF:

$$P(A|\underline{x} = x_i) = \frac{P(A, \underline{x} = x_i)}{P(\underline{x} = x_i)} \quad (1.49)$$

$$= \frac{P(A, \underline{x} = x_i)}{p_x(\underline{x} = x_i)}. \quad (1.50)$$

- **Total probability:**

$$P(A) = \sum_{i=1}^{\infty} P(A|\underline{x} = x_i) p(x_i). \quad (1.51)$$

- **Bayes rule:**

$$p(x|A) = \frac{P(A|\underline{x} = x) p(x)}{P(A)} \quad (1.52)$$

$$= \frac{P(A|\underline{x} = x) p(x)}{\sum_{i=1}^{\infty} P(A|\underline{x} = x_i) P(x_i)}. \quad (1.53)$$

When \underline{x} is continuous, it is not possible to divide by $P(\underline{x} = x) = 0$. Therefore, another definition for conditional PDF is needed.

Definition 1.4.7 (Conditional PDF). The *conditional PDF* of $A \in \mathcal{F}$ given $\underline{x} = x$ (assuming $f_x(x) \neq 0$) is

$$P(A|\underline{x} = x) = \frac{f_x(x|A)P(A)}{f_x(x)}. \quad (1.54)$$

The properties of such conditional PDF are

- **Total probability:** $P(A) = \int_{-\infty}^{\infty} P(A|\underline{x} = x) f_x(x) dx$
- **Bayes rule:**

$$f_x(x|A) = \frac{P(A|\underline{x} = x) f_x(x)}{P(A)} \quad (1.55)$$

$$= \frac{P(A|\underline{x} = x) f_x(x)}{\int_{-\infty}^{\infty} P(A|\underline{x} = x) f_x(x) dx}. \quad (1.56)$$

1.5 Transformed random variables

Say there exists a random variable $\underline{x} : S \rightarrow \mathbb{R}$ that is applied on the probability space (S, \mathcal{F}, P) . Further, say there exists a [deterministic] function $g : \mathbb{R} \rightarrow \mathbb{R}$. Then, there is a new RV \underline{y} defined as

$$\underline{y}(s) = g(\underline{x}(s)) \quad \forall s \in S. \quad (1.57)$$

Furthermore, say the distribution (i.e., CDF and PDF (or PMF)) of \underline{x} is known. Then what is the distribution of $\underline{y} = g(\underline{x})$?

The foundation is the following

$$P(\underline{y} \in A) = P(\{s \in S : \underline{y}(s) \in A\}) \quad (1.58)$$

$$= P(\{s \in S : g(\underline{x}(s)) \in A\}). \quad (1.59)$$

There are different methods to find the distribution of a transformed random variable (i.e., $\underline{y} = g(\underline{x})$) depending on the types of variables.

1. **Method of distributions** works when for *any* types of random variables (continuous or discrete). However, it is the most exhaustive method.
2. **Method of transformations** works for *continuous* random variables; i.e., both \underline{x} and \underline{y} are continuous.
3. **Discrete transformation** for *discrete* transformed variable \underline{y} and *continuous* or *discrete* domain random variable \underline{x} .

1.5.1 Method of distribution

The method will be motivated by an example.

Example 1.5.1. Consider the transformed random variable

$$\underline{y} = g(\underline{x}) \quad (1.60)$$

$$= \underline{x}^2. \quad (1.61)$$

Evaluate the CDF as follows

$$F_y(y) = P(\underline{y} \leq y) = P(g(\underline{x}) \leq y) = P(\underline{x}^2 \leq y) \quad (1.62)$$

$$= P(\{s \in S : \underline{x}(s)^2 \leq y\}). \quad (1.63)$$

There are three cases

1. $y < 0$: $P(\underline{x}^2 \leq y) = 0 \implies F_y(y) = 0$.
2. $y = 0$: $P(\underline{x}^2 \leq 0) = P(\underline{x}^2 = 0) = P(\underline{x} = 0) \implies F_y(0) = P(\underline{x} = 0)$.
3. $y > 0$: $P(\underline{x}^2 \leq y) = P(|\underline{x}| \leq \sqrt{y}) = P(-\sqrt{y} \leq \underline{x} \leq \sqrt{y}) \implies$

$$F_y(y) = \int_{-\sqrt{y}}^{\sqrt{y}} f_x(x) dx. \quad (1.64)$$

1.5.2 Method of transformation

The process is outlined by the following steps.

For any y (treat it as a deterministic variable):

1. Solve $g(x) = y$ for x . List *all* possible solutions x_1, x_2, \dots, x_n .
2. Calculate $g'(x) = \frac{dg(x)}{dx}$.
3. Evaluate

$$f_y(y) = \begin{cases} 0, & \text{no solutions to } g(x) = y, \\ \sum_{i=1}^n \frac{f_x(x_i)}{|g'(x_i)|}, & \text{otherwise.} \end{cases} \quad (1.65)$$

The method is derived in [2, Sec. 2.6.2].

Example 1.5.2 (Method of transformation). Consider

$$\underline{y} = g(\underline{x}) \quad (1.66)$$

$$= \underline{x}^2. \quad (1.67)$$

Carry the steps:

1. Solve for x in $y = g(x)$ for all possible cases of y :
 - 1.1. $y < 0$: No solutions.

1.2. $y = 0$: One solution $x_1 = 0$.

1.3. $y > 0$: Two solutions $x_1 = -\sqrt{y}$ and $x_2 = \sqrt{y}$.

2. Calculate the derivative of the function

$$g'(x) = 2x, \quad \forall x. \quad (1.68)$$

3. Calculate $f_y(y)$:

3.1. $y < 0$: $f_y(y) = 0$.

3.2. $y = 0$: $f_y(0) = 0$.

3.3. $y > 0$:

$$f_y(y) = \frac{f_x(-\sqrt{y})}{|-2\sqrt{y}|} + \frac{f_x(\sqrt{y})}{|2\sqrt{y}|} \quad (1.69)$$

$$= \frac{1}{2\sqrt{y}} [f_x(-\sqrt{y}) + f_x(\sqrt{y})]. \quad (1.70)$$

△

1.5.3 Discrete transformation

The method will be motivated by two examples:

1. continuous to discrete transformation, and
2. discrete to discrete transformation.

Example 1.5.3 (Continuous to discrete transformation). Let \underline{x} be continuous with $f(x) = 0, x < 0$ and

$$\underline{y} = g(\underline{x}) \quad (1.71)$$

$$= \lfloor \underline{x} \rfloor \quad (\text{floor of } \underline{x}). \quad (1.72)$$

Then the transformed variable \underline{y} is discrete. I.e., $\mathcal{R}_y = \{0, 1, 2, \dots\}$.

The PMF of \underline{y} is then given by

$$p_y(y) = P(\underline{y} = y) \quad (1.73)$$

$$= P(\lfloor \underline{x} \rfloor = y) \quad (1.74)$$

$$= P(y \leq \underline{x} < y + 1) \quad (1.75)$$

$$= \int_y^{y+1} f_x(x) dx, \quad y = 0, 1, 2, \dots \quad (1.76)$$

△

Example 1.5.4 (Discrete to discrete transformation). Let \underline{x} be discrete with $\mathcal{R}_x = \mathbb{Z}$, and

$$\underline{y} = g(\underline{x}) \quad (1.77)$$

$$= |\underline{x}|. \quad (1.78)$$

Then \underline{y} is also discrete random variable with $\mathcal{R}_y = \{0, 1, 2, \dots\}$.

The PMF of \underline{y} can then be given by

$$p_y(y) = P(\underline{y} = y) \quad (1.79)$$

$$= P(|\underline{x}| = y) \quad (1.80)$$

$$= P(\underline{x} = \pm y) \quad (1.81)$$

$$= \begin{cases} p_x(0), & y = 0, \\ p_x(-y) + p_x(y), & y = 1, 2, \dots \end{cases} \quad (1.82)$$

△

1.6 Statistics of random variables

Sometimes, information about the random variables such as the mean is sought after. Information about a random variables are referred to as the *statistics* of it.

Definition 1.6.1 (*Expectation operator*). The expectation operator $\mathbb{E}[\cdot]$ is a functional operator that operates on random variables. Specifically:

- For continuous random variables

$$\mathbb{E}[g(\underline{x})] = \int_{-\infty}^{\infty} g(x) f(x) dx. \quad (1.83)$$

- For discrete random variables

$$\mathbb{E}[g(\underline{x})] = \sum_{-\infty}^{\infty} g(x) p(x). \quad (1.84)$$

Definition 1.6.2 (*Mean of a random variable*). The mean of a random variable \underline{x} , μ_x , is defined as

$$\mu_x = \mathbb{E}[\underline{x}] \quad (1.85)$$

which

- for continuous random variable

$$\mu_x = \mathbb{E}[\underline{x}] \quad (1.86)$$

$$= \int_{-\infty}^{\infty} x f(x) dx, \quad (1.87)$$

- while for discrete random variable

$$\mu_x = \mathbb{E}[\underline{x}] \quad (1.88)$$

$$= \sum_{-\infty}^{\infty} x p(x). \quad (1.89)$$

and the conditional mean of \underline{x} given event M is

$$\mathbb{E}[\underline{x}|M] = \int_{-\infty}^{\infty} x f(x|M) dx. \quad (1.90)$$

Properties of the mean

- Mean of $\underline{y} = g(\underline{x})$ is

$$\mathbb{E}[\underline{y}] = \mathbb{E}[g(\underline{x})] \quad (1.91)$$

$$= \int_{-\infty}^{\infty} g(x) f_x(x) dx. \quad (1.92)$$

- Linear property:

$$\mathbb{E}\left[\sum_{i=1}^n \alpha_i g_i(\underline{x})\right] = \sum_{i=1}^n \alpha_i \mathbb{E}[g_i(\underline{x})]. \quad (1.93)$$

Note that in general,

$$\mathbb{E}[g(\underline{x})] \neq g(\mathbb{E}[\underline{x}]). \quad (1.94)$$

Definition 1.6.3 (*Variance and standard deviation*). Let \underline{x} be a random variable. Then the *variance* of \underline{x} , $\text{Var}[\underline{x}]$, is defined as

$$\text{Var}[\underline{x}] = \sigma_x^2 \quad (1.95)$$

$$= \mathbb{E}\left[(\underline{x} - \mathbb{E}[\underline{x}])^2\right] \quad (1.96)$$

$$= \mathbb{E}[\underline{x}^2] - \mathbb{E}[\underline{x}]^2. \quad (1.97)$$

The *standard deviation* of \underline{x} is $\sigma_x = \sqrt{\text{Var}[\underline{x}]}$.

A useful property of the variance is

$$\text{Var} [a\underline{x} + b] = a^2 \text{Var} [\underline{x}]. \quad (1.98)$$

Theorem 1.6.1 (Markov's inequality). Let \underline{x} be a non-negative random variable. That is, $P(\underline{x} < 0) = 0$. For any $\epsilon > 0$:

$$P(\underline{x} \geq \epsilon) \leq \frac{\mathbb{E}[\underline{x}]}{\epsilon}. \quad (1.99)$$

Theorem 1.6.2 (Chebyshev's inequality). Let \underline{x} be a random variable with mean μ and variance σ^2 . Then, for any $\epsilon > 0$:

$$P(|\underline{x} - \mu| \geq \epsilon) \leq \frac{\sigma^2}{\epsilon^2}. \quad (1.100)$$

Theorem 1.6.3 (Jensen's inequality). Let $g(\cdot)$ be a convex function. Then,

$$\mathbb{E}[g(\underline{x})] \geq g(\mathbb{E}[\underline{x}]). \quad (1.101)$$

1.6.1 Moments

Definition 1.6.4 (*Moments*). The different moments of a random variable are defined as follows.

- n -th moment of \underline{x} : $m_n = \mathbb{E}[\underline{x}^n]$.
- n -th central moment of \underline{x} : $\mu_n = \mathbb{E}[(\underline{x} - \mathbb{E}[\underline{x}])^n]$.
- n -th absolute moment of \underline{x} : $\mathbb{E}[|\underline{x}|^n]$.
- n -th absolute central moment of \underline{x} : $\mu_n = \mathbb{E}[|\underline{x} - \mathbb{E}[\underline{x}]|^n]$.

1.6.2 Characteristic functions

Definition 1.6.5 (*Characteristic function*). The *characteristic function* of a random variable \underline{x} is the Fourier transform of its PDF. Specifically,

$$\Phi_x(\omega) = \int_{-\infty}^{\infty} e^{j\omega x} f_x(x) dx, \quad (1.102)$$

where the inverse expression is given by

$$f_x(x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{-j\omega x} \Phi_x(\omega) d\omega. \quad (1.103)$$

Definition 1.6.6 (*Moment generating function*). The *moment generating function* of \underline{x} is the Laplace transform of its PDF. Specifically,

$$\bar{\Phi}_x(s) = \int_{-\infty}^{\infty} e^{Sx} f_x(x) dx. \quad (1.104)$$

Theorem 1.6.4 (Moment theorem).

$$\mathbb{E}[\underline{x}^n] = \left. \frac{d^n}{ds^n} \bar{\Phi}_x \right|_{s=0} = \frac{1}{j^n} \left. \frac{d^n}{d\omega^n} \Phi_x(\omega) \right|_{\omega=0}. \quad (1.105)$$

Chapter 2

Multivariate Probability

2.1 Multivariate random variables and distributions

In this chapter, we are dealing with multiple random variables.

Definition 2.1.1 (*Multiple random variable*). Let $\underline{x}_1, \underline{x}_2, \dots, \underline{x}_n$ be n random variables defined on the same probability space (S, \mathcal{F}, P) . Then, the mapping

$$s \in S \rightarrow [\underline{x}_1(s), \underline{x}_2(s), \dots, \underline{x}_n(s)]^T \in \mathbb{R}^n \quad (2.1)$$

defines an n -dimensional random variable (random vector or column matrix). Random vectors will be denoted by boldface underlined alphabet. Specifically,

$$\underline{\mathbf{x}} = \begin{bmatrix} \underline{x}_1 & \underline{x}_2 & \dots & \underline{x}_n \end{bmatrix}^T, \quad (2.2)$$

where $\underline{x}_i, i = 1, \dots, n$ are random variables.

2.1.1 Joint cumulative distribution function (JCDF)

Definition 2.1.2 (*Joint cumulative distribution function (JCDF)*). Let $\underline{x}_1, \underline{x}_2, \dots, \underline{x}_n$ be n random variables defined on the same probability space (S, \mathcal{F}, P) . Then, the *joint cumulative distribution function (JCDF)* of these random variables are

$$F_{x_1, x_2, \dots, x_n}(x_1, x_2, \dots, x_n) = P(\underline{x}_1 \leq x_1, \underline{x}_2 \leq x_2, \dots, \underline{x}_n \leq x_n) \quad (2.3)$$

$$= F_{\underline{\mathbf{x}}}(\underline{\mathbf{x}}) \quad (2.4)$$

$$= F(\underline{\mathbf{x}}) \quad (2.5)$$

$$= P(\underline{\mathbf{x}} \leq \underline{\mathbf{x}}). \quad (2.6)$$

Properties of JCDFs

1. $F(\underline{\mathbf{x}})$ is non-decreasing in each of its arguments.

2. $F(\mathbf{x})$ is right-continuous in each of its arguments.
3. For a fixed i , $\lim_{x_i \rightarrow -\infty} F(\mathbf{x}) = 0$.
4. When $\underline{x}_i \rightarrow \infty$ for all i , then $F(\mathbf{x}) \rightarrow 1$.

Definition 2.1.3 (*Marginal CDF*). Let $\underline{\mathbf{x}} = \begin{bmatrix} \underline{\mathbf{x}}_1^T & \underline{\mathbf{x}}_2^T \end{bmatrix}^T \in \mathbb{R}^n$ with joint CDF $F(\mathbf{x})$. Say we want to know the joint CDF of a subset of these random variables $\underline{\mathbf{x}}_1$. That is, we want to *marginalize out* the random variables $\underline{\mathbf{x}}_2$. The joint CDF of $\underline{\mathbf{x}}_1$ is computed by taking the limit

$$F(\mathbf{x}_1) = \lim_{\mathbf{x}_2 \rightarrow \infty} F(\mathbf{x}). \quad (2.7)$$

Here are the definitions:

- The process is called *marginalization*.
- The retained random variables $\underline{\mathbf{x}}_1$ are called *marginal random variables*.
- The joint CDF of the marginal random variables is called the *marginal CDF*.
- The discarded random variables $\underline{\mathbf{x}}_2$ is said to be *marginalized out*.

Example 2.1.1 (Marginalizing a CDF). Given the joint CDF $F(x_1, x_2, x_3)$, then the joint CDF $F(x_1, x_3)$ is given by

$$F(x_1, x_3) = \lim_{x_2 \rightarrow \infty} F(x_1, x_2, x_3). \quad (2.8)$$

△

2.1.2 Joint probability mass function (JPMF)

Definition 2.1.4 (*Joint probability mass function (JPMF)*). This is the analog to joint CDF for discrete random variables. The *joint probability mass function (JPMF)* of the discrete random variables $\underline{x}_1, \dots, \underline{x}_n$ is

$$p_{x_1, \dots, x_n}(x_1, \dots, x_n) = P(\underline{x}_1 = x_1, \dots, \underline{x}_n = x_n) \quad (2.9)$$

$$= p_{\mathbf{x}}(\mathbf{x}) \quad (2.10)$$

$$= P(\underline{\mathbf{x}} = \mathbf{x}). \quad (2.11)$$

Note that the subscripts may be omitted. That is, $p_{\mathbf{x}}(\mathbf{x}) = p(\mathbf{x})$.

Properties of JPMFs

1. **Total probability:** $0 \leq p(\mathbf{x}) \leq 1$.

2. **Normalization Property:**

$$\sum_{\mathbf{x} \in \mathcal{R}_{\mathbf{x}}} p(\mathbf{x}) = \sum_{x_1 \in \mathcal{R}_{x_1}} \dots \sum_{x_n \in \mathcal{R}_{x_n}} p(x_1, \dots, x_n) \quad (2.12)$$

$$= 1. \quad (2.13)$$

3. Probability that $\underline{\mathbf{x}}$ is in a region \mathcal{D} (i.e., $\underline{\mathbf{x}} \in \mathcal{D}$) is given by

$$P(\underline{\mathbf{x}} \in \mathcal{D}) = \sum_{\underline{\mathbf{x}} \in \mathcal{D}} p(\underline{\mathbf{x}}). \quad (2.14)$$

4. **Marginalization property:** Let $\underline{\mathbf{x}} = \begin{bmatrix} \underline{\mathbf{x}}_1^\top & \underline{\mathbf{x}}_2^\top \end{bmatrix}^\top$ be a discrete random variable with joint PMF $p(\mathbf{x})$. Then the joint PMF of $\underline{\mathbf{x}}_1$ is given by marginalization. Specifically,

$$p(\mathbf{x}_1) = \sum_{\mathbf{x}_2 \in \mathcal{R}_{\mathbf{x}_2}} p(\mathbf{x}). \quad (2.15)$$

2.1.3 Joint probability density function (JPDF)

Definition 2.1.5 (*Joint probability density function (JPDF)*). The *joint PDF* of random variables $\underline{\mathbf{x}} = \begin{bmatrix} x_1 & \dots & x_n \end{bmatrix}^\top$, $f(\mathbf{x})$, is given by

$$f(\mathbf{x}) = f(x_1, \dots, x_n) \quad (2.16)$$

$$= \frac{\partial^n F(x_1, \dots, x_n)}{\partial x_1 \dots \partial x_n}. \quad (2.17)$$

Often, subscripts will be dropped (i.e., $f_{\mathbf{x}}(\mathbf{x}) = f(\mathbf{x})$).

Properties of JPDEs

1. $f(\mathbf{x}) \geq 0$.

2. **Normalization:**

$$\int_{\mathbb{R}^n} f(\mathbf{x}) = \int_{\mathbb{R}^n} f(x_1, \dots, x_n) dx_1 \dots dx_n \quad (2.18)$$

$$= 1. \quad (2.19)$$

3. Probability that $\underline{\mathbf{x}}$ is in a region \mathcal{D} :

$$P(\underline{\mathbf{x}} \in \mathcal{D}) = \int_{\mathcal{D}} f(\mathbf{x}) d\mathbf{x}. \quad (2.20)$$

4. **Marginalization property:** Let $\underline{\mathbf{x}} = \begin{bmatrix} \underline{\mathbf{x}}_1^\top & \underline{\mathbf{x}}_n^\top \end{bmatrix}^\top$ be a continuous random variable with joint PDF

$f(\mathbf{x})$. Then the joint PDF of $\underline{\mathbf{x}}_1$ is

$$f(\mathbf{x}_1) = \int_{-\infty}^{\infty} f(\mathbf{x}) d\mathbf{x}. \quad (2.21)$$

2.1.4 Independence

Definition 2.1.6 (*Mutual independence*). The random variables $\underline{x}_1, \dots, \underline{x}_n$ are called *mutually independent* if the events $\{\underline{x}_1 \leq x_1\}, \dots, \{\underline{x}_n \leq x_n\}$ are independent for any $x_1, \dots, x_n \in \mathbb{R}^n$.

It follows that

$$F(x_1, \dots, x_n) = F(x_1) \cdots F(x_n), \quad (2.22)$$

$$f(x_1, \dots, x_n) = f(x_1) \cdots f(x_n), \quad (2.23)$$

$$p(x_1, \dots, x_n) = p(x_1) \cdots p(x_n). \quad (2.24)$$

Remark 2.1.1. Here some remarks about mutual independence.

- Any subset of the mutually independent random variables $\{\underline{x}_1, \dots, \underline{x}_n\}$ is a set of mutually independent random variables.
- If $\{\underline{x}_1, \dots, \underline{x}_n\}$ are independent *in pairs*, they are not necessarily independent.
- The random variables $\underline{y}_1 = g_1(\underline{x}_1), \dots, \underline{y}_n = g_n(\underline{x}_n)$ are independent if $\{\underline{x}_1, \dots, \underline{x}_n\}$ are independent. Note that the functions g_1, \dots, g_n need not be the same.
- Care must be taken when talking about independence. For instance, “independence” of a group of random variables may signify pairwise independence which does not imply mutual independence.

Definition 2.1.7 (*Group independence*). A group of random variables $\mathcal{G}_{\mathbf{x}} = \{\underline{x}_1, \dots, \underline{x}_n\}$ is *independent* of a group of random variables $\mathcal{G}_{\mathbf{y}} = \{\underline{y}_1, \dots, \underline{y}_m\}$ if

$$f(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) f(\mathbf{y}). \quad (2.25)$$

Definition 2.1.8 (*Independent and identically distributed (i.i.d.) random variables*). The random variables $\underline{x}_1, \dots, \underline{x}_n$ are called *independent and identically distributed (i.i.d.)* if

1. they are independent, and
2. $F_{x_1}(x) = F_{x_2}(x) = \cdots = F_{x_n}(x)$, or equivalently

$$f_{x_1}(x) = \cdots = f_{x_n}(x), \quad \text{or} \quad (2.26)$$

$$p_{x_1}(x) = \cdots = p_{x_n}(x). \quad (2.27)$$

2.1.4.1 Complex random variables

The statistics (PDF, CDF, etc.) of the n complex random variables

$$\begin{aligned} \underline{z}_1 &= \underline{x}_1 + j\underline{y}_1 \\ &\vdots \\ \underline{z}_n &= \underline{x}_n + j\underline{y}_n \end{aligned} \quad (2.28)$$

are determined by the joint PDF $f(x_1, \dots, x_n, y_1, \dots, y_n)$ of the $2n$ real random variables \mathbf{x}, \mathbf{y} .

Definition 2.1.9 (*Mutually independent complex random variables*). The complex random variables $\underline{z}_1, \dots, \underline{z}_n$ as defined in (2.28) are *mutually independent* if

$$f(\mathbf{x}, \mathbf{y}) = f(x_1, y_1) \cdots f(x_n, y_n). \quad (2.29)$$

2.2 Transformed random variables

Consider n random variables $\underline{x}_1, \underline{x}_2, \dots, \underline{x}_n$ with joint PDF $f(x_1, \dots, x_n)$. Further, given a function $g : \mathbb{R}^n \rightarrow \mathbb{R}^k$

$$g(\mathbf{x}) = \begin{bmatrix} g_1(\mathbf{x}) \\ \vdots \\ g_k(\mathbf{x}) \end{bmatrix}. \quad (2.30)$$

Then there are k new random variables $\underline{\mathbf{y}} = \begin{bmatrix} \underline{y}_1 & \cdots & \underline{y}_k \end{bmatrix}^T$. Specifically,

$$\underline{\mathbf{y}} = g(\mathbf{x}). \quad (2.31)$$

Basic concepts

If $\underline{\mathbf{x}}$ is continuous, then

$$F(\mathbf{y}) = F(y_1, \dots, y_k) \quad (2.32)$$

$$= P(\underline{y}_1 \leq y_1, \dots, \underline{y}_k \leq y_k) \quad (2.33)$$

$$= P(g_1(\mathbf{x}) \leq y_1, \dots, g_k(\mathbf{x}) \leq y_k) \quad (2.34)$$

$$= \int_{\mathcal{D}} f(\mathbf{x}) d\mathbf{x}, \quad (2.35)$$

where

$$\mathcal{D} = \{\mathbf{x} \in \mathbb{R}^n \text{ s.t. } g_1(\mathbf{x}) \leq y_1, \dots, g_k(\mathbf{x}) \leq y_k\} \quad (2.36)$$

$$= \{x_1, \dots, x_n \text{ s.t. } g_1(\cdot) \leq y_1, \dots, g_k(\cdot) \leq y_k\}. \quad (2.37)$$

That is, \mathcal{D} depends on \mathbf{y} .

If $\underline{\mathbf{x}}$ is discrete, then

$$p(\mathbf{y}) = P(y_1, \dots, y_k) \quad (2.38)$$

$$= \sum_{\mathbf{x} \in \mathcal{D}} P(x_1, \dots, x_n) \quad (2.39)$$

where

$$\mathcal{D} = \{\mathbf{x} \in \mathbb{R}^n \text{ s.t. } g_1(\mathbf{x}) = y_1, \dots, g_k(\mathbf{x}) = y_k\} \quad (2.40)$$

$$= \{x_1, \dots, x_n \text{ s.t. } g_1(\cdot) = y_1, \dots, g_k(\cdot) = y_k\}. \quad (2.41)$$

Again, \mathcal{D} depends on \mathbf{y} .

Transformed variables distribution

Given the distribution of $\underline{\mathbf{x}}$ (CDF, PDF, or PMF) and the function $g : \mathbb{R}^k \rightarrow \mathbb{R}^n$, then what is the distribution of $\underline{\mathbf{y}} = g(\underline{\mathbf{x}}) \in \mathbb{R}^k$?

Assumption: $k \leq n$. That is, the mapping g is surjective.

Just as is the case for a single transformed random variable (Section 1.5), there are different methods to find the distribution of a transformed random variable (i.e., $\underline{\mathbf{y}} = g(\underline{\mathbf{x}})$) depending on the types of variables.

1. **Method of distributions** works for *any* type of random variables (continuous or discrete). However, it is the most exhaustive method.
2. **Method of transformations** works for *continuous* random variables. Specifically, both $\underline{\mathbf{x}}$ and $\underline{\mathbf{y}}$ are continuous.

2.2.1 Method of distribution

Given the function $g : \mathbb{R}^n \rightarrow \mathbb{R}^k$, where $k \leq n$, then carry the following.

1. For all $\mathbf{y} \in \mathbb{R}^k$, find $\mathbf{x} \in D_y \subseteq \mathbb{R}^n$ such that

$$g(\mathbf{x}) \leq \mathbf{y}. \quad (2.42)$$

Note that D_y is a function of \mathbf{y} .

2. The CDF of $\underline{\mathbf{y}}$ is then

$$F(\mathbf{y}) = P(\underline{\mathbf{y}} \leq \mathbf{y}) \quad (2.43)$$

$$= P(\underline{\mathbf{x}} \in D_y) \quad (2.44)$$

$$= \int_{D_y} f(\mathbf{x}) d\mathbf{x}. \quad (2.45)$$

3. The PDF of \underline{y} is then

$$f(\underline{y}) = \frac{dF(\underline{y})}{d\underline{y}}. \quad (2.46)$$

Example 2.2.1. Consider the transformed random variable

$$\underline{y} = g(\underline{x}_1, \underline{x}_2) \quad (2.47)$$

$$= \underline{x}_1 + \underline{x}_2. \quad (2.48)$$

To find the PDF of \underline{y} , the steps outlined above will be follows.

1. Find D_y :

$$D_y = \{(x_1, x_2) \in \mathbb{R}^2 : x_2 \leq y - x_1\}. \quad (2.49)$$

2. CDF of \underline{y} :

$$F(y) = \int \int_{D_y} f(x_1, x_2) dx_1 dx_2 \quad (2.50)$$

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{y-x_1} f(x_1, x_2) dx_2 dx_1. \quad (2.51)$$

3. PDF of \underline{y} is

$$f(y) = \frac{dF(y)}{dy} \quad (2.52)$$

$$= \int_{-\infty}^{\infty} \frac{d}{dy} \int_{-\infty}^{y-x_1} f(x_1, x_2) dx_2 dx_1 \quad (2.53)$$

$$= \int_{-\infty}^{\infty} f(x_1, y - x_1) dx_1. \quad (2.54)$$

△

Theorem 2.2.1 (Transformed variable PDF of linear combination of random variables). Let \underline{x}_1 and \underline{x}_2 be jointly distributed random variables with joint PDF $f(x_1, x_2)$. The PDF of $\underline{y} = \underline{x}_1 + \underline{x}_2$ is given by

$$f(y) = \int_{-\infty}^{\infty} f(\lambda, y - \lambda) d\lambda. \quad (2.55)$$

Theorem 2.2.2 (Transformed variable PDF of linear combination of independent random variables). Let \underline{x}_1 and \underline{x}_2 be *independent* random variables with marginal PDFs $f_1(x_1)$ and $f_2(x_2)$, respectively. Then,

the PDF of $\underline{y} = \underline{x}_1 + \underline{x}_2$ is given by

$$f(y) = \int_{-\infty}^{\infty} f_1(\lambda) f_2(y - \lambda) d\lambda \quad (2.56)$$

$$= f_1(y) * f_2(y), \quad (2.57)$$

where $*$ denotes the convolution operator.

2.2.2 Method of transformation

Note that the method is derived in [2, Sec. 2.6.2].

2.2.2.1 Approach 1

Given the transformed random variable

$$\underline{y} = g(\underline{x}), \quad (2.58)$$

where $\underline{y} \in \mathbb{R}^k$ (i.g., $g : \mathbb{R}^n \rightarrow \mathbb{R}^k$) and given the PDF of the random variable \underline{x} . For now, **assume** $k = n$. The distribution of \underline{y} is then found as follows.

1. For all \underline{y} , solve the system of equations

$$\underline{y} = g(\underline{x}) \quad (2.59)$$

for \underline{x} . For a given \underline{y} , the system has either

- 1.1. m solutions

$$\underline{x}^{(1)}, \dots, \underline{x}^{(m)}, \quad (2.60)$$

- 1.2. or no solutions.

Note that the system will not have infinite solutions since it is assumed that $k = n$.

2. Evaluate the Jacobian of g

$$J(\underline{x}) = \frac{dg(\underline{x})}{d\underline{x}}. \quad (2.61)$$

3. Then, the PDF of \underline{y} is

$$f(\underline{y}) = \sum_{i=1}^m \frac{f(\underline{x}^{(i)})}{|\det(J(\underline{x}^{(i)}))|}, \quad (2.62)$$

where $\det(\cdot)$ is the determinant of a square matrix and $|\cdot|$ is the absolute value. **If for \underline{y} the system**

has no solution, then

$$f(\mathbf{y}) = \mathbf{0}. \quad (2.63)$$

2.2.2.2 Approach 2

Given the transformed random variable

$$\underline{\mathbf{y}} = g(\underline{\mathbf{x}}), \quad (2.64)$$

where $\underline{\mathbf{y}} \in \mathbb{R}^k$ (i.g., $g : \mathbb{R}^n \rightarrow \mathbb{R}^k$) and given the PDF of the random variable \mathbf{x} . For now, **assume** $k = n$. The distribution of $\underline{\mathbf{y}}$ is then found as follows.

1. Identical to in Step 1 of approach 1 (Section 2.2.2.1): for all \mathbf{y} , solve the system of equations

$$\mathbf{y} = g(\mathbf{x}) \quad (2.65)$$

for \mathbf{x} . For a given \mathbf{y} , the system has either

- 1.1. m solutions

$$\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}, \quad (2.66)$$

- 1.2. or no solutions.

Note that the system will not have infinite solutions since it is assumed that $k = n$.

2. Evaluate the Jacobian

$$\tilde{J}_i(\mathbf{y}) = \frac{d\mathbf{x}^{(i)}}{d\mathbf{y}} \quad (2.67)$$

$$= \begin{bmatrix} \frac{\partial x_1^{(i)}}{\partial y_1} & \dots & \frac{\partial x_1^{(i)}}{\partial y_n} \\ \vdots & & \vdots \\ \frac{\partial x_n^{(i)}}{\partial y_1} & \dots & \frac{\partial x_n^{(i)}}{\partial y_n} \end{bmatrix} \quad (2.68)$$

for each solution $i = 1, \dots, m$.

3. Then, the PDF of $\underline{\mathbf{y}}$ is

$$f(y) = \sum_{i=1}^m \left| \det \left(\tilde{J}_i(\mathbf{x}^{(i)}) \right) \right| f(\mathbf{x}^{(i)}), \quad (2.69)$$

where $\det(\cdot)$ is the determinant of a square matrix and $|\cdot|$ is the absolute value. **If for \mathbf{y} the system has no solution, then**

$$f(\mathbf{y}) = \mathbf{0}. \quad (2.70)$$

Example 2.2.2. Consider the transformed random variable $\underline{\mathbf{y}}$ given by

$$\underline{\mathbf{y}} = g(\underline{\mathbf{x}}) \quad (2.71)$$

$$= \begin{bmatrix} \underline{x}_1^2 \\ \underline{x}_1 + \underline{x}_2 \end{bmatrix}. \quad (2.72)$$

1. For all \mathbf{y} solve

$$\mathbf{y} = g(\mathbf{x}) \quad (2.73)$$

$$= \begin{bmatrix} \underline{x}_1^2 \\ \underline{x}_1 + \underline{x}_2 \end{bmatrix} \quad (2.74)$$

for \mathbf{x} .

Case I ($y_1 < 0$): No solutions.

Case II ($y_1 = 0$): One solution

$$\mathbf{x} = \begin{bmatrix} 0 \\ y_2 \end{bmatrix}. \quad (2.75)$$

Case III ($y_1 > 0$): two solutions:

$$\mathbf{x}^{(1)} = \begin{bmatrix} \sqrt{y_1} \\ y_2 - \sqrt{y_1} \end{bmatrix} \quad (2.76)$$

$$\mathbf{x}^{(2)} = \begin{bmatrix} -\sqrt{y_1} \\ y_2 + \sqrt{y_1} \end{bmatrix}. \quad (2.77)$$

2. For the i th solution (using approach 1):

$$\det(J_i(\mathbf{x})) = \det\left(\frac{d\mathbf{g}(\mathbf{x})}{d\mathbf{x}}\right) \quad (2.78)$$

$$= \det\left(\begin{bmatrix} \frac{\partial x_1^{(i)}}{\partial y_1} & \frac{\partial x_1^{(i)}}{\partial y_2} \\ \frac{\partial x_2^{(i)}}{\partial y_1} & \frac{\partial x_2^{(i)}}{\partial y_2} \end{bmatrix}\right) \quad (2.79)$$

$$= \frac{\partial x_1^{(i)}}{\partial y_1} \frac{\partial x_2^{(i)}}{\partial y_2} - \frac{\partial x_2^{(i)}}{\partial y_1} \frac{\partial x_1^{(i)}}{\partial y_2}. \quad (2.80)$$

Case I ($y_1 < 0$): no solutions.

Case II ($y_1 = 0$): one solution:

$$\det(J(\mathbf{x})) = 0. \quad (2.81)$$

Case III ($y_1 > 0$): two solutions:

$$J_1(\mathbf{x}) = \frac{1}{2\sqrt{y_1}} \quad (2.82)$$

$$J_2(\mathbf{x}) = -\frac{1}{2\sqrt{y_1}}. \quad (2.83)$$

3. The JPDP of $\underline{\mathbf{y}}$ is then

$$f_{\mathbf{y}}(\mathbf{y}) = \sum_i f_{\mathbf{x}}(\mathbf{x}) |\det(J_i)| \quad (2.84)$$

or $f_{\mathbf{y}}(\mathbf{y}) = \mathbf{0}$ if there are no solutions.

$$f_{\mathbf{y}}(\mathbf{y}) = f_{\mathbf{x}^{(1)}}(\mathbf{x}) |\det J_1| + f_{\mathbf{x}^{(2)}}(\mathbf{x}) |\det J_2| \quad (2.85)$$

$$= f_{\mathbf{x}}(\sqrt{y_1}, y_2 - \sqrt{y_1}) \left| \frac{1}{2\sqrt{y_1}} \right| + f_{\mathbf{x}}(-\sqrt{y_1}, y_2 + \sqrt{y_1}) \left| -\frac{1}{2\sqrt{y_1}} \right|. \quad (2.86)$$

△

So far, in this method, we assumed that $k = n$. But if $k < n$, then follow the following steps.

1. Introduce $n - k$ “dummy” variables^a:

$$\underline{y}_{k+1} = g_{k+1}(\underline{\mathbf{x}}) \quad (2.87)$$

$$\vdots \quad (2.88)$$

$$\underline{y}_n = g_n(\underline{\mathbf{x}}). \quad (2.89)$$

2. Evaluate $f(y_1, \dots, y_k, \dots, y_n)$ using the methods outlined before (now we have $k = n$).

3. Marginalize out $\underline{y}_{k+1}, \dots, \underline{y}_n$ to obtain $f(\underline{y}_1, \dots, \underline{y}_k)$.

^aI believe it's better if the introduced mapping is bijective.

2.2.3 Linear transformation

If $\mathbf{g}(\mathbf{x}) = \mathbf{A}\mathbf{x}$ (i.e., \mathbf{g} is linear), then the mean and the covariance are given by

$$\mathbb{E}[\underline{\mathbf{y}}] = \mathbf{A}\bar{\mathbf{x}}, \quad (2.90)$$

$$\text{Cov}[\underline{\mathbf{y}}] = \mathbb{E}[(\underline{\mathbf{y}} - \bar{\mathbf{y}})(\underline{\mathbf{y}} - \bar{\mathbf{y}})^{\top}] \quad (2.91)$$

$$= \mathbb{E}[\mathbf{A}(\underline{\mathbf{x}} - \bar{\mathbf{x}})(\mathbf{A}(\underline{\mathbf{x}} - \bar{\mathbf{x}}))^{\top}] \quad (2.92)$$

$$= \mathbf{A}\mathbb{E}[(\underline{\mathbf{x}} - \bar{\mathbf{x}})(\underline{\mathbf{x}} - \bar{\mathbf{x}})^{\top}] \mathbf{A}^{\top} \quad (2.93)$$

$$= \mathbf{A}\Sigma_x \mathbf{A}^{\top}. \quad (2.94)$$

Transform	Acronym	Benefits	Drawbacks
First order Taylor expansion	TT1	Simple and fast to compute.	Inaccurate for highly nonlinear functions.
Second order Taylor expansion	TT2	Better accuracy than TT1.	Requires computation of Hessian of the function.
Monte-Carlo simulation	MCT	Very accurate.	Requires many samples, thus computationally inefficient.
Unscented transform	UCT	Accuracy close to TT2	Doesn't require second order information (Hessian).

Table 2.1: Summary of transforms that approximate the covariances of propagated random variables.

Note that the relation is an exact relation, not an approximation.

2.2.4 Nonlinear transformation (covariance propagation)

For the nonlinear case, getting the exact transformed covariance is tedious and cumbersome. Therefore, tools that approximate the covariance are used that trade off between accuracy and computational efficiency. Most of the information from this section is obtained from [3]. Table 2.1 summarizes the results.

2.2.4.1 First order Taylor expansion

For a nonlinear transformation, a simple approximation is linearization. Thus,

$$\underline{\mathbf{y}} \approx \mathbf{g}(\bar{\mathbf{x}}) + \mathbf{g}'(\bar{\mathbf{x}})\delta\mathbf{x}. \quad (2.95)$$

The random variable becomes $\delta\mathbf{x} \sim \mathcal{N}(\mathbf{x} - \bar{\mathbf{x}}, \Sigma_x)$. The covariance of $\underline{\mathbf{y}}$ is then given by

$$\mathbb{E}[\underline{\mathbf{y}}] = \mathbf{g}(\bar{\mathbf{x}}), \quad (2.96)$$

$$\text{Cov}[\underline{\mathbf{y}}] = \mathbb{E}[(\underline{\mathbf{y}} - \bar{\mathbf{y}})(\underline{\mathbf{y}} - \bar{\mathbf{y}})^\top] \quad (2.97)$$

$$= \mathbb{E}[(\underline{\mathbf{y}} - \mathbf{g}(\bar{\mathbf{x}}))(\underline{\mathbf{y}} - \mathbf{g}(\bar{\mathbf{x}}))^\top] \quad (2.98)$$

$$= \mathbb{E}[\mathbf{g}'(\bar{\mathbf{x}})\delta\mathbf{x}\delta\mathbf{x}^\top\mathbf{g}'(\bar{\mathbf{x}})^\top] \quad (2.99)$$

$$= \mathbf{g}'(\bar{\mathbf{x}})\Sigma_x\mathbf{g}'(\bar{\mathbf{x}})^\top. \quad (2.100)$$

2.2.4.2 Monte-Carlo simulation

Monte-Carlo simulation consists of simply sampling a *large* number (say in the order of $N = 10^4$) of the domain random variable \mathbf{x} and passing it through the nonlinear function \mathbf{g} to get a large number of samples of $\underline{\mathbf{y}}, \mathbf{y}^{(i)}, i = 1, \dots, N$. The mean and covariance are then given by

$$\bar{\mathbf{y}} = \frac{1}{N} \sum_{i=1}^N \mathbf{y}^{(i)}, \quad (2.101)$$

$$\text{Cov}[\bar{\mathbf{y}}] \approx \frac{1}{N-1} \sum_{i=1}^N (\mathbf{y}^{(i)} - \bar{\mathbf{y}})(\mathbf{y}^{(i)} - \bar{\mathbf{y}})^\top, \quad (2.102)$$

where the $\frac{1}{N-1}$ term is to ensure that the estimator is *unbiased*.

2.2.4.3 Unscented transform

Let n_x be the length of the domain random variable (i.e., $\mathbf{x} \in \mathbb{R}^{n_x}$). Further, let $\mathbf{u}_i \in \mathbb{R}^{n_x}$ be the left singular vector of the singular value σ_i of the domain covariance matrix Σ_x . That is,

$$\Sigma_x = \mathbf{U}\Sigma\mathbf{U}^\top \quad (2.103)$$

$$= \sum_{i=1}^{n_x} \sigma_i^2 \mathbf{u}_i \mathbf{u}_i^\top, \quad (2.104)$$

where Σ is a diagonal matrix of the singular values.

Let $\mathbf{x}^{(i)}$ be the i -th realization of \mathbf{x} . For this version of the unscented transform (UT), the number of realizations will be $2n_x + 1$. Specifically, set

$$\mathbf{x}^{(0)} = \bar{\mathbf{x}}, \quad (2.105)$$

$$\mathbf{x}^{(\pm i)} = \bar{\mathbf{x}} \pm \sqrt{n_x + \lambda} \sigma_i \mathbf{u}_i, \quad i = 1, \dots, n_x, \quad (2.106)$$

where λ is a user-defined parameter to be discussed and $\bar{\mathbf{x}}$ is the mean of \mathbf{x} (in estimation problems, this is usually the predicted value). Further, let

$$\omega^{(0)} = \frac{\lambda}{n_x + \lambda}, \quad (2.107)$$

$$\omega^{(\pm i)} = \frac{1}{2(n_x + \lambda)}, \quad i = 1, \dots, n_x. \quad (2.108)$$

Next, set

$$\mathbf{y}^{(\pm i)} = \mathbf{g}(\mathbf{x}^{(\pm i)}), \quad i = 0, \dots, n_x. \quad (2.109)$$

Then, the stats of the propagated random variable \mathbf{y} are

$$\bar{\mathbf{y}} = \sum_{i=-n_x}^{n_x} \omega^{(i)} \mathbf{y}^{(i)}, \quad (2.110)$$

$$\Sigma_y = (1 - \alpha^2 + \beta) \left(\mathbf{y}^{(0)} - \bar{\mathbf{y}} \right) \left(\mathbf{y}^{(0)} - \bar{\mathbf{y}} \right)^\top \quad (2.111)$$

$$+ \sum_{i=-n_x}^{n_x} \omega^{(i)} \left(\mathbf{y}^{(i)} - \bar{\mathbf{y}} \right) \left(\mathbf{y}^{(i)} - \bar{\mathbf{y}} \right)^\top, \quad (2.112)$$

where α and β are user-defined parameters to be discussed. Further, $\omega^{(0)} + (1 - \alpha^2 + \beta)$ is often denoted $\omega_c^{(0)}$. The UT discussed above is the ‘mod’ version of the UT. There is another version denoted ‘std’ obtained by removing the first term in (2.111). I have no idea about the differences between the ‘std’ and ‘mod’ versions.

User-defined parameters

The authors of [3] discuss the values to set the user-defined parameters.

- α controls the spread of the sigma points and is suggested to be approximately 10^{-3} .

- β compensates for the distribution (not really sure what that means), and should be chosen as $\beta = 2$ when $\underline{\mathbf{x}}$ is Gaussian.
- λ is defined by $\lambda = \alpha^2 (n_x + \kappa) - n_x$, where κ is usually zero.
- $\omega^{(0)} = 1 - \frac{n_x}{3}$ for the ‘std’ version of the UT when $\underline{\mathbf{x}}$ is Gaussian.

2.2.4.4 Spherical cubature method

Given a function a random variable $\underline{\mathbf{x}} \in \mathbb{R}^{n_x}$ and a nonlinear function $\mathbf{g} : \mathbb{R}^{n_x} \rightarrow \mathbb{R}^{n_y}$, then the spherical cubature method¹ approximates the mean and the covariance of the transformed random variable $\underline{\mathbf{y}} = \mathbf{g}(\underline{\mathbf{x}})$. The method is summarized in Algorithm 1. More details are available in [4, Algorithm 6.8]

Algorithm 1 Nonlinear transformation: The spherical cubature

- 1: **Input:** The distribution on the random variable $\underline{\mathbf{x}} \sim \mathcal{N}(\underline{\boldsymbol{\mu}}_x, \underline{\boldsymbol{\Sigma}}_x)$ where $\mathbf{x} \in \mathbb{R}^{n_x}$, and the nonlinear function $\mathbf{g} : \mathbb{R}^{n_x} \rightarrow \mathbb{R}^{n_y}$.
- 2: Generate the unit sigma points using

$$\boldsymbol{\xi} = \sqrt{n_x} \begin{bmatrix} \mathbf{1} & -\mathbf{1} \end{bmatrix} \in \mathbb{R}^{n_x \times 2n_x}. \quad (2.113)$$

- 3: Generate the sigma points

$$\mathcal{X}^{(i)} = \underline{\boldsymbol{\mu}}_x + \sqrt{\underline{\boldsymbol{\Sigma}}_x} \boldsymbol{\xi}^{(i)}, \quad i = 1, \dots, 2n_x, \quad (2.114)$$

where $\boldsymbol{\xi}^{(i)} \in \mathbb{R}^{n_x}$ is the i -th column of $\boldsymbol{\xi}$, and $\underline{\boldsymbol{\Sigma}}_x = \sqrt{\underline{\boldsymbol{\Sigma}}_x} \sqrt{\underline{\boldsymbol{\Sigma}}_x}^\top$ is obtained using Cholesky decomposition.

- 4: Pass the sigma-points through the nonlinearity

$$\mathcal{Y}^{(i)} = \mathbf{g}(\mathcal{X}^{(i)}), \quad i = 1, \dots, 2n_x. \quad (2.115)$$

- 5: Compute

$$\underline{\boldsymbol{\mu}}_y = \frac{1}{2n_x} \sum_{i=1}^{2n_x} \mathcal{Y}^{(i)}, \quad (2.116)$$

$$\underline{\boldsymbol{\Sigma}}_y = \frac{1}{2n_x} \sum_{i=1}^{2n_x} (\mathcal{Y}^{(i)} - \underline{\boldsymbol{\mu}}_y) (\mathcal{Y}^{(i)} - \underline{\boldsymbol{\mu}}_y)^\top. \quad (2.117)$$

2.3 Statistics of random variables

2.3.1 Order statistics

Definition 2.3.1 (Order statistics). Let x_1, \dots, x_n be i.i.d. random variables with marginal CDF and

¹It's a type of a sigma-point transformation.

PDF $F(\cdot)$ and $f(\cdot)$, respectively. Consider the transformation

$$\underline{x}_{(1)} = \text{smallest value of } (\underline{x}_1, \dots, \underline{x}_n) \quad (2.118)$$

$$\underline{x}_{(2)} = \text{2nd smallest value of } (\underline{x}_1, \dots, \underline{x}_n) \quad (2.119)$$

$$\vdots \quad (2.120)$$

$$\underline{x}_{(n)} = \text{nth smallest (i.e., largest) value of } (\underline{x}_1, \dots, \underline{x}_n), \quad (2.121)$$

then $\underline{x}_{(1)}, \dots, \underline{x}_{(n)}$ are the *order statistics* of $\underline{x}_1, \dots, \underline{x}_n$. Furthermore, $\underline{x}_{(k)}$ is the *kth order statistic* of $\underline{x}_1, \dots, \underline{x}_n$.

Theorem 2.3.1 (CDF of order statistics). The CDF of the k th order statistic of $\underline{x}_1, \dots, \underline{x}_n$ is

$$F_{(k)}(x) = \sum_{i=k}^n \binom{n}{i} F(x)^i (1 - F(x))^{n-i}. \quad (2.122)$$

2.3.2 Mean

Theorem 2.3.2 (Mean of a function). The *mean* of a function $g : \mathbb{R}^n \rightarrow \mathbb{R}$ of the n real random variables $\underline{x}_1, \dots, \underline{x}_n$ is given by

$$\mathbb{E}[g(\mathbf{x})] = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} g(\mathbf{x}) f(\mathbf{x}) d\mathbf{x}, \quad (2.123)$$

where $\mathbb{E}[\cdot]$ is the multidimensional expectation operator (Definition 1.6.1).

Theorem 2.3.3 (Mean of a function of complex random variables). Mean of a function $g : \mathbb{R}^n \rightarrow \mathbb{R}$ of the n complex random variables $\underline{z}_i = \underline{x}_i + j\underline{y}_i, i = 1, \dots, n$ is given by

$$\mathbb{E}[g(\underline{z}_1, \dots, \underline{z}_n)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} g(x_1 + jy_1, \dots, x_n + jy_n) f(x_1, \dots, x_n, y_1, \dots, y_n) d\mathbf{x} d\mathbf{y}, \quad (2.124)$$

where $\mathbb{E}[\cdot]$ is the multidimensional expectation operator (Definition 1.6.1).

Definition 2.3.2 (Mean of a random vector). The *mean* of a random vector $\mathbf{x} = [\underline{x}_1 \ \cdots \ \underline{x}_n]^T$ is defined as

$$\mathbb{E}[\mathbf{x}] = [\mathbb{E}[\underline{x}_1] \ \cdots \ \mathbb{E}[\underline{x}_n]]^T. \quad (2.125)$$

Note that the mean of a random vector is also a vector of the same length. That is, $\mathbb{E}[\mathbf{x}] \in \mathbb{R}^n$.

Definition 2.3.3 (*Mean of a random matrix*). Let

$$\underline{\mathbf{X}} = \begin{bmatrix} \underline{\mathbf{x}}_1^\top \\ \vdots \\ \underline{\mathbf{x}}_m^\top \end{bmatrix}_{m \times n}, \quad (2.126)$$

be a random matrix, where $\underline{\mathbf{x}}_i \in \mathbb{R}^n, i = 1, \dots, m$. Then, the mean of $\underline{\mathbf{X}}$ is given by

$$\mathbb{E} [\underline{\mathbf{X}}] = \begin{bmatrix} \mathbb{E} [\underline{\mathbf{x}}_1]^\top \\ \vdots \\ \mathbb{E} [\underline{\mathbf{x}}_m]^\top \end{bmatrix}. \quad (2.127)$$

Properties of the expectation operator

The expectation operator is linear. That is

$$\mathbb{E} [\alpha_1 g_1(\underline{\mathbf{x}}) + \dots + \alpha_m g_m(\underline{\mathbf{x}})] = \alpha_1 \mathbb{E} [g_1(\underline{\mathbf{x}})] + \dots + \alpha_m \mathbb{E} [g_m(\underline{\mathbf{x}})] \quad (2.128)$$

for any $\underline{\mathbf{x}}$ (real or complex) and any *deterministic* $\alpha_1, \dots, \alpha_m$ (real or complex). As a result,

$$\mathbb{E} [\mathbf{A}\underline{\mathbf{x}}] = \mathbf{A}\mathbb{E} [\underline{\mathbf{x}}], \quad (2.129)$$

for any deterministic matrix \mathbf{A} .

2.3.3 Covariance and correlation

Definition 2.3.4 (*Correlation matrix*). The *correlation matrix* of the (complex) random variable $\underline{\mathbf{x}} = \begin{bmatrix} \underline{x}_1 & \dots & \underline{x}_n \end{bmatrix}^\top$ is defined as

$$\mathbf{R} = \mathbb{E} [\underline{\mathbf{x}}\underline{\mathbf{x}}^\mathbf{H}], \quad (2.130)$$

where $(\cdot)^\mathbf{H}$ is the Hermitian^a of a vector. The (i, j) th element of \mathbf{R} is the correlation of \underline{x}_i and \underline{x}_j .

^aHermitian: $\mathbf{x}^\mathbf{H} = (\mathbf{x}^*)^\top = (\mathbf{x}^\top)^*$.

Definition 2.3.5 (*Covariance matrix*). The *covariance matrix* of the random vector $\underline{\mathbf{x}} = \begin{bmatrix} \underline{x}_1 & \dots & \underline{x}_n \end{bmatrix}^\top$ is

$$\mathbf{C} = \mathbb{E} [(\underline{\mathbf{x}} - \mathbb{E} [\underline{\mathbf{x}}]) (\underline{\mathbf{x}} - \mathbb{E} [\underline{\mathbf{x}}])^\mathbf{H}]. \quad (2.131)$$

The (i, j) th element of \mathbf{C} is the covariance of \underline{x}_i and \underline{x}_j . That is,

$$C_{ij} = \mathbb{E} [(\underline{x}_i - \mathbb{E} [\underline{x}_i]) (\underline{x}_j - \mathbb{E} [\underline{x}_j])^*]. \quad (2.132)$$

The covariance matrix \mathbf{C} of $\underline{\mathbf{x}}$ is equal to the correlation matrix of the “centered” random vector $\underline{\mathbf{x}} - \mathbb{E}[\underline{\mathbf{x}}]$.

Definition 2.3.6 (*Mutually orthogonal random variables*). The random variables x_1, \dots, x_n are called *mutually orthogonal* if the correlation matrix \mathbf{R} is diagonal. That is, $R_{ij} = 0, i \neq j$.

Definition 2.3.7 (*Mutually uncorrelated random variables*). The random variables x_1, \dots, x_n are called *mutually uncorrelated* if the covariance matrix \mathbf{C} is diagonal. That is, $C_{ij} = 0, i \neq j$.

Notes: If the random variables x_1, \dots, x_n are independent then they are uncorrelated. The converse is not true in general.

Definition 2.3.8 (*Uncorrelated random vectors*). Two random vectors $\underline{\mathbf{x}}$ and $\underline{\mathbf{y}}$ are called *uncorrelated* if

$$\mathbb{E}[(\underline{\mathbf{x}} - \mathbb{E}[\underline{\mathbf{x}}])(\underline{\mathbf{y}} - \mathbb{E}[\underline{\mathbf{y}}])^H] = \mathbf{0}. \quad (2.133)$$

Or equivalently

$$\mathbb{E}[\underline{\mathbf{x}}\underline{\mathbf{y}}^H] = \mathbb{E}[\underline{\mathbf{x}}]\mathbb{E}[\underline{\mathbf{y}}]^H. \quad (2.134)$$

Definition 2.3.9 (*Orthogonal random vectors*). Two random vectors $\underline{\mathbf{x}}$ and $\underline{\mathbf{y}}$ are called *orthogonal* if

$$\mathbb{E}[\underline{\mathbf{x}}\underline{\mathbf{y}}^H] = \mathbf{0}. \quad (2.135)$$

Properties

1. Relationship between the correlation and covariance matrix

$$\mathbf{C} = \mathbf{R} - \mathbb{E}[\underline{\mathbf{x}}]\mathbb{E}[\underline{\mathbf{x}}]^H \quad (2.136)$$

$$= \mathbb{E}[\underline{\mathbf{x}}\underline{\mathbf{x}}^H] - \mathbb{E}[\underline{\mathbf{x}}]\mathbb{E}[\underline{\mathbf{x}}]^H. \quad (2.137)$$

2. The correlation and covariance matrix are positive semidefinite matrices.

Proof. For any vector \mathbf{a} ,

$$\mathbf{a}^H \mathbf{R} \mathbf{a} = \mathbf{a}^H \mathbb{E}[\underline{\mathbf{x}}\underline{\mathbf{x}}^H] \mathbf{a} \quad (2.138)$$

$$= \mathbb{E} \left[\left(\mathbf{a}^H \underline{\mathbf{x}} \right) \left(\mathbf{a}^H \underline{\mathbf{x}} \right)^H \right] \quad (2.139)$$

$$= \mathbb{E} \left[\left| \mathbf{a}^H \underline{\mathbf{x}} \right|^2 \right] \quad (2.140)$$

$$\geq 0. \quad (2.141)$$

□

3. The eigenvalues of the correlation and covariance matrix are real and non-negative.

Theorem 2.3.4. If \mathbf{R} is the correlation matrix of the n -dimensional random vector $\underline{\mathbf{x}}$. Then,

$$\mathbb{E} \left[\underline{\mathbf{x}}^H \mathbf{R}^{-1} \underline{\mathbf{x}} \right] = n. \quad (2.142)$$

Proof. Since $\underline{\mathbf{x}}^H \mathbf{R}^{-1} \underline{\mathbf{x}}$ is a scalar, we have

$$\mathbb{E} \left[\underline{\mathbf{x}}^H \mathbf{R}^{-1} \underline{\mathbf{x}} \right] = \mathbb{E} \left[\text{tr} \left(\underline{\mathbf{x}}^H \mathbf{R}^{-1} \underline{\mathbf{x}} \right) \right] \quad (2.143)$$

$$= \mathbb{E} \left[\text{tr} \left(\underline{\mathbf{x}} \underline{\mathbf{x}}^H \mathbf{R}^{-1} \right) \right] \quad (2.144)$$

$$= \text{tr} \mathbb{E} \left[\underline{\mathbf{x}} \underline{\mathbf{x}}^H \mathbf{R}^{-1} \right] \quad (2.145)$$

$$= \text{tr} \left(\mathbb{E} \left[\underline{\mathbf{x}} \underline{\mathbf{x}}^H \right] \mathbf{R}^{-1} \right) \quad (2.146)$$

$$= \text{tr} \left(\mathbf{R} \mathbf{R}^{-1} \right) \quad (2.147)$$

$$= \text{tr} \mathbf{1} \quad (2.148)$$

$$= n. \quad (2.149)$$

□

2.4 Moments and characteristic functions

Definition 2.4.1 (*Moments of random vectors*). Consider the random variables $\underline{x}_1, \dots, \underline{x}_N$ and $n_i = 0, 1, \dots$

- The (n_1, n_2, \dots, n_N) th *moment* of $\underline{x}_1, \dots, \underline{x}_N$ is

$$\mathbb{E} \left[\underline{x}_1^{n_1} \dots \underline{x}_N^{n_N} \right]. \quad (2.150)$$

Example ($N = 3$): (2,3,1)th moment is $\mathbb{E} \left[\underline{x}_1^2 \cdot \underline{x}_2^3 \cdot \underline{x}_3^1 \right]$.

- The (n_1, \dots, n_N) th *central moment* of $\underline{x}_1, \dots, \underline{x}_N$ is

$$\mathbb{E} \left[(\underline{x}_1 - \mathbb{E}[\underline{x}_1])^{n_1} \dots (\underline{x}_N - \mathbb{E}[\underline{x}_N])^{n_N} \right]. \quad (2.151)$$

- The (n_1, \dots, n_N) th *absolute moment* of $\underline{x}_1, \dots, \underline{x}_N$ is

$$\mathbb{E} \left[|\underline{x}_1|^{n_1} \dots |\underline{x}_N|^{n_N} \right]. \quad (2.152)$$

Definition 2.4.2 (*Characteristic function*). The *characteristic function* of a random vector \mathbf{x} is defined as

$$\Phi_{\mathbf{x}}(\Omega) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} e^{j(\omega_1 x_1 + \cdots + \omega_N x_N)} f(x_1, \dots, x_N) dx_1 \cdots dx_N \quad (2.153)$$

$$= \mathbb{E} \left[e^{j\Omega^T \mathbf{x}} \right] \quad (2.154)$$

where $\Omega = [\omega_1 \ \cdots \ \omega_N]^T \in \mathbb{R}^N$ and $\mathbf{x} = [x_1 \ \cdots \ x_N]^T$. The inverse relation is given by

$$f(x_1, \dots, x_N) = \frac{1}{(2\pi)^N} \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \Phi_{\mathbf{x}}(\Omega) e^{-j(\omega_1 x_1 + \cdots + \omega_N x_N)} d\omega_1 \cdots d\omega_N. \quad (2.155)$$

Note that it looks similar to the Fourier transform of the PDF but it's not exactly the same because of the sign in the exponential.

Definition 2.4.3 (*Moment generating function*). The *moment generating function* of a random vector \mathbf{x} is defined as

$$\text{Var} [\Phi]_{\mathbf{x}}(\mathbf{s}) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} e^{s_1 x_1 + \cdots + s_N x_N} f(x_1, \dots, x_N) dx_1 \cdots dx_N \quad (2.156)$$

$$= \mathbb{E} \left[e^{\mathbf{s}^T \mathbf{x}} \right] \quad (2.157)$$

where

$$\mathbf{s} = [s_1 \ \cdots \ s_N]^T. \quad (2.158)$$

Theorem 2.4.1 (*Characteristic functions*). The (n_1, \dots, n_N) th *moment* of \mathbf{x} is given by

$$m_{n_1, \dots, n_N} = \mathbb{E} \left[x_1^{n_1} \cdots x_N^{n_N} \right] \quad (2.159)$$

$$= \frac{\partial^{n_1 + \cdots + n_N}}{\partial s_1^{n_1} \cdots \partial s_N^{n_N}} \bar{\Phi}_{\mathbf{x}}(\mathbf{s}) \Big|_{s_1 = \cdots = s_N = 0} \quad (2.160)$$

$$= \frac{1}{j^{n_1 + \cdots + n_N}} \frac{\partial^{n_1 + \cdots + n_N}}{\partial \omega_1^{n_1} \cdots \partial \omega_N^{n_N}} \bar{\Phi}_{\mathbf{x}}(\Omega) \Big|_{\omega_1 = \cdots = \omega_N = 0}. \quad (2.161)$$

Theorem 2.4.2 (*PDF of independent random variables*). Let x_1, \dots, x_n be independent random variables with marginal PDFs $f_i(x), i = 1, \dots, n$, respectively. The PDF of $y = x_1 + \cdots + x_n$ is

$$f_y(y) = f_1(y) * \cdots * f_n(y), \quad (2.162)$$

where $*$ denotes the convolution.

2.5 Conditional distributions

Similar to the single random variable definition of conditional distributions (Definition 1.4.6), the following are extensions for such definitions to the multivariate case.

Definition 2.5.1 (*CDF of continuous RV conditioned on discrete RV*). Let \underline{x} and \underline{y} be two random variables with \underline{x} discrete. The *conditional CDF of \underline{y} given $\underline{x} = x$* is defined as

$$F(\underline{y} | \underline{x} = x) = P(\underline{y} \leq \underline{y} | \underline{x} = x). \quad (2.163)$$

Remark 2.5.1. Some remarks about CDFs.

- Conditions: $p(x) = P(\underline{x} = x) \neq 0$.
- Notation: $F(\underline{y} | \underline{x} = x) \rightarrow F_{y|x}(y|x)$, or $F(y|x)$.

Definition 2.5.2 (*PMF of discrete RV conditioned a discrete RV*). Let \underline{x} and \underline{y} be two joint discrete random variables. The *conditional PMF of \underline{y} given $\underline{x} = x$* is defined as

$$p(\underline{y} | \underline{x} = x) = P(\underline{y} = \underline{y} | \underline{x} = x) \quad (2.164)$$

$$= \frac{p_{xy}(x, y)}{p_x(x)}. \quad (2.165)$$

Notation: $p(\underline{y} | \underline{x} = x) \rightarrow p_{y|x}(y|x)$, or $p(y|x)$.

Definition 2.5.3 (*PDF of continuous RV conditioned on continuous RV*). Let \underline{x} and \underline{y} be two jointly continuous random variables. The *conditional PDF of \underline{y} given $\underline{x} = x$* is defined as

$$f(\underline{y} | \underline{x} = x) = \frac{f_{xy}(x, y)}{f_x(x)}. \quad (2.166)$$

Remark 2.5.2. Some remarks about PDFs.

- Note: if \underline{x} is continuous, then $P(\underline{x} = x) = 0$. Therefore, it is not possible to use a PDF definition similar to that of a PMF.
- Notation: $f(\underline{y} | \underline{x} = x) \rightarrow f_{y|x}(y|x)$, or $f(y|x)$.

Definition 2.5.4 (*CDF of continous RV conditioned on continuous RV*). Let $\underline{x}, \underline{y}$ be two jointly continuous random variables. The *conditional CDF of \underline{y} given $\underline{x} = x$* is defined as

$$F(\underline{y} | \underline{x} = x) = \int_{-\infty}^y f_{y|x}(\lambda|x) d\lambda. \quad (2.167)$$

- Notation: $F(\underline{y} | \underline{x} = x) \rightarrow F_{y|x}(y|x)$, or $F(y|x)$.

Remark 2.5.3. Some remarks about CDFs.

1. Note that the conditional CDF is derived from the conditional PDF, rather than the other way around (deriving PDF from CDF as was done in the unconditional case).
2. For a given x , $f(y|x)$ is a valid PDF. That is,
 - 2.1. $f(y|x) \geq 0$
 - 2.2. $\int_{-\infty}^{\infty} f(y|x) dy = 1$.
3. $f(x|y) = \frac{f(x,y)}{f(y)} = \frac{f(y|x)f(x)}{f(y)}$.
4. If \underline{x} and \underline{y} are independent, then $f(x,y) = f(x)f(y)$. Therefore,

$$f(y|x) = f(y), \quad (2.168)$$

$$f(x|y) = f(x). \quad (2.169)$$

Definition 2.5.5 (*Multivariate conditional PDF*). The multivariate conditional PDF of $\underline{x}_n, \dots, \underline{x}_{k+1}$ given $\underline{k} = x_k, \dots, \underline{x}_1 = x_1$ is

$$f(x_n, \dots, x_{k+1} | x_k, \dots, x_1) = \frac{f(x_n, \dots, x_{k+1}, x_k, \dots, x_1)}{f(x_k, \dots, x_1)} \quad (2.170)$$

Definition 2.5.6 (*Multivariate conditional CDF*). The multivariate conditional CDF of $\underline{x}_n, \dots, \underline{x}_{k+1}$ given $\underline{k} = x_k, \dots, \underline{x}_1 = x_1$ is

$$F(x_n, \dots, x_{k+1} | x_k, \dots, x_1) = \quad (2.171)$$

$$\int_{-\infty}^{x_n} \cdots \int_{-\infty}^{x_{k+1}} f_{x_n, \dots, x_{k+1} | x_k, \dots, x_1}(\lambda_n, \dots, \lambda_{k+1} | x_k, \dots, x_1) d\lambda_n \cdots d\lambda_{k+1}. \quad (2.172)$$

2.5.1 Properties of conditional probability

1. **Total probability:**

$$f(y) = \int_{-\infty}^{\infty} f(y|x) f(x) dx. \quad (2.173)$$

2. **Bayes Theorem:**

$$f(x|y) = \frac{f(y|x) f(x)}{\int_{-\infty}^{\infty} f(y|x) f(x) dx}. \quad (2.174)$$

3. **Chain rule:**

$$f(x_1, \dots, x_n) = f(x_n | x_{n-1}, \dots, x_1) \cdots f(x_3 | x_2, x_1) f(x_2 | x_1) f(x_1). \quad (2.175)$$

4. **Marginalization:** Removing any number of variables on the left of the conditional line can be done by integrating the PDF with respect to them. For example,

$$f(\mathbf{x}_1 | \mathbf{x}_3) = \int_{-\infty}^{\infty} f(\mathbf{x}_1, \mathbf{x}_2 | \mathbf{x}_3) d\mathbf{x}_2. \quad (2.176)$$

5. **Chapman-Kolmogoroff equation:** Removing any number of variables to the right of the conditional line can be done by multiplying by their conditional density given the remaining variables on the right, and then integrating the product. That is,

$$f(\mathbf{x}_1 | \mathbf{x}_3) = \int_{-\infty}^{\infty} f(\mathbf{x}_1 | \mathbf{x}_2, \mathbf{x}_3) f(\mathbf{x}_2 | \mathbf{x}_3) d\mathbf{x}_2. \quad (2.177)$$

2.5.2 Conditional covariance

These equations are obtained from [5, Eq. (2.53b)].

Given a the random variable

$$\underline{\mathbf{x}} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \sim \mathcal{N}(\bar{\mathbf{x}}, \Sigma), \quad (2.178)$$

where

$$\bar{\mathbf{x}} = \begin{bmatrix} \bar{\mathbf{x}}_1 \\ \bar{\mathbf{x}}_2 \end{bmatrix}, \quad (2.179)$$

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{12}^T & \Sigma_{22} \end{bmatrix}. \quad (2.180)$$

Then, what is the distribution of $\underline{\mathbf{x}}_1$ given $\underline{\mathbf{x}}_2 = \hat{\mathbf{x}}_2$?

The Schur's compliment can be used to compute this value. Specifically,

$$\mathbb{E}[\underline{\mathbf{x}}_1 | \underline{\mathbf{x}}_2 = \hat{\mathbf{x}}_2] = \bar{\mathbf{x}}_1 + \Sigma_{12} \Sigma_{22}^{-1} (\hat{\mathbf{x}}_2 - \bar{\mathbf{x}}_2), \quad (2.181)$$

$$\text{Cov}[\underline{\mathbf{x}}_1 | \underline{\mathbf{x}}_2 = \hat{\mathbf{x}}_2] = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{12}^T. \quad (2.182)$$

2.6 Conditional expectation and variance

Definition 2.6.1 (*Conditional mean*). Given n random variables $\underline{x}_1, \dots, \underline{x}_n$, the *conditional mean* of the random variable $h(\underline{x}_1, \dots, \underline{x}_k)$ conditioned on $\underline{x}_{k+1}, \dots, \underline{x}_n$ is

$$\mathbb{E}[h(\underline{x}_1, \dots, \underline{x}_k) | x_{k+1}, \dots, x_n] = \quad (2.183)$$

$$\int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} h(x_1, \dots, x_k) f(x_1, \dots, x_k | x_{k+1}, \dots, x_n) dx_1 \cdots dx_k. \quad (2.184)$$

Remark 2.6.1. Some remarks about conditional means.

- Special case:

$$\mathbb{E} [\underline{x}_1 | x_2, \dots, x_n] = \int_{-\infty}^{\infty} x_1 f(x_1 | x_2, \dots, x_n) dx_1. \quad (2.185)$$

- If $\underline{x}_1, \dots, \underline{x}_n$ are independent of $\underline{x}_{k+1}, \dots, \underline{x}_n$, then

$$\mathbb{E} [h(\underline{x}_1, \dots, \underline{x}_k) | x_{k+1}, \dots, x_n] = \mathbb{E} [h(\underline{x}_1, \dots, \underline{x}_k)]. \quad (2.186)$$

Theorem 2.6.1 (Nested expectation). Let $\underline{x}, \underline{y}$ be two jointly distributed RVs. Further, let

$$g(y) = \mathbb{E} [\underline{x} | \underline{y} = y] \quad (2.187)$$

$$= \int_{-\infty}^{\infty} x f(x | y) dx. \quad (2.188)$$

Then

$$\mathbb{E} [g(\underline{y})] = \mathbb{E} [\underline{x}]. \quad (2.189)$$

In other words,

$$\mathbb{E}_y [\mathbb{E}_{x|y} [\underline{x} | \underline{y}]] \rightarrow \mathbb{E} [\mathbb{E} [\underline{x} | \underline{y}]] \quad (2.190)$$

$$= \mathbb{E} [\underline{x}]. \quad (2.191)$$

Proof.

$$\mathbb{E} [g(\underline{y})] = \int_{-\infty}^{\infty} g(y) f(y) dy \quad (2.192)$$

$$= \int_{-\infty}^{\infty} \mathbb{E} [\underline{x} | y] f(y) dy \quad (2.193)$$

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x f(x | y) f(y) dx dy \quad (2.194)$$

$$= \int_{-\infty}^{\infty} x \int_{-\infty}^{\infty} f(x, y) dy dx \quad (2.195)$$

$$= \int_{-\infty}^{\infty} x f(x) dx \quad (2.196)$$

$$= \mathbb{E} [\underline{x}]. \quad (2.197)$$

□

Theorem 2.6.2 (Multivariate nested expectation). Let $\underline{x}_1, \dots, \underline{x}_n$ be n random variables. Also, let

$$g(x_2, \dots, x_n) = \mathbb{E} [\underline{x}_1 | x_2, \dots, x_n]. \quad (2.198)$$

Then,

$$\mathbb{E} [g(\underline{x}_2, \dots, \underline{x}_n)] = \mathbb{E} [\underline{x}_1]. \quad (2.199)$$

Notation:

$$\mathbb{E}_{x_2, \dots, x_n} [\mathbb{E}_{x_1 | x_2, \dots, x_n} [\underline{x}_1 | x_2, \dots, x_n]] \rightarrow \mathbb{E} [\mathbb{E} [\underline{x}_1 | x_2, \dots, x_n]]. \quad (2.200)$$

Theorem 2.6.3 (Multivariate nested expectation: general case). Let $\underline{x}_1, \dots, \underline{x}_n$ be n random variables and $h : \mathbb{R}^k \rightarrow \mathbb{R}$. Also, let

$$g(x_{k+1}, \dots, x_n) = \mathbb{E} [h(\underline{x}_1, \dots, \underline{x}_k) | x_{k+1}, \dots, x_n]. \quad (2.201)$$

Then,

$$\mathbb{E} [g(\underline{x}_{k+1}, \dots, \underline{x}_n)] = \mathbb{E} [\mathbb{E} [h(\underline{x}_1, \dots, \underline{x}_k) | \underline{x}_{k+1}, \dots, \underline{x}_n]] \quad (2.202)$$

$$= \mathbb{E} [h(\underline{x}_1, \dots, \underline{x}_k)]. \quad (2.203)$$

Definition 2.6.2 (Conditional variance). Let $\underline{x}, \underline{y}$ be two random variables. Then the *conditional variance of \underline{x} given $\underline{y} = y$* is

$$\text{Var} [\underline{x} | \underline{y} = y] = \mathbb{E} \left[(\underline{x} - \mathbb{E} [\underline{x} | \underline{y} = y])^2 | \underline{y} = y \right]. \quad (2.204)$$

2.6.1 Properties of conditional variance

1. $\text{Var} [\underline{x} | \underline{y} = y] = \mathbb{E} [\underline{x}^2 | \underline{y} = y] - \mathbb{E} [\underline{x} | \underline{y} = y]^2$.
2. $\text{Var} [\underline{x} | \underline{y} = y]$ is a function of y .
3. $g(y) = \text{Var} [\underline{x} | \underline{y} = y]$ is a random variable.

Theorem 2.6.4 (Nested conditional variance). Let $\underline{x}, \underline{y}$ be two random variables. Then

$$\text{Var} [\underline{x}] = \mathbb{E} [\text{Var} [\underline{x} | \underline{y}]] + \text{Var} [\mathbb{E} [\underline{x} | \underline{y}]]. \quad (2.205)$$

Proof.

$$\mathbb{E} [\text{Var} [\underline{x} | \underline{y}]] = \mathbb{E} [\mathbb{E} [\underline{x}^2 | \underline{y}] - \mathbb{E} [\underline{x} | \underline{y}]^2] \quad (2.206)$$

$$= \mathbb{E} [\mathbb{E} [\underline{x}^2 | \underline{y}]] - \mathbb{E} [\mathbb{E} [\underline{x} | \underline{y}]^2] \quad (2.207)$$

$$= \mathbb{E} [\underline{x}^2] - \mathbb{E} [\mathbb{E} [\underline{x} | \underline{y}]^2]. \quad (2.208)$$

Further,

$$\text{Var} [\mathbb{E} [\underline{x} | \underline{y}]] = \mathbb{E} [\mathbb{E} [\underline{x} | \underline{y}]^2] - \mathbb{E} [\mathbb{E} [\underline{x} | \underline{y}]]^2 \quad (2.209)$$

$$= \mathbb{E} [\mathbb{E} [\underline{x} | \underline{y}]^2] - \mathbb{E} [\underline{x}]^2. \quad (2.210)$$

Adding these two results proves the theorem. \square

2.7 Passing measurements through a function

Theorem 2.7.1 (Passing measurements through a function). Let $\underline{y} \in \mathbb{R}^m$ be a (vector) random variable that depends on $\underline{x} \in \mathbb{R}^n$ through

$$\underline{y} = \mathbf{g}(\underline{x}). \quad (2.211)$$

Furthermore, let $\underline{z} \in \mathbb{R}^p$ be a transformed random variable given by

$$\underline{z} = \mathbf{f}(\underline{y}). \quad (2.212)$$

Then

$$f(\underline{x} | \underline{y}, \underline{z}) = f(\underline{x} | \underline{y}). \quad (2.213)$$

Furthermore, if \mathbf{g} is invertible and \mathbf{f} is invertible, then

$$f(\underline{x} | \underline{y}, \underline{z}) = f(\underline{x} | \underline{z}). \quad (2.214)$$

Note that if \mathbf{g} is invertible but \mathbf{f} is non-invertible, then (2.214) does not hold in general!

Proof. 1. The first part.

$$f(\mathbf{x}|\mathbf{y}, \mathbf{z}) = \frac{f(\mathbf{y}, \mathbf{z}|\mathbf{x}) f(\mathbf{x})}{f(\mathbf{y}, \mathbf{z})} \quad (2.215)$$

$$= \frac{f(\mathbf{z}|\mathbf{y}, \mathbf{x}) f(\mathbf{y}|\mathbf{x}) f(\mathbf{x})}{f(\mathbf{z}|\mathbf{y}) f(\mathbf{y})} \quad (2.216)$$

$$= \frac{\cancel{f(\mathbf{z}|\mathbf{y})} f(\mathbf{y}|\mathbf{x}) f(\mathbf{x})}{\cancel{f(\mathbf{z}|\mathbf{y})} f(\mathbf{y})} \quad (2.217)$$

$$= \frac{f(\mathbf{y}|\mathbf{x}) f(\mathbf{x})}{f(\mathbf{y})} \quad (2.218)$$

$$= f(\mathbf{x}|\mathbf{y}). \quad (2.219)$$

2. The second part. Let's expand $f(\mathbf{x}|\mathbf{y}, \mathbf{z})$.

$$f(\mathbf{x}|\mathbf{y}, \mathbf{z}) = \frac{f(\mathbf{y}, \mathbf{z}|\mathbf{x}) f(\mathbf{x})}{f(\mathbf{y}, \mathbf{z})} \quad (2.220)$$

$$= \frac{f(\mathbf{y}|\mathbf{z}, \mathbf{x}) f(\mathbf{z}|\mathbf{x}) f(\mathbf{x})}{f(\mathbf{y}|\mathbf{z}) f(\mathbf{z})} \quad (2.221)$$

$$= \frac{f(\mathbf{y}|\mathbf{z}, \mathbf{x})}{f(\mathbf{y}|\mathbf{z})} f(\mathbf{x}|\mathbf{z}). \quad (2.222)$$

Thus, (2.214) satisfied if and only if

$$f(\mathbf{y}|\mathbf{z}, \mathbf{x}) = f(\mathbf{y}|\mathbf{z}). \quad (2.223)$$

That is, if \mathbf{y} is independent of \mathbf{x} given \mathbf{z} . More rigorously, if there exists $\mathbf{y}_1 \neq \mathbf{y}_2$ such that

$$\begin{aligned} \mathbf{g}(\mathbf{x}_1) = \mathbf{y}_1 \\ \mathbf{g}(\mathbf{x}_2) = \mathbf{y}_2 \end{aligned} \implies \mathbf{x}_1 \neq \mathbf{x}_2, \quad (2.224)$$

and

$$\mathbf{z} = \mathbf{f}(\mathbf{y}_1) = \mathbf{f}(\mathbf{y}_2). \quad (2.225)$$

Then knowing \mathbf{z} without knowing \mathbf{x} will result in a non-singleton sample space S for the random variable $\underline{\mathbf{y}}$. That is, a set with more than one element which implies that $\underline{\mathbf{y}}$ cannot be determined uniquely. In such case, the sample space is the set of \mathbf{y} that satisfy (2.225).

On the other hand, if in addition to knowing \mathbf{z} , \mathbf{x} is also known, then the sample space of \mathbf{y} is a singleton. That is, $\underline{\mathbf{y}}$ can be determined uniquely.

The condition (2.225) occurs only if the function \mathbf{f} is non-invertible.

□

Chapter 3

Distributions

3.1 Bernoulli distribution

From [1].

Definition 3.1.1 (*Bernoulli drandom variable*). A random variable \underline{x} is called Bernoulli with parameter p if there exists an event A with probability $p = P(A)$ such that

$$x(s) = \begin{cases} 1, & s \in A \\ 0, & s \notin A, \end{cases} \quad (3.1)$$

where $s \in S$.

Remark 3.1.1. Some remarks on Bernoulli random variables.

- \underline{x} is a *discrete* random variable with range $\mathcal{R}_x = \{0, 1\}$.
- PMF of \underline{x}

$$p(0) = P(\underline{x} = 0) = P(A^C) = 1 - p = q \quad (3.2)$$

$$p(1) = P(\underline{x} = 1) = P(A) = p \quad (3.3)$$

$$p(x) = 0 \quad \text{if } x \notin \{0, 1\}. \quad (3.4)$$

3.2 Geometric random variables

From [1].

Definition 3.2.1 (*Geometric random variable*). A geometric random variable is a discrete random variable with a countable infinite set of possible values. The set of possible values for the geometric random

variable \underline{x} is

$$\mathcal{R}_x = \{1, 2, \dots\} = \mathbb{N}. \quad (3.5)$$

Theorem 3.2.1 (PMF of geometric random variables). Let \underline{x} be a geometric random variable with parameter p . The PMF of \underline{x} is then

$$p(x) = \begin{cases} (1-p)^{x-1} p, & x = 1, 2, \dots, \\ 0, & \text{otherwise.} \end{cases} \quad (3.6)$$

3.3 Binomial random variable

Definition 3.3.1 (*Binomial random variable*). Consider a sequence of n identical and independent Bernoulli trials with probability of success p . The random variable \underline{x} defined by

$$\underline{x} = \text{number of successes in the } n \text{ trials} \quad (3.7)$$

is called binomial with parameters n and p , or simply $B(n, p)$.

Remark 3.3.1. Some remarks on Binomial random variables.

- Set of possible values of \underline{x} :

$$\mathcal{R}_x = \{0, 1, \dots, n\}. \quad (3.8)$$

- Notation:

$$\underline{x} \sim B(n, p). \quad (3.9)$$

- Basic examples

- Number of heads in a sequence of 10 independent flips of a coin.
- Number of defective IC chips in a production sample of size n .

Theorem 3.3.1 (PMF of a Binomial random variable). Let $\underline{x} \sim B(n, p)$. The PMF of \underline{x} is

$$p(x) = \begin{cases} \binom{n}{x} p^x (1-p)^{n-x}, & x = 0, 1, \dots, n \\ 0, & \text{otherwise.} \end{cases} \quad (3.10)$$

3.4 Poisson distribution

A discrete random variable \underline{x} is called Poisson with parameter $\lambda > 0$ if its PMF has the form

$$p(x) = \begin{cases} \frac{\lambda^x}{x!} e^{-\lambda}, & x \in \mathbb{N}_0 \\ 0, & \text{otherwise,} \end{cases} \quad (3.11)$$

where $\mathbb{N}_0 = \{0, 1, \dots\}$.

3.5 Uniform random variables

Definition 3.5.1 (*Uniform random variable*). A continuous random variable \underline{x} is called uniform over the interval (a, b) , written as $\underline{x} \sim U(a, b)$, where $a < b$ are real numbers, if its PDF is

$$f(x) = \begin{cases} \frac{1}{b-a}, & a < x < b, \\ 0, & \text{otherwise.} \end{cases} \quad (3.12)$$

Theorem 3.5.1 (Uniform random variable CDF). The CDF is given by

$$F(x) = \begin{cases} 0, & x \leq a, \\ \frac{x-a}{b-a}, & a \leq x \leq b, \\ 1, & x \geq b. \end{cases} \quad (3.13)$$

3.6 Exponential random variable

Definition 3.6.1 (*Exponential random variable*). A continuous random variable \underline{x} is called exponential with parameter $\lambda > 0$ if its PDF is given by

$$f(x) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0, \\ 0, & x < 0. \end{cases} \quad (3.14)$$

Theorem 3.6.1 (Exponential random variable CDF). The CDF is given by

$$F(x) = \begin{cases} 1 - e^{-\lambda x}, & x > 0, \\ 0, & x \leq 0. \end{cases} \quad (3.15)$$

3.7 Gaussian random variables

Definition 3.7.1 (*Gaussian (normal) distribution*). A random vector $\underline{\mathbf{x}}$ is called a *Gaussian (or normal) random vector* with mean $\underline{\boldsymbol{\mu}}$ and covariance matrix \mathbf{C} if it follows a *Gaussian (or normal) distribution*

given by

$$f(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^n \det \mathbf{C}}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^\top \mathbf{C}^{-1}(\mathbf{x}-\boldsymbol{\mu})} \quad (3.16)$$

$$= \mathcal{N}(\boldsymbol{\mu}, \mathbf{C}). \quad (3.17)$$

The n elements of \mathbf{x} , $\underline{x}_1, \dots, \underline{x}_n$ are called *jointly normal random variables*.

Definition 3.7.2 (Standard Normal Distribution). A *standard normal distribution* is a normal distribution with a special case of $\boldsymbol{\mu} = \mathbf{0}$ and $\mathbf{C} = \mathbf{1}$. That is,

$$f(\mathbf{x}) = \mathcal{N}(\mathbf{0}, \mathbf{1}). \quad (3.18)$$

3.7.1 Properties of Gaussian random variables

Theorem 3.7.1 (Linear combination of normal RVs). A random vector $\underline{\mathbf{x}}$ is called a Gaussian (or normal) random vector if and only if the random variable $\mathbf{a}^\top \underline{\mathbf{x}}$ is a Gaussian random variable for any vector \mathbf{a} .

Remark 3.7.1 (Properties of Gaussian random variables). Here are some properties of Gaussian random variables.

1. A Gaussian random vector consists of marginally Gaussian random variables.
2. However, marginally Gaussian random variables are *not* necessarily jointly Gaussian.^a

^aThat is, if you take 3 Gaussian random variables, then they are not necessarily jointly Gaussian. Unless they are independent.

Definition 3.7.3 (Complex Gaussian random vector). A complex random vector

$$\underline{\mathbf{z}} = \begin{bmatrix} \underline{x}_1 + jy_1 \\ \vdots \\ \underline{x}_n + jy_n \end{bmatrix} \quad (3.19)$$

is a *complex Gaussian random vector* if and only if the real random variables $\underline{x}_1, \underline{y}_1, \dots, \underline{x}_n, \underline{y}_n$ are jointly normal.

Theorem 3.7.2 (Characteristic function of normal random variables). If the random variables $\underline{x}_1, \dots, \underline{x}_n$ are jointly Normal with mean 0 and covariance matrix \mathbf{C} . Then, their characteristic function is

$$\Phi_{\mathbf{x}}(\boldsymbol{\Omega}) = e^{-\frac{1}{2}\boldsymbol{\Omega}^\top \mathbf{C} \boldsymbol{\Omega}}. \quad (3.20)$$

3.8 Chi-squared distribution

Definition 3.8.1 (*Chi-squared distribution*). If $\underline{x}_1, \dots, \underline{x}_n$ are **independent** standard normal variables (check Definition 3.7.2), then the sum of their squares

$$\underline{y} = \sum_{k=1}^n \underline{x}_k^2 \in \mathbb{R} \quad (3.21)$$

is distributed according to the Chi-square distribution with n degrees of freedom, denoted by χ_n^2 . That is,

$$\underline{y} \sim \chi_n^2. \quad (3.22)$$

Remark 3.8.1. If we have a random vector $\underline{\mathbf{x}} \in \mathbb{R}^n$ with mean $\underline{\mu}_x = \mathbf{0}$ and a unit covariance $\mathbf{C} = \mathbf{1}$, then the sum of squares

$$\underline{y} = \underline{\mathbf{x}}^T \underline{\mathbf{x}} \quad (3.23)$$

$$(3.24)$$

follows a Chi-squared distribution with n degrees of freedom. That is

$$\underline{y} \sim \chi_n^2. \quad (3.25)$$

Remark 3.8.2. MATLAB has a built-in inverse CDF function called `chi2inv(p, nu)` that evaluates the inverse CDF of a Chi-squared distribution at the probability value $p \in [0, 1]$ and the degrees of freedom ν (which we denote as n in Definition 3.8.1).

3.8.1 Relation to the Mahalanobis distance

In estimation problems, the Mahalanobis distance is often encountered. Especially when testing for consistency of the estimates.

Definition 3.8.2 (*Mahalanobis distance*). The Mahalanobis distance $D_M \in \mathbb{R}$ is defined as

$$D_M \triangleq \sqrt{\mathbf{r}^T \mathbf{\Sigma}^{-1} \mathbf{r}}, \quad (3.26)$$

where $\mathbf{r} \in \mathbb{R}^n$ is a sample of

$$\underline{\mathbf{r}} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma}). \quad (3.27)$$

It is a *Euclidean distance* normalized by *uncertainty*.

Example 3.8.1. Consider a random variable $\underline{\mathbf{x}} \in \mathbb{R}^n$ distributed according to $\underline{\mathbf{x}} \sim \mathcal{N}(\underline{\mu}_x, \mathbf{\Sigma}_x)$. Say a

realization $\mathbf{x} \in \mathbb{R}^n$ is sampled from the distribution. Then the Mahalanobis distance

$$D_M = \sqrt{(\mathbf{x} - \boldsymbol{\mu}_x)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})} \quad (3.28)$$

is a measure of how far off the realization \mathbf{x} is from the mean $\boldsymbol{\mu}_x$. \triangle

Theorem 3.8.1. The squared Mahalanobis distance follows a Chi-squared distribution. That is,

$$\underline{D}_M^2 \sim \chi_n^2. \quad (3.29)$$

Proof. Given a random variable $\mathbf{r} \in \mathbb{R}^n$ that follows the distribution $\mathbf{r} \sim (\mathbf{0}, \boldsymbol{\Sigma})$, then the Mahalanobis distance squared is a random variable. Specifically,

$$\underline{D}_M^2 = \mathbf{r}^\top \boldsymbol{\Sigma}^{-1} \mathbf{r}. \quad (3.30)$$

Since the covariance matrix $\boldsymbol{\Sigma}$ is symmetric positive definite, then its eigenvectors are orthonormal and the eigenvalues are positive. Therefore, (3.30) can be written as

$$\underline{D}_M^2 = \mathbf{r}^\top \boldsymbol{\Sigma}^{-1} \mathbf{r} \quad (3.31)$$

$$= \mathbf{r}^\top (\mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^\top)^{-1} \mathbf{r} \quad (3.32)$$

$$= \mathbf{r}^\top \mathbf{U} \boldsymbol{\Lambda}^{-1} \mathbf{U}^\top \mathbf{r} \quad (3.33)$$

$$= \mathbf{r}^\top \mathbf{U} \boldsymbol{\Lambda}^{-\frac{1}{2}} \boldsymbol{\Lambda}^{-\frac{1}{2}} \mathbf{U}^\top \mathbf{r} \quad (3.34)$$

$$= (\boldsymbol{\Lambda}^{-\frac{1}{2}} \mathbf{U}^\top \mathbf{r})^\top (\boldsymbol{\Lambda}^{-\frac{1}{2}} \mathbf{U}^\top \mathbf{r}) \quad (3.35)$$

$$= \underline{\mathbf{y}}^\top \underline{\mathbf{y}}, \quad (3.36)$$

where

$$\underline{\mathbf{y}} = \boldsymbol{\Lambda}^{-\frac{1}{2}} \mathbf{U}^\top \mathbf{r} \in \mathbb{R}^n. \quad (3.37)$$

The next step is to check whether $\underline{\mathbf{y}}$ is a standard normal variable. First, let's check the mean.

$$\boldsymbol{\mu}_y = \mathbb{E}[\underline{\mathbf{y}}] \quad (3.38)$$

$$= \boldsymbol{\Lambda}^{-\frac{1}{2}} \mathbf{U}^\top \underbrace{\mathbb{E}[\mathbf{r}]}_{\mathbf{0}} \quad (3.39)$$

$$= \mathbf{0}. \quad (3.40)$$

The second condition to check is whether the covariance of $\underline{\mathbf{y}}$ is the identity matrix or not.

$$\text{Cov} [\underline{\mathbf{y}}] = \mathbb{E} [\underline{\mathbf{y}}\underline{\mathbf{y}}^T] \quad (3.41)$$

$$= \mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{U}^T \underbrace{\mathbb{E} [\underline{\mathbf{r}}\underline{\mathbf{r}}^T]}_{\mathbf{\Sigma}} \mathbf{U} \mathbf{\Lambda}^{-\frac{1}{2}} \quad (3.42)$$

$$= \mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{U}^T \underbrace{\mathbf{\Sigma}}_{\mathbf{U} \mathbf{\Lambda} \mathbf{U}^T} \mathbf{U} \mathbf{\Lambda}^{-\frac{1}{2}} \quad (3.43)$$

$$= \mathbf{\Lambda}^{-\frac{1}{2}} \cancel{\mathbf{U}^T} \mathbf{1} \cancel{\mathbf{U} \mathbf{\Lambda} \mathbf{U}^T} \mathbf{1} \cancel{\mathbf{U}} \mathbf{\Lambda}^{-\frac{1}{2}} \quad (3.44)$$

$$= \mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{\Lambda} \mathbf{\Lambda}^{-\frac{1}{2}} \quad (3.45)$$

$$= \mathbf{1}. \quad (3.46)$$

Thus, $\underline{\mathbf{y}}$ is a standard normal random variable. Therefore, according to Definition 3.8.1, \underline{D}_M^2 follows a Chi-squared distribution with n degrees of freedom. \square

Chapter 4

Graphs in probability

4.1 Graph definitions and terminology

Definition 4.1.1 (*Directed graphs*). A directed graph consists of a set of vertices and a set of directed edges that connect pairs of vertices.

In probability, the vertices are random variables. Figure 4.1 shows an example of a directed graph.

Definition 4.1.2 (*Directed acyclic graph*). A directed acyclic graph (DAG) is a directed graph that contains no cycles. That is, no directed paths from a vertex back to itself.

In case of probability, a DAG is called a *Bayesian network* (BN).

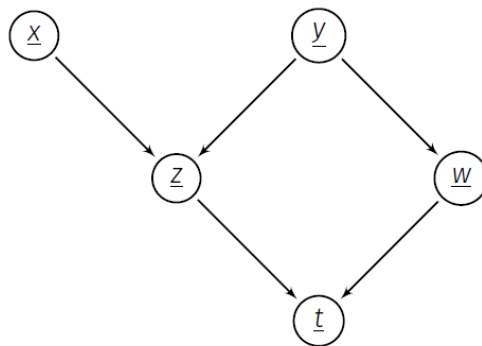


Figure 4.1: Illustration of a simple directed graph in probability.

The terms are illustrated in the figures below.

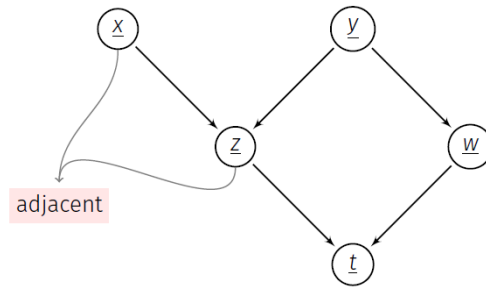


Figure 4.2: Directed graph terminology: adjacent vertices.

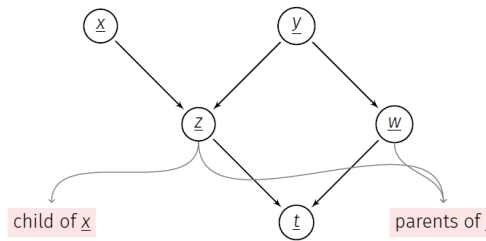


Figure 4.3: Directed graph terminology: child and parent vertices.

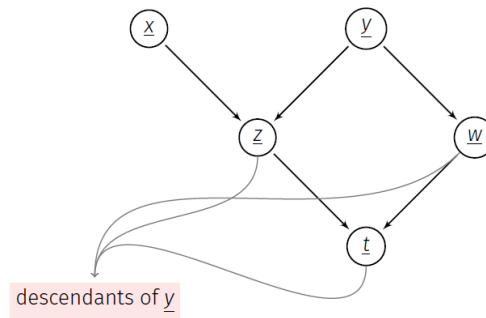


Figure 4.4: Directed graph terminology: descendants vertices.

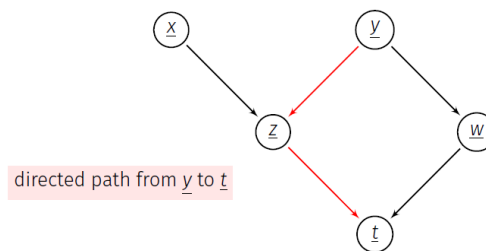


Figure 4.5: Directed graph terminology: directed path.

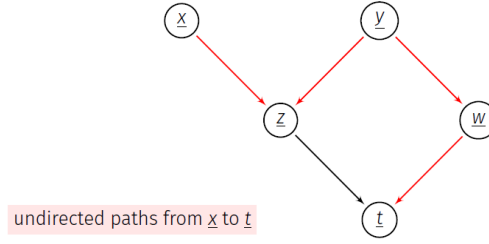


Figure 4.6: Directed graph terminology: undirected path.

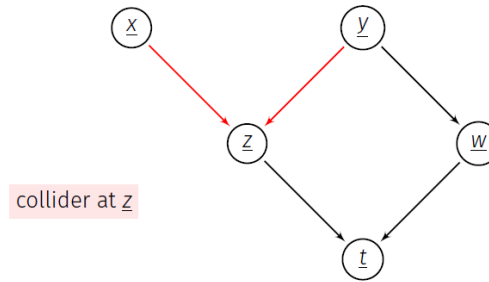


Figure 4.7: Directed graph terminology: collider.

4.2 Graphs in Probability

Definition 4.2.1 (*Bayesian networks*). A Bayesian network is a directed acyclic graph (DAG) (check Def. 4.1.2) that has random variables as vertices.

Theorem 4.2.1 (Representing JPDPs using DAGs). Representing JPDPs using DAGs Consider n random variables $\underline{x}_1, \dots, \underline{x}_n$ with joint pdf (JPDP) $f(x_1, \dots, x_n)$. Let G be a DAG with n vertices representing the random variables. Then, G represents the JPDP f if

$$f(x_1, \dots, x_n) = \prod_{i=1}^n f(x_i | \text{parents of } x_i). \quad (4.1)$$

Theorem 4.2.2 (Markov condition in DAG). The DAG G represents the JPDP f if and only if the Markov condition holds, which states the following:

$$\forall \underline{x}_i : \quad \underline{x}_i \perp \bar{\mathcal{X}}_i | \text{parents of } \underline{x}_i, \quad (4.2)$$

where $\bar{\mathcal{X}}_i$ are all other random variables but the parents and descendants of \underline{x}_i .

Definition 4.2.2 (*D-donnection and d-separation*). Consider a set of n random variables $\underline{x}_1, \dots, \underline{x}_n$ whose JPDP is represented by G . Furthermore, let \mathcal{A} , \mathcal{B} , and \mathcal{C} be mutually exclusive sets with $\mathcal{A}, \mathcal{C} \neq \emptyset$. The sets \mathcal{A} and \mathcal{C} are called *d-connected* given \mathcal{B} if there exists an *undirected* path between some random

variable in \mathcal{A} and some random variable in \mathcal{C} such that

1. every collider along the path is in \mathcal{B} , or has a descendent in \mathcal{B} ;
2. no non-collider along the path belongs in \mathcal{B} .

If \mathcal{A} and \mathcal{C} are not d-connected given \mathcal{B} , then they are called *d-separated given \mathcal{B}* .

Theorem 4.2.3 (D-connection and d-separation). Consider a set of n random variables $\underline{x}_1, \dots, \underline{x}_n$ whose JPDP is represented by G . Furthermore, let \mathcal{A} , \mathcal{B} , and \mathcal{C} be mutually exclusive sets with $\mathcal{A}, \mathcal{C} \neq \emptyset$. Then, the sets \mathcal{A} and \mathcal{C} are *d-separated given \mathcal{B}* if and only if

$$\mathcal{A} \perp \mathcal{C} | \mathcal{B}. \tag{4.3}$$

Part II

Statistical Inference

Chapter 5

Parameter Estimation

5.1 Motivation

Let \underline{x} be a random variable whose CDF, $F(x)$, is *unknown*. The goal is estimate $F(x)$ from n observations of \underline{x} , x_1, \dots, x_n .

Definition 5.1.1 (Empirical CDF). The empirical CDF of \underline{x} from the observations x_1, \dots, x_n is

$$\hat{F}(x) = \frac{1}{n} \times (\# \text{ of observations } x_i \leq x). \quad (5.1)$$

The empirical CDF $\hat{F}(x)$ is the CDF of a discrete random variable with PMF

$$\hat{p}_n(x) = \begin{cases} \frac{1}{n}, & x = x_i, \quad i = 1, \dots, n, \\ 0 & \text{Otherwise.} \end{cases} \quad (5.2)$$

The PDF is then given by

$$\hat{f}_n(x) = \sum_{i=1}^n \frac{1}{n} \delta(x - x_i). \quad (5.3)$$

The empirical PDF $\hat{f}_n(x)$ in (5.3) assigns constant weights to every observation. The empirical PDF can be generalized by giving different weights to different observations. Specifically,

$$\hat{f}_n(x) = \sum_{i=1}^n \frac{1}{n} k(x - x_i), \quad (5.4)$$

where $k(\cdot)$ is the kernel function such that it is

1. symmetric around 0, and
2. $\int_{-\infty}^{\infty} k(x) dx = 1$.

A common choice is the Gaussian kernel. That is, a PDF of $\mathcal{N}(0, 1)$.

5.2 Problem statement

Let $\underline{\mathbf{x}}$ be a random variable. Furthermore, let $f(\underline{\mathbf{x}}; \boldsymbol{\theta})$ be the PDF of $\underline{\mathbf{x}}$ of *known form* (e.g., Gaussian distribution) but depending on *unknown* parameters $\boldsymbol{\theta}$ (e.g., mean and covariance of a Gaussian distribution).

Example 5.2.1. A Gaussian random variable has a PDF parameterized w.r.t. to the mean μ and the variance σ^2 . Thus, the parameters could be $\boldsymbol{\theta} = [\mu, \sigma^2]$.

The goal is to estimate $\boldsymbol{\theta}$ from n observations of $\underline{\mathbf{x}}$: $\mathbf{x}_1, \dots, \mathbf{x}_n$.

Define the data vector

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix}. \quad (5.5)$$

The data vector \mathbf{X} is a realization of the random vector

$$\underline{\mathbf{X}} = \begin{bmatrix} \underline{\mathbf{x}}_1 \\ \vdots \\ \underline{\mathbf{x}}_n \end{bmatrix}, \quad (5.6)$$

where $\underline{\mathbf{x}}_i, i = 1, \dots, n$ are i.i.d.

Approaches

There are two approaches to parameter estimation.

1. *Point estimation*: identify a function $g(\cdot)$ of the sample such that $g(\underline{\mathbf{x}}_1, \dots, \underline{\mathbf{x}}_n)$ is a good (in some sense) approximation of $\boldsymbol{\theta}$.
2. *Interval estimation*: identify two functions $g_1(\cdot)$ and $g_2(\cdot)$ of the sample, such that

$$P(g_1(\underline{\mathbf{x}}_1, \dots, \underline{\mathbf{x}}_n) < \boldsymbol{\theta} < g_2(\underline{\mathbf{x}}_1, \dots, \underline{\mathbf{x}}_n)) = \gamma, \quad (5.7)$$

where γ is the confidence interval.

5.3 Definitions and terminology

Definition 5.3.1 (Estimator). The estimator of $\boldsymbol{\theta}$ is given in the form

$$\hat{\boldsymbol{\theta}} = g(\underline{\mathbf{x}}_1, \dots, \underline{\mathbf{x}}_n). \quad (5.8)$$

The estimator models the behaviour of the estimation method over all possible samples.

Definition 5.3.2 (*Estimate of θ*). The deterministic quantity $g(\mathbf{x}_1, \dots, \mathbf{x}_n)$. That is, the the output of $\hat{\theta}$ for the realizations $\mathbf{x}_1, \dots, \mathbf{x}_n$.

Definition 5.3.3 (*Statistic*). Any function of $\mathbf{x}_1, \dots, \mathbf{x}_n$.

Definition 5.3.4 (*Bias*). The estimator $g(\mathbf{x}_1, \dots, \mathbf{x}_n)$ is called

- *unbiased* if

$$\mathbb{E} [g(\mathbf{x}_1, \dots, \mathbf{x}_n)] = \theta. \quad (5.9)$$

- Otherwise, it's called *biased* with bias

$$\mathbb{E} [g(\mathbf{x}_1, \dots, \mathbf{x}_n)] - \theta. \quad (5.10)$$

Definition 5.3.5 (*Asymptic unbiased*). A biased estimator for which

$$\mathbb{E} [g(\mathbf{x}_1, \dots, \mathbf{x}_n)] \xrightarrow{N \rightarrow \infty} \theta \quad (5.11)$$

is called *asymptotically unbiased*.

Definition 5.3.6 (*Mean-squared best estimator*). For a fixed n , the estimator that minimizes the mean-squared (MS) error

$$\mathbf{e}_n = \mathbb{E} \left[\left\| g(\mathbf{x}_1, \dots, \mathbf{x}_n)^\top g(\mathbf{x}_1, \dots, \mathbf{x}_n) \right\| \right] \quad (5.12)$$

is called the *best* estimator (in the mean-squared sense).

Definition 5.3.7 (*Consistency*). The estimator $g(\mathbf{x}_1, \dots, \mathbf{x}_n)$ is called *consistent* if

$$\lim_{n \rightarrow \infty} P(|g(\mathbf{x}_1, \dots, \mathbf{x}_n)| > \epsilon) = 0, \forall \epsilon > 0, \quad (5.13)$$

or

$$\lim_{n \rightarrow \infty} P(|g(\mathbf{x}_1, \dots, \mathbf{x}_n)| < \epsilon) = 1, \forall \epsilon > 0. \quad (5.14)$$

Theorem 5.3.1 (Sufficient condition for consistency). Let $g(\mathbf{x}_1, \dots, \mathbf{x}_n)$ be an estimator of θ . Then, if

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\left\| g(\mathbf{x}_1, \dots, \mathbf{x}_n)^\top g(\mathbf{x}_1, \dots, \mathbf{x}_n) \right\| \right] = 0, \quad (5.15)$$

then $g(\mathbf{x}_1, \dots, \mathbf{x}_n)$ is a *consistent estimator* of θ .

The converse is not true in general. That is, this is a sufficient but not necessary condition.

Theorem 5.3.2 (Sample mean estimator). Let $\underline{\mathbf{x}}_1, \dots, \underline{\mathbf{x}}_n$ be a random sample of size n and let $\underline{\mathbf{x}}_1, \dots, \underline{\mathbf{x}}_n$ be *uncorrelated*. The *sample mean estimator*

$$\hat{\underline{\mu}}_n = \frac{1}{n} \sum_{i=1}^n \underline{\mathbf{x}}_i \quad (5.16)$$

is

- an *unbiased estimator* of $\underline{\mu}$;
- a *consistent estimator* of $\underline{\mu}$.

Theorem 5.3.3 (Variance estimaton). Consider the sample average estimators of the mean and variance of Gaussian random variable \underline{x} from a random sample $\underline{x}_1, \dots, \underline{x}_n$ of size n :

$$\hat{\underline{\mu}}_n = \frac{1}{n} \sum_{i=1}^n \underline{x}_i, \quad (5.17)$$

$$\hat{\underline{\sigma}}_n^2 = \frac{1}{n} \sum_{i=1}^n \left(\underline{x}_i - \hat{\underline{\mu}}_n \right)^2. \quad (5.18)$$

The random variables $\hat{\underline{\mu}}_n$ and $\hat{\underline{\sigma}}_n^2$ are independent.

5.4 Method of moments

Recall that moments are defined in Def. 1.6.4 and Def. 2.4.1. The m -th moment estimator is given by

$$\hat{m}_m(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n x_i^m. \quad (5.19)$$

Let $\underline{\theta} \in \mathbb{R}^k$ be of size k and let

$$\mathbf{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \quad (5.20)$$

be the samples of the random variable \underline{x} . Then, the idea of method of moments is to estimate the first k moments, thus getting k equations, and then solving the k equations for $\underline{\theta}$.

Example 5.4.1. Estimate the mean μ and variance σ^2 of a random variable \underline{x} from a random sample $\mathbf{X} = \begin{bmatrix} \underline{x}_1 & \dots & \underline{x}_n \end{bmatrix}^T$.

The moments are given by

$$m_1(\mu, \sigma^2) = \mathbb{E}[\underline{x}] = \mu, \quad (5.21)$$

$$m_2(\mu, \sigma^2) = \mathbb{E}[\underline{x}^2] = \mu^2 + \sigma^2. \quad (5.22)$$

The sample average estimates of the moments are given by

$$\hat{m}_1(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n x_i, \quad (5.23)$$

$$\hat{x}_2(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n x_i^2. \quad (5.24)$$

Thus, the parameter estimates can be computed by solving the system of equations

$$\mu = \hat{m}_1(\mathbf{X}) \quad (5.25)$$

$$\mu^2 + \sigma^2 = \hat{m}_2(\mathbf{X}). \quad (5.26)$$

The solution is given by

$$\hat{\mu} = \hat{m}_1(\mathbf{X}) \quad (5.27)$$

$$= \frac{1}{n} \sum_{i=1}^n x_i, \quad (5.28)$$

$$\hat{\sigma}^2 = \hat{m}_2(\mathbf{X}) - \hat{m}_1^2(\mathbf{X}) \quad (5.29)$$

$$= \frac{1}{n} \sum_{i=1}^n x_i^2 - \left(\frac{1}{n} \sum_{i=1}^n x_i \right)^2. \quad (5.30)$$

5.5 Maximum likelihood (ML) estimator

Definition 5.5.1 (*Likelihood*). Given the observations $\underline{\mathbf{x}}_1 = \mathbf{x}'_1, \dots, \underline{\mathbf{x}}_n = \mathbf{x}'_n$, the *likelihood* of $\boldsymbol{\theta}$ is the joint PDF of $\underline{\mathbf{x}}_1, \dots, \underline{\mathbf{x}}_n$ evaluated at the observed points

$$L(\boldsymbol{\theta}; \mathbf{x}'_1, \dots, \mathbf{x}'_n) = f(\mathbf{x}'_1, \dots, \mathbf{x}'_n; \boldsymbol{\theta}) \quad (5.31)$$

$$= f_{\mathbf{x}}(\mathbf{x}'_1; \boldsymbol{\theta}) \cdots f_{\mathbf{x}}(\mathbf{x}'_n; \boldsymbol{\theta}). \quad (5.32)$$

Remark 5.5.1. • The likelihood $f(\mathbf{x}'_1, \dots, \mathbf{x}'_n; \boldsymbol{\theta})$ is treated as a function of $\boldsymbol{\theta}$.

• It is not a PDF. That is,

$$\int_{-\infty}^{\infty} f(\mathbf{x}'_1, \dots, \mathbf{x}'_n; \boldsymbol{\theta}) d\boldsymbol{\theta} \neq 1. \quad (5.33)$$

Definition 5.5.2 (*Maximum-likelihood estimator*). The maximum-likelihood (ML) estimate of θ , $\hat{\theta}_{\text{ML}}$, is the value of θ that maximizes $f(\mathbf{x}'_1, \dots, \mathbf{x}'_n; \theta)$. That is,

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta \in \mathbb{R}^k} L(\theta; \mathbf{x}'_1, \dots, \mathbf{x}'_N) \quad (5.34)$$

$$= \arg \max_{\theta \in \mathbb{R}^k} f(\mathbf{x}'_1, \dots, \mathbf{x}'_n; \theta). \quad (5.35)$$

Remark 5.5.2. Some remarks regarding the ML estimator.

- Log-likelihood function

$$\ell(\mathbf{x}'_1, \dots, \mathbf{x}'_n; \theta) = \log L(\theta; \mathbf{x}'_1, \dots, \mathbf{x}'_N) \quad (5.36)$$

$$= \log f(\mathbf{x}'_1, \dots, \mathbf{x}'_n; \theta). \quad (5.37)$$

- The maximizer of the log-likelihood function is the same maximizer of the likelihood function. That is,

$$\arg \max_{\theta \in \mathbb{R}^k} L(\theta; \mathbf{x}'_1, \dots, \mathbf{x}'_N) = \arg \max_{\theta \in \mathbb{R}^k} \ell(\theta; \mathbf{x}'_1, \dots, \mathbf{x}'_N). \quad (5.38)$$

- The ML estimate is

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta \in \mathbb{R}^k} f(\mathbf{x}'_1, \dots, \mathbf{x}'_n; \theta) \quad (5.39)$$

$$= \arg \max_{\theta \in \mathbb{R}^k} \log f(\mathbf{x}'_1, \dots, \mathbf{x}'_n; \theta) \quad (5.40)$$

$$= \arg \max_{\theta \in \mathbb{R}^k} L(\mathbf{x}'_1, \dots, \mathbf{x}'_n; \theta). \quad (5.41)$$

Remark 5.5.3. Some remarks on ML estimation.

1. Let $\hat{\theta}_{\text{ML}}$ be the ML estimator of θ and let $g(\theta)$ be a function of θ . Then, the ML estimator of $g(\theta)$ is $g(\hat{\theta}_{\text{ML}})$.
2. Limitations of the ML estimation:
 - 2.1. The ML estimate may not exist.
 - 2.2. The ML estimate may not be unique.
3. Maximum likelihood estimate can be finding by minimizing the negative likelihood or negative log-likelihood.

Example 5.5.1. Let $\underline{x} \sim U[0, \theta]$ be a uniformly distributed random variable in the form

$$f(x; \theta) = \begin{cases} \frac{1}{\theta}, & x \in [0, \theta], \\ 0, & \text{otherwise,} \end{cases} \quad (5.42)$$

where θ is some unknown parameter. Find the ML estimator of θ from a random sample of size N .

1. First, get the joint PDF of all the random samples

$$f(x_1, \dots, x_N; \theta) = f(x_1; \theta) \cdots f(x_N; \theta) \quad (5.43)$$

$$= \begin{cases} \frac{1}{\theta^N}, & x_i \in [0, \theta], i = 1, \dots, N, \\ 0, & \text{otherwise,} \end{cases} \quad (5.44)$$

where the fact that the random samples \underline{x}_i for $i = 1, \dots, N$ are i.i.d.

2. The likelihood of θ is then given as a function of the samples of $\underline{x}_i = x'_i$. Specifically, the likelihood is

$$L(\theta; x'_1, \dots, x'_N) = f(\underline{x}_1 = x'_1, \dots, \underline{x}_N = x'_N; \theta) \quad (5.45)$$

$$= \begin{cases} \frac{1}{\theta^N}, & x'_i \in [0, \theta], i = 1, \dots, N, \\ 0, & \text{otherwise} \end{cases} \quad (5.46)$$

$$= \begin{cases} \frac{1}{\theta^N}, & \max_i(x_i) \leq \theta, \\ 0, & \text{otherwise.} \end{cases} \quad (5.47)$$

Note that it was not necessary to specify the lower bound $\min_i(x_i) \geq 0$ since the samples are drawn from a uniform distribution $U[0, \theta]$ which is bounded from below by 0. Thus, $P(\underline{x} < 0) = 0$. Therefore, it is impossible to have a sample $x'_i < 0$.

3. Maximize the likelihood function. Maximizing (5.47) is a *constrained* optimization problem that entails maximizing $\frac{1}{\theta^N}$.

Maximizing $\frac{1}{\theta^N}$ is the same as minimizing θ . However, since θ is bounded from below by the constraint in (5.47), then the minimum value of θ is $\max_i(x_i)$. Thus,

$$\hat{\theta}_{\text{ML}} = \max_i(x_i). \quad (5.48)$$

△

5.6 Bayesian parameter estimation

In the previous sections, the approaches used were the *classical* or *frequentist* approaches. The classical approach

1. treats the unknown parameter as deterministic parameter, and

2. assumes no prior information on the unknown parameter was available.

For Bayesian approach on the other hand

1. assumes that unknown parameter is random with known pdf $f(\theta)$, called the *prior*;
2. the JPDF of random sample given $\underline{\theta} = \theta$ (likelihood) is given by

$$f(\mathbf{x}_1, \dots, \mathbf{x}_n | \theta); \quad (5.49)$$

3. The estimation problem becomes a prediction problem; predict $\underline{\theta}$ from a correlated (with $\underline{\theta}$) sample $\underline{\mathbf{X}}$.

The mean-squared (or Bayes) estimate $\hat{\theta}$ of $\underline{\theta}$ given $\underline{\mathbf{X}} = \mathbf{X}$ is

$$\hat{\theta} = \mathbb{E}[\underline{\theta} | \mathbf{X}] \quad (5.50)$$

$$= \int_{-\infty}^{\infty} \theta f(\theta | \mathbf{X}) d\theta, \quad (5.51)$$

where

$$f(\theta | \mathbf{X}) = \frac{f(\mathbf{X} | \theta) f(\theta)}{f(\mathbf{X})} \quad (5.52)$$

$$= \frac{f(\mathbf{X} | \theta) f(\theta)}{\int_{-\infty}^{\infty} f(\mathbf{X} | \theta) f(\theta) d\theta}. \quad (5.53)$$

Terminology

- The mean-squared estimator is also referred to as the Bayes estimator.
- $f(\theta)$: *prior* pdf of $\underline{\theta}$ (before observing the sample).
- $f(\theta | \mathbf{X})$: *posterior* pdf of θ (after observing the sample).
- In [5], the term $f(\mathbf{X} | \theta)$ is sometimes referred to as the likelihood of \mathbf{X} given θ . This is not to be confused with the likelihood from the ML estimator. In estimation, the term $f(\mathbf{X} | \theta)$ can also be referred to as an *observation model*.

Remark 5.6.1. Some remarks about Bayesian estimation.

- $\mathbb{E}[\underline{\theta} | \mathbf{X}]$ from (5.50) is difficult to solve for two reasons. First, the denominator of the posterior pdf involves evaluating an integral over the entire sample space of θ . Second, evaluating the expectation $\mathbb{E}[\cdot]$ requires *another* integral to be evaluated over the entire sample space.
- Maximum a-posteriori (MAP) estimate maximizes the *posterior* pdf $f(\theta | \mathbf{X})$ with respect to θ .

Since the denominator of the prior is constant, then it can be omitted. That is,

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} f(\theta | \mathbf{X}) \quad (5.54)$$

$$= \arg \max_{\theta} \frac{f(\mathbf{X} | \theta) f(\theta)}{f(\mathbf{X})} \quad (5.55)$$

$$= \arg \max_{\theta} f(\mathbf{X} | \theta) f(\theta). \quad (5.56)$$

5.7 Maximum a-posteriori (MAP) estimator

This estimator is an *approximation* to the Bayes estimator discussed in the previous section. It is easier (less computationally intensive) to compute.

Definition 5.7.1 (*Maximum a-posteriori estimator*). Maximum a-posteriori (MAP) estimate maximizes the *posterior* pdf $f(\theta | \mathbf{X})$ with respect to θ . Since the denominator of the prior is constant, then it can be omitted. That is,

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} f(\theta | \mathbf{X}) \quad (5.57)$$

$$= \arg \max_{\theta} \frac{f(\mathbf{X} | \theta) f(\theta)}{f(\mathbf{X})} \quad (5.58)$$

$$= \arg \max_{\theta} f(\mathbf{X} | \theta) f(\theta). \quad (5.59)$$

5.7.1 MAP estimator of a Markov normally distributed random variable

In this section, the MAP estimator for normally distributed variables that follow a Markov chain will be derived. For each problem, a process model is provided in the form

$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k), \quad (5.60)$$

where $\mathbf{w}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_k)$ is the process noise and the covariance matrix \mathbf{Q}_k is known. A prior is distributed according to $\hat{\mathbf{x}}_0 \sim \mathcal{N}(\hat{\mathbf{x}}_1, \check{\mathbf{P}}_0)$.

Furthermore, a measurement model is given in the form

$$\mathbf{y}_k = \mathbf{g}(\mathbf{x}_k, \mathbf{v}_k), \quad (5.61)$$

where $\mathbf{v}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_k)$ is the Gaussian measurement noise, where the covariance matrix \mathbf{R}_k is known.

The goal is to estimate states' $\mathbf{x}_{1:K}$ parameters using the realizations of the prior \mathbf{x}_0 , the realizations of the interoceptive measurements \mathbf{u}_k and realizations of the exteroceptive measurements \mathbf{y}_k . If the RVs are Gaussian, then they can be parametrized by a mean and a covariance (i.e., $\mathbf{x}_k \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$).

Remark 5.7.1. The interoceptive measurement \mathbf{u}_k is a random variable that can be assumed to be dis-

tributed according to

$$\underline{\mathbf{u}}_k \sim \mathcal{N}(\underline{\mathbf{u}}_k, \mathbf{Q}_k^u). \quad (5.62)$$

The measurement is assumed to be random due to measurement noise. However, the measurement noise can be embedded into the process noise \mathbf{Q} . Therefore, it is possible to assume that the interoceptive measurements are not random, but the interoceptive measurement process noise is indeed random and will be embedded with the process noise $\underline{\mathbf{w}}_k$.

Remark 5.7.2. The states will be assumed to be normally distributed. That is, $\underline{\mathbf{x}}_k \sim \mathcal{N}(\underline{\mathbf{x}}_k, \mathbf{P}_k)$. Therefore, the parameters to be estimated should

$$\tilde{\boldsymbol{\theta}} = \{\mathbf{x}_0, \dots, \mathbf{x}_N, \mathbf{P}_0, \dots, \mathbf{P}_N\}. \quad (5.63)$$

However, to make things easier, the covariances $\mathbf{P}_{0:N}$ will not be estimated (they will be computed after the MAP estimate is computed). Thus, the parameters that'll be estimated is given by

$$\boldsymbol{\theta} = \{\mathbf{x}_0, \dots, \mathbf{x}_N\}. \quad (5.64)$$

The realizations are then

$$\mathbf{X} = \{\check{\mathbf{x}}_0, \mathbf{u}_{0:N-1}, \mathbf{y}_{1:N}\}. \quad (5.65)$$

For the MAP estimator, the PDF of interest is

$$f(\boldsymbol{\theta} | \mathbf{X}) = f(\mathbf{x}_0, \dots, \mathbf{x}_N | \check{\mathbf{x}}_0, \mathbf{u}_{0:N-1}, \mathbf{y}_{1:N}) \quad (5.66)$$

$$= \frac{f(\mathbf{y}_{1:N} | \mathbf{x}_{0:N}, \check{\mathbf{x}}_0, \mathbf{u}_{0:N-1}) f(\mathbf{x}_{0:N} | \check{\mathbf{x}}_0, \mathbf{u}_{0:N-1})}{f(\mathbf{y}_{1:N} | \check{\mathbf{x}}_0, \mathbf{u}_{0:N-1})} \quad (5.67)$$

$$= \eta f(\mathbf{y}_1 | \mathbf{x}_1) \cdots f(\mathbf{y}_N | \mathbf{x}_N) f(\mathbf{x}_N | \mathbf{x}_{N-1}, \mathbf{u}_{N-1}) \cdots f(\mathbf{x}_1 | \mathbf{x}_0, \mathbf{u}_0) f(\mathbf{x}_0 | \check{\mathbf{x}}_0) \quad (5.68)$$

$$= \eta \prod_{k=1}^N f(\mathbf{y}_k | \mathbf{x}_k) \prod_{k=1}^N f(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{u}_{k-1}) f(\mathbf{x}_0 | \check{\mathbf{x}}_0), \quad (5.69)$$

where η is a normalizing parameter that is independent of the design variables $(\mathbf{x}_{0:N})$ and the following assumptions were used

1. The measurements \mathbf{y}_k depends only on \mathbf{x}_k (and measurement noise $\underline{\mathbf{n}}_k$) as shown in the measurement model (5.61). Thus, the random variable $\underline{\mathbf{y}}_k$ is independent of all other random variables.
2. The state $\underline{\mathbf{x}}_k$ depends only on $\underline{\mathbf{x}}_{k-1}$ and \mathbf{u}_{k-1} (and process noise $\underline{\mathbf{w}}_{k-1}$) as shown in (5.60). Therefore, the random variables $\mathbf{x}_{0:N}$ follow the **Markov sequence**: the future states $\mathbf{x}_{k+1:N}$ are independent of the past states $\mathbf{x}_{0:k-1}$ given the present state \mathbf{x}_k . The Markov sequence is dictated by the process model (5.60).

If $\underline{\mathbf{x}}_{0:N}$ are assumed to be [marginally] Gaussian, then it is easier to minimize the negative-log of (5.66) than maximize (5.66) directly.

Taking the negative-log of (5.69) gives

$$\begin{aligned}
 -\log f(\mathbf{x}_{0:N} | \check{\mathbf{x}}_0, \mathbf{u}_{0:N-1}, \mathbf{y}_{1:N}) &= \frac{1}{2} \|\mathbf{x}_0 - \check{\mathbf{x}}_0\|_{\mathbf{P}_0}^2 + \frac{1}{2} \sum_{k=1}^N \|\mathbf{y}_k - \mathbf{g}(\mathbf{x}_k, \mathbf{0})\|_{\mathbf{R}_k}^2 + \\
 &\quad \frac{1}{2} \sum_{k=1}^N \|\mathbf{x}_k - \mathbf{f}(\mathbf{x}_{k-1}, \mathbf{u}_{k-1}, \mathbf{0})\| + \gamma,
 \end{aligned} \tag{5.70}$$

where

$$\|\mathbf{z}\|_{\Sigma}^2 \triangleq \mathbf{z}^\top \Sigma^{-1} \mathbf{z}, \tag{5.71}$$

and γ are constant terms that are independent of the design variables thus they will not affect the optimization.

Therefore, the MAP estimate of $\hat{\mathbf{x}}_{0:N}^{\text{MAP}}$ is the solution to

$$\hat{\mathbf{x}}_{0:N}^{\text{MAP}} = \arg \min_{\mathbf{x}_{0:N} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{x}_0 - \check{\mathbf{x}}_0\|_{\mathbf{P}_0}^2 + \frac{1}{2} \sum_{k=1}^N \|\mathbf{y}_k - \mathbf{g}(\mathbf{x}_k, \mathbf{0})\|_{\mathbf{R}_k}^2 + \frac{1}{2} \sum_{k=1}^N \|\mathbf{x}_k - \mathbf{f}(\mathbf{x}_{k-1}, \mathbf{u}_{k-1}, \mathbf{0})\|, = \arg \min_{\mathbf{x}_{0:N} \in \mathbb{R}^n} J \tag{5.72}$$

which is a (nonlinear) weighted least squares problem.

Remark 5.7.3. In (5.70), it was assumed that $\mathbf{w}_k = \mathbf{n}_k = \mathbf{0}$. I'm not not sure why this assumption is done. The way I think of it is that the noise \mathbf{w}_k is embedded into the interoceptive measurement \mathbf{u}_k . Similarly, the noise in \mathbf{n}_k is embedded in the exteroceptive measurement \mathbf{y}_k .

5.7.2 Expressing the nonlinear least squares problem in matrix form

The objective function can be expressed as

$$J(\mathbf{x}_{0:N}) = \frac{1}{2} \|\mathbf{x}_0 - \check{\mathbf{x}}_0\|_{\mathbf{P}_0}^2 + \frac{1}{2} \sum_{k=1}^N \|\mathbf{x}_k - \mathbf{f}(\mathbf{x}_{k-1}, \mathbf{u}_{k-1}, \mathbf{0})\| + \frac{1}{2} \sum_{k=1}^N \|\mathbf{y}_k - \mathbf{g}(\mathbf{x}_k, \mathbf{0})\|_{\mathbf{R}_k}^2 \tag{5.73}$$

$$= \frac{1}{2} \mathbf{e}(\mathbf{x}_{0:N})^\top \mathbf{W} \mathbf{e}(\mathbf{x}_{0:N}), \tag{5.74}$$

where

$$\mathbf{e}(\mathbf{x}_{0:N}) = \begin{bmatrix} \mathbf{x}_0 - \check{\mathbf{x}}_0 \\ \mathbf{x}_1 - \mathbf{f}(\mathbf{x}_0, \mathbf{u}_0, \mathbf{0}) \\ \vdots \\ \mathbf{x}_N - \mathbf{f}(\mathbf{x}_{N-1}, \mathbf{u}_{N-1}, \mathbf{0}) \\ \mathbf{y}_1 - \mathbf{g}(\mathbf{x}_1, \mathbf{0}) \\ \vdots \\ \mathbf{y}_N - \mathbf{g}(\mathbf{x}_N, \mathbf{0}) \end{bmatrix} \tag{5.75}$$

is the *error function*, and

$$\mathbf{W}^{-1} = \begin{bmatrix} \check{\mathbf{P}}_0 & & & & \\ & \mathbf{Q}_0 & & & \\ & & \ddots & & \\ & & & \mathbf{Q}_N & \\ & & & & \mathbf{R}_1 \\ & & & & & \ddots \\ & & & & & & \mathbf{R}_N \end{bmatrix} \quad (5.76)$$

is the *weight matrix*.

5.7.3 Covariance on MAP estimate

In the optimization problem above, the covariance was not estimated. However, it can be estimated by using the MAP state estimates $\hat{\mathbf{x}}_{0:N}^{\text{MAP}}$. The *joint* covariance can be approximated by

$$\text{Cov} [\hat{\mathbf{x}}_{0:N}^{\text{MAP}}] = \left(\mathbf{J} (\hat{\mathbf{x}}_{0:N}^{\text{MAP}})^T \mathbf{W} \mathbf{J} (\hat{\mathbf{x}}_{0:N}^{\text{MAP}}) \right)^{-1}, \quad (5.77)$$

where

$$\mathbf{J} (\hat{\mathbf{x}}_{0:N}^{\text{MAP}}) = \left. \frac{d\mathbf{e}(\mathbf{x}_{0:N})}{d\mathbf{x}_{0:N}} \right|_{\mathbf{x}_{0:N}=\hat{\mathbf{x}}_{0:N}^{\text{MAP}}}. \quad (5.78)$$

is the Jacobian of the error function (5.75) w.r.t. the design variables, evaluated at the MAP estimate $\hat{\mathbf{x}}_{0:N}^{\text{MAP}}$.

Appendices

Appendix A

Linear algebra

A.1 Schur complement

Theorem A.1.1 (Schur's complement). Let

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \quad (\text{A.1})$$

be a square matrix. If \mathbf{A}_{11} is nonsingular, then the Schur complement of the block \mathbf{A}_{11} of the matrix \mathbf{A} is the matrix

$$\mathbf{A}/\mathbf{A}_{11} := \mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12}. \quad (\text{A.2})$$

Similarly, if \mathbf{A}_{22} is nonsingular, then the Schur complement of the block \mathbf{A}_{22} of the matrix \mathbf{A} is the matrix

$$\mathbf{A}/\mathbf{A}_{22} := \mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21}. \quad (\text{A.3})$$

A.1.1 Application to solving linear equations

Given a system of equations

$$\begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix}. \quad (\text{A.4})$$

Then what is the solution \mathbf{x}_1 without computing \mathbf{x}_2 (or without solving the big system of equations all at once)?

The solution is given by using the Schur complement of the block \mathbf{A}_{22} of the matrix \mathbf{A} . Specifically, pre-multiply the left and right sides of (A.4) by the matrix

$$\mathbf{L} = \begin{bmatrix} \mathbf{1} & -\mathbf{A}_{12}\mathbf{A}_{22}^{-1} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}. \quad (\text{A.5})$$

Then, the result is

$$\begin{bmatrix} \mathbf{1} & -\mathbf{A}_{12}\mathbf{A}_{22}^{-1} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{1} & -\mathbf{A}_{12}\mathbf{A}_{22}^{-1} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix} \quad (\text{A.6})$$

$$\begin{bmatrix} \mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21} & \mathbf{0} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{b}_1 - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{b}_2 \\ \mathbf{b}_2 \end{bmatrix}. \quad (\text{A.7})$$

The solution \mathbf{x}_1 is then given by solving

$$(\mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21}) \mathbf{x}_1 = \mathbf{b}_1 - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{b}_2. \quad (\text{A.8})$$

A.1.2 Application in probability

A.1.2.1 Conditioning using covariance matrix

Given a the random variable

$$\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \bar{\mathbf{x}}_1 \\ \bar{\mathbf{x}}_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{12}^T & \Sigma_{22} \end{bmatrix} \right). \quad (\text{A.9})$$

Then, the marginal covariance of \mathbf{x}_1 is Σ_{11} . However, the covariance of \mathbf{x}_1 *given* $\mathbf{x}_2 = \hat{\mathbf{x}}_2$ is given by the Schur complement. Specifically,

$$\mathbb{E} [\mathbf{x}_1 | \mathbf{x}_2 = \hat{\mathbf{x}}_2] = \bar{\mathbf{x}}_1 + \Sigma_{12}\Sigma_{22}^{-1} (\hat{\mathbf{x}}_2 - \bar{\mathbf{x}}_2), \quad (\text{A.10})$$

$$\text{Cov} [\mathbf{x}_1 | \mathbf{x}_2 = \hat{\mathbf{x}}_2] = \Sigma / \Sigma_{22} \quad (\text{A.11})$$

$$= \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{12}^T. \quad (\text{A.12})$$

Note that this is equivalent to solving the linear system of equations

$$\begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{12}^T & \Sigma_{22} \end{bmatrix} \begin{bmatrix} \mathbb{E} [\mathbf{x}_1 | \mathbf{x}_2 = \hat{\mathbf{x}}_2] \\ \mathbf{x}_2 \end{bmatrix} = \begin{bmatrix} \bar{\mathbf{x}}_1 \\ \bar{\mathbf{x}}_2 - \hat{\mathbf{x}}_2 \end{bmatrix}. \quad (\text{A.13})$$

I'm not sure what's the relation between (A.13) and (A.10).

A.1.2.2 Conditioning using information matrix

On the other hand, if we are given the same random variable

$$\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \bar{\mathbf{x}}_1 \\ \bar{\mathbf{x}}_2 \end{bmatrix}, \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{12}^T & \mathbf{A}_{22} \end{bmatrix}^{-1} \right), \quad (\text{A.14})$$

where

$$\mathbf{A} = \Sigma^{-1} \quad (\text{A.15})$$

is the *information matrix*. Then, conditioning is easily done by partitioning the system. That is,

$$\text{Cov} [\underline{\mathbf{x}}_1 | \underline{\mathbf{x}}_2 = \hat{\mathbf{x}}_2] = \mathbf{A}_{11}^{-1}. \quad (\text{A.16})$$

On the other hand, marginalization is more “difficult” and is given by the Schur complement. Specifically,

$$\text{Cov} [\underline{\mathbf{x}}_1] = (\mathbf{A} / \mathbf{A}_{22})^{-1} \quad (\text{A.17})$$

$$= \left(\mathbf{A}_{11} - \mathbf{A}_{12} \mathbf{A}_{22}^{-1} \mathbf{A}_{12}^T \right)^{-1}. \quad (\text{A.18})$$

Appendix B

Numerically sampling from a normal distribution

Say we want to sample realizations of a normal random variable $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_x, \boldsymbol{\Sigma})$, given the mean $\boldsymbol{\mu}_x$ and covariance $\boldsymbol{\Sigma}_x$. How can this be done using MATLAB's `randn` function?

B.1 MATLAB's `randn` function

The `randn` function produces a realization of a normal random variable with zero mean and unit variance. How to compute the variance of such variable? Well, that depends on how the realizations are stored in MATLAB. That is, the realizations can be stored as columns or as rows.

Say the random variable has n degrees of freedom and we want m realizations. Thus, if the realizations are stored as columns, then the matrix is of size $n \times m$. On the other hand, if the realizations are stored as rows, then the matrix is of size $m \times n$. Specifically,

$$\mathbf{x}_{\text{unit,col}} = \text{randn}(n, m) \in \mathbb{R}^{n \times m}, \quad (\text{B.1})$$

$$\mathbf{x}_{\text{unit,row}} = \text{randn}(m, n) \in \mathbb{R}^{m \times n}. \quad (\text{B.2})$$

The covariance of \mathbf{x}_{unit} can be approximated numerically as

$$\text{Cov}[\mathbf{x}_{\text{unit}}] = \boldsymbol{\Sigma}_{\text{unit}} \quad (\text{B.3})$$

$$= \mathbf{1}_{n \times n} \quad (\text{B.4})$$

$$\approx \frac{1}{m-1} \mathbf{x}_{\text{unit,col}} \mathbf{x}_{\text{unit,col}}^T \quad (\text{B.5})$$

$$= \frac{1}{m-1} \mathbf{x}_{\text{row,unit}}^T \mathbf{x}_{\text{row,unit}}. \quad (\text{B.6})$$

Well, what if we want a realization of a normal random variable with zero mean but non-unit variance $\boldsymbol{\Sigma}_x$?

The covariance is then

$$\text{Cov} [\mathbf{x}] = \Sigma_x \quad (\text{B.7})$$

$$= \mathbf{L} \underbrace{\Sigma_{\text{unit}}}_{\mathbf{I}} \mathbf{L}^T \quad (\text{B.8})$$

$$= \mathbf{R}^T \underbrace{\Sigma_{\text{unit}}}_{\mathbf{I}} \mathbf{R}, \quad (\text{B.9})$$

where \mathbf{L} and \mathbf{R} are the lower- and upper-triangular matrices obtained from the Cholesky factorization. So which of these two should we use? Well, the answer depends on how the realizations of $\underline{\mathbf{x}}$ are stored.

If the realizations are stored in columns, then from (B.5),

$$\text{Cov} [\mathbf{x}] = \Sigma_x \quad (\text{B.10})$$

$$= \mathbf{L} \Sigma_{\text{unit}} \mathbf{L}^T \quad (\text{B.11})$$

$$\approx \mathbf{L} \mathbf{x}_{\text{unit,col}} \mathbf{x}_{\text{unit,col}}^T \mathbf{L}^T \quad (\text{B.12})$$

$$= \underbrace{(\mathbf{L} \mathbf{x}_{\text{unit,col}})}_{\mathbf{x}_{\text{col}}} (\mathbf{L} \mathbf{x}_{\text{unit,col}})^T \quad (\text{B.13})$$

$$= \mathbf{x}_{\text{col}} \mathbf{x}_{\text{col}}^T, \quad (\text{B.14})$$

where

$$\mathbf{x}_{\text{col}} = \mathbf{L} \mathbf{x}_{\text{unit,col}} \quad (\text{B.15})$$

$$= \mathbf{L} \text{randn}(n, m). \quad (\text{B.16})$$

It should be noted that if \mathbf{R} is used in place of \mathbf{L} , then the results would not be the same in general. The upper triangular matrix \mathbf{R} can still be used by setting

$$\mathbf{L} = \mathbf{R}^T. \quad (\text{B.17})$$

Specifically,

$$\text{Cov} [\mathbf{x}] = \Sigma_x \quad (\text{B.18})$$

$$= \mathbf{R}^T \Sigma_{\text{unit}} \mathbf{R} \quad (\text{B.19})$$

$$\approx \mathbf{R}^T \mathbf{x}_{\text{unit,row}}^T \mathbf{x}_{\text{unit,row}} \mathbf{R} \quad (\text{B.20})$$

$$= (\mathbf{x}_{\text{unit,row}} \mathbf{R})^T \underbrace{(\mathbf{x}_{\text{unit,row}} \mathbf{R})}_{\mathbf{x}_{\text{row}}} \quad (\text{B.21})$$

$$= \mathbf{x}_{\text{row}}^T \mathbf{x}_{\text{row}}, \quad (\text{B.22})$$

where

$$\mathbf{x}_{\text{row}} = \mathbf{x}_{\text{unit,row}} \mathbf{R} \quad (\text{B.23})$$

$$= \text{randn}(m, n) \mathbf{R}. \quad (\text{B.24})$$

Summary:

$$\Sigma_x = \mathbf{R}^T \mathbf{R} \quad (\text{B.25})$$

$$= \mathbf{L} \mathbf{L}^T, \quad (\text{B.26})$$

$$\mathbf{x}_{\text{unit,col}} = \text{randn}(n, m) \quad (\text{B.27})$$

$$= \mathbf{x}_{\text{unit,row}}^T, \quad (\text{B.28})$$

$$\mathbf{x}_{\text{col}} = \mathbf{L} \mathbf{x}_{\text{unit,col}}, \quad (\text{B.29})$$

$$= \mathbf{x}_{\text{row}}^T, \quad (\text{B.30})$$

$$\mathbf{x}_{\text{row}} = \mathbf{x}_{\text{unit,row}} \mathbf{R}. \quad (\text{B.31})$$

B.2 Sampling (another derivation)

Assume we have N realizations $\mathbf{x} \in \mathbb{R}^{n \times N}$ of the random variable $\underline{\mathbf{x}} \in \mathbb{R}^n$. For simplicity, we'll assume the random variable has zero mean. Let's assume that there exists a matrix \mathbf{R} such that

$$\underline{\mathbf{x}} = \mathbf{R} \underline{\mathbf{x}}_0, \quad (\text{B.32})$$

where $\underline{\mathbf{x}}_0 \sim N(\mathbf{0}, \mathbf{1})$. Then,

$$\Sigma = \mathbb{E}[\underline{\mathbf{x}} \underline{\mathbf{x}}^T] \quad (\text{B.33})$$

$$= \mathbf{R} \mathbb{E}[\underline{\mathbf{x}}_0 \underline{\mathbf{x}}_0^T] \mathbf{R}^T \quad (\text{B.34})$$

$$= \mathbf{R} \mathbf{R}^T. \quad (\text{B.35})$$

well, how can we construct \mathbf{R} ? Decompose the covariance matrix using eigendecomposition. Specifically,

$$\Sigma = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T \quad (\text{B.36})$$

$$= \underbrace{\mathbf{U} \sqrt{\mathbf{\Lambda}}}_{\mathbf{R}} \underbrace{\mathbf{1} \sqrt{\mathbf{\Lambda}} \mathbf{U}^T}_{\mathbf{R}^T} \quad (\text{B.37})$$

$$= \mathbf{R} \mathbf{R}^T, \quad (\text{B.38})$$

where \mathbf{U} and $\mathbf{\Lambda}$ are the eigenvector and the eigenvalue matrices, respectively. Thus,

$$\mathbf{R} = \mathbf{U} \sqrt{\mathbf{\Lambda}}. \quad (\text{B.39})$$

Therefore, the random variable $\underline{\mathbf{x}}$ can be sampled by sampling from the unit normal random variable $\underline{\mathbf{x}}_0$ and scaling it using

$$\mathbf{x} = \mathbf{U} \sqrt{\mathbf{\Lambda}} \mathbf{x}_0. \quad (\text{B.40})$$

B.3 Covariance ellipses

To plot covariance ellipses (in 2D), eigendecomposition is used. Specifically, given a covariance matrix Σ , then decompose

$$\Sigma = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T, \quad (\text{B.41})$$

where $\mathbf{U}^{-1} = \mathbf{U}^T$ because $\Sigma = \Sigma^T$ is symmetric. The standard deviation ellipses¹ are then computed as

$$\boldsymbol{\sigma} = \mathbf{U}\sqrt{\mathbf{\Lambda}} \begin{bmatrix} \cos(\theta) \\ \sin(\theta) \end{bmatrix}, \quad \theta \in [0, 2\pi). \quad (\text{B.42})$$

Figure B.1 shows an example of such ellipses.

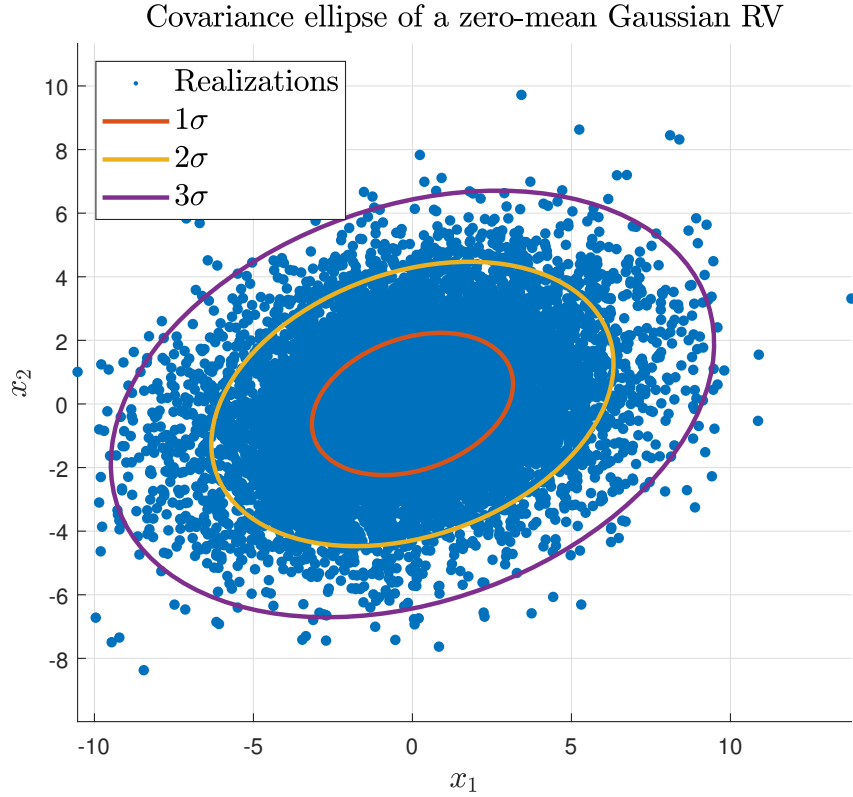


Figure B.1: Covariance ellipses (actually, standard deviation ellipses) around $N = 10^4$ samples.

¹I'm not sure about the naming here. But the ellipse computed represents the 1σ bounds.

Bibliography

- [1] I. Psaromiligkos, “Slides from ECSE 509: Probability and random signals 2 taught in fall 2019 at McGill university.”
- [2] K. P. Murphy, *Machine learning: a probabilistic perspective*.
- [3] F. Gustafsson and G. Hendeby, “On nonlinear transformations of stochastic variables and its application to nonlinear filtering,” in *2008 IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, pp. 3617–3620, ISSN: 1520-6149. [Online]. Available: <http://ieeexplore.ieee.org/document/4518435/>
- [4] S. Sarkka, *Bayesian Filtering and Smoothing*. Cambridge University Press. [Online]. Available: <http://ebooks.cambridge.org/ref/id/CBO9781139344203>
- [5] T. D. Barfoot, *State Estimation for Robotics: A Matrix Lie Group Approach*. Cambridge University Press.