

PROBABILISTIC MODELS – PART 2: FUNDAMENTAL CONCEPTS OF PROBABILITY

Gerhard Widmer

Institute of Computational Perception
Johannes Kepler University
Linz, Austria

gerhard.widmer@jku.at
www.cp.jku.at/people/widmer



October 2, 2025

Many thanks to Daphne Koller, Nir Friedman, Stuart Russell, and Peter Norvig
for making these available
(pgm.stanford.edu; aima.cs.berkeley.edu).

Do not distribute!

Goals of this Lecture

- ▶ Establish terminology and mathematical notation for this class
- ▶ Remind you of fundamental concepts of probability
- ▶ Introduce the probabilistic query
- ▶ Show that the answer to any probabilistic query can be computed from the full joint distribution via *Inference by Enumeration*
- ▶ Introduce the central concept for this entire class: *Independencies*
- ▶ Motivate you to internalise these concepts
(to be able to follow the rest of the class)

Basics of Probability: Random Variables

Definition

A **Random Variable** X is a variable with a fixed (finite or infinite) **domain** (i.e., set of possible values) $Val(X)$, which represents some aspect of the system's world (or of the system's internal state).

Different types:

- ▶ **Boolean variables** take values $\in \{false, true\}$ (or $\{0, 1\}$)
- ▶ **Discrete (categorical) variables** have a finite domain of symbolic values.
Example: *Season* with $Val(Season) = \{spring, summer, fall, winter\}$
- ▶ **Continuous (real-valued) variables** take a numeric value $\in \mathbb{R}$
Most real-world sensors (e.g., temperature, sonar, radar, ...) give real-valued readings.

👉 **Basic atomic building blocks of our models and world representation.**

A Simple Example

Consider a system for medical diagnosis

in a very simple world where there are only

- ▶ two diseases (not directly observable): influenza (flu) and hayfever;
- ▶ two symptoms (possibly observable): congestion and muscle pain; and
- ▶ four seasons (observable): spring, summer, fall, and winter.

and where at a given moment, there is exactly one patient to be diagnosed.

World model using 5 random variables:

<i>Flu</i>	$\in \{true, false\}$
<i>Hayfever</i>	$\in \{true, false\}$
<i>Congestion</i>	$\in \{true, false\}$
<i>MusclePain</i>	$\in \{true, false\}$
<i>Season</i>	$\in \{spring, summer, fall, winter\}$

Note the use of boolean variables for diseases and symptoms – to allow for the case where a patient may have both symptoms and/or both diseases!

Basics of Probability: Events

Definition

An **Event** is a fixed assignment of values to some or all of the variables in the system's world.

Definition

An **Atomic Event** is an event where **all** random variables in the system's world have a specific value assigned.

Notes:

- ▶ An *atomic event* corresponds to one particular possible state of the world
- ▶ An *event* corresponds to a *set* of possible states of the world (it comprises all those states that share the specific values fixed in the event)
- ▶ Events (possible value assignments to variables) are what the system will try to compute probabilities for.

A Simple Example

Simple medical diagnosis world:

<i>Flu</i>	$\in \{true, false\}$
<i>Hayfever</i>	$\in \{true, false\}$
<i>Congestion</i>	$\in \{true, false\}$
<i>MusclePain</i>	$\in \{true, false\}$
<i>Season</i>	$\in \{spring, summer, fall, winter\}$

Some events:

- ▶ $Season = spring \wedge MusclePain = true$
denotes all situations where it is spring and the patient has a muscle pain (regardless of whether or not she has the flu, hayfever, or congestion)
- ▶ $Congestion = false$
- ▶ $Season = fall \wedge MusclePain = true \wedge Congestion = false \wedge Hayfever = false \wedge Flu = false$ (an atomic event)

👉 **Total number of possible *atomic* events:** $2 \times 2 \times 2 \times 2 \times 4 = 64$

Some Notational Conventions (1)

NOTATION

Examples	Type of object
\mathcal{X}	Complete set of variables that a model is defined over
A, X_i, E_j	Single random variable (e.g., $X_i \in \mathcal{X}$)
$\mathbf{X}, \mathbf{Y}, \mathbf{Z}, \mathbf{E}, \mathbf{U}, \dots$	Sets of random variables (e.g., $\mathbf{X} = \{X_1, X_2, X_3, X_4\}$)
$Val(X)$	Domain of a variable X (set of values that X can take)
x^0, x^1, \dots	$x^i = i^{th}$ value in the domain $Val(X)$ of variable X
x^0, x^1	Special convention for boolean variables: $x^0 = false, x^1 = true$
$\mathbf{x}_i, \mathbf{e}, \dots$	Set/list of specific values of a variable (e.g., $\mathbf{E} = e$)
$Val(\mathbf{X})$	Set of possible value assignments to the variables in \mathbf{X} (= Cartesian product $Val(X_1) \times Val(X_2) \times \dots$)
$P(X = x^i)$	Probability (probability mass function)
$P(x^i), P(x^i, y^j)$	Shorthand notation for $P(X = x^i), P(X = x^i \wedge Y = y^j)$
$P(X, Y)$	Probability distribution over discrete variables X, Y

Probabilities and Distributions

Definition

Given: A set \mathcal{S} of (mutually exclusive) atomic events.

A **Probability Distribution** over \mathcal{S} is a function $P : \mathcal{S} \mapsto \mathbb{R}$ that satisfies the following:

- 1 $P(\alpha) \geq 0$ for all $\alpha \in \mathcal{S}$
- 2 $P(\Omega) = 1$, where $\Omega = \bigcup_{\alpha_i \in \mathcal{S}} \alpha_i$ (the disjunction of all possible events)^a
- 3 If $\alpha, \beta \in \mathcal{S}$ and $\alpha \cap \beta = \emptyset$, then $P(\alpha \cup \beta) = P(\alpha) + P(\beta)$.

A **Probability** $P(\alpha)$ is the value that the probability distribution P assigns to the specific event α .

^aWe will use \cap and \cup to denote logical AND and OR, respectively, between events

- 👉 **Look at the argument of an expression $P(\cdot)$** to decide whether the P denotes a *probability* (a single value) or a *distribution* (a list of values):
- ▶ If the argument is an event, it is a probability.
 - ▶ If the argument is a (set of) random variable(s), it is a distribution.

Interpretations of Probability

Defining what probabilities “*mean*” is a deep philosophical problem (with a long history ...)

Intuitively:

“The probability $P(\alpha)$ of an event α quantifies the degree of confidence that α will occur.”

- ▶ $P(\alpha) = 1$: we are certain that α will occur
- ▶ $P(\alpha) = 0$: we are certain that α will *not* occur
- ▶ But what does $P(\alpha) = x$ mean, if $0 < x < 1$?

Interpretations of Probability

Two main philosophical stances:

'Frequentist interpretation':

The probability of an event is the **proportion of times** that the event α would occur if we repeated the experiment an *infinite number of times*.

Problem: How does this apply to, for example, “*the probability of rain this afternoon is 0.7*”?

'Subjectivist interpretation':

The probability of an event expresses a **subjective degree of belief** that the event α will happen (or the degree to which α is supported by available evidence).

In this class:

- ▶ Will adopt the subjectivist (Bayesian) interpretation
- ▶ $P(x)$ represents the **system's degree of belief** that x is true in the world.
- ▶ (Of course, will have to resort to frequentist view when we want to learn probabilities from observations ...)

Joint and Marginal Distributions

Definition

In a world defined by a set of random variables \mathcal{X} , the probability distribution over *all atomic events* possible over \mathcal{X} is called **Full Joint Distribution** over \mathcal{X} . It assigns probabilities to all possible value combinations of the \mathcal{X} .

Definition

A probability distribution defined over the events induced by a subset $X \subset \mathcal{X}$ of variables is called **Marginal Distribution** over X .
It assigns probabilities to all value combinations of the selected X .

Definition

Special Case:

A probability distribution defined over the values of a *single* variable $X \in \mathcal{X}$ is called the **Marginal Distribution** of variable X .

Some More Notational Conventions

NOTATION

- **Probability of conjunction of events:** Write

$$\boxed{P(X = x, Y = y)} \quad \text{or} \quad \boxed{P(x, y)} \quad \text{instead of} \quad P((X = x) \cap (Y = y))$$

- **Joint distribution over sets of variables:** Write

$$P(\mathbf{X}) = P(X_1, X_2, \dots, X_k)$$

to denote joint distribution over variables $\mathbf{X} = \{X_1, X_2, \dots, X_k\}$

- **Joint distribution over several sets of variables:** Write

$$P(\mathbf{X}, \mathbf{Y}) = P(X_1, \dots, X_k, Y_1, \dots, Y_l)$$

to denote joint distribution over variable set $X \cup Y$

- **Conditional distributions:** Write

$$P(\mathbf{X} \mid \mathbf{Y}) = P(X_1, X_2, \dots, X_k \mid Y_1, Y_2, \dots, Y_l)$$

to denote the joint distributions over \mathbf{X} , conditioned on values of \mathbf{Y} .

Joint and Marginal Distributions

Example: Diagnosis world with 3 (Boolean) random variables

$$\mathcal{X} = \{Flu, MusclePain(MP), Fever\}$$

Full Joint Distribution:

$P(Flu, MP, Fever)$	$musclepain$		$\neg musclepain$	
	$fever$	$\neg fever$	$fever$	$\neg fever$
flu	0.15	0.1	0.1	0.03
$\neg flu$	0.03	0.04	0.05	0.5

Note:

- ▶ Sum over all entries is 1.0 (a proper distribution)

Joint and Marginal Distributions

$P(Flu, MP, Fever)$	$musclepain$		$\neg musclepain$	
	$fever$	$\neg fever$	$fever$	$\neg fever$
flu	0.15	0.1	0.1	0.03
$\neg flu$	0.03	0.04	0.05	0.5

Marginal distributions

can be computed by summing over the full joint distribution:

$$\begin{aligned}
P(flu) &= P(flu, mp, fever) + P(flu, mp, \neg fever) + \\
&\quad P(flu, \neg mp, fever) + P(flu, \neg mp, \neg fever) \\
&= \sum_{x,y} P(flu, MP = x, Fever = y) \\
&= 0.15 + 0.1 + 0.1 + 0.03 = \underline{0.38} \\
P(\neg flu) &= 0.03 + 0.04 + 0.05 + 0.5 = 0.62
\end{aligned}$$

Marginal distribution $P(Flu) = \{0.62, 0.38\}$

$P(Flu, MP, Fever)$	$musclepain$		$\neg musclepain$	
	$fever$	$\neg fever$	$fever$	$\neg fever$
flu	0.15	0.1	0.1	0.03
$\neg flu$	0.03	0.04	0.05	0.5

Likewise:

$$P(\neg musclepain) = 0.1 + 0.03 + 0.05 + 0.5 = \underline{0.68}$$

$$P(\text{musclepain}) = 0.15 + 0.1 + 0.03 + 0.04 = \underline{0.32}$$

Marginal distribution $P(MusclePain) = \{0.68, 0.32\}$

$$P(\neg fever) = 0.1 + 0.04 + 0.03 + 0.5 = \underline{0.67}$$

$$P(fever) = 0.15 + 0.03 + 0.1 + 0.05 = \underline{0.33}$$

Marginal distribution $P(Fever) = \{0.67, 0.33\}$

We are computing the probability of the disjunction (\cup) of all atomic events that

 **Inference by Enumeration** (see below)

Conditional Probabilities

Conditional probabilities

- ▶ will be the main way of modelling questions about a given world:

“Given that I have observed some event β , what is the probability that some other (unobservable) event α of interest to me is true in the world?”

Examples:

- ▶ $\underline{P(\text{Hayfever} = \text{true} \mid \text{Congestion} = \text{true}) = ?}$
Given that my patient suffers from a sinus congestion, what is the probability that she has a hay fever?
- ▶ $\underline{P(\text{Hayfever} = \text{true} \mid \text{Congestion} = \text{true}, \text{MP} = \text{true}) = ?}$
If I learn, in addition, that she has muscle pain, how does that change my degree of belief in hay fever?
- ▶ $\underline{P(\text{Location}, \text{Velocity} \mid \text{GPS} = \langle ., ., . \rangle, \text{Radar} = \langle ., ., . \rangle) = ?}$
Given the readings of a robot's GPS and radar sensors, what is the probability distribution over the possible current locations and velocities of the robot?

Conditional Probabilities

Definition

For two events α, β , the **Conditional Probability** of α , given that we know that β is true, is

$$P(\alpha \mid \beta) = \frac{P(\alpha \cap \beta)}{P(\beta)}$$

Note:

- ▶ For a fixed β and some variables \mathbf{X} , the **conditional distribution** $P(\mathbf{X} \mid \beta)$ satisfies all the properties of a proper probability distribution (specifically: sums to 1.0 over all value combinations of \mathbf{X})
- ▶ $P(\mathbf{X} \mid \beta)$ is a legitimate probability distribution in its own right
- ▶ **Conditioning** is an operation that takes one distribution $P(\mathbf{X})$ and returns another distribution $P(\mathbf{X} \mid \beta)$.

The Chain Rule

Consequences of the above definition:

$$P(\alpha \mid \beta) = \frac{P(\alpha \cap \beta)}{P(\beta)}$$

implies:

$$P(\alpha \cap \beta) = P(\alpha \mid \beta)P(\beta)$$

and

$$P(\alpha \cap \beta) = P(\beta \cap \alpha) = P(\beta \mid \alpha)P(\alpha) = P(\alpha)P(\beta \mid \alpha)$$

Definition

The generalisation of the above to n events is known as the **Chain Rule**:

$$\begin{aligned} P(\alpha_1, \dots, \alpha_n) &= P(\alpha_1 \mid \alpha_2, \dots, \alpha_n)P(\alpha_2, \dots, \alpha_n) \\ &= P(\alpha_1 \mid \alpha_2, \dots, \alpha_n)P(\alpha_2 \mid \alpha_3, \dots, \alpha_n)P(\alpha_3, \dots, \alpha_n) \\ &\dots \\ &= P(\alpha_1 \mid \alpha_2, \dots, \alpha_n)P(\alpha_2 \mid \alpha_3, \dots, \alpha_n) \cdots P(\alpha_{n-1} \mid \alpha_n)P(\alpha_n) \end{aligned}$$

The Chain Rule

Chain Rule

$$P(\alpha_1, \dots, \alpha_n) = P(\alpha_1 \mid \alpha_2, \dots, \alpha_n)P(\alpha_2 \mid \alpha_3, \dots, \alpha_n) \cdots P(\alpha_{n-1} \mid \alpha_n)P(\alpha_n)$$

Notes:

- ▶ Because conjunction is commutative and associative, this expression can be expanded using **any order of events**, for example:

$$P(\alpha_1, \dots, \alpha_n) = P(\alpha_1)P(\alpha_2 \mid \alpha_1) \cdots P(\alpha_n \mid \alpha_1, \dots, \alpha_{n-1})$$

- ▶ ... and also **partially**, e.g.,

$$P(\alpha_1, \dots, \alpha_n) = P(\alpha_1, \dots, \alpha_k)P(\alpha_{k+1}, \dots, \alpha_n \mid \alpha_1, \dots, \alpha_k)$$

- ▶ The chain rule can also be written for whole **distributions**, e.g.:

$$P(X_1, X_2, \dots, X_N) = P(X_1, \dots, X_k)P(X_{k+1}, \dots, X_N \mid X_1, \dots, X_k)$$

The Law of Total Probability

Definition

The **Law of Total Probability**:

In a world defined over random variables \mathcal{X} , for any event $(\mathbf{X} = \mathbf{x})$, $\mathbf{X} \subset \mathcal{X}$, and any subset of variables $\mathbf{Y} \subset \mathcal{X}$, the following holds:

$$P(\mathbf{x}) = \sum_{\mathbf{y} \in \text{Val}(\mathbf{Y})} P(\mathbf{x}, \mathbf{y})$$

or, alternatively, following from the Chain Rule:

$$P(\mathbf{x}) = \sum_{\mathbf{y} \in \text{Val}(\mathbf{Y})} P(\mathbf{x} \mid \mathbf{y})P(\mathbf{y}) = \sum_{\mathbf{y} \in \text{Val}(\mathbf{Y})} P(\mathbf{x})P(\mathbf{y} \mid \mathbf{x})$$

Proof:

► Obvious. \square

The Law of Total Probability

Generalises to:

► **Marginal distributions:**

$$P(X) = \sum_y P(X, y) = \sum_y P(X | y)P(y)$$

► **Conditional distributions:**

$$P(X | Z) = \sum_y P(X, y | Z) = \sum_y P(X | y, Z)P(y | Z)$$



The Law of Total Probability justifies our marginalisation and “Inference by Enumeration” algorithms.

Bayes' Rule

Another view on conditional probability:

$$P(\alpha \mid \beta) = \frac{P(\alpha \cap \beta)}{P(\beta)} \quad \text{and} \quad P(\alpha \cap \beta) = P(\beta \mid \alpha)P(\alpha)$$

imply:

Definition

BAYES' RULE:

$$P(\alpha \mid \beta) = \frac{P(\beta \mid \alpha)P(\alpha)}{P(\beta)}$$

Notes:

- ▶ Bayes' rule allows us to compute a conditional probability $P(\alpha \mid \beta)$ from the “inverse” conditional probability $P(\beta \mid \alpha)$
- ▶ Bayes' rule is one of the most fundamental tools in scientific data analysis and probabilistic reasoning.

The Reverend Thomas Bayes



Quote from the Encyclopedia Britannica (quote and image taken from the *International Society for Bayesian Analysis* [<http://bayesian.org/bayes>]):

Bayes, Thomas (b. 1702, London - d. 1761, Tunbridge Wells, Kent), mathematician who first used probability inductively and established a mathematical basis for probability inference (a means of calculating, from the number of times an event has not occurred, the probability that it will occur in future trials).

He set down his findings on probability in *“Essay Towards Solving a Problem in the Doctrine of Chances”* (1763), published posthumously in the Philosophical Transactions of the Royal Society of London.

The only works he is known to have published in his lifetime are *“Divine Benevolence, or an Attempt to Prove That the Principal End of the Divine Providence and Government is the Happiness of His Creatures”* (1731) and *“An Introduction to the Doctrine of Fluxions, and a Defence of the Mathematicians Against the Objections of the Author of the Analyst”* (1736) which countered attacks by Bishop Berkeley on the logical foundations of Newton’s calculus.

The Reverend Thomas Bayes



Quote and image taken from the *International Society for Bayesian Analysis* [<http://bayesian.org/bayes>]):

Bayes is buried in Bunhill Fields in the heart of the City of London. The cemetery was used for the burial of nonconformists in the 18th century, but is now a public park maintained by the Corporation of London. Also buried in Bunhill Fields is Bayes's friend Richard Price, a pioneer of insurance, who presented Bayes's famous paper on probability to the Royal Society in 1763, two years after Bayes's death

A Simple Example

Assume we know (or believe, or estimate) from experience that

- ▶ $\underline{P(mp \mid flu) = 0.5}$ (50% of the flu patients tend to have a muscle pain)
- ▶ $\underline{P(flu) = 0.004}$ (0.4% of the population at any time tend to have the flu)
Or: the probability of a random person having the flu (if we know nothing else about the person's condition) is 0.004
- ▶ $\underline{P(mp) = 0.1}$ (10% of the population tend to suffer from a muscle pain).

Question:

- ▶ Given that you meet a person who suffers from a muscle pain, what is the probability that the person has the flu?

By Bayes' rule:

$$P(flu \mid mp) = \frac{P(mp \mid flu)P(flu)}{P(mp)} = (0.5 \times 0.004)/0.1 = \underline{0.02}$$

Bayes' Rule in Diagnostic Settings

Know about:

- ▶ Possible observations (e.g., symptoms, lab tests, ...)
- ▶ Problems or diseases that would explain the symptoms
- ▶ General probabilities $P(Problem)$ of diseases/problems to occur
- ▶ Probabilities with which problems produce specific symptoms:
 $P(Symptoms \mid Problem)$.

Observe:

- ▶ A patient/situation with specific *Symptoms*

Want to determine:

- ▶ Probability of the underlying problem, given the symptoms:

$$P(Problem \mid Symptoms) = \frac{P(Symptoms \mid Problem)P(Problem)}{P(Symptoms)}$$

👉 **Diagnosis is a very common reasoning pattern (not only in medicine)!**

Bayes' Rule in Diagnostic Settings

Bayes' Rule in Diagnosis Settings:

$$P(\textit{Problem} \mid \textit{Symptoms}) = \frac{P(\textit{Symptoms} \mid \textit{Problem})P(\textit{Problem})}{P(\textit{Symptoms})}$$

Names of the components:

$P(\textit{Problem})$

Prior Probability (degree of belief in *Problem* **before** we have observed anything else about the current case)

$P(\textit{Problem} \mid \textit{Symptoms})$

Posterior Probability
(degree of belief in *Problem* **after** we have observed *Symptoms*)

$P(\textit{Symptoms} \mid \textit{Problem})$

Likelihood of *Problem* in the presence of *Symptoms*
= probability with which *Problem* produces *Symptoms*

$P(\textit{Symptoms})$

“Evidence” (probability of these *Symptoms* occurring at all).

Note: By the Law of Total Probability, $P(\textit{Symptoms})$ can be calculated as

$$P(\textit{Symptoms}) = \sum_{p \in \textit{Problems}} P(\textit{Symptoms} \mid p)P(p)$$

Probabilistic Inference and Reasoning

Scenario:

- ▶ The **full joint distribution** over some set of variables \mathcal{X} will be the system's **Knowledge Base (KB)**, its **model of the world**.
- ▶ Gives a probability to each possible state of the world (= atomic event): how likely is it to encounter this specific combination of values?
- ▶ The system's task will be to compute, from this knowledge base and the known (observed) values of some variables, the probability of any hypothesis of interest.

Definition

The process of deriving the truth or probability of a hypothesis from a knowledge base and from specific observations is called **INFERENCE**, or **REASONING**.



In our context, inference will mainly mean *computing (conditional) probability distributions over variables of interest*.

Probabilistic Inference and Reasoning

$P(Flu, MP, Fever)$	$musclepain$		$\neg musclepain$	
	$fever$	$\neg fever$	$fever$	$\neg fever$
flu	0.15	0.1	0.1	0.03
$\neg flu$	0.03	0.04	0.05	0.5

Possible queries we might like to answer:

► $P(fever \wedge \neg musclepain) = ?$

What is the probability that someone has a fever, but no muscle pain?

(☞ an *event probability*)

► $P(Flu) = ?$

What is the probability that a random person has the flu (or not)?

(☞ the *marginal distribution* over variable Flu)

► $P(Flu \mid fever, \neg musclepain) = ?$

What is the probability that someone does or does not have the flu, if they have a fever, but no muscle pain?

(☞ a *conditional distribution*)

► and many more ...

Inference by Enumeration

$P(Flu, MP, Fever)$	$musclepain$		$\neg musclepain$	
	$fever$	$\neg fever$	$fever$	$\neg fever$
flu	0.15	0.1	0.1	0.03
$\neg flu$	0.03	0.04	0.05	0.5

Remember how we computed marginal probabilities from the full joint:

$$\begin{aligned} P(flu) &= P(flu, mp, fever) + P(flu, mp, \neg fever) + \\ &\quad P(flu, \neg mp, fever) + P(flu, \neg mp, \neg fever) \\ &= \sum_{x,y} P(flu, MP = x, Fever = y) \\ &= 0.15 + 0.1 + 0.1 + 0.03 = 0.38 \end{aligned}$$

- ▶ Can compute probability of any event α as the sum over all atomic events that ‘match’ α
- ▶ This is a direct application of the *Law of Total Probability*
- ▶ Will permit us to compute the probability of any query we may want to ask of a distribution.

 **Very general algorithm!**

$P(Flu, MP, Fever)$	$musclepain$		$\neg musclepain$	
	$fever$	$\neg fever$	$fever$	$\neg fever$
flu	0.15	0.1	0.1	0.03
$\neg flu$	0.03	0.04	0.05	0.5

$$\begin{aligned} P(\neg fever) &= \sum_{x,y} P(\neg fever, MP = x, Flu = y) \\ &= 0.1 + 0.03 + 0.04 + 0.5 = 0.67 \end{aligned}$$

$$\begin{aligned} P(fever \wedge \neg musclepain) &= \sum_{x \in Val(Flu)} P(fever, \neg musclepain, Flu = x) \\ &= 0.1 + 0.05 = 0.15 \end{aligned}$$

$P(Flu, MP, Fever)$	$musclepain$		$\neg musclepain$	
	$fever$	$\neg fever$	$fever$	$\neg fever$
flu	0.15	0.1	0.1	0.03
$\neg flu$	0.03	0.04	0.05	0.5

$$\begin{aligned}
 P(flu \mid \neg fever) &= \frac{P(flu \wedge \neg fever)}{P(\neg fever)} && \text{(by definition of cond. prob.)} \\
 &= \frac{\sum_x P(flu, \neg fever, MP = x)}{\sum_{x,y} P(Flu = x, \neg fever, MP = y)} \\
 &= \frac{0.1 + 0.03}{0.1 + 0.03 + 0.04 + 0.5} \\
 &= \frac{0.13}{0.67} \approx 0.194
 \end{aligned}$$

Example 3

$P(Flu, MP, Fever)$	$musclepain$		$\neg musclepain$	
	$fever$	$\neg fever$	$fever$	$\neg fever$
flu	0.15	0.1	0.1	0.03
$\neg flu$	0.03	0.04	0.05	0.5

Computing Conditional Distributions:

$$\begin{aligned} P(flu \mid \neg fever) &= \frac{P(flu \wedge \neg fever)}{P(\neg fever)} \\ &= \frac{0.1 + 0.03}{0.1 + 0.03 + 0.04 + 0.5} = \frac{0.13}{0.67} \approx 0.194 \end{aligned}$$

$$\begin{aligned} P(\neg flu \mid \neg fever) &= \frac{P(\neg flu \wedge \neg fever)}{P(\neg fever)} \\ &= \frac{0.04 + 0.5}{0.1 + 0.03 + 0.04 + 0.5} = \frac{0.54}{0.67} \approx 0.806 \end{aligned}$$

$$\Rightarrow P(Flu \mid \neg fever) = \{0.806, 0.194\}$$

Example 3 (revisited)

$$P(flu \mid \neg fever) = \frac{P(flu \wedge \neg fever)}{P(\neg fever)} = \frac{0.13}{0.67} \approx 0.194$$

$$P(\neg flu \mid \neg fever) = \frac{P(\neg flu \wedge \neg fever)}{P(\neg fever)} = \frac{0.54}{0.67} \approx 0.806$$

$$\Rightarrow P(Flu \mid \neg fever) = \{0.194, 0.806\}$$

Notes:

- ▶ The denominator $P(\neg fever) = 0.67$ is the same in both cases (because $P(\neg fever)$ does not depend on the value of Flu)
- ▶ In general: In computing a distribution $P(\mathbf{X} \mid \mathbf{y})$, the term $P(\mathbf{y})$ will be computed $|Val(\mathbf{X})|$ times!
- ▶ This constant term $P(\mathbf{y})$ serves as a **normalisation factor** Z that makes the resulting conditional distribution $P(\mathbf{X} \mid \mathbf{y})$ sum to 1.0
- ▶ $Z = P(\mathbf{y}) = \sum_{\mathbf{x}} P(\mathbf{x}, \mathbf{y})$ can be computed afterwards, in a simple way (without having to sum out over the full joint distribution) ...

- Compute $P(\mathbf{X} \mid \mathbf{y})$
for some sets of variables $\mathbf{X}, \mathbf{Y} \subseteq \mathcal{X}$ and some specific value assignment $\mathbf{Y} = \mathbf{y}$.

$$P(\mathbf{X} \mid \mathbf{y}) = \frac{P(\mathbf{X}, \mathbf{y})}{P(\mathbf{y})} \quad \text{(definition of conditional probability)}$$

- 1 For each possible value combination $x \in \text{Val}(\mathbf{X})$, compute $P(x, \mathbf{y})$ via enumeration
- 2 Compute denominator = normalisation constant $P(\mathbf{y})$ from these $P(x, \mathbf{y})$ (simply by summing them up (Law of Total Probability))
- 3 Normalise $P(\mathbf{X}, \mathbf{y})$ to obtain conditional distribution $P(\mathbf{X} \mid \mathbf{y})$.

The General Inference-by-Enumeration Algorithm

Algorithm

- 1 Compute $P(\mathbf{X}, y)$:

$$P(\mathbf{X}, \mathbf{y}) = \sum_{\mathbf{z} \in \text{Val}(\mathbf{Z})} P(\mathbf{X}, \mathbf{y}, \mathbf{z})$$

(Inference by Enumeration, $\mathbf{Z} = \mathcal{X} - \mathbf{X} - \mathbf{Y}$)

- 2 Compute normalisation constant Z as

$$Z = P(\mathbf{y}) = \sum_{\mathbf{x}} P(\mathbf{x}, \mathbf{y})$$

(simply sum up all the $P(x, y)$ computed in previous step)

- ### 3 Conditioning via renormalisation:

$$P(\mathbf{X} \mid \mathbf{y}) = \frac{1}{Z} P(\mathbf{X}, \mathbf{y})$$

$P(Flu, MP, Fever)$	$musclepain$		$\neg musclepain$	
	$fever$	$\neg fever$	$fever$	$\neg fever$
flu	0.15	0.1	0.1	0.03
$\neg flu$	0.03	0.04	0.05	0.5

$$\begin{aligned} \boxed{P(Flu, \neg fever)} &= \sum_{x \in \{mp, \neg mp\}} P(Flu, \neg fever, MP = x) \\ &= \begin{bmatrix} 0.1 + 0.03 \\ 0.04 + 0.5 \end{bmatrix} = \begin{bmatrix} 0.13 \\ 0.54 \end{bmatrix} \quad \begin{array}{l} \leftarrow P(flu, \neg fever) \\ \leftarrow P(\neg flu, \neg fever) \end{array} \end{aligned}$$

$$\boxed{Z} = \sum_{x \in \{flu, \neg flu\}} P(x, \neg fever) = 0.13 + 0.54 = 0.67$$

$$\boxed{P(Flu \mid \neg fever)} = \frac{1}{Z} \times \begin{bmatrix} 0.13 \\ 0.54 \end{bmatrix} = \begin{bmatrix} \mathbf{0.194} \\ \mathbf{0.806} \end{bmatrix} \quad \begin{array}{l} \leftarrow P(flu \mid \neg fever) \\ \leftarrow P(\neg flu \mid \neg fever) \end{array}$$

Independence Between Events

Definition

Two events α and β are said to be **independent** in a joint distribution P , denoted as

$$(\alpha \perp \beta)$$

if^a

$$P(\alpha \mid \beta) = P(\alpha) \quad \text{or, equivalently,}$$

$$P(\beta \mid \alpha) = P(\beta) \quad \text{or, equivalently,}$$

$$P(\alpha \cap \beta) = P(\alpha)P(\beta)$$

^aExercise: Prove that these three conditions are equivalent!

Note that independence is **symmetric**:

$$(\alpha \perp \beta) \Leftrightarrow (\beta \perp \alpha)$$

Independence Between Random Variables

Definition

Two **random variables** X and Y are **independent**, denoted as

$$(X \perp Y)$$

if for all values $x \in Val(X), y \in Val(Y)$:

$$P(x \mid y) = P(x) \quad \text{and} \quad P(y \mid x) = P(y)$$

or, in short:

$$P(X \mid Y) = P(X) \quad \text{and} \quad P(Y \mid X) = P(Y)$$

Consequence:

- ▶ If X and Y are independent, then

$$P(X, Y) = P(X)P(Y)$$

Independence Between Random Variables

Generalises to *Sets* of Random Variables:

Definition

Two **sets of random variables** X and Y are **independent**, denoted as

$$(\mathbf{X} \perp \mathbf{Y})$$

if for all value combinations $x \in Val(\mathbf{X}), y \in Val(\mathbf{Y})$:

$$P(\boldsymbol{x} \mid \boldsymbol{y}) = P(\boldsymbol{x}) \quad \text{and} \quad P(\boldsymbol{y} \mid \boldsymbol{x}) = P(\boldsymbol{y})$$

or, in short:

$$P(\mathbf{X} \mid \mathbf{Y}) = P(\mathbf{X}) \quad \text{and} \quad P(\mathbf{Y} \mid \mathbf{X}) = P(\mathbf{Y})$$

Consequence:

- If X and Y are independent, then

$$P(\mathbf{X}, \mathbf{Y}) = P(\mathbf{X})P(\mathbf{Y})$$

The Importance of Independence

If X and Y are independent, then $P(X, Y) = P(X)P(Y)$

Consequences:

- ▶ If two variable sets \mathbf{X}, \mathbf{Y} are independent, their joint distribution $P(\mathbf{X}, \mathbf{Y})$ is equal to the product $P(\mathbf{X})P(\mathbf{Y})$
- ▶ Each entry $P(x, y)$ in the joint distribution $P(\mathbf{X}, \mathbf{Y})$ can be reconstructed from the lower-dimensional distributions $P(\mathbf{X})$ and $P(\mathbf{Y})$:

$$P(\boldsymbol{x}, \boldsymbol{y}) = P(\boldsymbol{x})P(\boldsymbol{y})$$

- ▶ Need to model (store) joint distributions only over smaller subsets X and Y rather than over full set $X \cup Y$
- ▶ Remember: Number of entries in joint distribution grows **exponentially** with number of variables!

👉 **Effect: Reduction of Space Complexity.**

An Extreme Example

Consider N independent coins (binary variables $\{X_1, \dots, X_N\}$)



- Can assume that all X_i, X_j are pairwise independent:

$$P(X_1, \dots, X_N) = P(X_1)P(X_2) \cdots P(X_N)$$

- ▶ Can reconstruct joint distribution over $\{X_1, \dots, X_N\}$ from one-dimensional distributions $P(X_1), P(X_2) \dots$
- ▶ Joint distribution $P(X_1, \dots, X_N)$ over all variables has 2^N entries
- ▶ Each marginal distribution $P(X_i)$ has only 2 entries $P(heads), P(tails)^1$
- ▶ Need to store only $2N$ probabilities instead of 2^N

 **Exponential Reduction of Complexity from $O(2^N)$ to $O(N)$!**

¹Actually, only one is needed: $P(tails) = 1 - P(heads)$

👉 **α and β are dependent:** Learning that Smith feels a muscle pain increases our expectation that she will also run a fever (because these often go together):

$$P(\alpha \mid \beta) > P(\alpha)$$

... then learning about β will probably not change our degree of belief in α (because α is already “explained” by γ):

$$P(\alpha \mid \beta \cap \gamma) = P(\alpha \mid \gamma)$$

 β tells us nothing new about α , if we already know γ ...

Conditional Independence

Definition

Two events α and β are said to be **conditionally independent** given a third event γ , denoted as

$$(\alpha \perp \beta \mid \gamma)$$

 if^a

$$P(\alpha \mid \beta \cap \gamma) = P(\alpha \mid \gamma) \quad \text{or, equivalently,}$$

$$P(\beta \mid \alpha \cap \gamma) = P(\beta \mid \gamma) \quad \text{or, equivalently,}$$

$$P(\alpha \cap \beta \mid \gamma) = P(\alpha \mid \gamma)P(\beta \mid \gamma)$$

^aExercise: Prove that these three conditions are equivalent!

Note that conditional independence is **symmetric**:

$$(\alpha \perp \beta \mid \gamma) \Leftrightarrow (\beta \perp \alpha \mid \gamma)$$

100%

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

$$P(\mathbf{x} \mid \mathbf{y}, \mathbf{z}) = P(\mathbf{x} \mid \mathbf{z}) \quad \text{and} \quad P(\mathbf{y} \mid \mathbf{x}, \mathbf{z}) = P(\mathbf{y} \mid \mathbf{z})$$

$$P(\mathbf{Y} \mid \mathbf{V}, \mathbf{Z}) = P(\mathbf{Y} \mid \mathbf{Z}) \quad \text{and} \quad P(\mathbf{V} \mid \mathbf{Y}, \mathbf{Z}) = P(\mathbf{V} \mid \mathbf{Z})$$

$$P(\mathbf{Y} \ \mathbf{V} \mid \mathbf{Z}) = P(\mathbf{Y} \mid \mathbf{Z}) P(\mathbf{V} \mid \mathbf{Z})$$

The Importance of Conditional Independence

If X and Y are conditionally independent given Z , then

$$P(\mathbf{X} \mid \mathbf{Y}, \mathbf{Z}) = P(\mathbf{X} \mid \mathbf{Z})$$

Consequence:

- ▶ Joint distribution $P(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$ over all variables can be decomposed (“factorised”) into product of simpler (conditional) distributions:

$$P(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) = P(\mathbf{X}|\mathbf{Y}, \mathbf{Z})P(\mathbf{Y}, \mathbf{Z}) \quad (1)$$

$$= P(\mathbf{X}|\mathbf{Y}, \mathbf{Z})P(\mathbf{Y}|\mathbf{Z})P(\mathbf{Z}) \quad (2)$$

$$= P(\mathbf{X}|\mathbf{Z}) \cdot P(\mathbf{Y}|\mathbf{Z}) \cdot P(\mathbf{Z}) \quad (3)$$

(1) *chain rule*

(2) *chain rule*

(3) *conditional independence*

👉 Most fundamental and important concept in Graphical Models

 **Be sure to understand this!**

Test Yourself²

$P(X_1, X_2)$	$X_2 = 1$	$X_2 = 2$
$X_1 = 1$	0.3	0.2
$X_1 = 2$	0.1	0.4

1. Find the following quantities:

- ▶ Marginals: $P(X_1), P(X_2)$
- ▶ Conditionals: $P(X_1 \mid X_2), P(X_2 \mid X_1)$
- ▶ Joint vs. Posterior: $P(X_1, X_2 = 2)$ vs. $P(X_1 \mid X_2 = 2)$
- ▶ Evidence (relative to previous item): $P(X_2 = 2)$
- ▶ Normalisation constant $P(\Omega) = P(\text{true})$
- ▶ Max: $P(x_1^*) = \max_{x_1} P(x_1 \mid X_2 = 1)$
- ▶ Mode: $x_1^* = \arg \max_{x_1} P(x_1 \mid X_2 = 1)$
- ▶ Max-marginal: $\max_{x_1} P(x_1, X_2)$

2. Are X_1 and X_2 independent?

²Test motivated by Ali Cemgil's ISMIR 2006 Tutorial (www-sigproc.eng.cam.ac.uk/~atc27/).

Answers

$P(X_1, X_2)$	$X_2 = 1$	$X_2 = 2$
$X_1 = 1$	0.3	0.2
$X_1 = 2$	0.1	0.4

Marginals:

$P(X_1)$	
$X_1 = 1$	0.5
$X_1 = 2$	0.5

$P(X_2)$	$X_2 = 1$	$X_2 = 2$
	0.4	0.6

Conditionals:

$P(X_1 \mid X_2)$	$X_2 = 1$	$X_2 = 2$
$X_1 = 1$	0.75	0.33
$X_1 = 2$	0.25	0.67

$P(X_2 \mid X_1)$	$X_2 = 1$	$X_2 = 2$
$X_1 = 1$	0.6	0.4
$X_1 = 2$	0.2	0.8

Answers

$P(X_1, X_2)$	$X_2 = 1$	$X_2 = 2$
$X_1 = 1$	0.3	0.2
$X_1 = 2$	0.1	0.4

Joint vs. Posterior:

$P(X_1, X_2 = 2)$	$X_2 = 2$
$X_1 = 1$	0.2
$X_1 = 2$	0.4

$P(X_1 \mid X_2 = 2)$	$X_2 = 2$
$X_1 = 1$	0.33
$X_1 = 2$	0.67

Evidence: (relative to previous item):

$$P(X_2 = 2) = \sum_{x_1} P(x_1, X_2 = 2) = 0.6$$

Normalisation constant:

$$P(\Omega) = P(true) = \sum_{x_1} \sum_{x_2} P(x_1, x_2) = 1$$

Answers

$P(X_1, X_2)$	$X_2 = 1$	$X_2 = 2$
$X_1 = 1$	0.3	0.2
$X_1 = 2$	0.1	0.4

Max:

$$P(x_1^*) = \max_{x_1} P(x_1 \mid X_2 = 1) = 0.75$$

Mode:

$$x_1^* = \arg \max_{x_1} P(x_1 \mid X_2 = 1) = 1$$

Max-marginal (get the “skyline”):

$\max_{x_1} P(x_1, X_2)$	$X_2 = 1$	$X_2 = 2$
	0.3	0.4

$P(X_1, X_2)$	$X_2 = 1$	$X_2 = 2$
$X_1 = 1$	0.3	0.2
$X_1 = 2$	0.1	0.4

2. Are X_1 and X_2 independent?

Remember: Two random variables X, Y are (marginally) independent if

$$P(X \mid Y) = P(X) \quad (\text{or} \quad P(Y \mid X) = P(Y))$$

👉 easy to see that X_1 and X_2 are

- ▶ So far, have considered only *discrete* variables with *finite* domains $Val(X)$
- ▶ But: In real-world applications, many (most) of the relevant variables will be **continuous** (real-valued) !
- ▶ Examples: Sensors in an autonomous robot – range sensors, temperature and pressure sensors, GPS, camera, ...

- ▶ Probability of each specific value (from an infinite set) is zero!

- ▶ Assume continuous variable X with $Val(X) = [0, 1] \subset \mathbb{R}$
 - ▶ Assume we want to assign the same probability to each number in this range (uniform distribution)
- ⇒ infinite number of possible values x
- ⇒ $P(X = x)$ must be 0 for all x (otherwise $\sum_x P(x)$ would be infinite)!

(This problem also appears for other, non-uniform distributions)

Probabilities from Densities

The PDF defines a set of **distributions** over X as follows:

Definition

For any $a \in \mathbb{R}$,

$$P(X \leq a) = \int_{-\infty}^a p(x)dx$$

This is called the **Cumulative Distribution Function (CDF)** for X .

Intuitively: PDF $p(x)$ is the incremental amount that x adds to the cumulative distribution in the integration process.

Consequence (follows from rules of probability)

$$P(a \leq X \leq b) = P(X \leq b) - P(X \leq a) = \int_a^b p(x)dx$$

👉 PDFs permit us to compute probabilities for (arbitrarily small) *intervals* of values.

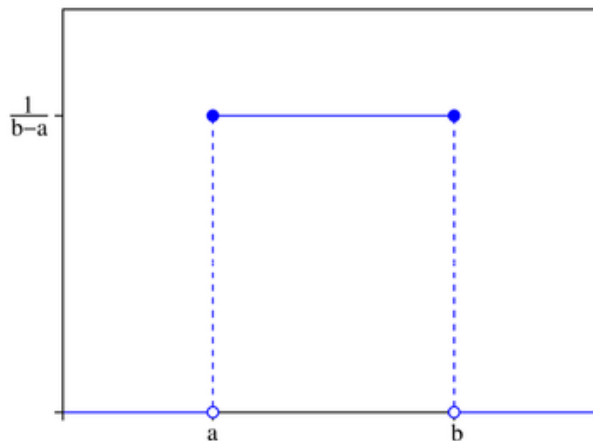
Example 1: The Uniform Distribution

Definition

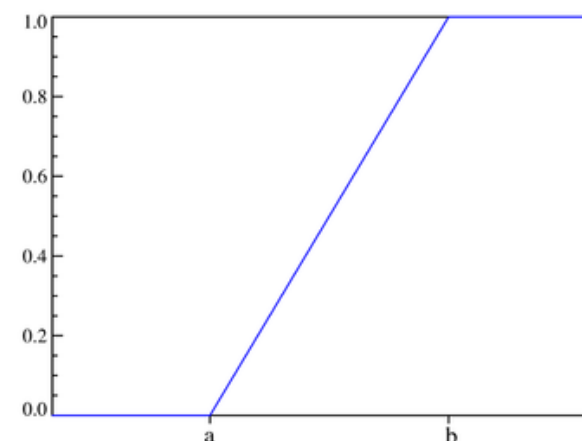
A real variable X has a **Uniform Distribution** over a range $[a, b]$, denoted $X \sim Unif[a, b]$, if it has the PDF

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

PDF:



CDF:



Note:

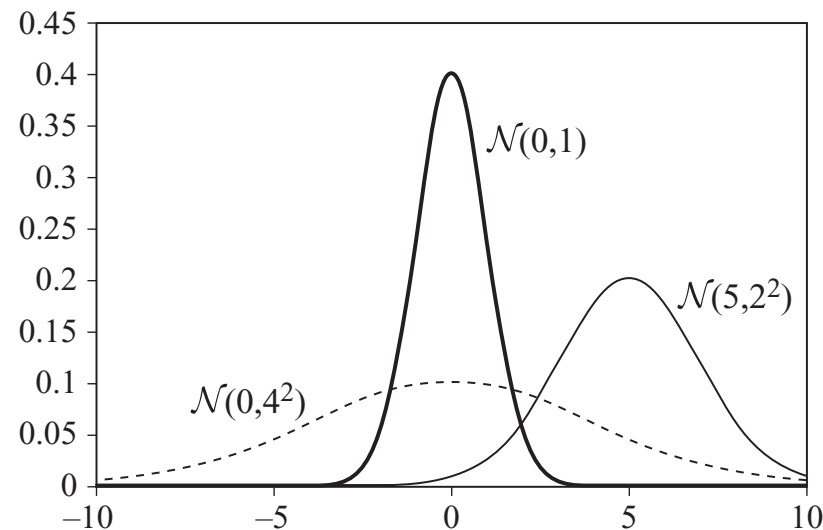
- ▶ If $(b - a) < 1$, the PDF $p(x)$ is > 1 !
- ▶ Can happen in any legal PDF, as long as total area under the PDF is 1.0.

Example 2a: The Normal (Gaussian) Distribution

Definition

A variable X has a **Normal Distribution** with mean μ and variance σ^2 , denoted $X \sim \mathcal{N}(\mu, \sigma^2)$, if it has the **Gaussian PDF**

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



Example 2b: The Multivariate Gaussian

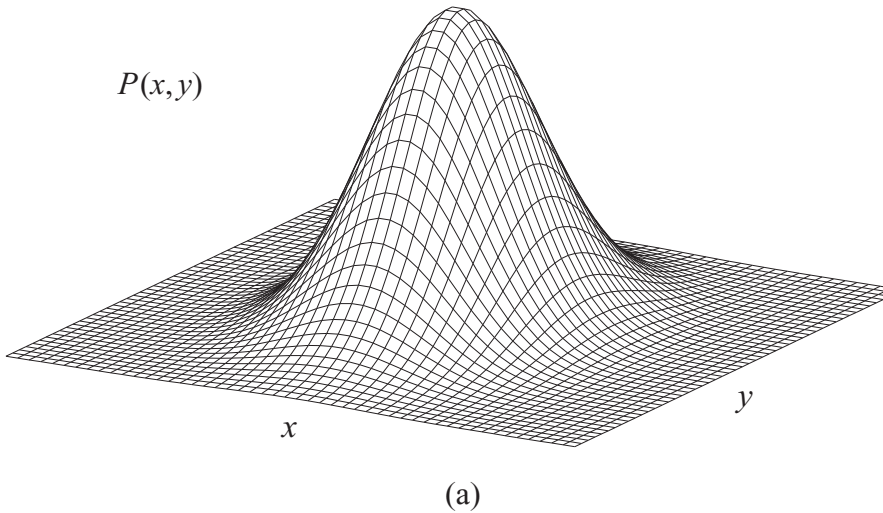
Definition

N random variables $\mathbf{X} = \{X_1, \dots, X_N\}$ have a **Joint (Multivariate) Normal Distribution** with mean vector $\boldsymbol{\mu}$ and covariance Σ , denoted $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$, if they have the **Multivariate Gaussian PDF**

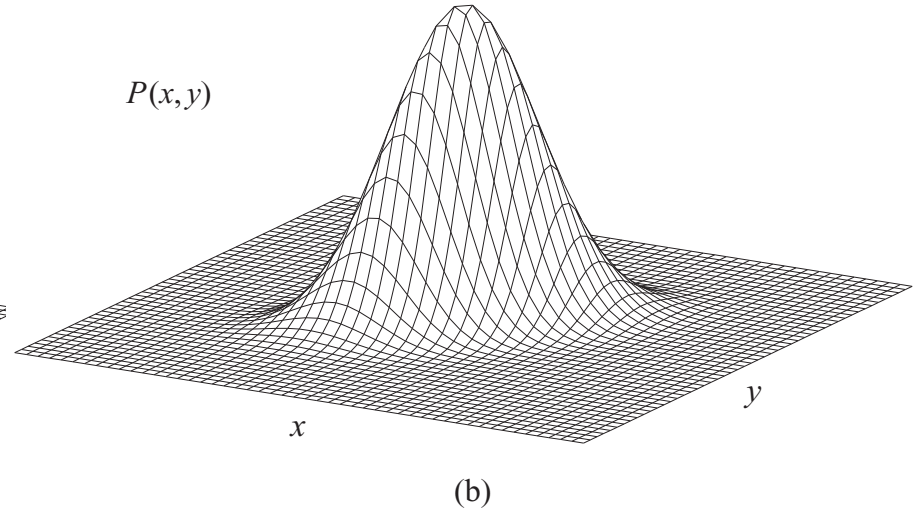
$$p(\mathbf{x}) = \frac{1}{(2\pi)^{N/2} |\mathbf{\Sigma}|^{1/2}} e^{[-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \mathbf{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})]}$$

where

- ▶ $\boldsymbol{\mu}$ is a *mean vector* $\boldsymbol{\mu} = [\mu_{X_1}, \dots, \mu_{X_N}]$
- ▶ $\boldsymbol{\Sigma}$ is the $N \times N$ *covariance matrix* (with $\Sigma_{ij} = \sigma_{ij}$ the *covariance* (= *pairwise correlation*) between variables X_i and X_j)
- ▶ $|\boldsymbol{\Sigma}|$ is the determinant of $\boldsymbol{\Sigma}$



X and Y are uncorrelated ($\sigma_{X,Y} = 0$)



X and Y are (positively) correlated
($\sigma_{X,Y} > 0$)

Notational Conventions: Summary

NOTATION

Examples	Type of object
\mathcal{X}	Complete set of variables that a model is defined over
A, X_i, E_j	Single random variable (e.g., $X_i \in \mathcal{X}$)
$\mathbf{X}, \mathbf{Y}, \mathbf{Z}, \mathbf{E}, \mathbf{U}, \dots$	Sets of random variables (e.g., $\mathbf{X} = \{X_1, X_2, X_3, X_4\}$)
$Val(X)$	Domain of variable X (set of values that X can take)
x^0, x^1, \dots	$x^i = i^{th}$ value in the domain $Val(X)$ of variable X
x^0, x^1	Special convention for boolean variables: $x^0 = false, x^1 = true$
$\mathbf{x}_i, \mathbf{e}, \dots$	Set/list of specific values of a variable (e.g., $\mathbf{E} = e$)
$Val(\mathbf{X})$	Set of possible value assignments to the variables in \mathbf{X} (= Cartesian product $Val(X_1) \times Val(X_2) \times \dots$)
$\mathbf{x}\langle\mathbf{Y}\rangle$	Instance \mathbf{x} reduced to its values for variables $\mathbf{Y} \subset \mathbf{X}$
$X_i^{(0)}, X_i^{(1)}, X_i^{(t)}$	Instantiation of a template variable X_i at time t (see Chapter 6)
θ	Set of parameters of a model
$P(X = x^i)$	Probability (probability mass function)
$P(x^i), P(x^i, y^j)$	Shorthand notation for $P(X = x^i), P(X = x^i, Y = y^j)$
$P(X, Y)$	Probability distribution over discrete variables X, Y
$p(x)$	Probability Density Function (PDF) over continuous variable X

