

Introduction to Probability

SECOND EDITION

Dimitri P. Bertsekas and John N. Tsitsiklis

Massachusetts Institute of Technology

Selected Summary Material – All Rights Reserved

WWW site for book information and orders

<http://www.athenasc.com>



Athena Scientific, Belmont, Massachusetts

1

Sample Space and Probability

Excerpts from Introduction to Probability: Second Edition

by Dimitri P. Bertsekas and John N. Tsitsiklis ©

Massachusetts Institute of Technology

Contents

1.1. Sets	p. 3
1.2. Probabilistic Models	p. 6
1.3. Conditional Probability	p. 18
1.4. Total Probability Theorem and Bayes' Rule	p. 28
1.5. Independence	p. 34
1.6. Counting	p. 44
1.7. Summary and Discussion	p. 51
Problems	p. 53

sample space

- discrete
 - finite : roll a dice
 - countable : $P(n) = \frac{1}{2^n}$
 - (infinite)
- continuous, infinite : $0 \leq x, y \leq 1$

2

From Introduction to Probability, by Bertsekas and Tsitsiklis

Chap. 1

1.1 SETS

1.2 PROBABILISTIC MODELS

Elements of a Probabilistic Model

- The **sample space** Ω , which is the set of all possible **outcomes** of an experiment.
- The **probability law**, which assigns to a set A of possible outcomes (also called an **event**) a nonnegative number $\mathbf{P}(A)$ (called the **probability** of A) that encodes our knowledge or belief about the collective “likelihood” of the elements of A . The probability law must satisfy certain properties to be introduced shortly.

$$\begin{cases} A \cap B : \text{intersection} \\ A \cup B : \text{union} \end{cases}$$

• A and B are disjoint:

$$A \cap B = \emptyset$$

• A^c : complement of set A

Probability Axioms

1. (**Nonnegativity**) $\mathbf{P}(A) \geq 0$, for every event A .
2. (**Additivity**) If A and B are two disjoint events, then the probability of their union satisfies

$$\mathbf{P}(A \cup B) = \mathbf{P}(A) + \mathbf{P}(B).$$

More generally, if the sample space has an infinite number of elements and A_1, A_2, \dots is a sequence of disjoint events, then the probability of their union satisfies A_1, A_2, \dots should be countable

$$\mathbf{P}(A_1 \cup A_2 \cup \dots) = \mathbf{P}(A_1) + \mathbf{P}(A_2) + \dots$$

3. (**Normalization**) The probability of the entire sample space Ω is equal to 1, that is, $\mathbf{P}(\Omega) = 1$.

Consequences

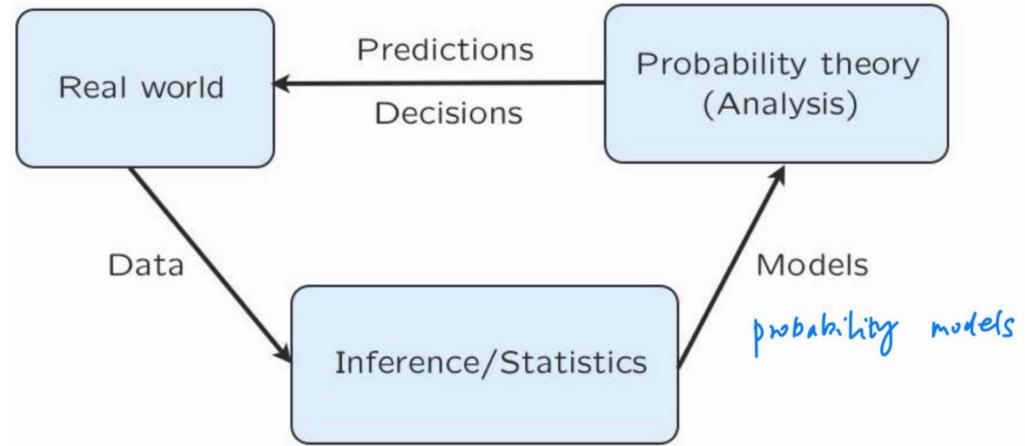
$$\mathbf{P}(A) \leq 1$$

$$\mathbf{P}(\emptyset) = 0$$

$$\mathbf{P}(A) + \mathbf{P}(A^c) = 1$$

if A_1, \dots, A_k disjoint, \Rightarrow

$$\mathbf{P}(A_1 \cup \dots \cup A_k) = \sum_{i=1}^k \mathbf{P}(A_i)$$



Discrete Probability Law

If the sample space consists of a finite number of possible outcomes, then the probability law is specified by the probabilities of the events that consist of a single element. In particular, the probability of any event $\{s_1, s_2, \dots, s_n\}$ is the sum of the probabilities of its elements:

$$\mathbf{P}(\{s_1, s_2, \dots, s_n\}) = \mathbf{P}(s_1) + \mathbf{P}(s_2) + \dots + \mathbf{P}(s_n).$$

Discrete Uniform Probability Law

If the sample space consists of n possible outcomes which are equally likely (i.e., all single-element events have the same probability), then the probability of any event A is given by

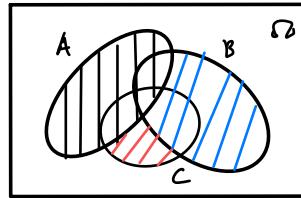
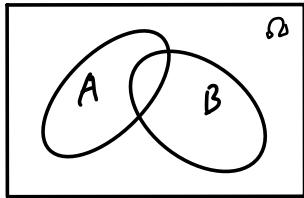
$$\mathbf{P}(A) = \frac{\text{number of elements of } A}{n} = \frac{k}{n}$$

- **Uniform probability law: Probability = Area**

Some Properties of Probability Laws

Consider a probability law, and let A , B , and C be events.

- If $A \subset B$, then $\mathbf{P}(A) \leq \mathbf{P}(B)$.
- $\mathbf{P}(A \cup B) = \mathbf{P}(A) + \mathbf{P}(B) - \mathbf{P}(A \cap B)$.
- $\mathbf{P}(A \cup B) \leq \mathbf{P}(A) + \mathbf{P}(B)$.
- $\mathbf{P}(A \cup B \cup C) = \mathbf{P}(A) + \mathbf{P}(A^c \cap B) + \mathbf{P}(A^c \cap B^c \cap C)$.



Probability calculation steps

- Specify the sample space *possible outcomes*
- Specify a probability law *(somewhat arbitrary choice)*
- Identify an event of interest *useful to describe in a picture*)
- Calculate... •

1.3 CONDITIONAL PROBABILITY

Conditional probability
satisfies 3 probability
axioms

$$\left\{ \begin{array}{l} P(A|B) \geq 0 \\ P(B|B) = 1 \\ \text{Countable additivity} \end{array} \right.$$

\Rightarrow all theorems for ordinary probability apply to conditional probability

Properties of Conditional Probability

- The conditional probability of an event A , given an event B with $P(B) > 0$, is defined by

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

and specifies a new (conditional) probability law on the same sample space Ω . In particular, all properties of probability laws remain valid for conditional probability laws.

- Conditional probabilities can also be viewed as a probability law on a new universe B , because all of the conditional probability is concentrated on B .
- If the possible outcomes are finitely many and equally likely, then

$$P(A|B) = \frac{\text{number of elements of } A \cap B}{\text{number of elements of } B}.$$

1.4 TOTAL PROBABILITY THEOREM AND BAYES' RULE

Total Probability Theorem

Let A_1, \dots, A_n be disjoint events that form a partition of the sample space (each possible outcome is included in exactly one of the events A_1, \dots, A_n) and assume that $P(A_i) > 0$, for all i . Then, for any event B , we have

$$\begin{aligned} P(B) &= P(A_1 \cap B) + \dots + P(A_n \cap B) = \sum_i P(A_i) P(B|A_i) \\ &= P(A_1)P(B|A_1) + \dots + P(A_n)P(B|A_n). \end{aligned}$$

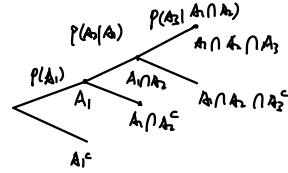
Weighted average of $P(B|A_i)$

Multiplication Rule:

Probability of leaf = product of probability of branches (given in conditional probabilities)

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

$$P(A_1 \cap \dots \cap A_n) = P(A_1) P(A_2|A_1) P(A_3|A_1 \cap A_2) \dots$$



Bayes' Rule:

initial "beliefs" of A_i : $P(A_i)$, have $P(B|A_i)$

\Rightarrow revised "beliefs". give B occurred:

$$P(A_i|B) = \frac{P(A_i) P(B|A_i)}{\sum_j P(A_j) P(B|A_j)}$$

1.5 INDEPENDENCE

Independence

- if A and B are disjoint,
 A and B are NOT
independent!

- Two events A and B are said to be **independent** if

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

If in addition, $\mathbf{P}(B) > 0$, independence is equivalent to the condition

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

- If A and B are independent, so are A and B^c .
- Two events A and B are said to be **conditionally independent**, given another event C with $\mathbf{P}(C) > 0$, if

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

If in addition, $\mathbf{P}(B \cap C) > 0$, conditional independence is equivalent to the condition

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

$$\begin{aligned} \mathbf{P}(A | B \cap C) &= \frac{\mathbf{P}(A \cap (B \cap C))}{\mathbf{P}(B \cap C)} \\ &= \frac{\mathbf{P}((A \cap B) \cap C)}{\mathbf{P}(B \cap C)} \end{aligned}$$

- Independence does not imply conditional independence, and vice versa.
- conditioning may affect independence

Definition of Independence of Several Events

We say that the events A_1, A_2, \dots, A_n are **independent** if

$$\mathbf{P}\left(\bigcap_{i \in S} A_i\right) = \prod_{i \in S} \mathbf{P}(A_i), \quad \text{for every subset } S \text{ of } \{1, 2, \dots, n\}.$$

$$\begin{aligned} &= \frac{\mathbf{P}(C) \mathbf{P}(A \cap B | C)}{\mathbf{P}(C) \mathbf{P}(B | C)} \\ &= \frac{\mathbf{P}(A | C) \mathbf{P}(B | C)}{\mathbf{P}(B | C)} \\ &= \mathbf{P}(A | C) \end{aligned}$$

- pairwise independence vs. independence
(is different!)
- Modeling: Be careful to state the assumptions
before making a probability model

Union \rightarrow intersection
 • De Morgan's law
 $(U_1 \cup U_2 \cup U_3)^c = F_1 \cap F_2 \cap F_3,$
 $F_i = U_i^c$

1.6 COUNTING

The Counting Principle

Consider a process that consists of r stages. Suppose that:

- (a) There are n_1 possible results at the first stage.
- (b) For every possible result at the first stage, there are n_2 possible results at the second stage.
- (c) More generally, for any sequence of possible results at the first $i - 1$ stages, there are n_i possible results at the i th stage.

Then, the total number of possible results of the r -stage process is

$$n_1 n_2 \cdots n_r.$$

Summary of Counting Results

- **Permutations** of n objects: $n!$.
- **k -permutations** of n objects: $n!/(n - k)!$.
- **Combinations** of k out of n objects: $\binom{n}{k} = \frac{n!}{k!(n - k)!}$. (Binomial coefficient)
- **Partitions** of n objects into r groups, with the i th group having n_i objects:

$$\binom{n}{n_1, n_2, \dots, n_r} = \frac{n!}{n_1! n_2! \cdots n_r!}$$
. (Multinomial coefficient)

- number of subsets of $\{1, \dots, n\}$: 2^n (go by each element)

$$\sum_{k=0}^n \binom{n}{k} = 2^n$$

1.7 SUMMARY AND DISCUSSION

- **Binomial probability**: $P(K \text{ heads}) = \binom{n}{k} p^k (1-p)^{n-k}$

- **Multinomial probability**: $P(\text{get type } (n_1, n_2, \dots, n_r)) = \frac{n!}{n_1! n_2! \cdots n_r!} p_1^{n_1} p_2^{n_2} \cdots p_r^{n_r}$

2

Discrete Random Variables

Excerpts from Introduction to Probability: Second Edition

by Dimitri P. Bertsekas and John N. Tsitsiklis ©

Massachusetts Institute of Technology

Contents

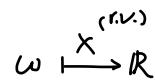
2.1. Basic Concepts	p. 72
2.2. Probability Mass Functions	p. 74
2.3. Functions of Random Variables	p. 80
2.4. Expectation, Mean, and Variance	p. 81
2.5. Joint PMFs of Multiple Random Variables	p. 92
2.6. Conditioning	p. 97
2.7. Independence	p. 109
2.8. Summary and Discussion	p. 115
Problems	p. 119

2.1 BASIC CONCEPTS

Main Concepts Related to Random Variables

Starting with a probabilistic model of an experiment:

- A **random variable** is a real-valued function of the outcome of the experiment.
- A function of a random variable defines another random variable.
- We can associate with each random variable certain “averages” of interest, such as the **mean** and the **variance**.
- A random variable can be **conditioned** on an event or on another random variable.
- There is a notion of **independence** of a random variable from an event or from another random variable.



*pay attention to the meaning of a func of r.v.
e.g. $Y = g(X)$*

• Notation:

$$p_{X=x} = P(X=x)$$

• Properties:

$$p_{X=x} \geq 0, \quad \sum_x p_{X=x} = 1$$

Concepts Related to Discrete Random Variables

Starting with a probabilistic model of an experiment:

- A **discrete random variable** is a real-valued function of the outcome of the experiment that can take a finite or countably infinite number of values.
- A discrete random variable has an associated **probability mass function (PMF)**, which gives the probability of each numerical value that the random variable can take.
- A function of a discrete random variable defines another discrete random variable, whose PMF can be obtained from the PMF of the original random variable.

- In general: $E[X]$ is the center of gravity
- PMF is symmetric: $E[X]$ will be the center of symmetry

2.2 PROBABILITY MASS FUNCTIONS

Calculation of the PMF of a Random Variable X

For each possible value x of X :

1. Collect all the possible outcomes that give rise to the event $\{X = x\}$.
2. Add their probabilities to obtain $p_X(x)$.

2.3 FUNCTIONS OF RANDOM VARIABLES

2.4 EXPECTATION, MEAN, AND VARIANCE

• basic properties:

$$\bullet X \geq 0 \Rightarrow E[X] \geq 0$$

$$\bullet a \leq X \leq b \Rightarrow a \leq E[X] \leq b$$

$$\bullet E[c] = c. \quad (c \text{ is const.})$$

Expectation

We define the **expected value** (also called the **expectation** or the **mean**) of a random variable X , with PMF p_X , by

$$E[X] = \sum_x x p_X(x).$$

Expected Value Rule for Functions of Random Variables

Let X be a random variable with PMF p_X , and let $g(X)$ be a function of X . Then, the expected value of the random variable $g(X)$ is given by

$$E[g(X)] = \sum_x g(x) p_X(x).$$

Variance

The variance $\text{var}(X)$ of a random variable X is defined by

$$\text{var}(X) = \mathbf{E} \left[(X - \mathbf{E}[X])^2 \right],$$

and can be calculated as $= \mathbf{E}[(X-\mu)^2]$ distance from the mean

$$\text{var}(X) = \sum_x (x - \mathbf{E}[X])^2 p_X(x).$$

It is always nonnegative. Its square root is denoted by σ_X and is called the **standard deviation**.

Variance → measure of the amount of randomness / uncertainty in a r.v.

Mean and Variance of a Linear Function of a Random Variable

Let X be a random variable and let

$$Y = aX + b,$$

where a and b are given scalars. Then,

$$\mathbf{E}[Y] = a\mathbf{E}[X] + b, \quad \text{var}(Y) = a^2 \text{var}(X).$$

(Linearity of Expectation)

Variance in Terms of Moments Expression

$$\text{var}(X) = \mathbf{E}[X^2] - (\mathbf{E}[X])^2.$$

2.5 JOINT PMFS OF MULTIPLE RANDOM VARIABLES

Summary of Facts About Joint PMFs

Let X and Y be random variables associated with the same experiment.

- The **joint PMF** $p_{X,Y}$ of X and Y is defined by

$$\sum_x \sum_y p_{X,Y}(x,y) = 1$$

$$p_{X,Y}(x,y) = \mathbf{P}(X = x, Y = y).$$

- The **marginal PMFs** of X and Y can be obtained from the joint PMF, using the formulas

$$p_X(x) = \sum_y p_{X,Y}(x,y), \quad p_Y(y) = \sum_x p_{X,Y}(x,y).$$

- A function $g(X, Y)$ of X and Y defines another random variable, and

$$\mathbf{E}[g(X, Y)] = \sum_x \sum_y g(x, y) p_{X,Y}(x, y).$$

If g is linear, of the form $aX + bY + c$, we have

$$(\text{Linearity of Expectation}) \quad \mathbf{E}[aX + bY + c] = a\mathbf{E}[X] + b\mathbf{E}[Y] + c.$$

always true

- The above have natural extensions to the case where more than two random variables are involved.

$$p_{X,Y,Z}(x,y,z) = \mathbf{P}(X = x \text{ and } Y = y \text{ and } Z = z)$$

2.6 CONDITIONING

Summary of Facts About Conditional PMFs

Let X and Y be random variables associated with the same experiment.

- Conditional PMFs are similar to ordinary PMFs, but pertain to a universe where the conditioning event is known to have occurred.
- The conditional PMF of X given an event A with $\mathbf{P}(A) > 0$, is defined by

$$p_{X|A}(x) = \mathbf{P}(X = x | A)$$

and satisfies

$$\sum_x p_{X|A}(x) = 1.$$

- If A_1, \dots, A_n are disjoint events that form a partition of the sample space, with $\mathbf{P}(A_i) > 0$ for all i , then

$$p_X(x) = \sum_{i=1}^n \mathbf{P}(A_i) p_{X|A_i}(x).$$

(This is a special case of the total probability theorem.) Furthermore, for any event B , with $\mathbf{P}(A_i \cap B) > 0$ for all i , we have

$$p_{X|B}(x) = \sum_{i=1}^n \mathbf{P}(A_i | B) p_{X|A_i \cap B}(x).$$

conditional PMF defined as:
 $p_{X|Y}(x|y) = \frac{p_{X,Y}(x,y)}{p_Y(y)}$,

$p_Y(y) > 0$

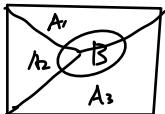
- The conditional PMF of X given $Y = y$ is related to the joint PMF by

$$p_{X,Y}(x,y) = p_Y(y) p_{X|Y}(x|y). = p_X(x) p_{Y|X}(y|x)$$

- The conditional PMF of X given Y can be used to calculate the marginal PMF of X through the formula

$$p_X(x) = \sum_y p_Y(y) p_{X|Y}(x|y). \quad y \text{ finite or countable}$$

- There are natural extensions of the above involving more than two random variables.



$$\begin{aligned}
 \bullet \quad p_{X|B}(x) &= \frac{p(X=x \text{ and } B)}{p(B)} & \bullet \quad p_{X,Y,Z}(x,y,z) &= \frac{p_{X,Y,Z}(x,y,z)}{p_{Y,Z}(y,z)} \\
 &= \frac{\sum_i p(X=x \text{ and } A_i \cap B)}{p(B)} & p_{X,Y,Z}(x,y,z) &= p_X(x) p_{Y|X}(y|x) p_{Z|Y,X}(z|y,x) \\
 &= \frac{\sum_i p(A_i \cap B) p(X=x | A_i \cap B)}{p(B)} & &
 \end{aligned}$$

Summary of Facts About Conditional Expectations

Let X and Y be random variables associated with the same experiment.

- The conditional expectation of X given an event A with $\mathbf{P}(A) > 0$, is defined by

$$\mathbf{E}[X | A] = \sum_x x p_{X|A}(x).$$

For a function $g(X)$, we have

$$\mathbf{E}[g(X) | A] = \sum_x g(x) p_{X|A}(x).$$

- The conditional expectation of X given a value y of Y is defined by

$$\mathbf{E}[g(x) | Y=y] = \sum_x g(x) p_{X|Y}(x | y) \quad \mathbf{E}[X | Y=y] = \sum_x x p_{X|Y}(x | y).$$

- If A_1, \dots, A_n be disjoint events that form a partition of the sample space, with $\mathbf{P}(A_i) > 0$ for all i , then

$$\mathbf{E}[X] = \sum_{i=1}^n \mathbf{P}(A_i) \mathbf{E}[X | A_i]. \quad (\text{total expectation theorem})$$

Furthermore, for any event B with $\mathbf{P}(A_i \cap B) > 0$ for all i , we have

$$\mathbf{E}[X | B] = \sum_{i=1}^n \mathbf{P}(A_i | B) \mathbf{E}[X | A_i \cap B].$$

- We have

$$\mathbf{E}[X] = \sum_y p_Y(y) \mathbf{E}[X | Y=y]. \quad y \text{ finite or countable}$$

2.7 INDEPENDENCE

Summary of Facts About Independent Random Variables

Let A be an event, with $\mathbf{P}(A) > 0$, and let X and Y be random variables associated with the same experiment.

- X is independent of the event A if

$$p_{X|A}(x) = p_X(x), \quad \text{for all } x,$$

that is, if for all x , the events $\{X = x\}$ and A are independent.

- X and Y are independent if for all pairs (x, y) , the events $\{X = x\}$ and $\{Y = y\}$ are independent, or equivalently

$$p_{X,Y}(x, y) = p_X(x)p_Y(y), \quad \text{for all } x, y.$$

- If X and Y are independent random variables, then

$$\mathbf{E}[XY] = \mathbf{E}[X]\mathbf{E}[Y].$$

Furthermore, for any functions g and h , the random variables $g(X)$ and $h(Y)$ are independent, and we have

$$\mathbf{E}[g(X)h(Y)] = \mathbf{E}[g(X)]\mathbf{E}[h(Y)].$$

- If X and Y are independent, then

$$\text{var}(X + Y) = \text{var}(X) + \text{var}(Y).$$

X, Y, Z are independent if:

$$p_{X,Y,Z}(x, y, z) = p_X(x)p_Y(y)p_Z(z), \text{ for all } x, y, z$$

2.8 SUMMARY AND DISCUSSION

Summary of Results for Special Random Variables

Discrete Uniform over $[a, b]$:

$a = b$. deterministic
r.v.

I_A : Indicator r.v. of event A
event \rightarrow r.v.

$$P_{I_A}(1) = P(I_A = 1) = P(A) = p$$

$$E[I_A] = P(A)$$

$$E[X] = E[X_1] + \dots + E[X_n]$$

X_i: indicator of ith trial

$$E[X_i] = p$$

$$\text{var}(x) = \text{var}(x_1) + \dots + \text{var}(x_n)$$

(independence)

$$\text{var}(x_i) = p(1-p)$$

$$p_X(k) = \begin{cases} \frac{1}{b-a+1}, & \text{if } k = a, a+1, \dots, b, \\ 0, & \text{otherwise,} \end{cases}$$

$$E[X] = \frac{a+b}{2}, \quad \text{var}(X) = \frac{(b-a)(b-a+2)}{12}.$$

Bernoulli with Parameter p : (Describes the success or failure in a single trial.)

$$p_X(k) = \begin{cases} p, & \text{if } k = 1, \\ 1-p, & \text{if } k = 0, \end{cases}$$

$$E[X] = p, \quad \text{var}(X) = p(1-p).$$

Binomial with Parameters p and n : (Describes the number of successes in n independent Bernoulli trials.)

$$p_X(k) = \binom{n}{k} p^k (1-p)^{n-k}, \quad k = 0, 1, \dots, n, \quad P(\underbrace{k \text{ heads}}_{\text{in } n \text{ trials}})$$

$$E[X] = np, \quad \text{var}(X) = np(1-p).$$

Geometric with Parameter p : (Describes the number of trials until the first success, in a sequence of independent Bernoulli trials.)

$$p_X(k) = (1-p)^{k-1} p, \quad k = 1, 2, \dots, \quad P(\underbrace{\text{Tails}}_{k-1} \text{ followed by Head})$$

$$E[X] = \frac{1}{p}, \quad \text{var}(X) = \frac{1-p}{p^2}.$$

Poisson with Parameter λ : (Approximates the binomial PMF when n is large, p is small, and $\lambda = np$.)

$$p_X(k) = e^{-\lambda} \frac{\lambda^k}{k!}, \quad k = 0, 1, \dots,$$

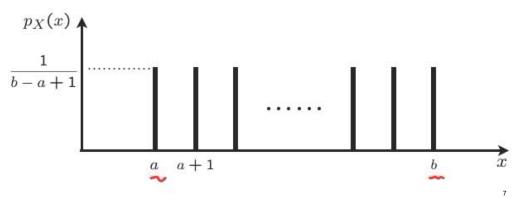
$$E[X] = \lambda, \quad \text{var}(X) = \lambda.$$

memorylessness:

$$P_{X|X>n}(k) = P_{X-n|X>n}(k)$$

Memorylessness:

Number of remaining coin tosses, conditioned on Tails in the first toss, is Geometric, with parameter p



Conditioned on $X > n$, $X - n$ is geometric with parameter p

3

General Random Variables

Excerpts from **Introduction to Probability: Second Edition**

by Dimitri P. Bertsekas and John N. Tsitsiklis ©

Massachusetts Institute of Technology

Contents

3.1. Continuous Random Variables and PDFs	p. 140
3.2. Cumulative Distribution Functions	p. 148
3.3. Normal Random Variables	p. 153
3.4. Joint PDFs of Multiple Random Variables	p. 158
3.5. Conditioning	p. 164
3.6. The Continuous Bayes' Rule	p. 178
3.7. Summary and Discussion	p. 182
Problems	p. 184

3.1 CONTINUOUS RANDOM VARIABLES AND PDFS

Summary of PDF Properties

Let X be a continuous random variable with PDF f_X .

- $f_X(x) \geq 0$ for all x .
- $\int_{-\infty}^{\infty} f_X(x) dx = 1$.
- If δ is very small, then $\mathbf{P}([x, x + \delta]) \approx f_X(x) \cdot \delta$. PDFs are densities:
probability per unit length
- For any subset B of the real line,

$$\mathbf{P}(X \in B) = \int_B f_X(x) dx. \quad \mathbf{P}(a \leq X \leq b) = \int_a^b f_X(x) dx$$

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

Expectation of a Continuous Random Variable and its Properties

Let X be a continuous random variable with PDF f_X .

- The expectation of X is defined by

$$\mathbf{E}[X] = \int_{-\infty}^{\infty} x f_X(x) dx.$$

- The expected value rule for a function $g(X)$ has the form

$$\mathbf{E}[g(X)] = \int_{-\infty}^{\infty} g(x) f_X(x) dx.$$

- The variance of X is defined by

$$\text{var}(X) = \mathbf{E}[(X - \mathbf{E}[X])^2] = \int_{-\infty}^{\infty} (x - \mathbf{E}[X])^2 f_X(x) dx.$$

- We have

$$0 \leq \text{var}(X) = \mathbf{E}[X^2] - (\mathbf{E}[X])^2.$$

- If $Y = aX + b$, where a and b are given scalars, then

$$\mathbf{E}[Y] = a\mathbf{E}[X] + b, \quad \text{var}(Y) = a^2\text{var}(X).$$

linearity

3.2 CUMULATIVE DISTRIBUTION FUNCTIONS

Properties of a CDF

- CDF of a r.v. : enough info for anything of that r.v.

The CDF F_X of a random variable X is defined by

$$F_X(x) = \mathbf{P}(X \leq x), \quad \text{for all } x,$$

and has the following properties.

- F_X is monotonically nondecreasing:

$$\text{if } x \leq y, \text{ then } F_X(x) \leq F_X(y).$$

- $F_X(x)$ tends to 0 as $x \rightarrow -\infty$, and to 1 as $x \rightarrow \infty$.
- If X is discrete, then $F_X(x)$ is a piecewise constant function of x .
- If X is continuous, then $F_X(x)$ is a continuous function of x .
- If X is discrete and takes integer values, the PMF and the CDF can be obtained from each other by summing or differencing:

$$F_X(k) = \sum_{i=-\infty}^k p_X(i),$$

$$p_X(k) = \mathbf{P}(X \leq k) - \mathbf{P}(X \leq k-1) = F_X(k) - F_X(k-1),$$

for all integers k .

- If X is continuous, the PDF and the CDF can be obtained from each other by integration or differentiation:

$$F_X(x) = \int_{-\infty}^x f_X(t) dt, \quad f_X(x) = \frac{dF_X}{dx}(x).$$

(The second equality is valid for those x at which the PDF is continuous.)

3.3 NORMAL RANDOM VARIABLES

Normality is Preserved by Linear Transformations

If X is a normal random variable with mean μ and variance σ^2 , and if $a \neq 0$, b are scalars, then the random variable

$$Y = aX + b$$

is also normal, with mean and variance

$$\mathbf{E}[Y] = a\mu + b, \quad \text{var}(Y) = a^2\sigma^2.$$

CDF Calculation for a Normal Random Variable

For a normal random variable X with mean μ and variance σ^2 , we use a two-step procedure.

- (a) “Standardize” X , i.e., subtract μ and divide by σ to obtain a standard normal random variable Y .
- (b) Read the CDF value from the standard normal table:

$$\mathbf{P}(X \leq x) = \mathbf{P}\left(\frac{X - \mu}{\sigma} \leq \frac{x - \mu}{\sigma}\right) = \mathbf{P}\left(Y \leq \frac{x - \mu}{\sigma}\right) = \Phi\left(\frac{x - \mu}{\sigma}\right).$$

$$X \sim N(\mu, \sigma^2)$$

$$Y = \frac{X - \mu}{\sigma}$$

$$Y \sim N(0, 1)$$

$$X = \mu + \sigma Y$$

	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
0.0	.5000	.5040	.5080	.5120	.5160	.5199	.5239	.5279	.5319	.5359
0.1	.5398	.5438	.5478	.5517	.5557	.5596	.5636	.5675	.5714	.5753
0.2	.5793	.5832	.5871	.5910	.5948	.5987	.6026	.6064	.6103	.6141
0.3	.6179	.6217	.6255	.6293	.6331	.6368	.6406	.6443	.6480	.6517
0.4	.6554	.6591	.6628	.6664	.6700	.6736	.6772	.6808	.6844	.6879
0.5	.6915	.6950	.6985	.7019	.7054	.7088	.7123	.7157	.7190	.7224
0.6	.7257	.7291	.7324	.7357	.7389	.7422	.7454	.7486	.7517	.7549
0.7	.7580	.7611	.7642	.7673	.7704	.7734	.7764	.7794	.7823	.7852
0.8	.7881	.7910	.7939	.7967	.7995	.8023	.8051	.8078	.8106	.8133
0.9	.8159	.8186	.8212	.8238	.8264	.8289	.8315	.8340	.8365	.8389
1.0	.8413	.8438	.8461	.8485	.8508	.8531	.8554	.8577	.8599	.8621
1.1	.8643	.8665	.8686	.8708	.8729	.8749	.8770	.8790	.8810	.8830
1.2	.8849	.8869	.8888	.8907	.8925	.8944	.8962	.8980	.8997	.9015
1.3	.9032	.9049	.9066	.9082	.9099	.9115	.9131	.9147	.9162	.9177
1.4	.9192	.9207	.9222	.9236	.9251	.9265	.9279	.9292	.9306	.9319
1.5	.9332	.9345	.9357	.9370	.9382	.9394	.9406	.9418	.9429	.9441
1.6	.9452	.9463	.9474	.9484	.9495	.9505	.9515	.9525	.9535	.9545
1.7	.9554	.9564	.9573	.9582	.9591	.9599	.9608	.9616	.9625	.9633
1.8	.9641	.9649	.9656	.9664	.9671	.9678	.9686	.9693	.9699	.9706
1.9	.9713	.9719	.9726	.9732	.9738	.9744	.9750	.9756	.9761	.9767
2.0	.9772	.9778	.9783	.9788	.9793	.9798	.9803	.9808	.9812	.9817
2.1	.9821	.9826	.9830	.9834	.9838	.9842	.9846	.9850	.9854	.9857
2.2	.9861	.9864	.9868	.9871	.9875	.9878	.9881	.9884	.9887	.9890
2.3	.9893	.9896	.9898	.9901	.9904	.9906	.9909	.9911	.9913	.9916
2.4	.9918	.9920	.9922	.9925	.9927	.9929	.9931	.9932	.9934	.9936
2.5	.9938	.9940	.9941	.9943	.9945	.9946	.9948	.9949	.9951	.9952
2.6	.9953	.9955	.9956	.9957	.9959	.9960	.9961	.9962	.9963	.9964
2.7	.9965	.9966	.9967	.9968	.9969	.9970	.9971	.9972	.9973	.9974
2.8	.9974	.9975	.9976	.9977	.9977	.9978	.9979	.9979	.9980	.9981
2.9	.9981	.9982	.9982	.9983	.9984	.9984	.9985	.9985	.9986	.9986
3.0	.9987	.9987	.9987	.9988	.9988	.9989	.9989	.9989	.9990	.9990
3.1	.9990	.9991	.9991	.9991	.9992	.9992	.9992	.9992	.9993	.9993
3.2	.9993	.9993	.9994	.9994	.9994	.9994	.9994	.9995	.9995	.9995
3.3	.9995	.9995	.9995	.9996	.9996	.9996	.9996	.9996	.9996	.9997
3.4	.9997	.9997	.9997	.9997	.9997	.9997	.9997	.9997	.9997	.9998

The standard normal table. The entries in this table provide the numerical values of $\Phi(y) = \mathbf{P}(Y \leq y)$, where Y is a standard normal random variable, for y between 0 and 3.49. For example, to find $\Phi(1.71)$, we look at the row corresponding to 1.7 and the column corresponding to 0.01, so that $\Phi(1.71) = .9564$. When y is negative, the value of $\Phi(y)$ can be found using the formula $\Phi(y) = 1 - \Phi(-y)$.

3.4 JOINT PDFS OF MULTIPLE RANDOM VARIABLES

Summary of Facts about Joint PDFs

Let X and Y be jointly continuous random variables with joint PDF $f_{X,Y}$.

- The **joint PDF** is used to calculate probabilities:

$$\mathbf{P}((X, Y) \in B) = \int_{(x,y) \in B} \int f_{X,Y}(x, y) dx dy.$$

- The **marginal PDFs** of X and Y can be obtained from the joint PDF, using the formulas

$$f_X(x) = \int_{-\infty}^{\infty} f_{X,Y}(x, y) dy, \quad f_Y(y) = \int_{-\infty}^{\infty} f_{X,Y}(x, y) dx.$$

- The **joint CDF** is defined by $F_{X,Y}(x, y) = \mathbf{P}(X \leq x, Y \leq y)$, and determines the joint PDF through the formula

$$f_{X,Y}(x, y) = \frac{\partial^2 F_{X,Y}}{\partial x \partial y}(x, y),$$

for every (x, y) at which the joint PDF is continuous.

- A function $g(X, Y)$ of X and Y defines a new random variable, and

$$\mathbf{E}[g(X, Y)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y) f_{X,Y}(x, y) dx dy.$$

If g is linear, of the form $aX + bY + c$, we have

$$\mathbf{E}[aX + bY + c] = a\mathbf{E}[X] + b\mathbf{E}[Y] + c.$$

- The above have natural extensions to the case where more than two random variables are involved.

Definition: Two random variables are **jointly continuous** if they can be described by a joint PDF

joint PDF: $\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{X,Y}(x, y) dx dy = 1$

$f_{X,Y}(x, y) \geq 0$

3.5 CONDITIONING

Conditional PDF Given an Event

- The conditional PDF $f_{X|A}$ of a continuous random variable X , given an event A with $\mathbf{P}(A) > 0$, satisfies

$$\mathbf{P}(X \in B | A) = \int_B f_{X|A}(x) dx.$$

- If A is a subset of the real line with $\mathbf{P}(X \in A) > 0$, then

$$f_{X|\{X \in A\}}(x) = \begin{cases} \frac{f_X(x)}{\mathbf{P}(X \in A)}, & \text{if } x \in A, \\ 0, & \text{otherwise.} \end{cases}$$

- Let A_1, A_2, \dots, A_n be disjoint events that form a partition of the sample space, and assume that $\mathbf{P}(A_i) > 0$ for all i . Then,

$$f_X(x) = \sum_{i=1}^n \mathbf{P}(A_i) f_{X|A_i}(x)$$

(a version of the total probability theorem).

Conditional PDF Given a Random Variable

Let X and Y be jointly continuous random variables with joint PDF $f_{X,Y}$.

- The joint, marginal, and conditional PDFs are related to each other by the formulas

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

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

$$f_X(x) = \int_{-\infty}^{\infty} f_Y(y) f_{X|Y}(x|y) dy. \quad \text{total probability theorem}$$

The conditional PDF $f_{X|Y}(x|y)$ is defined only for those y for which $f_Y(y) > 0$.

- We have

$$\mathbf{P}(X \in A | Y = y) = \int_A f_{X|Y}(x|y) dx.$$

Summary of Facts About Conditional Expectations

Let X and Y be jointly continuous random variables, and let A be an event with $\mathbf{P}(A) > 0$.

- **Definitions:** The conditional expectation of X given the event A is defined by

$$\mathbf{E}[X | A] = \int_{-\infty}^{\infty} x f_{X|A}(x) dx.$$

The conditional expectation of X given that $Y = y$ is defined by

$$\mathbf{E}[X | Y = y] = \int_{-\infty}^{\infty} x f_{X|Y}(x | y) dx.$$

- **The expected value rule:** For a function $g(X)$, we have

$$\mathbf{E}[g(X) | A] = \int_{-\infty}^{\infty} g(x) f_{X|A}(x) dx,$$

and

$$\mathbf{E}[g(X) | Y = y] = \int_{-\infty}^{\infty} g(x) f_{X|Y}(x | y) dx.$$

- **Total expectation theorem:** Let A_1, A_2, \dots, A_n be disjoint events that form a partition of the sample space, and assume that $\mathbf{P}(A_i) > 0$ for all i . Then,

$$\mathbf{E}[X] = \sum_{i=1}^n \mathbf{P}(A_i) \mathbf{E}[X | A_i].$$

Similarly,

$$\mathbf{E}[X] = \int_{-\infty}^{\infty} \mathbf{E}[X | Y = y] f_Y(y) dy.$$

- There are natural analogs for the case of functions of several random variables. For example,

$$\mathbf{E}[g(X, Y) | Y = y] = \int g(x, y) f_{X|Y}(x | y) dx,$$

and

$$\mathbf{E}[g(X, Y)] = \int \mathbf{E}[g(X, Y) | Y = y] f_Y(y) dy.$$

Mixed random variable:

$$X = \begin{cases} \text{discrete w.p. } p \\ \text{continuous w.p. } 1-p \end{cases}$$

describe by CDF

$$F_X(x) = p F_Y(x) + (1-p) F_Z(x)$$

Independence of Continuous Random Variables

Let X and Y be jointly continuous random variables.

- X and Y are **independent** if

$$f_{X,Y}(x,y) = f_X(x)f_Y(y), \quad \text{for all } x,y.$$

- If X and Y are independent, then

$$\mathbf{E}[XY] = \mathbf{E}[X]\mathbf{E}[Y].$$

Furthermore, for any functions g and h , the random variables $g(X)$ and $h(Y)$ are independent, and we have

$$\mathbf{E}[g(X)h(Y)] = \mathbf{E}[g(X)]\mathbf{E}[h(Y)].$$

- If X and Y are independent, then

$$\text{var}(X + Y) = \text{var}(X) + \text{var}(Y).$$

$$f_{X|Y}(x|y) = f_X(x)$$

$$f_{Y|X}(y|x) = f_Y(y)$$

3.6 BAYES' RULE AND APPLICATIONS IN INFERENCE

Bayes' Rule Relations for Random Variables

Let X and Y be two random variables.

- If X and Y are discrete, we have for all x, y with $p_X(x) \neq 0, p_Y(y) \neq 0$,

$$p_X(x)p_{Y|X}(y|x) = p_Y(y)p_{X|Y}(x|y),$$

and the terms on the two sides in this relation are both equal to

$$p_{X,Y}(x,y).$$

- If X is discrete and Y is continuous, we have for all x, y with $p_X(x) \neq 0, f_Y(y) \neq 0$,

$$p_X(x)f_{Y|X}(y|x) = f_Y(y)p_{X|Y}(x|y),$$

and the terms on the two sides in this relation are both equal to

$$\lim_{\delta \rightarrow 0} \frac{\mathbf{P}(X = x, y \leq Y \leq y + \delta)}{\delta}.$$

- If X and Y are continuous, we have for all x, y with $f_X(x) \neq 0, f_Y(y) \neq 0$,

$$f_X(x)f_{Y|X}(y|x) = f_Y(y)f_{X|Y}(x|y),$$

and the terms on the two sides in this relation are both equal to

$$\lim_{\delta \rightarrow 0} \frac{\mathbf{P}(x \leq X \leq x + \delta, y \leq Y \leq y + \delta)}{\delta^2}.$$

3.7 SUMMARY AND DISCUSSION

Summary of Results for Special Random Variables

Continuous Uniform Over $[a, b]$:

$$f_X(x) = \begin{cases} \frac{1}{b-a}, & \text{if } a \leq x \leq b, \\ 0, & \text{otherwise,} \end{cases}$$

$$\mathbb{E}[X] = \frac{a+b}{2}, \quad \text{var}(X) = \frac{(b-a)^2}{12}.$$

Exponential with Parameter λ : time we have to wait until

$$f_X(x) = \begin{cases} \lambda e^{-\lambda x}, & \text{if } x \geq 0, \\ 0, & \text{otherwise,} \end{cases} \quad F_X(x) = \begin{cases} 1 - e^{-\lambda x}, & \text{if } x \geq 0, \\ 0, & \text{otherwise,} \end{cases}$$

$$\mathbb{E}[X] = \frac{1}{\lambda}, \quad \text{var}(X) = \frac{1}{\lambda^2}.$$

Normal with Parameters μ and $\sigma^2 > 0$: $N(\mu, \sigma^2)$

$$f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2},$$

$$\mathbb{E}[X] = \mu, \quad \text{var}(X) = \sigma^2.$$

Surprising happens.
Compare to geometric

Standard normal

 $N(0, 1)$

$$P(X \geq a) = e^{-\lambda a} \quad (a \geq 0)$$

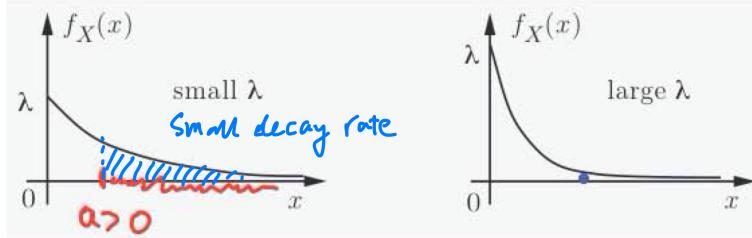
Memorylessness:

$$P(T > x) = e^{-\lambda x} \quad x \geq 0.$$

$$P(X > x | T > t) = e^{-\lambda x}, \quad x \geq 0$$

$$X = T - t$$

T: light bulb lifetime

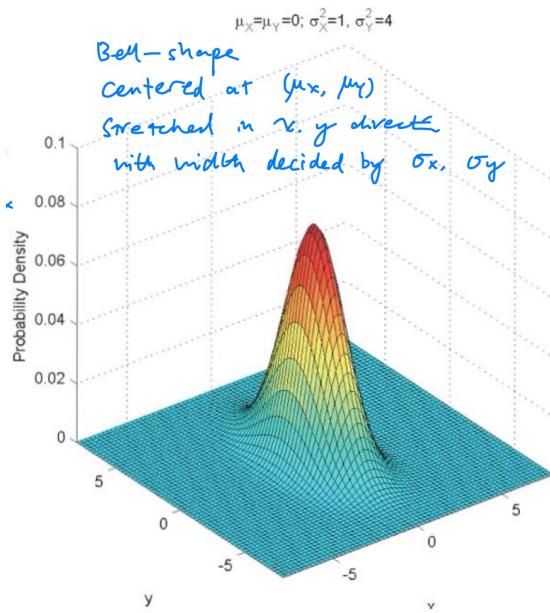
X: remaining lifetime, given that
light bulb has been on for
time t

Independent standard normal:

$$f_{X,Y}(x, y) = \frac{1}{2\pi} \exp\left\{-\frac{1}{2}(x^2 + y^2)\right\}$$

Independent normal:

$$f_{X,Y}(x, y) = \frac{1}{2\pi\sigma_x\sigma_y} \exp\left\{-\frac{(x-\mu_x)^2}{2\sigma_x^2} - \frac{(y-\mu_y)^2}{2\sigma_y^2}\right\}$$



4

Further Topics on Random Variables

Excerpts from Introduction to Probability: Second Edition

by Dimitri P. Bertsekas and John N. Tsitsiklis ©

Massachusetts Institute of Technology

Contents

4.1. Derived Distributions	p. 202
4.2. Covariance and Correlation	p. 217
4.3. Conditional Expectation and Variance Revisited	p. 222
4.4. Transforms	p. 229
4.5. Sum of a Random Number of Independent Random Variables	p. 240
4.6. Summary and Discussion	p. 244
Problems	p. 246

In general, X is r.v. with PDF $f_X(x)$.

PDF of $Y = g(X)$ is

$$f_Y(y) = \begin{cases} f_X(g^{-1}(y))|g'(g^{-1}(y))|, & y \geq 0 \\ 0, & y < 0 \end{cases}$$

can prove by using the ³⁰CDF of y . From Introduction to Probability, by Bertsekas and Tsitsiklis Chap. 4

4.1 DERIVED DISTRIBUTIONS

multiple r.v. $Z = g(X, Y)$

same method.

① CDF of Z

② $\frac{d}{dz} F_Z(z)$

Calculation of the PDF of a Function $Y = g(X)$ of a Continuous Random Variable X

- Calculate the CDF F_Y of Y using the formula

$$F_Y(y) = P(g(X) \leq y) = \int_{\{x \mid g(x) \leq y\}} f_X(x) dx.$$

- Differentiate to obtain the PDF of Y :

$$f_Y(y) = \frac{dF_Y}{dy}(y).$$

The PDF of a Linear Function of a Random Variable

Let X be a continuous random variable with PDF f_X , and let

$$Y = aX + b,$$

where a and b are scalars, with $a \neq 0$. Then,

$$f_Y(y) = \frac{1}{|a|} f_X\left(\frac{y-b}{a}\right).$$

- If $X \sim \mathcal{N}(\mu, \sigma^2)$, then $aX + b \sim \mathcal{N}(a\mu + b, a^2\sigma^2)$



- discrete case:

$$P_Y(y) = P(g(X) = y) = \sum_{x: g(x)=y} p_X(x)$$

$$\text{linear: } P_Y(y) = P(aX+b=y) = P(X = \frac{y-b}{a}) = p_X(\frac{y-b}{a})$$

- convolution formula:

$$Z = X + Y, \quad X, Y \text{ independent}$$

discrete:

$$p_Z(z) = \sum_x p_X(x) p_Y(z-x)$$

continuous:

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(x) f_Y(z-x) dx$$

- $X \sim N(\mu_X, \sigma_X^2), \quad Y \sim N(\mu_Y, \sigma_Y^2) \Rightarrow Z = X + Y \sim N(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2)$



The sum of finitely many independent normals is normal

PDF Formula for a Strictly Monotonic Function of a Continuous Random Variable

Suppose that g is strictly monotonic and that for some function h and all x in the range of X we have

$$y = g(x) \quad \text{if and only if} \quad x = h(y).$$

Assume that h is differentiable. Then, the PDF of Y in the region where $f_Y(y) > 0$ is given by

$$f_Y(y) = f_X(h(y)) \left| \frac{dh}{dy}(y) \right|.$$

- Correlation often reflects underlying, common, hidden factor

4.2 COVARIANCE AND CORRELATION

Covariance and Correlation

- covariance describes if X and Y tend to move in the same direction ($X > E[X]$, $Y > E[Y]$)
- independent $\Rightarrow \text{cov}(X, Y) = 0$
 $\rho = 0$
(converse not true)
- $\rho(X, X) = 1$
- $\rho(X, -X) = -1$
- $\rho(aX+b, Y)$
= $\text{Sign}(a) \cdot \rho(X, Y)$

- The covariance of X and Y is given by

$$\text{cov}(X, Y) = \mathbf{E}[(X - E[X])(Y - E[Y])] = E[XY] - E[X]E[Y].$$

(also $\rho = 0$)

- If $\text{cov}(X, Y) = 0$, we say that X and Y are **uncorrelated**.
- If X and Y are independent, they are uncorrelated. The converse is not always true.
- We have

$$\text{var}(X + Y) = \text{var}(X) + \text{var}(Y) + 2\text{cov}(X, Y).$$

- The **correlation coefficient** $\rho(X, Y)$ of two random variables X and Y with positive variances is defined by

$$\rho(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X)\text{var}(Y)}}, \quad = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$$

and satisfies

$$-1 \leq \rho(X, Y) \leq 1.$$

- $|\rho| = 1 \Leftrightarrow$ linearly related : $X - E[X] = C(Y - E[Y])$

- $\text{cov}(X, X) = \text{var}(X)$

$$\text{cov}(aX+b, Y) = a \cdot \text{cov}(X, Y)$$

$$\text{cov}(X, Y+Z) = \text{cov}(X, Y) + \text{cov}(X, Z)$$



$$\text{var}(X_1 + \dots + X_n) = \sum_{i=1}^n \text{var}(X_i) + \sum_{\{(i,j): i \neq j\}} \text{cov}(X_i, X_j)$$

4.3 CONDITIONAL EXPECTATION AND VARIANCE REVISITED

Law of Iterated Expectations: $\mathbf{E}[\mathbf{E}[X | Y]] = \mathbf{E}[X].$

Law of Total Variance: $\text{var}(X) = \mathbf{E}[\text{var}(X | Y)] + \text{var}(\mathbf{E}[X | Y]).$

Properties of the Conditional Expectation and Variance

- $\mathbf{E}[X | Y = y]$ is a number whose value depends on y .
- $\mathbf{E}[X | Y]$ is a function of the random variable Y , hence a random variable. Its value is $\mathbf{E}[X | Y = y]$ whenever the value of Y is y .
- $\mathbf{E}[\mathbf{E}[X | Y]] = \mathbf{E}[X]$ (law of iterated expectations). \Leftrightarrow total expectation theorem
- $\mathbf{E}[X | Y = y]$ may be viewed as an estimate of X given $Y = y$. The corresponding error $\mathbf{E}[X | Y] - X$ is a zero mean random variable that is uncorrelated with $\mathbf{E}[X | Y]$. (?)
- $\text{var}(X | Y)$ is a random variable whose value is $\text{var}(X | Y = y)$ whenever the value of Y is y .
- $\text{var}(X) = \mathbf{E}[\text{var}(X | Y)] + \text{var}(\mathbf{E}[X | Y])$ (law of total variance).

$$\text{var}(X) = (\text{average variability within sections}) + (\text{variability between sections})$$

See Lecture 13 example

- sum of a random number of independent r.v.s

4.4 TRANSFORMS

Summary of Transforms and their Properties

- The transform associated with a random variable X is given by

$$M_X(s) = \mathbf{E}[e^{sX}] = \begin{cases} \sum_x e^{sx} p_X(x), & X \text{ discrete}, \\ \int_{-\infty}^{\infty} e^{sx} f_X(x) dx, & X \text{ continuous}. \end{cases}$$

- The distribution of a random variable is completely determined by the corresponding transform.
- Moment generating properties:

$$M_X(0) = 1, \quad \left. \frac{d}{ds} M_X(s) \right|_{s=0} = \mathbf{E}[X], \quad \left. \frac{d^n}{ds^n} M_X(s) \right|_{s=0} = \mathbf{E}[X^n].$$

- If $Y = aX + b$, then $M_Y(s) = e^{sb} M_X(as)$.
- If X and Y are independent, then $M_{X+Y}(s) = M_X(s)M_Y(s)$.

Transforms for Common Discrete Random Variables

Bernoulli(p) ($k = 0, 1$)

$$p_X(k) = \begin{cases} p, & \text{if } k = 1, \\ 1 - p, & \text{if } k = 0, \end{cases} \quad M_X(s) = 1 - p + pe^s.$$

Binomial(n, p) ($k = 0, 1, \dots, n$)

$$p_X(k) = \binom{n}{k} p^k (1-p)^{n-k}, \quad M_X(s) = (1 - p + pe^s)^n.$$

Geometric(p) ($k = 1, 2, \dots$)

$$p_X(k) = p(1 - p)^{k-1}, \quad M_X(s) = \frac{pe^s}{1 - (1 - p)e^s}.$$

Poisson(λ) ($k = 0, 1, \dots$)

$$p_X(k) = \frac{e^{-\lambda} \lambda^k}{k!}, \quad M_X(s) = e^{\lambda(e^s - 1)}.$$

Uniform(a, b) ($k = a, a+1, \dots, b$)

$$p_X(k) = \frac{1}{b-a+1}, \quad M_X(s) = \frac{e^{sa}(e^{s(b-a+1)} - 1)}{(b-a+1)(e^s - 1)}.$$

Transforms for Common Continuous Random Variables

Uniform(a, b) ($a \leq x \leq b$)

$$f_X(x) = \frac{1}{b-a}, \quad M_X(s) = \frac{e^{sb} - e^{sa}}{s(b-a)}.$$

Exponential(λ) ($x \geq 0$)

$$f_X(x) = \lambda e^{-\lambda x}, \quad M_X(s) = \frac{\lambda}{\lambda - s}, \quad (s < \lambda).$$

Normal(μ, σ^2) ($-\infty < x < \infty$)

$$f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2}, \quad M_X(s) = e^{(\sigma^2 s^2/2) + \mu s}.$$

4.5 SUM OF A RANDOM NUMBER OF INDEPENDENT RANDOM VARIABLES

Properties of the Sum of a Random Number of Independent Random Variables

Let X_1, X_2, \dots be identically distributed random variables with mean $\mathbf{E}[X]$ and variance $\text{var}(X)$. Let N be a random variable that takes nonnegative integer values. We assume that all of these random variables are independent, and we consider the sum

$$Y = X_1 + \dots + X_N.$$

Then:

- $\mathbf{E}[Y] = \mathbf{E}[N] \mathbf{E}[X]$.
- $\text{var}(Y) = \mathbf{E}[N] \text{var}(X) + (\mathbf{E}[X])^2 \text{var}(N)$.
- We have

$$M_Y(s) = M_N(\log M_X(s)).$$

Equivalently, the transform $M_Y(s)$ is found by starting with the transform $M_N(s)$ and replacing each occurrence of e^s with $M_X(s)$.

4.6 SUMMARY AND DISCUSSION

5

Limit Theorems

Excerpts from Introduction to Probability: Second Edition

by Dimitri P. Bertsekas and John N. Tsitsiklis ©

Massachusetts Institute of Technology

Contents

5.1. Markov and Chebyshev Inequalities	p. 265
5.2. The Weak Law of Large Numbers	p. 269
5.3. Convergence in Probability	p. 271
5.4. The Central Limit Theorem	p. 273
5.5. The Strong Law of Large Numbers	p. 280
5.6. Summary and Discussion	p. 282
Problems	p. 284

5.1 MARKOV AND CHEBYSHEV INEQUALITIES

Markov Inequality

If a random variable X can only take nonnegative values, then

$$\mathbf{P}(X \geq a) \leq \frac{\mathbf{E}[X]}{a}, \quad \text{for all } a > 0.$$

Chebyshev Inequality

If X is a random variable with mean μ and variance σ^2 , then

$$\mathbf{P}(|X - \mu| \geq c) \leq \frac{\sigma^2}{c^2}, \quad \text{for all } c > 0.$$

5.2 THE WEAK LAW OF LARGE NUMBERS

• sample mean

$$M_n = \frac{x_1 + \dots + x_n}{n} \quad \text{rv}$$

true mean

$$\mu = E[X_i] \quad \text{number}$$

The Weak Law of Large Numbers

Let X_1, X_2, \dots be independent identically distributed random variables with mean μ . For every $\epsilon > 0$, we have

$$\mathbf{P}(|M_n - \mu| \geq \epsilon) = \mathbf{P}\left(\left|\frac{X_1 + \dots + X_n}{n} - \mu\right| \geq \epsilon\right) \rightarrow 0, \quad \text{as } n \rightarrow \infty.$$

• empirical frequency of event A (with $p = \mathbf{P}(A)$)

is close to p (the probability of event A)

• $\mathbf{P}(|M_n - \mu| \geq \epsilon) \leq \frac{\sigma^2}{n\epsilon^2}$ (fixed $\epsilon > 0$)

• M_n converges in probability to μ

- properties : $x_n \rightarrow a$, $y_n \rightarrow b$,

$$g(x_n) \rightarrow g(a).$$

$$x_n + y_n \rightarrow a+b$$

but $E[X_n]$ need not converge to a

5.3 CONVERGENCE IN PROBABILITY

Convergence of a Deterministic Sequence

Let a_1, a_2, \dots be a sequence of real numbers, and let a be another real number. We say that the sequence a_n converges to a , or $\lim_{n \rightarrow \infty} a_n = a$, if for every $\epsilon > 0$ there exists some n_0 such that

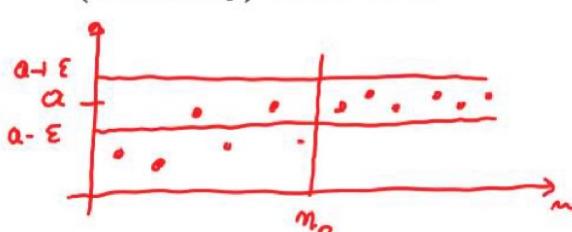
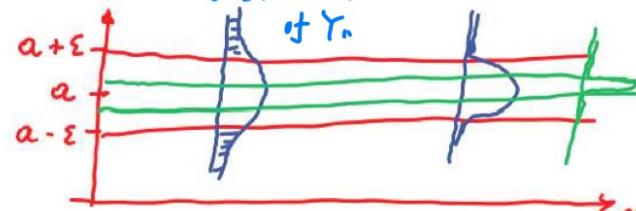
$$|a_n - a| \leq \epsilon, \quad \text{for all } n \geq n_0.$$

Convergence in Probability

Let Y_1, Y_2, \dots be a sequence of random variables (not necessarily independent), and let a be a real number. We say that the sequence Y_n **converges to a in probability**, if for every $\epsilon > 0$, we have

$$\lim_{n \rightarrow \infty} \mathbf{P}(|Y_n - a| \geq \epsilon) = 0.$$

Understanding convergence "in probability"

- Ordinary convergence
 - Sequence a_n ; number a
 $a_n \rightarrow a$
 " a_n eventually gets and stays (arbitrarily) close to a "
- Convergence in probability
 - Sequence Y_n ; number a
 $Y_n \rightarrow a$
 • for any $\epsilon > 0$, $\mathbf{P}(|Y_n - a| \geq \epsilon) \rightarrow 0$
distribution of Y_n

"(almost all) of the PMF/PDF of Y_n eventually gets concentrated (arbitrarily) close to a ".

5.4 THE CENTRAL LIMIT THEOREM

The Central Limit Theorem

Let X_1, X_2, \dots be a sequence of independent identically distributed random variables with common mean μ and variance σ^2 , and define

$$Z_n = \frac{X_1 + \dots + X_n - n\mu}{\sigma\sqrt{n}} = \frac{S_n - n\mu}{\sigma\sqrt{n}} \sim N(0, 1)$$

Then, the CDF of Z_n converges to the standard normal CDF

$$\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-x^2/2} dx, \quad = P(Z \leq z)$$

in the sense that

$$\lim_{n \rightarrow \infty} P(Z_n \leq z) = \Phi(z), \quad \text{for every } z.$$

- $E[S_n] = n\mu$
 - $\text{Var}(S_n) = n\sigma^2$
 - treat Z_n as if normal
 - for sample mean M_n :
- $$\frac{S_n - n\mu}{\sigma\sqrt{n}} = \frac{M_n - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$$

• $P(S_n \leq a) \approx b$

give 2 parameters, find the other one

Normal Approximation Based on the Central Limit Theorem

Let $S_n = X_1 + \dots + X_n$, where the X_i are independent identically distributed random variables with mean μ and variance σ^2 . If n is large, the probability $P(S_n \leq c)$ can be approximated by treating S_n as if it were normal, according to the following procedure.

- Calculate the mean $n\mu$ and the variance $n\sigma^2$ of S_n .
- Calculate the normalized value $z = (c - n\mu)/\sigma\sqrt{n}$.
- Use the approximation

$$P(S_n \leq c) \approx \Phi(z),$$

where $\Phi(z)$ is available from standard normal CDF tables.

De Moivre-Laplace Approximation to the Binomial (1/2 correction)

If S_n is a binomial random variable with parameters n and p , n is large, and k, l are nonnegative integers, then

$$\mathbf{P}(k \leq S_n \leq l) \approx \Phi\left(\frac{l + \frac{1}{2} - np}{\sqrt{np(1-p)}}\right) - \Phi\left(\frac{k - \frac{1}{2} - np}{\sqrt{np(1-p)}}\right).$$

$l = k \Rightarrow$ approximate binomial PMF

5.5 THE STRONG LAW OF LARGE NUMBERS

The Strong Law of Large Numbers

Let X_1, X_2, \dots be a sequence of independent identically distributed random variables with mean μ . Then, the sequence of sample means $M_n = (X_1 + \dots + X_n)/n$ converges to μ , **with probability 1**, in the sense that

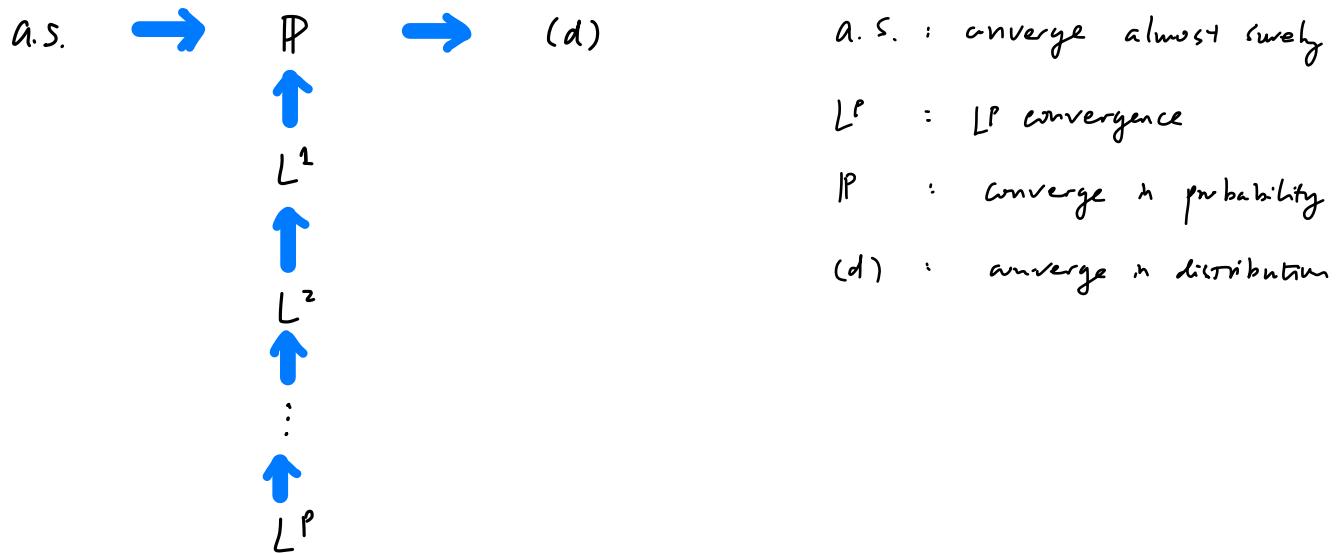
$$\mathbf{P}\left(\lim_{n \rightarrow \infty} \frac{X_1 + \dots + X_n}{n} = \mu\right) = 1.$$

Convergence with Probability 1

Let Y_1, Y_2, \dots be a sequence of random variables (not necessarily independent). Let c be a real number. We say that Y_n converges to c **with probability 1** (or **almost surely**) if

$$\mathbf{P}\left(\lim_{n \rightarrow \infty} Y_n = c\right) = 1.$$

5.6 SUMMARY AND DISCUSSION



Remember this diagram!

6

The Bernoulli and Poisson Processes

Excerpts from Introduction to Probability: Second Edition

by Dimitri P. Bertsekas and John N. Tsitsiklis ©

Massachusetts Institute of Technology

Contents

6.1. The Bernoulli Process	p. 297
6.2. The Poisson Process	p. 309
6.3. Summary and Discussion	p. 324
Problems	p. 326

• Stochastic processes:

{ infinite sequence of r.v. X_1, X_2, \dots
 Sample space Ω : set of infinite sequences of 0s and 1s (for Bernoulli process)

6.1 THE BERNOULLI PROCESS

Some Random Variables Associated with the Bernoulli Process and their Properties

$$S = X_1 + \dots + X_n$$

- The binomial with parameters p and n . This is the number S of successes in n independent trials. Its PMF, mean, and variance are

arrivals

$$P(S=k) = p_S(k) = \binom{n}{k} p^k (1-p)^{n-k}, \quad k = 0, 1, \dots, n,$$

$$E[S] = np, \quad \text{var}(S) = np(1-p).$$

- The geometric with parameter p . This is the number T of trials up to (and including) the first success. Its PMF, mean, and variance are

arrival

$$P(T=t) = p_T(t) = (1-p)^{t-1} p, \quad t = 1, 2, \dots,$$

$$E[T] = \frac{1}{p}, \quad \text{var}(T) = \frac{1-p}{p^2}.$$

Independence Properties of the Bernoulli Process

- For any given time n , the sequence of random variables X_{n+1}, X_{n+2}, \dots (the future of the process) is also a Bernoulli process, and is independent from X_1, \dots, X_n (the past of the process).
- Let n be a given time and let \bar{T} be the time of the first success after time n . Then, $\bar{T} - n$ has a geometric distribution with parameter p , and is independent of the random variables X_1, \dots, X_n .

as long as n is determined
"causally"

Alternative Description of the Bernoulli Process

1. Start with a sequence of independent geometric random variables T_1, T_2, \dots , with common parameter p , and let these stand for the interarrival times.
2. Record a success (or arrival) at times $T_1, T_1 + T_2, T_1 + T_2 + T_3$, etc.

- Bernoulli process key assumptions:

{ independence
 time-homogeneity: p is the same for all trials

- r.v.: time length of the first busy period
 $\text{Geo}(1-p)$

Properties of the k th Arrival Time

- The k th arrival time is equal to the sum of the first k interarrival times

$$Y_k = T_1 + T_2 + \cdots + T_k,$$

T_i are i.i.d. Geometric(p)

and the latter are independent geometric random variables with common parameter p .

- The mean and variance of Y_k are given by

$$\mathbf{E}[Y_k] = \mathbf{E}[T_1] + \cdots + \mathbf{E}[T_k] = \frac{k}{p},$$

$$\text{var}(Y_k) = \text{var}(T_1) + \cdots + \text{var}(T_k) = \frac{k(1-p)}{p^2}.$$

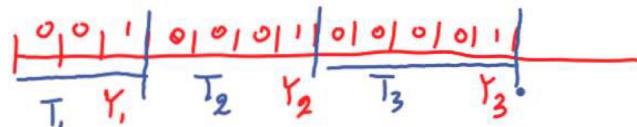
- $P_{Y_k}(t) = P(Y_k = t)$

$$= P(k-1 \text{ arrivals in time } t-1) \cdot P(\text{arrival at time } t)$$

- The PMF of Y_k is given by

$$p_{Y_k}(t) = \binom{t-1}{k-1} p^k (1-p)^{t-k}, \quad t = k, k+1, \dots,$$

and is known as the **Pascal PMF of order k** .



- Y_k = time of k th arrival $\quad Y_k = T_1 + \cdots + T_k$
- T_k = k th inter-arrival time $= Y_k - Y_{k-1}$ ($k \geq 2$)
Ti: time until the first arrival. Geo(p)
- The process starts fresh after time T_1

Merging of independent Bernoulli processes

$$\begin{array}{c} \text{Bernoulli}(p) \\ \downarrow \\ \text{Independence} \\ \left\{ \begin{array}{l} P(\text{success}) = 1 - (1-p)(1-p) \\ = p + p - p^2 \end{array} \right. \\ \text{Bernoulli}(p+q-p^2) \end{array}$$

Splitting of a Bernoulli process (upon coin flip q)

$$\begin{array}{c} \text{Bernoulli}(pq) \\ \uparrow q \\ \text{Bernoulli}(p) \\ \sqrt{1-q} \\ \text{Bernoulli}(p(1-q)) \end{array}$$

Poisson Approximation to the Binomial

- A Poisson random variable Z with parameter λ takes nonnegative integer values and is described by the PMF

$$p_Z(k) = e^{-\lambda} \frac{\lambda^k}{k!}, \quad k = 0, 1, 2, \dots$$

Its mean and variance are given by

$$\mathbf{E}[Z] = \lambda, \quad \text{var}(Z) = \lambda.$$

- For any fixed nonnegative integer k , the binomial probability

$$p_S(k) = \frac{n!}{(n-k)! k!} \cdot p^k (1-p)^{n-k}$$

converges to $p_Z(k)$, when we take the limit as $n \rightarrow \infty$ and $p = \lambda/n$, while keeping λ constant.

- In general, the Poisson PMF is a good approximation to the binomial as long as $\lambda = np$, n is very large, and p is very small.

6.2 THE POISSON PROCESS

(continuous time version of Bernoulli process)

Definition of the Poisson Process

An arrival process is called a Poisson process with rate λ if it has the following properties:

- (a) **(Time-homogeneity)** The probability $P(k, \tau)$ of k arrivals is the same for all intervals of the same length τ .
- (b) **(Independence)** The number of arrivals during a particular interval is independent of the history of arrivals outside this interval.
- (c) **(Small interval probabilities)** The probabilities $P(k, \tau)$ satisfy
(see below)

$$P(0, \tau) = 1 - \lambda\tau + o(\tau),$$

$$P(1, \tau) = \lambda\tau + o_1(\tau),$$

$$P(k, \tau) = o_k(\tau), \quad \text{for } k = 2, 3, \dots$$

Here, $o(\tau)$ and $o_k(\tau)$ are functions of τ that satisfy

$$\lim_{\tau \rightarrow 0} \frac{o(\tau)}{\tau} = 0, \quad \lim_{\tau \rightarrow 0} \frac{o_k(\tau)}{\tau} = 0.$$

Random Variables Associated with the Poisson Process and their Properties

- **The Poisson with parameter $\lambda\tau$.** This is the number N_τ of arrivals in a Poisson process with rate λ , over an interval of length τ . Its PMF, mean, and variance are

$$p_{N_\tau}(k) = P(k, \tau) = e^{-\lambda\tau} \frac{(\lambda\tau)^k}{k!}, \quad k = 0, 1, \dots,$$

$$\mathbf{E}[N_\tau] = \lambda\tau, \quad \text{var}(N_\tau) = \lambda\tau.$$

- **The exponential with parameter λ .** This is the time T until the first arrival. Its PDF, mean, and variance are

$$f_T(t) = \lambda e^{-\lambda t}, \quad t \geq 0, \quad \mathbf{E}[T] = \frac{1}{\lambda}, \quad \text{var}(T) = \frac{1}{\lambda^2}.$$

- **Memorylessness:** conditioned on $T_1 > t$,
the PDF of $T_1 - t$ is again exponential

small interval probabilities:

$$P(k, \delta) = \begin{cases} 1 - \lambda\delta + O(\delta^2) & \text{if } k = 0 \\ \lambda\delta + O(\delta^2) & \text{if } k = 1 \\ 0 + O(\delta^2) & \text{if } k > 1 \end{cases}$$

$\frac{O(\delta^2)}{\delta} \xrightarrow{\delta \rightarrow 0} 0$

- "Random Incidence": sampling method matters!
- Be careful about what you choose to sample
When you choose at random, what is it exactly that you're choosing at random

Independence Properties of the Poisson Process

- For any given time $t > 0$, the history of the process after time t is also a Poisson process, and is independent from the history of the process until time t .
- Let t be a given time and let \bar{T} be the time of the first arrival after time t . Then, $\bar{T} - t$ has an exponential distribution with parameter λ , and is independent of the history of the process until time t .

Alternative Description of the Poisson Process

1. Start with a sequence of independent exponential random variables T_1, T_2, \dots , with common parameter λ , and let these represent the interarrival times.
2. Record an arrival at times $T_1, T_1 + T_2, T_1 + T_2 + T_3$, etc.

Properties of the k th Arrival Time

- The k th arrival time is equal to the sum of the first k interarrival times

$$Y_k = T_1 + T_2 + \dots + T_k, \quad T_i \text{ are i.i.d. } \text{Exp}(\lambda)$$

and the latter are independent exponential random variables with common parameter λ .

- The mean and variance of Y_k are given by

$$\mathbf{E}[Y_k] = \mathbf{E}[T_1] + \dots + \mathbf{E}[T_k] = \frac{k}{\lambda},$$

$$\text{var}(Y_k) = \text{var}(T_1) + \dots + \text{var}(T_k) = \frac{k}{\lambda^2}.$$

- The PDF of Y_k is given by

$$f_{Y_k}(y) = \frac{\lambda^k y^{k-1} e^{-\lambda y}}{(k-1)!}, \quad y \geq 0,$$

and is known as the **Erlang PDF of order k** .

$$f_{Y_k}(y) \delta \approx P(k-1, y) \lambda \delta$$

- Sum of independent Poisson r.v.

The sum of independent Poisson random variables, with means/parameters μ and ν , is Poisson with mean/parameter $\mu + \nu$

Special property: Poisson, normal, ...

- Merge independent Poisson processes

Poisson (λ_1), Poisson (λ_2) \rightarrow Poisson ($\lambda_1 + \lambda_2$)

- Splitting a Poisson process with com flp q

Resulting streams are Poisson, rates λq , $\lambda(1-q)$

Properties of Sums of a Random Number of Random Variables

Let N, X_1, X_2, \dots be independent random variables, where N takes nonnegative integer values. Let $Y = X_1 + \dots + X_N$ for positive values of N , and let $Y = 0$ when $N = 0$.

- If X_i is Bernoulli with parameter p , and N is binomial with parameters m and q , then Y is binomial with parameters m and pq .
- If X_i is Bernoulli with parameter p , and N is Poisson with parameter λ , then Y is Poisson with parameter λp .
- If X_i is geometric with parameter p , and N is geometric with parameter q , then Y is geometric with parameter pq .
- If X_i is exponential with parameter λ , and N is geometric with parameter q , then Y is exponential with parameter λq .

6.3 SUMMARY AND DISCUSSION

Bernoulli/Poisson relation



$$n = \tau / \delta,$$

$$p = \lambda \delta$$

$$np = \lambda \tau$$

	POISSON	BERNOULLI
Times of Arrival	Continuous	Discrete
Arrival Rate	$\lambda/\text{unit time}$	$p/\text{per trial}$
PMF of # of Arrivals	Poisson	Binomial
Interarrival Time Distr.	Exponential	Geometric
Time to k -th arrival	Erlang	Pascal

poisson process is
a limiting case of
the Bernoulli process

7

Markov Chains

Excerpts from Introduction to Probability: Second Edition

by Dimitri P. Bertsekas and John N. Tsitsiklis ©

Massachusetts Institute of Technology

Contents

7.1. Discrete-Time Markov Chains	p. 340
7.2. Classification of States	p. 346
7.3. Steady-State Behavior	p. 352
7.4. Absorption Probabilities and Expected Time to Absorption . .	p. 362
7.5. Continuous-Time Markov Chains	p. 369
7.6. Summary and Discussion	p. 378
Problems	p. 380

7.1 DISCRETE-TIME MARKOV CHAINS

Specification of Markov Models

- A Markov chain model is specified by identifying:
 - (a) the set of states $\S = \{1, \dots, m\}$,
 - (b) the set of possible transitions, namely, those pairs (i, j) for which $p_{ij} > 0$, and,
 - (c) the numerical values of those p_{ij} that are positive.
- The Markov chain specified by this model is a sequence of random variables X_0, X_1, X_2, \dots , that take values in \S , and which satisfy

$$\mathbf{P}(X_{n+1} = j | X_n = i, X_{n-1} = i_{n-1}, \dots, X_0 = i_0) = p_{ij},$$

for all times n , all states $i, j \in \S$, and all possible sequences i_0, \dots, i_{n-1} of earlier states.

Chapman-Kolmogorov Equation for the n -Step Transition Probabilities

The n -step transition probabilities can be generated by the recursive formula

$$r_{ij}(n) = \sum_{k=1}^m r_{ik}(n-1)p_{kj}, \quad \text{for } n > 1, \text{ and all } i, j,$$

starting with

$$r_{ij}(1) = p_{ij}.$$

7.2 CLASSIFICATION OF STATES

Markov Chain Decomposition

- A Markov chain can be decomposed into one or more recurrent classes, plus possibly some transient states.
- A recurrent state is accessible from all states in its class, but is not accessible from recurrent states in other classes.
- A transient state is not accessible from any recurrent state.
- At least one, possibly more, recurrent states are accessible from a given transient state.

Periodicity

Consider a recurrent class R .

- The class is called **periodic** if its states can be grouped in $d > 1$ disjoint subsets S_1, \dots, S_d , so that all transitions from S_k lead to S_{k+1} (or to S_1 if $k = d$).
- The class is **aperiodic** (not periodic) if and only if there exists a time n such that $r_{ij}(n) > 0$, for all $i, j \in R$.

7.3 STEADY-STATE BEHAVIOR

Steady-State Convergence Theorem

Consider a Markov chain with a single recurrent class, which is aperiodic. Then, the states j are associated with steady-state probabilities π_j that have the following properties.

- (a) For each j , we have

$$\lim_{n \rightarrow \infty} r_{ij}(n) = \pi_j, \quad \text{for all } i.$$

- (b) The π_j are the unique solution to the system of equations below:

$$\begin{aligned}\pi_j &= \sum_{k=1}^m \pi_k p_{kj}, \quad j = 1, \dots, m, \\ 1 &= \sum_{k=1}^m \pi_k.\end{aligned}$$

- (c) We have

$$\begin{aligned}\pi_j &= 0, \quad \text{for all transient states } j, \\ \pi_j &> 0, \quad \text{for all recurrent states } j.\end{aligned}$$

Steady-State Probabilities as Expected State Frequencies

For a Markov chain with a single class which is aperiodic, the steady-state probabilities π_j satisfy

$$\pi_j = \lim_{n \rightarrow \infty} \frac{v_{ij}(n)}{n},$$

where $v_{ij}(n)$ is the expected value of the number of visits to state j within the first n transitions, starting from state i .

Expected Frequency of a Particular Transition

Consider n transitions of a Markov chain with a single class which is aperiodic, starting from a given initial state. Let $q_{jk}(n)$ be the expected number of such transitions that take the state from j to k . Then, regardless of the initial state, we have

$$\lim_{n \rightarrow \infty} \frac{q_{jk}(n)}{n} = \pi_j p_{jk}.$$

7.4 ABSORPTION PROBABILITIES AND EXPECTED TIME TO ABSORPTION

Absorption Probability Equations

Consider a Markov chain where each state is either transient or absorbing, and fix a particular absorbing state s . Then, the probabilities a_i of eventually reaching state s , starting from i , are the unique solution to the equations

$$\begin{aligned} a_s &= 1, \\ a_i &= 0, \quad \text{for all absorbing } i \neq s, \\ a_i &= \sum_{j=1}^m p_{ij} a_j, \quad \text{for all transient } i. \end{aligned}$$

Equations for the Expected Times to Absorption

Consider a Markov chain where all states are transient, except for a single absorbing state. The expected times to absorption, μ_1, \dots, μ_m , are the unique solution to the equations

$$\begin{aligned}\mu_i &= 0, && \text{if } i \text{ is the absorbing state,} \\ \mu_i &= 1 + \sum_{j=1}^m p_{ij} \mu_j, && \text{if } i \text{ is transient.}\end{aligned}$$

Equations for Mean First Passage and Recurrence Times

Consider a Markov chain with a single recurrent class, and let s be a particular recurrent state.

- The mean first passage times μ_i to reach state s starting from i , are the unique solution to the system of equations

$$\mu_s = 0, \quad \mu_i = 1 + \sum_{j=1}^m p_{ij} \mu_j, \quad \text{for all } i \neq s.$$

- The mean recurrence time μ_s^* of state s is given by

$$\mu_s^* = 1 + \sum_{j=1}^m p_{sj} \mu_j.$$

7.5 CONTINUOUS-TIME MARKOV CHAINS**Continuous-Time Markov Chain Assumptions**

- If the current state is i , the time until the next transition is exponentially distributed with a given parameter ν_i , independent of the past history of the process and of the next state.
- If the current state is i , the next state will be j with a given probability p_{ij} , independent of the past history of the process and of the time until the next transition.

Alternative Description of a Continuous-Time Markov Chain

Given the current state i of a continuous-time Markov chain, and for any $j \neq i$, the state δ time units later is equal to j with probability

$$q_{ij}\delta + o(\delta),$$

independent of the past history of the process.

Steady-State Convergence Theorem

Consider a continuous-time Markov chain with a single recurrent class. Then, the states j are associated with steady-state probabilities π_j that have the following properties.

- (a) For each j , we have

$$\lim_{t \rightarrow \infty} \mathbf{P}(X(t) = j \mid X(0) = i) = \pi_j, \quad \text{for all } i.$$

- (b) The π_j are the unique solution to the system of equations below:

$$\begin{aligned} \pi_j \sum_{k \neq j} q_{jk} &= \sum_{k \neq j} \pi_k q_{kj}, \quad j = 1, \dots, m, \\ 1 &= \sum_{k=1}^m \pi_k. \end{aligned}$$

- (c) We have

$$\begin{aligned} \pi_j &= 0, && \text{for all transient states } j, \\ \pi_j &> 0, && \text{for all recurrent states } j. \end{aligned}$$

7.6 SUMMARY AND DISCUSSION

Bayesian Statistical Inference

Excerpts from **Introduction to Probability: Second Edition**

by Dimitri P. Bertsekas and John N. Tsitsiklis ©

Massachusetts Institute of Technology

Contents

8.1. Bayesian Inference and the Posterior Distribution	p. 412
8.2. Point Estimation, Hypothesis Testing, and the MAP Rule	p. 420
8.3. Bayesian Least Mean Squares Estimation	p. 430
8.4. Bayesian Linear Least Mean Squares Estimation	p. 437
8.5. Summary and Discussion	p. 444
Problems	p. 446

Major Terms, Problems, and Methods in this Chapter

- Bayesian statistics treats unknown parameters as random variables with known prior distributions.
- In parameter estimation, we want to generate estimates that are close to the true values of the parameters in some probabilistic sense.
- In hypothesis testing, the unknown parameter takes one of a finite number of values, corresponding to competing hypotheses; we want to choose one of the hypotheses, aiming to achieve a small probability of error.
- Principal Bayesian inference methods:
 - (a) Maximum a posteriori probability (MAP) rule: Out of the possible parameter values/hypotheses, select one with maximum conditional/posterior probability given the data (Section 8.2).
 - (b) Least mean squares (LMS) estimation: Select an estimator/function of the data that minimizes the mean squared error between the parameter and its estimate (Section 8.3).
 - (c) Linear least mean squares estimation: Select an estimator which is a linear function of the data and minimizes the mean squared error between the parameter and its estimate (Section 8.4). This may result in higher mean squared error, but requires simple calculations, based only on the means, variances, and covariances of the random variables involved.

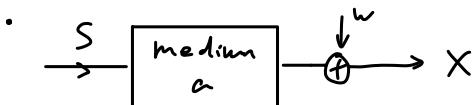
(instead of unknown constant)

numeric unknown $\in \mathbb{R}$

8.1 BAYESIAN INFERENCE AND THE POSTERIOR DISTRIBUTION

Summary of Bayesian Inference

- We start with a prior distribution p_{Θ} or f_{Θ} for the unknown random variable Θ .
- We have a model $p_{X|\Theta}$ or $f_{X|\Theta}$ of the observation vector X .
- After observing the value x of X , we form the posterior distribution of Θ , using the appropriate version of Bayes' rule.



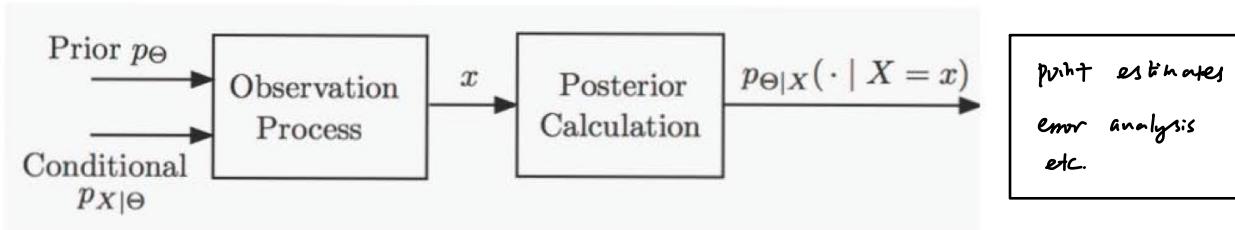
$$x = a s + w$$

model building: infer a
(know s, x)

variable estimation: infer s
(know a, x)

* The Bayesian inference framework: treat the unknown Θ as r.v. (instead of unknown constant)

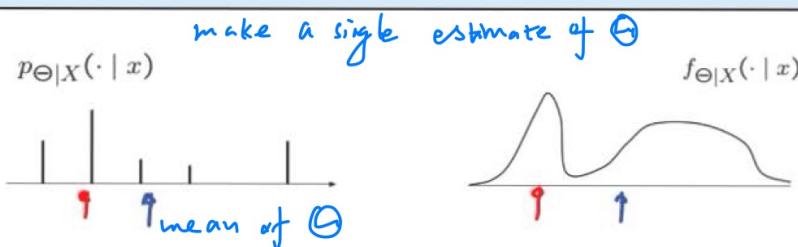
- Unknown Θ
 - treated as a random variable
 - prior distribution p_Θ or f_Θ
What we believe about Θ before we get any data
- Observation X r.v.
aug. data
- observation model $p_{X|\Theta}$ or $f_{X|\Theta}$
- Use appropriate version of the Bayes rule to find $p_{\Theta|X}(\cdot | X = x)$ or $f_{\Theta|X}(\cdot | X = x)$



12

Point estimates in Bayesian inference

The complete answer is a posterior distribution:
PMF $p_{\Theta|X}(\cdot | x)$ or PDF $f_{\Theta|X}(\cdot | x)$



- Maximum a posteriori probability (MAP):

$$p_{\Theta|X}(\theta^* | x) = \max_{\theta} p_{\Theta|X}(\theta | x),$$

$$f_{\Theta|X}(\theta^* | x) = \max_{\theta} f_{\Theta|X}(\theta | x).$$

- Conditional expectation: $E[\Theta | X = x]$ (LMS: Least Mean Squares)

the way to process the data
↓ ↓ data

estimate: $\hat{\theta} = g(x)$

(number) r.v. r.v.

estimator: $\hat{\Theta} = g(X)$

(random variable)

g is also called an estimator

The Four Versions of Bayes' Rule

- Θ discrete, X discrete:

$$p_{\Theta|X}(\theta|x) = \frac{p_\Theta(\theta)p_{X|\Theta}(x|\theta)}{\sum_{\theta'} p_\Theta(\theta')p_{X|\Theta}(x|\theta')}.$$

• example:

discrete signal Θ

random noise $w \sim N(0, \sigma^2)$ in dep. of Θ

$$X = \Theta + w$$

- Θ discrete, X continuous:

$$p_{\Theta|X}(\theta|x) = \frac{p_\Theta(\theta)f_{X|\Theta}(x|\theta)}{\sum_{\theta'} p_\Theta(\theta')f_{X|\Theta}(x|\theta')}.$$

- Θ continuous, X discrete:

$$f_{\Theta|X}(\theta|x) = \frac{f_\Theta(\theta)p_{X|\Theta}(x|\theta)}{\int f_\Theta(\theta')p_{X|\Theta}(x|\theta') d\theta'}.$$

- Θ continuous, X continuous:

$$f_{\Theta|X}(\theta|x) = \frac{f_\Theta(\theta)f_{X|\Theta}(x|\theta)}{\int f_\Theta(\theta')f_{X|\Theta}(x|\theta') d\theta'}.$$

• example:

noisy signal

Θ : signal > independent normals

w : noise

$$X = \Theta + w$$

- Performance evaluation of an estimator $\hat{\Theta}$

$$P(\hat{\Theta} \neq \Theta | X=x) \quad E[(\hat{\Theta} - \Theta)^2 | X=x] \quad \text{inc } p_{\Theta|X}(\cdot|x)$$

$$P(\hat{\Theta} \neq \Theta)$$

hypothesis testing

$$E[(\hat{\Theta} - \Theta)^2]$$

estimate

total probability, then
exp

8.2 POINT ESTIMATION, HYPOTHESIS TESTING, AND THE MAP RULE

The Maximum a Posteriori Probability (MAP) Rule

- Given the observation value x , the MAP rule selects a value $\hat{\theta}$ that maximizes over θ the posterior distribution $p_{\Theta|X}(\theta|x)$ (if Θ is discrete) or $f_{\Theta|X}(\theta|x)$ (if Θ is continuous).
 - Equivalently, it selects $\hat{\theta}$ that maximizes over θ : *see formula in last page,
the denominator is only a factor of X*
- $$p_{\Theta}(\theta)p_{X|\Theta}(x|\theta) \quad (\text{if } \Theta \text{ and } X \text{ are discrete}),$$
- $$p_{\Theta}(\theta)f_{X|\Theta}(x|\theta) \quad (\text{if } \Theta \text{ is discrete and } X \text{ is continuous}),$$
- $$f_{\Theta}(\theta)p_{X|\Theta}(x|\theta) \quad (\text{if } \Theta \text{ is continuous and } X \text{ is discrete}),$$
- $$f_{\Theta}(\theta)f_{X|\Theta}(x|\theta) \quad (\text{if } \Theta \text{ and } X \text{ are continuous}).$$
- If Θ takes only a finite number of values, the MAP rule minimizes (over all decision rules) the probability of selecting an incorrect hypothesis. This is true for both the unconditional probability of error and the conditional one, given any observation value x .

- conditional probability of error:*

$P(\hat{\theta} \neq \theta | X=x)$ is smallest under the MAP rule

overall probability of error:

$$P(\hat{\theta} \neq \theta) = \sum_x P(\hat{\theta} \neq \theta | X=x) p_X(x) \quad \text{also smallest under the MAP rule}$$

t $g(x)$

Point Estimates

- An **estimator** is a random variable of the form $\hat{\Theta} = g(X)$, for some function g . Different choices of g correspond to different estimators.
- An **estimate** is the value $\hat{\theta}$ of an estimator, as determined by the realized value x of the observation X .
- Once the value x of X is observed, the **Maximum a Posteriori Probability (MAP)** estimator, sets the estimate $\hat{\theta}$ to a value that maximizes the posterior distribution over all possible values of θ .
- Once the value x of X is observed, the **Conditional Expectation (LMS)** estimator sets the estimate $\hat{\theta}$ to $\mathbf{E}[\Theta | X = x]$.

The MAP Rule for Hypothesis Testing

- Given the observation value x , the MAP rule selects a hypothesis H_i for which the value of the posterior probability $\mathbf{P}(\Theta = \theta_i | X = x)$ is largest.
- Equivalently, it selects a hypothesis H_i for which $p_{\Theta}(\theta_i)p_{X|\Theta}(x | \theta_i)$ (if X is discrete) or $p_{\Theta}(\theta_i)f_{X|\Theta}(x | \theta_i)$ (if X is continuous) is largest.
- The MAP rule minimizes the probability of selecting an incorrect hypothesis for any observation value x , as well as the probability of error over all decision rules.

8.3 BAYESIAN LEAST MEAN SQUARES ESTIMATION

LMS relevant for estimation, not hypothesis testing

• Mean Squared Error

(MSE):

$$E[(\Theta - \hat{\theta})^2]$$

• optimal mean squared error:

$$E[(\Theta - E[\Theta])^2] = \text{var}(\Theta)$$

Key Facts About Least Mean Squares Estimation

- In the absence of any observations, $E[(\Theta - \hat{\theta})^2]$ is minimized when $\hat{\theta} = E[\Theta]$:

$$E[(\Theta - E[\Theta])^2] \leq E[(\Theta - \hat{\theta})^2], \quad \text{for all } \hat{\theta}.$$

- For any given value x of X , $E[(\Theta - \hat{\theta})^2 | X = x]$ is minimized when $\hat{\theta} = E[\Theta | X = x]$: LMS estimator: $\hat{\Theta}_{LMS} = E[\Theta | X]$

$$E[(\Theta - E[\Theta | X = x])^2 | X = x] \leq E[(\Theta - \hat{\theta})^2 | X = x], \quad \text{for all } \hat{\theta}.$$

- Out of all estimators $g(X)$ of Θ based on X , the mean squared estimation error $E[(\Theta - g(X))^2]$ is minimized when $g(X) = E[\Theta | X]$: $\hat{\Theta}_{LMS}$

$$E[(\Theta - E[\Theta | X])^2] \leq E[(\Theta - g(X))^2], \quad \text{for all estimators } g(X).$$

- Expected performance, once we have a measurement:

$$\text{MSE} = E[(\Theta - E[\Theta | X = x])^2 | X = x] = \text{var}(\Theta | X = x)$$

- Expected performance of the design:

$$\text{MSE} = E[(\Theta - E[\Theta | X])^2] = E[\text{var}(\Theta | X)]$$

- LMS same as MAP for posterior unimodal, symmetric around the mean



- multidimensional Θ, X :

$$\Theta = (\Theta_1, \dots, \Theta_n), \quad X = (X_1, \dots, X_n)$$

$$\hat{\Theta}_j = E[\Theta_j | X_1 = x_1, \dots, X_n = x_n]$$

Properties of the Estimation Error $\tilde{\Theta} = \hat{\Theta} - \Theta$

- The estimation error $\tilde{\Theta}$ is **unbiased**, i.e., it has zero unconditional and conditional mean:

$$\mathbf{E}[\tilde{\Theta}] = 0, \quad \mathbf{E}[\tilde{\Theta} | X = x] = 0, \quad \text{for all } x.$$

- The estimation error $\tilde{\Theta}$ is uncorrelated with the estimate $\hat{\Theta}$:

$$\text{cov}(\hat{\Theta}, \tilde{\Theta}) = 0.$$

- The variance of Θ can be decomposed as

$$\text{var}(\Theta) = \text{var}(\hat{\Theta}) + \text{var}(\tilde{\Theta}).$$

8.4 BAYESIAN LINEAR LEAST MEAN SQUARES ESTIMATION

• Estimators of form

$$\hat{\Theta} = aX + b$$

$$\begin{aligned} & \text{minimize}_{\text{w.r.t. } a, b} E[(\Theta - aX - b)^2] \\ & \text{only means, variances, covariances matter} \end{aligned}$$

$$\cdot |p| = 1. \quad \hat{\Theta}_L = \Theta$$

Linear LMS Estimation Formulas

- The linear LMS estimator $\hat{\Theta}$ of Θ based on X is

$$\hat{\Theta} = \mathbf{E}[\Theta] + \frac{\text{cov}(\Theta, X)}{\text{var}(X)}(X - \mathbf{E}[X]) = \mathbf{E}[\Theta] + \rho \frac{\sigma_\Theta}{\sigma_X} (X - \mathbf{E}[X]),$$

↑ baseline ↑ correction

where

$$\rho = \frac{\text{cov}(\Theta, X)}{\sigma_\Theta \sigma_X}$$

is the correlation coefficient.

- The resulting mean squared estimation error is equal to

$$E[(\hat{\Theta}_L - \Theta)^2] = (1 - \rho^2)\sigma_\Theta^2.$$

- If $\mathbf{E}[\Theta | X]$ is linear in X , then $\hat{\Theta}_{\text{LMS}} = \hat{\Theta}_{\text{LLMS}}$

8.5 SUMMARY AND DISCUSSION

LLMS with multiple observations

- Unknown Θ ; observations $X = (\mathbf{x}_1, \dots, \mathbf{x}_n)$
 $\xrightarrow{\text{linear}}$
- Consider estimators of the form: $\hat{\Theta} = a_1X_1 + \dots + a_nX_n + b$

• Depending on what we know about $\Theta \sim X$, can consider

$$\hat{\Theta} = a_1x_1 + a_2x_2^2 + a_3x_3^3 + b, \quad \hat{\Theta} = a_1x_1 + a_2e^{x_2} + a_3\ln x_3 + b, \dots$$

they're still LLMS estimators! (linearity in the coefficients is what matters)

Linear model with normal noise

- $X \sim N(\mu, \sigma^2)$.

$$f_X(x) = C e^{-(\alpha x^2 + \beta x + \gamma)}, \quad \alpha > 0 \quad \sim N(-\frac{\beta}{2\alpha}, \frac{1}{2\alpha})$$

- estimate a normal r.v., with additive normal noise
multiple observations

$$X_1 = \Theta + w_1 \quad \Theta \sim N(\mu_0, \sigma_0^2), \quad w_i \sim N(0, \sigma_i^2)$$

⋮

$$X_n = \Theta + w_n \quad \Theta, w_1, \dots, w_n \text{ independent}$$

- posterior $f(\Theta | X = \theta | x)$ is normal

$$\hat{\theta}_{\text{MAP}} = \hat{\theta}_{\text{LMS}} = E[\Theta | X = x] = \frac{\sum_{i=0}^n \frac{x_i}{\sigma_i^2}}{\sum_{i=0}^n \frac{1}{\sigma_i^2}}$$

- $\hat{\theta}$ linear, $\hat{\theta} = \alpha_0 + \alpha_1 x_1 + \dots + \alpha_n x_n$
weighted average of μ_0 (prior mean) and x_i (observations)
- mean squared error

$$E[(\Theta - \hat{\Theta})^2 | X = \underline{x}] \underset{\bullet}{=} E[(\Theta - \hat{\Theta})^2] = 1 / \sum_{i=0}^n \frac{1}{\sigma_i^2}$$

- Bayesian confidence interval
 $P(\text{True position } \in \text{interval} | \text{data}) = 0.95$



linear functions + normal r.v.s

Linear normal models •

underlying

- Θ_j and X_i are linear functions of independent normal random variables
- $f_{\Theta|X}(\theta|x) = c(x) \exp \left\{ -\text{quadratic}(\theta_1, \dots, \theta_m) \right\}$ linear regression
- MAP estimate: maximize over $(\theta_1, \dots, \theta_m)$; linear equations
(minimize quadratic function)
- $\widehat{\Theta}_{\text{MAP},j}$: linear function of $X = (X_1, \dots, X_n)$
- Facts:
 - $\widehat{\Theta}_{\text{MAP},j} = E[\Theta_j | X]$
 - marginal posterior PDF of Θ_j : $f_{\Theta_j|X}(\theta_j | x)$, is normal
 - MAP estimate based on the joint posterior PDF:
same as MAP estimate based on the marginal posterior PDF
 - $E[(\widehat{\Theta}_{i,\text{MAP}} - \Theta_i)^2 | X = x]$: same for all x

9

Classical Statistical Inference

Excerpts from **Introduction to Probability: Second Edition**

by Dimitri P. Bertsekas and John N. Tsitsiklis ©

Massachusetts Institute of Technology

Contents

9.1. Classical Parameter Estimation	p. 460
9.2. Linear Regression	p. 475
9.3. Binary Hypothesis Testing	p. 485
9.4. Significance Testing	p. 495
9.5. Summary and Discussion	p. 506
Problems	p. 507

Major Terms, Problems, and Methods in this Chapter

- $P_X(x|\theta)$
 θ is a parameter

- **Classical statistics** treats unknown parameters as constants to be determined. A separate probabilistic model is assumed for each possible value of the unknown parameter.
- In **parameter estimation**, we want to generate estimates that are nearly correct under any possible value of the unknown parameter.
- In **hypothesis testing**, the unknown parameter takes a finite number m of values ($m \geq 2$), corresponding to competing hypotheses; we want to choose one of the hypotheses, aiming to achieve a small probability of error under any of the possible hypotheses.
- In **significance testing**, we want to accept or reject a single hypothesis, while keeping the probability of false rejection suitably small.
- Principal classical inference methods in this chapter:
 - (a) **Maximum likelihood (ML) estimation:** Select the parameter that makes the observed data “most likely,” i.e., maximizes the probability of obtaining the data at hand (Section 9.1).
 - (b) **Linear regression:** Find the linear relation that matches best a set of data pairs, in the sense that it minimizes the sum of the squares of the discrepancies between the model and the data (Section 9.2).
 - (c) **Likelihood ratio test:** Given two hypotheses, select one based on the ratio of their “likelihoods,” so that certain error probabilities are suitably small (Section 9.3).
 - (d) **Significance testing:** Given a hypothesis, reject it if and only if the observed data falls within a certain rejection region. This region is specially designed to keep the probability of false rejection below some threshold (Section 9.4).

9.1 CLASSICAL PARAMETER ESTIMATION

- mean squared error (MSE):
 $E[(\hat{\theta} - \theta)^2]$
- bias: whether an estimator is systematically above or below the true value
- $E[(\hat{\theta} - \theta)^2] = \text{Var}(\hat{\theta}) + (\text{bias})^2$
- $\sqrt{\text{Var}(\hat{\theta})}$: standard error

Terminology Regarding Estimators

Let $\hat{\Theta}_n$ be an **estimator** of an unknown parameter θ , that is, a function of n observations X_1, \dots, X_n whose distribution depends on θ .

- The **estimation error**, denoted by $\tilde{\Theta}_n$, is defined by $\tilde{\Theta}_n = \hat{\Theta}_n - \theta$.
- The **bias** of the estimator, denoted by $b_\theta(\hat{\Theta}_n)$, is the expected value of the estimation error:

$$b_\theta(\hat{\Theta}_n) = E_\theta[\hat{\Theta}_n] - \theta.$$
- The expected value, the variance, and the bias of $\hat{\Theta}_n$ depend on θ , while the estimation error depends in addition on the observations X_1, \dots, X_n .
- We call $\hat{\Theta}_n$ **unbiased** if $E_\theta[\hat{\Theta}_n] = \theta$, for every possible value of θ .
- We call $\hat{\Theta}_n$ **asymptotically unbiased** if $\lim_{n \rightarrow \infty} E_\theta[\hat{\Theta}_n] = \theta$, for every possible value of θ .
- We call $\hat{\Theta}_n$ **consistent** if the sequence $\hat{\Theta}_n$ converges to the true value of the parameter θ , in probability, for every possible value of θ .

Maximum Likelihood Estimation

- We are given the realization $x = (x_1, \dots, x_n)$ of a random vector $X = (X_1, \dots, X_n)$, distributed according to a PMF $p_X(x; \theta)$ or PDF $f_X(x; \theta)$.
- The maximum likelihood (ML) estimate is a value of θ that maximizes the likelihood function, $p_X(x; \theta)$ or $f_X(x; \theta)$, over all θ .
- The ML estimate of a one-to-one function $h(\theta)$ of θ is $h(\hat{\theta}_n)$, where $\hat{\theta}_n$ is the ML estimate of θ (the invariance principle).
- When the random variables X_i are i.i.d., and under some mild additional assumptions, each component of the ML estimator is consistent and asymptotically normal.

$$\hat{\theta}_{\text{ML}} = \underset{\theta}{\operatorname{argmax}} p_X(x; \theta)$$

Estimates of the Mean and Variance of a Random Variable $= \sigma^2$

Let the observations X_1, \dots, X_n be i.i.d., with mean θ and variance v that are unknown.

- The sample mean

$$M_n = \frac{X_1 + \dots + X_n}{n}$$

is an unbiased estimator of θ , and its mean squared error is $v/n = \frac{\sigma^2}{n}$

- Two variance estimators are

$$\bar{S}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - M_n)^2, \quad \hat{S}_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - M_n)^2.$$

- The estimator \bar{S}_n^2 coincides with the ML estimator if the X_i are normal. It is biased but asymptotically unbiased. The estimator \hat{S}_n^2 is unbiased. For large n , the two variance estimators essentially coincide.

$$\hat{\Theta} = \frac{1}{n} \sum_{i=1}^n g(X_i)$$

Confidence Intervals

- A **confidence interval** for a scalar unknown parameter θ is an interval whose endpoints $\hat{\Theta}_n^-$ and $\hat{\Theta}_n^+$ bracket θ with a given high probability.
- $\hat{\Theta}_n^-$ and $\hat{\Theta}_n^+$ are random variables that depend on the observations X_1, \dots, X_n .
- A $1 - \alpha$ confidence interval is one that satisfies

$$\mathbf{P}_\theta(\hat{\Theta}_n^- \leq \theta \leq \hat{\Theta}_n^+) \geq 1 - \alpha,$$

for all possible values of θ .

• interpretation: 95% of the time, your method will capture the true value θ

or

in one experiment, the probability that θ falls in the confidence interval is 95%

- CI for the estimation of the mean:

$$\mathbf{P}\left(\hat{\Theta}_n - \frac{1.96\sigma}{\sqrt{n}} \leq \theta \leq \hat{\Theta}_n + \frac{1.96\sigma}{\sqrt{n}}\right) \approx 0.95$$

Θ^- Θ^+

- for σ unknown:

{ we upper bound on σ
 ad hoc estimate for σ
 e.g. $\hat{\sigma} = \sqrt{\hat{\Theta}((1-\hat{\Theta}))}$, $\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \hat{\Theta})^2$
 sample mean estimate for σ^2
 $\sigma^2 = E[(X_i - \theta)^2] \rightarrow \frac{1}{n-1} \sum_{i=1}^n (X_i - \hat{\Theta})^2 \rightarrow \hat{\sigma}^2$

9.2 LINEAR REGRESSION**Linear Regression**

Given n data pairs (x_i, y_i) , the estimates that minimize the sum of the squared residuals are given by

$$\hat{\theta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad \hat{\theta}_0 = \bar{y} - \hat{\theta}_1 \bar{x},$$

where

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i.$$

Bayesian Linear Regression

- **Model:**

- (a) We assume a linear relation $Y_i = \Theta_0 + \Theta_1 x_i + W_i$.
- (b) The x_i are modeled as known constants.
- (c) The random variables $\Theta_0, \Theta_1, W_1, \dots, W_n$ are normal and independent.
- (d) The random variables Θ_0 and Θ_1 have mean zero and variances σ_0^2, σ_1^2 , respectively.
- (e) The random variables W_i have mean zero and variance σ^2 .

- **Estimation Formulas:**

Given the data pairs (x_i, y_i) , the MAP estimates of Θ_0 and Θ_1 are

$$\hat{\theta}_1 = \frac{\sigma_1^2}{\sigma^2 + \sigma_1^2 \sum_{i=1}^n (x_i - \bar{x})^2} \cdot \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}),$$

$$\hat{\theta}_0 = \frac{n\sigma_0^2}{\sigma^2 + n\sigma_0^2} (\bar{y} - \hat{\theta}_1 \bar{x}),$$

where

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i.$$

9.3 BINARY HYPOTHESIS TESTING

Likelihood Ratio Test (LRT)

- Start with a target value α for the false rejection probability.
- Choose a value for ξ such that the false rejection probability is equal to α :

$$\mathbf{P}(L(X) > \xi; H_0) = \alpha.$$

- Once the value x of X is observed, reject H_0 if $L(x) > \xi$.

Neyman-Pearson Lemma

Consider a particular choice of ξ in the LRT, which results in error probabilities

$$\mathbf{P}(L(X) > \xi; H_0) = \alpha, \quad \mathbf{P}(L(X) \leq \xi; H_1) = \beta.$$

Suppose that some other test, with rejection region R , achieves a smaller or equal false rejection probability:

$$\mathbf{P}(X \in R; H_0) \leq \alpha.$$

Then,

$$\mathbf{P}(X \notin R; H_1) \geq \beta,$$

with strict inequality $\mathbf{P}(X \notin R; H_1) > \beta$ when $\mathbf{P}(X \in R; H_0) < \alpha$.

9.4 SIGNIFICANCE TESTING

Significance Testing Methodology

A statistical test of a hypothesis H_0 is to be performed, based on the observations X_1, \dots, X_n .

- The following steps are carried out before the data are observed.
 - (a) Choose a **statistic** S , that is, a scalar random variable that will summarize the data to be obtained. Mathematically, this involves the choice of a function $h : \Re^n \rightarrow \Re$, resulting in the statistic $S = h(X_1, \dots, X_n)$.
 - (b) Determine the **shape of the rejection region** by specifying the set of values of S for which H_0 will be rejected as a function of a yet undetermined critical value ξ .
 - (c) Choose the **significance level**, i.e., the desired probability α of a false rejection of H_0 .
 - (d) Choose the **critical value** ξ so that the probability of false rejection is equal (or approximately equal) to α . At this point, the rejection region is completely determined.
- Once the values x_1, \dots, x_n of X_1, \dots, X_n are observed:
 - (i) Calculate the value $s = h(x_1, \dots, x_n)$ of the statistic S .
 - (ii) Reject the hypothesis H_0 if s belongs to the rejection region.

The Chi-Square Test:

- Use the statistic

$$S = \sum_{k=1}^m N_k \log \left(\frac{N_k}{n\theta_k^*} \right)$$

(or possibly the related statistic T) and a rejection region of the form

reject H_0 if $2S > \gamma$

(or $T > \gamma$, respectively).

- The critical value γ is determined from the CDF tables for the χ^2 distribution with $m - 1$ degrees of freedom so that

$$\mathbf{P}(2S > \gamma; H_0) = \alpha,$$

where α is a given significance level.

9.5 SUMMARY AND DISCUSSION

MIT OpenCourseWare

<https://ocw.mit.edu>

Resource: Introduction to Probability

John Tsitsiklis and Patrick Jaillet

The following may not correspond to a particular course on MIT OpenCourseWare, but has been provided by the author as an individual learning resource.

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>