

Neural Networks

A microscopic image of neurons, showing several cell bodies (soma) in shades of blue and purple, interconnected by a dense network of thin, cyan-colored axons. The background is dark, making the glowing structures stand out.

image: Ardy Rahman, UCI Research

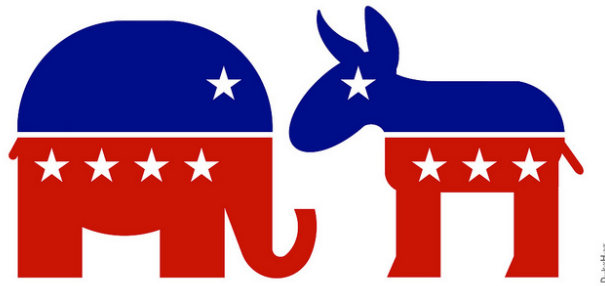
Will Monroe
August 16, 2017

with materials by
Mehran Sahami
and Chris Piech

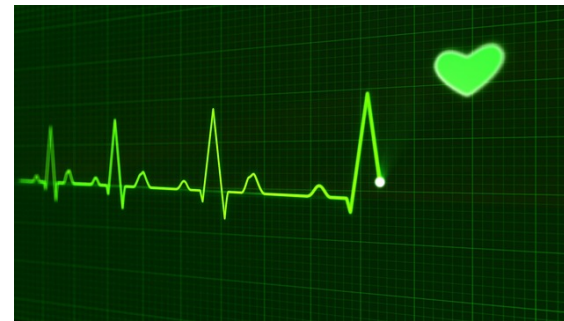
Announcement: Problem Set #6

Due **today!**

That's all, folks!



Congressional voting



Heart disease
diagnosis

Announcements: Final exam



This Saturday, August 19, 12:15-3:15pm
in NVIDIA Auditorium
(pending maintenance)

Two pages (both sides) of notes

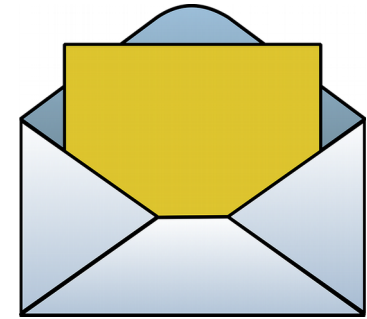
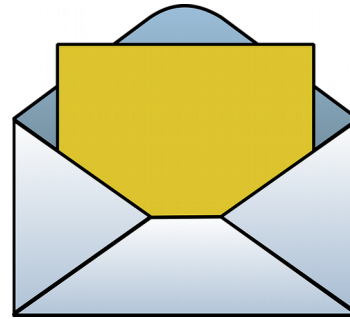
All material in the class through Monday

Review session:

Today after lecture, 2:30-3:20 in **Huang 18**

Two envelopes: A resolution

“I’m trying to think: how likely is it that you would have put \$40 in an envelope?”



$Y = y$: amount in envelope chosen

$$E[W|Y = y, \text{stay}] = y$$

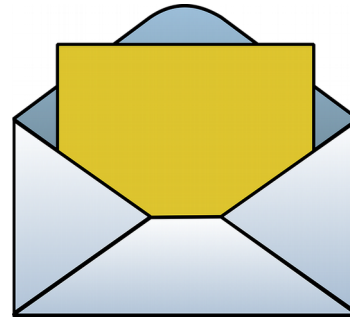
$$E[W|Y = y, \text{switch}] = \frac{y}{2} P(X = \frac{y}{2}|Y = y) + 2y P(X = y|Y = y)$$

not necessarily 0.5!

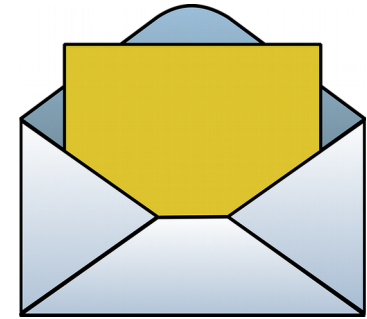
$$P(X = y|Y = y) = \frac{P(X = y)}{P(X = y) + P(X = y/2)}$$

Two envelopes: A resolution

“I’m trying to think: how likely is it that you would have put \$40 in an envelope?”



\$X



\$2X

$Y = y$: amount in envelope chosen

$$E[W | Y = y, \text{stay}] = y$$

$$E[W | Y = y, \text{switch}] = \frac{y}{2} P(X = \frac{y}{2} | Y = y) + 2y P(X = y | Y = y)$$

not necessarily 0.5!

$$P(X = y | Y = y) = \frac{P(X = y)}{P(X = y) + P(X = y/2)}$$

prior: if all equally likely, then this will be 0.5

$$P(X = y) = C?$$

$$\sum_y P(X = y) = \sum_y C = 1$$

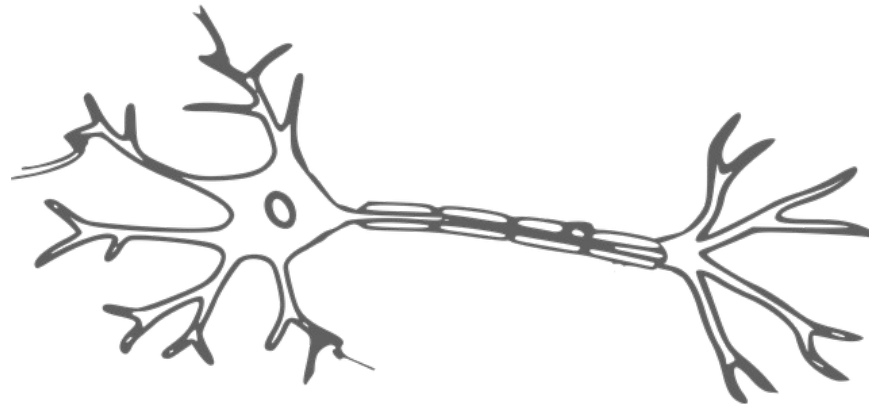
$$\infty \cdot C = 1???$$

Logistic regression

A classification algorithm using the assumption that **log odds** are a linear function of the features.



$$\hat{y} = \frac{1}{1 + e^{-\vec{\theta}^T \vec{x}}}$$



Review: The logistic function

$$z = \log o_f$$

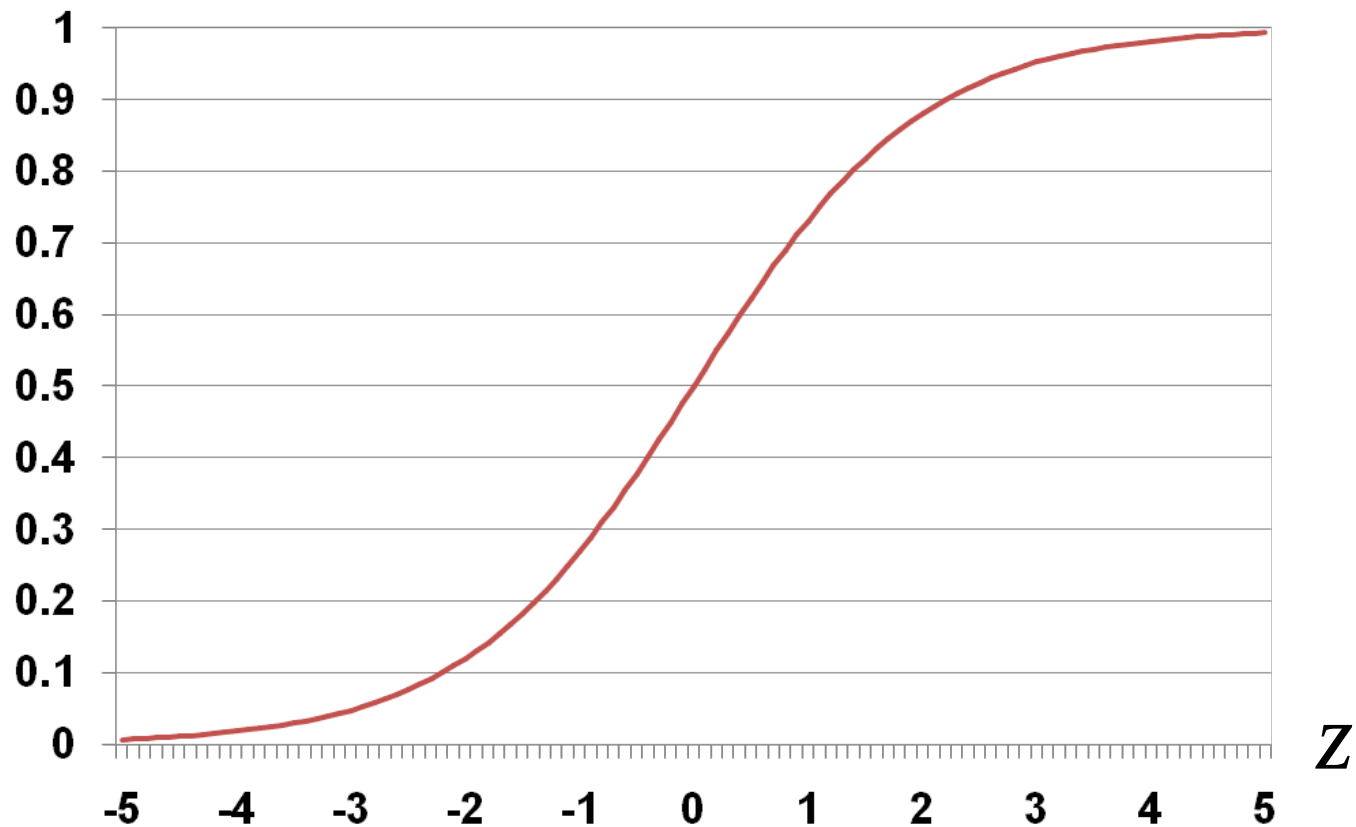
$$p = \frac{o_f}{o_f + 1} = \frac{1}{1 + \frac{1}{o_f}}$$

$$= \frac{1}{1 + e^{-\log(o_f)}}$$

$$= \frac{1}{1 + e^{-z}}$$

$$= \sigma(z)$$

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



Review: Logistic regression assumption

$$P(Y = 1 | \vec{X} = \vec{x}) = \sigma(\vec{\theta}^T \vec{x}) = \frac{1}{1 + e^{-\vec{\theta}^T \vec{x}}}$$

or in other words:

$$p = \sigma(z)$$

$$z = \log o_f$$

$$\vec{\theta}^T \vec{x} = \log o_f(Y = 1 | \vec{X} = \vec{x})$$

$$\vec{\theta}^T \vec{x} = \vec{\theta} \cdot \vec{x} = \theta_0 \cdot 1 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_m x_m$$

$$= \sum_{i=0}^m \theta_i x_i$$

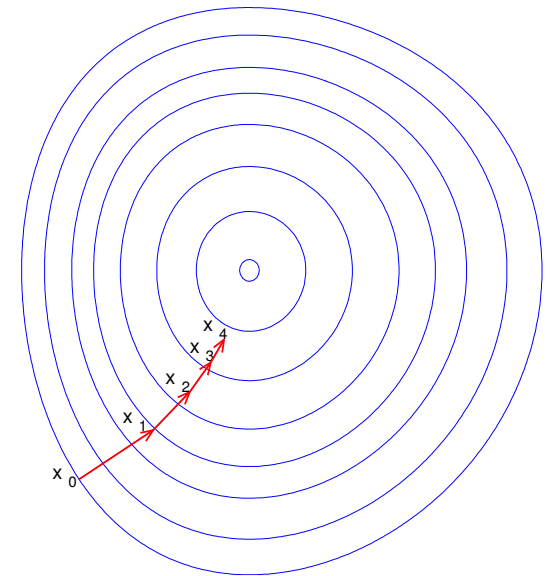
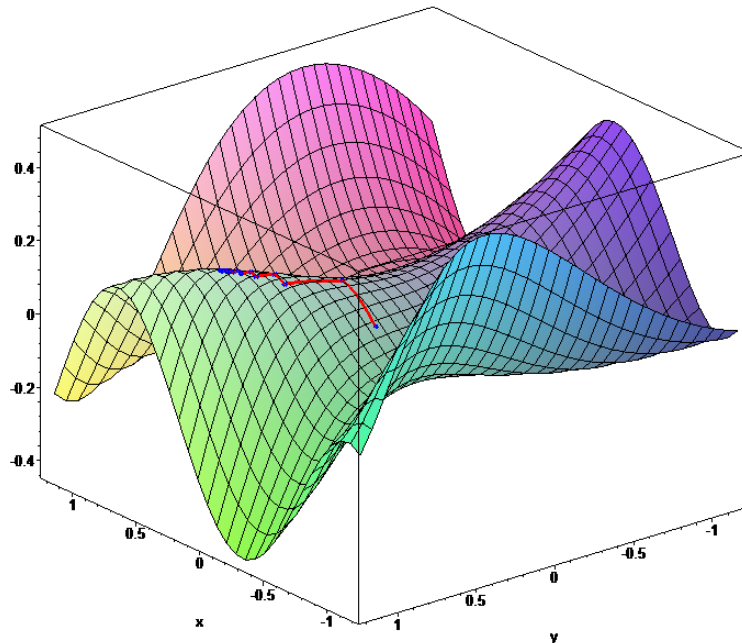
$$(x_0 = 1)$$

Review: Gradient ascent

An algorithm for computing an **arg max** by taking small steps **uphill** (i.e., in the direction of the **gradient** of the function).



$$\vec{\theta} \leftarrow \vec{\theta} + \eta \cdot \nabla_{\vec{\theta}} f(\vec{\theta})$$



Review: Logistic regression algorithm

initialize: $\theta = [0, 0, \dots, 0]$ (m elements)

repeat many times:

 gradient = $[0, 0, \dots, 0]$ (m elements)

for each training example $(\vec{x}^{(i)}, y^{(i)})$:

for j = 0 **to** m:

 gradient[j] += $[y^{(i)} - \sigma(\vec{\theta}^T \vec{x}^{(i)})] \vec{x}_j^{(i)}$

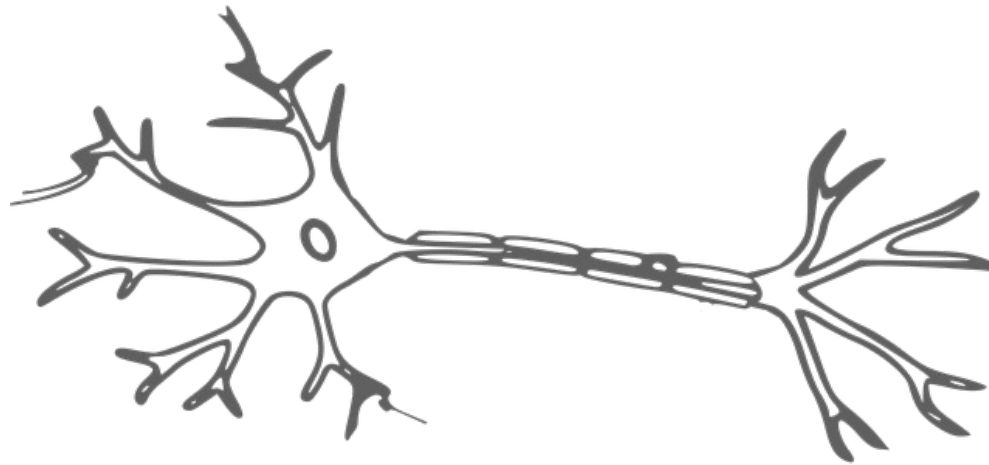
for j = 0 **to** m:

$\theta[j] += \eta * \text{gradient}[j]$

return θ

Your brain on logistic regression

$$p = \sigma(z)$$
$$z = \vec{\theta}^T \vec{x}^{(i)}$$



dendrites:

take a weighted sum
of incoming stimuli
with electric potential

axon:

carries outgoing
pulse if potential
exceeds a threshold

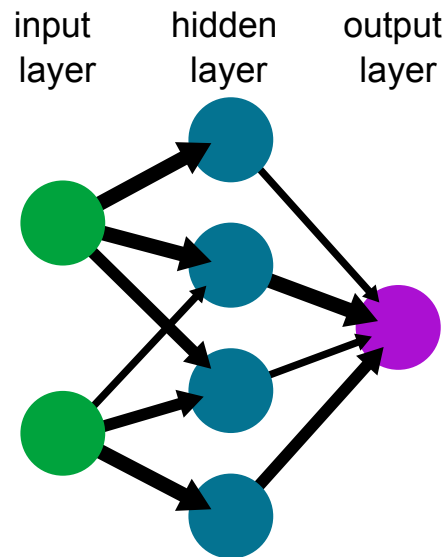
Caution: Just a (greatly simplified)
model! All models are wrong—but
some are useful...

Feedforward neural network

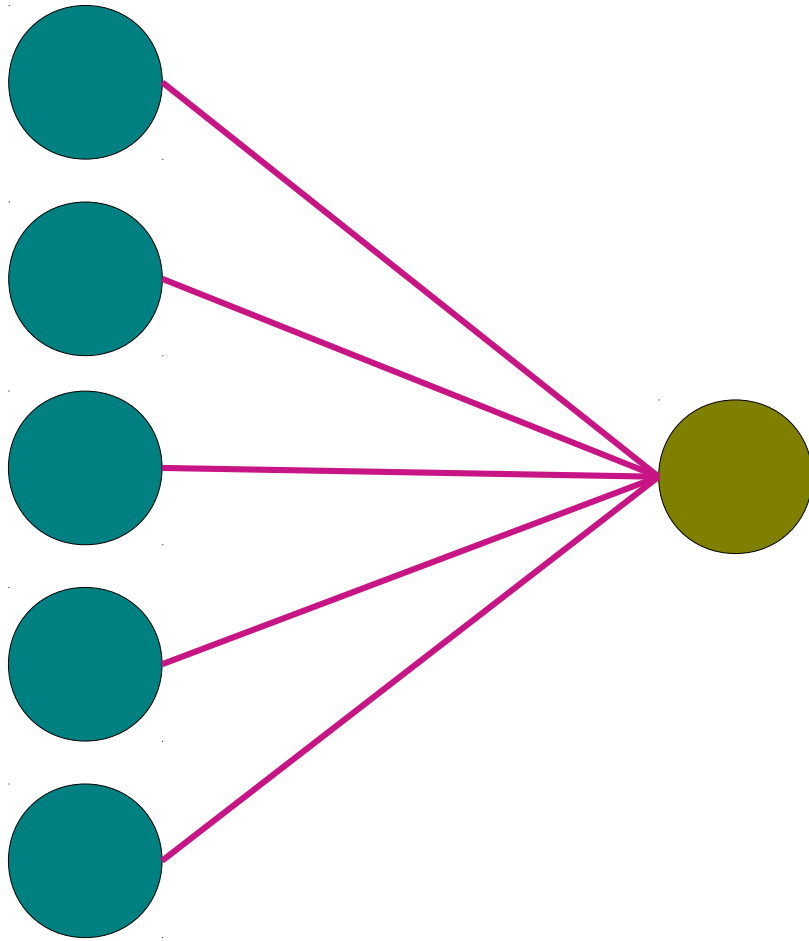
An algorithm for classification or regression that uses **layers of logistic regressions** to discover its own features.



$$\hat{y} = \sigma \left(\theta^{(\hat{y})} \sigma \left(\theta^{(h)} \vec{x} \right) \right)$$



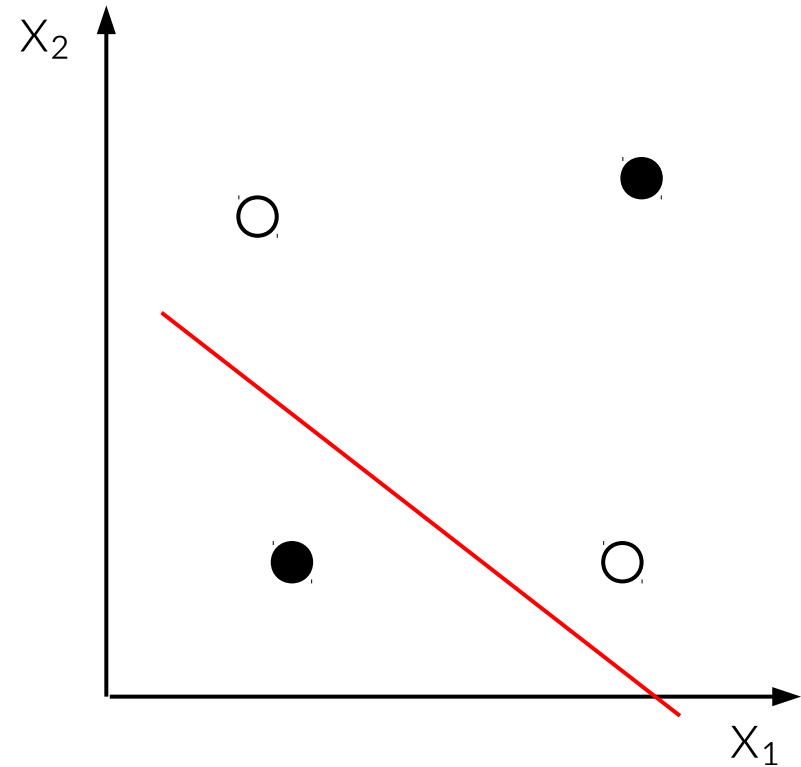
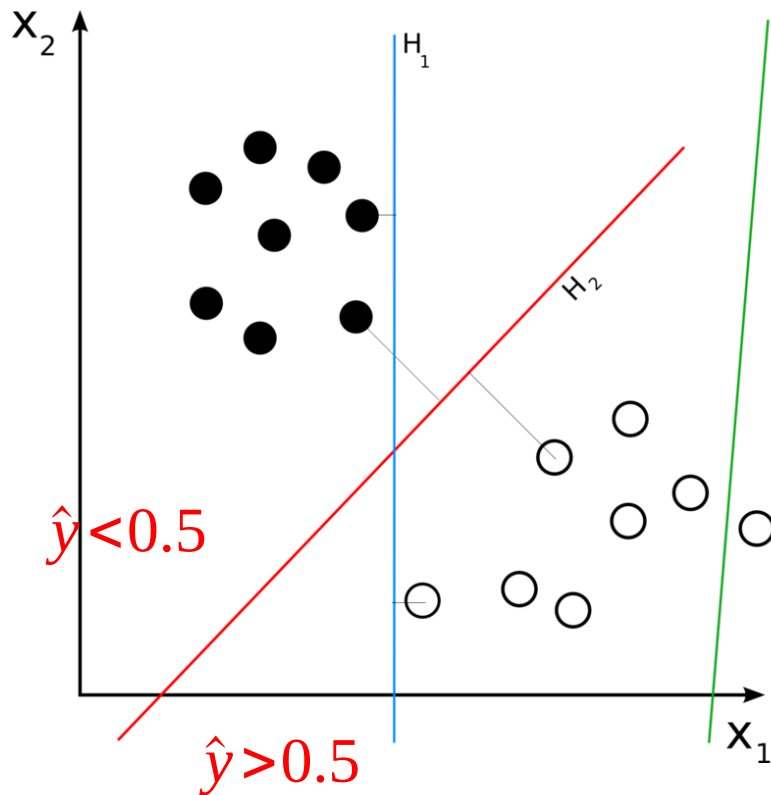
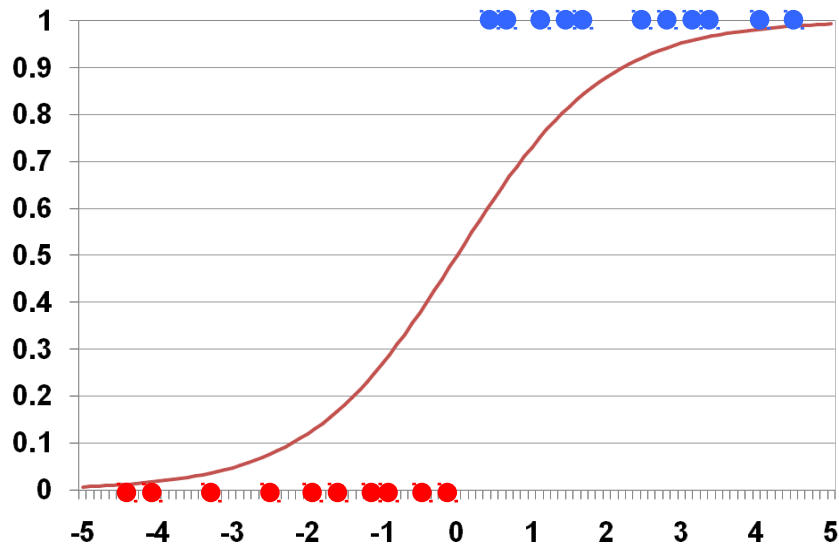
A cartoon of logistic regression



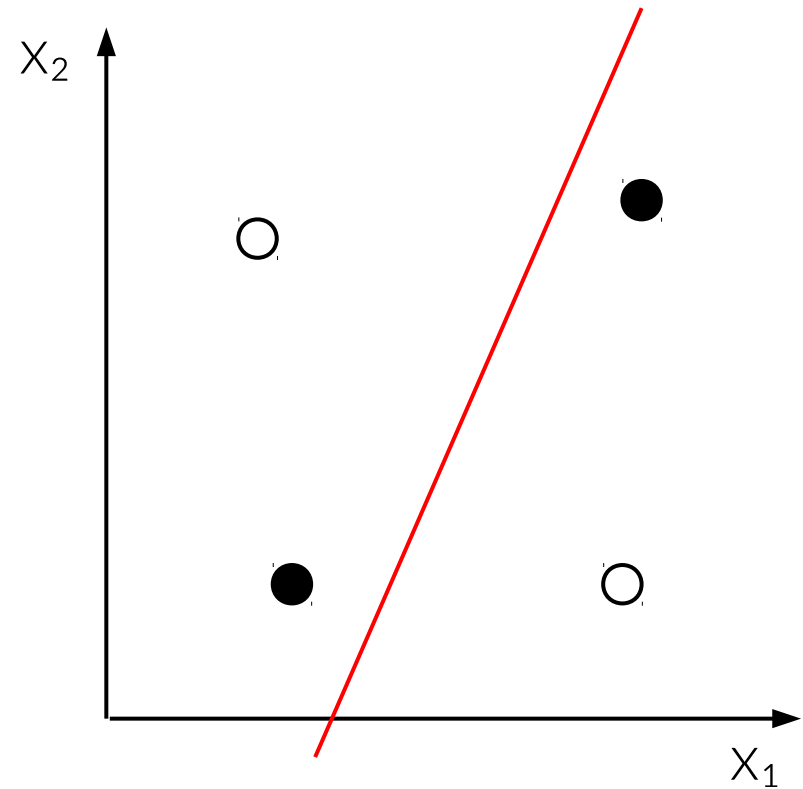
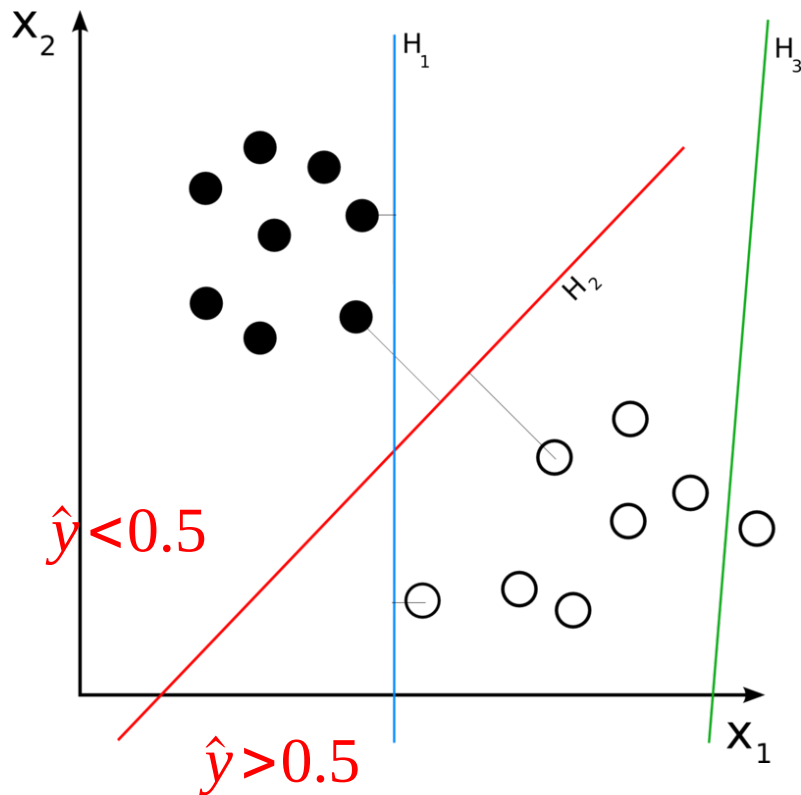
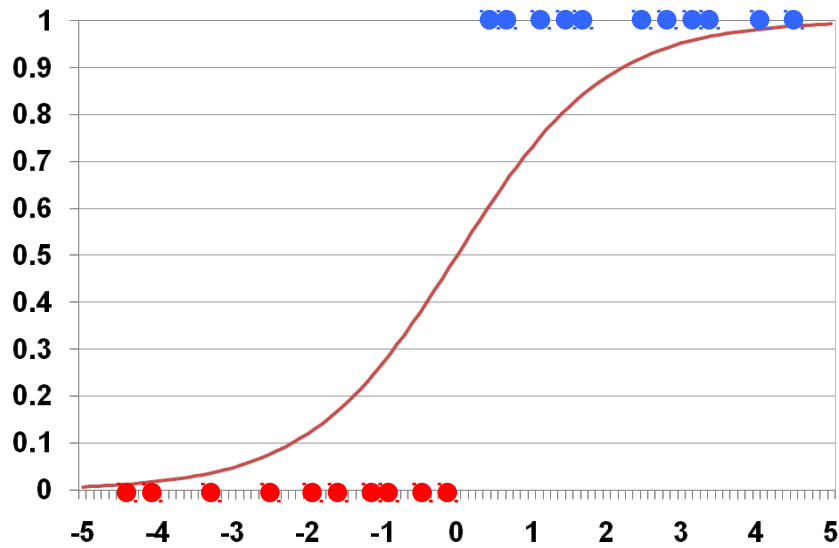
\vec{x}

$$\hat{y} = \sigma(\vec{\theta}^T \vec{x})$$

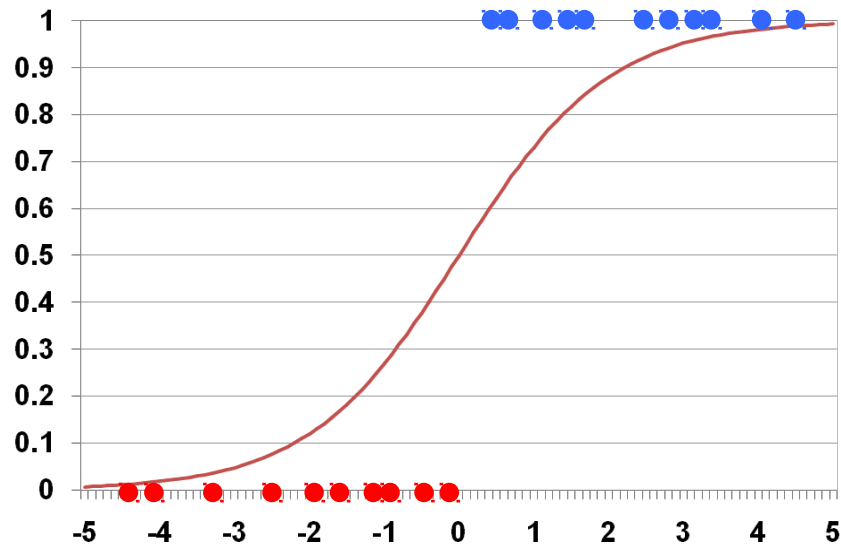
Logistic regression is linear



Logistic regression is linear



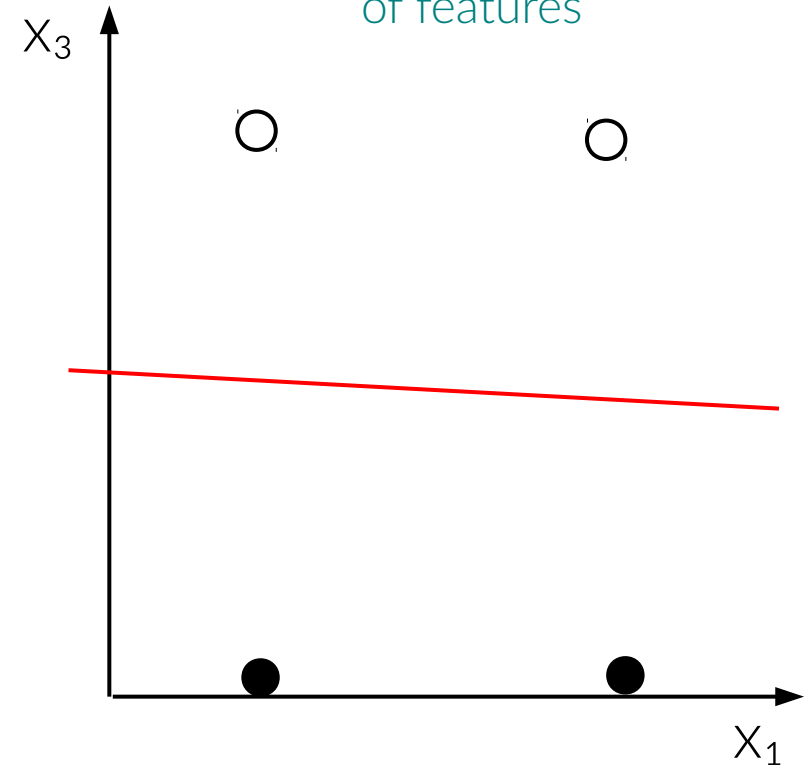
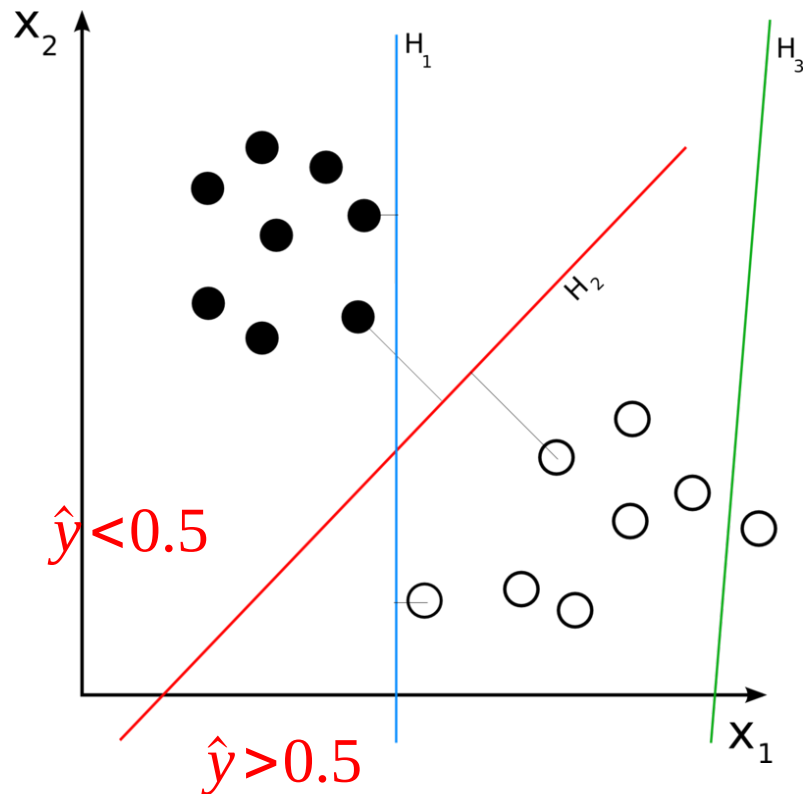
Logistic regression is linear



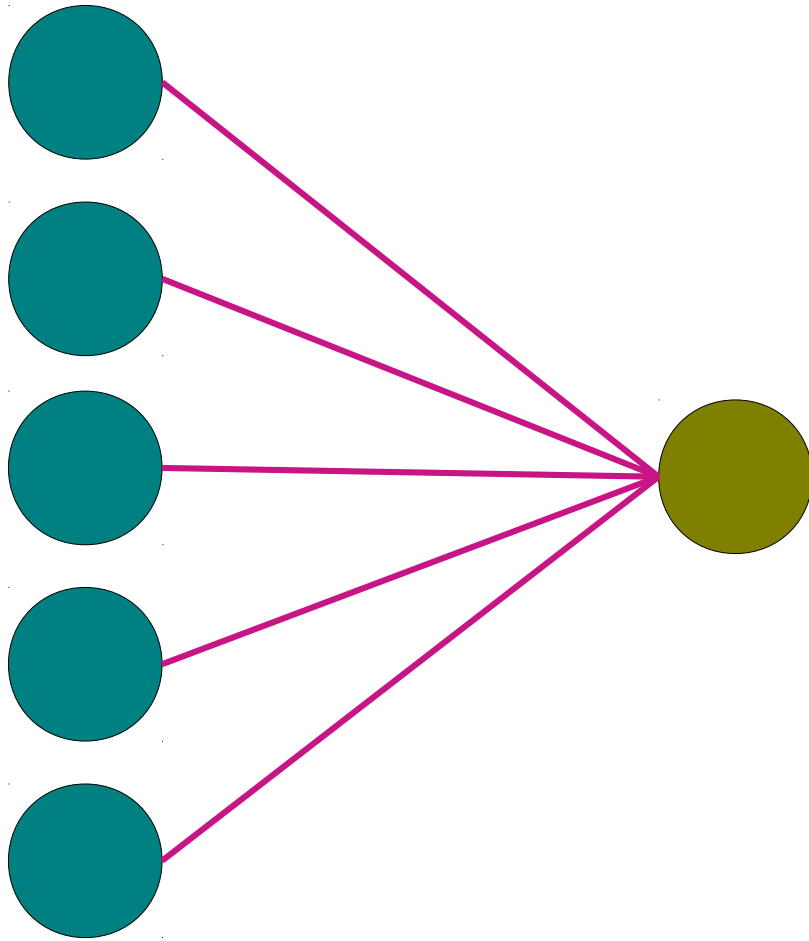
non-linear function

$$X_3 = (X_1 - X_2)^2$$

linear combination
of features



A cartoon of logistic regression



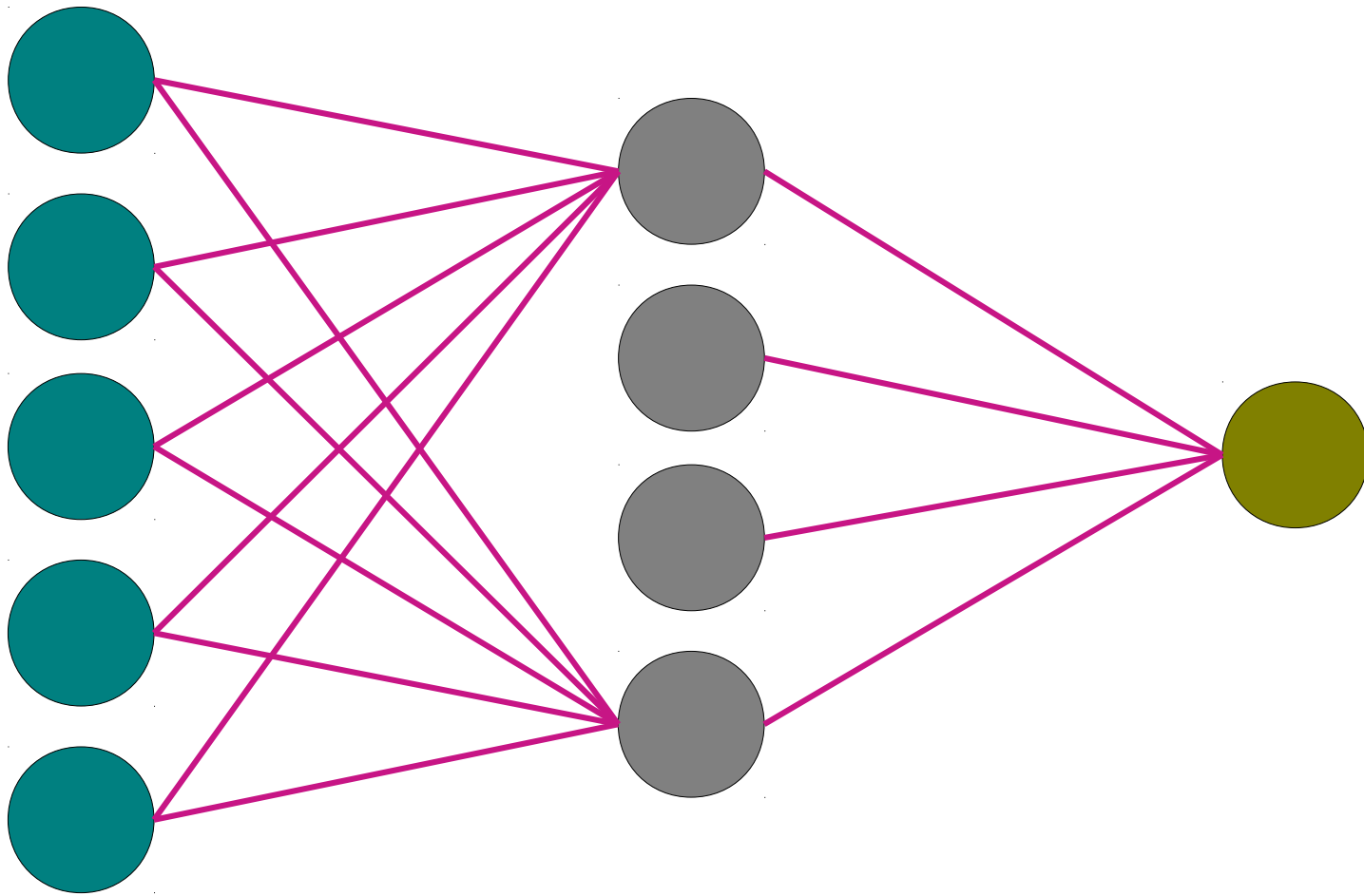
\vec{x}

non-linear function

$$\hat{y} = \sigma \left(\vec{\theta}^T \vec{x} \right)$$

linear combination of features

Stacking logistic regression

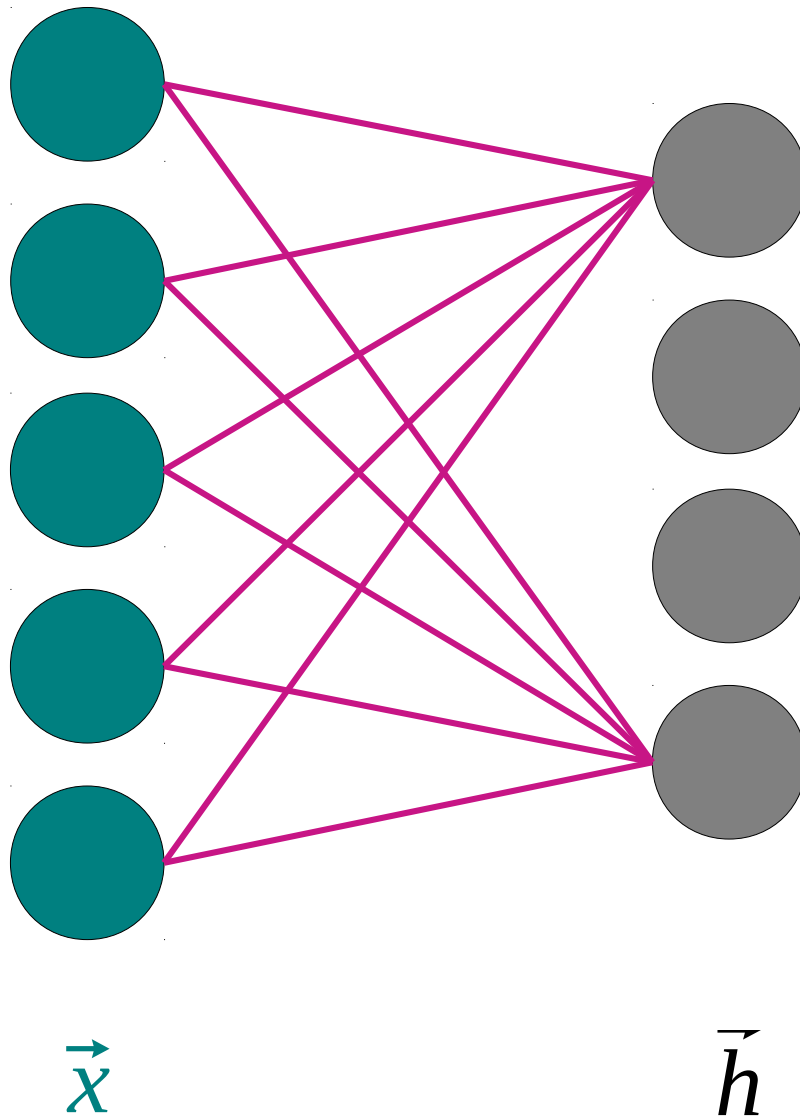


$$\vec{x}$$

$$\vec{h} = \sigma(\vec{\theta}^{(h)} \vec{x})$$

$$\hat{y} = \sigma(\vec{\theta}^{(\hat{y})T} \vec{h})$$

Unpacking the linear algebra



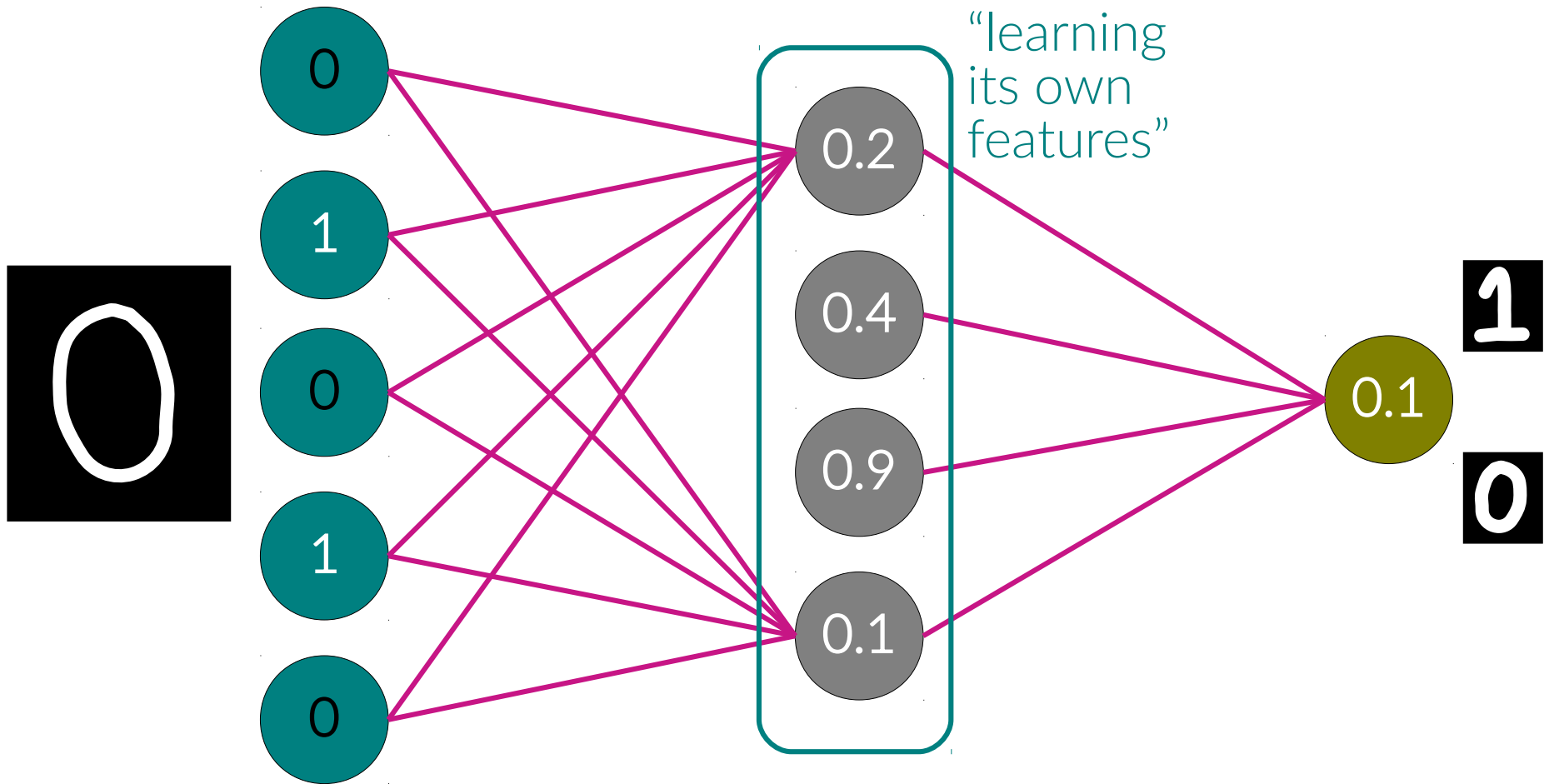
matrix $\theta^{(h)}$ vector \vec{x}

$$\vec{h} = \sigma \left(\underbrace{\theta^{(h)}}_{\text{matrix}} \underbrace{\vec{x}}_{\text{vector}} \right)$$

vector

$$h_i = \sigma \left(\sum_{j=0}^m \theta_{i,j}^{(h)} x_j \right)$$

Stacking logistic regression

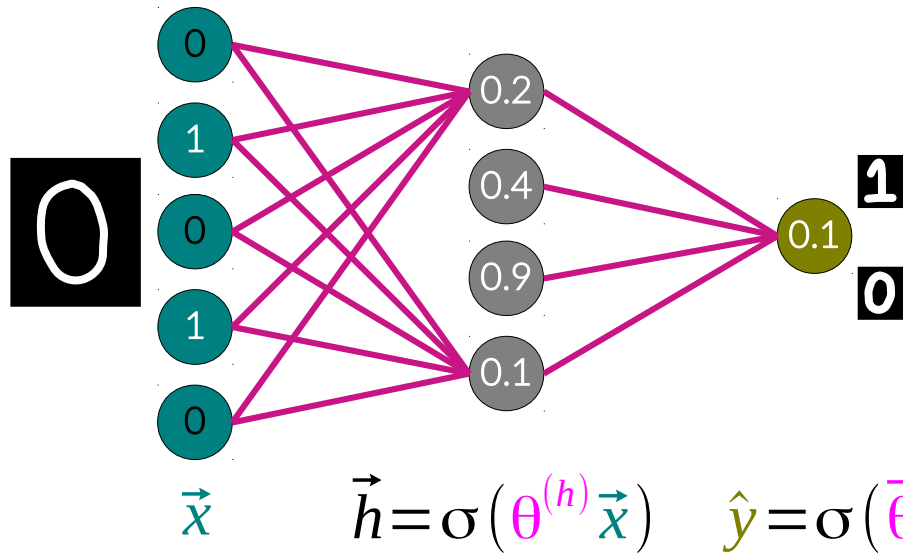


\vec{x}
input
features

$\vec{h} = \sigma(\theta^{(h)} \vec{x})$
hidden
representation

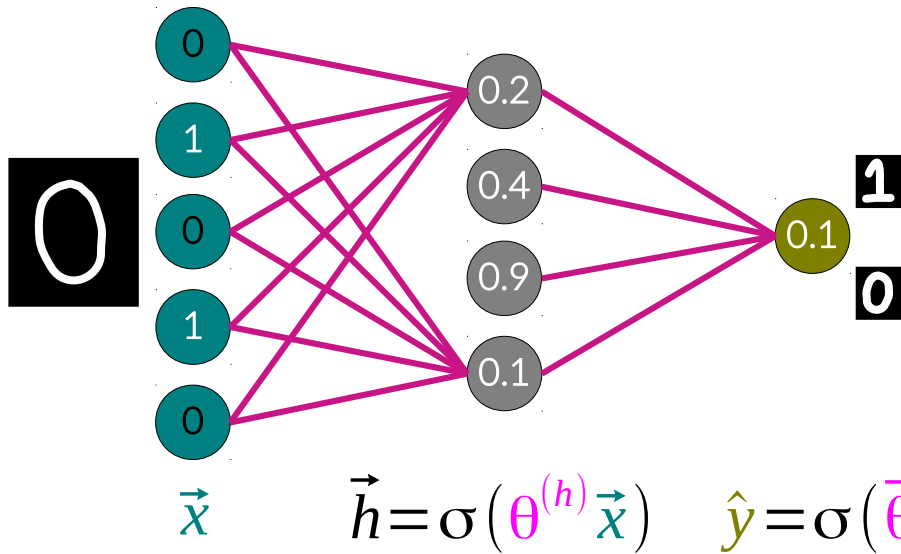
$\hat{y} = \sigma(\vec{\theta}^{(\hat{y})T} \vec{h})$
output
class label

Maximum likelihood with neural nets



$$\begin{aligned} L(\theta) &= P(X^{(1)}, \dots, X^{(n)}, Y^{(1)}, \dots, Y^{(n)} | \theta) \\ &= \prod_{i=1}^n P(X^{(i)}, Y^{(i)} | \theta) \\ &= \prod_{i=1}^n \hat{y}^{y^{(i)}} (1 - \hat{y})^{1 - y^{(i)}} P(X^{(i)}) \end{aligned}$$

Maximum likelihood with neural nets



$$L(\theta) = \prod_{i=1}^n \hat{y}^{y^{(i)}} (1 - \hat{y})^{1 - y^{(i)}} P(\mathbf{X}^{(i)})$$

$$LL(\theta) = \sum_{i=1}^n \left[y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log (1 - \hat{y}^{(i)}) + \log P(\mathbf{X}^{(i)}) \right]$$

Maximum likelihood with neural nets

$$L(\boldsymbol{\theta}) = \prod_{i=1}^n \hat{y}^{y^{(i)}} (1 - \hat{y})^{1 - y^{(i)}} P(\mathbf{X}^{(i)})$$

$$LL(\boldsymbol{\theta}) = \sum_{i=1}^n \left[y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log (1 - \hat{y}^{(i)}) + \log P(\mathbf{X}^{(i)}) \right]$$

$$\frac{\partial}{\partial \boldsymbol{\theta}_j^{(\hat{y})}} LL(\boldsymbol{\theta}) = \sum_{i=1}^n \frac{\partial}{\partial \boldsymbol{\theta}_j^{(\hat{y})}} \left[y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log (1 - \hat{y}^{(i)}) \right]$$

$$= \sum_{i=1}^n \left[\frac{y^{(i)}}{\hat{y}^{(i)}} - \frac{(1 - y^{(i)})}{(1 - \hat{y}^{(i)})} \right] \frac{\partial \hat{y}^{(i)}}{\partial \boldsymbol{\theta}_j^{(\hat{y})}}$$

$$= \sum_{i=1}^n \left[\frac{y^{(i)}}{\hat{y}^{(i)}} - \frac{(1 - y^{(i)})}{(1 - \hat{y}^{(i)})} \right] \frac{\partial}{\partial \boldsymbol{\theta}_j^{(\hat{y})}} \sigma(\hat{\boldsymbol{\theta}}^{(\hat{y})T} \vec{h})$$

$$= \sum_{i=1}^n \left[\frac{y^{(i)}}{\hat{y}^{(i)}} - \frac{(1 - y^{(i)})}{(1 - \hat{y}^{(i)})} \right] \hat{y}^{(i)} (1 - \hat{y}^{(i)}) h_j$$

Maximum likelihood with neural nets

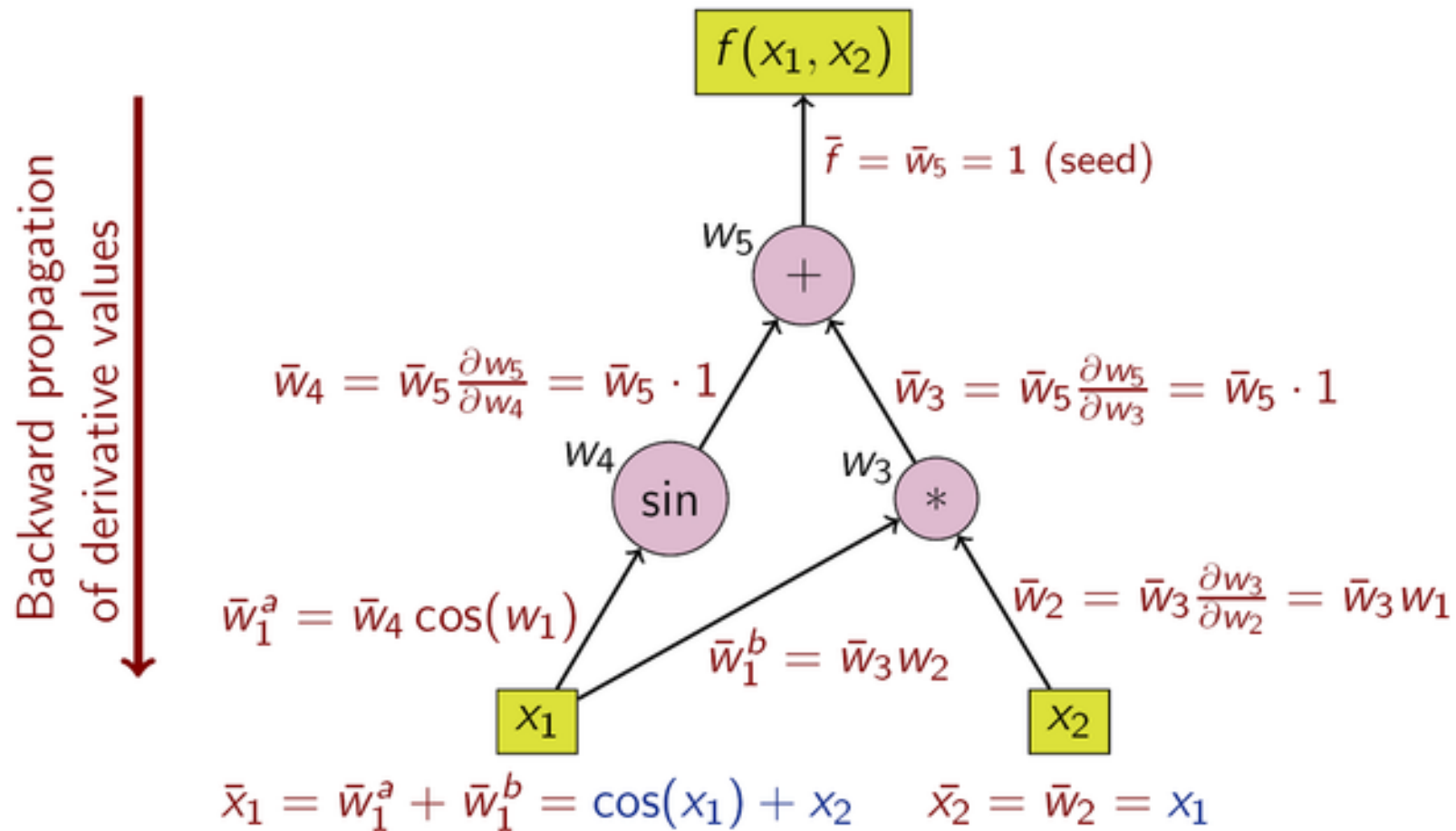
$$L(\boldsymbol{\theta}) = \prod_{i=1}^n \hat{y}^{y^{(i)}} (1 - \hat{y})^{1 - y^{(i)}} P(\mathbf{X}^{(i)})$$

$$LL(\boldsymbol{\theta}) = \sum_{i=1}^n \left[y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log (1 - \hat{y}^{(i)}) + \log P(\mathbf{X}^{(i)}) \right]$$

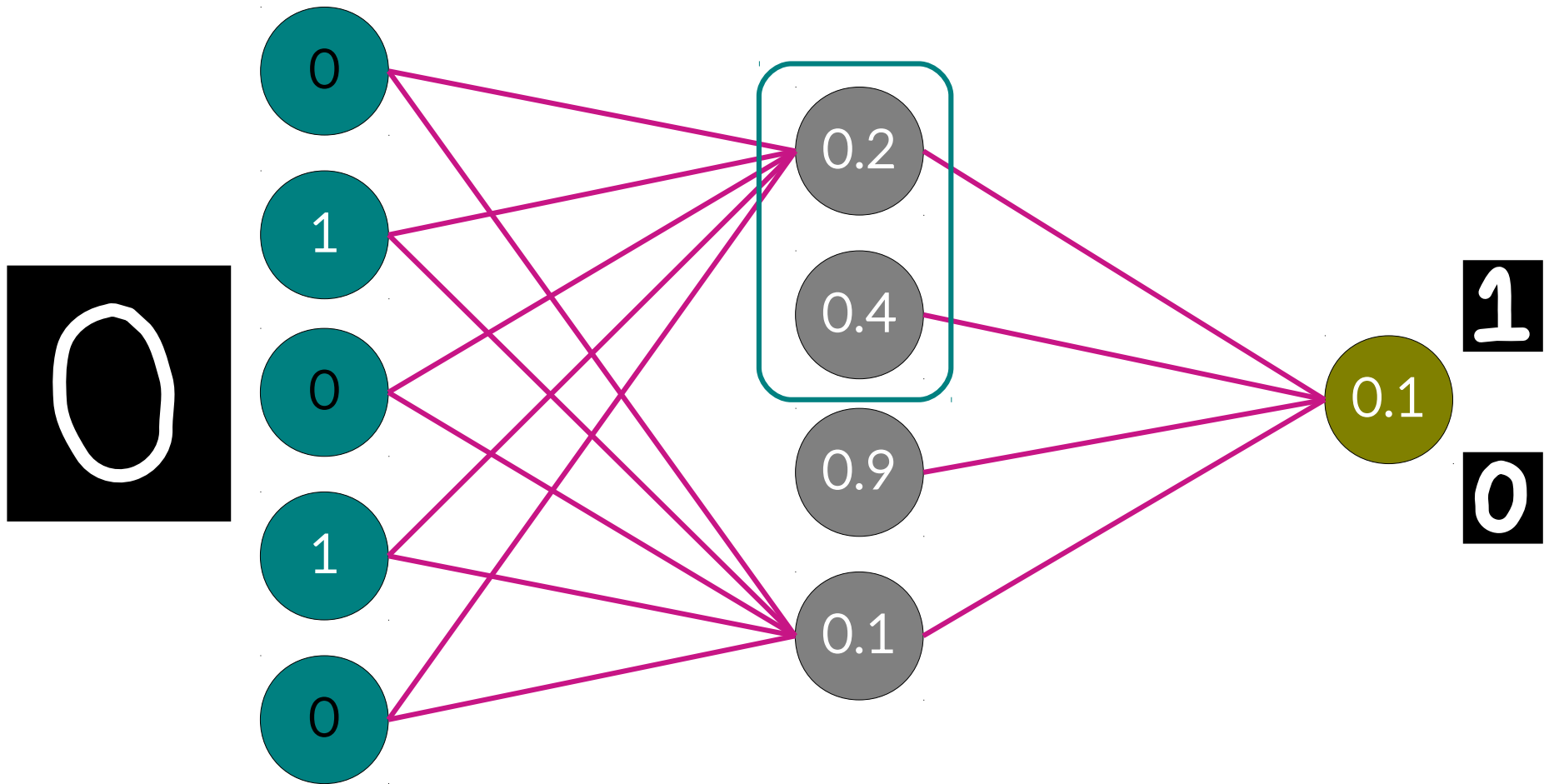
$$\frac{\partial}{\partial \boldsymbol{\theta}_j^{(\hat{y})}} LL(\boldsymbol{\theta}) = \sum_{i=1}^n \left[\underbrace{\frac{y^{(i)}}{\hat{y}^{(i)}} - \frac{(1 - y^{(i)})}{(1 - \hat{y}^{(i)})}}_{\frac{\partial}{\partial \hat{y}^{(i)}} LL^{(i)}(\boldsymbol{\theta})} \underbrace{\hat{y}^{(i)} (1 - \hat{y}^{(i)})}_{\frac{\partial \hat{y}^{(i)}}{\partial z^{(i)}}} \underbrace{h_j}_{\frac{\partial z^{(i)}}{\partial \boldsymbol{\theta}_j^{(\hat{y})}}}$$

$$\frac{\partial}{\partial \boldsymbol{\theta}_{j,k}^{(h)}} LL(\boldsymbol{\theta}) = \sum_{i=1}^n \left[\underbrace{\frac{y^{(i)}}{\hat{y}^{(i)}} - \frac{(1 - y^{(i)})}{(1 - \hat{y}^{(i)})}}_{\frac{\partial}{\partial \hat{y}^{(i)}} LL^{(i)}(\boldsymbol{\theta})} \underbrace{\hat{y}^{(i)} (1 - \hat{y}^{(i)})}_{\frac{\partial \hat{y}^{(i)}}{\partial z^{(i)}}} \underbrace{\boldsymbol{\theta}_j^{(\hat{y})}}_{\frac{\partial z^{(i)}}{dh_j^{(i)}}} \underbrace{h_j (1 - h_j) x_k}_{\frac{dh_j^{(i)}}{\partial \boldsymbol{\theta}_{j,k}^{(h)}}$$

Automatic differentiation



Breaking the symmetry

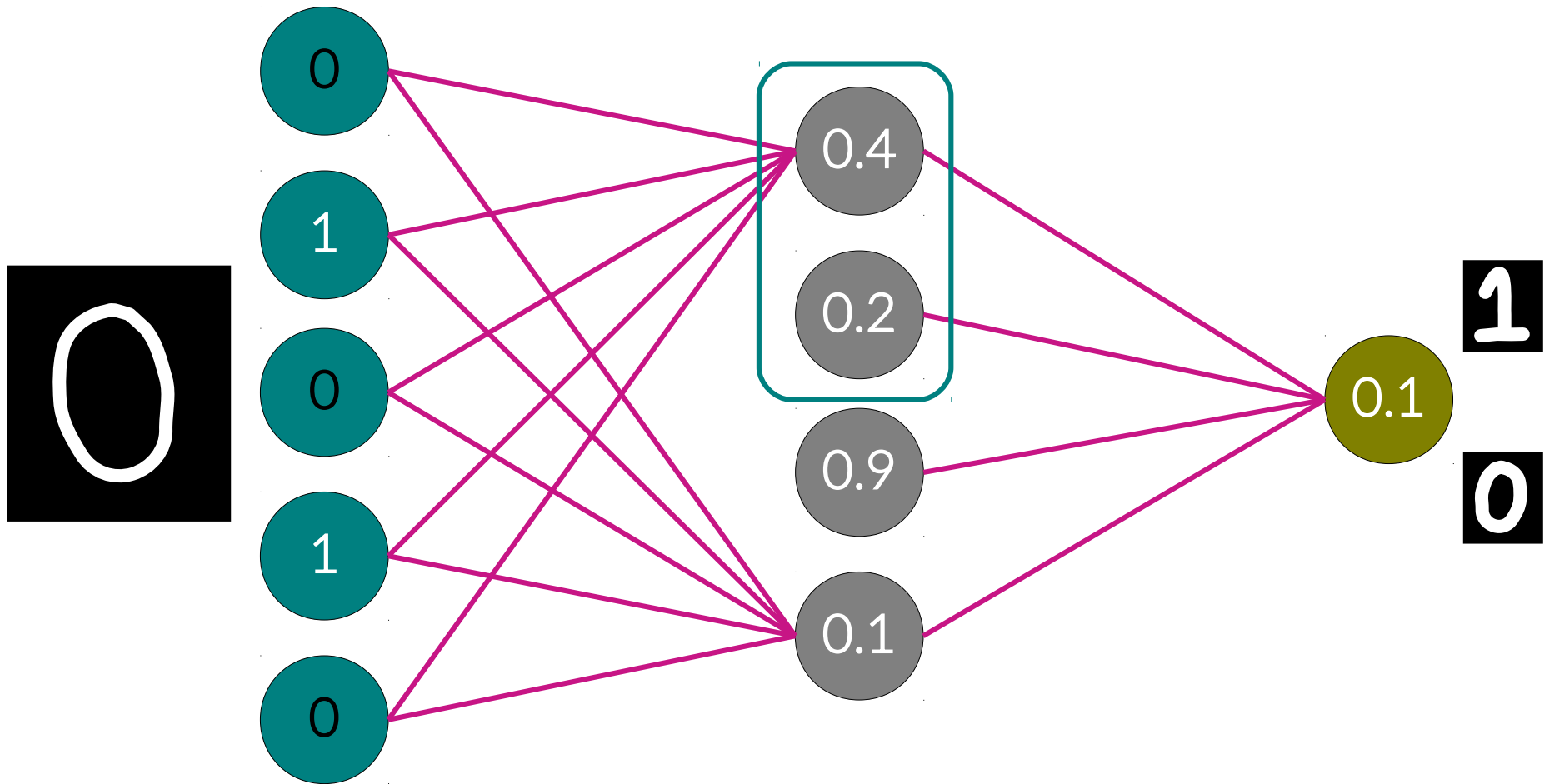


\vec{x}
input
features

$\vec{h} = \sigma(\theta^{(h)} \vec{x})$
hidden
representation

$\hat{y} = \sigma(\vec{\theta}^{(\hat{y})T} \vec{h})$
output
class label

Breaking the symmetry

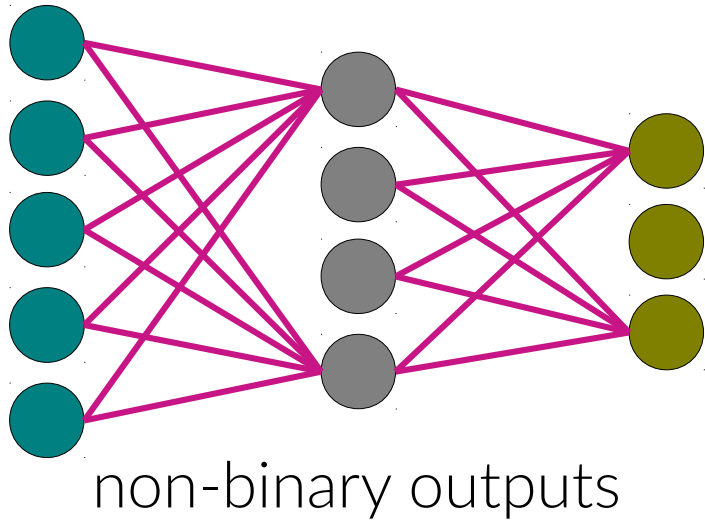


\vec{x}
input
features

$\vec{h} = \sigma(\theta^{(h)} \vec{x})$
hidden
representation

$\hat{y} = \sigma(\vec{\theta}^{(\hat{y})T} \vec{h})$
output
class label

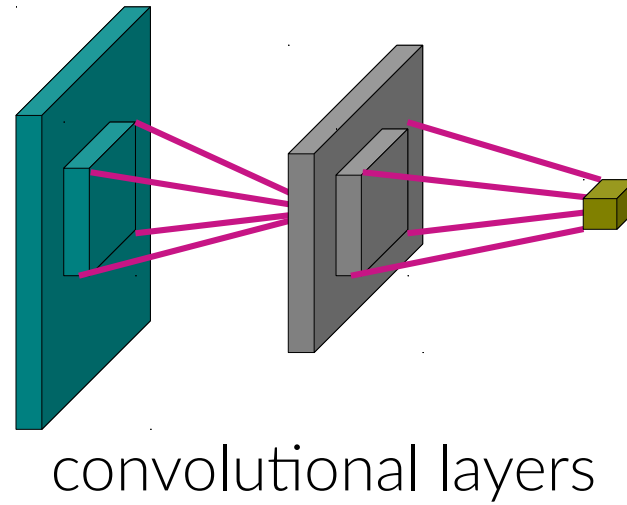
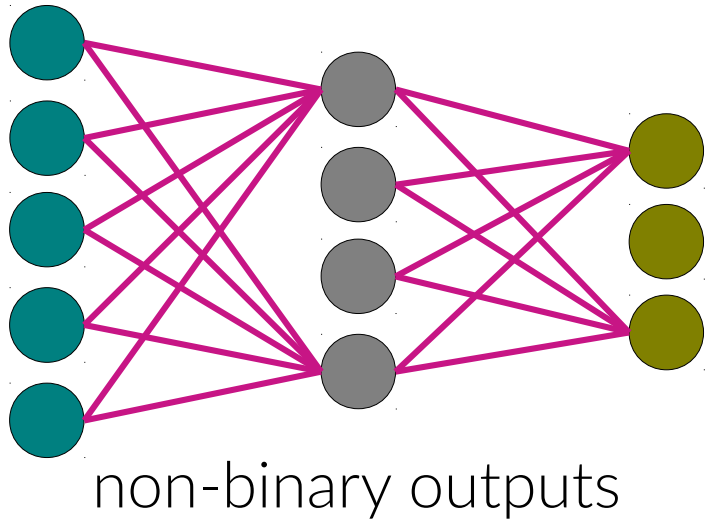
Expanding the toolbox



Applications: Image recognition

0 0 0 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9 9 9 9 9

Expanding the toolbox



Applications: Image recognition

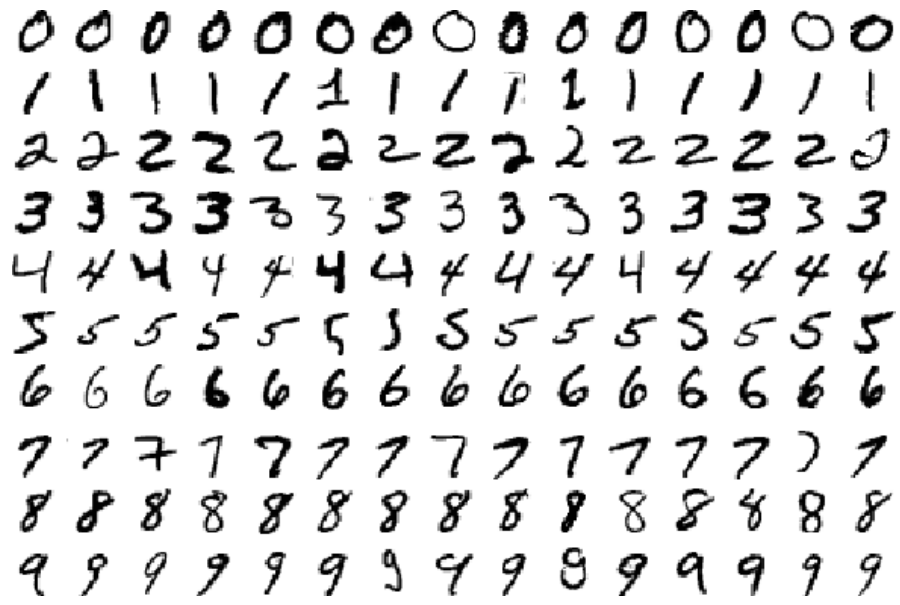
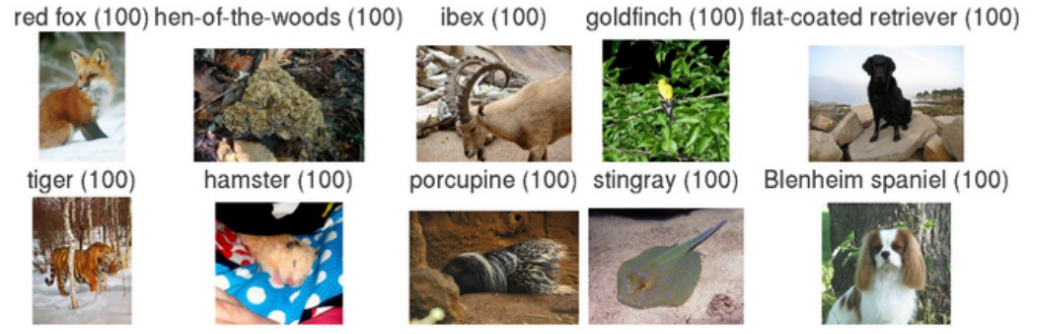


Image classification

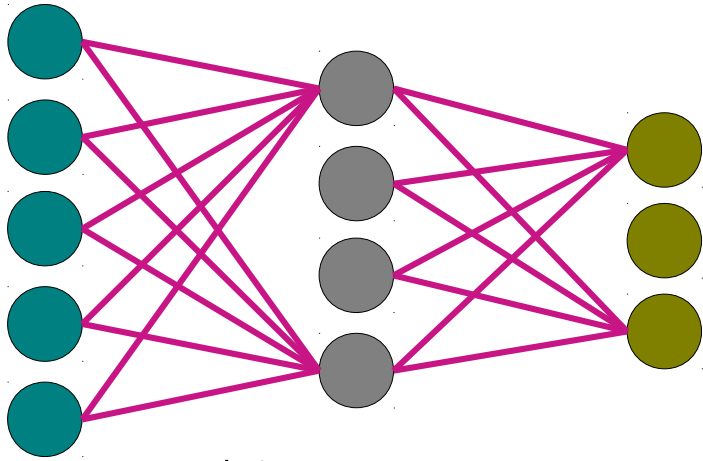
Easiest classes



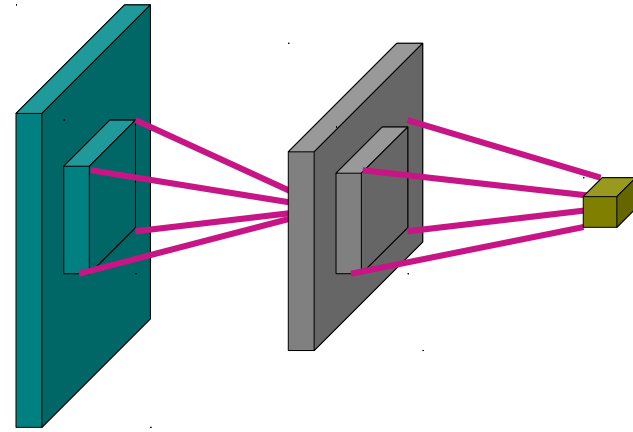
Hardest classes



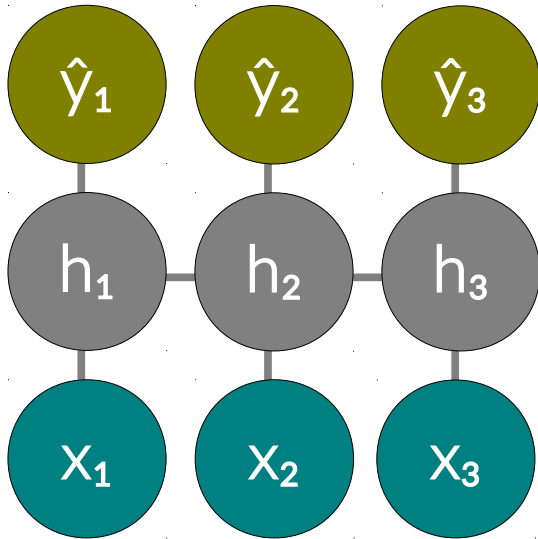
Expanding the toolbox



non-binary outputs



convolutional layers



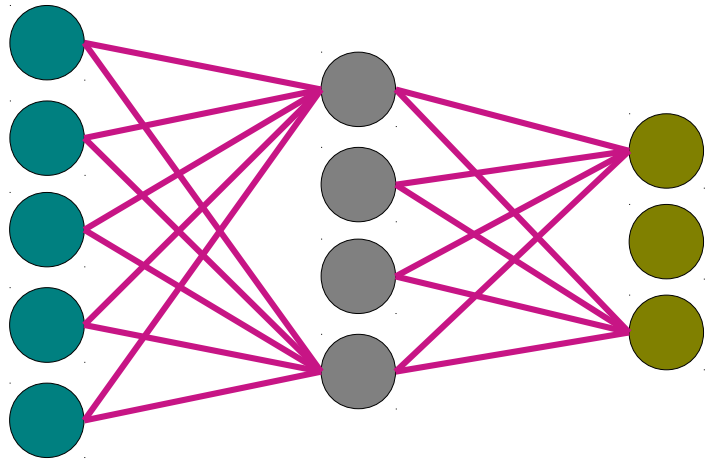
recurrent connections

Applications: Speech recognition

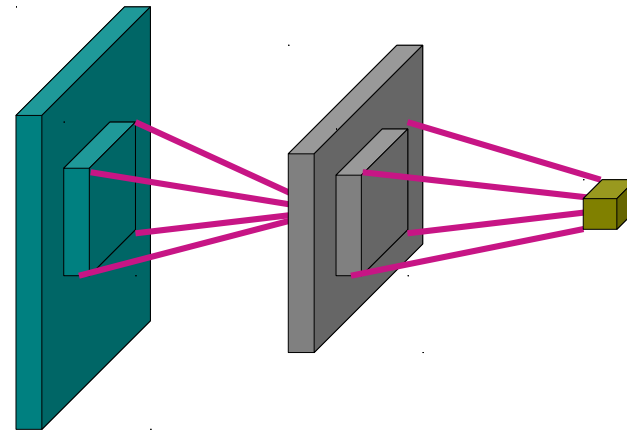


who is the current president of France?

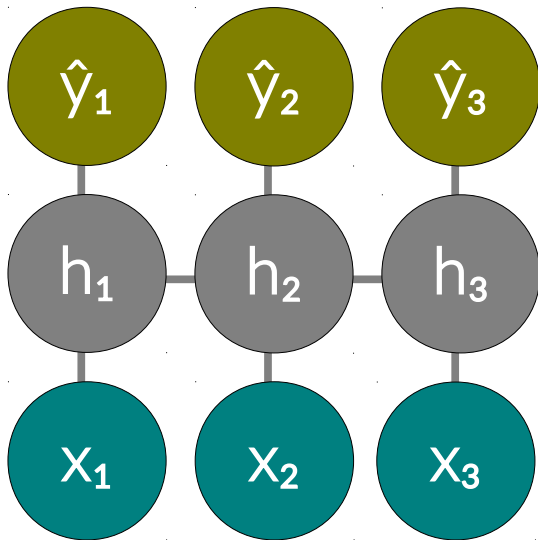
Expanding the toolbox



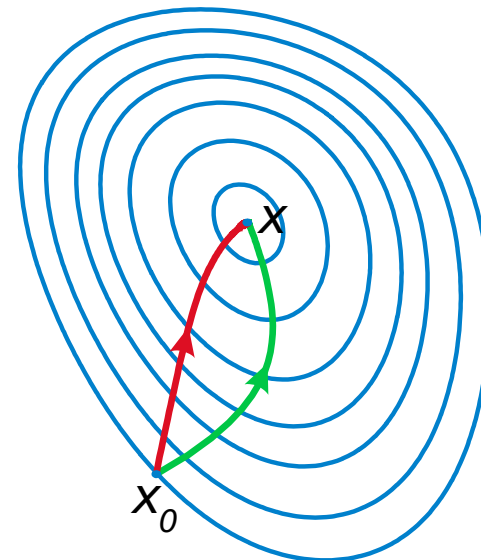
non-binary outputs



convolutional layers



recurrent connections



fancy maximization

Break time!

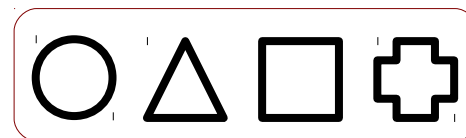
General principle of counting

An experiment consisting of two or more **separate parts** has a number of outcomes equal to the **product** of the number of outcomes of each part.

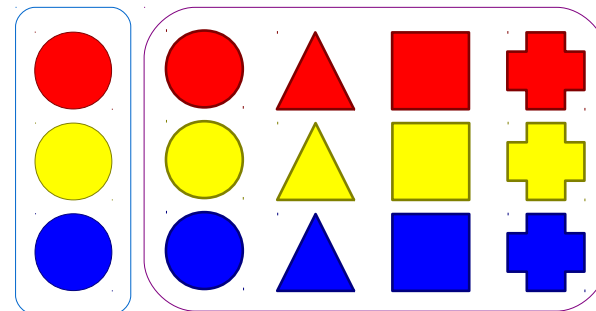


$$|A_1 \times A_2 \times \cdots \times A_n| = \prod_i |A_i|$$

shapes: 4



colors: 3



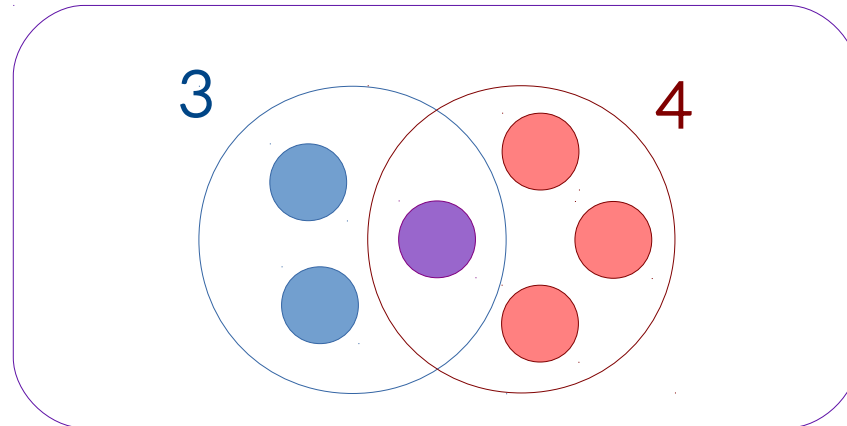
total:
 $4 \cdot 3 = 12$

Principle of Inclusion/Exclusion

The **total number of elements** in two sets is the sum of the number of elements of each set, **minus** the number of elements in both sets.



$$|A \cup B| = |A| + |B| - |A \cap B|$$



$$3 + 4 - 1 = 6$$

Inclusion/exclusion with more than two sets

size of the union

add or subtract (based on size)

size of intersections

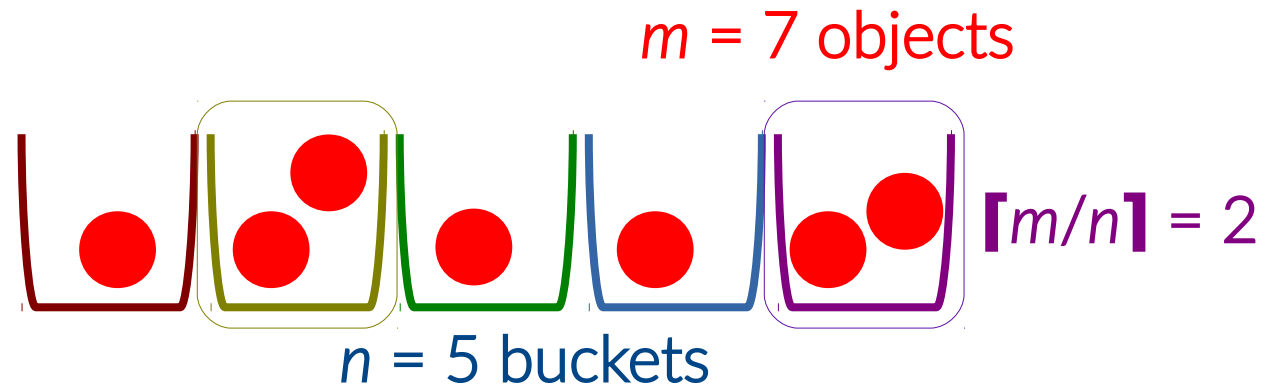
$$\left| \bigcup_{i=1}^n E_i \right| = \sum_{r=1}^n (-1)^{(r+1)} \sum_{i_1 < \dots < i_r} \left| \bigcap_{j=1}^r E_{i_j} \right|$$

sum over subset sizes

sum over all **subsets** of that size

General Pigeonhole Principle

If m objects are placed in n buckets, then **at least one bucket** must contain at least $\lceil m/n \rceil$ objects.



Permutations

The number of **ways** of ordering n **distinguishable** objects.



$$n! = 1 \cdot 2 \cdot 3 \cdot \dots \cdot n = \prod_{i=1}^n i$$

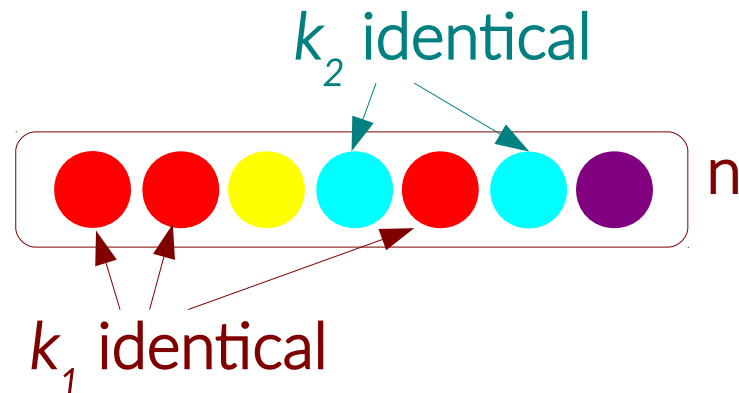


Permutations with indistinct elements

The number of **ways of ordering** n . objects, where some groups are **indistinguishable**.



$$\binom{n}{k_1, k_2, \dots, k_m} = \frac{n!}{k_1! k_2! \dots k_m!}$$

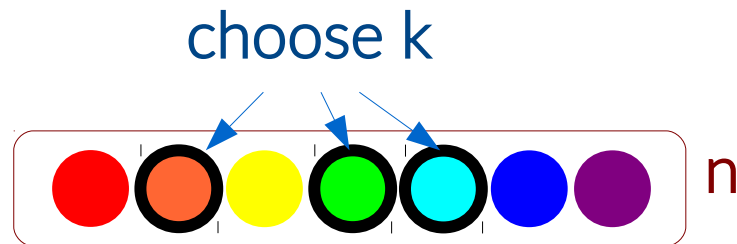


Combinations

The number of unique **subsets** of size k from a larger set of size n .
(objects are distinguishable, unordered)



$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$



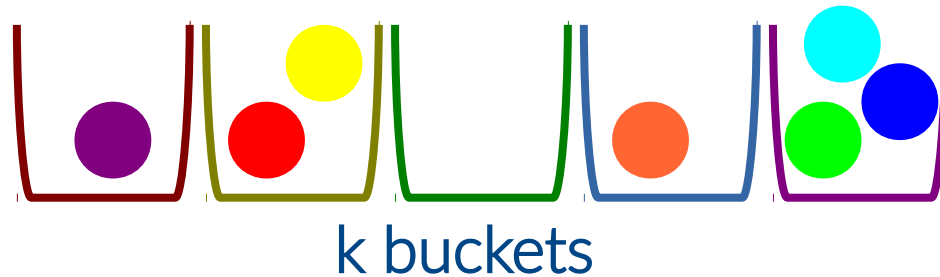
Bucketing

The number of ways of assigning n distinguishable objects to a fixed set of k buckets or labels.



$$k^n$$

n objects

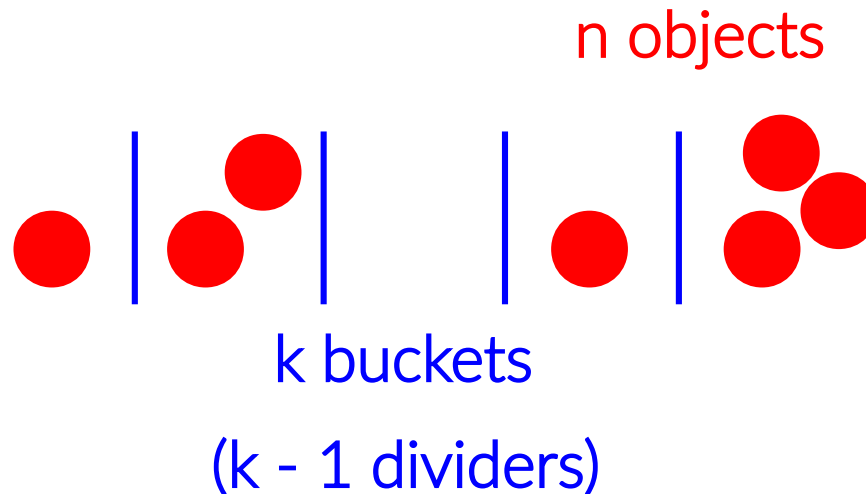


Divider method

The number of ways of assigning n indistinguishable objects to a fixed set of k buckets or labels.



$$\binom{n + (k - 1)}{n}$$



A grid of ways of counting

Ordering

Subsets

Bucketing

All
distinct

$$n!$$

$$\binom{n}{k}$$

$$k^n$$

Some
indistinct

$$\frac{n!}{k_1! k_2! \dots k_m!}$$

Creativity!

- Split into cases
- Use inclusion/exclusion
- Reframe the problem

All
indistinct

$$1$$

$$1$$

$$\binom{n + (k-1)}{n}$$

Axioms of probability

(1) $0 \leq P(E) \leq 1$

(2) $P(S) = 1$

(3) If $E \cap F = \emptyset$, then
$$P(E \cup F) = P(E) + P(F)$$

(Sum rule, but with probabilities!)



How do I get started?



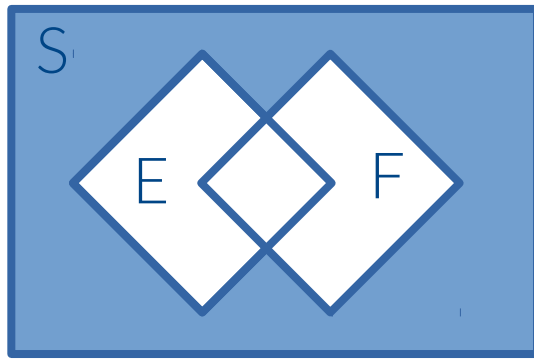
For word problems involving probability, start by defining **events**!

Getting rid of ORs

Finding the probability of an OR of events can be nasty. Try **using De Morgan's laws** to turn it into an AND!



$$P(A \cup B \cup \dots \cup Z) = 1 - P(A^c B^c \dots Z^c)$$

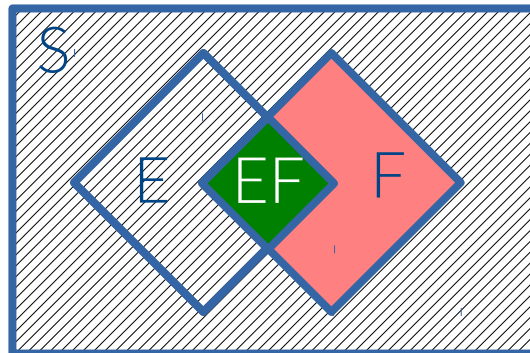


Definition of conditional probability

The conditional probability $P(E | F)$ is the probability that E happens, **given** that F has happened. F is the new sample space.



$$P(E | F) = \frac{P(EF)}{P(F)}$$

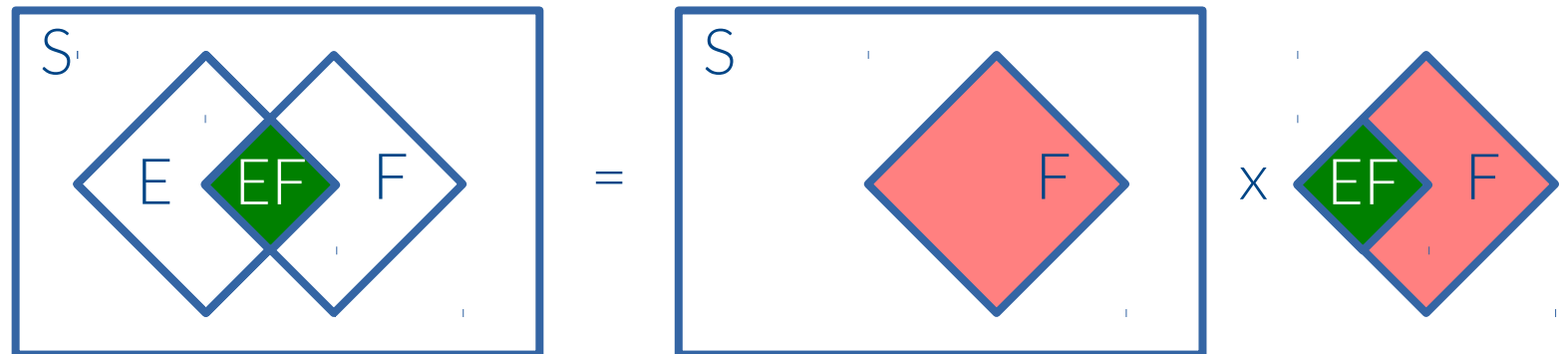


Chain rule of probability

The probability of **both** events happening is the probability of **one happening** times the probability of **the other happening given the first one**.



$$P(EF) = P(F)P(E|F)$$



General chain rule of probability

The probability of **all** events happening is the probability of **the first** happening times the prob. of **the second given the first** times the prob. of **the third given the first two** ...etc.



$$P(EFG\dots) = P(E)P(F|E)P(G|EF)\dots$$

Law of total probability

You can compute an overall probability by adding up the case when an event **happens** and when it **doesn't happen**.

$$P(F) = P(EF) + P(E^C F)$$

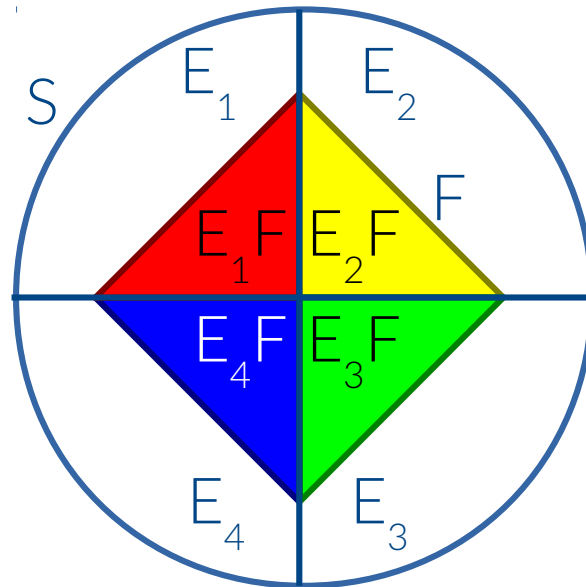
$$= P(E)P(F|E) + P(E^C)P(F|E^C)$$



General law of total probability

You can compute an overall probability by summing over **mutually exclusive** and **exhaustive** sub-cases.

$$\begin{aligned} P(F) &= \sum_i P(E_i F) \\ &= \sum_i P(E_i) P(F|E_i) \end{aligned}$$



Bayes' theorem


You can “flip” a conditional probability if you multiply by the probability of the **hypothesis** and divide by the probability of the **observation**.

$$P(E|F) = \frac{P(F|E)P(E)}{P(F)}$$



Finding the denominator

If you don't know $P(F)$ on the bottom, try using the **law of total probability**.


$$P(E|F) = \frac{P(F|E)P(E)}{P(F|E)P(E) + P(F|E^c)P(E^c)}$$

$$P(E|F) = \frac{P(F|E)P(E)}{\sum_i P(F|E_i)P(E_i)}$$

Independence

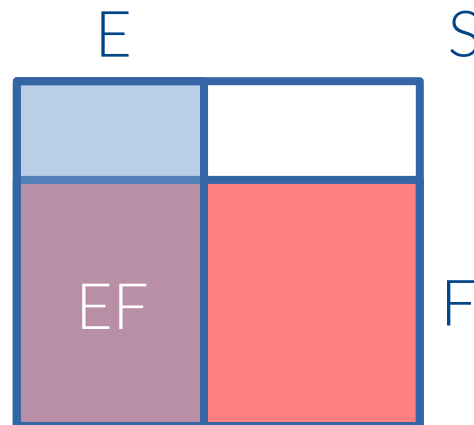
Two events are **independent** if you can **multiply** their probabilities to get the probability of **both** happening.

$$P(EF) = P(E)P(F)$$



$$E \perp F$$

← (“independent of”)



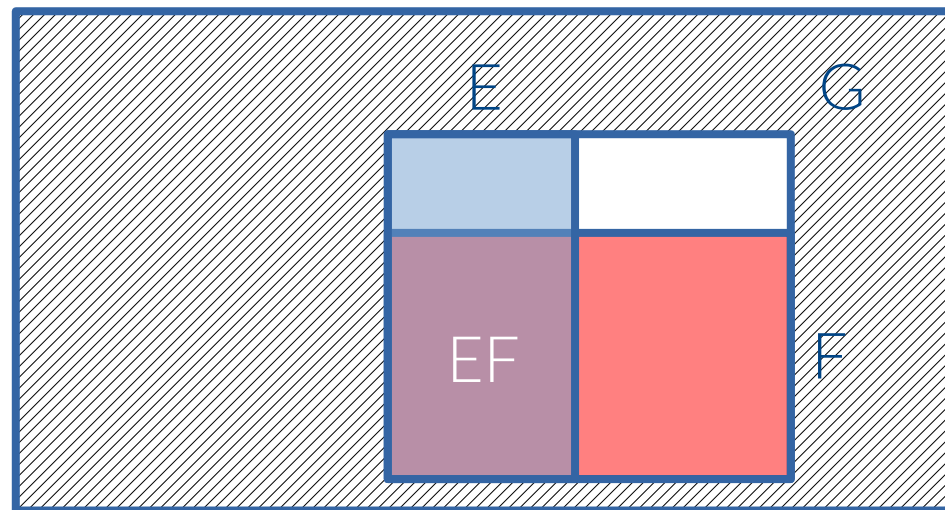
Conditional independence

Two events are **conditionally independent** if you can **multiply** their conditional probabilities to get the conditional probability of **both** happening.

$$P(EF|G) = P(E|G)P(F|G)$$



$$(E \perp F) | G$$

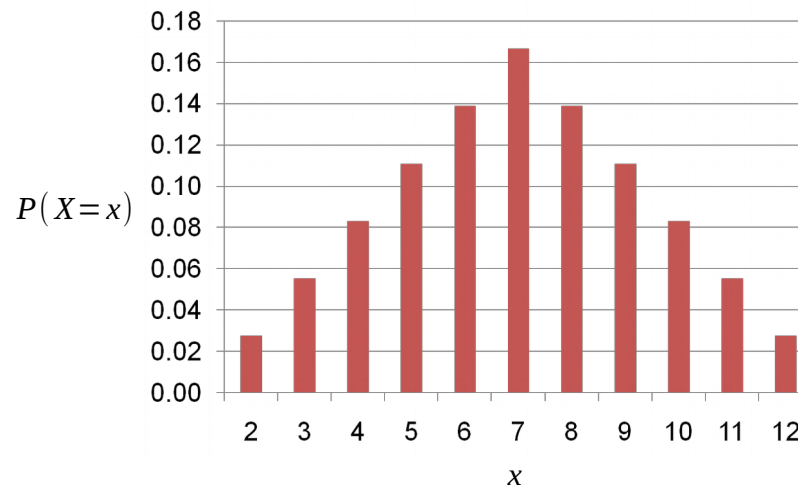


Random variables

A **random variable** takes on values probabilistically.



$$P(X=2) = \frac{1}{36}$$



How do I get started?



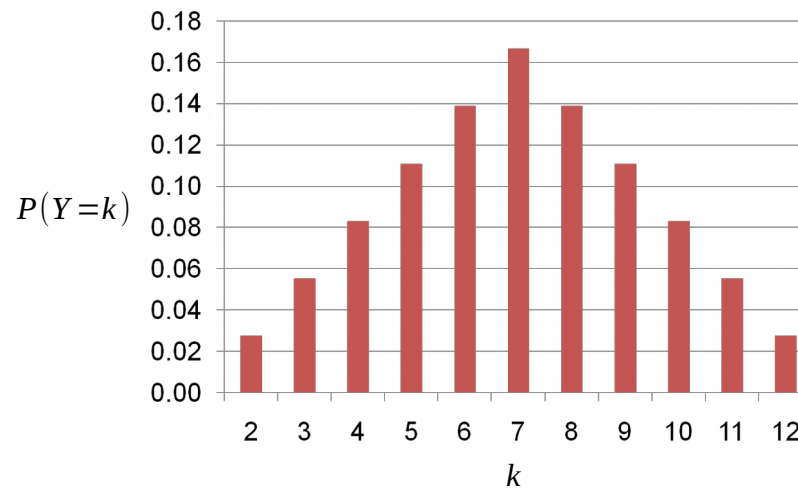
For word problems involving probability, start by defining **events** and **random variables**!

Probability mass function

The **probability mass function** (PMF) of a random variable is a function from values of the variable to probabilities.



$$p_Y(k) = P(Y = k)$$

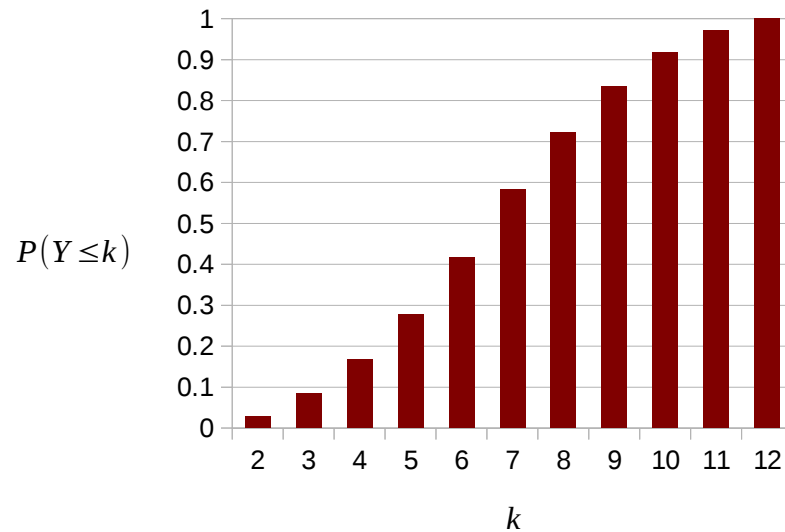


Cumulative distribution function

The **cumulative distribution function** (CDF) of a random variable is a function giving the probability that the random variable is **less than or equal to** a value.



$$F_Y(k) = P(Y \leq k)$$

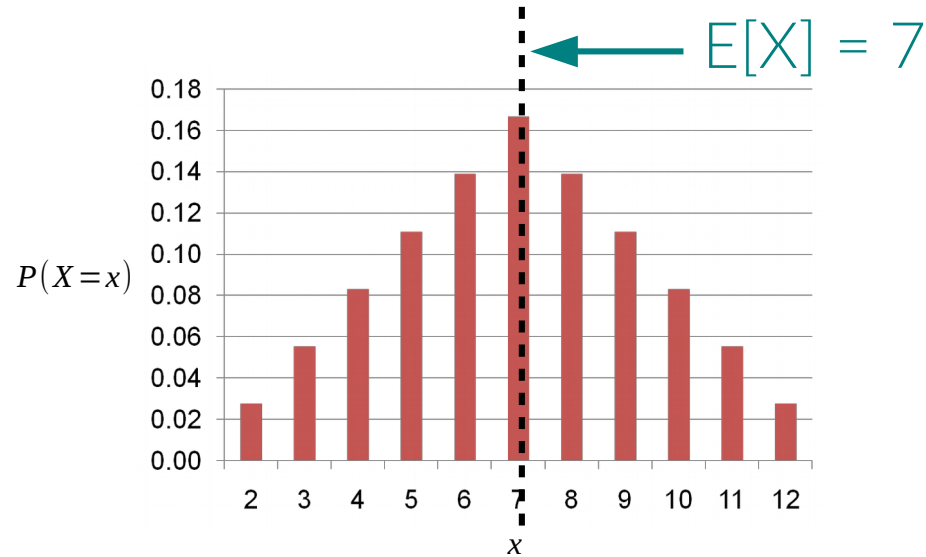


Expectation

The **expectation** of a random variable is the “**average**” value of the variable (weighted by probability).



$$E[X] = \sum_{x: p(x) > 0} p(x) \cdot x$$

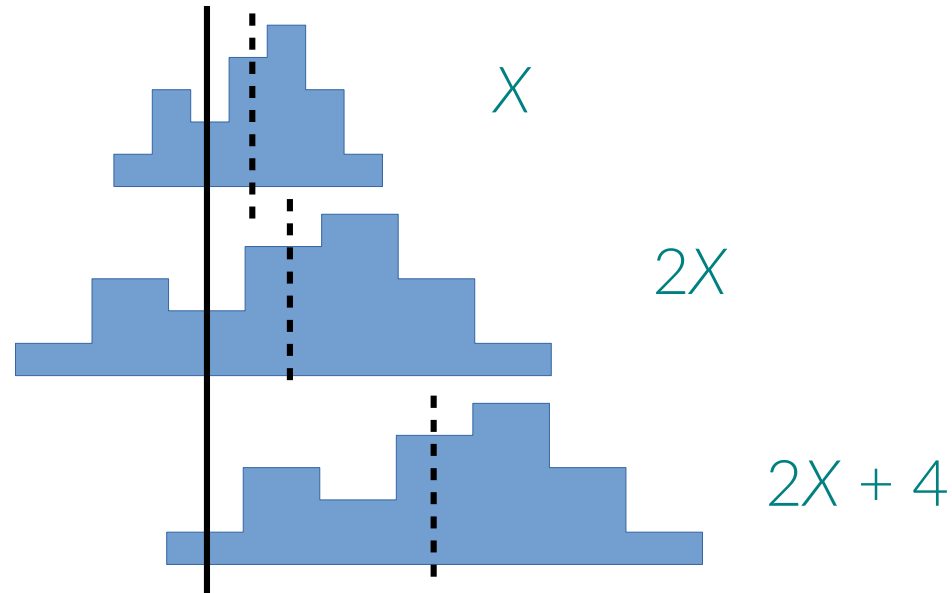


Linearity of expectation

Adding random variables or constants? **Add** the expectations.
Multiplying by a constant? **Multiply** the expectation by the constant.



$$E[aX + bY + c] = aE[X] + bE[Y] + c$$



Indicator variable

An **indicator variable** is a “Boolean” variable, which takes values 0 or 1 corresponding to whether an event takes place.



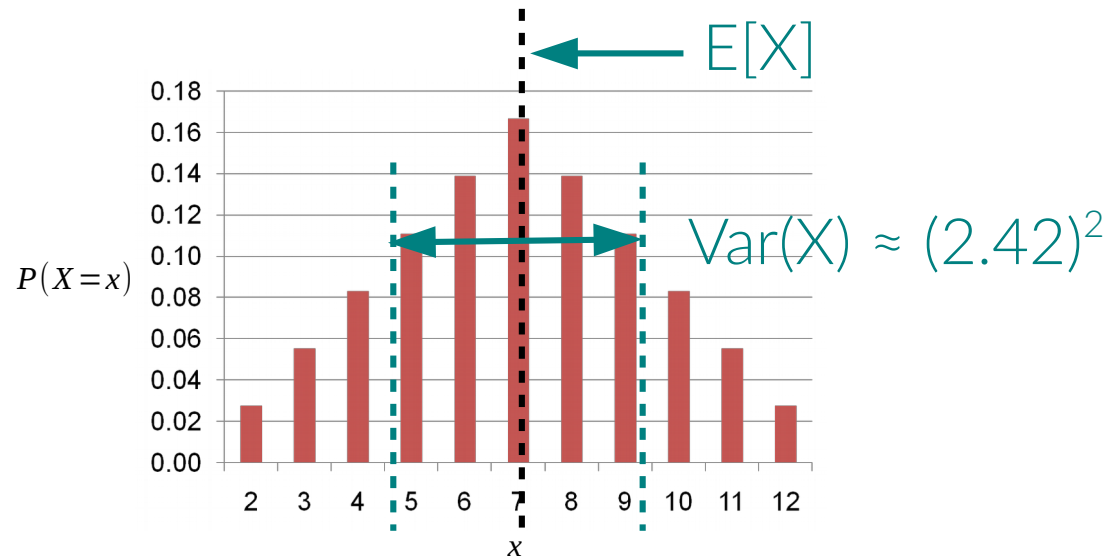
$$I = \mathbb{1}[A] = \begin{cases} 1 & \text{if event } A \text{ occurs} \\ 0 & \text{otherwise} \end{cases}$$



Variance

Variance is the average **square** of the **distance** of a variable from the expectation. Variance measures the “**spread**” of the variable.

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

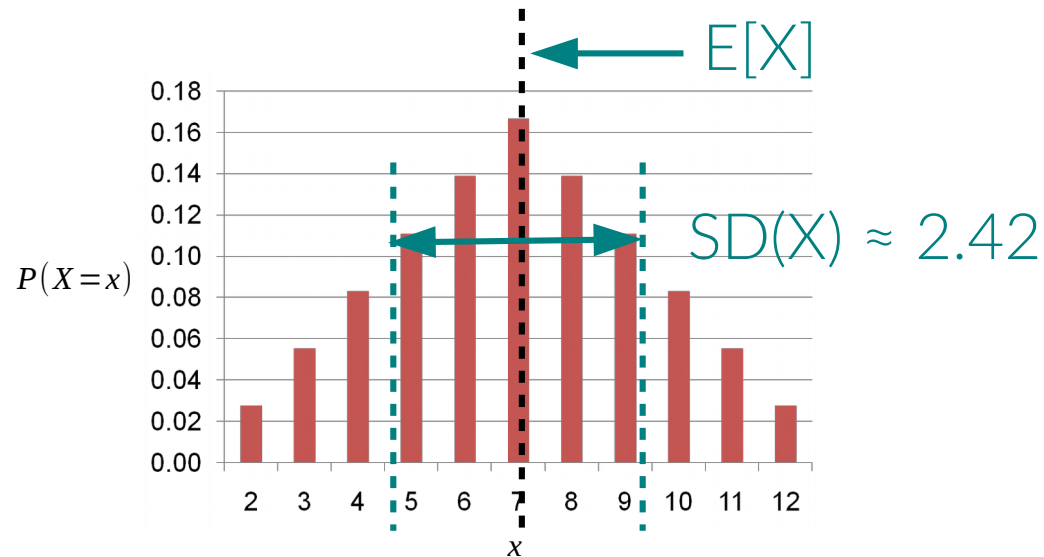


Standard deviation

Standard deviation is the (“root-mean-square”) average of the **distance** of a variable from the expectation.



$$\text{SD}(X) = \sqrt{\text{Var}(X)} = \sqrt{E[(X - E[X])^2]}$$



Variance of a linear function

Adding a constant? Variance **doesn't** change.
Multiplying by a constant? **Multiply** the variance by the **square** of the constant.

$$\begin{aligned}\text{Var}(aX + b) &= E[(aX + b)^2] - (E[aX + b])^2 \\ &= E[a^2 X^2 + 2abX + b^2] - (aE[X] + b)^2 \\ &= a^2 E[X^2] + 2ab E[X] + b^2 \\ &\quad - [a^2 (E[X])^2 + 2abE[X] + b^2] \\ &= a^2 E[X^2] - a^2 (E[X])^2 \\ &= a^2 [E[X^2] - (E[X])^2] \\ &= a^2 \text{Var}(X)\end{aligned}$$

Bernoulli random variable

An indicator variable (a possibly biased coin flip) obeys a **Bernoulli distribution**. Bernoulli random variables can be 0 or 1.



$$X \sim \text{Ber}(p)$$

$$p_X(1) = p$$

$$p_X(0) = 1 - p \quad (0 \text{ elsewhere})$$



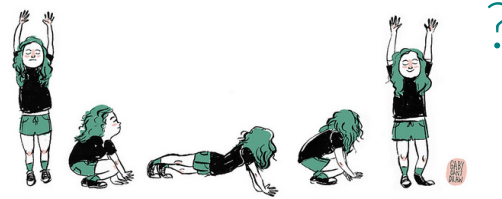
Bernoulli: Fact sheet



$$X \sim \text{Ber}(p)$$



probability of “success” (heads, ad click, ...)



PMF:

$$p_X(1) = p$$

$$p_X(0) = 1 - p \quad (0 \text{ elsewhere})$$

expectation:

$$E[X] = p$$

variance:

$$\text{Var}(X) = p(1 - p)$$

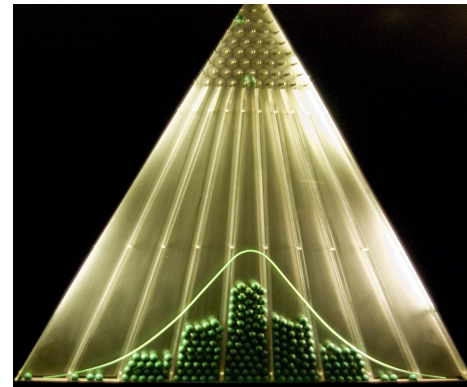
Binomial random variable

The **number of heads** on n (possibly biased) coin flips obeys a **binomial distribution**.



$$X \sim \text{Bin}(n, p)$$

$$p_X(k) = \begin{cases} \binom{n}{k} p^k (1-p)^{n-k} & \text{if } k \in \mathbb{N}, 0 \leq k \leq n \\ 0 & \text{otherwise} \end{cases}$$



Binomial: Fact sheet

number of trials (flips, program runs, ...)



$$X \sim \text{Bin}(n, p)$$



probability of "success" (heads, crash, ...)



PMF:

$$p_X(k) = \begin{cases} \binom{n}{k} p^k (1-p)^{n-k} & \text{if } k \in \mathbb{N}, 0 \leq k \leq n \\ 0 & \text{otherwise} \end{cases}$$

expectation:

$$E[X] = np$$

variance:

$$\text{Var}(X) = np(1-p)$$

note: $\text{Ber}(p) = \text{Bin}(1, p)$

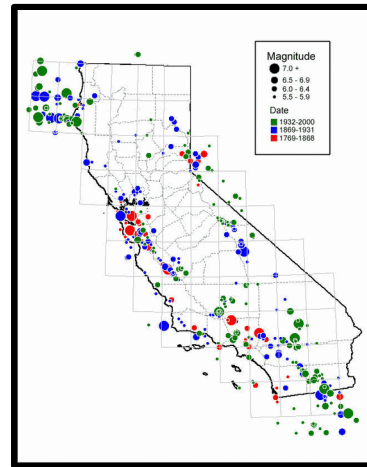
Poisson random variable

The **number of occurrences** of an event that occurs with **constant rate** λ (per unit time), in 1 unit of time, obeys a **Poisson distribution**.



$$X \sim \text{Poi}(\lambda)$$

$$p_X(k) = \begin{cases} e^{-\lambda} \frac{\lambda^k}{k!} & \text{if } x \in \mathbb{Z}, x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$



Poisson: Fact sheet



$$X \sim \text{Poi}(\lambda)$$



rate of events (requests, earthquakes,
chocolate chips, ...)
per unit time (hour, year, cookie, ...)

PMF:

$$p_X(k) = \begin{cases} e^{-\lambda} \frac{\lambda^k}{k!} & \text{if } k \in \mathbb{Z}, k \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

expectation:

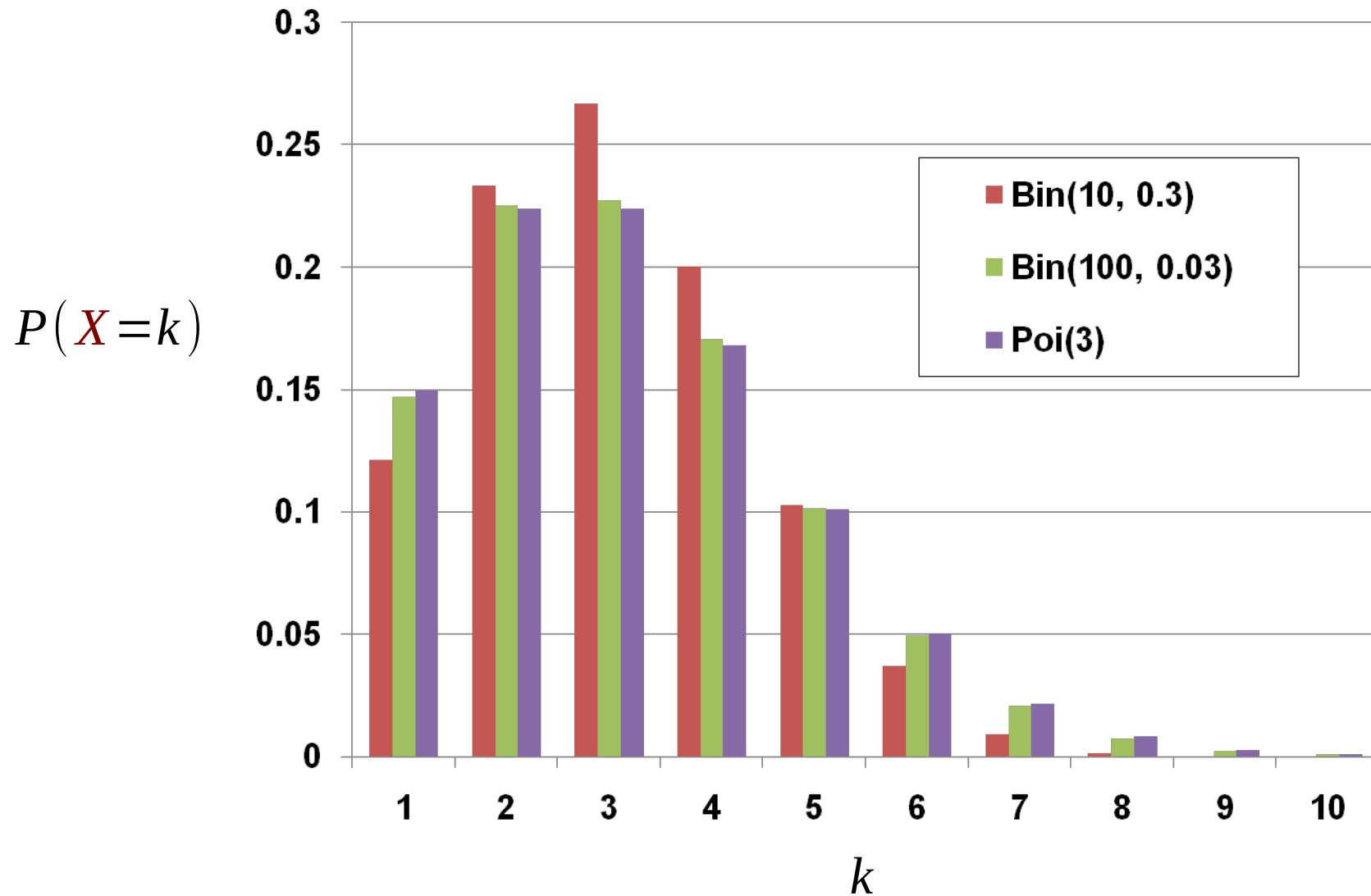
$$E[X] = \lambda$$

variance:

$$\text{Var}(X) = \lambda$$

Poisson approximation of binomial

$$\lambda = np$$



Geometric random variable

The **number of trials** it takes to get **one success**, if successes occur independently with probability p , obeys a **geometric distribution**.



$$X \sim \text{Geo}(p)$$

$$p_X(k) = \begin{cases} (1-p)^{k-1} \cdot p & \text{if } k \in \mathbb{Z}, k \geq 1 \\ 0 & \text{otherwise} \end{cases}$$



Geometric: Fact sheet



$$X \sim \text{Geo}(p)$$



probability of “success” (catch, heads, crash, ...)

PMF:

$$p_X(k) = \begin{cases} (1-p)^{k-1} \cdot p & \text{if } k \in \mathbb{Z}, k \geq 1 \\ 0 & \text{otherwise} \end{cases}$$

CDF:

$$F_X(k) = \begin{cases} 1 - (1-p)^k & \text{if } k \in \mathbb{Z}, k \geq 1 \\ 0 & \text{otherwise} \end{cases}$$

expectation:

$$E[X] = \frac{1}{p}$$

variance:

$$\text{Var}(X) = \frac{1-p}{p^2}$$

Negative binomial random variable

The **number of trials** it takes to get r successes, if successes occur independently with probability p , obeys a **negative binomial distribution**.



$$X \sim \text{NegBin}(r, p)$$

$$p_X(n) = \begin{cases} \binom{n-1}{r-1} p^r (1-p)^{n-r} & \text{if } n \in \mathbb{Z}, n \geq r \\ 0 & \text{otherwise} \end{cases}$$



Negative binomial: Fact sheet

number of **successes** (heads, crash, ...)

$$X \sim \text{NegBin}(r, p)$$

number of **trials** (flips, program runs, ...)

probability of “success”

$$\text{PMF: } p_X(n) = \begin{cases} \binom{n-1}{r-1} p^r (1-p)^{n-r} & \text{if } n \in \mathbb{Z}, n \geq r \\ 0 & \text{otherwise} \end{cases}$$

expectation: $E[X] = \frac{r}{p}$

variance: $\text{Var}(X) = \frac{r(1-p)}{p^2}$

note:

$$\text{Geo}(p) = \text{NegBin}(1, p)$$

Continuous random variables

A **continuous** random variable has a value that's a **real number** (not necessarily an integer).

Replace sums with **integrals**!



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

$$F_X(a) = \int_{x=-\infty}^a dx f_X(x)$$

Probability density function

The **probability density function** (PDF) of a continuous random variable represents the relative likelihood of various values.

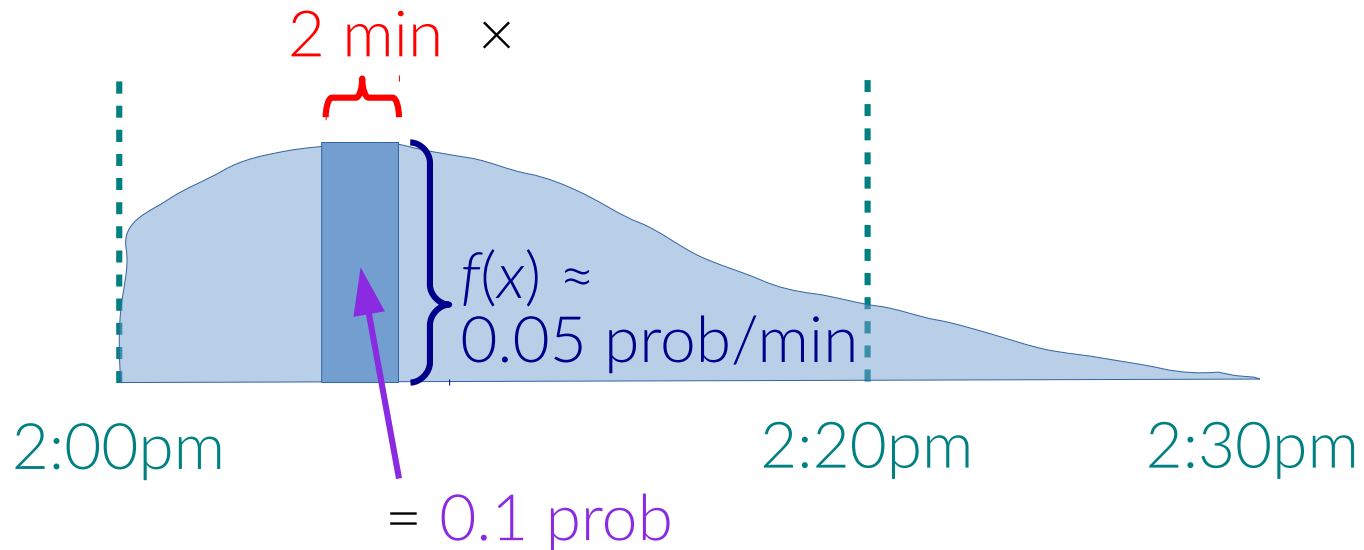
Units of probability *divided by units of X*.
Integrate it to get probabilities!



$$P(a < X \leq b) = \int_{x=a}^b dx \boxed{f_X(x)}$$

$f(x)$ is not a probability

Rather, it has “units” of probability
divided by units of X .



Uniform random variable

A **uniform** random variable is **equally likely** to be any value in a single real number interval.

$$X \sim \text{Uni}(\alpha, \beta)$$

$$f_X(x) = \begin{cases} \frac{1}{\beta - \alpha} & \text{if } x \in [\alpha, \beta] \\ 0 & \text{otherwise} \end{cases}$$



Uniform: Fact sheet



minimum value
↓
 $X \sim \text{Uni}(\alpha, \beta)$
↑
maximum value

PDF: $f_X(x) = \begin{cases} \frac{1}{\beta - \alpha} & \text{if } x \in [\alpha, \beta] \\ 0 & \text{otherwise} \end{cases}$

CDF: $F_X(x) = \begin{cases} \frac{x - \alpha}{\beta - \alpha} & \text{if } x \in [\alpha, \beta] \\ 1 & \text{if } x > \beta \\ 0 & \text{otherwise} \end{cases}$

expectation: $E[X] = \frac{\alpha + \beta}{2}$

variance: $\text{Var}(X) = \frac{(\beta - \alpha)^2}{12}$

Exponential random variable

An **exponential** random variable is the **amount of time until the first event** when events occur as in the Poisson distribution.



$$X \sim \text{Exp}(\lambda)$$

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



Exponential: Fact sheet



rate of events per unit time

$$X \sim \text{Exp}(\lambda)$$

time until first event



$$\text{PDF: } f_X(x) = \begin{cases} \lambda e^{-\lambda x} & \text{if } x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

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

$$\text{expectation: } E[X] = \frac{1}{\lambda}$$

$$\text{variance: } \text{Var}(X) = \frac{1}{\lambda^2}$$

A grid of random variables

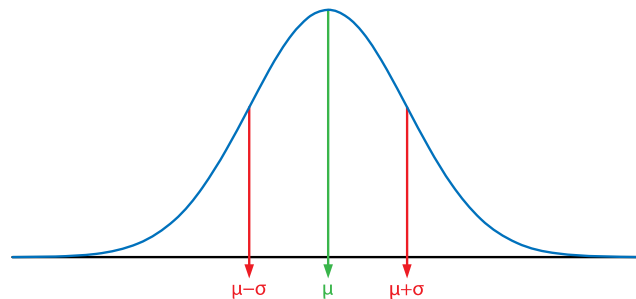
	number of successes	time to get successes	
One trial	$X \sim \text{Ber}(p)$  $n = 1$	$X \sim \text{Geo}(p)$  $r = 1$	One success
Several trials	$X \sim \text{Bin}(n, p)$	$X \sim \text{NegBin}(r, p)$	Several successes
Interval of time	$X \sim \text{Poi}(\lambda)$	$X \sim \text{Exp}(\lambda)$	One success after interval of time

Normal random variable

An **normal** (= **Gaussian**) random variable is a good approximation to many other distributions. It often results from **sums or averages** of independent random variables.



$$X \sim N(\mu, \sigma^2)$$
$$f_X(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma} \right)^2}$$



Normal: Fact sheet



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

mean



variance (σ = standard deviation)

PDF: $f_X(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$

CDF: $F_X(x) = \Phi\left(\frac{x-\mu}{\sigma}\right) = \int_{-\infty}^x dx f_X(x)$
(no closed form)

expectation: $E[X] = \mu$

variance: $\text{Var}(X) = \sigma^2$

The Standard Normal

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

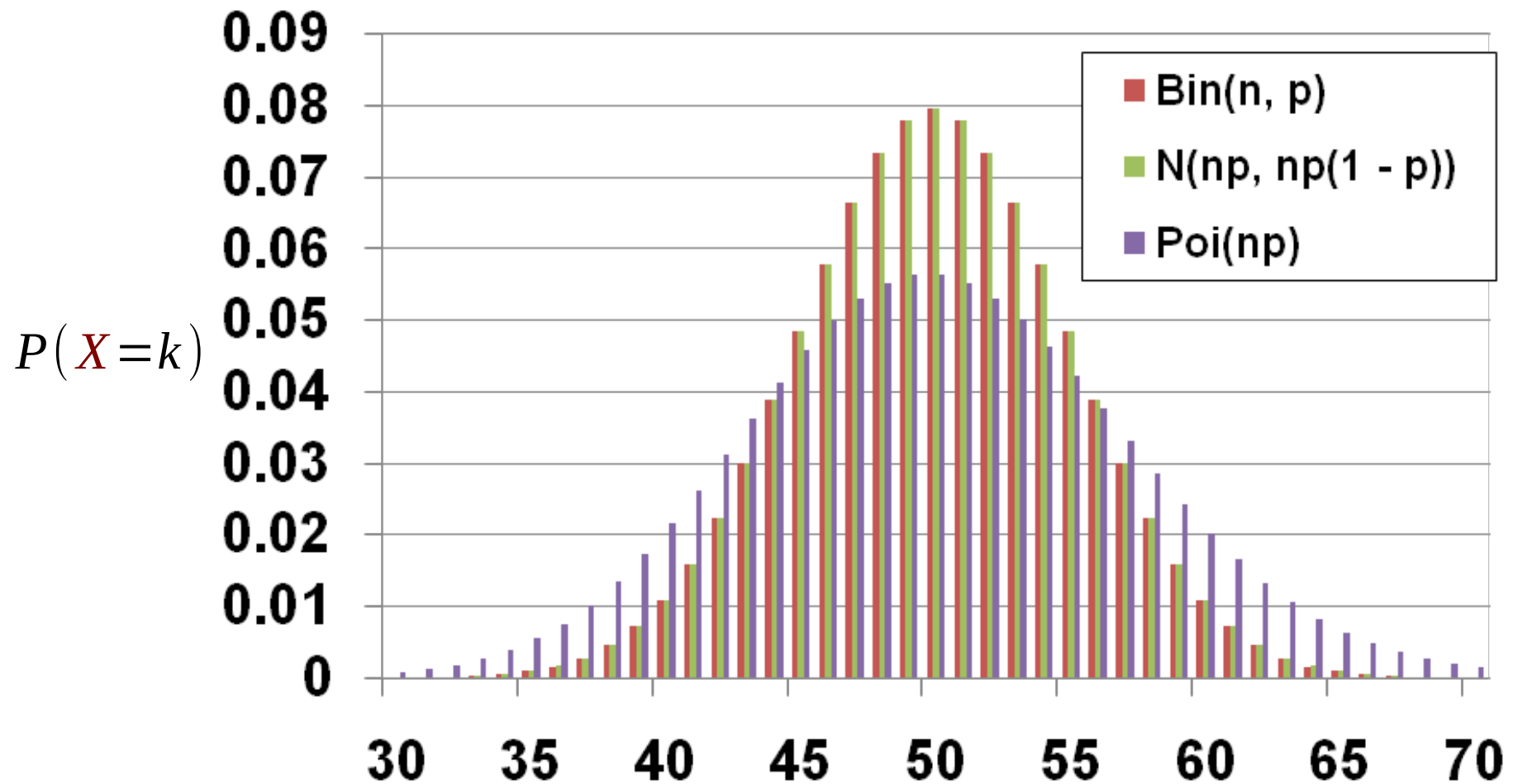
↑ ↑
μ σ²

$$X \sim N(\mu, \sigma^2) \rightarrow X = \sigma Z + \mu$$
$$Z = \frac{X - \mu}{\sigma}$$

$$\Phi(z) = F_Z(z) = P(Z \leq z)$$

Normal approximation to binomial

large n , medium p



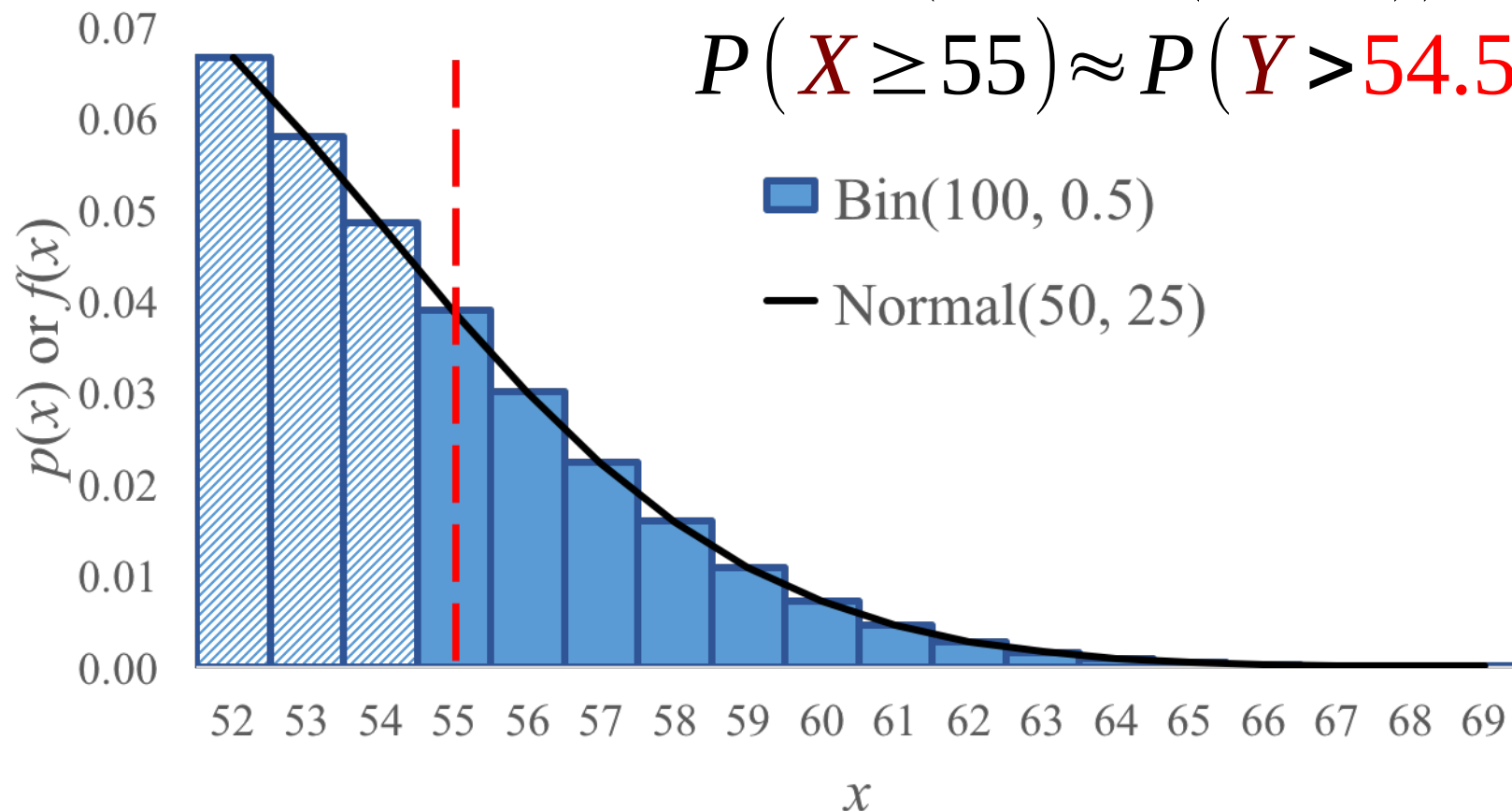
$$\text{Bin}(n, p) \approx N(\mu, \sigma^2) \quad k$$

Continuity correction

$$X \sim \text{Bin}(n, p)$$

$$Y \sim N(np, np(1-p))$$

$$P(X \geq 55) \approx P(Y > 54.5)$$



When approximating a **discrete** distribution with a **continuous** distribution, adjust the bounds by 0.5 to account for the missing half-bar.

Joint distributions

A **joint distribution** combines multiple random variables. Its PDF or PMF gives the probability or relative likelihood of **both** random variables taking on specific values.



$$p_{X,Y}(a,b) = P(X=a, Y=b)$$

Joint probability mass function

A joint probability mass function gives the probability of **more than one** discrete random variable **each** taking on a specific value (an AND of the 2+ values).



$$p_{X,Y}(a,b) = P(X=a, Y=b)$$

		Y		
		0	1	2
X	0	0.05	0.20	0.10
	1	0.10	0.10	0.10
	2	0.05	0.10	0.20

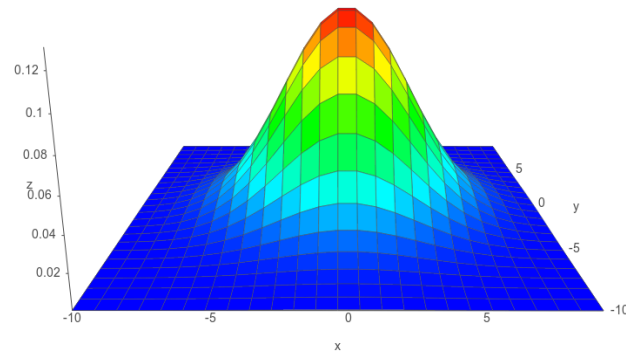
Joint probability density function

A joint probability density function gives the relative likelihood of **more than one** continuous random variable **each** taking on a specific value.



$$P(a_1 < X \leq a_2, b_1 < Y \leq b_2) =$$

$$\int_{a_1}^{a_2} dx \int_{b_1}^{b_2} dy f_{X,Y}(x, y)$$



Marginalization

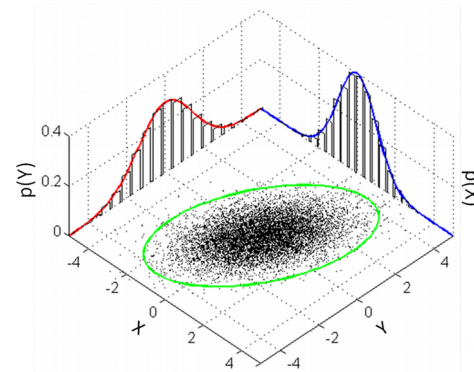
Marginal probabilities give the distribution of **a subset of the variables** (often, just one) of a joint distribution.

Sum/integrate over the variables you don't care about.



$$p_X(a) = \sum_y p_{X,Y}(a, y)$$

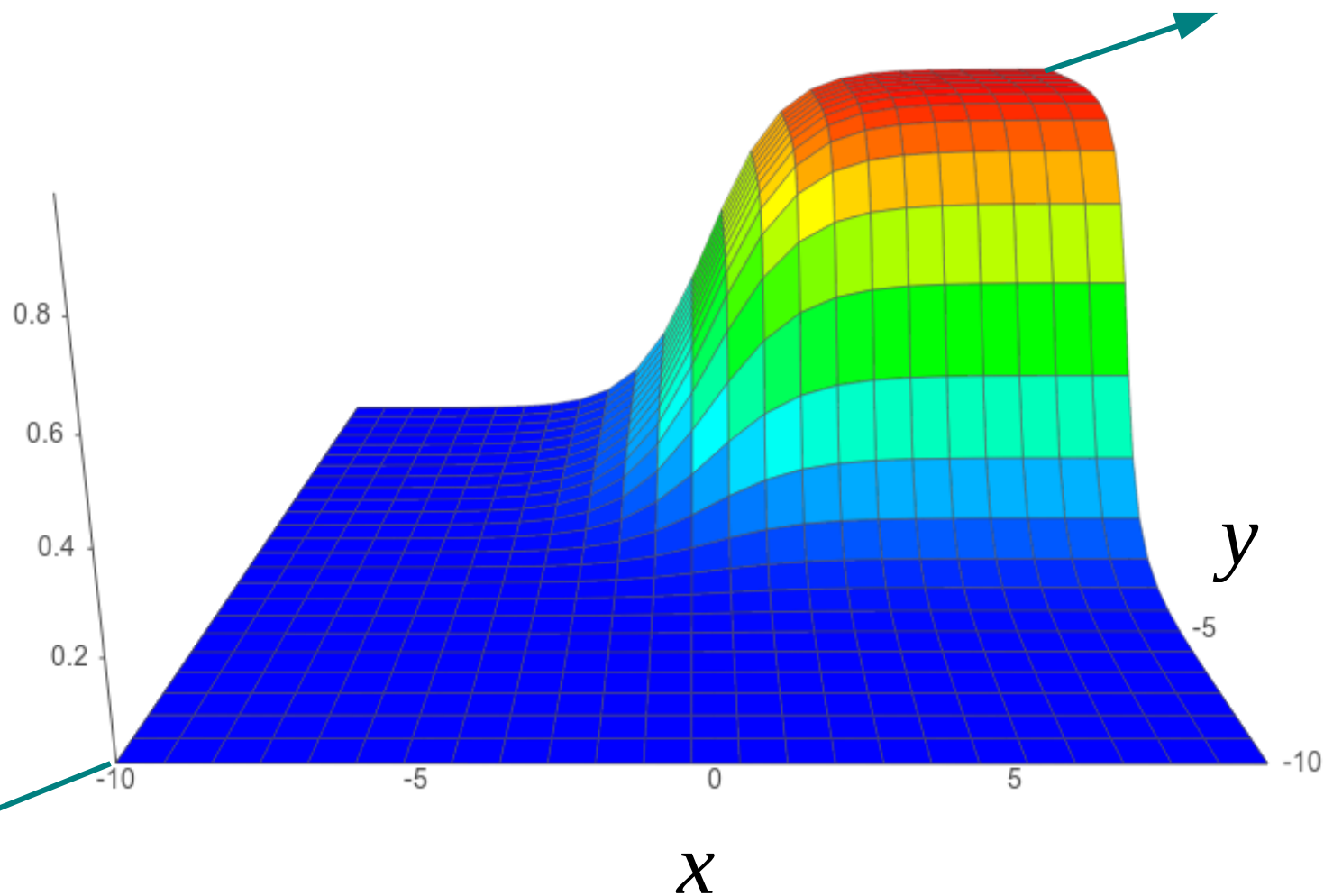
$$f_X(a) = \int_{-\infty}^{\infty} dy f_{X,Y}(a, y)$$



Joint cumulative distribution function

$$F_{X,Y}(x,y) = P(X \leq x, Y \leq y)$$

to 1 as
 $x \rightarrow +\infty,$
 $y \rightarrow +\infty$



to 0 as
 $x \rightarrow -\infty,$
 $y \rightarrow -\infty$

Multinomial random variable

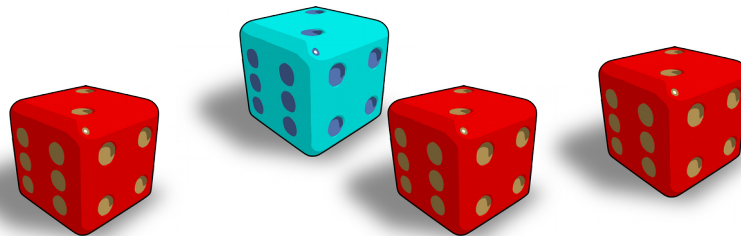
An **multinomial** random variable records the number of times each outcome occurs, when an experiment with multiple outcomes (e.g. die roll) is run multiple times.



vector!

$$X_1, \dots, X_m \sim \text{MN}(n, p_1, p_2, \dots, p_m)$$

$$P(X_1 = c_1, X_2 = c_2, \dots, X_m = c_m) = \binom{n}{c_1, c_2, \dots, c_m} p_1^{c_1} p_2^{c_2} \dots p_m^{c_m}$$



Independence of discrete random variables

Two random variables are **independent** if knowing the value of one tells you nothing about the value of the other (for **all** values!).



$X \perp Y$ iff $\forall x, y$:

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

- or -

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

Independence of continuous random variables

Two random variables are **independent** if knowing the value of one tells you nothing about the value of the other (for **all** values!).



$X \perp Y$ iff $\forall x, y$:

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

- or -

$$f_{X,Y}(x, y) = g(x) h(y)$$

- or -

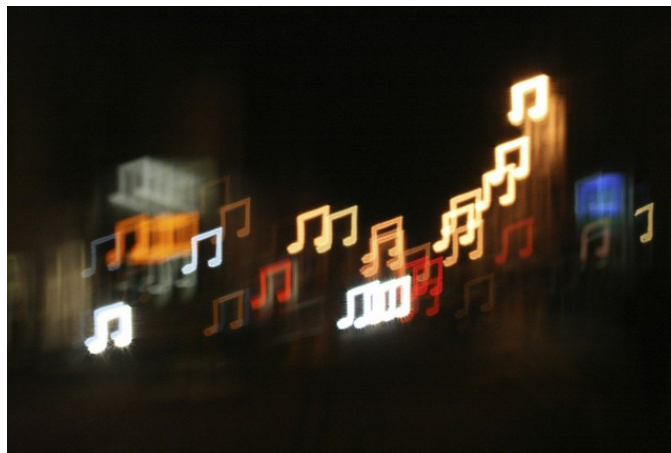
$$F_{X,Y}(x, y) = F_X(x) F_Y(y)$$

Convolution

A **convolution** is the distribution of the **sum** of two independent random variables.



$$f_{X+Y}(a) = \int_{-\infty}^{\infty} dy f_X(a-y) f_Y(y)$$



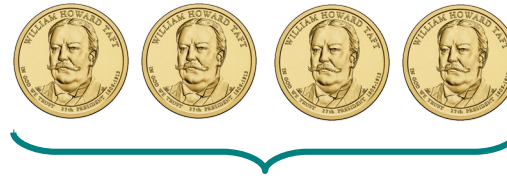
Sum of independent binomials



n flips

X : number of heads
in first n flips

$$X \sim \text{Bin}(n, p)$$



m flips

Y : number of heads
in next m flips

$$Y \sim \text{Bin}(m, p)$$

$$X + Y \sim \text{Bin}(n + m, p)$$

More generally:

$$\begin{array}{l} X_i \sim \text{Bin}(n_i, p) \\ \text{all } X_i \text{ independent} \end{array} \Rightarrow \sum_{i=1}^N X_i \sim \text{Bin}\left(\sum_{i=1}^N n_i, p\right)$$

Sum of independent Poissons



λ_1 chips/cookie

X : number of chips
in first cookie

$$X \sim \text{Poi}(\lambda_1)$$



λ_2 chips/cookie

Y : number of chips
in second cookie

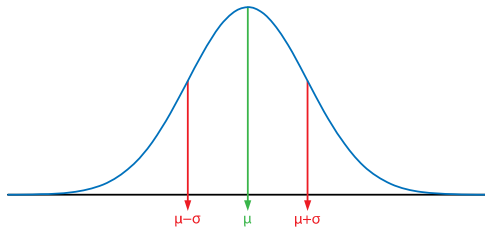
$$Y \sim \text{Poi}(\lambda_2)$$

$$X + Y \sim \text{Poi}(\lambda_1 + \lambda_2)$$

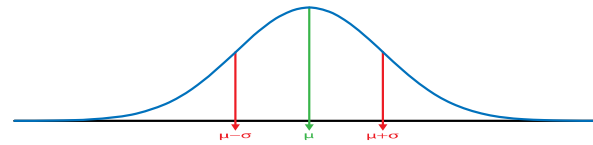
More generally:

$$\begin{array}{l} X_i \sim \text{Poi}(\lambda_i) \Rightarrow \\ \text{all } X_i \text{ independent} \end{array} \quad \sum_{i=1}^N X_i \sim \text{Poi}\left(\sum_{i=1}^N \lambda_i\right)$$

Sum of independent normals



$$X \sim N(\mu_1, \sigma_1^2)$$



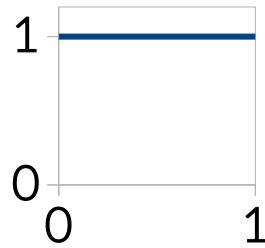
$$Y \sim N(\mu_2, \sigma_2^2)$$

$$X + Y \sim N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$$

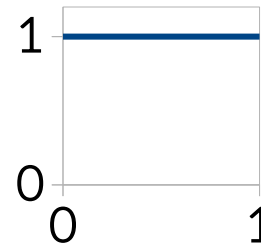
More generally:

$$\begin{array}{l} X_i \sim N(\mu_i, \sigma_i^2) \\ \text{all } X_i \text{ independent} \end{array} \Rightarrow \sum_{i=1}^N X_i \sim N\left(\sum_{i=1}^N \mu_i, \sum_{i=1}^N \sigma_i^2\right)$$

Sum of independent uniforms



$$X \sim \text{Uni}(0, 1)$$



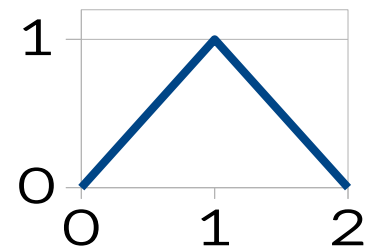
$$Y \sim \text{Uni}(0, 1)$$

$$\begin{aligned} f_{X+Y}(a) &= \int_{-\infty}^{\infty} dy f_X(a-y) f_Y(y) \\ &= \int_0^1 dy f_X(a-y) f_Y(y) \end{aligned}$$

Case 1: if $0 \leq a \leq 1$, then we need $0 \leq y \leq a$ (for $a - y$ to be in $[0, 1]$)

Case 2: if $1 \leq a \leq 2$, then we need $a - 1 \leq y \leq 1$

$$= \begin{cases} \int_0^a dy \cdot 1 = a & 0 \leq a \leq 1 \\ \int_{a-1}^1 dy \cdot 1 = 2 - a & 1 \leq a \leq 2 \\ 0 & \text{otherwise} \end{cases}$$



Discrete conditional distributions

The value of a random variable, conditioned on the value of some other random variable, has a probability distribution.



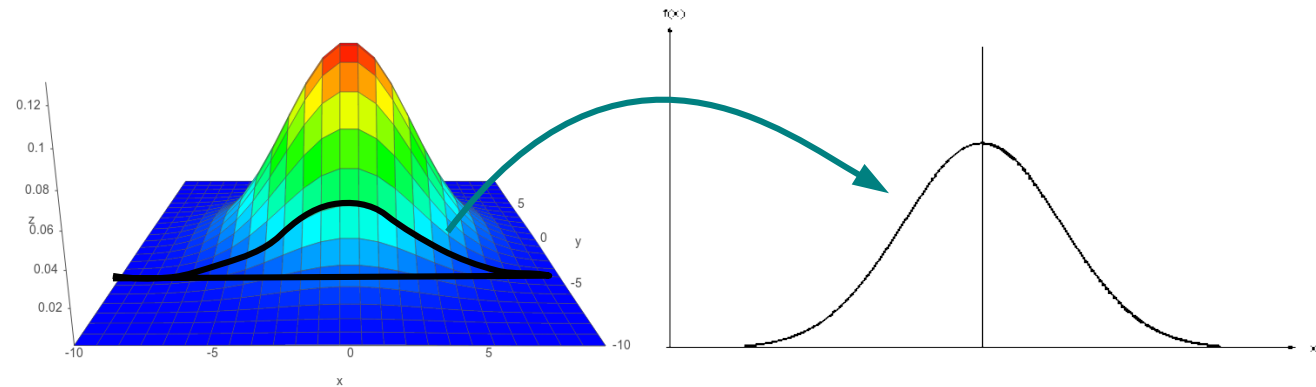
$$p_{X|Y}(x, y) = \frac{P(X=x, Y=y)}{P(Y=y)} \\ = \frac{p_{X,Y}(x, y)}{p_Y(y)}$$

PDF	Single	In a relationship	It's complicated / Other	TOTALS
Freshman	0.00	0.00	0.00	0.00
Sophomore	0.06	0.00	0.00	0.06
Junior	0.19	0.19	0.13	0.50
Senior	0.00	0.00	0.00	0.00
Grad student / Other	0.38	0.06	0.00	0.44
TOTALS	0.63	0.25	0.13	1.00

Continuous conditional distributions

The value of a random variable, **conditioned on the value of some other random variable**, has a probability distribution.

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



Ratios of continuous probabilities

The probability of an exact value for a continuous random variable is 0.

But ratios of these probabilities are still well-defined!

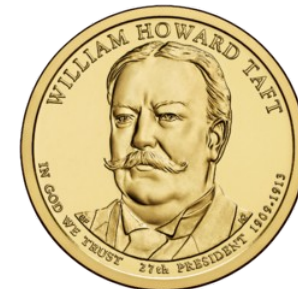
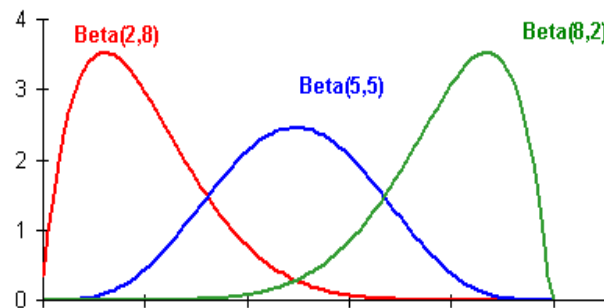
$$\frac{P(X=a)}{P(X=b)} = \frac{f_X(a)}{f_X(b)}$$

Beta random variable

An **beta** random variable models the **probability** of a trial's success, given previous trials. The PDF/CDF let you compute **probabilities of probabilities!**

$$X \sim \text{Beta}(a, b)$$

$$f_X(x) = \begin{cases} C x^{a-1} (1-x)^{b-1} & \text{if } 0 < x < 1 \\ 0 & \text{otherwise} \end{cases}$$



Beta: Fact sheet



$$X \sim \text{Beta}(a, b)$$

number of successes + 1

probability of success

number of failures + 1

$$\text{PDF: } f_X(x) = \begin{cases} C x^{a-1} (1-x)^{b-1} & \text{if } 0 < x < 1 \\ 0 & \text{otherwise} \end{cases}$$

expectation: $E[X] = \frac{a}{a+b}$

variance: $\text{Var}(X) = \frac{ab}{(a+b)^2(a+b+1)}$

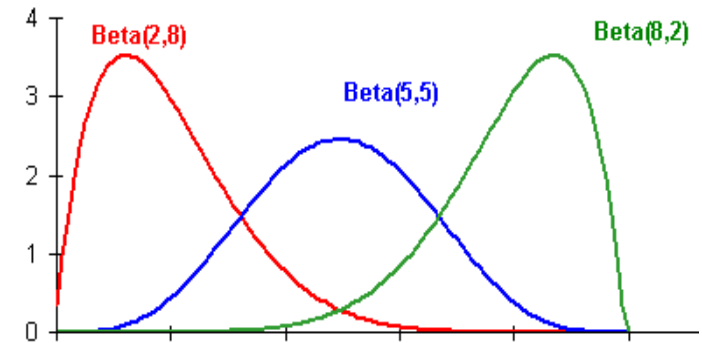
Subjective priors

$X | A \sim \text{Beta}(a + 1, N - a + 1)$
“posterior”

$X \sim \text{Beta}(1, 1)$
“prior”

$$f_{X|A}(x|a) = \frac{P(A=a|X=x) f_X(x)}{P(A=a)}$$

How did we decide on
Beta(1, 1) for the prior?



Beta(1, 1): “we haven’t seen any rolls yet.”

Beta(4, 1): “we’ve seen 3 sixes and 0 non-sixes.”

Beta(2, 6): “we’ve seen 1 six and 5 non-sixes.”

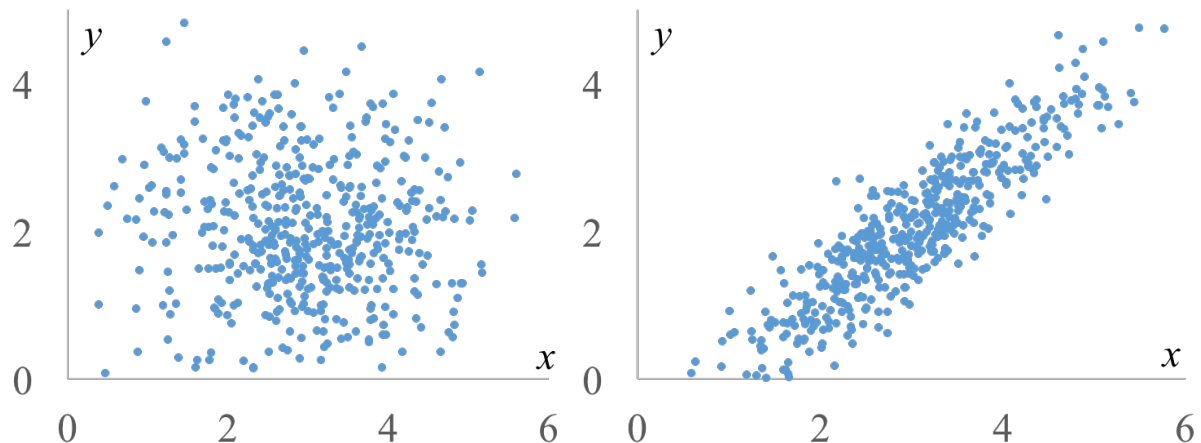
Beta prior = “imaginary” previous trials

Covariance

The **covariance** of two variables is a measure of how much they **vary together**.



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



Expectation of a product

If two random variables are **independent**, then the **expectation of their product** equals the **product of their expectations**.



$$X \perp Y \Rightarrow$$

$$E[XY] = E[X]E[Y]$$

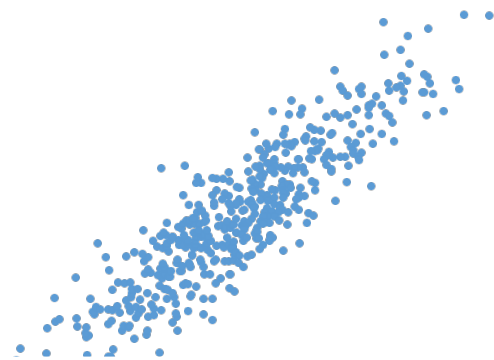
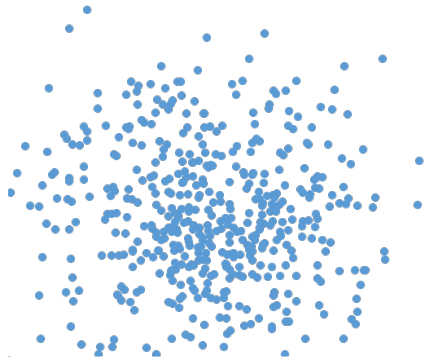
$$E[g(X)h(Y)] = E[g(X)]E[h(Y)]$$

Correlation

The **correlation** of two variables is a measure of the **linear dependence** between them, scaled to always take on values between -1 and 1.



$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}}$$



Conditional expectation

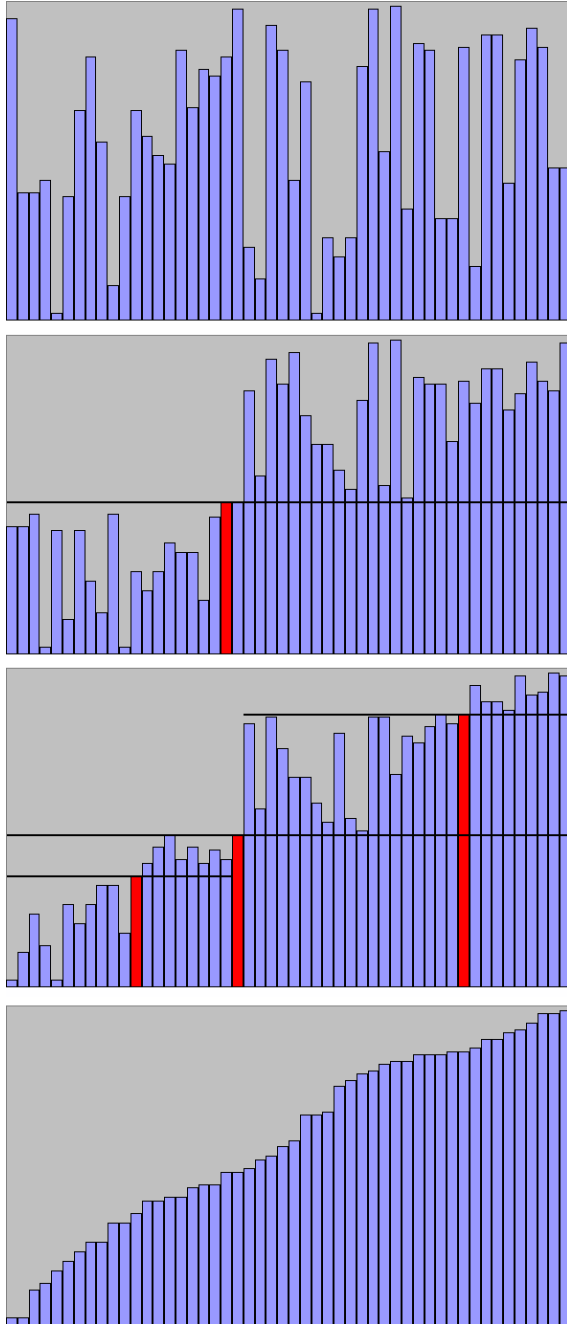
One can compute the **expectation** of a random variable while **conditioning** on the values of other random variables.



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

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

Quicksort



You've been told Quicksort is $O(n \log n)$, "average case".

Now you get to find out why!

Quicksort's ordinary life

Let X = number of comparisons to the pivot.

What is $E[X]$? expected number of events = indicator variables!

1	2	3	4	5	6	7	8
---	---	---	---	---	---	---	---

Y_1 Y_2 ... Y_n

Define $Y_1 \dots Y_n$ = elements in sorted order.

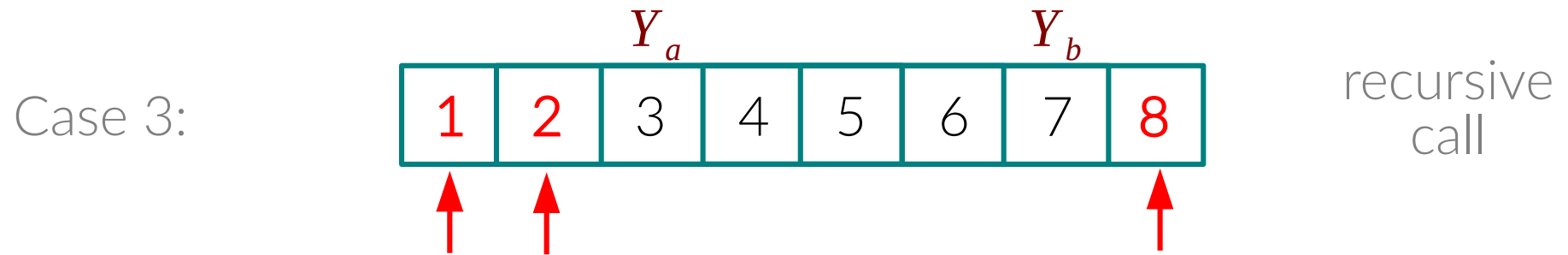
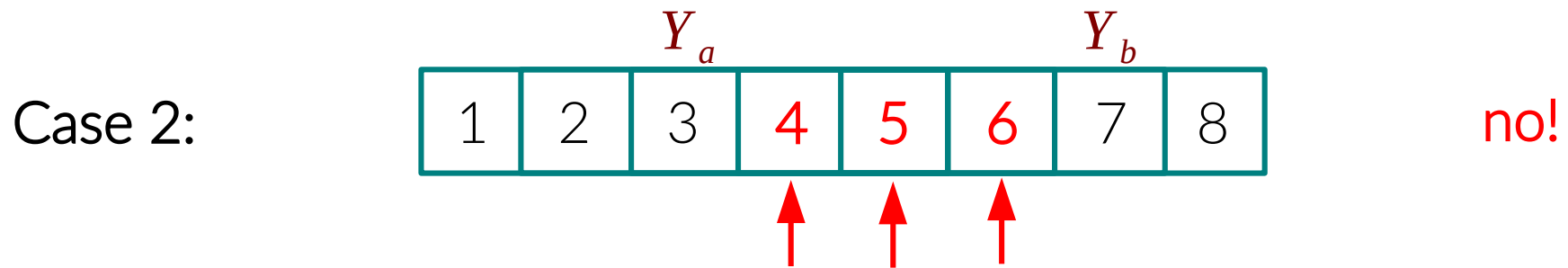
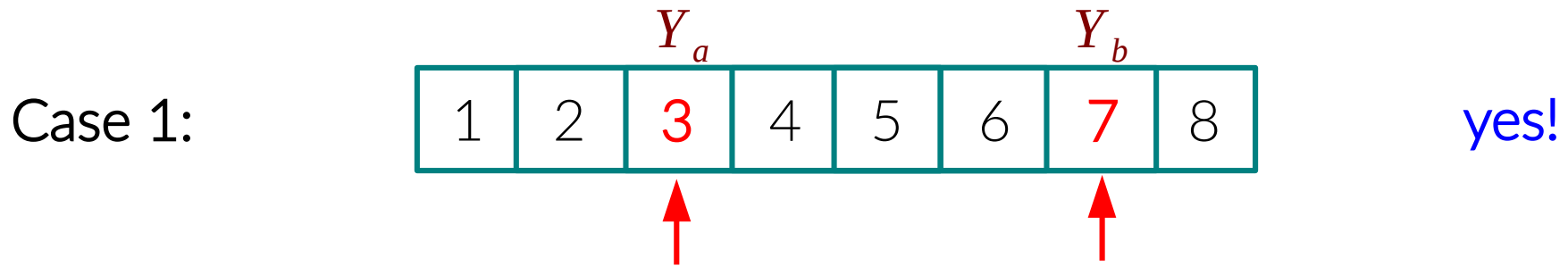
Indicator variables $I_{ab} = 1$ if Y_a and Y_b are ever compared.

$$\begin{aligned} E[X] &= E \left[\sum_{a=1}^{n-1} \sum_{b=a+1}^n I_{ab} \right] = \sum_{a=1}^{n-1} \sum_{b=a+1}^n E[I_{ab}] \\ &= \sum_{a=1}^{n-1} \sum_{b=a+1}^n P(Y_a \text{ and } Y_b \text{ ever compared}) \end{aligned}$$

unique pairs

Shall I compare thee to...

$P(Y_a \text{ and } Y_b \text{ ever compared}) = ?$



$$\therefore P(Y_a \text{ and } Y_b \text{ ever compared}) = \frac{2}{b-a+1}$$

The home stretch

$$E[X] = \sum_{a=1}^{n-1} \sum_{b=a+1}^n P(Y_a \text{ and } Y_b \text{ ever compared})$$

$$= \sum_{a=1}^{n-1} \sum_{b=a+1}^n \frac{2}{b-a+1}$$

$$\approx \sum_{a=1}^{n-1} 2 \ln(n-a+1)$$

$$\approx \int_{a=1}^{n-1} da 2 \ln(n-a+1)$$

$$\sum_{b=a+1}^n \frac{2}{b-a+1} \approx \int_{b=a+1}^n db \frac{2}{b-a+1}$$

$$= [2 \ln(b-a+1)]_{b=a+1}^n$$

$$= 2 \ln(n-a+1) - 2 \ln 2$$

$$\approx 2 \ln(n-a+1) \quad \text{for large } n$$

$$= -2 \int_{y=n}^2 dy \ln y$$

$$= -2 [y \ln y - y]_{y=n}^2$$

$$= -2 [(\cancel{2 \ln 2} - 2) - (n \ln n - \cancel{n})]$$

constants

$$u = \ln y \quad v = y$$

$$du = \frac{1}{y} dy \quad dv = dy$$

$$\int u dv = uv - \int v du$$

$$\int \ln y dy = y \ln y - \int y \frac{1}{y} dy$$

$$= y \ln y - y + C$$

lower-order term

$$= O(n \ln n)$$



Variance of a sum

The **variance of a sum** of random variables is equal to the **sum of pairwise covariances** (*including* variances and double-counted pairs).



$$\begin{aligned}\text{Var}\left(\sum_{i=1}^n X_i\right) &= \text{Cov}\left(\sum_{i=1}^n X_i, \sum_{j=1}^n X_j\right) \\ &= \sum_{i=1}^n \text{Var}(X_i) + 2 \sum_{i=1}^n \sum_{j=i+1}^n \text{Cov}(X_i, X_j)\end{aligned}$$

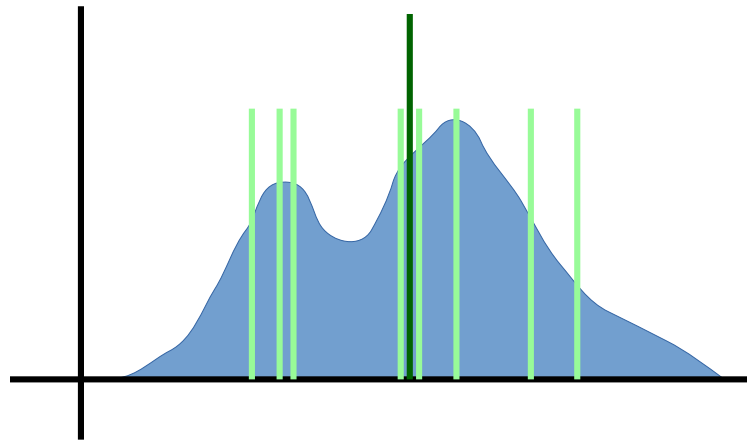
note: independent \Rightarrow Cov = 0

Sample mean

A **sample mean** is an **average** of random variables drawn (usually independently) from the **same distribution**.



$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$



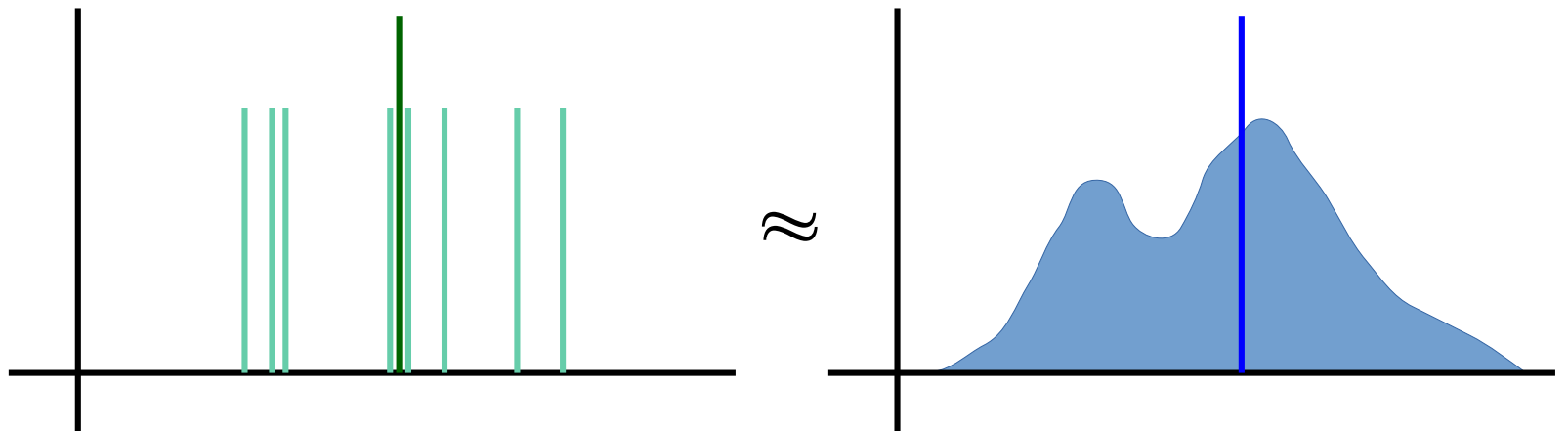
Parameter estimation

Sometimes we **don't know** things like the expectation and variance of a distribution; we have to **estimate** them from incomplete information.



$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \quad S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

$$\hat{\theta} = \arg \max_{\theta} \text{LL}(\theta)$$

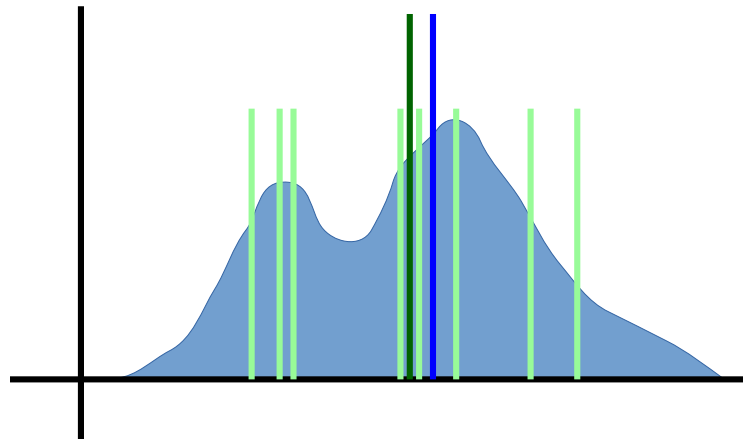


Unbiased estimator

An **unbiased estimator** is a random variable that has **expectation** equal to the quantity you are estimating.



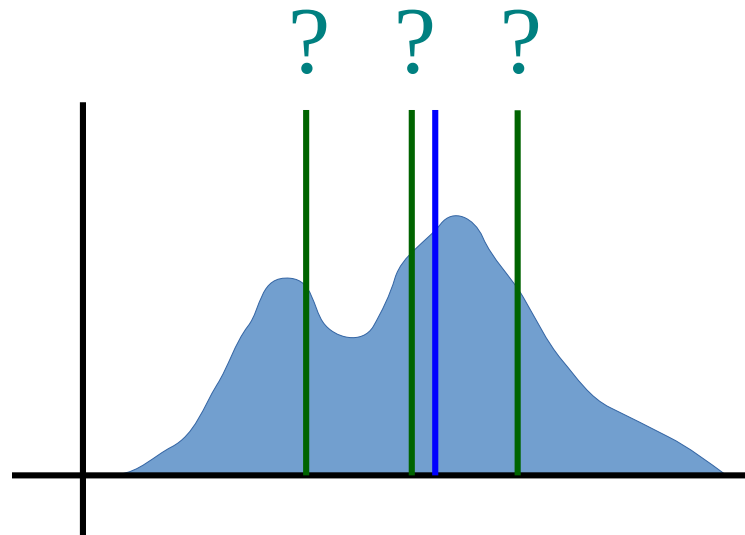
$$E[\bar{X}] = \mu = E[X_i]$$



Variance of the sample mean

The **sample mean** is a random variable; it can differ among samples. That means it has a **variance**.

$$\text{Var}(\bar{X}) = \frac{\sigma^2}{n}$$

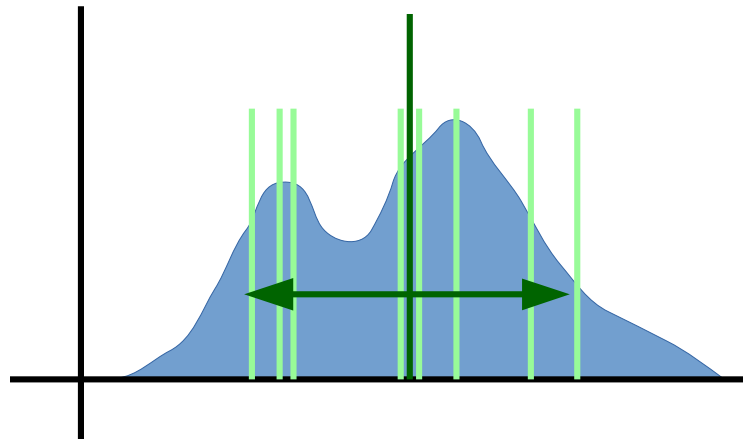


Sample variance

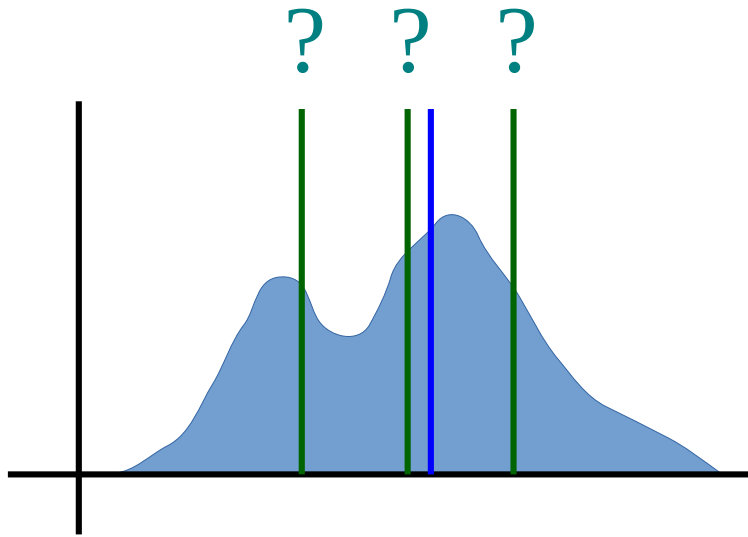
Samples can be used to **estimate the variance** of the original distribution.



$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$



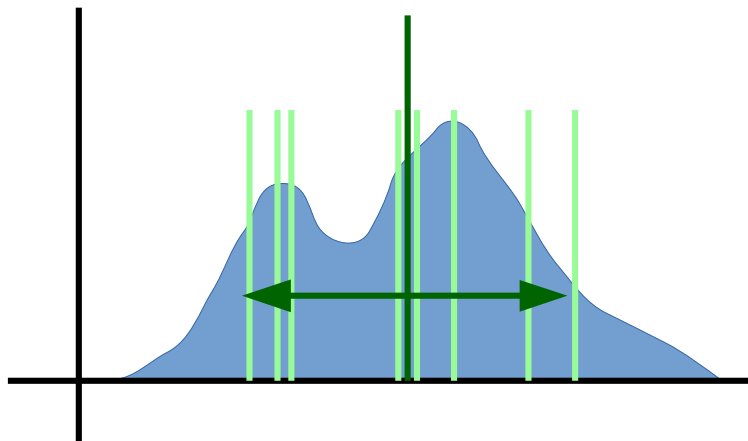
Variance of the sample mean



- Is a single number
- Shrinks with number of samples $\left(= \frac{\sigma^2}{n} \right)$
- Measures the stability of an estimate

vs.

Sample variance



- Is a random variable
- Constant with number of samples $\left(\approx \sigma^2 \right)$
- Is an estimate (of a variance) itself

p-values

A **p-value** gives the probability of an extreme result, assuming that any extremeness is due to chance.



$$p = P(|\bar{X} - \mu| > d | H_0)$$



Bootstrapping



Bootstrapping allows you to compute complicated statistics from samples using simulation.

Bootstrap for p-values

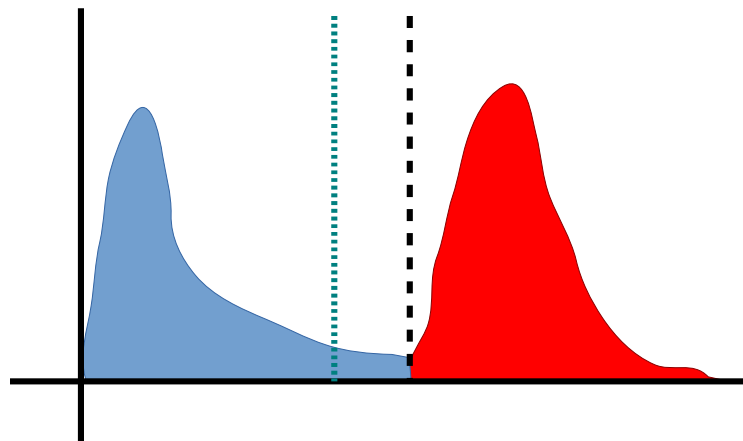
```
def pvalue_bootstrap(sample1, sample2):  
    n = len(sample1)  
    m = len(sample2)  
    observed_diff = abs(np.mean(sample2) -  
                        np.mean(sample1))  
    universal_pmf = sample1 + sample2  
    count_extreme = 0  
    for i in range(10000):  
        resample1 = np.random.choice(universal_pmf, n)  
        resample2 = np.random.choice(universal_pmf, m)  
        new_diff = abs(np.mean(resample2) -  
                        np.mean(resample1))  
        if new_diff >= observed_diff:  
            count_extreme += 1  
    return count_extreme / 10000.
```

Markov's inequality

Knowing the **expectation** of a **non-negative** random variable lets you bound the probability of **high** values for that variable.



$$X \geq 0 \Rightarrow P(X \geq a) \leq \frac{E[X]}{a}$$

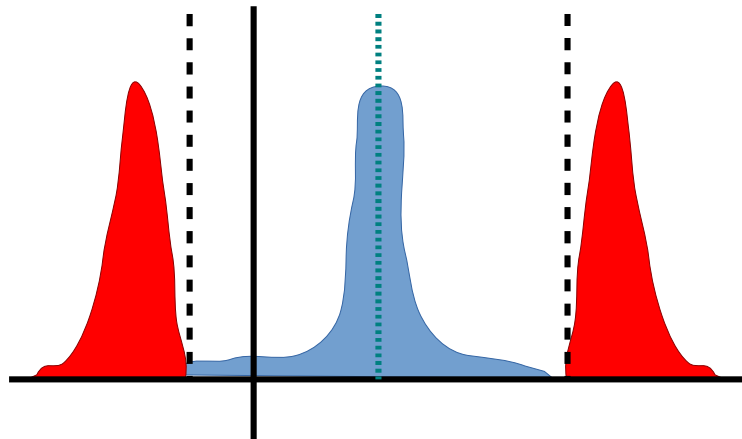


Chebyshev's inequality

Knowing the expectation **and variance** of a random variable lets you bound the probability of **extreme** values for that variable.



$$P(|X - \mu| \geq k) \leq \frac{\sigma^2}{k^2}$$



One-sided Chebyshev's inequality

$$P(X \geq \mu + a) \leq \frac{\sigma^2}{\sigma^2 + a^2}$$

$$P(X \leq \mu - a) \leq \frac{\sigma^2}{\sigma^2 + a^2}$$

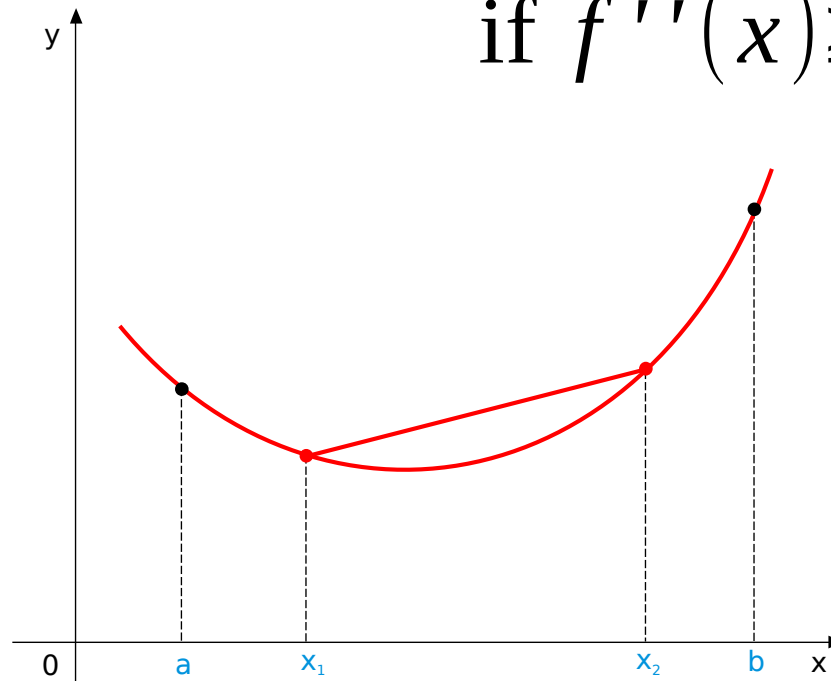
Jensen's inequality

The expectation of a **convex function** of a random variable can't be less than the value of the function applied to the expectation.



$$E[f(X)] \geq f(E[X])$$

$$\text{if } f''(x) \geq 0$$



Law of large numbers

A sample mean will converge to the true mean if you take a large enough sample.



$$\lim_{n \rightarrow \infty} P(|\bar{X} - \mu| \geq \varepsilon) = 0$$

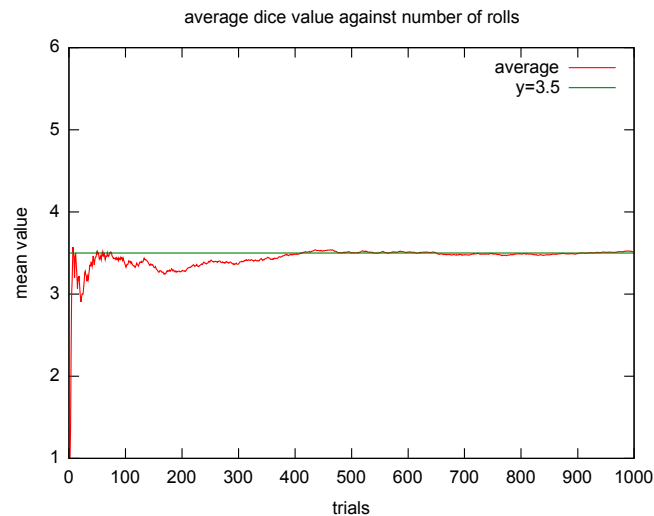
$$P\left(\lim_{n \rightarrow \infty} (\bar{X}) = \mu\right) = 1$$

Consistent estimator

An **consistent estimator** is a random variable that has a **limit** (as number of samples gets large) equal to the quantity you are estimating.



$$\lim_{n \rightarrow \infty} P(|\hat{\theta} - \theta| < \varepsilon) = 1$$



Review: Central limit theorem

Sums and averages of IID random variables are normally distributed.



$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

$$Y = n \bar{X} = \sum_{i=1}^n X_i \sim N(n\mu, n\sigma^2)$$

Easily-confused principles

Constant multiple
of a normal

Sum of identical
normals

CLT

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

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

$$X_i \sim ???$$

(independent
& identical)

(independent
& identical)



$$nX \sim N(n\mu, n^2\sigma^2)$$

$$\sum_{i=1}^n X_i \sim N(n\mu, n\sigma^2)$$

$$\sum_{i=1}^n X_i \sim N(n\mu, n\sigma^2)$$

(exactly)

(approximately,
for large n)

Parameters

θ

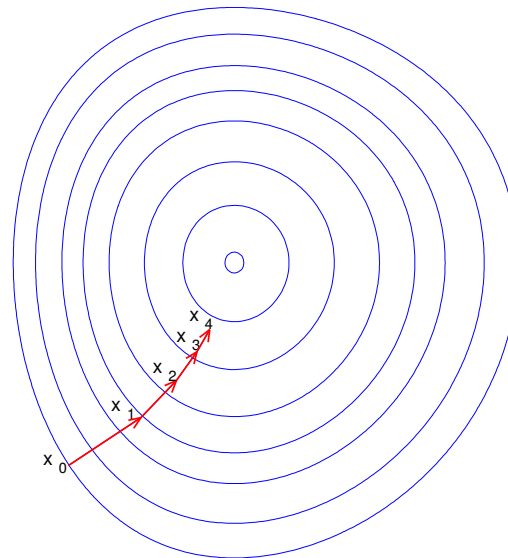
$X \sim$	Ber(p)	$\theta = p$
	Poi(λ)	$\theta = \lambda$
	Uni(a, b)	$\theta = [a, b]$
	N(μ, σ^2)	$\theta = [\mu, \sigma^2]$

Maximum likelihood estimation

Choose parameters that **maximize** the likelihood (**joint probability given parameters**) of the example data.



$$\hat{\theta} = \arg \max_{\theta} LL(\theta)$$



How to: MLE

1. Compute the likelihood.

$$L(\theta) = P(X_1, \dots, X_m | \theta)$$

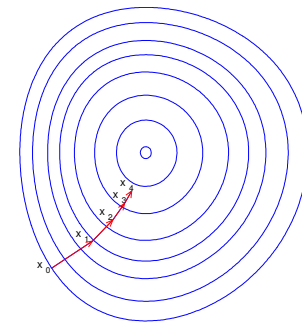


2. Take its log.

$$LL(\theta) = \log L(\theta)$$

3. Maximize this as a function of the parameters.

$$\frac{d}{d\theta} LL(\theta) = 0$$



Maximum likelihood for Bernoulli

The maximum likelihood p for Bernoulli random variables is the sample mean.



$$\hat{p} = \frac{1}{m} \sum_{i=1}^m X_i$$

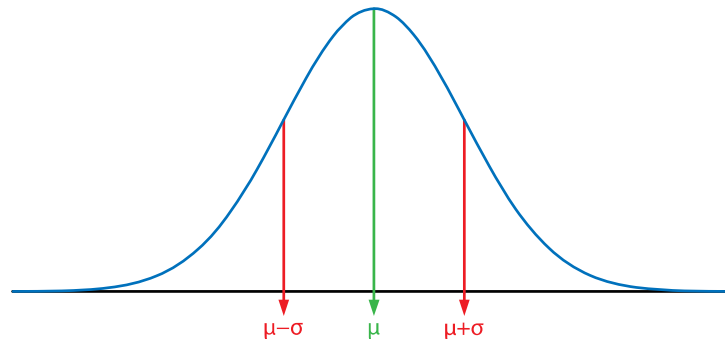


Maximum likelihood for normal

The maximum likelihood μ for normal random variables is the **sample mean**, and the maximum likelihood σ^2 is the “uncorrected” **mean square deviation**.



$$\hat{\mu} = \frac{1}{m} \sum_{i=1}^m X_i \quad \hat{\sigma}^2 = \frac{1}{m} \sum_{i=1}^m (X_i - \hat{\mu})^2$$



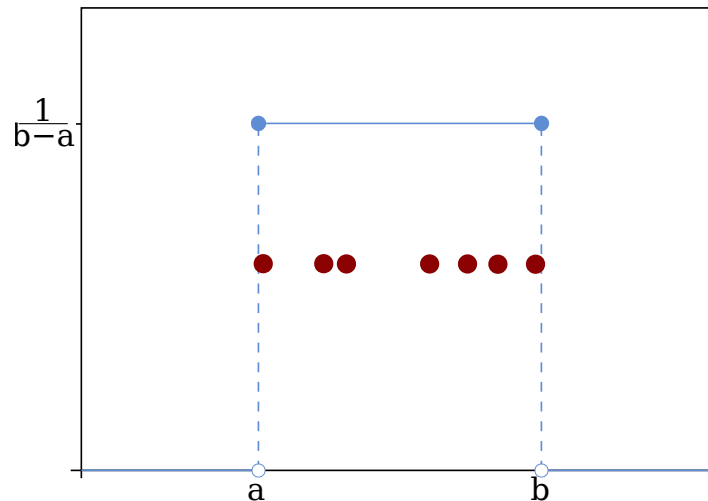
Maximum likelihood for uniform

The maximum likelihood a and b for **uniform** random variables are the **minimum and maximum** of the data.



$$\hat{a} = \min_i X_i$$

$$\hat{b} = \max_i X_i$$



Maximum a posteriori estimation

Choose the **most likely** parameters given the **example data**. You'll need a **prior probability** over the parameters.



$$\begin{aligned}\hat{\theta} &= \arg \max_{\theta} P(\theta | X_1, \dots, X_n) \\ &= \arg \max_{\theta} [LL(\theta) + \log P(\theta)]\end{aligned}$$

Laplace smoothing

Also known as **add-one** smoothing:
assume you've seen one "imaginary"
occurrence of each possible outcome.

$$p_i = \frac{\#(X=i) + 1}{n + m}$$

or: "add- k " smoothing
(if you believe equally
likely is more plausible)

$$p_i = \frac{\#(X=i) + k}{n + mk}$$



Parameter priors

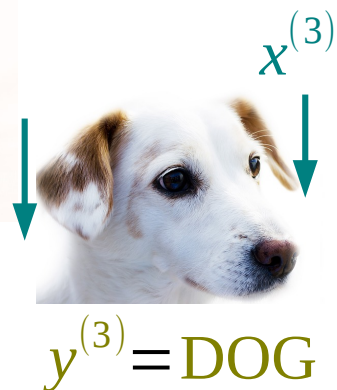
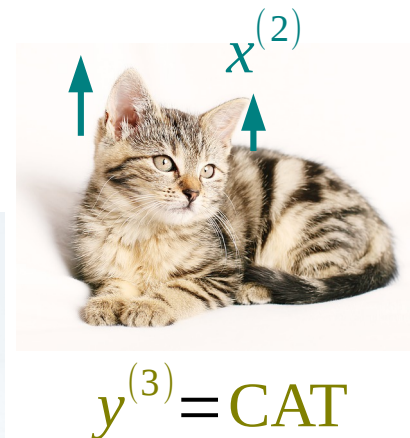
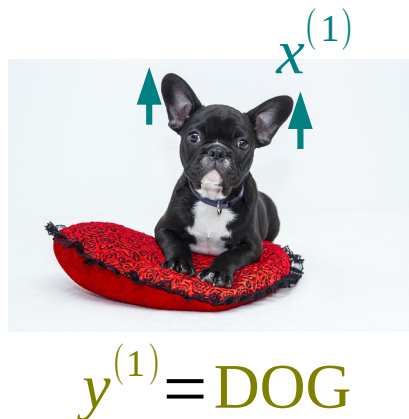
$X \sim$	Ber(p)	$p \sim \text{Beta}(a, b)$
	Bin(n, p)	$p \sim \text{Beta}(a, b)$
	MN(p)	$p \sim \text{Dir}(a)$
	Poi(λ)	$\lambda \sim \text{Gamma}(k, \theta)$
	Exp(λ)	$\lambda \sim \text{Gamma}(k, \theta)$
	N(μ, σ^2)	$\mu \sim \text{N}(\mu', \sigma'^2)$ $\sigma^2 \sim \text{InvGamma}(\alpha, \beta)$

Classification

The most basic machine learning task:
predict a **label** from a vector of **features**.



$$\hat{y} = \arg \max_y P(\mathbf{Y} = y | \vec{\mathbf{X}} = \vec{x})$$

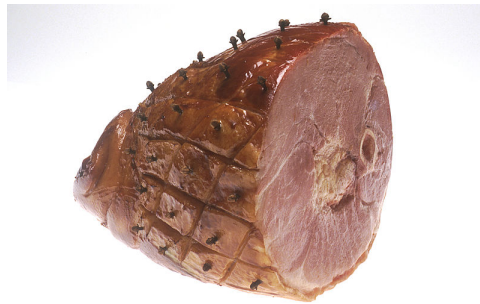


Naïve Bayes

A classification algorithm using the assumption that features are **conditionally independent** given the label.



$$\hat{y} = \arg \max_y \hat{P}(Y = y) \prod_j \hat{P}(X_j = x_j | Y = y)$$



Three secret ingredients

1. Maximum likelihood or maximum a posteriori for conditional probabilities.

$$\hat{P}(X_j = x_j | Y = y) = \frac{\#(X_j = x_j, Y = y)[+1]}{\#(Y = y)[+2]}$$

2. “Naïve Bayes assumption”: features are independent conditioned on the label.

$$\hat{P}(\vec{X} = \vec{x} | Y = y) = \prod_j \hat{P}(X_j = x_j | Y = y)$$

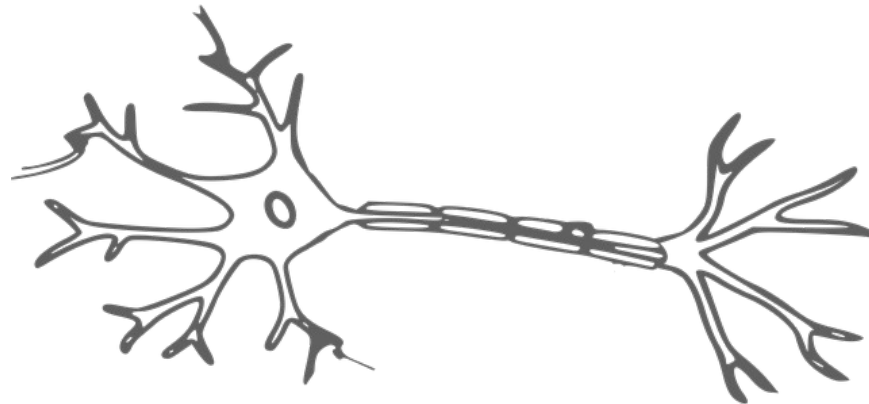
3. (Take logs for numerical stability.)

Logistic regression

A classification algorithm using the assumption that **log odds** are a linear function of the features.

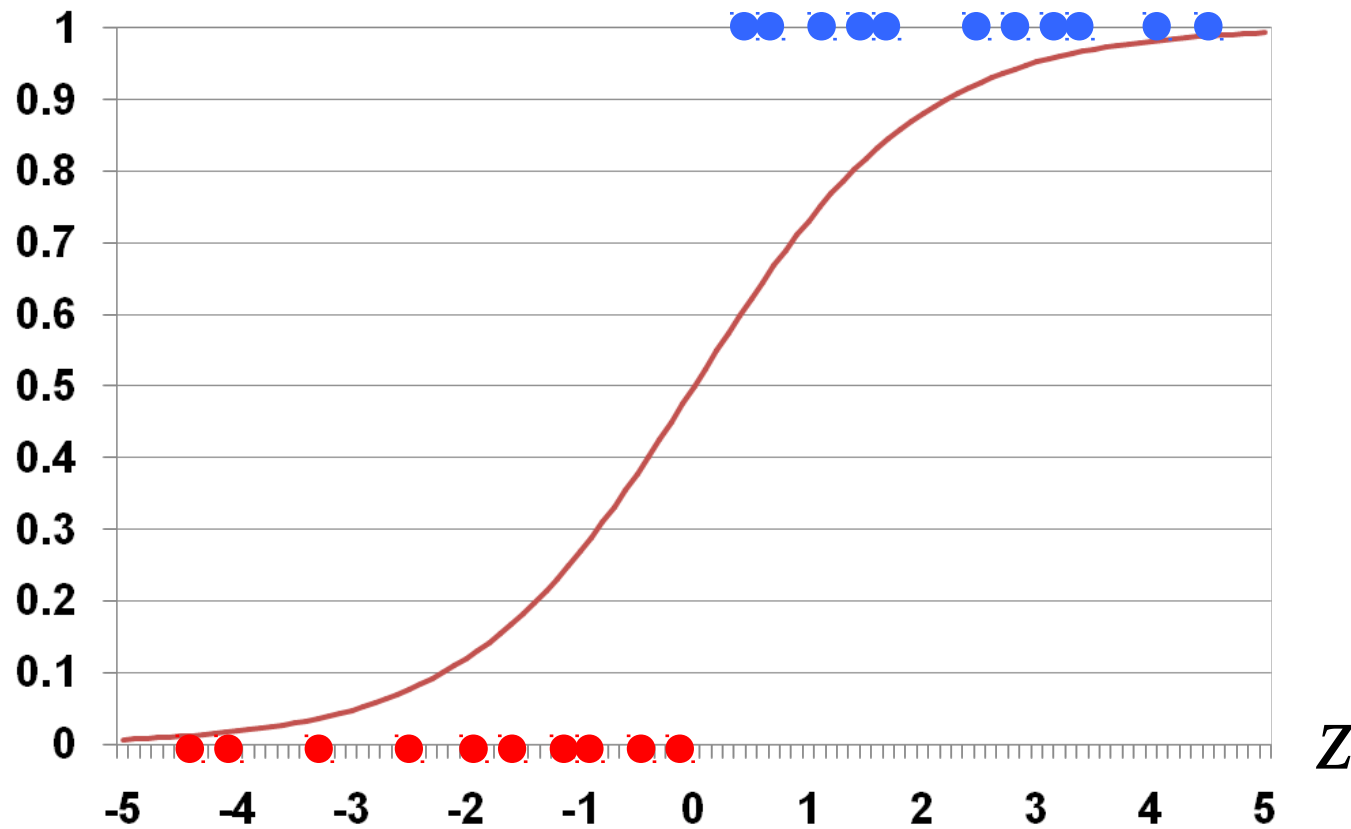


$$\hat{y} = \arg \max_y \frac{1}{1 + e^{-\vec{\theta}^T \vec{x}}}$$



Predicting 0/1 with the logistic

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



Logistic regression: Pseudocode

initialize: $\theta = [0, 0, \dots, 0]$ (m elements)

repeat many times:

 gradient = $[0, 0, \dots, 0]$ (m elements)

for each training example $(\vec{x}^{(i)}, y^{(i)})$:

for j = 0 **to** m:

 gradient[j] += $[y^{(i)} - \sigma(\vec{\theta}^T \vec{x}^{(i)})] x_j^{(i)}$

for j = 0 **to** m:

$\theta[j] += \eta * \text{gradient}[j]$

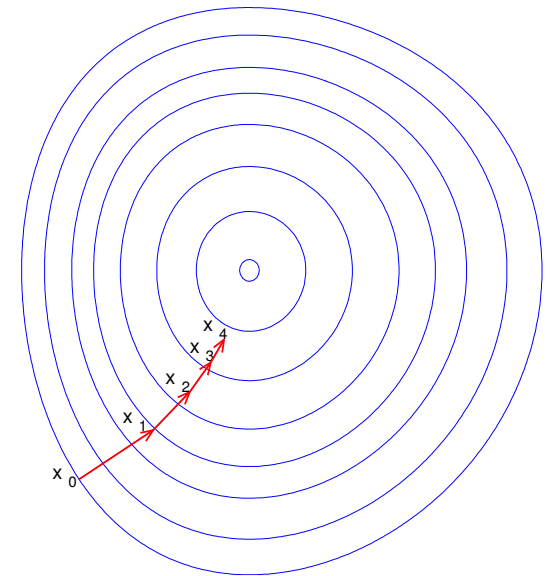
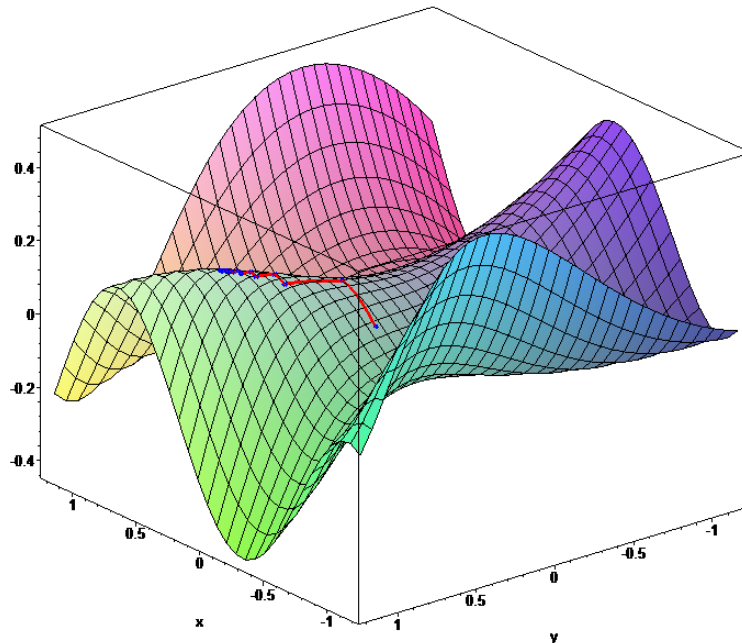
return θ

Gradient ascent

An algorithm for computing an **arg max** by taking small steps uphill (i.e., in the direction of the **gradient** of the function).



$$\vec{\theta} \leftarrow \vec{\theta} + \eta \cdot \nabla_{\vec{\theta}} f(\vec{\theta})$$

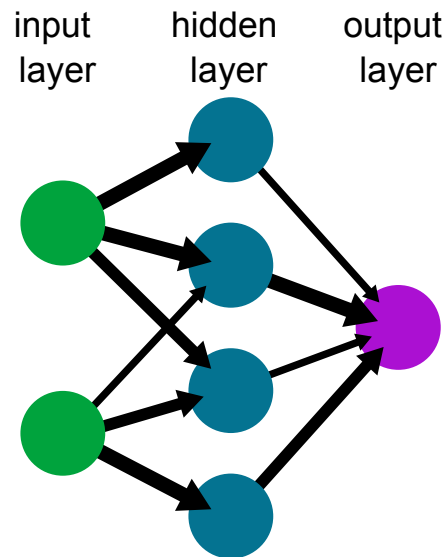


Feedforward neural network

An algorithm for classification or regression that uses **layers of logistic regressions** to discover its own features.



$$\hat{y} = \sigma \left(\theta^{(\hat{y})} \sigma \left(\theta^{(h)} \vec{x} \right) \right)$$





Keep in touch!