



PROBABILITY THEORY REVIEW

CMPUT 466/566

Martha White



Fall, 2019

REMINDERS

- Assignment 1 is due on September 26
 - 566 has to do extra question (bonus for 466)
- Thought questions 1 are due on September 19
 - Preface and Chapters 1-4, about 40 pages
 - If you are printing, don't print all the notes yet
- Start thinking about datasets for your mini-project
- I do not expect you to know formulas, like pdfs
- I will use a combination of slides and writing on the board (where you should take notes)
- Any questions?

BACKGROUND FOR COURSE

- Need to know calculus, mostly derivatives
 - Will almost never integrate
 - I will teach you multivariate calculus
- Need to know linear algebra
 - I assume you know about vectors, matrices and dot products
 - I will teach you more advanced topics, like singular value decompositions
- Need to have learned about probability

WHY DO WE NEED PROBABILITIES?

- We could just assume a deterministic world
- I see an input x , I can produce the output $y = f(x)$
 - Example: in a game (e.g., chess), you take an action, and the outcome is deterministic
- But, even in a deterministic world, we have a problem: partial observability
- Outcomes look random because we don't have enough information
 - Example: Imagine a (high-tech) gumball machine, where $f(x = \text{has candy, battery charged}) = \text{output candy}$
 - You can only see if it has candy

WHY DO YOU NEED TO KNOW PROBABILITIES?

- Once we derive the algorithms, you may not have needed to know about probabilities
 - isn't it all just linear algebra?
- BUT, actually, many models in machine learning use distributions explicitly, including
 - Graphical models
 - Boltzmann machine
 - Bayesian strategies, like Bayesian linear regression
 - Variational auto-encoders
 - ...
- Having a good grasp of probabilities (formally) really makes understanding these models easier

STRUCTURE FOR COURSE

- In Chapter 4, we will see that optimal regression and classification models model distribution $p(y | x)$
- In Chapter 3, we talk about how to estimate distributions, like $p(y)$ and $p(y | x)$ (maximum likelihood)
- Alternative order: motivate why we care about $p(y|x)$ first and then talk about how to estimate $p(y|x)$
- Problem: you are not yet sufficiently comfortable with distributions to jump right to Chapter 4
- Instead: we learn about MLE first, both because (a) we will need it later and (b) as practice to get comfortable with probability

SPACE OF OUTCOMES AND EVENTS

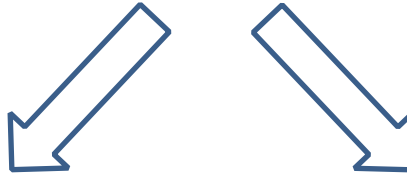
Ω = sample space, all outcomes of the experiment

\mathcal{E} = event space, set of subsets of Ω

Ω and \mathcal{E} must be non-empty

SAMPLE SPACES

Ω



discrete (countable)

continuous (uncountable)

$$\Omega = \{1, 2, 3, 4, 5, 6\}$$

$$\Omega = [0, 1]$$

$$\Omega = \mathbb{N}$$

$$\Omega = \mathbb{R}$$

e.g., $\mathcal{E} = \{\emptyset, \{1, 2\}, \{3, 4, 5, 6\}, \{1, 2, 3, 4, 5, 6\}\}$

e.g., $\mathcal{E} = \{\emptyset, [0, 0.5], (0.5, 1.0], [0, 1]\}$

$$\Omega = [0, 1] \cup \{2\} = \text{mixed space}$$

(MEASURABLE) SPACE OF OUTCOMES AND EVENTS

Ω = sample space, all outcomes of the experiment

\mathcal{E} = event space, set of subsets of Ω

Ω and \mathcal{E} must be non-empty

If the following conditions hold:

$$1. A \in \mathcal{E} \Rightarrow A^c \in \mathcal{E}$$

$$2. A_1, A_2, \dots \in \mathcal{E} \Rightarrow \bigcup_{i=1}^{\infty} A_i \in \mathcal{E}$$

\mathcal{E} is an event space

Note: terminology sigma field sounds technical, but it just means this event space

(Ω, \mathcal{E}) = a measurable space

WHY IS THIS THE DEFINITION?

Intuitively,

1. A collection of outcomes is an event (e.g., either a 1 or 6 was rolled)
2. If we can measure two events separately, then their union should also be a measurable event
3. If we can measure an event, then we should be able to measure that that event did not occur (the complement)

Ω = sample space, all outcomes of the experiment

\mathcal{E} = event space, set of subsets of Ω

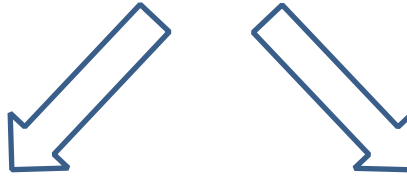
If the following conditions hold:

$$1. A \in \mathcal{E} \Rightarrow A^c \in \mathcal{E}$$

$$2. A_1, A_2, \dots \in \mathcal{E} \Rightarrow \bigcup_{i=1}^{\infty} A_i \in \mathcal{E}$$

SAMPLE SPACES

Ω



discrete (countable)

continuous (uncountable)

$$\Omega = \{1, 2, 3, 4, 5, 6\}$$

$$\Omega = [0, 1]$$

$$\Omega = \mathbb{N}$$

$$\Omega = \mathbb{R}$$

e.g., $\mathcal{E} = \{\emptyset, \{1, 2\}, \{3, 4, 5, 6\}, \{1, 2, 3, 4, 5, 6\}\}$

e.g., $\mathcal{E} = \{\emptyset, [0, 0.5], (0.5, 1.0], [0, 1]\}$

Typically: $\mathcal{E} = \mathcal{P}(\Omega)$

Typically: $\mathcal{E} = \mathcal{B}(\Omega)$



Power set



Borel field

$$\Omega = [0, 1] \cup \{2\} = \text{mixed space}$$

A FEW COMMENTS ON TERMINOLOGY

- A few new terms, including countable, closure
 - only a small amount of terminology used, can google these terms and learn on your own
 - notation sheet in notes
- Countable: integers, $\{0.1, 2.0, 3.6\}, \dots$
- Uncountable: real numbers, intervals, ...
- Interchangeably I use (though its somewhat loose)
 - discrete and countable
 - continuous and uncountable

AXIOMS OF PROBABILITY

$(\Omega, \mathcal{E}) =$ a measurable space

Any function $P : \mathcal{E} \rightarrow [0, 1]$ such that

1. (unit measure) $P(\Omega) = 1$
2. (σ -additivity) Any countable sequence of disjoint events $A_1, A_2, \dots \in \mathcal{E}$ satisfies $P(\cup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} P(A_i)$

is called a probability measure (probability distribution)

$(\Omega, \mathcal{E}, P) =$ a probability space

FINDING PROBABILITY DISTRIBUTIONS

$(\Omega, \mathcal{E}) =$ a measurable space

Do you recognize this distribution?

Example: $\Omega = \{0, 1\}$
 $\mathcal{E} = \{\emptyset, \{0\}, \{1\}, \Omega\}$

$$P(A) = \begin{cases} 1 - \alpha & A = \{0\} \\ \alpha & A = \{1\} \\ 0 & A = \emptyset \\ 1 & A = \Omega \end{cases} \quad \alpha \in [0, 1]$$

How can we choose P in practice?

Clearly, we cannot do it arbitrarily.

How can we satisfy all constraints?

PROBABILITY MASS FUNCTIONS

Ω = discrete sample space

$\mathcal{E} = \mathcal{P}(\Omega)$

Probability mass function:

1. $p : \Omega \rightarrow [0, 1]$

2. $\sum_{\omega \in \Omega} p(\omega) = 1$

The probability of any event $A \in \mathcal{E}$ is defined as

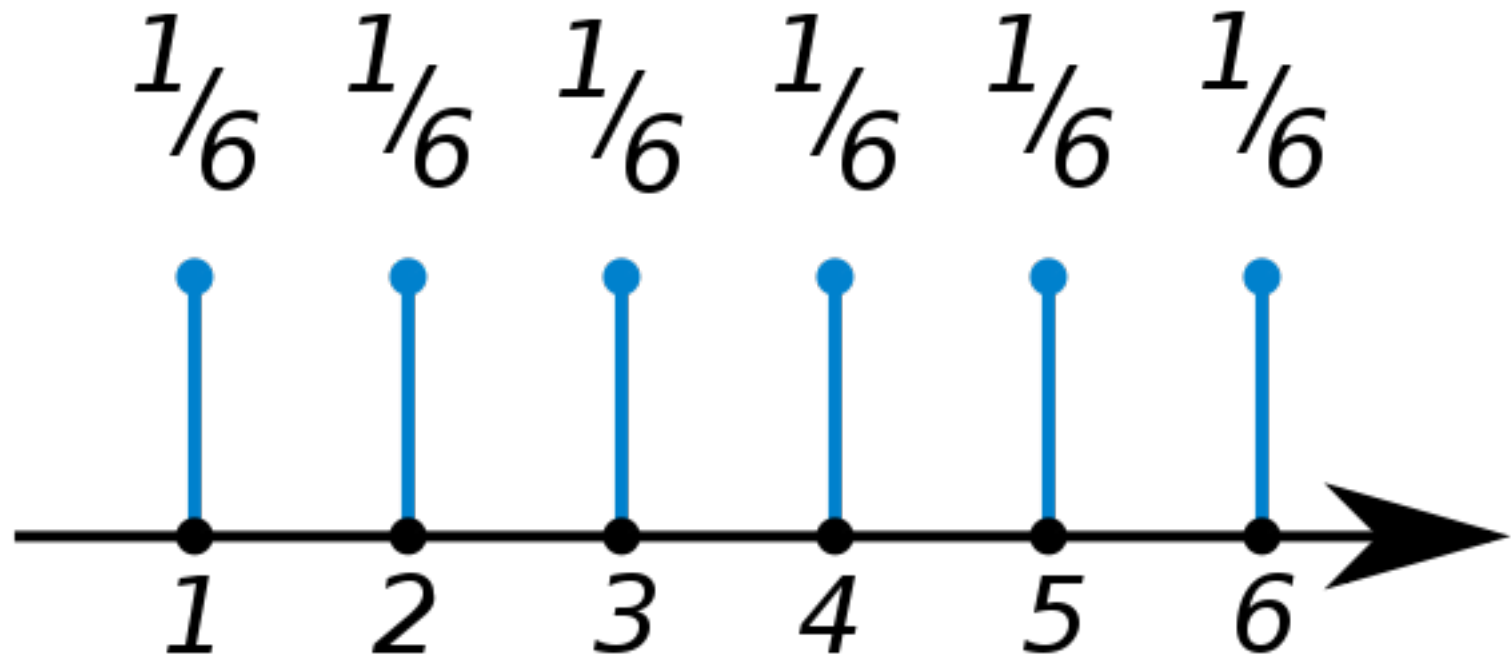
$$P(A) = \sum_{\omega \in A} p(\omega)$$

ARBITRARY PMFs

e.g. PMF for a fair die (table of values)

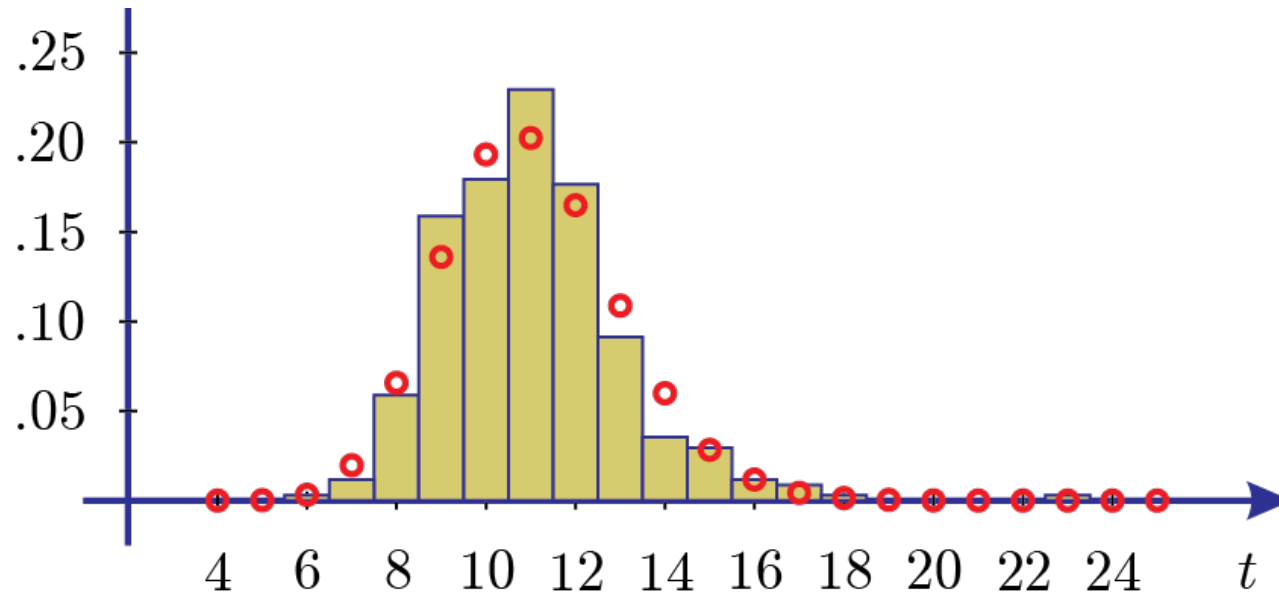
$$\Omega = \{1, 2, 3, 4, 5, 6\}$$

$$p(\omega) = 1/6 \quad \forall \omega \in \Omega$$



EXERCISE: HOW ARE PMFs USEFUL AS A MODEL?

- Recall we wanted to model commute times
- We could use a probability table for minutes: count number of times $t = 1, 2, 3, \dots$ occurs and then normalize probabilities by # samples
- Pick t with the largest $p(t)$



USEFUL PMFs

Bernoulli distribution:

$$\Omega = \{S, F\} \quad \alpha \in (0, 1)$$

$$p(\omega) = \begin{cases} \alpha & \omega = S \\ 1 - \alpha & \omega = F \end{cases}$$

Alternatively, $\Omega = \{0, 1\}$

$$p(k) = \alpha^k \cdot (1 - \alpha)^{1-k} \quad \forall k \in \Omega$$

USEFUL PMFs

Poisson distribution:

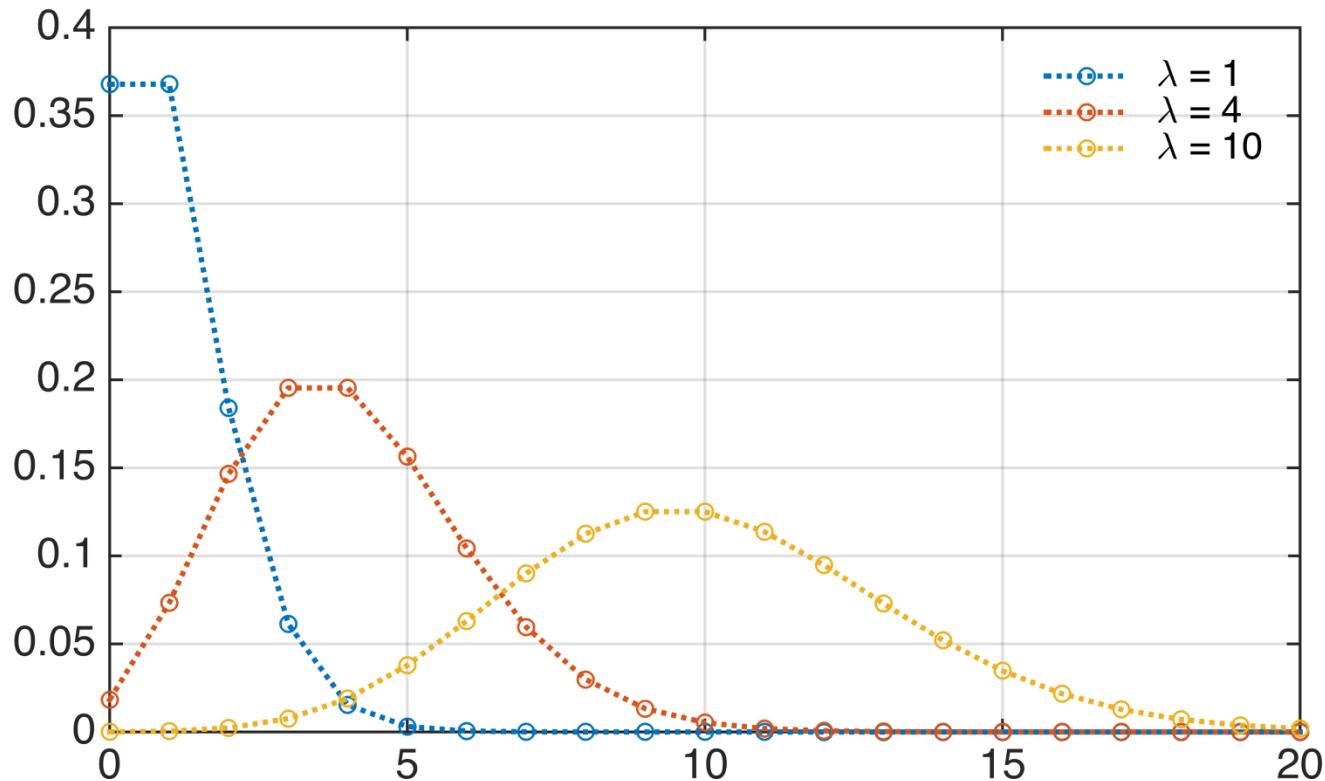
$$\Omega = \{0, 1, \dots\} \quad \lambda \in (0, \infty)$$

e.g., amount of mail received in a day

number of calls received by call center in an hour

$$p(k) = \frac{\lambda^k e^{-\lambda}}{k!}$$

$$\forall k \in \Omega$$



USEFUL PMFs

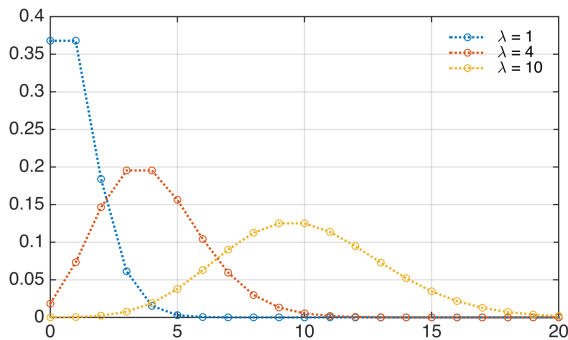
Poisson distribution:

e.g., amount of mail received in a day
number of calls received by call center in an hour

$$\Omega = \{0, 1, \dots\} \quad \lambda \in (0, \infty)$$

$$p(k) = \frac{\lambda^k e^{-\lambda}}{k!} \quad \forall k \in \Omega$$

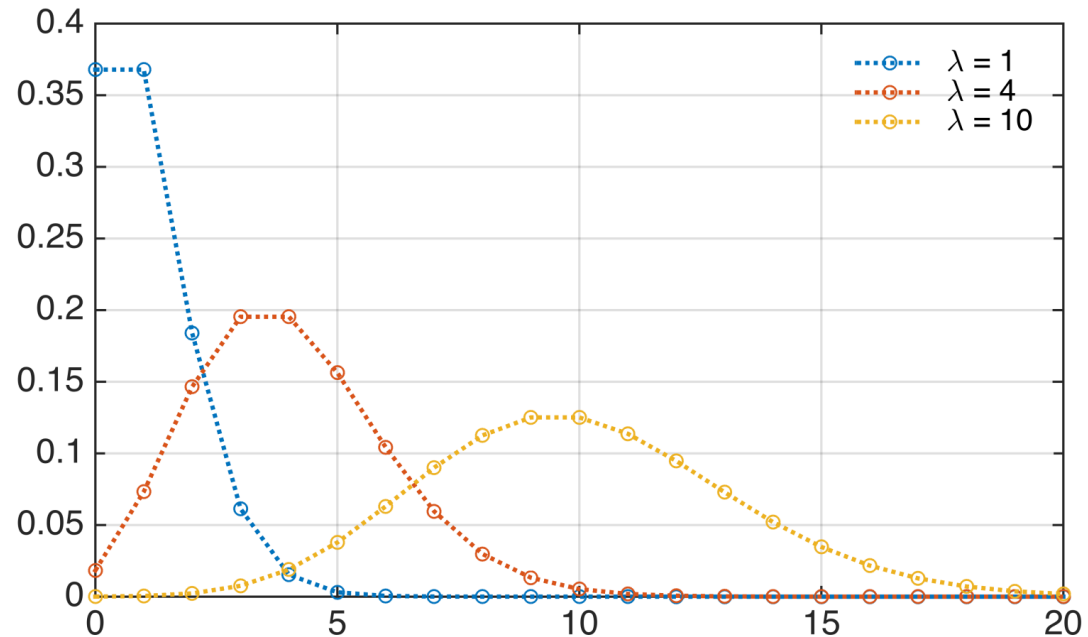
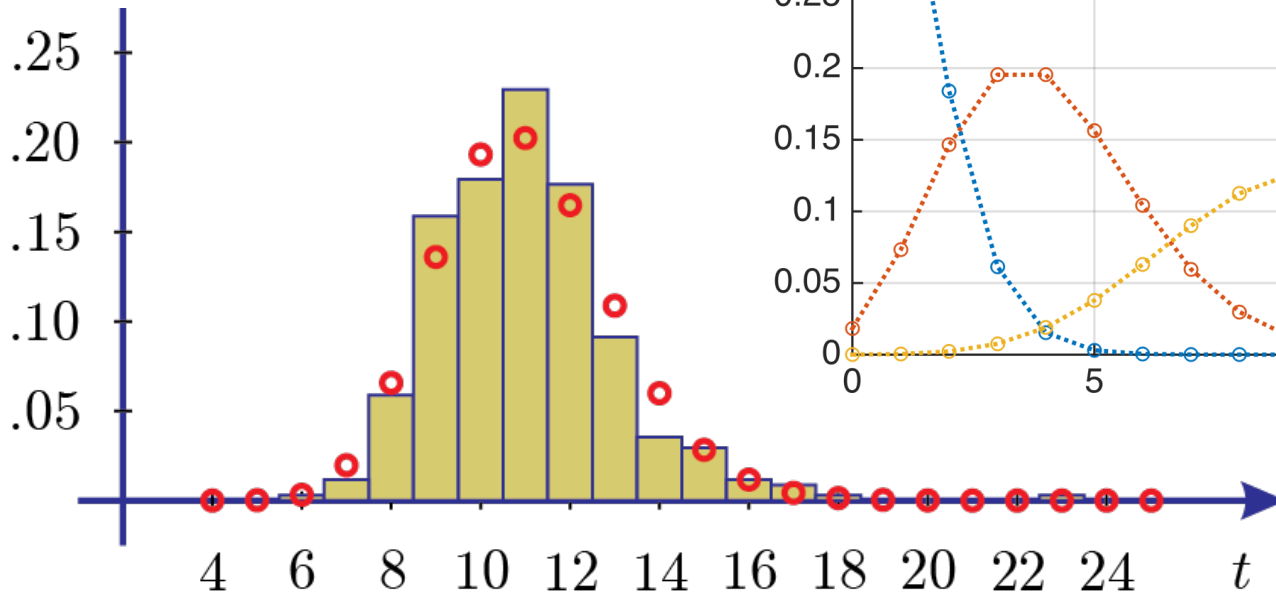
1. Can we use a table for this?
2. How do we know this is a valid pmf? How can you check?



EXERCISE: CAN WE USE A POISSON FOR COMMUTE TIMES?

- Used a probability table (histogram) for minutes: count number of times $t = 1, 2, 3, \dots$ occurs and then normalize probabilities by # samples
- Can we use a Poisson?

$$p(k) = \frac{\lambda^k e^{-\lambda}}{k!}$$



PROBABILITY DENSITY FUNCTIONS

Ω = continuous sample space

$$\mathcal{E} = \mathcal{B}(\Omega)$$

Probability density function:

1. $p : \Omega \rightarrow [0, \infty)$

2. $\int_{\Omega} p(\omega) d\omega = 1$

The probability of any event $A \in \mathcal{E}$ is defined as

$$P(A) = \int_A p(\omega) d\omega.$$

PMFs vs. PDFs

Ω = discrete sample space

Consider a singleton event $\{\omega\} \in \mathcal{E}$, where $\omega \in \Omega$

$$P(\{\omega\}) = p(\omega)$$

Ω = continuous sample space

Example:

- Stopping time of a car, in interval $[3, 15]$. What is the probability of seeing a stopping time of exactly 3.141596? (How much mass in $[3, 15]$?)
- More reasonable to ask the probability of stopping between 3 to 3.5 seconds.

PMFs vs. PDFs

$\Omega =$ discrete sample space

Consider a singleton event $\{\omega\} \in \mathcal{F}$, where $\omega \in \Omega$

$$P(\{\omega\}) = p(\omega)$$

$\Omega =$ continuous sample space

Consider an interval event $A = [x, x + \Delta x]$, where Δ is small

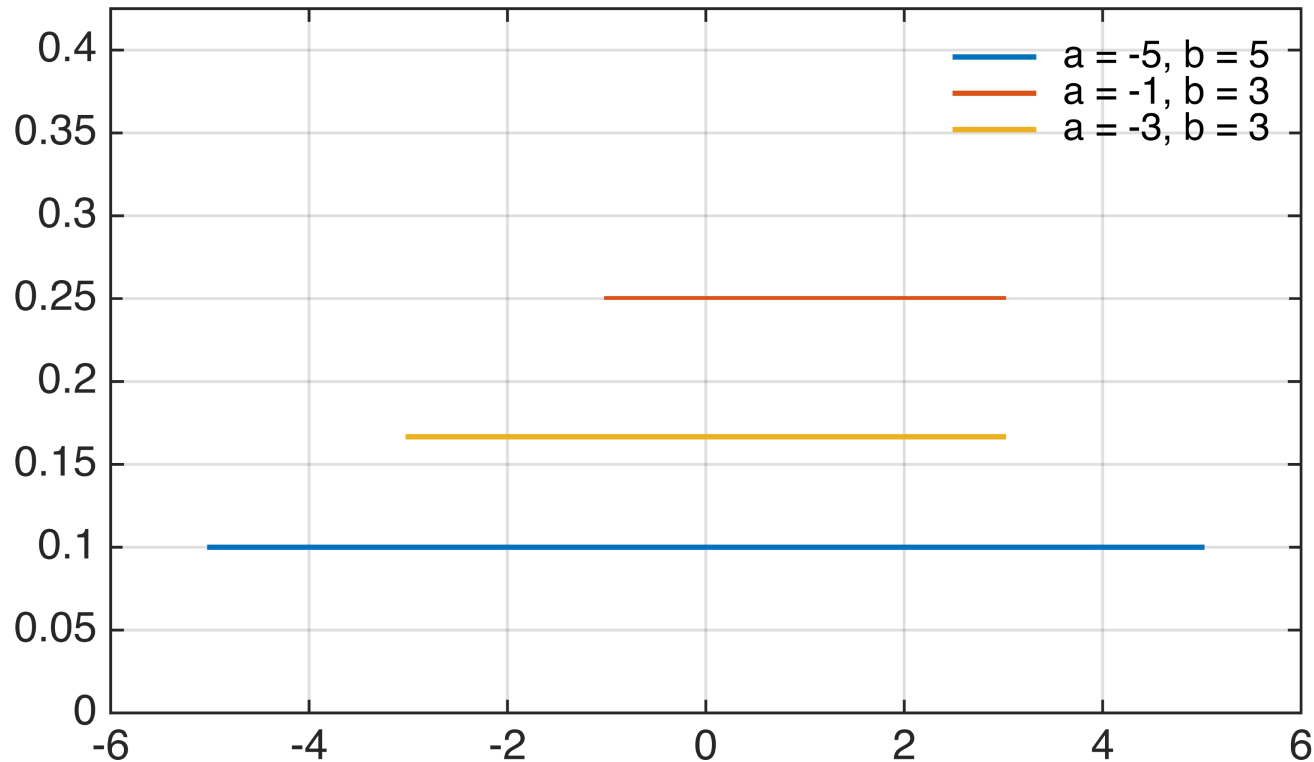
$$\begin{aligned} P(A) &= \int_x^{x+\Delta x} p(\omega) d\omega \\ &\approx p(x) \Delta x \end{aligned}$$

USEFUL PDFs

Uniform distribution: $\Omega = [a, b]$

$$p(\omega) = \frac{1}{b - a}$$

$\forall \omega \in [a, b]$

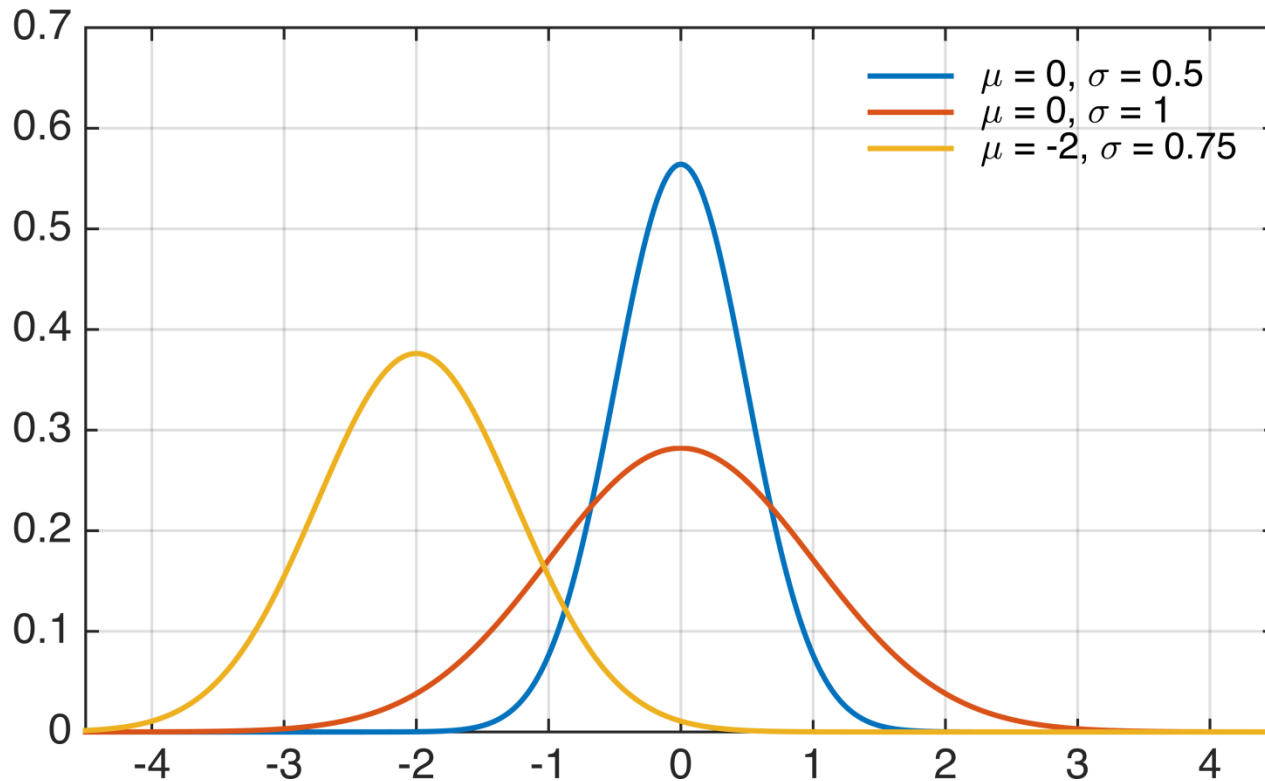


USEFUL PDFs

Gaussian distribution:

$$\Omega = \mathbb{R} \quad \mu \in \mathbb{R}, \sigma \in \mathbb{R}^+$$

$$p(\omega) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(\omega-\mu)^2} \quad \forall \omega \in \mathbb{R}$$



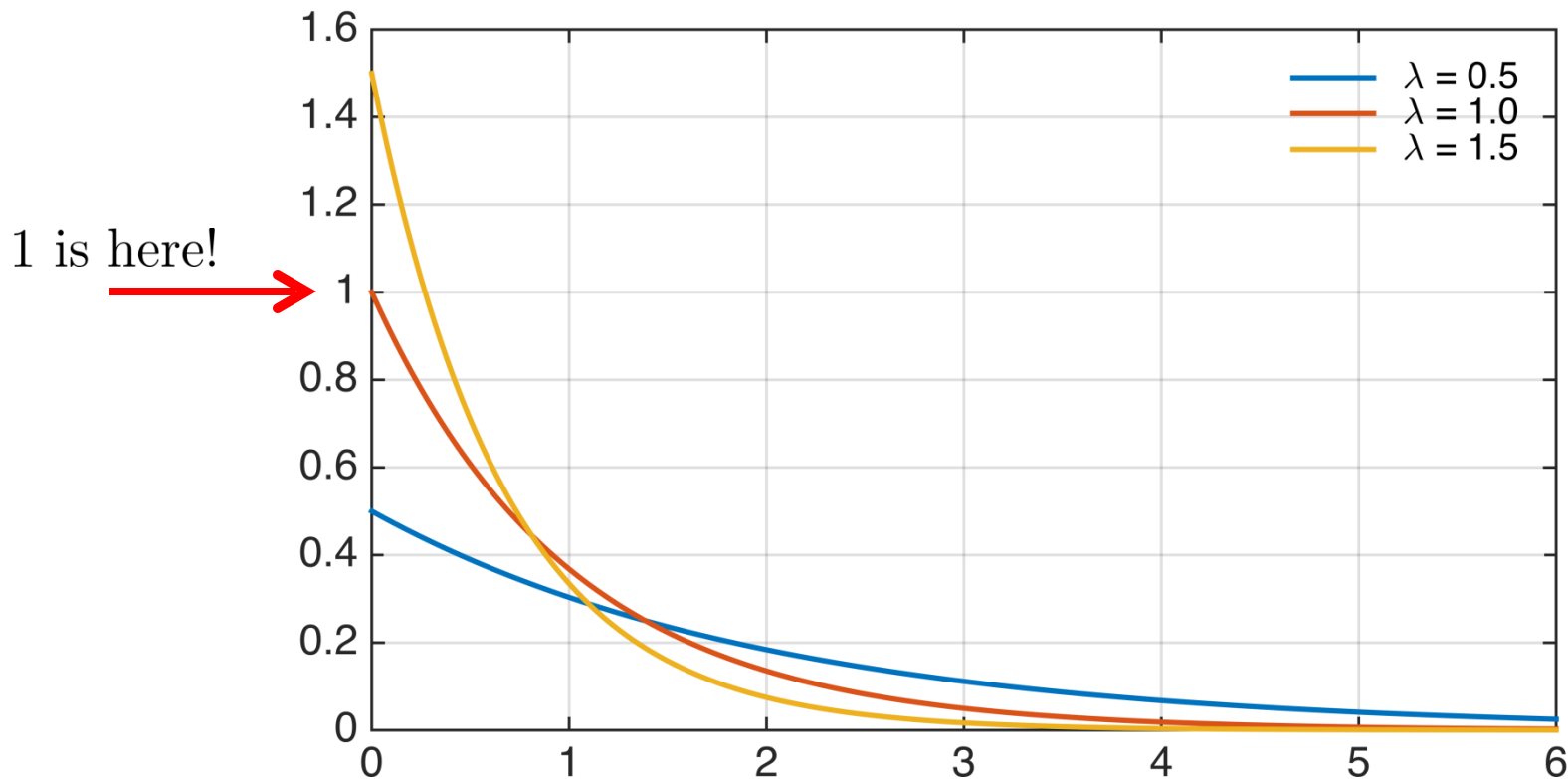
USEFUL PDFs

Exponential distribution:

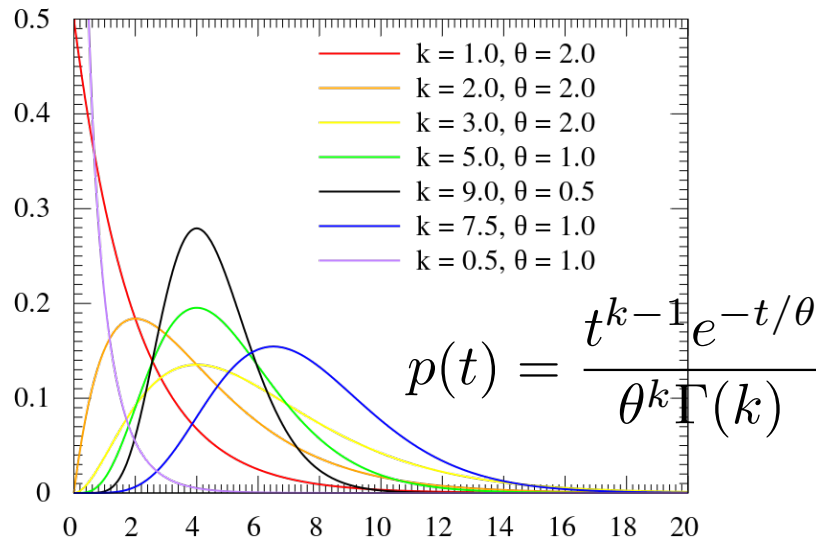
$$\Omega = [0, \infty) \quad \lambda > 0$$

$$p(\omega) = \lambda e^{-\lambda\omega}$$

$$\forall \omega \geq 0$$

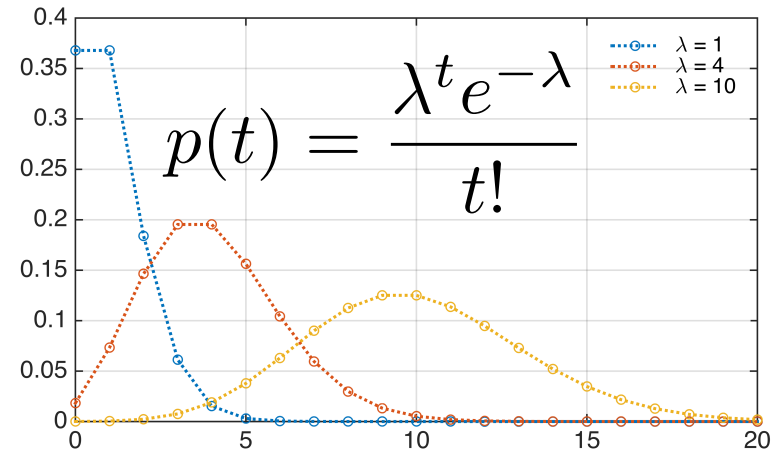
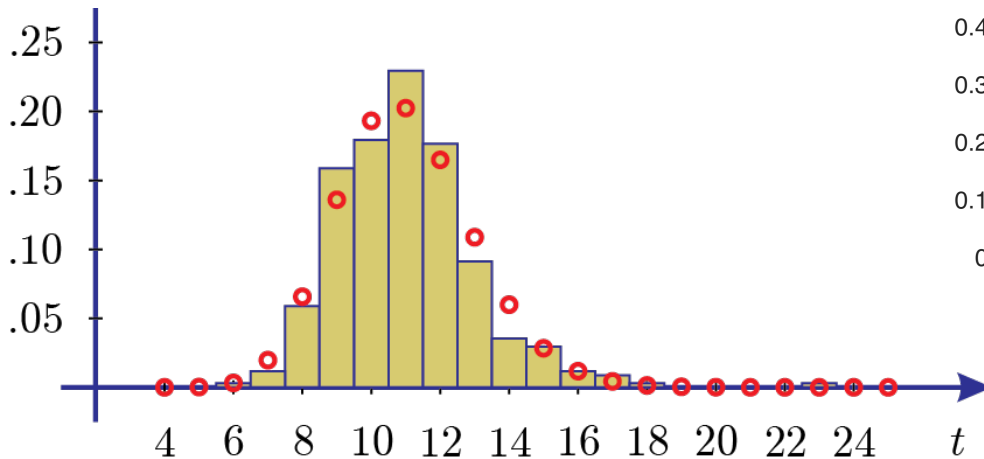


EXERCISE: MODELING COMMUTE TIMES

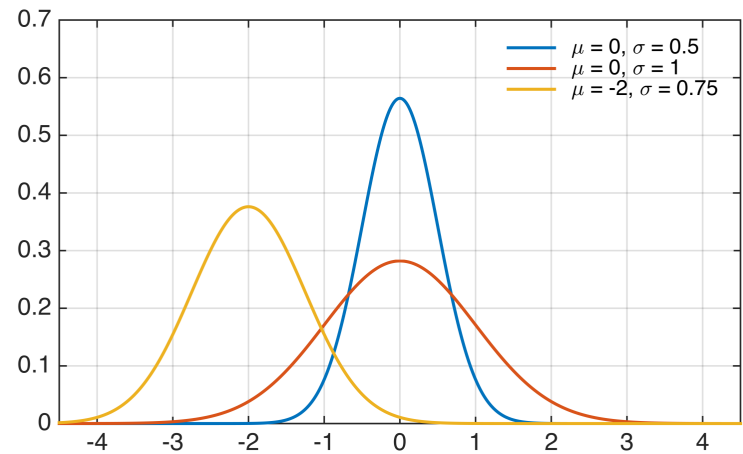


Gamma

Which might you choose?



Poisson



Gaussian

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

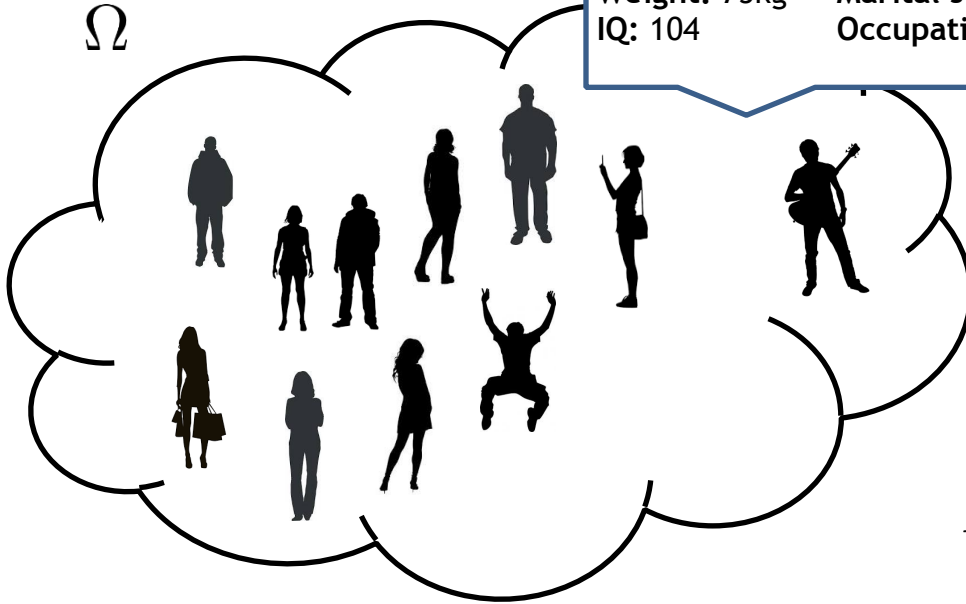
RANDOM VARIABLES

(Ω, \mathcal{E}, P)

Ω

Age: 35 Likes sports: Yes
Height: 1.85m Smokes: No
Weight: 75kg Marital st.: Single
IQ: 104 Occupation: Musician

Age: 26 Likes sports: Yes
Height: 1.75m Smokes: No
Weight: 79kg Marital st.: Divorced
IQ: 103 Occupation: Athlete



$$A = \{\omega \in \Omega : \text{Musician}(\omega) = \text{yes}\}$$

Musician is a random variable (a function)

A is a new event, let's call it 1; Not-A is 0

Omega is $\{0, 1\}$

Can ask $P(M = 0)$ and $P(M = 1)$

WE INSTINCTIVELY CREATE THIS TRANSFORMATION

Assume Ω is a set of people.

Compute the probability that a randomly selected person $\omega \in \Omega$ has a cold.

Define event $A = \{\omega \in \Omega : \text{Disease}(\omega) = \text{cold}\}$.

Disease is our new random variable, $P(\text{Disease} = \text{cold})$

Disease is a function that maps outcome space to new outcome space $\{\text{cold}, \text{not cold}\}$

Disease is a function, which is neither a variable nor random
BUT, this term is still a good one since we treat Disease as a variable
And assume it can take on different values
(randomly according to some distribution)

RANDOM VARIABLE: FORMAL DEFINITION

(Ω, \mathcal{E}, P) = a probability space

Random variable:

1. $X : \Omega \rightarrow \Omega_X$

2. $\forall A \in \mathcal{B}(\Omega_X)$ it holds that $\{\omega : X(\omega) \in A\} \in \mathcal{E}$

It follows that: $P_X(A) = P(\{\omega : X(\omega) \in A\})$

Example $X : \Omega \rightarrow [0, \infty)$

Ω is set of (measured) people in population

with associated measurements such as height and weight

$X(\omega)$ = height

A = interval = $[5'1'', 5'2'']$

$P(X \in A) = P(5'1'' \leq X \leq 5'2'') = P(\{\omega : X(\omega) \in A\})$

5 MINUTE BREAK AND EXERCISE

- Let X be a random variable that corresponds to the ratio of hard-to-easy problems on an assignment. Assume it takes values in $\{0.1, 0.25, 0.7\}$.
 - Is this discrete or continuous? Does it have a PMF or PDF?
 - Further, where could the variability come from? i.e., why is this a random variable?
- Let X be the stopping time of a car, taking values in $[3,5]$ union $[7,9]$. Is this discrete or continuous?
- Think of an example of a discrete random variable (RV) and a continuous RV

WHAT IF WE HAVE MORE THAN TWO VARIABLES...

- So far, we have considered scalar random variables
- Axioms of probability defined abstractly, apply to vector random variables

Ω = sample space, all outcomes of the experiment

\mathcal{E} = event space, set of subsets of Ω

$$\Omega = \mathbb{R}^2, \text{ e.g., } \omega = [-0.5, 10]$$

$$\Omega = [0, 1] \times [2, 5], \text{ e.g., } \omega = [0.2, 3.5]$$

But, when defining probabilities, we will want to consider how the variables interact

TWO DISCRETE RANDOM VARIABLES

Random variables X and Y

Outcome spaces \mathcal{X} and \mathcal{Y}

$$p(x, y) = P(X = x, Y = y)$$

$$\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) = 1.$$

$\mathcal{X} = \{\text{young, old}\}$ and $\mathcal{Y} = \{\text{no arthritis, arthritis}\}$.

		0	1
X	0	$P(X=0, Y=0) = 1/2$	$P(X=0, Y=1) = 1/100$
	1	$P(X=1, Y=0) = 1/10$	$P(X=1, Y=1) = 39/100$

Table 1.1: A joint probability table for random variables X and Y .

* these numbers are completely made up

SOME QUESTIONS WE MIGHT ASK NOW THAT WE HAVE TWO RANDOM VARIABLES

$\mathcal{X} = \{\text{young, old}\}$ and $\mathcal{Y} = \{\text{no arthritis, arthritis}\}$.

		Y	
		0	1
X	0	1/2	1/100
	1	1/10	39/100

Are these two variables related?

Or do they change completely independently of each other?

Given this joint distribution, can we determine just the distribution over arthritis? i.e., $P(Y = 1)$? (Marginal distribution)

If we knew something about one of the variables, say that the person is young, do we now know the distribution over Y? (Conditional distribution)

EXAMPLE: MARGINAL AND CONDITIONAL DISTRIBUTION

$\mathcal{X} = \{\text{young, old}\}$ and $\mathcal{Y} = \{\text{no arthritis, arthritis}\}$.

		Y	
		0	1
X	0	1/2	1/100
	1	1/10	39/100

$$P(Y = 1) = P(Y = 1, X = 0) + P(Y = 1, X = 1) = 40/100$$

What is $P(Y = 0)$?

$P(X = 1) = 49/100$. So what is $P(X = 0)$?

$P(Y = 1 \mid X = 0) = ?$

Is it $1/100$, where the table tells us $P(Y = 1, X=0)$?

No

CONDITIONAL DISTRIBUTIONS

Conditional probability distribution:

$$p(y|x) = \frac{p(x, y)}{p(x)}$$

The probability of an event A , given that $X = x$, is:

$$P(Y \in A|X = x) = \begin{cases} \sum_{y \in A} p(y|x) & Y : \text{discrete} \\ \int_A p(y|x) dy & Y : \text{continuous} \end{cases}$$

EXERCISE: CONDITIONAL DISTRIBUTION

$\mathcal{X} = \{\text{young, old}\}$ and $\mathcal{Y} = \{\text{no arthritis, arthritis}\}$.

		Y	
		0	1
X	0	1/2	1/100
	1	1/10	39/100

$$p(y|x) = \frac{p(x, y)}{p(x)}$$

$P(Y = 1 \mid X = 0) = ?$

What is $P(Y = 0 \mid X = 0)$?

Should $P(Y = 1 \mid X = 0) + P(Y = 0 \mid X = 0) = 1$?

JOINT DISTRIBUTIONS FOR MANY VARIABLES

In general, we can consider d -dimensional random variable $\mathbf{X} = (X_1, X_2, \dots, X_d)$ with vector-valued outcomes $\mathbf{x} = (x_1, x_2, \dots, x_d)$, such that each x_i is chosen from some \mathcal{X}_i . Then, for the discrete case, any function $p : \mathcal{X}_1 \times \mathcal{X}_2 \times \dots \times \mathcal{X}_d \rightarrow [0, 1]$ is called a multidimensional probability mass function if

$$\sum_{x_1 \in \mathcal{X}_1} \sum_{x_2 \in \mathcal{X}_2} \cdots \sum_{x_d \in \mathcal{X}_d} p(x_1, x_2, \dots, x_d) = 1.$$

or, for the continuous case, $p : \mathcal{X}_1 \times \mathcal{X}_2 \times \dots \times \mathcal{X}_d \rightarrow [0, \infty]$ is a multidimensional probability density function if

$$\int_{\mathcal{X}_1} \int_{\mathcal{X}_2} \cdots \int_{\mathcal{X}_d} p(x_1, x_2, \dots, x_d) dx_1 dx_2 \dots dx_d = 1.$$

MARGINAL DISTRIBUTIONS

A *marginal distribution* is defined for a subset of $\mathbf{X} = (X_1, X_2, \dots, X_d)$ by summing or integrating over the remaining variables. For the discrete case, the marginal distribution $p(x_i)$ is defined as

$$p(x_i) = \sum_{x_1 \in \mathcal{X}_1} \cdots \sum_{x_{i-1} \in \mathcal{X}_{i-1}} \sum_{x_{i+1} \in \mathcal{X}_{i+1}} \cdots \sum_{x_d \in \mathcal{X}_d} p(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_d),$$

where the variable x_i is fixed to some value and we sum over all possible values of the other variables. Similarly, for the continuous case, the marginal distribution $p(x_i)$ is defined as

$$p(x_i) = \int_{\mathcal{X}_1} \cdots \int_{\mathcal{X}_{i-1}} \int_{\mathcal{X}_{i+1}} \cdots \int_{\mathcal{X}_d} p(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_d) dx_1 \dots dx_{i-1} dx_{i+1} \dots dx_d.$$

Natural question: Why do you use p for $p(x_i)$ and for $p(x_1, \dots, x_d)$?
They have different domains, they can't be the same function!

DROPPING SUBSCRIPTS

Instead of:

$$p_{Y|X}(y|x) = \frac{p_{XY}(x, y)}{p_X(x)}$$

We will write:

$$p(y|x) = \frac{p(x, y)}{p(x)}$$

ANOTHER EXAMPLE FOR CONDITIONAL DISTRIBUTIONS

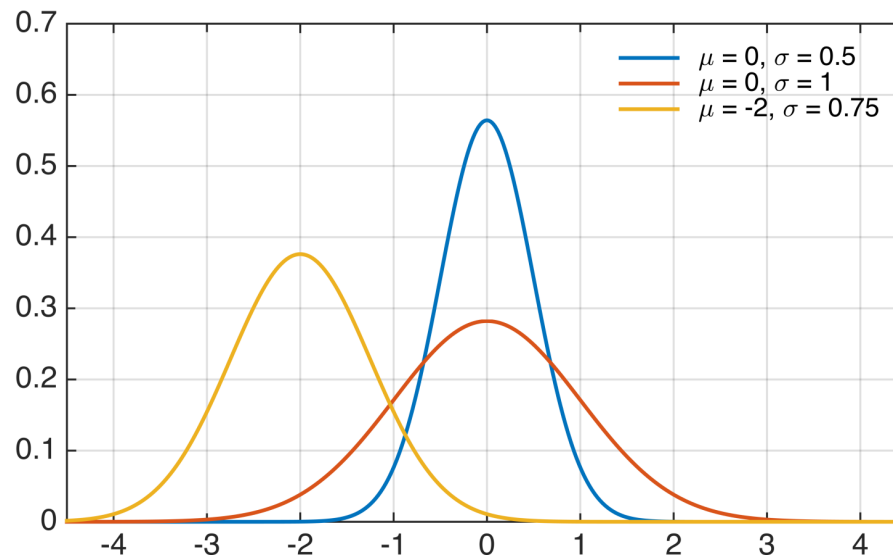
- Let **X** be a Bernoulli random variable (i.e., 0 or 1 with probability α)
- Let **Y** be a random variable in $\{10, 11, \dots, 1000\}$
- $p(y \mid X = 0)$ and $p(y \mid X = 1)$ are different distributions
- Two **types of books**: fiction ($X=0$) and non-fiction ($X=1$)
- Let **Y** correspond to **number of pages**
- Distribution over number of pages different for fiction and non-fiction books (e.g., average different)

EXAMPLE CONTINUED

- Two types of books: fiction ($X=0$) and non-fiction ($X=1$)
- Y corresponds to number of pages
- If most books are non-fiction, $p(X = 0, y)$ is small even if y is a likely number of pages for a fiction book
- $p(X = 0)$ accounts for the fact that joint probability small if $p(X = 0)$ is small
 - $p(y | X = 0) = p(X = 0, y)/p(X = 0)$
 - $p(X = 0, y)$ = probability that a book is fiction and has y pages (imagine randomly sampling a book)
 - $p(X = 0)$ = probability that a book is fiction

ANOTHER EXAMPLE

- Two types of books: fiction ($X=0$) and non-fiction ($X=1$)
- Let Y be a random variable over the reals, which corresponds to amount of money made
- $p(y | X = 0)$ and $p(y | X = 1)$ are different distributions
- e.g., even if both $p(y | X = 0)$ and $p(y | X = 1)$ are Gaussian, they likely have different means and variances



Review so far

- PMFs (discrete) and PDFs (continuous)
- Joint probabilities
- Marginals
- Conditional probabilities
- Chain rule (and Bayes rule)

CHAIN RULE

Conditional probability distribution:

$$p(x_k | x_1, \dots, x_{k-1}) = \frac{p(x_1, \dots, x_k)}{p(x_1, \dots, x_{k-1})}$$

This leads to:

$$\begin{aligned} p(x_1, \dots, x_k) &= p(x_k) \prod_{i=1}^{k-1} p(x_i | x_{i+1}, \dots, x_k) \\ &= p(x_1) \prod_{i=2}^k p(x_i | x_1, \dots, x_{i-1}) \end{aligned}$$

Two variable example $p(x, y) = p(x|y)p(y) = p(y|x)p(x)$

CHAIN RULE

Conditional probability distribution:

$$p(x_k | x_1, \dots, x_{k-1}) = \frac{p(x_1, \dots, x_k)}{p(x_1, \dots, x_{k-1})}$$

This leads to:

$$p(x_1, \dots, x_k) = p(x_k) \prod_{i=1}^{k-1} p(x_i | x_{i+1}, \dots, x_k)$$

Three variable example

$$\begin{aligned} p(x, y, z) &= p(y|x, z)p(x, z) = p(y|x, z)p(x|z)p(z) \\ &= p(x|y, z)p(y|z)p(z) \\ &= p(x|y, z)p(z|y)p(y) \\ &\vdots \end{aligned}$$

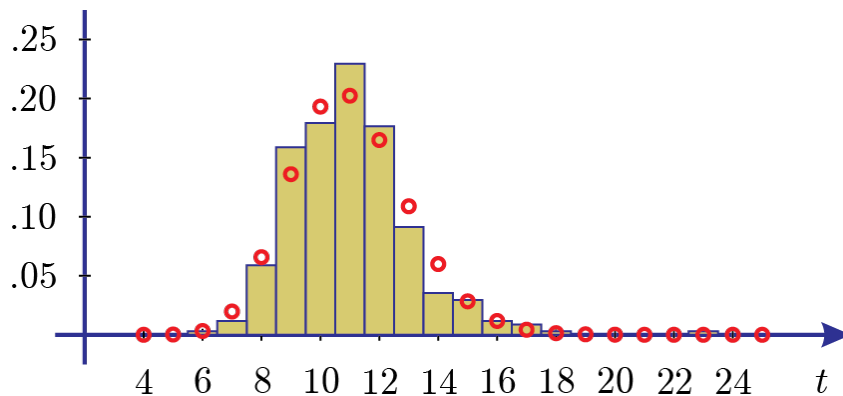
HOW DO WE GET BAYES RULE?

Recall chain rule: $p(x, y) = p(x|y)p(y) = p(y|x)p(x)$

Bayes rule:
$$p(y|x) = \frac{p(x|y)p(y)}{p(x)}$$

EXERCISE: CONDITIONAL PROBABILITIES

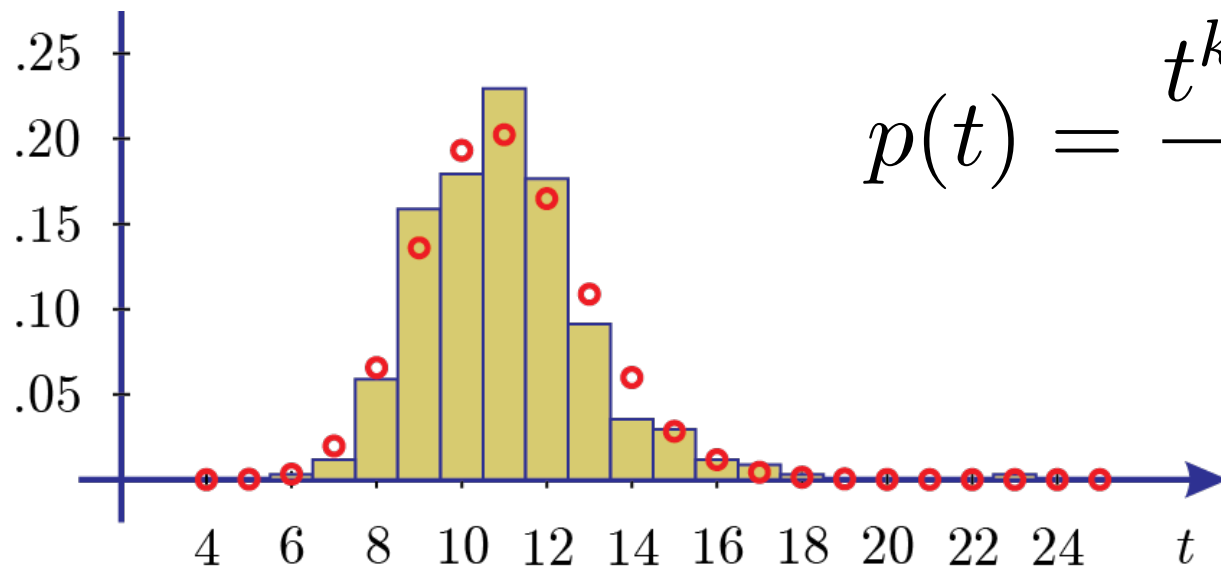
- Using conditional probabilities, we can incorporate other external information (features)
- Let y be the commute time, x the day of the year
- Array of conditional probability values $\rightarrow p(y | x)$
 - $y = 1, 2, \dots$ and $x = 1, 2, \dots, 365$
- What are some issues with this choice for x ?
- What other x could we use feasibly?



EXERCISE: ADDING IN AUXILIARY INFORMATION

- Gamma distribution for commute times extrapolates between recorded time in minutes
- Can incorporate external information (features) by modeling $\theta = \text{function}(\text{features})$

$$\theta = \sum_{i=1}^d w_i x_i$$



$$p(t) = \frac{t^{k-1} e^{-t/\theta}}{\theta^k \Gamma(k)}$$

REMINDERS: SEPTEMBER 10, 2019

- Assignment 1 is due on September 26
 - 566 has to do extra question (bonus for 466)
- Thought questions 1 are due on September 19
 - Preface and Chapters 1-4, about 40 pages
- I have Office Hours today, from 2 - 4 p.m.
- TAs are scheduling office hours
- Any questions?

INDEPENDENCE OF RANDOM VARIABLES

X and Y are **independent** if:

$$p(x, y) = p(x)p(y)$$

X and Y are **conditionally independent** given Z if:

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

CONDITIONAL INDEPENDENCE EXAMPLES

EXAMPLE 7 IN THE NOTES

- Imagine you have a biased coin (does not flip 50% heads and 50% tails, but skewed towards one)
- Let Z = bias of a coin (say outcomes are 0.3, 0.5, 0.8 with associated probabilities 0.7, 0.2, 0.1)
 - what other outcome space could we consider?
 - what kinds of distributions?
- Let X and Y be consecutive flips of the coin
- Are X and Y independent?
- Are X and Y conditionally independent, given Z ?

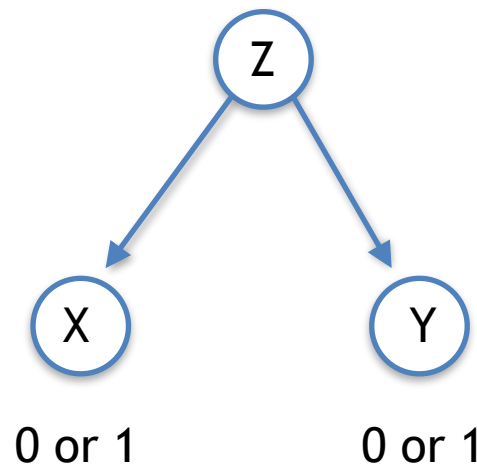
** (Basic example about an important issue in ML: hidden variables)

CONDITIONAL INDEPENDENCE EXAMPLES

EXAMPLE 7 IN THE NOTES

- Are X and Y independent? (don't know Z) $p(X, Y) = p(X)p(Y)$?
- Are X and Y conditionally independent, given Z ?

$$p(X, Y|Z) = p(X|Z)p(Y|Z)?$$



z	0.3	0.5	0.8
$p(z)$	0.7	0.2	0.1

bias

probability of that bias

Imagine don't know Z and flip two 0s. Does that tell you anything about Z ?

** (Basic example about an important issue in ML: hidden variables)

CONDITIONAL INDEPENDENCE EXAMPLES

EXAMPLE 7 IN THE NOTES

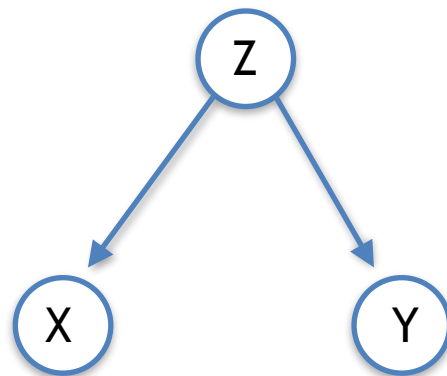
- What is $p(X, Y)$? $p(X, Y) = \sum_z p(X, Y, Z)$ Called the Law of Total Probability

$$= \sum_z p(X, Y|Z)p(Z)$$

- What is $p(X | Z)$? Its a Bernoulli

- Conditional independence: $p(X, Y|Z) = p(X|Z)p(Y|Z)$

$$P(X, Y) \neq P(X)P(Y)$$



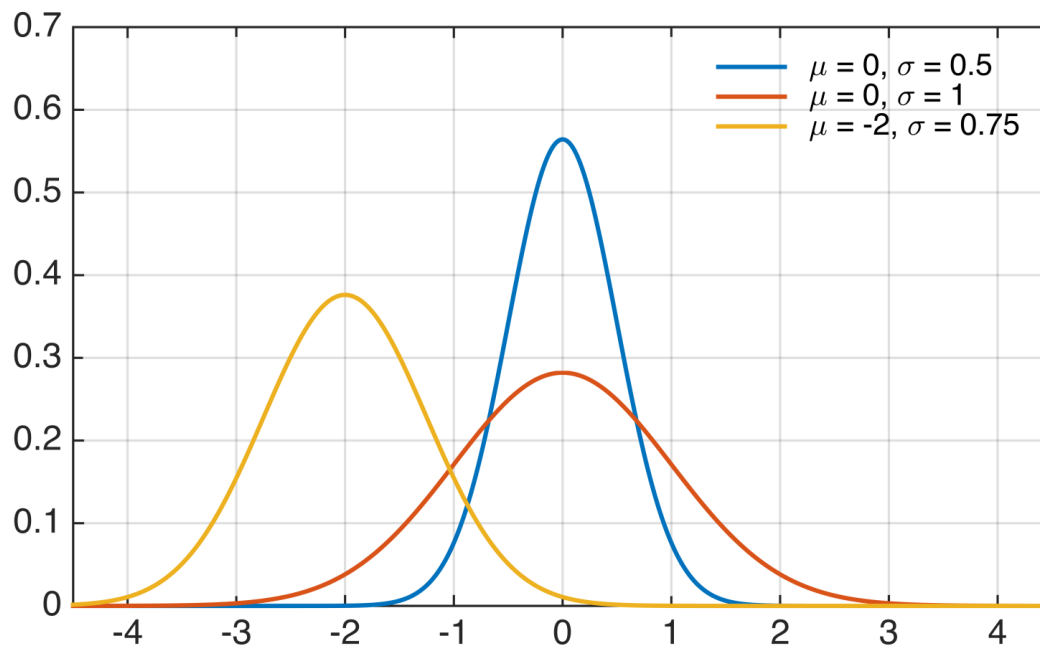
z	0.3	0.5	0.8
$p(z)$	0.7	0.2	0.1

bias

probability of that bias

EXPECTED VALUE (MEAN)

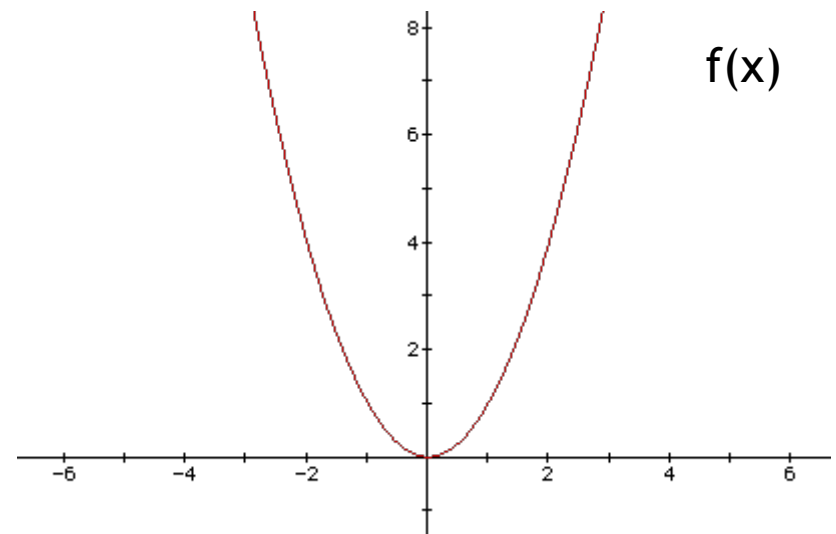
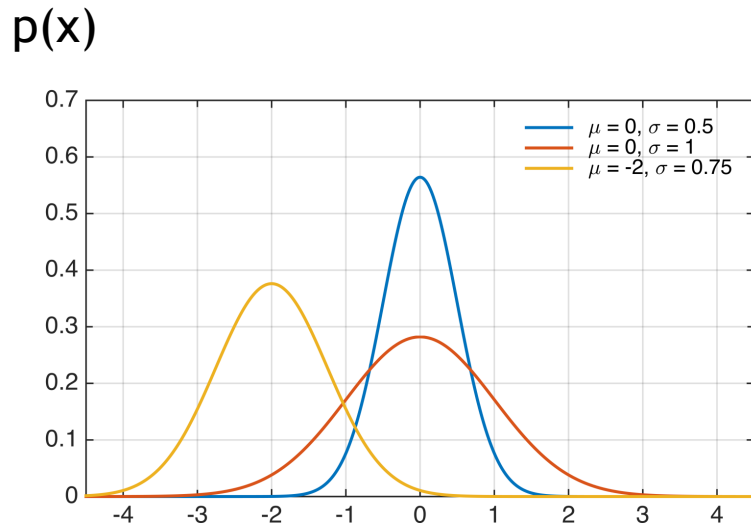
$$\mathbb{E}[X] = \begin{cases} \sum_{x \in \mathcal{X}} xp(x) & X : \text{discrete} \\ \int_{\mathcal{X}} xp(x)dx & X : \text{continuous} \end{cases}$$



EXPECTATIONS WITH FUNCTIONS

$$f : \mathcal{X} \rightarrow \mathbb{R}$$

$$\mathbb{E} [f(X)] = \begin{cases} \sum_{x \in \mathcal{X}} f(x)p(x) & X : \text{discrete} \\ \int_{\mathcal{X}} f(x)p(x)dx & X : \text{continuous} \end{cases}$$



CONDITIONAL EXPECTATIONS

$$\mathbb{E}[Y|X = x] = \begin{cases} \sum_{y \in \mathcal{Y}} yp(y|x) & Y : \text{discrete} \\ \int_{\mathcal{Y}} yp(y|x)dy & Y : \text{continuous} \end{cases}$$

Different expected value, depending on which x is observed

Example: $p(y | x) = \text{Gaussian with } N(\mu = x, \sigma^2 = 10)$

What is the $E[Y | x]$?

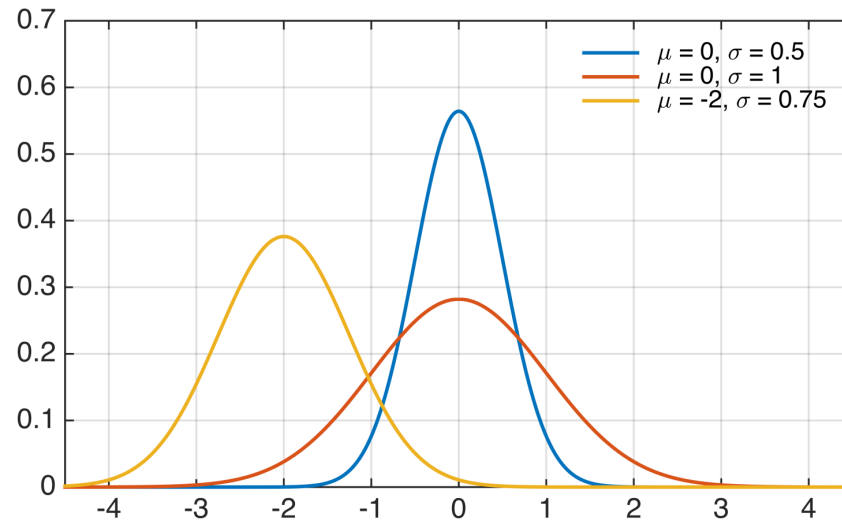
Example: $p(y | x) = \text{Gaussian with } N(\mu = f(x), \sigma^2 = 0.1)$

What is the $E[Y | x]$?

VARIANCE

$$\begin{aligned}\text{Variance}(X) &= \mathbb{E}[(X - \mathbb{E}[X])^2] \\ &= \mathbb{E}[X^2] - \mathbb{E}[X]^2\end{aligned}$$

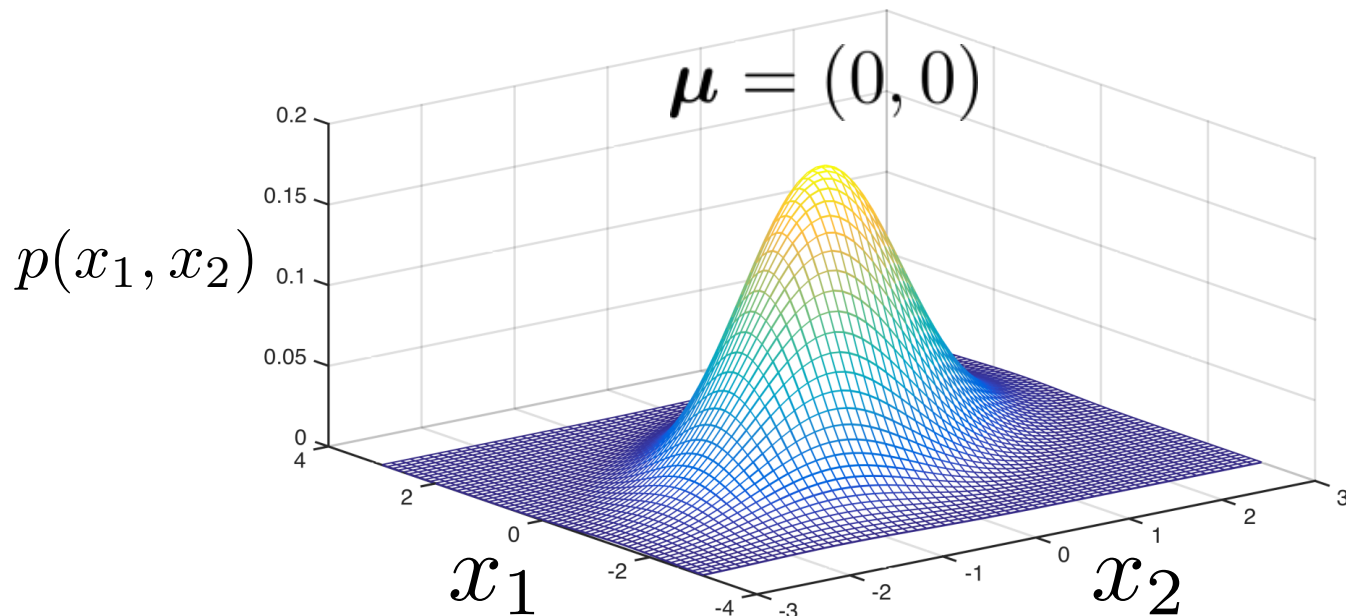
Why? See if you can get this formula



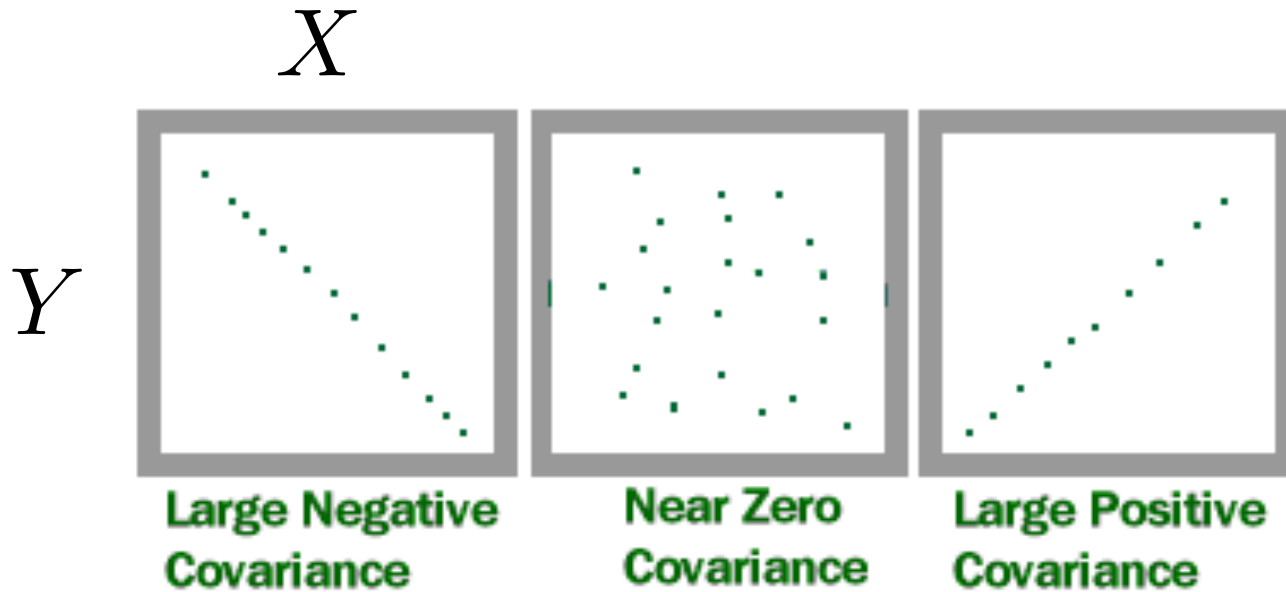
EXPECTED VALUE FOR MULTIVARIATE

$$\mathbb{E} [\mathbf{X}] = \begin{cases} \sum_{\mathbf{x} \in \mathcal{X}} \mathbf{x} p(\mathbf{x}) & \mathbf{X} : \text{discrete} \\ \int_{\mathcal{X}} \mathbf{x} p(\mathbf{x}) d\mathbf{x} & \mathbf{X} : \text{continuous} \end{cases}$$

Each instance \mathbf{x} is a vector, p is a function on these vectors



COVARIANCE



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

$$\text{Corr}[X, Y] = \frac{\text{Cov}[X, Y]}{\sqrt{V[X] \cdot V[Y]}}$$

COVARIANCE FOR MORE THAN TWO DIMENSIONS

$$\mathbf{X} = [X_1, \dots, X_d]$$

$$\begin{aligned}\Sigma_{ij} &= \text{Cov}[X_i, X_j] \\ &= \mathbb{E}[(X_i - \mathbb{E}[X_i])(X_j - \mathbb{E}[X_j])]\end{aligned}$$

$$\begin{aligned}\Sigma &= \text{Cov}[\mathbf{X}, \mathbf{X}] \in \mathbb{R}^{d \times d} \\ &= \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])^\top] \\ &= \mathbb{E}[\mathbf{X}\mathbf{X}^\top] - \mathbb{E}[\mathbf{X}]\mathbb{E}[\mathbf{X}]^\top.\end{aligned}$$

COVARIANCE FOR MORE THAN TWO DIMENSIONS

$$\begin{aligned}\mathbf{X} &= [X_1, \dots, X_d] & \boldsymbol{\Sigma} &= \text{Cov}[\mathbf{X}, \mathbf{X}] \in \mathbb{R}^{d \times d} \\ & & &= \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])^\top] \\ & & &= \mathbb{E}[\mathbf{X}\mathbf{X}^\top] - \mathbb{E}[\mathbf{X}]\mathbb{E}[\mathbf{X}]^\top.\end{aligned}$$

$$\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$$

Dot product

$$\mathbf{x}^\top \mathbf{y} = \sum_{i=1}^d x_i y_i$$

Outer product

$$\mathbf{xy}^\top = \begin{bmatrix} x_1 y_1 & x_1 y_2 & \dots & x_1 y_d \\ x_2 y_1 & x_2 y_2 & \dots & x_2 y_d \\ \vdots & \vdots & & \vdots \\ x_d y_1 & x_d y_2 & \dots & x_d y_d \end{bmatrix}$$

COVARIANCE FOR MORE THAN TWO DIMENSIONS

$$\begin{aligned}\mathbf{X} &= [X_1, \dots, X_d] & \boldsymbol{\Sigma} &= \text{Cov}[\mathbf{X}, \mathbf{X}] \in \mathbb{R}^{d \times d} \\ & & &= \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])^\top] \\ & & &= \mathbb{E}[\mathbf{X}\mathbf{X}^\top] - \mathbb{E}[\mathbf{X}]\mathbb{E}[\mathbf{X}]^\top.\end{aligned}$$

$$\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$$

Example:

$$\mathbb{E} \begin{bmatrix} X_1^2 & X_1 X_2 \\ X_2 X_1 & X_2^2 \end{bmatrix} - \begin{bmatrix} \mathbb{E}[X_1]^2 & \mathbb{E}[X_1]\mathbb{E}[X_2] \\ \mathbb{E}[X_2]\mathbb{E}[X_1] & \mathbb{E}[X_2]^2 \end{bmatrix}$$

SOME USEFUL PROPERTIES

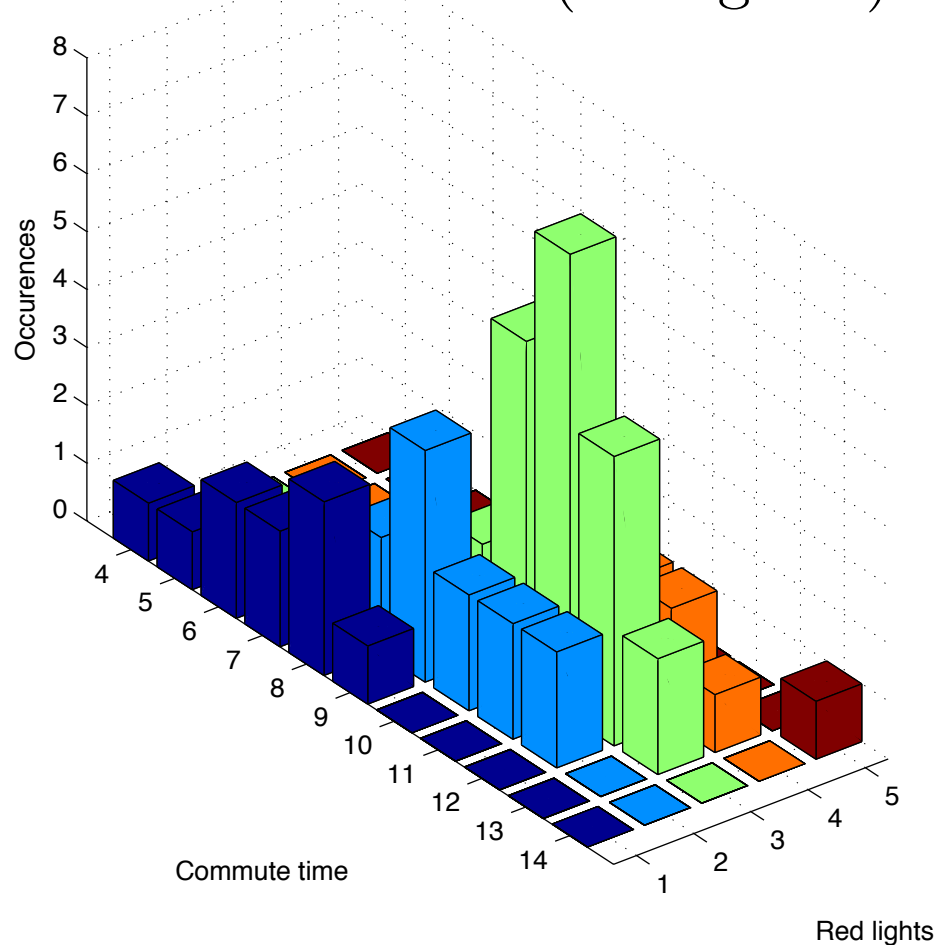
1. $\mathbb{E}[c\mathbf{X}] = c\mathbb{E}[\mathbf{X}]$
2. $\mathbb{E}[\mathbf{X} + \mathbf{Y}] = \mathbb{E}[\mathbf{X}] + \mathbb{E}[\mathbf{Y}]$
3. $V[c] = 0$ \triangleright the variance of a constant is zero
4. $V[\mathbf{X}] \succeq 0$ (i.e., is positive semi-definite), where for $d = 1$, $V[\mathbf{X}] \geq 0$
 $V[\mathbf{X}]$ is shorthand for $\text{Cov}[\mathbf{X}, \mathbf{X}]$.
5. $V[c\mathbf{X}] = c^2V[\mathbf{X}]$.
6. $\text{Cov}[\mathbf{X}, \mathbf{Y}] = \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{Y} - \mathbb{E}[\mathbf{Y}])^\top] = \mathbb{E}[\mathbf{X}\mathbf{Y}^\top] - \mathbb{E}[\mathbf{X}]\mathbb{E}[\mathbf{Y}]^\top$
7. $V[\mathbf{X} + \mathbf{Y}] = V[\mathbf{X}] + V[\mathbf{Y}] + 2\text{Cov}[\mathbf{X}, \mathbf{Y}]$

MULTIDIMENSIONAL PMF

Now record both commute time and number red lights

$$\Omega = \{4, \dots, 14\} \times \{1, 2, 3, 4, 5\}$$

PMF is normalized 2-d table (histogram) of occurrences



MULTIDIMENSIONAL GAUSSIAN

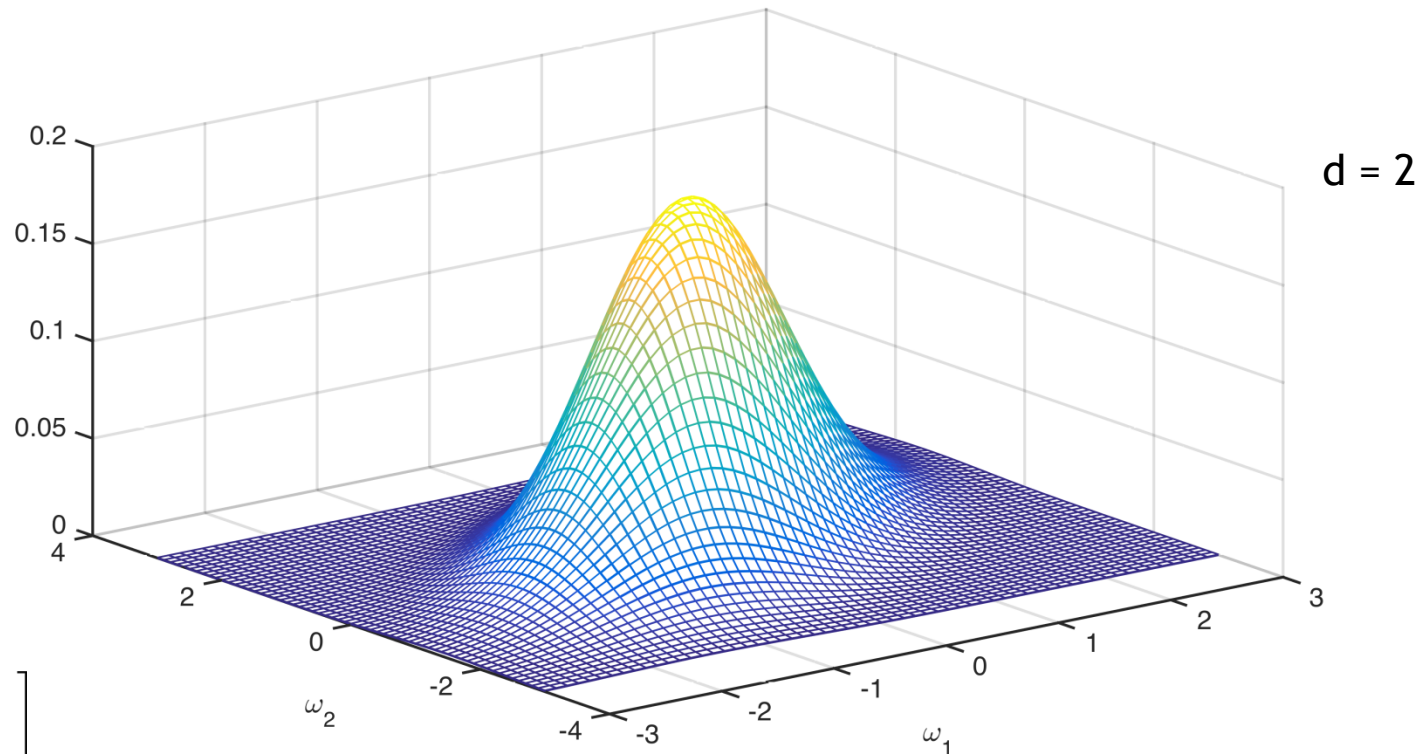
$$\Omega = \mathbb{R}^d$$
$$\mathcal{F} = \mathcal{B}(\mathbb{R})^d$$

$$\boldsymbol{\mu} \in \mathbb{R}^d$$

$\boldsymbol{\Sigma}$ = positive definite $d \times d$ matrix

$|\boldsymbol{\Sigma}|$ = determinant of $\boldsymbol{\Sigma}$

$$p(\boldsymbol{\omega}) = \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2}(\boldsymbol{\omega} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{\omega} - \boldsymbol{\mu})\right)$$

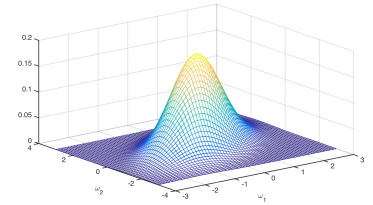


$$\boldsymbol{\mu} = (0, 0)$$

$$\boldsymbol{\Sigma} = \begin{bmatrix} 1 & .75 \\ .75 & 1 \end{bmatrix}$$

Example of multivariate Gaussian

$$p(\boldsymbol{\omega}) = \frac{1}{\sqrt{(2\pi)^k |\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2}(\boldsymbol{\omega} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{\omega} - \boldsymbol{\mu})\right)$$



$$\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \quad \boldsymbol{\Sigma} = \begin{bmatrix} 10 & 0 \\ 0 & 2 \end{bmatrix} \quad \boldsymbol{\Sigma}^{-1} = \begin{bmatrix} \frac{1}{10} & 0 \\ 0 & \frac{1}{2} \end{bmatrix}$$

$$\boldsymbol{\omega} - \boldsymbol{\mu} = \begin{bmatrix} \omega_1 - \mu_1 \\ \omega_2 - \mu_2 \end{bmatrix}$$

$$\begin{bmatrix} \omega_1 - \mu_1 \\ \omega_2 - \mu_2 \end{bmatrix} \begin{bmatrix} \frac{1}{10} & 0 \\ 0 & \frac{1}{2} \end{bmatrix} = \begin{bmatrix} \frac{1}{10}(\omega_1 - \mu_1) \\ \frac{1}{2}(\omega_2 - \mu_2) \end{bmatrix}$$

$$\begin{bmatrix} \frac{1}{10}(\omega_1 - \mu_1) \\ \frac{1}{2}(\omega_2 - \mu_2) \end{bmatrix}^T \begin{bmatrix} \omega_1 - \mu_1 \\ \omega_2 - \mu_2 \end{bmatrix} = \frac{1}{10}(\omega_1 - \mu_1)^2 + \frac{1}{2}(\omega_2 - \mu_2)^2$$

MIXTURES OF DISTRIBUTIONS

Mixture model:

A set of m probability distributions, $\{p_i(x)\}_{i=1}^m$

$$p(x) = \sum_{i=1}^m w_i p_i(x)$$

where $\mathbf{w} = (w_1, w_2, \dots, w_m)$ and non-negative and

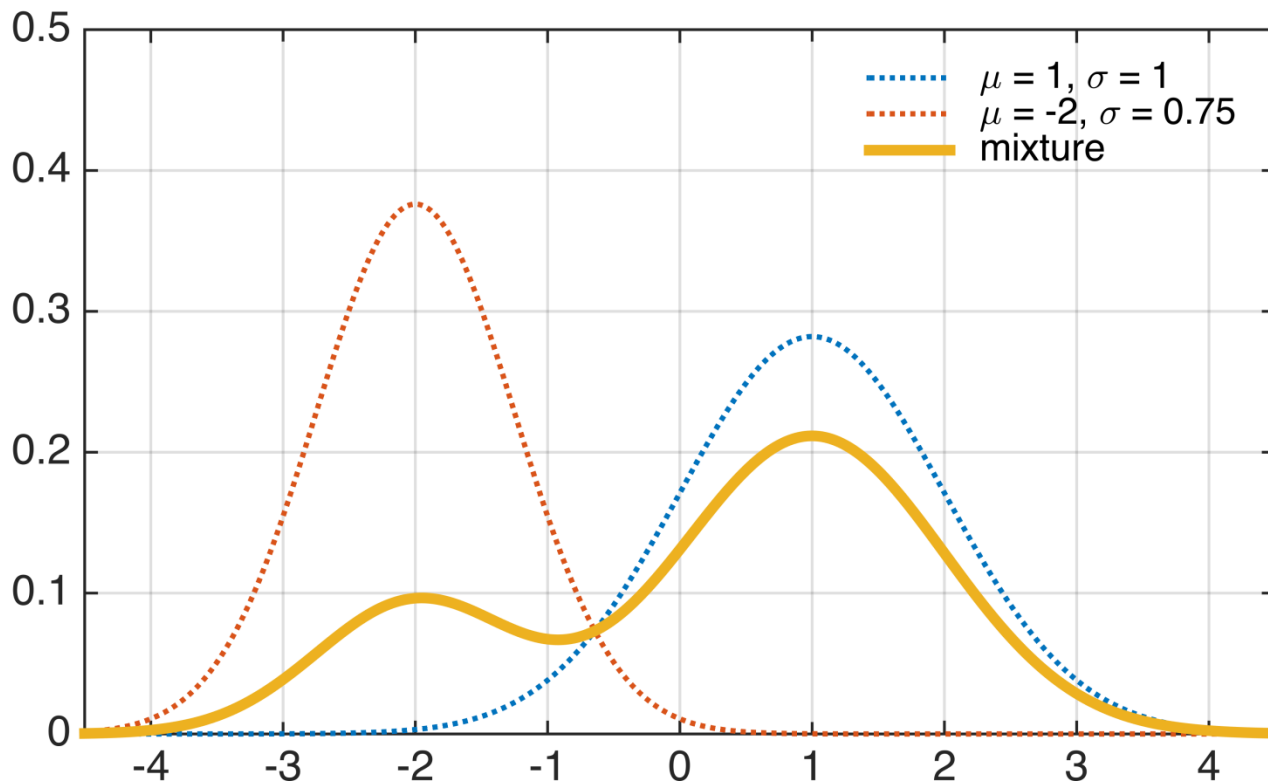
$$\sum_{i=1}^m w_i = 1$$

MIXTURES OF GAUSSIANS

Mixture of $m = 2$ Gaussian distributions:

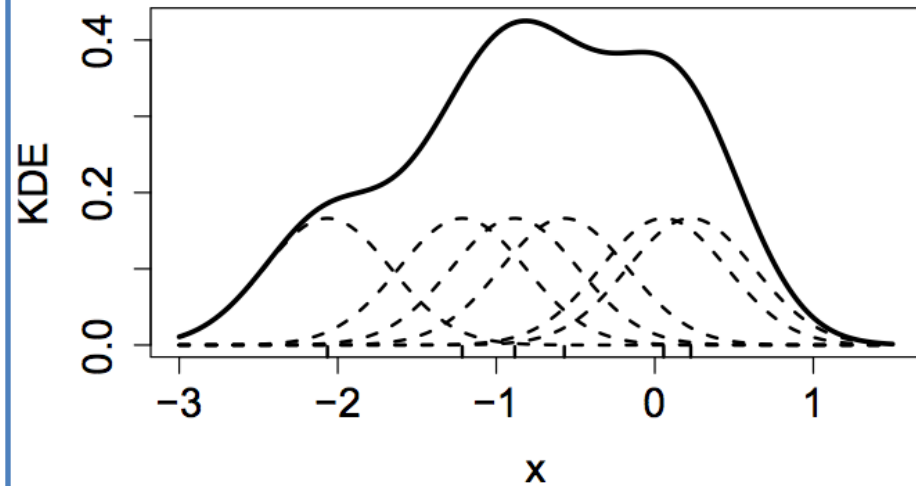
$$w_1 = 0.75, w_2 = 0.25$$

$$p(x) = \sum_{i=1}^m w_i p_i(x)$$

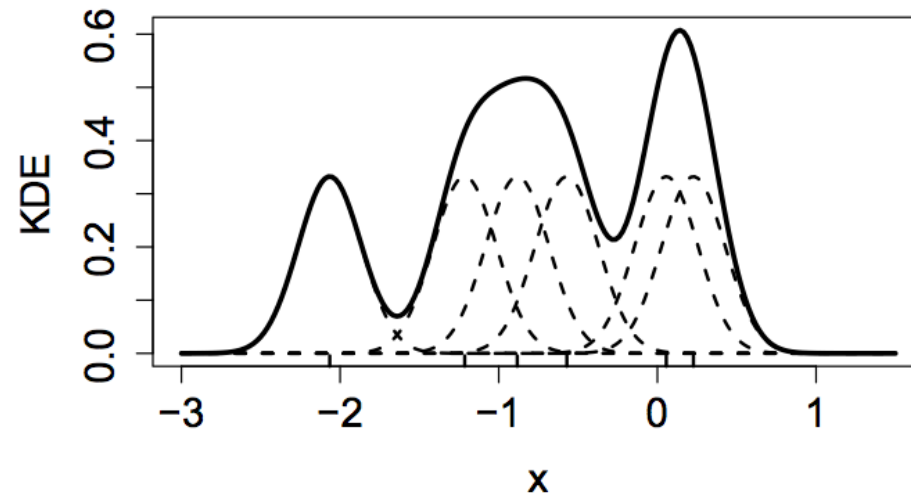


MIXTURES CAN PRODUCE COMPLEX DISTRIBUTIONS

b = 0.4



b = 0.2



* Image from <https://people.ucsc.edu/~ealdrich/Teaching/Econ114/LectureNotes/kde.html>

EXERCISE: SHOULD WE DISCRETIZE

- Notice that moving to continuous RVs puts more restrictions on the distributions we can define
 - mixtures provide more generality, but still restricted
- For discrete RVs, distributions are table of probabilities and so are highly flexible
 - we can define any possible distribution
- So, why not just discretize our variables?
- Example: imagine have d -dimensional random vectors, each entry in range $[0,1]$
- You discretize each dimension into 3 bins
- How many variables do you need for a PMF?
- What if had instead modelled it as a Gaussian? Or a Mixture?

EXERCISE: SAMPLE AVERAGE IS AN UNBIASED ESTIMATOR

Obtain instances x_1, \dots, x_n

What can we say about the sample average?

This sample is random, so we consider i.i.d. random variables X_1, \dots, X_n

Reflects that we could have seen a different set of instances x_i

$$\begin{aligned}\mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n X_i \right] &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i] \\ &= \frac{1}{n} \sum_{i=1}^n \mu \\ &= \mu\end{aligned}$$

For any one sample x_1, \dots, x_n , unlikely that $\frac{1}{n} \sum_{i=1}^n x_i = \mu$