



Computer
Science

CSC696H: Advanced Topics in Probabilistic Graphical Models

Probability and Statistics : Review

Prof. Jason Pacheco

Outline

- Random Variables and Discrete Probability
- Fundamental Rules of Probability
- Expected Value and Moments
- Continuous Probability
- Bayesian Inference

Outline

- **Random Variables and Discrete Probability**
- Fundamental Rules of Probability
- Expected Value and Moments
- Continuous Probability
- Bayesian Inference

Random Variables

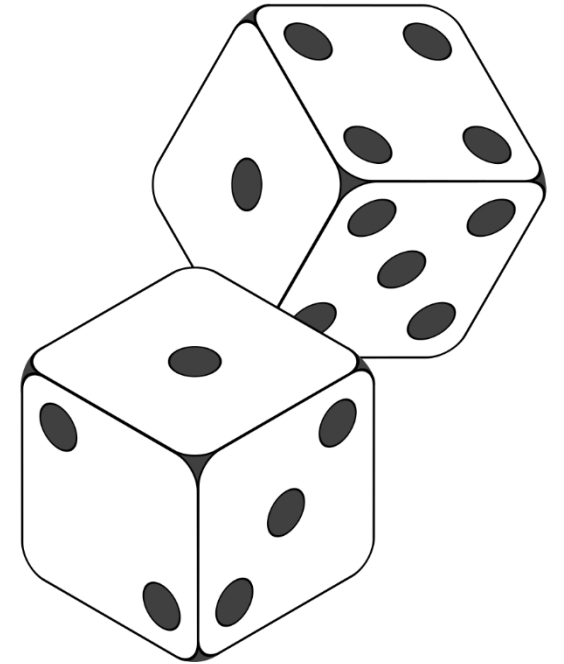
(Informally) A random variable is an unknown quantity whose value depends on the outcome of a random process

Example Roll 2 dice and let random variable X represent their sum. It takes values,

$$X \in \{2, 3, 4, \dots, 12\}$$

Example Flip a coin and let random variable Y represent the outcome,

$$Y \in \{\text{Heads}, \text{Tails}\}$$



Discrete vs. Continuous Probability

Discrete RVs take on a finite or countably infinite set of values

Continuous RVs take an uncountably infinite set of values

- Representing / interpreting / computing probabilities becomes more complicated in the continuous setting
- We will focus on discrete RVs for now...

Random Variables and Probability

Capitol letters represent
random variables

Lowercase letters are
realized *values*

$$X = x$$

$X = x$ is the **event** that X takes the value x

Example Let X be the random variable (RV) representing the sum of two dice with values,

$$X \in \{2, 3, 4, \dots, 12\}$$

$X=5$ is the *event* that the dice sum to 5.

Probability Mass Function

A function $p(X)$ is a **probability mass function (PMF)** of a discrete random variable if the following conditions hold:

(a) It is nonnegative for all values in the support,

$$p(X = x) \geq 0$$

(b) The sum over all values in the support is 1,

$$\sum_x p(X = x) = 1$$

Intuition Probability mass is conserved, just as in physical mass. Reducing probability mass of one event must increase probability mass of other events so that the definition holds...

Probability Mass Function

Example Let X be the outcome of a single fair die. It has the PMF,

$$p(X = x) = \frac{1}{6} \quad \text{for } x = 1, \dots, 6 \quad \text{Uniform Distribution}$$

Example We can often represent the PMF as a vector. Let S be an RV that is the *sum of two fair dice*. The PMF is then,

Observe that S does not follow a uniform distribution

$$p(S) = \begin{pmatrix} p(S = 2) \\ p(S = 3) \\ p(S = 4) \\ \vdots \\ p(S = 12) \end{pmatrix} = \begin{pmatrix} 1/36 \\ 1/18 \\ 1/2 \\ \vdots \\ 1/36 \end{pmatrix}$$

Functions of Random Variables

Any function $f(X)$ of a random variable X is also a random variable and it has a probability distribution

Example Let X_1 be an RV that represents the result of a fair die, and let X_2 be the result of another fair die. Then,

$$S = X_1 + X_2$$

Is an RV that is the *sum of two fair dice* with PMF $p(S)$.

NOTE Even if we know the PMF $p(X)$ and we know that the PMF $p(f(X))$ exists, it is not always easy to calculate!

PMF Notation

- We use $p(X)$ to refer to the probability mass *function* (i.e. a function of the RV X)
- We use $p(X=x)$ to refer to the probability of the *outcome* $X=x$ (also called an “event”)
- We will often use $p(x)$ as shorthand for $p(X=x)$

Outline

- Random Variables and Discrete Probability
- **Fundamental Rules of Probability**
- Expected Value and Moments
- Continuous Probability
- Bayesian Inference

Joint Probability

Definition Two (discrete) RVs X and Y have a *joint PMF* denoted by $p(X, Y)$ and the probability of the event $X=x$ and $Y=y$ denoted by $p(X = x, Y = y)$ where,

(a) It is nonnegative for all values in the support,

$$p(X = x, Y = y) \geq 0$$

(b) The sum over all values in the support is 1,

$$\sum_x \sum_y p(X = x, Y = y) = 1$$

Joint Probability

Let X and Y be *binary RVs*. We can represent the joint PMF $p(X, Y)$ as a 2x2 array (table):

		Y	
		0	1
X	0	0.04	0.36
	1	0.30	0.30

All values are nonnegative

Joint Probability

Let X and Y be *binary RVs*. We can represent the joint PMF $p(X,Y)$ as a 2x2 array (table):

		Y	
		0	1
X	0	0.04	0.36
	1	0.30	0.30

**The sum over all values is 1:
 $0.04 + 0.36 + 0.30 + 0.30 = 1$**

Joint Probability

Let X and Y be *binary RVs*. We can represent the joint PMF $p(X, Y)$ as a 2x2 array (table):

		Y	
		0	1
X	0	0.04	0.36
	1	0.30	0.30

$$P(X=1, Y=0) = 0.30$$

Fundamental Rules of Probability

Given two RVs X and Y the **conditional distribution** is:

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

Multiply both sides by $p(Y)$ to obtain the **probability chain rule**:

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

For N RVs X_1, X_2, \dots, X_N :

$$p(X_1, X_2, \dots, X_N) = p(X_1)p(X_2 | X_1) \dots p(X_N | X_{N-1}, \dots, X_1)$$

Chain rule valid
for any ordering

$$= p(X_1) \prod_{i=2}^N p(X_i | X_{i-1}, \dots, X_1)$$

Fundamental Rules of Probability

Law of total probability

$$p(Y) = \sum_x p(Y, X = x)$$

- $P(y)$ is a **marginal** distribution
- This is called **marginalization**

Proof

$$\begin{aligned} \sum_x p(Y, X = x) &= \sum_x p(Y) p(X = x | Y) && \text{(chain rule)} \\ &= p(Y) \sum_x p(X = x | Y) && \text{(distributive property)} \\ &= p(Y) && \text{(PMF sums to 1)} \end{aligned}$$

Generalization for conditionals:

$$p(Y | Z) = \sum_x p(Y, X = x | Z)$$

Tabular Method

Let X, Y be binary RVs with the joint probability table

For Binomial use K-by-K probability table.

		Y	
		y_1	y_2
X	x_1	0.04	0.36
	x_2	0.30	0.30

$P(y_1) = P(x_1, y_1) + P(x_2, y_1)$
 $P(y_2) = P(x_1, y_2) + P(x_2, y_2)$
[i.e., sum down columns]

$P(y)$

0.34 0.66

$P(y_1)$ $P(y_2)$

$P(x_1) = P(x_1, y_1) + P(x_1, y_2)$
 $P(x_2) = P(x_2, y_1) + P(x_2, y_2)$
[i.e., sum across rows]

Tabular Method

We don't care about event $Y=y_2$

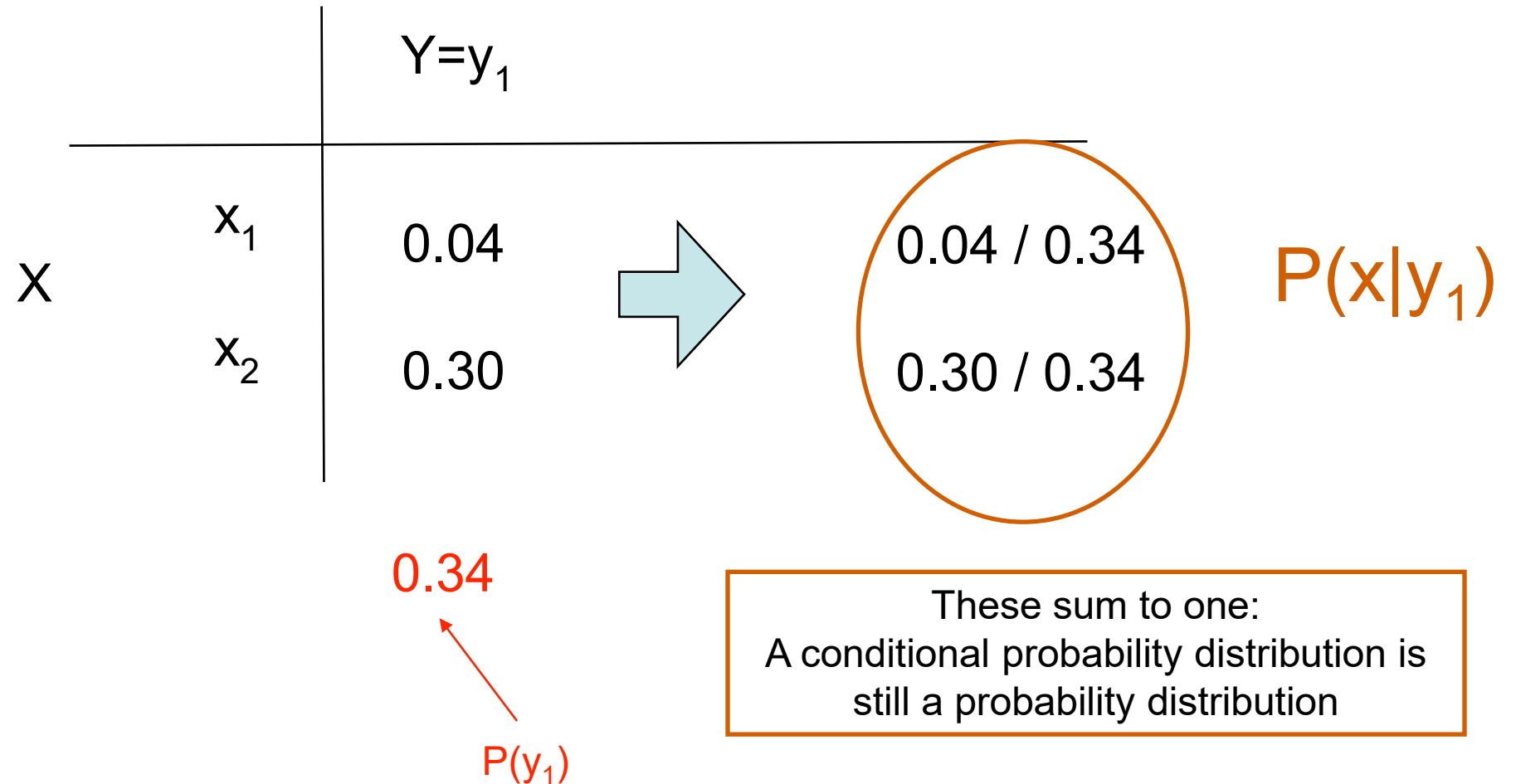
		Y	
		y_1	y_2
X	x_1	0.04	Censored!
	x_2	0.30	

$P(x|y_1)=?$

0.34

$P(y_1)$

Tabular Method



Summary

- A **random variable** is an unknown quantity whose value depends on the outcome a random process (informal definition)
- $X = x$ is an event with probability mass $p(X = x)$

- $p(X)$ is a **probability mass function (PMF)** satisfying

$$p(X = x) \geq 0 \qquad \sum_x p(X = x) = 1$$

- Some fundamental rules of probability:

- Conditional: $p(X | Y) = \frac{p(X, Y)}{p(Y)} = \frac{p(X, Y)}{\sum_x p(X=x, Y)}$
- Law of total probability: $p(Y) = \sum_x p(Y, X = x)$
- Probability chain rule: $p(X, Y) = p(Y)p(X | Y)$

Outline

- Random Variables and Discrete Probability
- Fundamental Rules of Probability
- **Expected Value and Moments**
- Continuous Probability
- Bayesian Inference

Moments of RVs

Definition The expectation of a discrete RV X , denoted by $\mathbf{E}[X]$, is:

$$\mathbf{E}[X] = \sum_x x p(X = x)$$

Summation over all values in domain of X

Example Let X be the sum of two fair dice, then:

$$\mathbf{E}[X] = \frac{1}{36} \cdot 2 + \frac{1}{18} \cdot 3 + \dots + \frac{1}{36} \cdot 12 = 7$$

Theorem (Linearity of Expectations) For any finite collection of discrete RVs X_1, X_2, \dots, X_N with finite expectations,

Corollary For any constant c
 $\mathbf{E}[cX] = c\mathbf{E}[X]$

$$\mathbf{E} \left[\sum_{i=1}^N X_i \right] = \sum_{i=1}^N \mathbf{E}[X_i]$$

E.g. for two RVs X and Y
 $\mathbf{E}[X + Y] = \mathbf{E}[X] + \mathbf{E}[Y]$

Moments of RVs

Law of Total Expectation *Let X and Y be discrete RVs with finite expectations, then:*

$$\mathbf{E}[X] = \mathbf{E}_Y[\mathbf{E}_X[X | Y]]$$

Proof

$$\begin{aligned} \mathbf{E}_Y[\mathbf{E}_X[X | Y]] &= \mathbf{E}_Y \left[\sum_x x \cdot p(x | Y) \right] \\ &= \sum_y \left[\sum_x x \cdot p(x | y) \right] \cdot p(y) && \text{(Definition of expectation)} \\ &= \sum_y \sum_x x \cdot p(x, y) && \text{(Probability chain rule)} \\ &= \sum_x x \sum_y p(x, y) && \text{(Linearity of expectations)} \\ &= \sum_x x \cdot p(x) = \mathbf{E}[X] && \text{(Law of total probability)} \end{aligned}$$

Moments of RVs

Definition The conditional expectation of a discrete RV X , given Y is:

$$\mathbf{E}[X \mid Y = y] = \sum_x x p(X = x \mid Y = y)$$

Example Roll two standard six-sided dice and let X be the result of the first die and let Y be the sum of both dice, then:

$$\begin{aligned} \mathbf{E}[X_1 \mid Y = 5] &= \sum_{x=1}^4 x p(X_1 = x \mid Y = 5) \\ &= \sum_{x=1}^4 x \frac{p(X_1 = x, Y = 5)}{p(Y = 5)} = \sum_{x=1}^4 x \frac{1/36}{4/36} = \frac{5}{2} \end{aligned}$$

Conditional expectation follows properties of expectation (linearity, etc.)

Moments of RVs

Definition The variance of a RV X is defined as,

$$\mathbf{Var}[X] = \mathbf{E}[(X - \mathbf{E}[X])^2] \quad \boxed{\text{(X-units)}^2}$$

The standard deviation is $\sigma[X] = \sqrt{\mathbf{Var}[X]}$. (X-units)

Lemma An equivalent form of variance is:

$$\mathbf{Var}[X] = \mathbf{E}[X^2] - (\mathbf{E}[X])^2$$

Proof Keep in mind that $E[X]$ is a constant,

$$\begin{aligned} \mathbf{E}[(X - \mathbf{E}[X])^2] &= \mathbf{E}[X^2 - 2X\mathbf{E}[X] + \mathbf{E}[X]^2] && \text{(Distributive property)} \\ &= \mathbf{E}[X^2] - 2\mathbf{E}[X]\mathbf{E}[X] + \mathbf{E}[X]^2 && \text{(Linearity of expectations)} \\ &= \mathbf{E}[X^2] - \mathbf{E}[X]^2 && \text{(Algebra)} \end{aligned}$$

Moments of RVs

Definition The covariance of two RVs X and Y is defined as,

$$\mathbf{Cov}(X, Y) = \mathbf{E}[(X - \mathbf{E}[X])(Y - \mathbf{E}[Y])]$$

Lemma For any two RVs X and Y ,

$$\mathbf{Var}[X + Y] = \mathbf{Var}[X] + \mathbf{Var}[Y] + 2\mathbf{Cov}(X, Y)$$

e.g. variance is not a linear operator.

Proof $\mathbf{Var}[X + Y] = \mathbf{E}[(X + Y - \mathbf{E}[X + Y])^2]$

(Linearity of expectation) $= \mathbf{E}[(X + Y - \mathbf{E}[X] - \mathbf{E}[Y])^2]$

(Distributive property) $= \mathbf{E}[(X - \mathbf{E}[X])^2 + (Y - \mathbf{E}[Y])^2 + 2(X - \mathbf{E}[X])(Y - \mathbf{E}[Y])]$

(Linearity of expectation) $= \mathbf{E}[(X - \mathbf{E}[X])^2] + \mathbf{E}[(Y - \mathbf{E}[Y])^2] + 2\mathbf{E}[(X - \mathbf{E}[X])(Y - \mathbf{E}[Y])]$

(Definition of Var / Cov) $= \mathbf{Var}[X] + \mathbf{Var}[Y] + 2\mathbf{Cov}(X, Y)$

Summary

Moments and Expected Value

- Expected value of a discrete RV:

$$\mathbf{E}[X] = \sum_x x p(X = x)$$

- Expectation is a linear operator

$$\mathbf{E} \left[\sum_{i=1}^N X_i \right] = \sum_{i=1}^N \mathbf{E}[X_i]$$

- Variance of a RV:

$$\mathbf{Var}[X] = \mathbf{E}[(X - \mathbf{E}[X])^2]$$

- Variance is **not** a linear operator

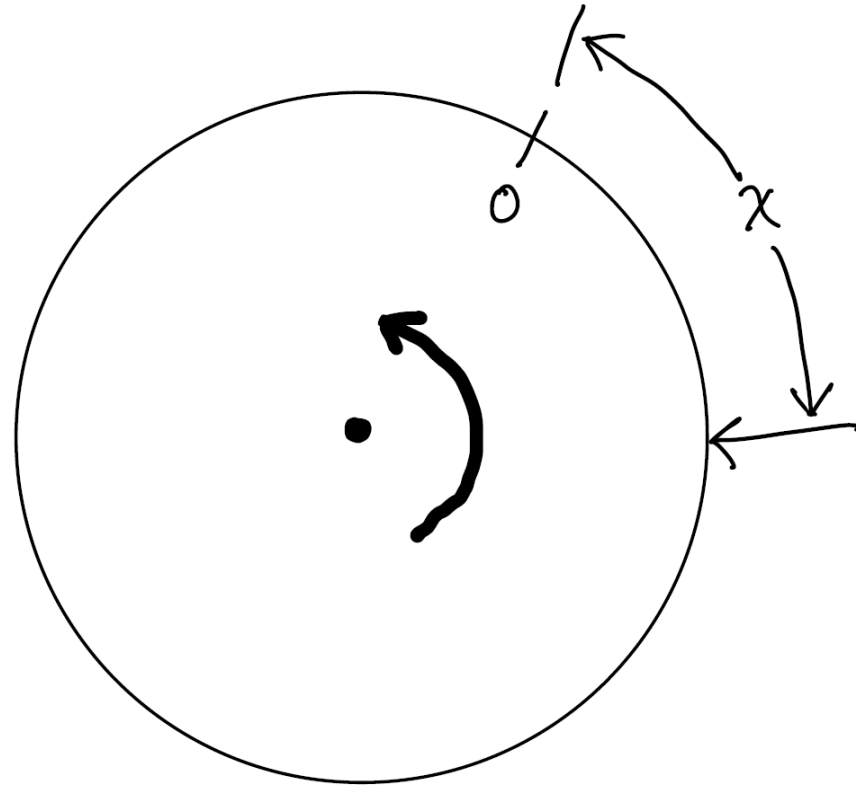
$$\mathbf{Var}[X + Y] = \mathbf{Var}[X] + \mathbf{Var}[Y] + 2\mathbf{Cov}(X, Y)$$

Outline

- Random Variables and Discrete Probability
- Fundamental Rules of Probability
- Expected Value and Moments
- **Continuous Probability**
- Bayesian Inference

Continuous Probability

Experiment Spin continuous wheel and measure X displacement from 0



Question Assuming uniform probability, what is $p(X = x)$?

Continuous Probability

➤ Let $p(X = x) = \pi$ be the probability of any single outcome

➤ Let $S(k)$ be set of any k *distinct* points in $[0, 1)$ then,

$$P(x \in S(k)) = k\pi$$

➤ Since $0 < P(x \in S(k)) < 1$ by axioms of probability, $k\pi < 1$ for any k

➤ Therefore: $\pi = 0$ and $P(x \in S(k)) = p(X = x) = 0$

Continuous Probability

- We have a well-defined event that x takes a value in set $x \in S(k)$
- Clearly this event can happen... i.e. **it is possible**
- But we have shown it has zero probability of occurring,

$$P(x \in S(k)) = 0$$

- By the axioms of probability, the probability that it **doesn't happen** is,

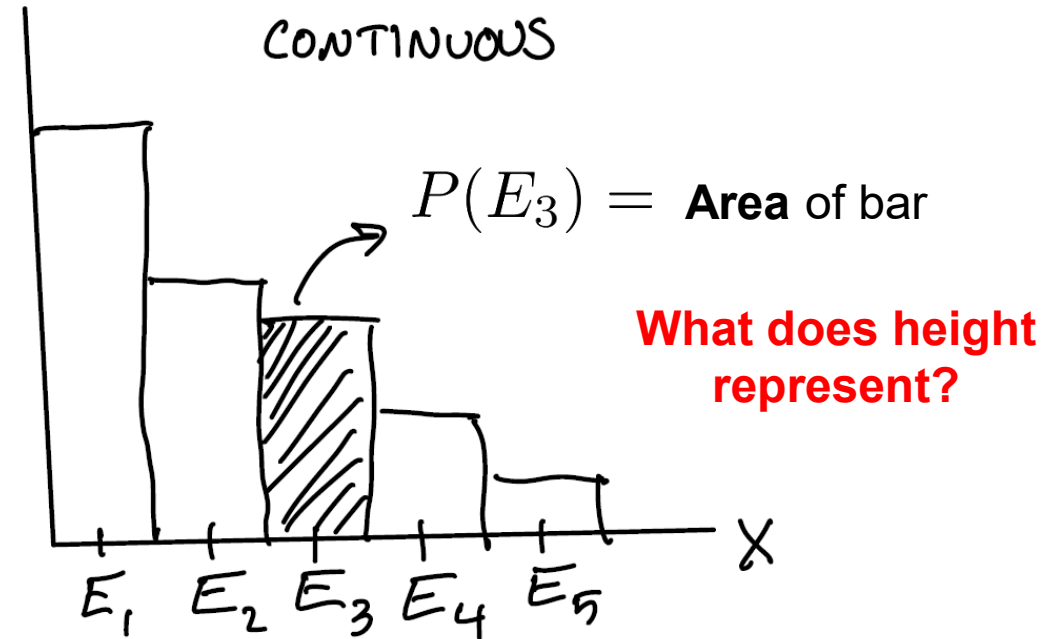
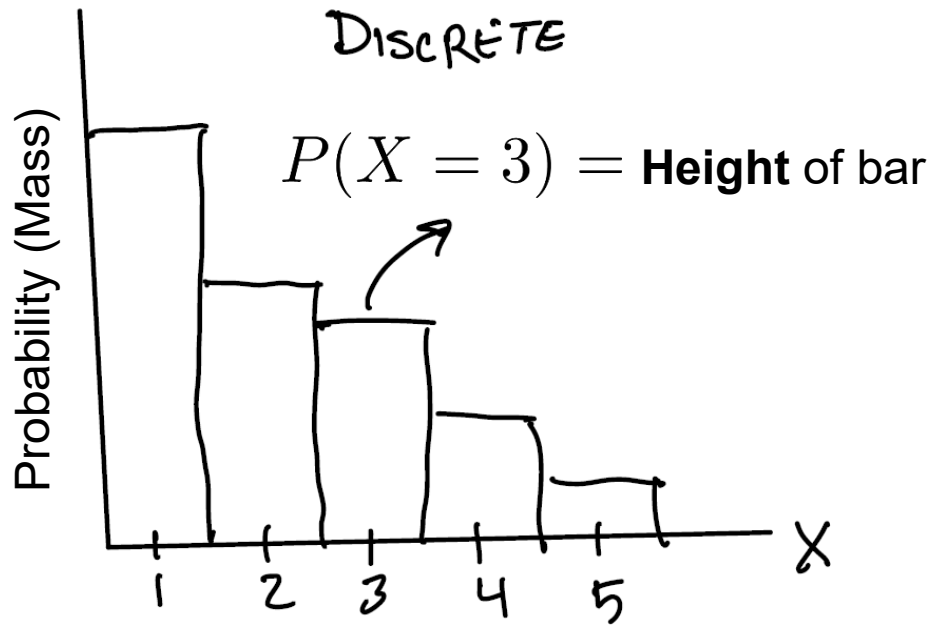
$$P(x \notin S(k)) = 1 - P(x \in S(k)) = 1$$

We seem to have a paradox!

Solution Rethink how we interpret probability in continuous setting

- Define events as *intervals* instead of discrete values
- Assign probability to those intervals

Continuous Probability



Probability

Δx

Height = $\frac{\text{Probability}}{\Delta x}$

Height represents *probability per unit* in the x-direction

We call this a **probability density** (as opposed to probability mass)

Continuous Probability

➤ We denote the **probability density function** (PDF) as, $p(X)$

➤ An event E corresponds to an *interval* $a \leq X < b$

➤ The probability of an interval is given by the *area under the PDF*,

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

➤ Specific outcomes have zero probability $P(X = 0) = P(x \leq X < x) = 0$

➤ But may have nonzero *probability density* $p(X = x)$

Continuous Probability Measures

Definition The cumulative distribution function (CDF) of a real-valued continuous RV X is the function given by,

$$P(x) = P(X \leq x)$$

Different ways to represent probability of interval, CDF is just a convention.

➤ Can easily measure probability of closed intervals,

$$P(a \leq X < b) = P(b) - P(a)$$

➤ If X is *absolutely continuous* (i.e. differentiable) then,

Fundamental Theorem of Calculus

$$p(x) = \frac{dP(x)}{dx} \quad \text{and} \quad P(t) = \int_{-\infty}^t p(x) dx$$

Where $p(x)$ is the *probability density function* (PDF)

Fundamental Laws of Probability (Continuous)

Most definitions for discrete RVs hold, replacing PMF with PDF/CDF...

Probability chain rule,

Shorthand: $P(x) = P(X \leq x)$

$$p(x, y) = p(x)p(y | x) \quad \text{and} \quad P(x, y) = P(x)P(y | x)$$

...and by replacing summation with integration...

Law of Total Probability for continuous distributions,

$$p(x) = \int_{\mathcal{Y}} p(x, y) dy$$

Expectation of a continuous random variable,

$$\mathbf{E}[X] = \int_{\mathcal{X}} x \cdot p(x) dx$$

Outline

- Random Variables and Discrete Probability
- Fundamental Rules of Probability
- Expected Value and Moments
- Continuous Probability
- **Bayesian Inference**

What is Probability?

What does it mean that the probability of heads is $\frac{1}{2}$?



Two schools of thought...

Frequentist Perspective

Proportion of successes (heads) in repeated trials (coin tosses)

Bayesian Perspective

Belief of outcomes based on assumptions about nature and the physics of coin flips

Neither is better/worse, but we can compare interpretations...

Frequentist & Bayesian Modeling

We will use the following notation throughout:

θ - Unknown (e.g. coin bias)

y - Data

Frequentist

(Conditional Model)

$$p(y; \theta)$$

- θ is a non-random unknown parameter
- $p(y; \theta)$ is the *sampling / data generating distribution*

Bayesian

(Generative Model)

Prior Belief $\rightarrow p(\theta)p(y | \theta) \leftarrow$ Likelihood

- θ is a random variable (latent)
- Requires specifying $p(\theta)$ the prior belief

Bayes' Rule

Posterior represents all uncertainty after observing data...

The diagram shows the Bayes' Rule equation with four labels and arrows pointing to the corresponding parts of the equation:

- prior** probability: points to $p(\theta)$
- likelihood** function for the parameters: points to $p(y | \theta)$
- posterior** probability: points to $p(\theta | y)$
- marginal likelihood** or: **evidence** or: **partition function** or: **normalizer**: points to $p(y)$

$$p(\theta | y) = \frac{p(\theta)p(y | \theta)}{p(y)}$$

Bayes' Rule : Marginal Likelihood

$$p(\theta | y) = \frac{p(\theta)p(y | \theta)}{p(y)} \propto \underbrace{p(\theta)p(y | \theta)}$$

Often hard to calculate

Often know this (the model)

Marginal likelihood integrates (marginalizes) over unknown θ :

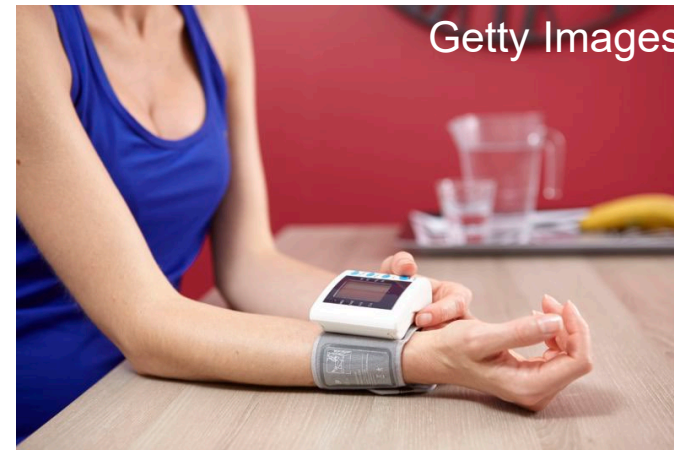
$$p(y) = \int p(\theta)p(y | \theta) d\theta$$

Marginal likelihood is less problematic in discrete models (not always)

This integral often lacks a closed form and cannot be computed...

Bayesian Inference Example

About **29%** of American adults have high blood pressure (BP). Home test has **30% false positive** rate and **no false negative error**.



A recent home test states that you have high BP. Should you start medication?

An Assessment of the Accuracy of Home Blood Pressure Monitors When Used in Device Owners

Jennifer S. Ringrose,¹ Gina Polley,¹ Donna McLean,²⁻⁴ Ann Thompson,^{1,5} Fraulein Morales,¹ and Raj Padwal^{1,4,6}

Bayesian Inference Example

About **29%** of American adults have high blood pressure (BP). Home test has **30% false positive** rate and **no false negative error**.



- Latent quantity of interest is hypertension: $\theta \in \{true, false\}$
- Measurement of hypertension: $y \in \{true, false\}$
- Prior: $p(\theta = true) = 0.29$
- Likelihood: $p(y = true \mid \theta = false) = 0.30$
 $p(y = true \mid \theta = true) = 1.00$

Bayesian Inference Example

About **29%** of American adults have high blood pressure (BP). Home test has **30% false positive** rate and **no false negative error**.



Suppose we get a positive measurement, then posterior is:

$$\begin{aligned} p(\theta = \text{true} \mid y = \text{true}) &= \frac{p(\theta = \text{true})p(y = \text{true} \mid \theta = \text{true})}{p(y = \text{true})} \\ &= \frac{0.29 * 1.00}{0.29 * 1.00 + 0.71 * 0.30} \approx 0.58 \end{aligned}$$

What conclusions can be drawn from this calculation?

Aside : Proportionality

Recall PMF / PDF must sum / integrate to 1,

$$\begin{array}{cc} \text{PMF} & \text{PDF} \\ \sum_x p(x) = 1 & \int p(x) dx = 1 \end{array}$$

May only know distribution constant that does not depend on RV x ,

$$\int \tilde{p}(x) dx = \mathcal{Z} \quad \text{so} \quad p(x) \propto \tilde{p}(x)$$

Properly normalized distribution by dividing our normalization constant:

$$\int p(x) dx = \int \frac{1}{\mathcal{Z}} \tilde{p}(x) dx = \frac{1}{\int \tilde{p}(x) dx} \int \tilde{p}(x) dx = 1$$

Aside : Proportionality

Example Let X be a Bernoulli RV (coinflip) with probabilities *proportional to*:

$$\tilde{p}(X = 0) = 0.5$$

$$\tilde{p}(X = 1) = 1.5$$

Greater than 1, but
It is an *unnormalized*
probability

Compute normalization constant,

$$\mathcal{Z} = \tilde{p}(X = 0) + \tilde{p}(X = 1) = 2.0$$

Normalize probability distribution,

$$p(X) = \frac{1}{\mathcal{Z}} \tilde{p}(X) = \begin{pmatrix} 1/4 \\ 3/4 \end{pmatrix}$$

Sums to 1

Frequentist vs. Bayesian Inference

We have data X_1, \dots, X_N and want to infer unknown parameter θ

Frequentist Inference

The data *uniquely determines* θ , e.g. by the likelihood:

Not a distribution on parameter $p(X_1, \dots, X_N; \theta)$ **How well it explains the data**

Bayesian Inference

The data *updates our belief* about θ , which is random:

$$p(\theta \mid X_1, \dots, X_N) \propto p(\theta \mid X_1, \dots, X_{N-1})p(X_N \mid \theta)$$

Our belief changes with more data

Minimum Mean Squared Error (MMSE)

Posterior mean minimizes squared error,

$$\hat{\theta}^{\text{MMSE}} = \arg \min \mathbb{E}[(\hat{\theta} - \theta)^2 \mid y] = E[\theta \mid y]$$

- Minimizes error conditioned on observed data
- MMSE is an **unbiased estimator**
- MMSE is **asymptotically unbiased** and **asymptotically normal**,

$$\sqrt{N}(\hat{\theta}^{\text{MMSE}} - \theta) \rightarrow \mathcal{N}(0, \sigma^2)$$

Bayes Estimators

Minimizes expected loss function,

$$\hat{\theta} = \arg \min_{\hat{\theta}} \mathbf{E} \left[L(\theta, \hat{\theta}) \mid y \right]$$

Expected loss referred to as *Bayes risk*.

MMSE minimizes squared-error loss $L(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$

Minimum absolute error (MAE) is posterior *median*,

$$\arg \min \mathbf{E}[|\hat{\theta} - \theta| \mid y] = \text{median}(\theta \mid y)$$

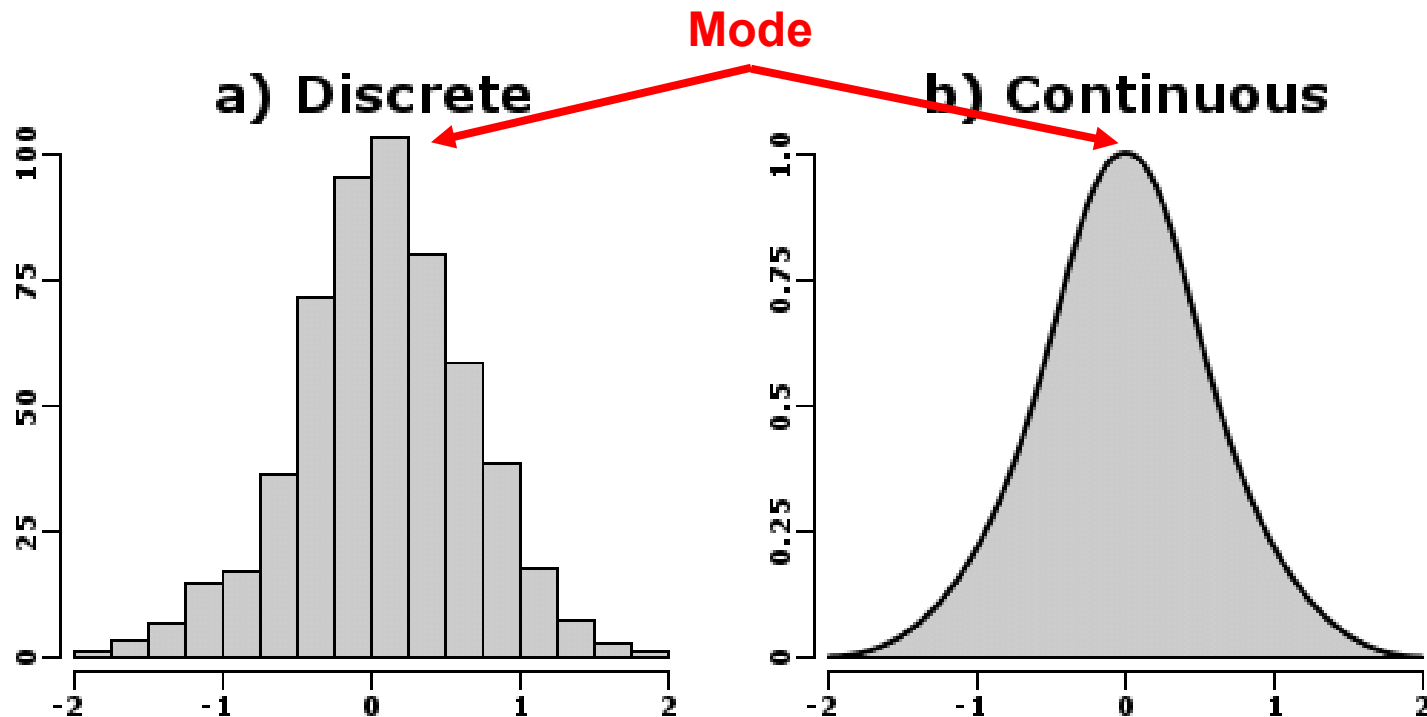
Note: Same answer for linear function: $L(\theta, \hat{\theta}) = c|\hat{\theta} - \theta|$

Maximum a Posteriori (MAP)

Very common to produce maximum probability estimates,

$$\hat{\theta}^{\text{MAP}} = \arg \max p(\theta | y)$$

*MAP is the **mode** (highest probability outcome) of the posterior*



Maximum a Posteriori (MAP)

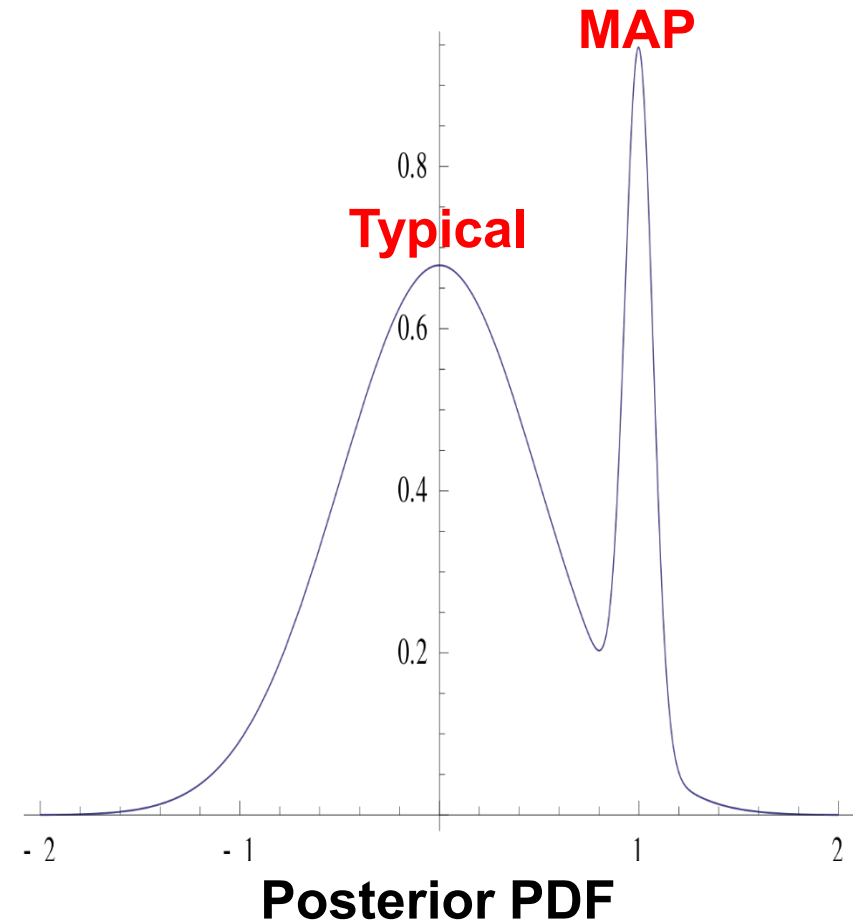
MAP (mode) may not be representative of typical outcomes

Also, not a Bayes estimator (unless discrete),

$$\lim_{c \rightarrow 0} L(\theta, \hat{\theta}) = \begin{cases} 0, & \text{if } |\hat{\theta} - \theta| < c \\ 1, & \text{otherwise} \end{cases}$$

Degenerate loss function

Despite its issues, MAP is frequently used in “Bayesian” inference and estimation



Example: Beta-Bernoulli MAP

Let $X_1, \dots, X_N \sim \text{Bernoulli}(\pi)$ and $\pi \sim \text{Beta}(\alpha, \beta)$ then posterior is,

$$p(\pi | X_1^N) = \text{Beta}(\alpha + \underbrace{\text{number of heads}}_{N_H}, \beta + \text{number of tails})$$

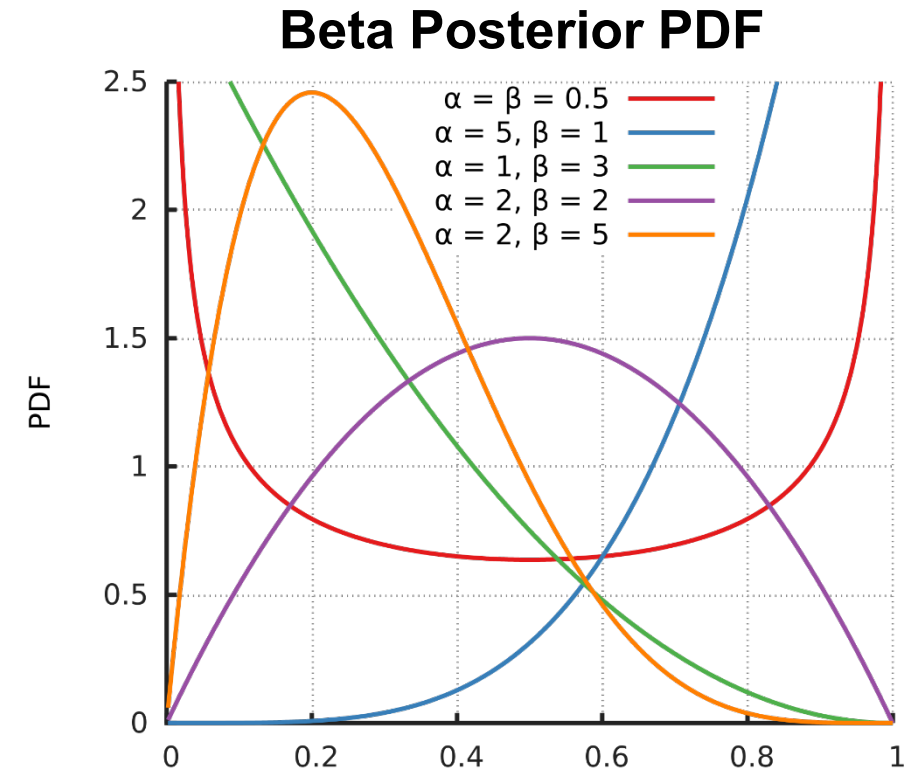
Highest probability (mode) of Beta given by,

Take derivative,
set to zero, solve.

$$\hat{\pi}^{\text{MAP}} = \frac{\alpha + N_H - 1}{\alpha + \beta + N - 2}$$

Beta distribution is not always convex!

- MAP is any value for $\alpha = \beta = 1$
- Two modes (bimodal) for $\alpha, \beta < 1$



Maximum a Posteriori (MAP)

Equivalent to maximizing joint probability,

$$\arg \max_{\theta} p(\theta | y) = \arg \max_{\theta} \frac{p(\theta, y)}{p(y)} = \arg \max_{\theta} p(\theta, y)$$

Constant

For iid y_1, \dots, y_N solve in log-domain (like *maximum likelihood est.*),

$$\hat{\theta}^{\text{MAP}} = \arg \max_{\theta} \log p(\theta, y_1, \dots, y_N) = \underbrace{\sum_i \log p(y_i | \theta)}_{\text{Log-Likelihood (how well it fits data)}} + \underbrace{\log p(\theta)}_{\text{Log-Prior (how well it agrees with prior)}}$$

Intuition MAP is like MLE but with a “penalty” term (log-prior)

Summary

- Bayesian statistics interprets probability differently than classical stats
 - Frequentist: Probability \rightarrow Long run odds in repeated trials
 - Bayesian: Probability \rightarrow Belief of outcome that captures all uncertainty
- Bayesian models treat unknown parameter as random, with a prior
- Bayesian inference via the *posterior distribution* using Bayes' rule

$$p(\theta | y) = \frac{p(\theta)p(y | \theta)}{p(y)}$$

- Bayesian estimators minimize expected risk (e.g. MMSE)
- Maximum a posteriori (MAP) estimate maximizes posterior probability

