

Foundations

COMP3314 — Lecture 2

Lingpeng Kong

Department of Computer Science, The University of Hong Kong

Based on: Probabilistic Machine Learning by Kevin Murphy

Slides from: Saw Shier Nee with special thanks!

What is probability?



$$P(\text{head}) = 50\%$$

Frequentist

Bayesian



events (repeated trials)



uncertainty (ignorance)

Probability

Joint probability:

$$\Pr(A \wedge B) = \Pr(A, B)$$

Probability of a union of two events:

$$\Pr(A \vee B) = \Pr(A) + \Pr(B) - \Pr(A \wedge B)$$

Conditional probability:

$$\Pr(B|A) \triangleq \frac{\Pr(A, B)}{\Pr(A)}$$

If A and B are independent events:

$$\Pr(A, B) = \Pr(A) \Pr(B)$$

If the events are mutually exclusive:

$$\Pr(A \vee B) = \Pr(A) + \Pr(B)$$

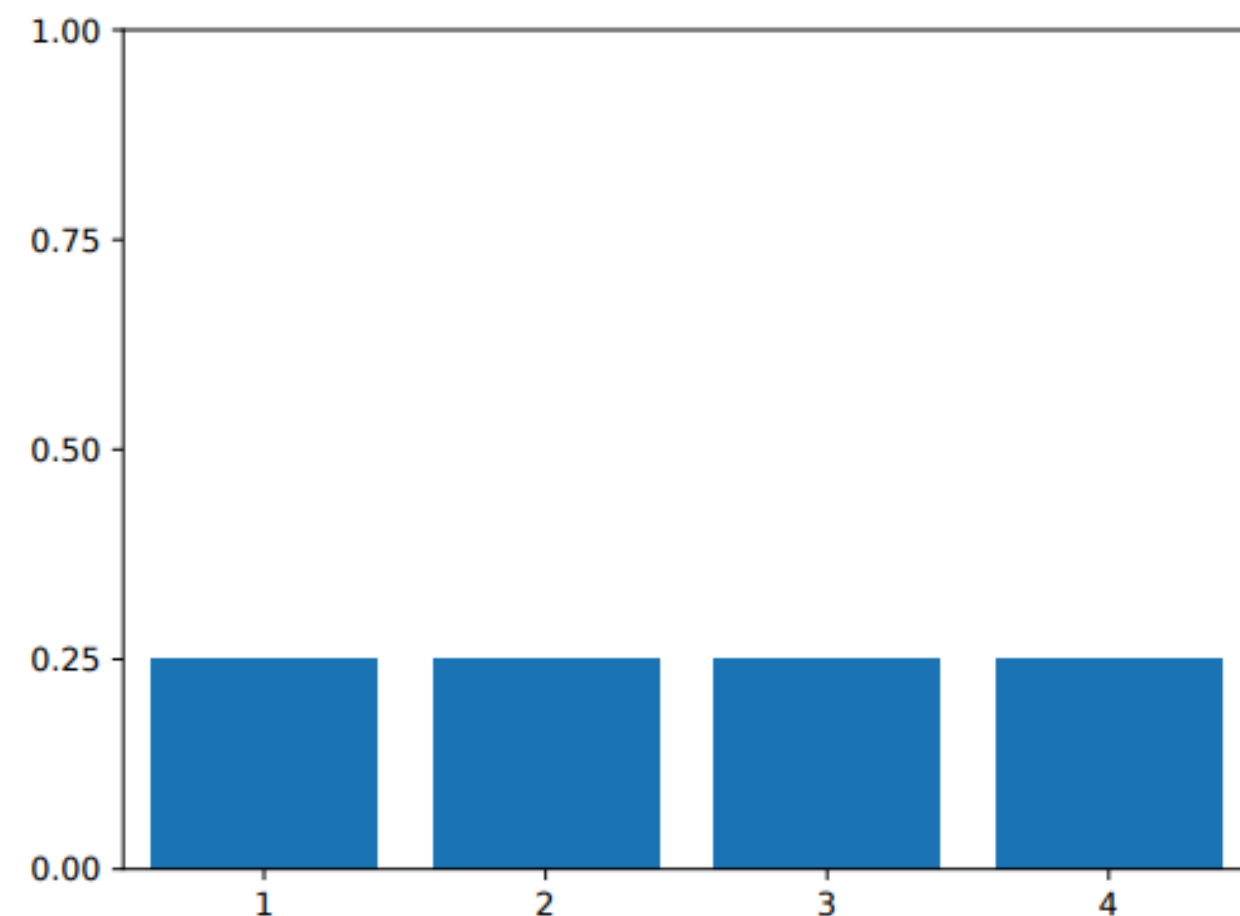
Conditional independence:

$$\Pr(A, B|C) = \Pr(A|C) \Pr(B|C)$$

Random Variables

X represents some unknown quantity of interest

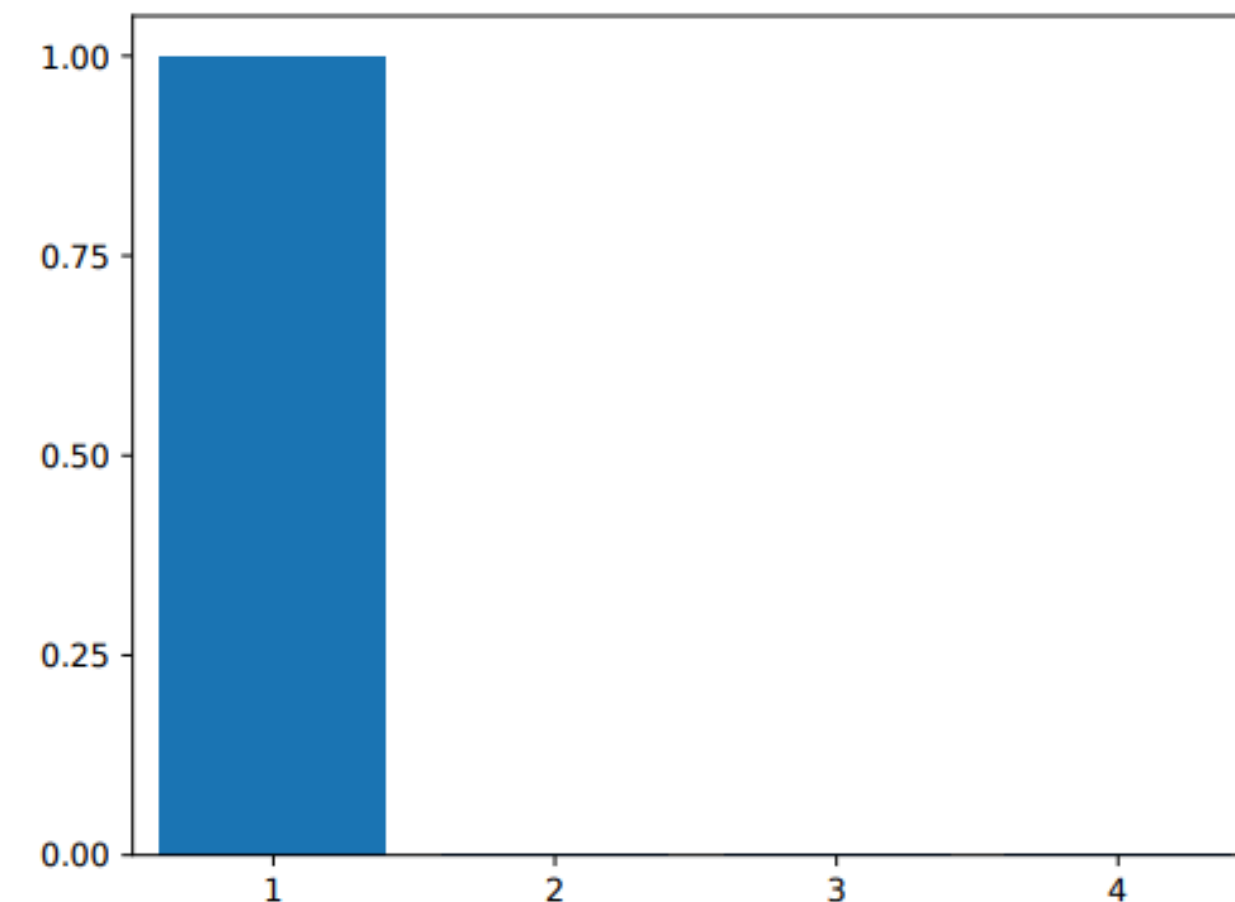
We call it a **random variable** if the value of X is unknown and/or could change.



(a)

uniform distribution

probability mass function (pmf):

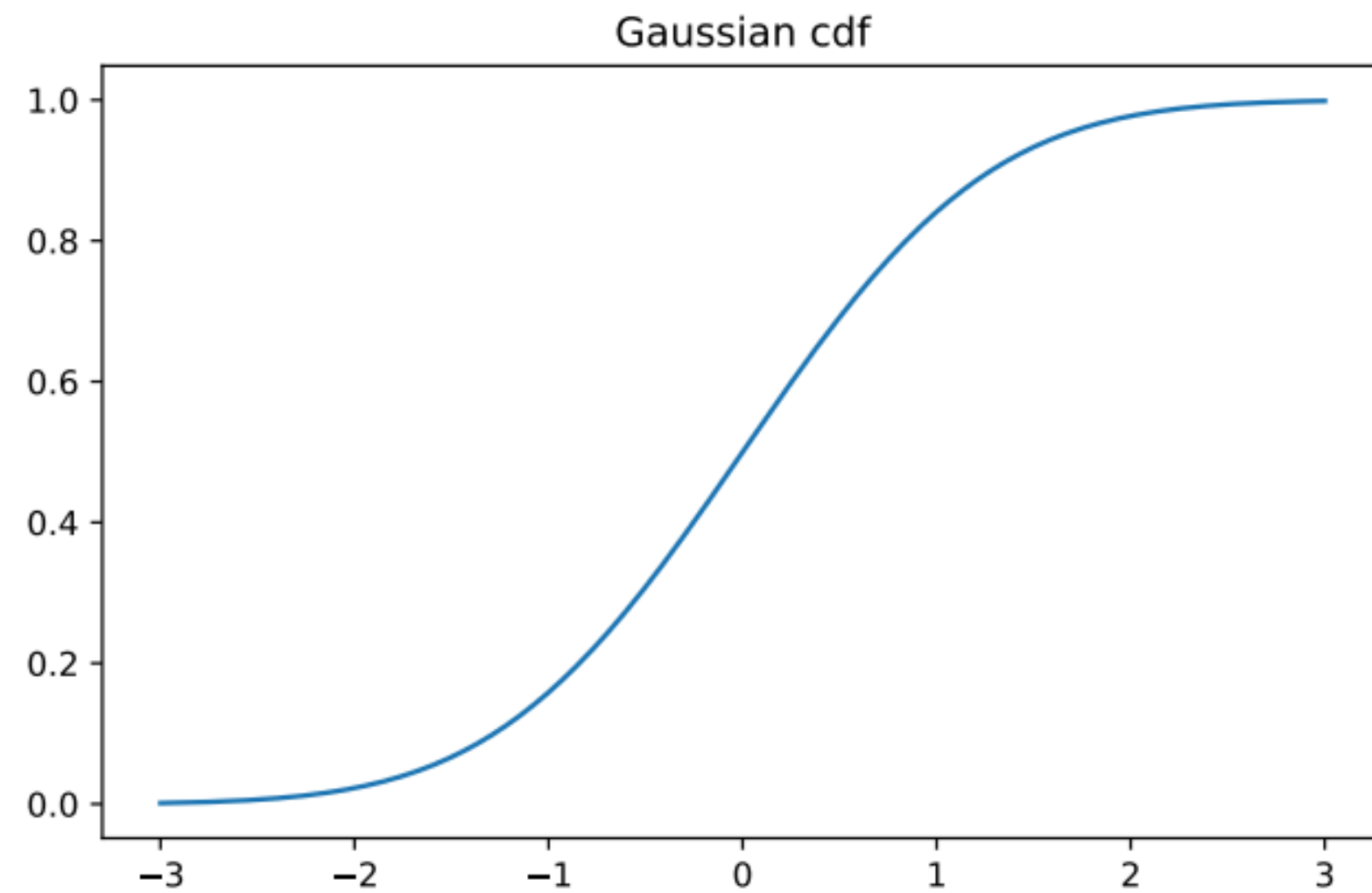


(b)

degenerate distribution
(delta function)

$$p(x) \triangleq \Pr(X = x)$$

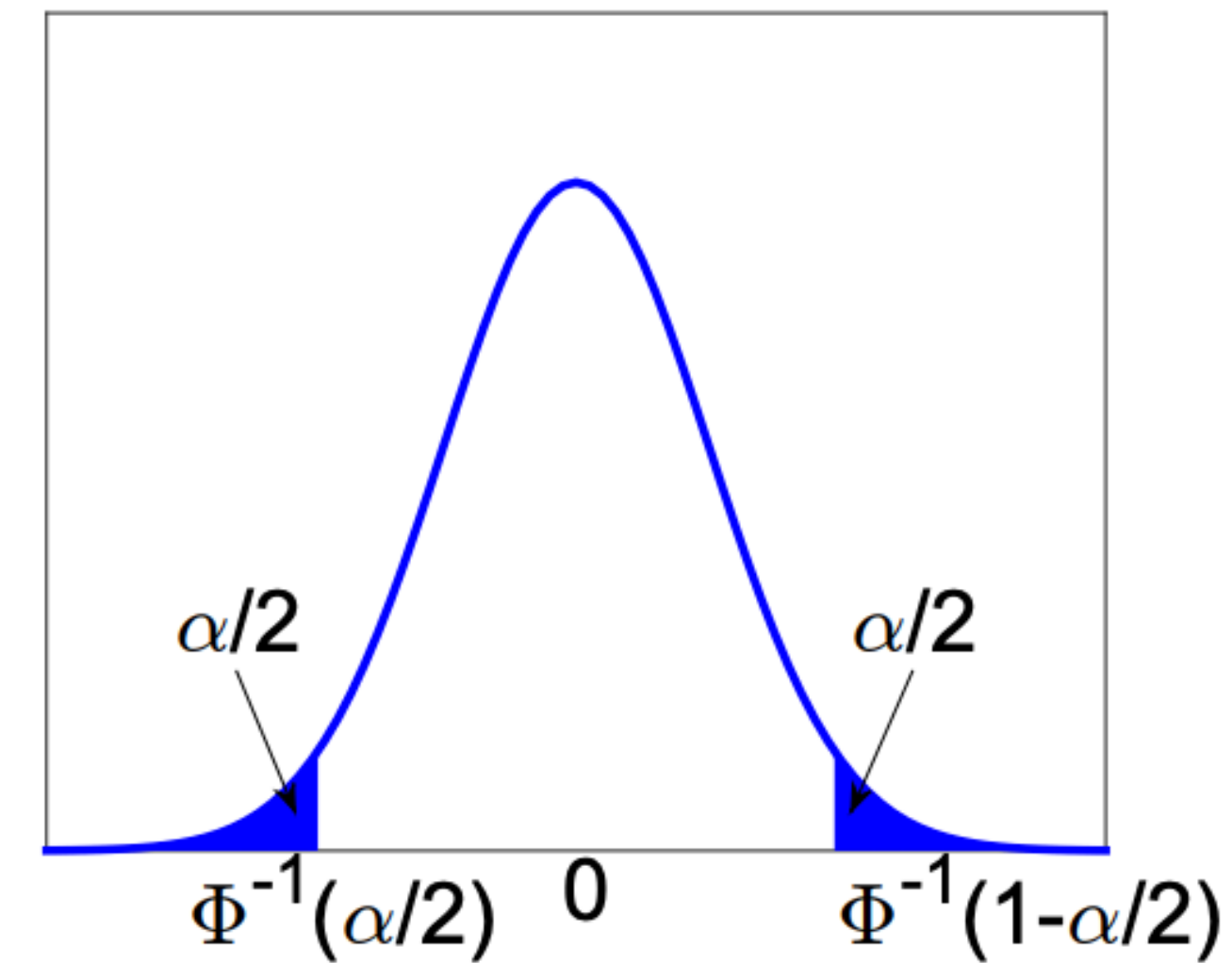
Continuous Random Variables



(a)

cumulative distribution function (cdf)

$$P(x) \triangleq \Pr(X \leq x)$$



(b)

probability density function (pdf)

$$p(x) \triangleq \frac{d}{dx} P(x)$$

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

$$\Pr(x \leq X \leq x + dx) \approx p(x) dx$$

Sets of Related Random Variables

$p(X, Y)$	$Y = 0$	$Y = 1$
$X = 0$	0.2	0.3
$X = 1$	0.3	0.2

marginal distribution

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

conditional distribution

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

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

chain rule

$$p(\mathbf{x}_{1:D}) = p(x_1)p(x_2|x_1)p(x_3|x_1, x_2)p(x_4|x_1, x_2, x_3) \dots p(x_D|\mathbf{x}_{1:D-1})$$

Moments of a Distribution

Mean of a distribution

$$\mathbb{E}[X] \triangleq \int_{\mathcal{X}} x p(x) dx$$

$$\mathbb{E}[X] \triangleq \sum_{x \in \mathcal{X}} x p(x)$$

Linearity of expectation

$$\mathbb{E}[aX + b] = a\mathbb{E}[X] + b$$

Variance of a distribution

$$\begin{aligned} \mathbb{V}[X] &\triangleq \mathbb{E}[(X - \mu)^2] = \int (x - \mu)^2 p(x) dx \\ &= \int x^2 p(x) dx + \mu^2 \int p(x) dx - 2\mu \int x p(x) dx = \mathbb{E}[X^2] - \mu^2 \end{aligned}$$

$$\mathbb{E}[X^2] = \sigma^2 + \mu^2 \quad \text{assuming some random variable with mean} = 0$$

$$\mathbb{V}[aX + b] = a^2 \mathbb{V}[X]$$

Moments of a Distribution

Mode of a distribution

$$\boldsymbol{x}^* = \operatorname{argmax}_{\boldsymbol{x}} p(\boldsymbol{x})$$

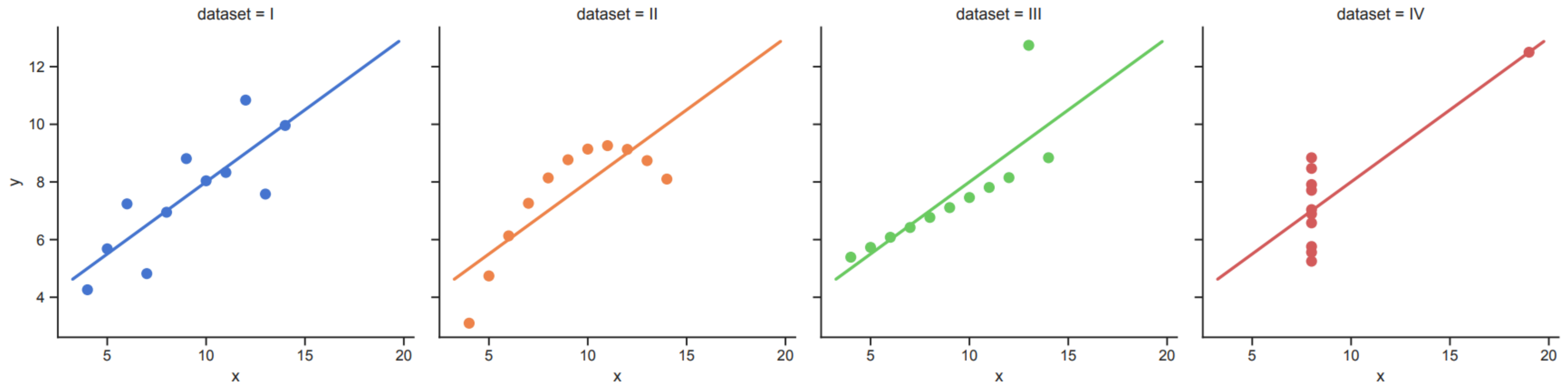
conditional moments

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

$$\begin{aligned}\mathbb{E}_Y[\mathbb{E}[X|Y]] &= \mathbb{E}_Y \left[\sum_x x p(X = x|Y) \right] \\ &= \sum_y \left[\sum_x x p(X = x|Y) \right] p(Y = y) = \sum_{x,y} xp(X = x, Y = y) = \mathbb{E}[X]\end{aligned}$$

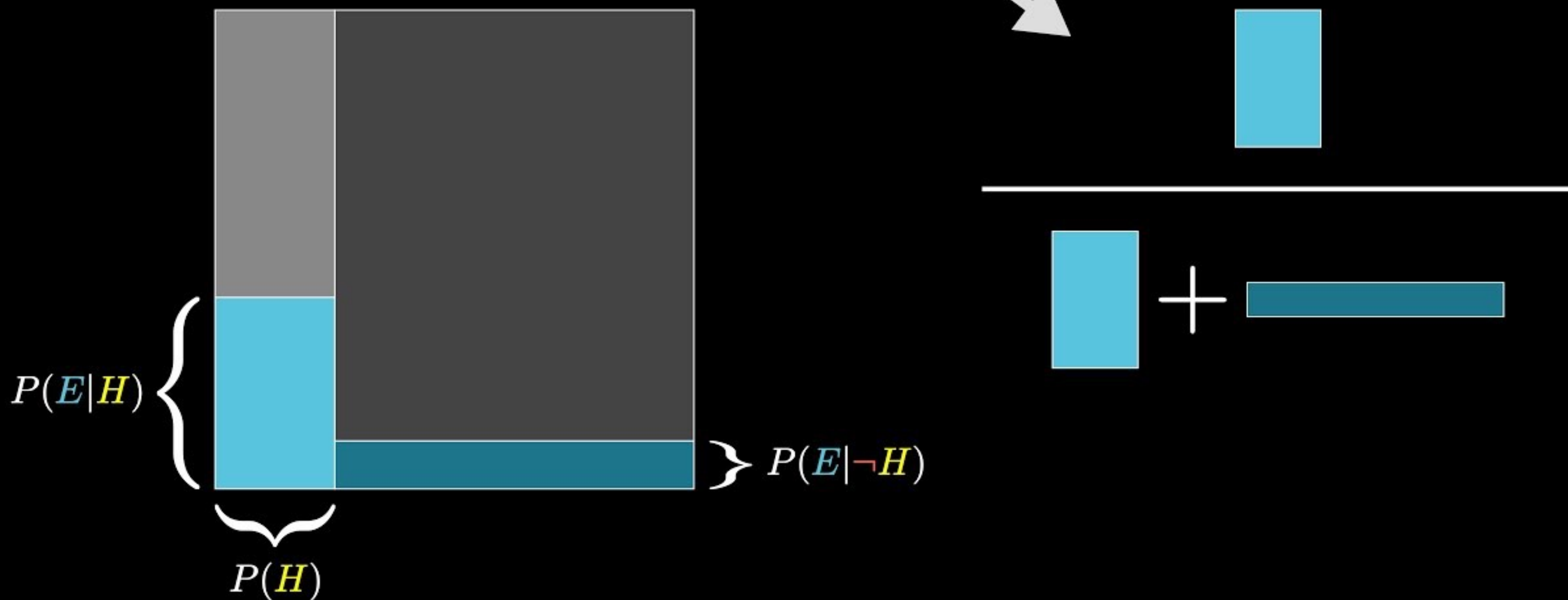
law of total (iterated) expectation

Limitations of Summary Statistics



Anscombe's Quartet: $E[x] = 9$, $V[x] = 11$, $E[y] = 7.50$, $V[y] = 4.125$

This is Bayes' rule



Bayes' rule

$$p(H = h|Y = y) = \frac{p(H = h)p(Y = y|H = h)}{p(Y = y)}$$

prior distribution — what we know about possible values of H before we see any data

$$p(H)$$

observation distribution — possible outcomes Y we expect to see if H =h

$$p(Y|H = h)$$

likelihood — evaluate the observation distribution at a point corresponding to the actual observations, y

$$p(Y = y|H = h)$$

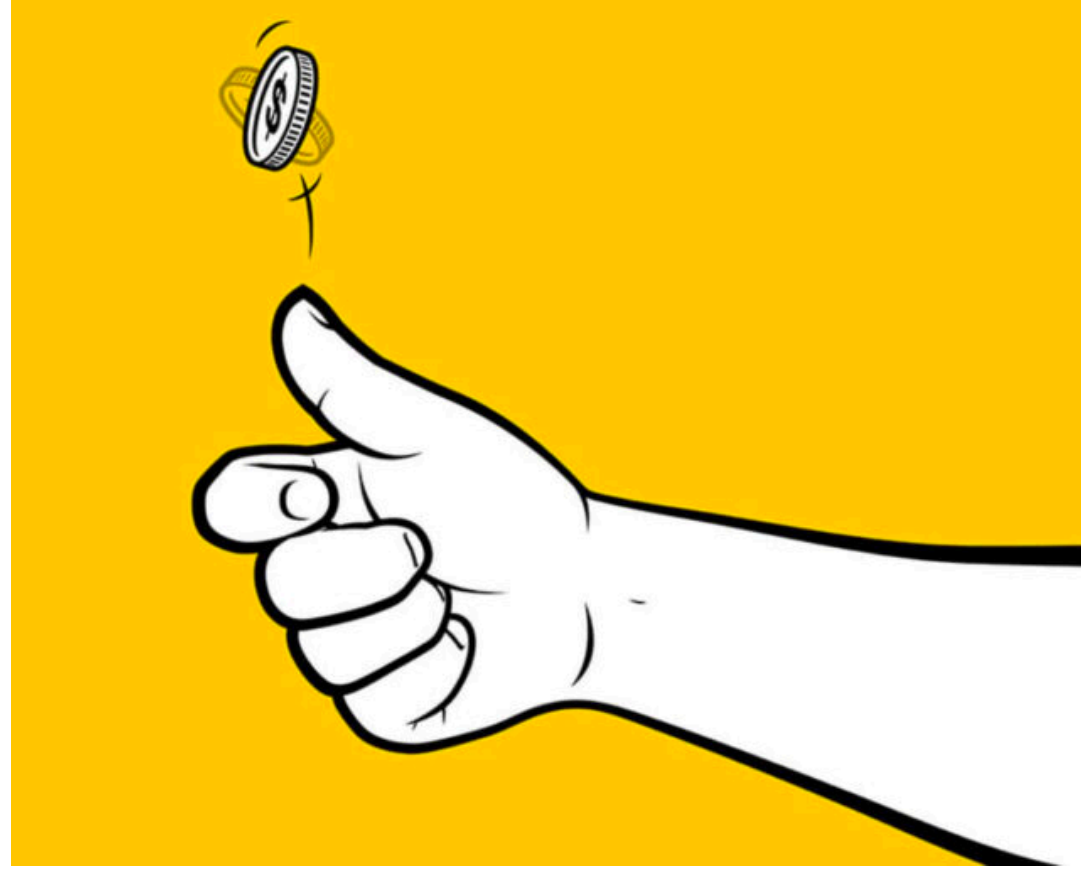
marginal likelihood

$$p(Y = y) = \sum_{h' \in \mathcal{H}} p(H = h')p(Y = y|H = h') = \sum_{h' \in \mathcal{H}} p(H = h', Y = y)$$

posterior — our new belief state about the possible value of H

$$p(H = h|Y = y)$$

Bernoulli Distribution

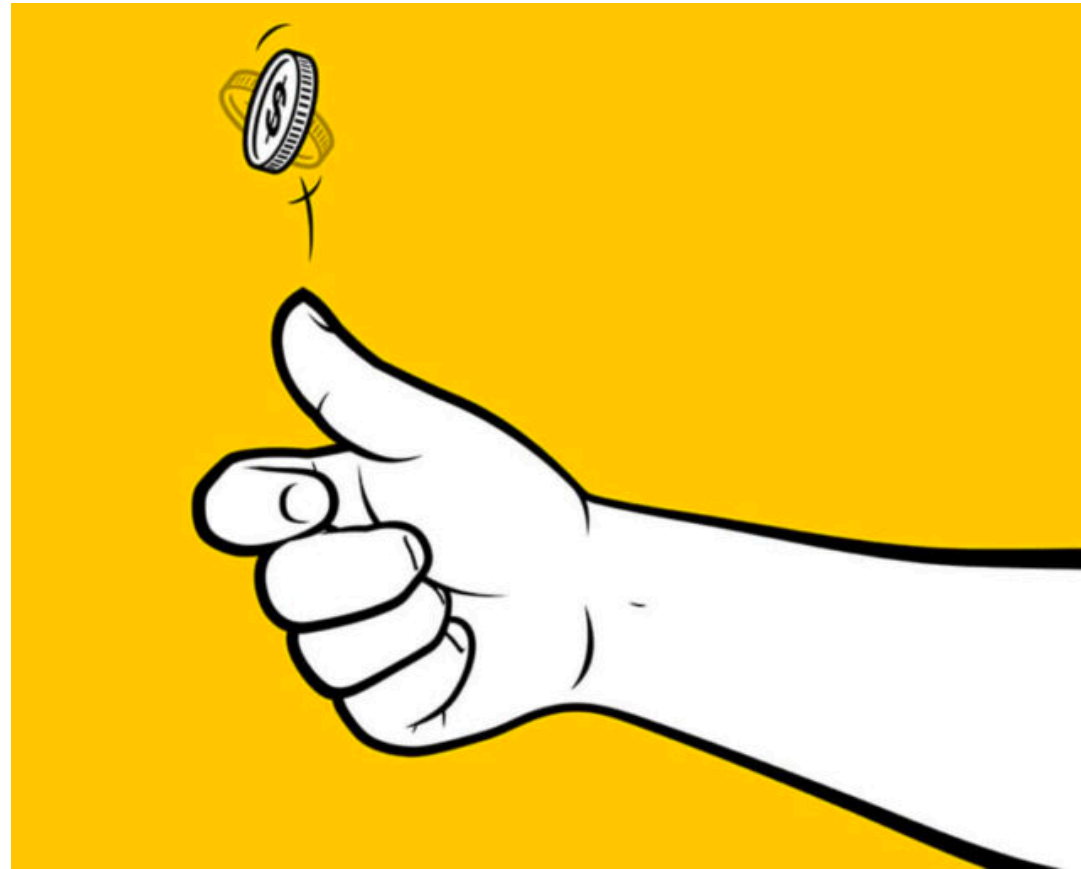


$$Y \sim \text{Ber}(\theta)$$

$$\text{Ber}(y|\theta) = \begin{cases} 1 - \theta & \text{if } y = 0 \\ \theta & \text{if } y = 1 \end{cases}$$

$$\text{Ber}(y|\theta) \triangleq \theta^y (1 - \theta)^{1-y}$$

Binomial Distribution



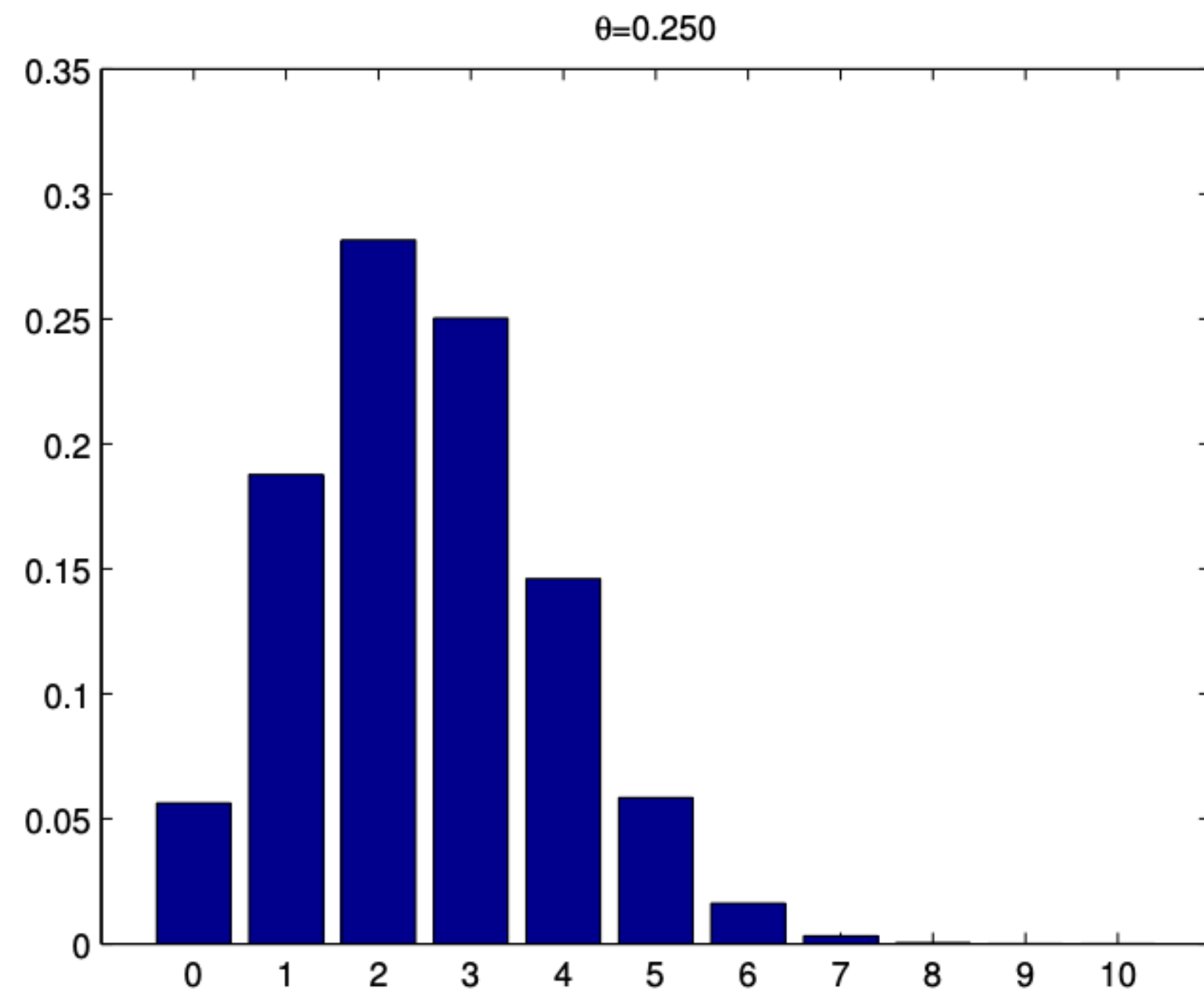
$N = 5$ (N Bernoulli trials), $s = 4$ (total number of heads)

$$\text{Bin}(s|N, \theta) \triangleq \binom{N}{s} \theta^s (1 - \theta)^{N-s}$$

$$\binom{N}{k} \triangleq \frac{N!}{(N-k)!k!}$$

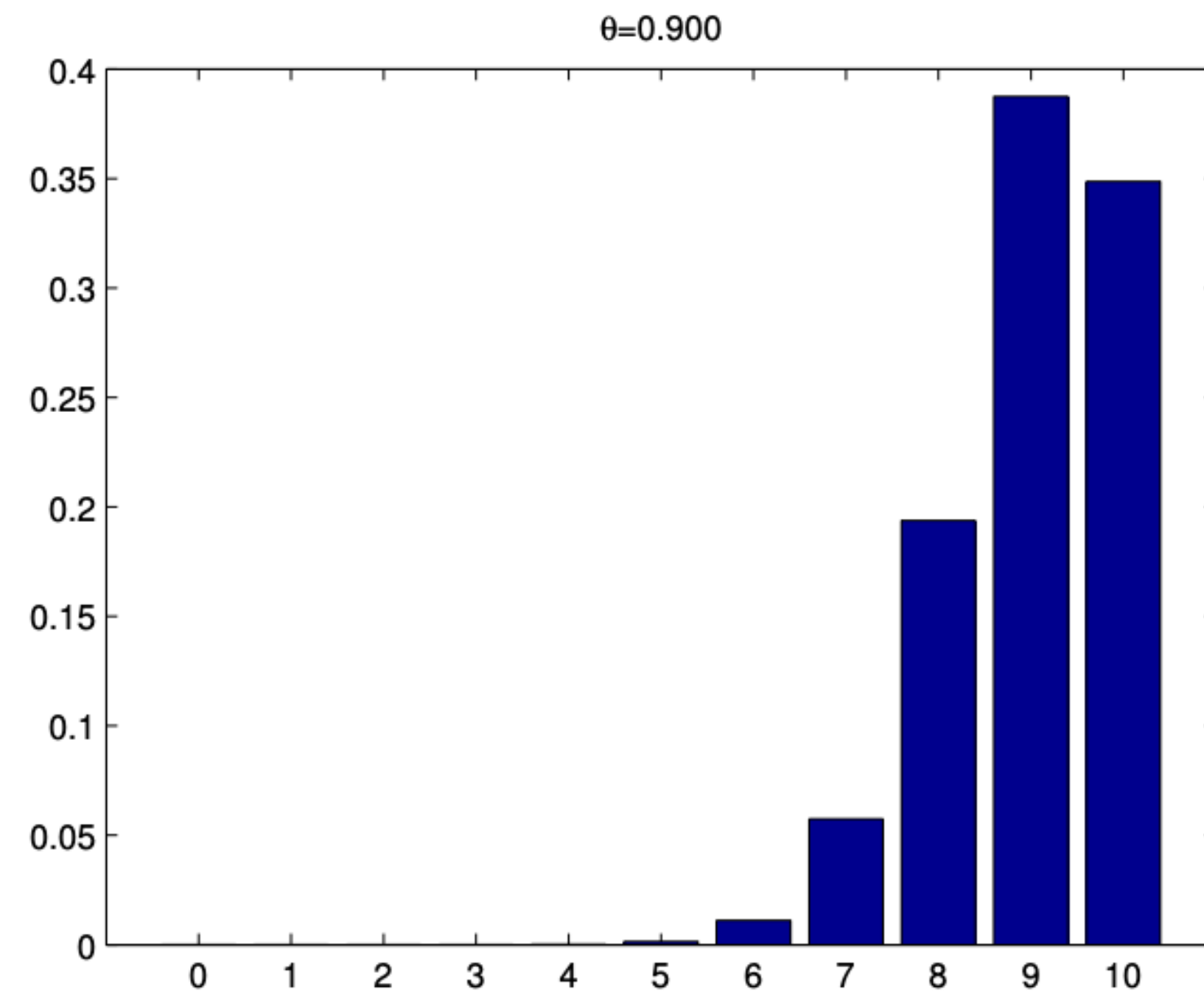
Binomial Distribution

$$\text{Bin}(s|N, \theta) \triangleq \binom{N}{s} \theta^s (1 - \theta)^{N-s}$$



(a)

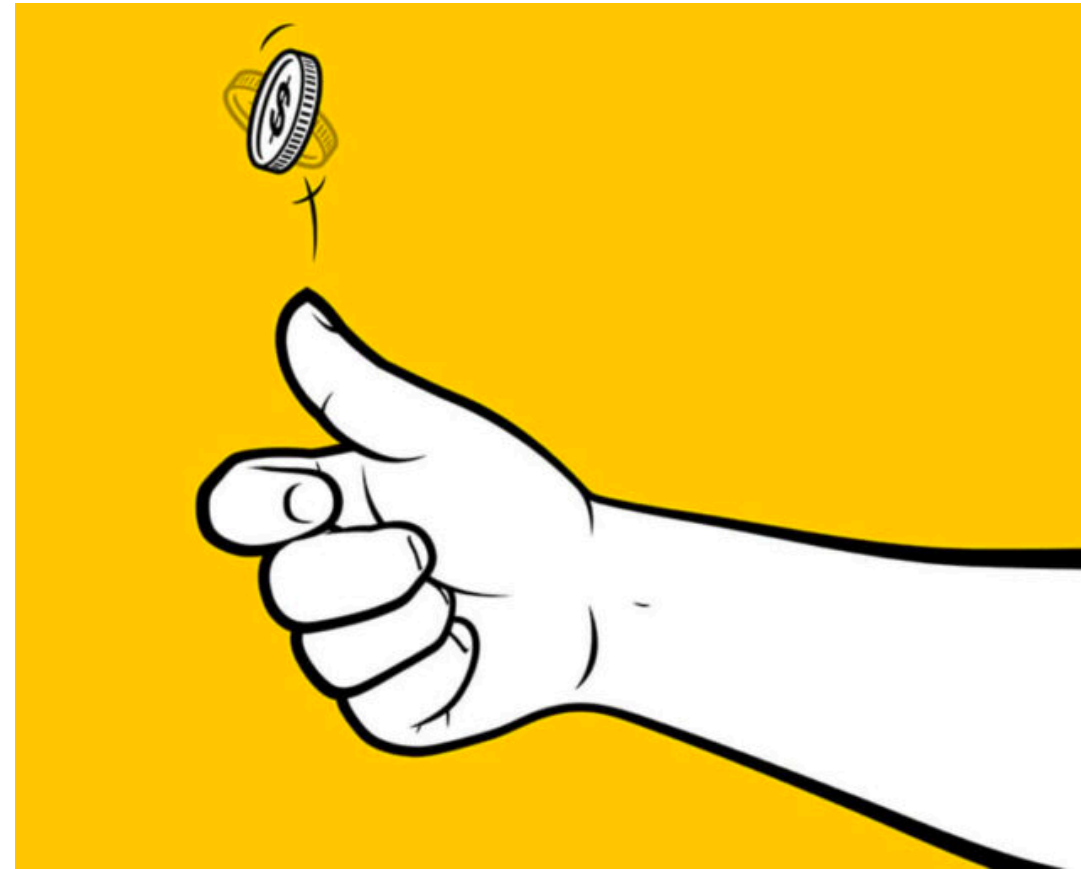
$N = 10, \theta = 0.25$



(b)

$N = 10, \theta = 0.9$

A Simple Question?



What's the probability of head up?

$P(\text{Heads})$

0.8?

A Simple Question?

$$P(\text{Heads}) = \theta$$

$$P(\text{Tails}) = 1 - \theta$$

$$D = \{\text{head, head, head, head, tail}\}$$

α_H total number of heads

α_T total number of tails

$$P(D|\theta) = P(\alpha_H, \alpha_T | \theta) = \theta^{\alpha_H} (1 - \theta)^{\alpha_T}$$

Maximum Likelihood Estimation

MLE: Choose θ that maximizes the probability of observed data

$$\begin{aligned}\hat{\theta} &= \arg \max_{\theta} P(D | \theta) \\ &= \arg \max_{\theta} \log P(D | \theta)\end{aligned}$$

Set derivative to zero

$$\frac{\partial \log P(D | \theta)}{\partial \theta} = 0$$

Maximum Likelihood Estimation

Set derivative to zero

$$\frac{\partial \log P(D | \theta)}{\partial \theta} = 0$$

$$\frac{\partial \log \theta^{\alpha_H} (1 - \theta)^{\alpha_T}}{\partial \theta} = 0$$

$$\hat{\theta}_{\text{MLE}} = \frac{\alpha_H}{\alpha_H + \alpha_T} = \frac{4}{4 + 1} = 0.8$$

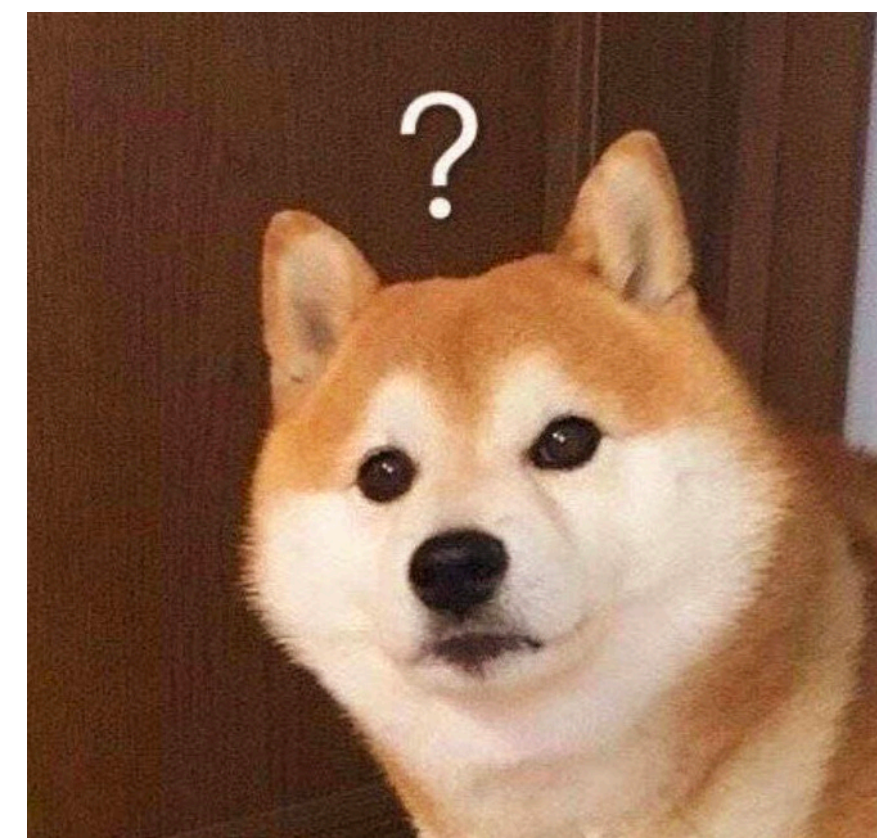
Maximum Likelihood Estimation

If you flip the coin 5 times, get 4 times head

$$\hat{\theta}_{\text{MLE}} = \frac{\alpha_H}{\alpha_H + \alpha_T} = \frac{4}{4 + 1} = 0.8$$

If you flip the coin 5000 times, get 4000 times head

$$\hat{\theta}_{\text{MLE}} = \frac{\alpha_H}{\alpha_H + \alpha_T} = \frac{4000}{4000 + 1000} = 0.8$$



Maximum Likelihood Estimation

MLE: Choose θ that maximizes the probability of observed data

$$\begin{aligned}\hat{\theta} &= \arg \max_{\theta} P(D | \theta) \\ &= \arg \max_{\theta} \log P(D | \theta)\end{aligned}$$

Set derivative to zero

$$\frac{\partial \log P(D | \theta)}{\partial \theta} = 0$$

$$\hat{\theta}_{\text{MLE}} = \frac{\alpha_H}{\alpha_H + \alpha_T}$$

Bayesian Learning

$$P(\theta | D) = \frac{P(D | \theta)P(\theta)}{P(D)}$$

posterior

$$\arg \max_{\theta} P(\theta | D) = \arg \max_{\theta} \frac{P(D | \theta)P(\theta)}{P(D)}$$

$$= \arg \max_{\theta} P(D | \theta)P(\theta)$$



$P(D | \theta)$

likelihood

$P(\theta)$

prior

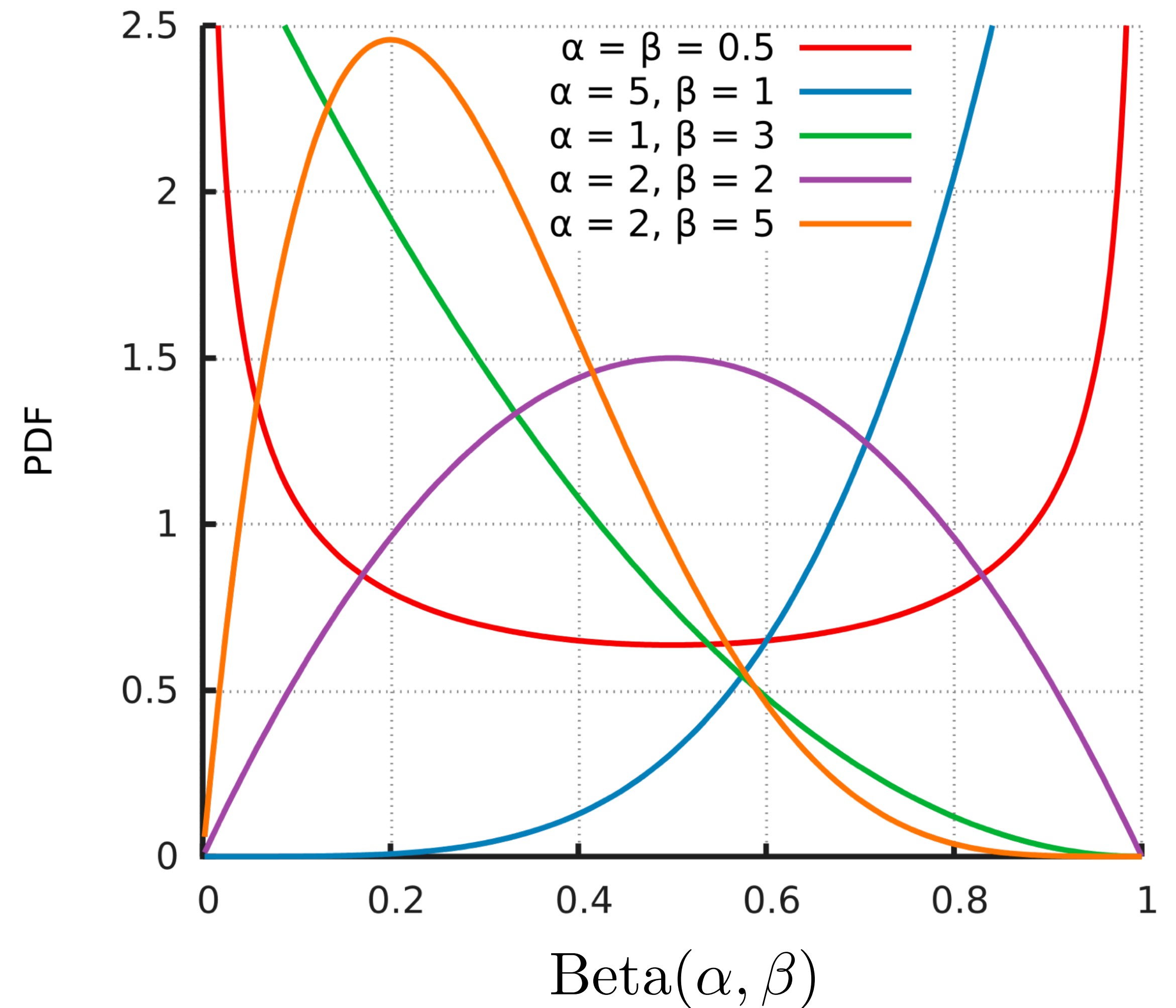
Bayesian Learning

prior distribution

$$P(\theta) \sim \text{Beta}(\beta_H, \beta_T)$$

Beta Distribution

$$P(\theta) = \frac{\theta^{\beta_H-1} (1-\theta)^{\beta_T-1}}{B(\beta_H, \beta_T)}$$



Bayesian Learning

prior distribution

$$P(\theta) \sim \text{Beta}(\beta_H, \beta_T)$$

$$P(\theta) = \frac{\theta^{\beta_H-1} (1-\theta)^{\beta_T-1}}{B(\beta_H, \beta_T)}$$

likelihood function

$$P(D|\theta) = P(\alpha_H, \alpha_T | \theta) = \theta^{\alpha_H} (1-\theta)^{\alpha_T}$$

posterior distribution

$$P(\theta | D) \propto \theta^{\alpha_H+\beta_H-1} (1-\theta)^{\alpha_T+\beta_T-1}$$

$$P(\theta | D) \sim \beta(\beta_H + \alpha_H, \beta_T + \alpha_T)$$

Estimating Parameters

MLE: Maximum Likelihood Estimate, choose θ that maximizes the probability of observed data

$$\hat{\theta} = \arg \max_{\theta} P(D | \theta)$$

MAP: Maximum a Posteriori, choose θ that is most probable given prior probability and the data

$$\hat{\theta} = \arg \max_{\theta} P(\theta | D) = \arg \max_{\theta} \frac{P(D | \theta)P(\theta)}{P(D)}$$

Sigmoid (logistic) function

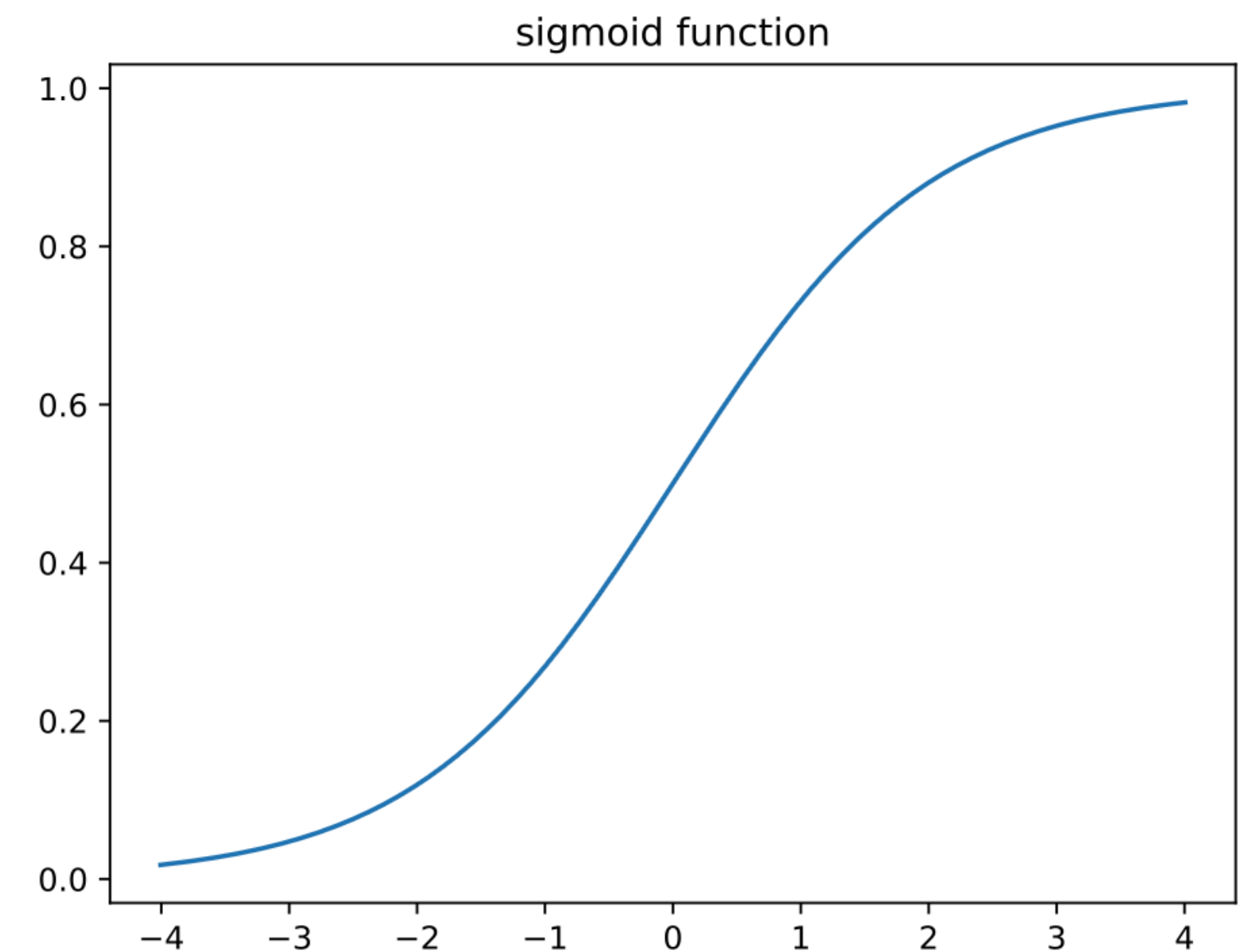
$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \text{Ber}(y|f(\mathbf{x}; \boldsymbol{\theta})) \dots\dots\dots \text{Ber}(y|\theta) \triangleq \theta^y (1 - \theta)^{1-y}$$



We want this to be an unconstrained function

sigmoid (logistic) function $\sigma()$ helps:

$$\sigma(a) \triangleq \frac{1}{1 + e^{-a}} \quad a = f(\mathbf{x}; \boldsymbol{\theta})$$



one type of “activation” in neural nets!

Categorical and Multinomial Distributions

categorical distribution

$$\mathcal{M}(y|\boldsymbol{\theta}) \triangleq \prod_{c=1}^C \theta_c^{\mathbb{I}(y=c)}$$

$$\mathcal{M}(\mathbf{y}|\boldsymbol{\theta}) \triangleq \prod_{c=1}^C \theta_c^{y_c} \quad C = 3 \quad (1, 0, 0)(0, 1, 0)(0, 0, 1)$$

one-hot vector

multinomial distribution

rolling a C-sided dice N times, \mathbf{s} to be a vector that counts the number of times each face shows up

$$\mathcal{M}(\mathbf{s}|N, \boldsymbol{\theta}) \triangleq \binom{N}{s_1 \dots s_C} \prod_{c=1}^C \theta_c^{s_c} \quad \binom{N}{s_1 \dots s_C} \triangleq \frac{N!}{s_1! s_2! \dots s_C!}$$

Softmax function

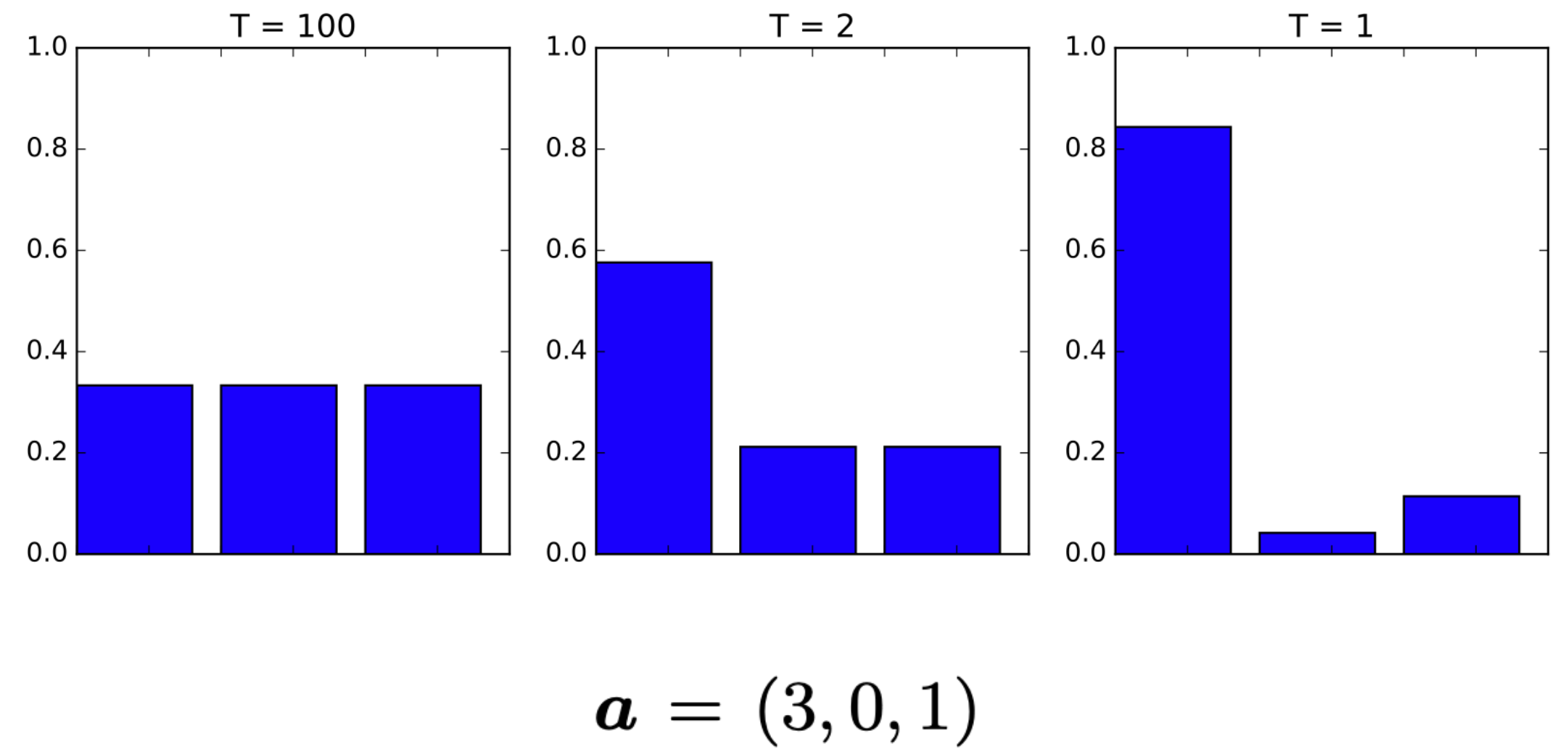
$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \mathcal{M}(y|f(\mathbf{x}; \boldsymbol{\theta})) \dots\dots\dots \mathcal{M}(y|\boldsymbol{\theta}) \triangleq \prod_{c=1}^C \theta_c^{\mathbb{I}(y=c)}$$

We want this to be an unconstrained function

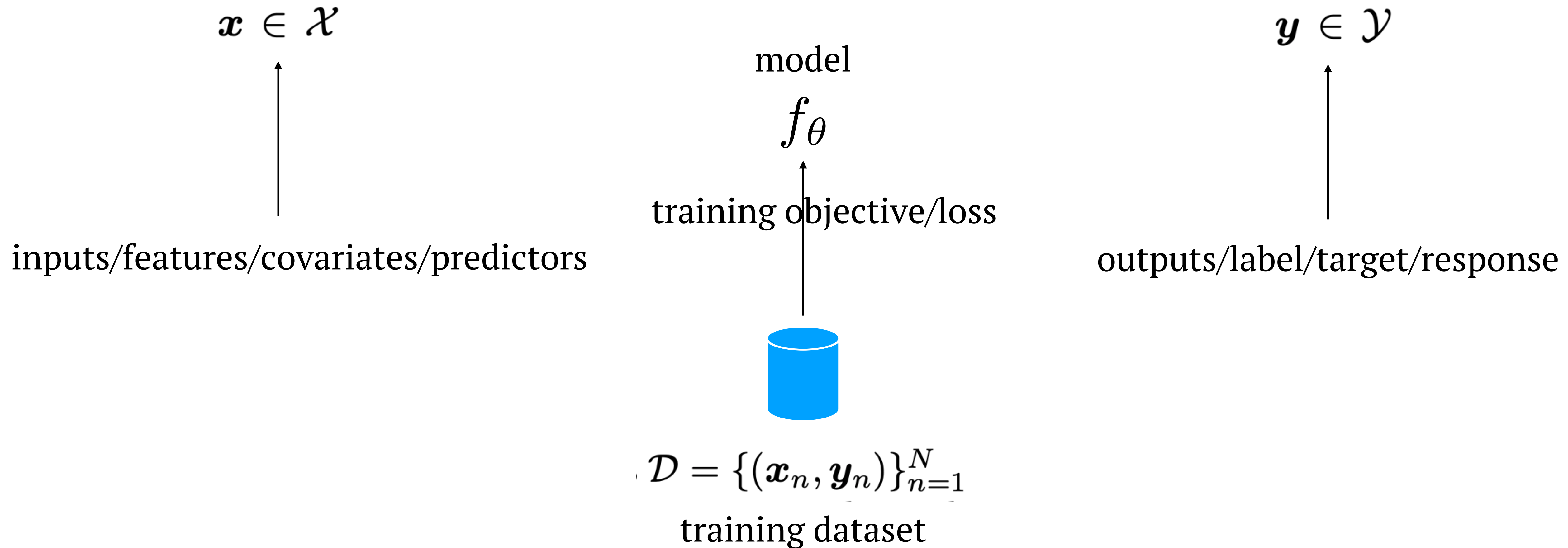
softmax function helps:

$$\mathcal{S}(\mathbf{a}) \triangleq \left[\frac{e^{a_1}}{\sum_{c'=1}^C e^{a_{c'}}}, \dots, \frac{e^{a_C}}{\sum_{c'=1}^C e^{a_{c'}}} \right]$$

$\mathcal{S}(\mathbf{a}/T)$ T is the temperature for the softmax

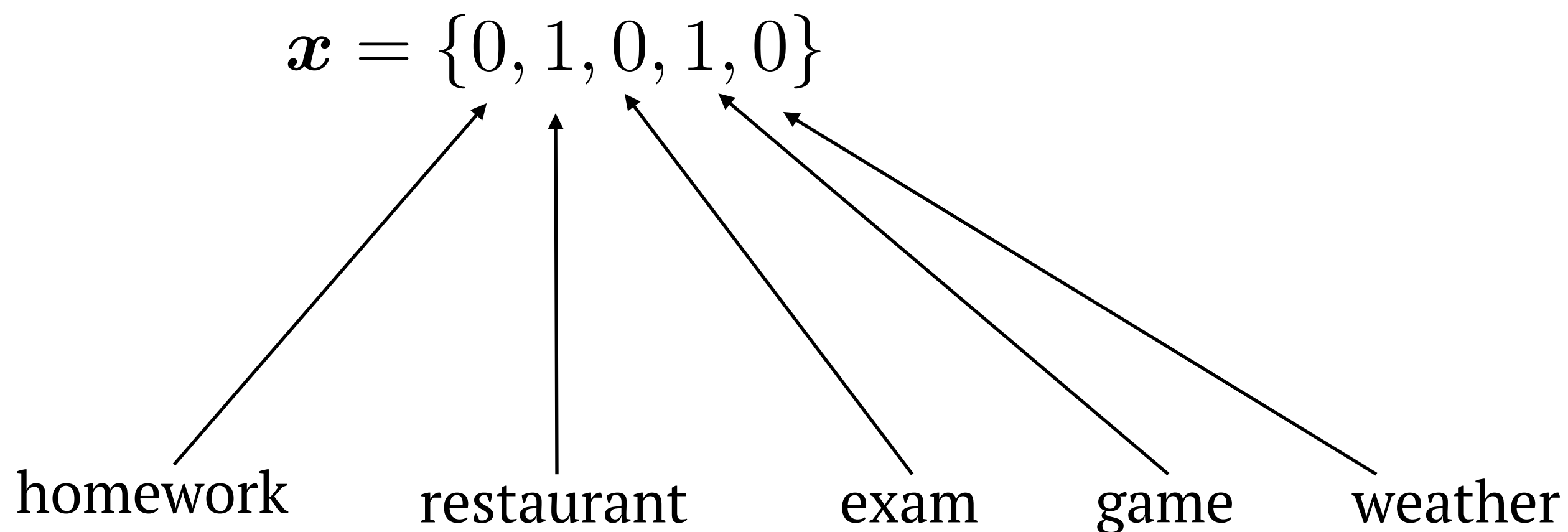


Recap: Supervised Learning in a Nutshell

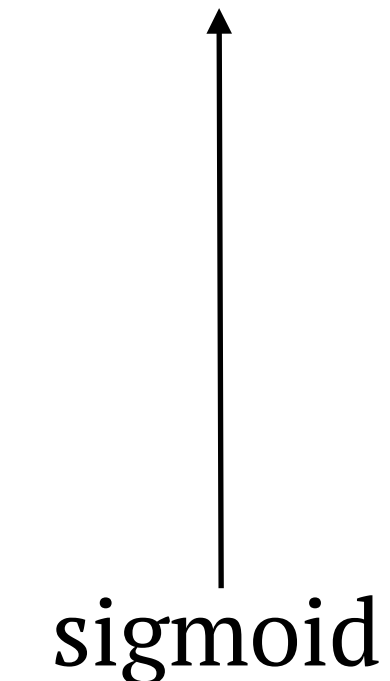


Recap: Supervised Learning in a Nutshell

Happy or sad? $y = 0, 1$



$$p(y | \mathbf{x}, \boldsymbol{\theta}) = \text{Ber}(y | \sigma(f(\mathbf{x}; \boldsymbol{\theta})))$$



	$\mathbf{x} = \{0, 0, 1, 1, 1\}$	$y = 1$
	$\mathbf{x} = \{0, 1, 1, 1, 1\}$	$y = 1$
D	$\mathbf{x} = \{0, 1, 1, 0, 1\}$	$y = 0$
	\vdots	\vdots
	$\mathbf{x} = \{1, 1, 1, 0, 1\}$	$y = 0$

$$\begin{aligned} P(D | \boldsymbol{\theta}) &= \prod_{i=0}^{|D|} p(y_i | \mathbf{x}_i, \boldsymbol{\theta}) \\ &= \prod_{i=0}^{|D|} \text{Ber}(y_i | \sigma(f(\mathbf{x}_i; \boldsymbol{\theta}))) \end{aligned}$$

Recap: Supervised Learning in a Nutshell

$$P(D | \boldsymbol{\theta}) = \prod_{i=0}^{|D|} \text{Ber}(y_i | \sigma(f(\mathbf{x}_i; \boldsymbol{\theta})))$$

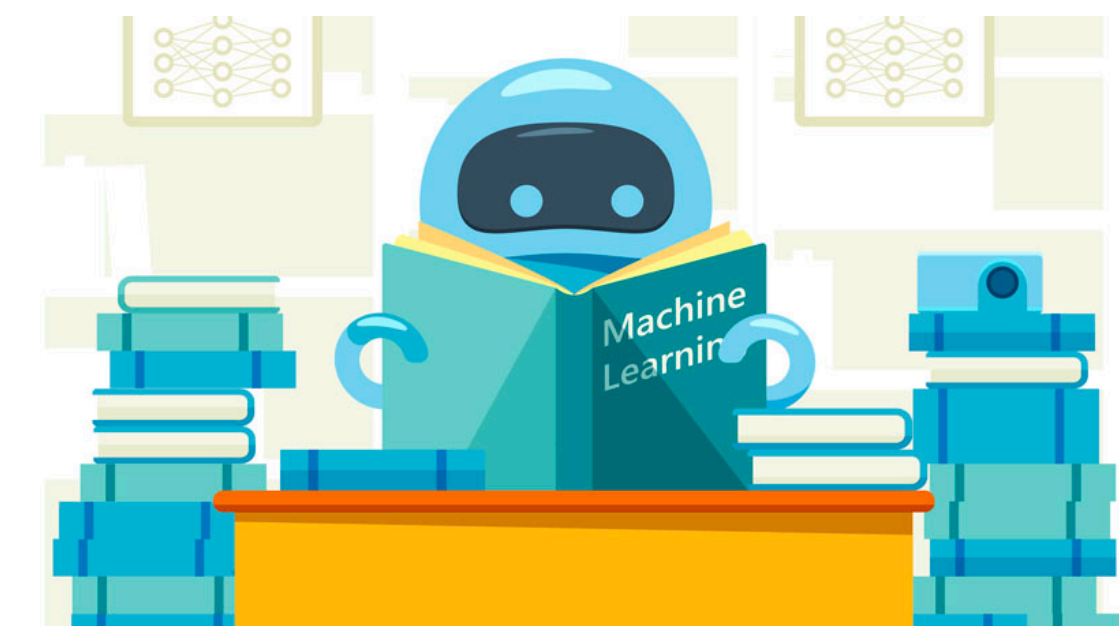
$$\sigma(a) \triangleq \frac{1}{1 + e^{-a}}$$

$$f(\mathbf{x}; \boldsymbol{\theta}) = w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + w_5x_5 + b$$

$$\boldsymbol{\theta} = \{w_1, w_2, w_3, w_4, w_5, b\} \quad \text{parameters}$$

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} P(D | \boldsymbol{\theta})$$

The machine, is learning!

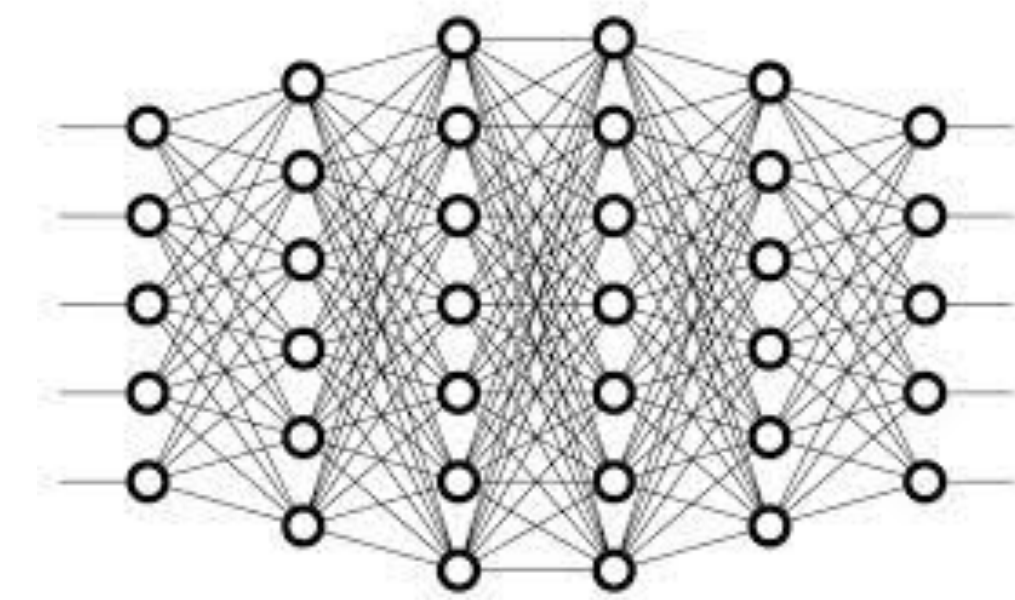


Recap: Supervised Learning in a Nutshell

$$P(D | \boldsymbol{\theta}) = \prod_{i=0}^{|D|} \text{Ber}(y_i | \sigma(f(\mathbf{x}_i; \boldsymbol{\theta}))) \quad \sigma(a) \triangleq \frac{1}{1 + e^{-a}} \quad \text{multinomial, HMMs...}$$

$$f(\mathbf{x}; \boldsymbol{\theta}) = w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + w_5 x_5 + b$$

$$\boldsymbol{\theta} = \{w_1, w_2, w_3, w_4, w_5, b\} \quad \text{parameters}$$



$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} P(D | \boldsymbol{\theta})$$

MLE, MAP ...

Gradient Descent, SGD, Adam ...

Recap: Supervised Learning in a Nutshell

First Step!

