

10-701

Probability and MLE

(brief) intro to probability

Basic notations

- Random variable
 - referring to an element / event whose status is unknown:
 $A = \text{"it will rain tomorrow"}$
- Domain (usually denoted by Ω)
 - The set of values a random variable can take:
 - " $A = \text{The stock market will go up this year}$ ": Binary
 - " $A = \text{Number of Steelers wins in 2019}$ ": Discrete
 - " $A = \text{\% change in Google stock in 2019}$ ": Continuous

Axioms of probability (Kolmogorov's axioms)

A variety of useful facts can be derived from just three axioms:

1. $0 \leq P(A) \leq 1$
2. $P(\text{true}) = 1, P(\text{false}) = 0$
3. $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

There have been several other attempts to provide a foundation for probability theory. Kolmogorov's axioms are the most widely used.

Priors

Degree of belief
in an event in the
absence of any
other information

No rain



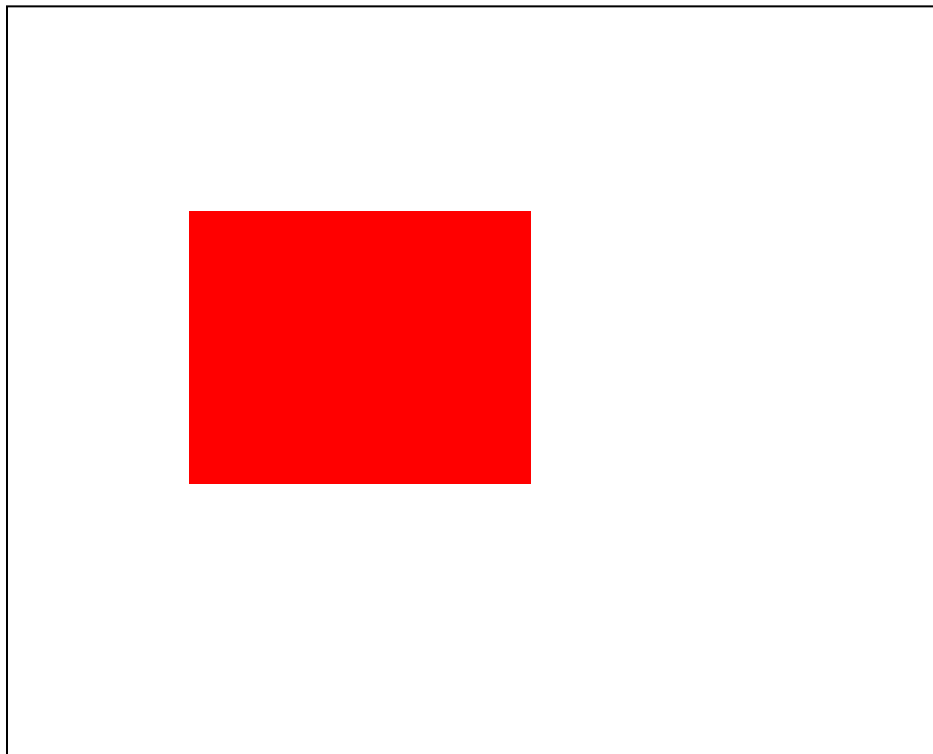
$$P(\text{rain tomorrow}) = 0.2$$

$$P(\text{no rain tomorrow}) = 0.8$$

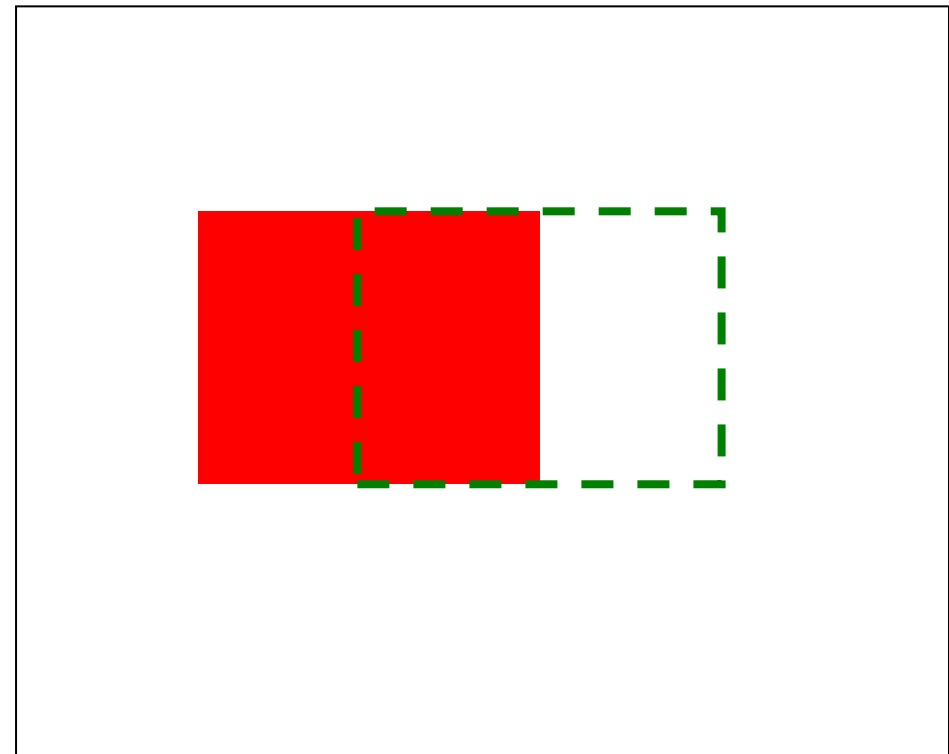
Conditional probability

- $P(A = 1 \mid B = 1)$: The fraction of cases where A is true if B is true

$$P(A = 0.2)$$



$$P(A|B = 0.5)$$



Conditional probability

- In some cases, given knowledge of one or more random variables we can improve upon our prior belief of another random variable
- For example:

$$p(\text{slept in movie}) = 0.5$$

$$p(\text{slept in movie} \mid \text{liked movie}) = 1/4$$

$$p(\text{didn't sleep in movie} \mid \text{liked movie}) = 3/4$$

Slept	Liked
1	0
0	1
1	1
1	0
0	0
1	0
0	1
0	1

Joint distributions

- The probability that a *set* of random variables will take a specific value is their joint distribution.
- Notation: $P(A \wedge B)$ or $P(A,B)$
- Example: $P(\text{liked movie, slept})$

If we assume independence then

$$P(A,B)=P(A)P(B)$$

However, in many cases such an assumption may be too strong
(more later in the class)

Joint distribution (cont)

$P(\text{class size} > 20) = 0.6$

$P(\text{summer}) = 0.4$

$P(\text{class size} > 20, \text{summer}) = ?$

Evaluation of classes

Size	Time	Eval
30	R	2
70	R	1
12	S	2
8	S	3
56	R	1
24	S	2
10	S	3
23	R	3
9	R	2
45	R	1

Joint distribution (cont)

$P(\text{class size} > 20) = 0.6$

$P(\text{summer}) = 0.4$

$P(\text{class size} > 20, \text{summer}) = 0.1$

Evaluation of classes

Size	Time	Eval
30	R	2
70	R	1
12	S	2
8	S	3
56	R	1
24	S	2
10	S	3
23	R	3
9	R	2
45	R	1

Joint distribution (cont)

$P(\text{class size} > 20) = 0.6$

$P(\text{eval} = 1) = 0.3$

$P(\text{class size} > 20, \text{eval} = 1) = 0.3$

Size	Time	Eval
30	R	2
70	R	1
12	S	2
8	S	3
56	R	1
24	S	2
10	S	3
23	R	3
9	R	2
45	R	1

Joint distribution (cont)

$P(\text{class size} > 20) = 0.6$

$P(\text{eval} = 1) = 0.3$

$P(\text{class size} > 20, \text{eval} = 1) = 0.3$

Evaluation of classes

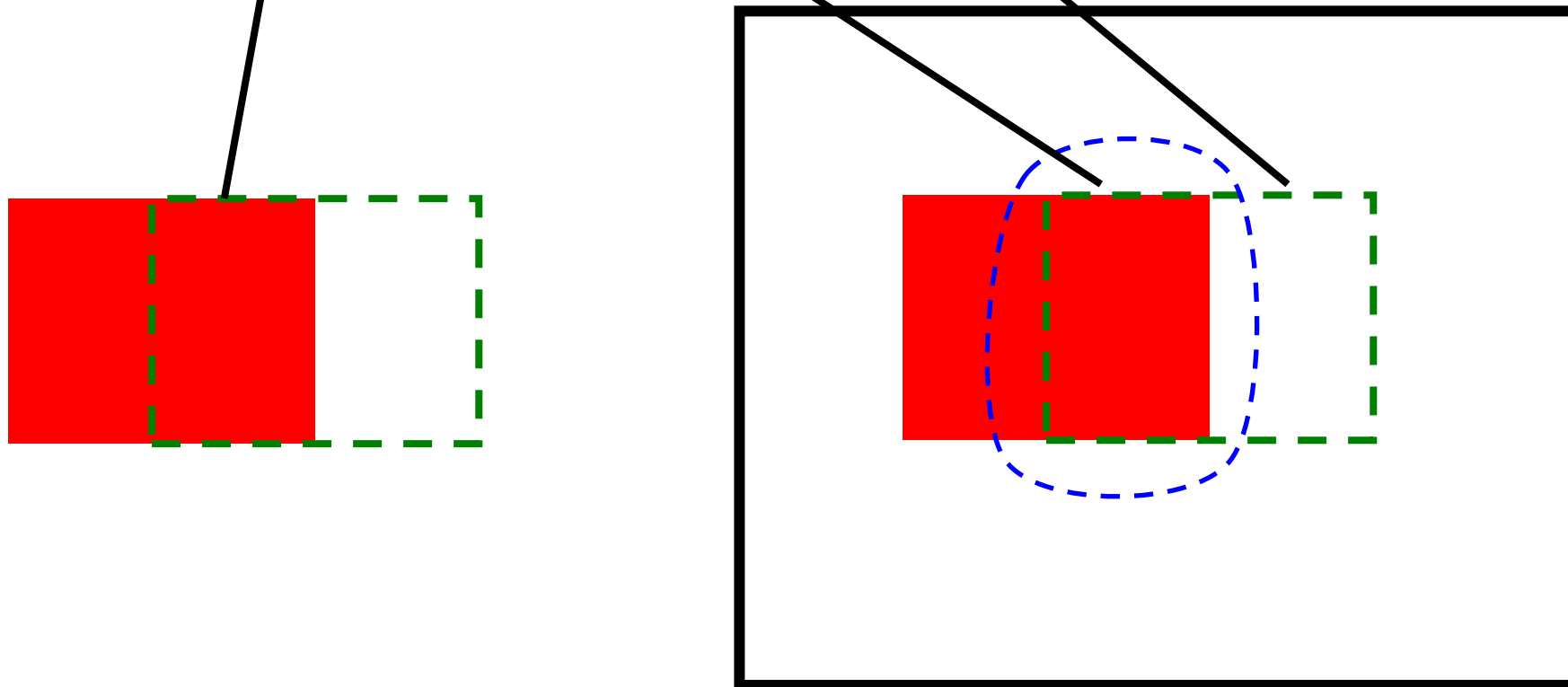
Size	Time	Eval
30	R	2
70	R	1
12	S	2
8	S	3
56	R	1
24	S	2
10	S	3
23	R	3
9	R	2
45	R	1

Chain rule

- The joint distribution can be specified in terms of conditional probability:

$$P(A,B) = P(A|B) \cdot P(B)$$

- Together with Bayes rule (which is actually derived from it) this is one of the most powerful rules in probabilistic reasoning



Bayes rule

- One of the most important rules for this class.
- Derived from the chain rule:

$$P(A,B) = P(A | B)P(B) = P(B | A)P(A)$$

- Thus,

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$



Thomas Bayes was an English clergyman who set out his theory of probability in 1764.

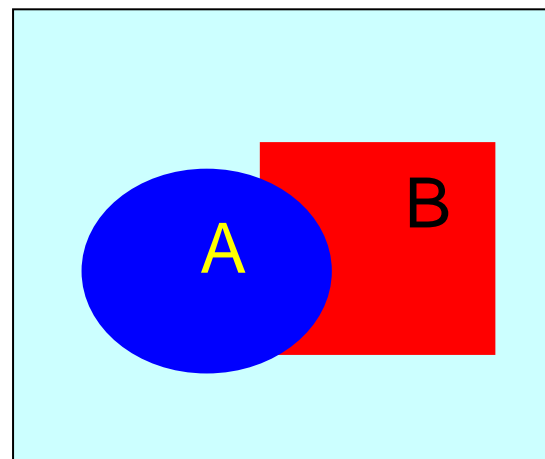
Bayes rule (cont)

Often it would be useful to derive the rule a bit further:

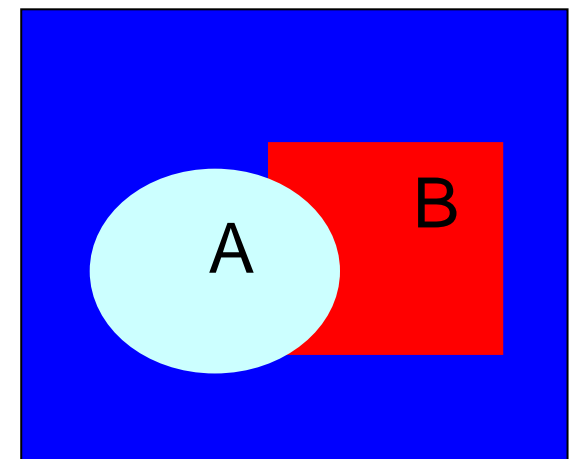
$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} = \frac{P(B|A)P(A)}{\sum_A P(B|A)P(A)}$$

This results from:
 $P(B) = \sum_A P(B,A)$

$P(B,A=1)$



$P(B,A=0)$



Bayes Rule for Continuous Distributions

- Standard form:

$$f(x|y) = \frac{f(y|x)f(x)}{f(y)}$$

- Replacing the bottom:

$$f(x|y) = \frac{f(y|x)f(x)}{\int f(y|x)f(x)dx}$$

AIDS test (Bayes rule)

Data

- Approximately 0.1% are infected
- Test detects all infections
- Test reports positive for 1% healthy people

AIDS test (Bayes rule)

Data

- Approximately 0.1% are infected
- Test detects all infections
- Test reports positive for 1% healthy people

Probability of having AIDS if test is positive:



AIDS test (Bayes rule)

Data

- Approximately 0.1% are infected
- Test detects all infections
- Test reports positive for 1% healthy people

Probability of having AIDS if test is positive:

$$\begin{aligned} P(a = 1|t = 1) &= \frac{P(t = 1|a = 1)P(a = 1)}{P(t = 1)} \\ &= \frac{P(t = 1|a = 1)P(a = 1)}{P(t = 1|a = 1)P(a = 1) + P(t = 1|a = 0)P(a = 0)} \\ &= \frac{1 \cdot 0.001}{1 \cdot 0.001 + 0.01 \cdot 0.999} = 0.091 \end{aligned}$$

Only 9%!...

Continuous distributions

Statistical Models

- Statistical models attempt to characterize properties of the population of interest
- For example, we might believe that repeated measurements follow a normal (Gaussian) distribution with some mean μ and variance σ^2 , $x \sim N(\mu, \sigma^2)$

where

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

and $\Theta=(\mu, \sigma^2)$ defines the parameters (mean and variance) of the model.

How much do grad students sleep?

- Lets try to estimate the distribution of the time students spend sleeping (outside class).

Possible statistics

- **X**

Sleep time

- **Mean of X:**

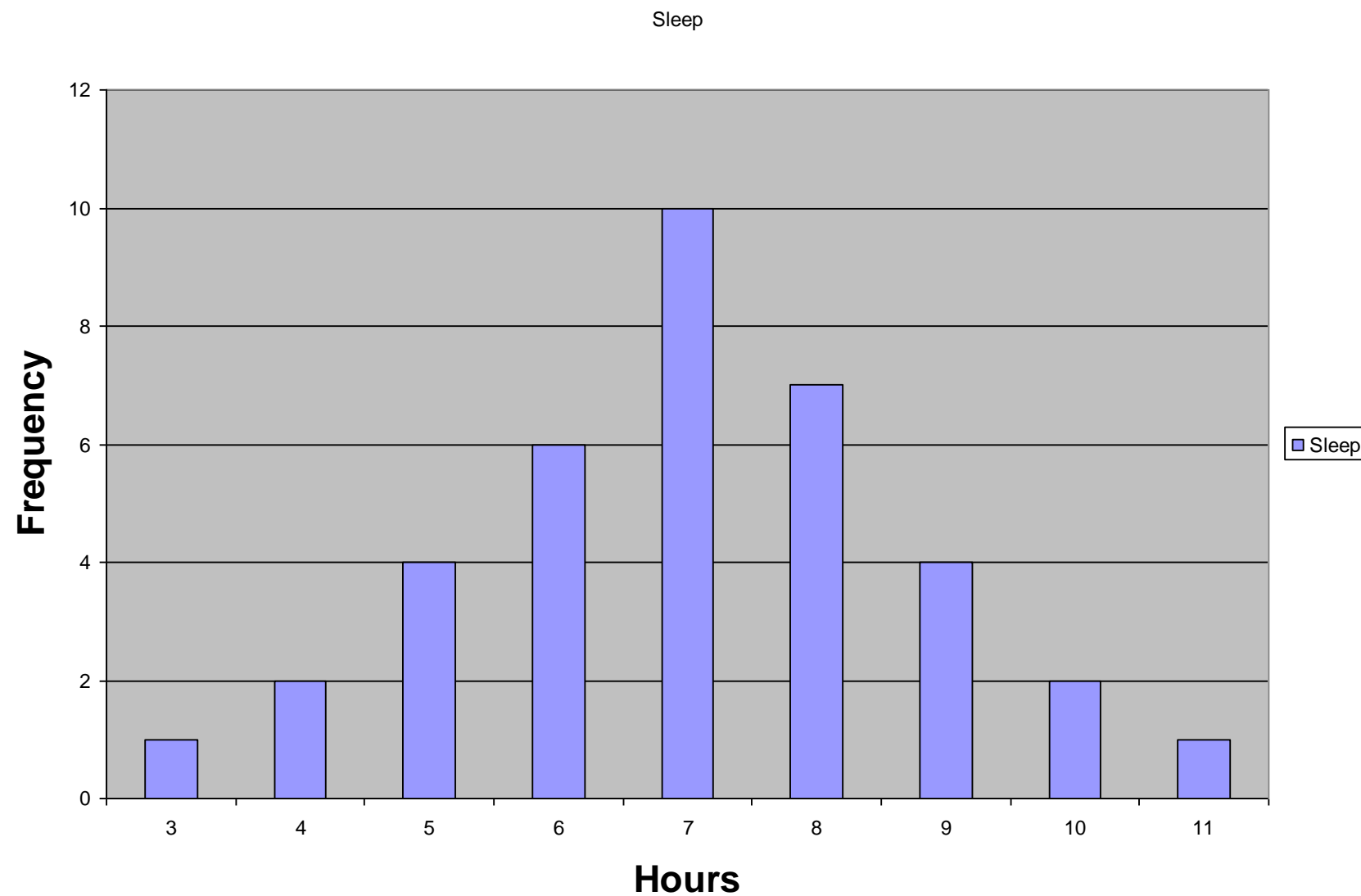
$$E\{X\}$$

7.03

- **Variance of X:**

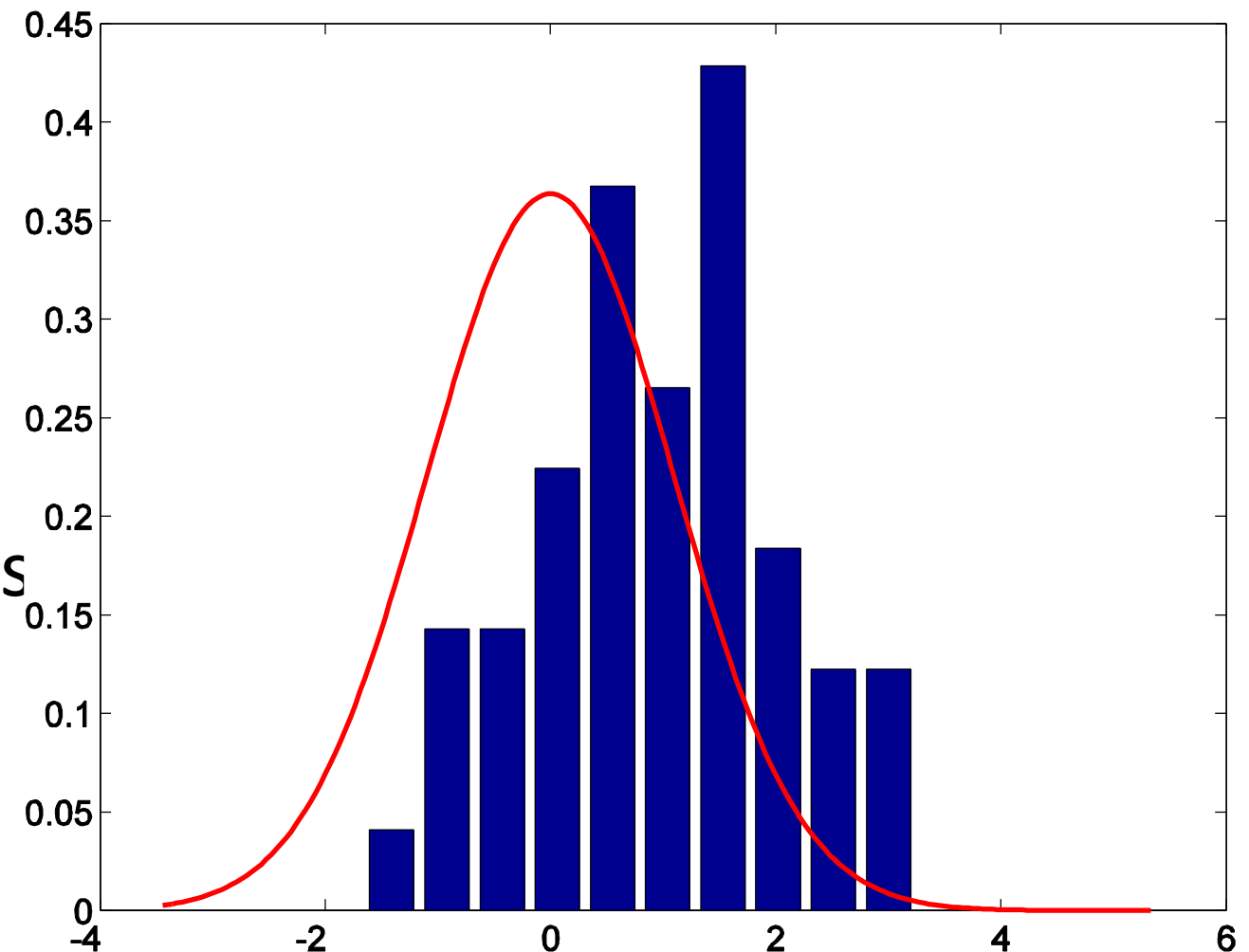
$$\text{Var}\{X\} = E\{(X - E\{X\})^2\}$$

3.05



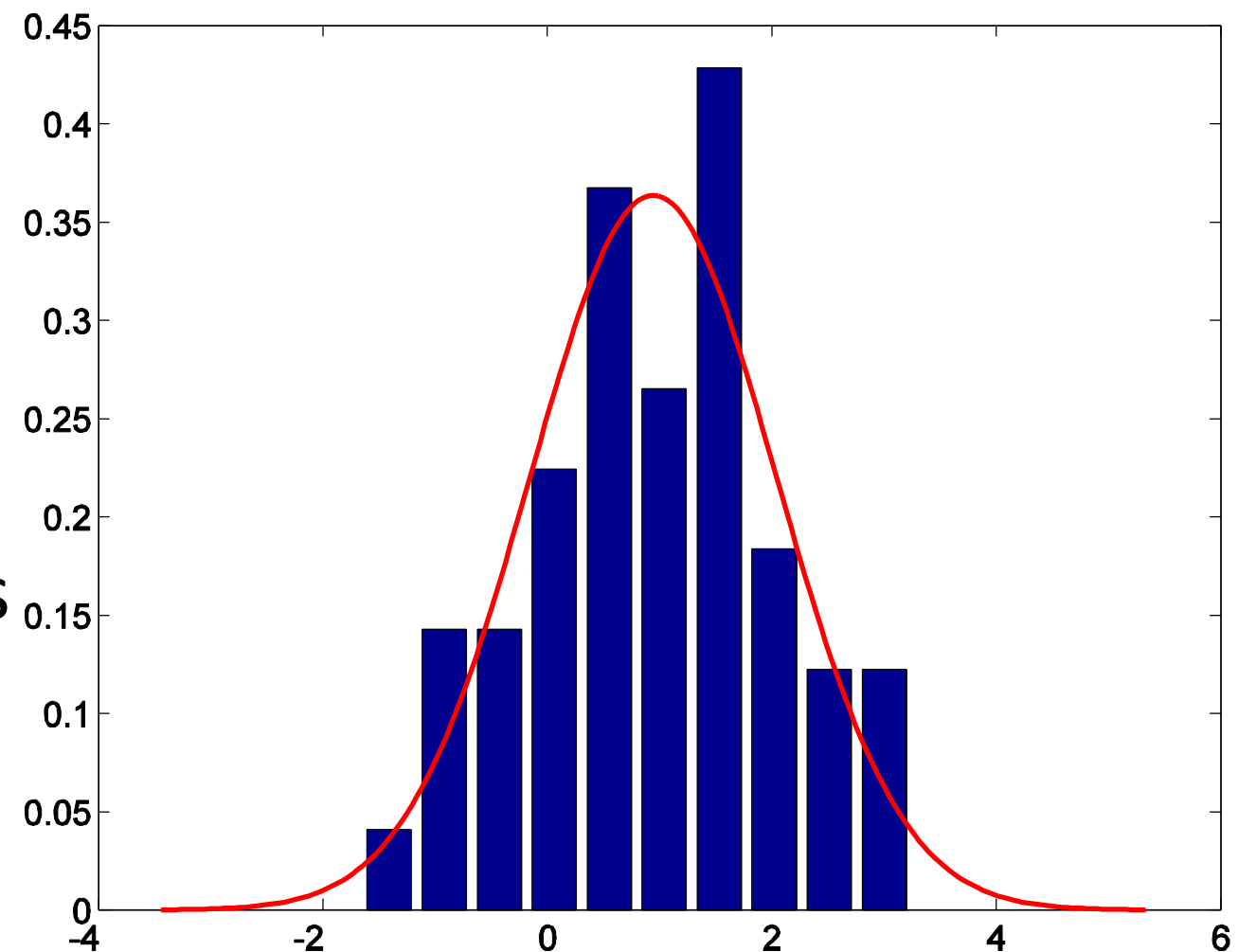
The Parameters of Our Model

- A statistical model is a **collection** of distributions; the **parameters** specify individual distributions $x \sim N(\mu, \sigma^2)$
- We need to adjust the parameters so that the resulting distribution **fits** the data well



The Parameters of Our Model

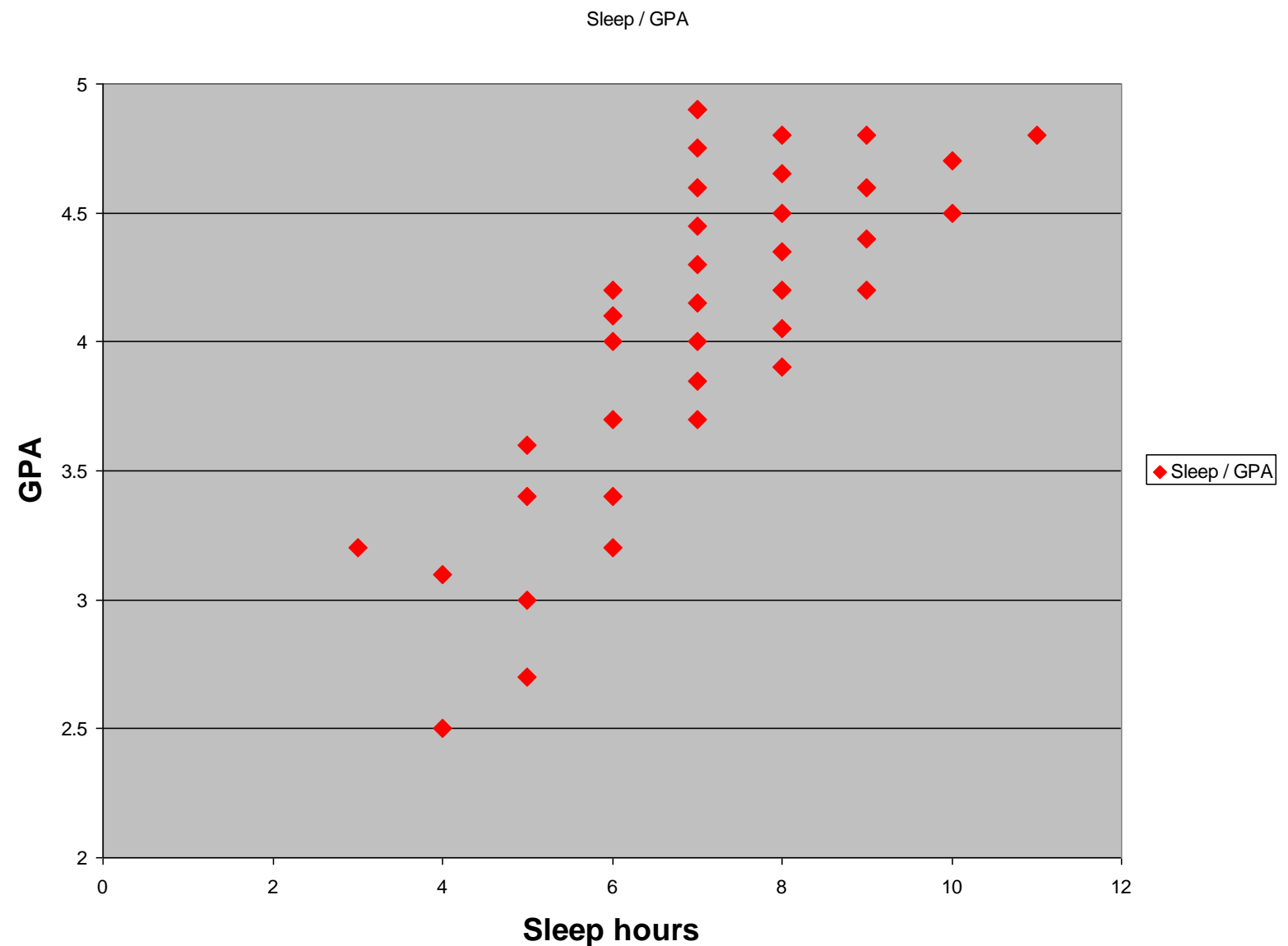
- A statistical model is a **collection** of distributions; the **parameters** specify individual distributions $x \sim N(\mu, \sigma^2)$
- We need to adjust the parameters so that the resulting distribution **fits** the data well



Covariance: Sleep vs. GPA

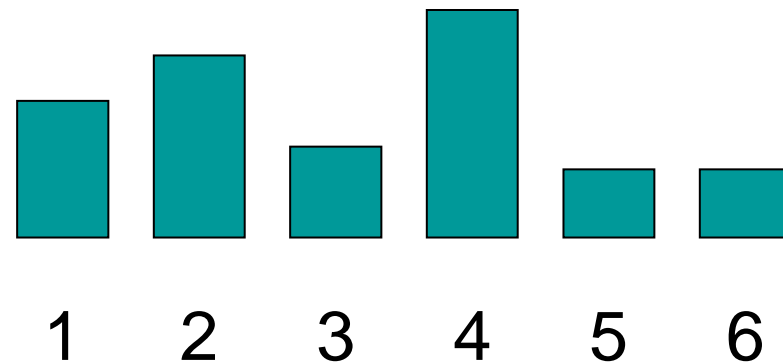
- **Co-Variance of X1, X2:**

$$\begin{aligned} \text{Covariance}\{X1, X2\} &= \\ E\{(X1 - E\{X1\})(X2 - E\{X2\})\} \\ &= 0.88 \end{aligned}$$



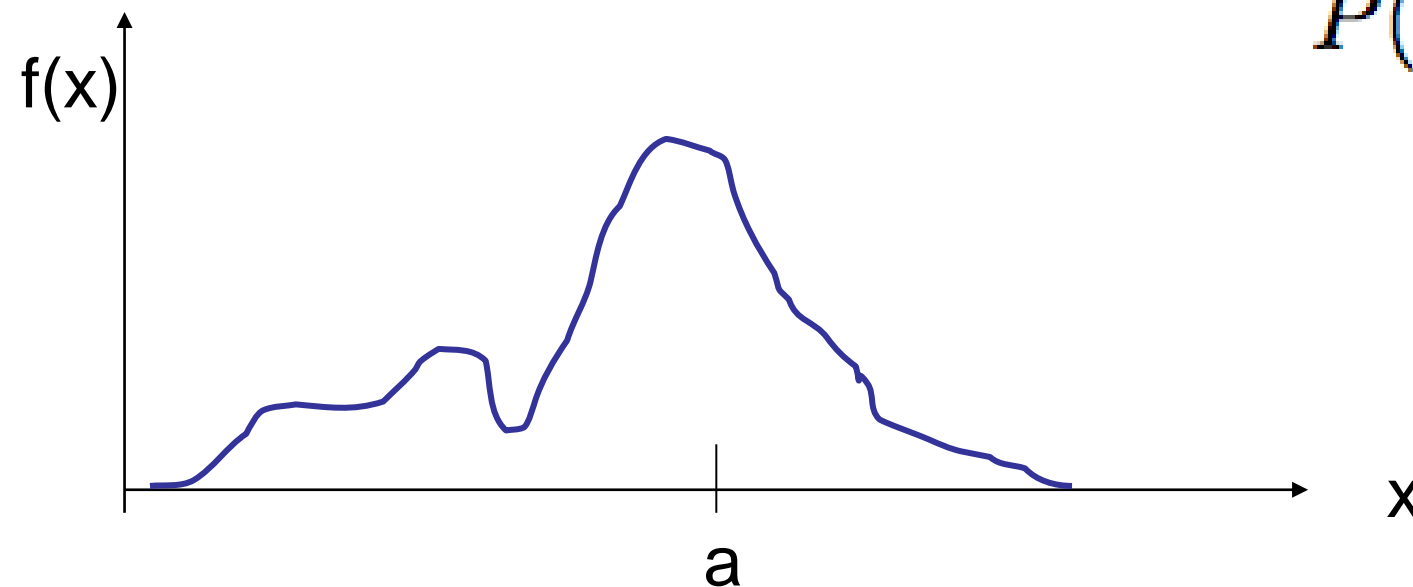
Probability Density Function

- Discrete distributions



$$\sum_i P(X = x_i) = 1$$

- Continuous: Cumulative Density Function (CDF): $F(a)$



$$P(x \leq a) = \int_{-\infty}^a f(\tau) d\tau$$

Cumulative Density Functions

- Total probability
$$P(\Omega) = \int_{-\infty}^{\infty} f(x)dx = 1$$

- Probability Density Function (PDF)
$$\frac{d}{dx}F(x) = f(x)$$

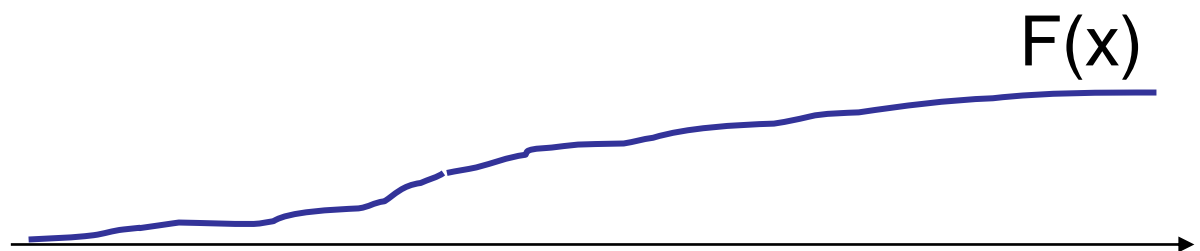
- Properties:

$$P(a \leq x \leq b) = \int_a^b f(x)dx = F(b) - F(a)$$

$$\lim_{x \rightarrow -\infty} F(x) = 0$$

$$\lim_{x \rightarrow \infty} F(x) = 1$$

$$F(a) \geq F(b) \quad \forall a \geq b$$



Density estimation: The Bayesian way

Your first consulting job

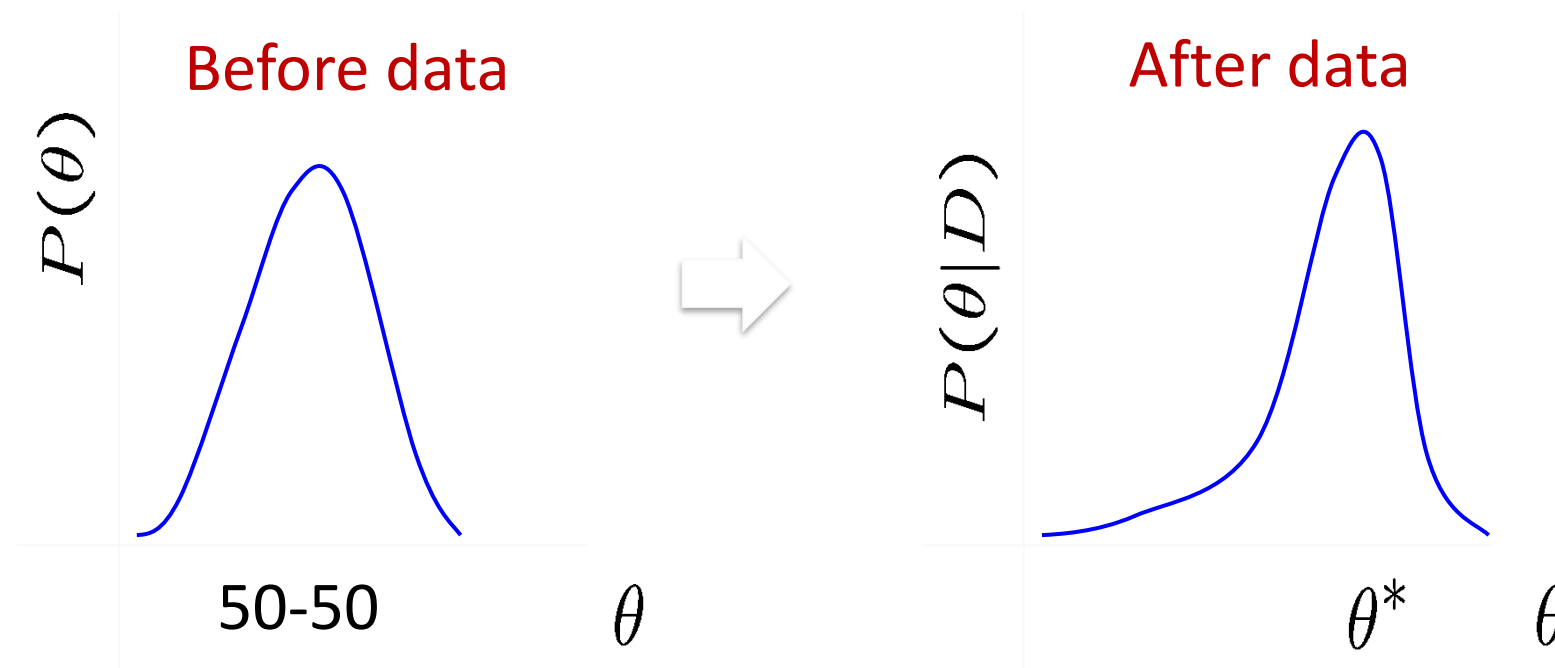
- A billionaire from the suburbs of Seattle asks you a question:
 - He says: I have a coin, if I flip it, what's the probability it will fall with the head up?
 - You say: Please flip it a few times:



- You say: The probability is: **3/5** because... frequency of heads in all flips
- **He says: But can I put money on this estimate?**
- You say: ummm.... Maybe not.
 - Not enough flips (less than sample complexity)

What about prior knowledge?

- Billionaire says: Wait, I know that the coin is “close” to 50-50. What can you do for me now?
- **You say: I can learn it the Bayesian way...**
- Rather than estimating a single θ , we obtain a distribution over possible values of θ



Bayesian Learning

- Use Bayes rule:

$$P(\theta \mid \mathcal{D}) = \frac{P(\mathcal{D} \mid \theta)P(\theta)}{P(\mathcal{D})}$$

- Or equivalently:

$$P(\theta \mid \mathcal{D}) \propto P(\mathcal{D} \mid \theta)P(\theta)$$

posterior likelihood prior



Prior distribution

- From where do we get the prior?
 - Represents expert knowledge (philosophical approach)
 - Simple posterior form (engineer's approach)
- Uninformative priors:
 - Uniform distribution
- Conjugate priors:
 - Closed-form representation of posterior
 - $P(q)$ and $P(q|D)$ have the same algebraic form as a function of θ

Conjugate Prior

- $P(q)$ and $P(q|D)$ have the same form as a function of θ

Eg. 1 Coin flip problem

Likelihood given Bernoulli model:

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

If prior is Beta distribution,

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

Then posterior is Beta distribution

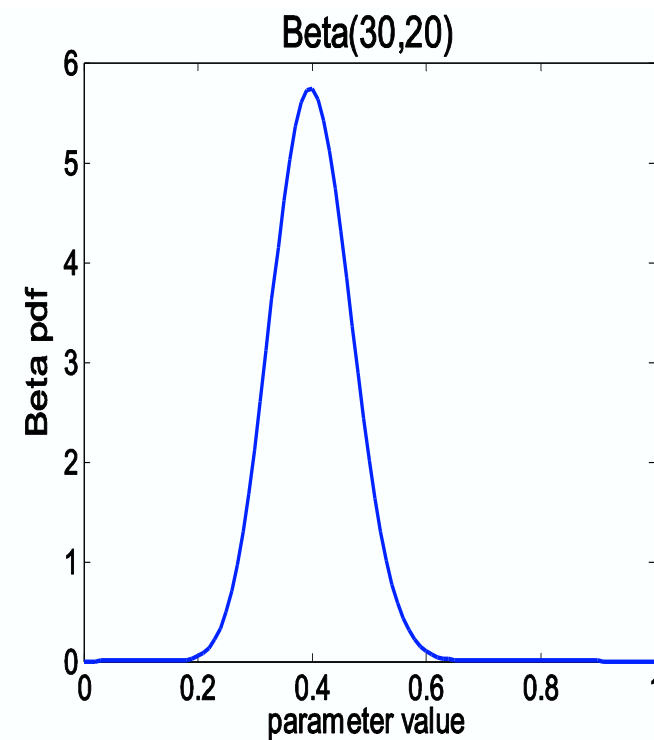
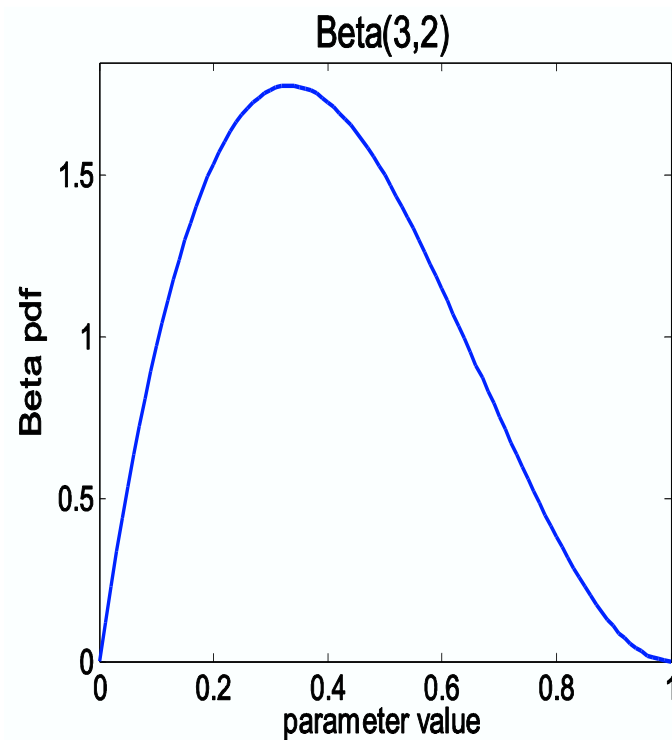
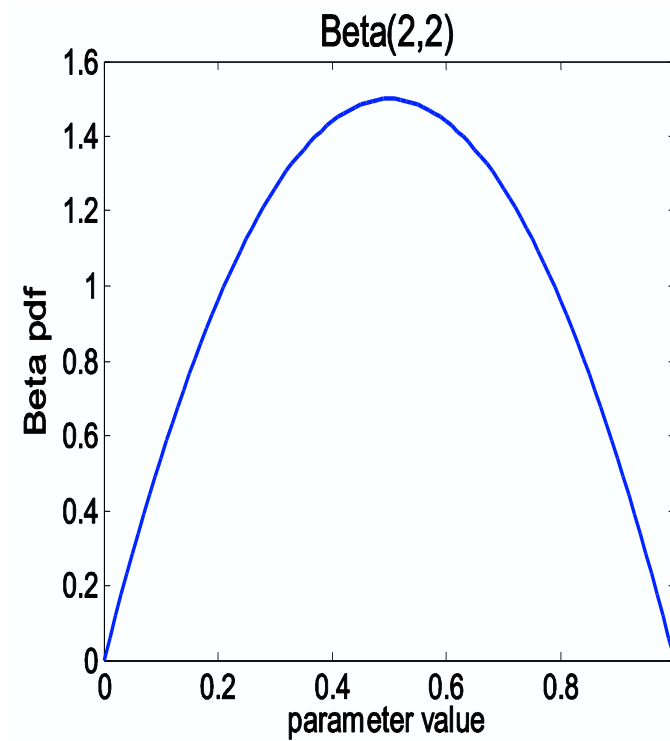
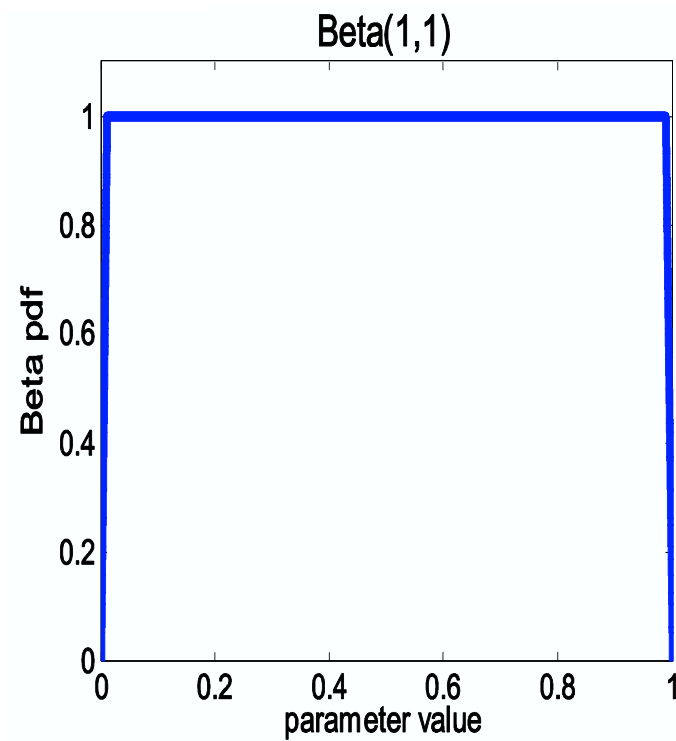
$$P(\theta|D) \sim \text{Beta}(\beta_H + \alpha_H, \beta_T + \alpha_T)$$



Beta distribution

$Beta(\beta_H, \beta_T)$

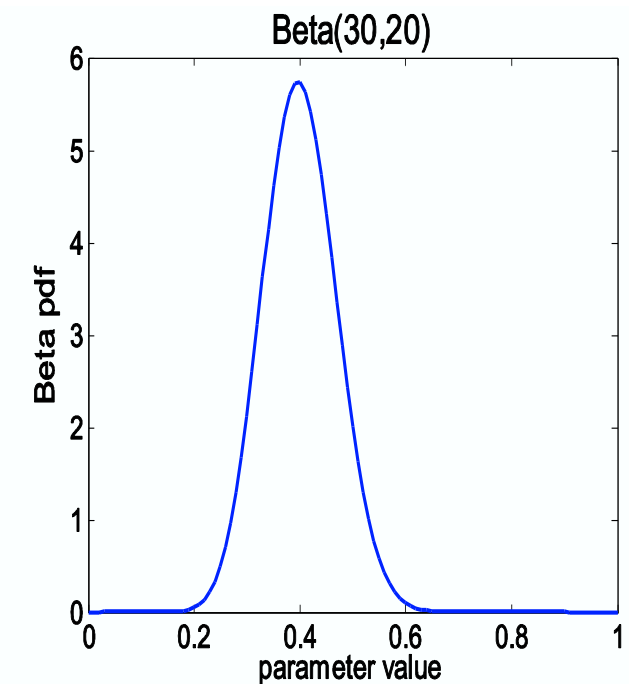
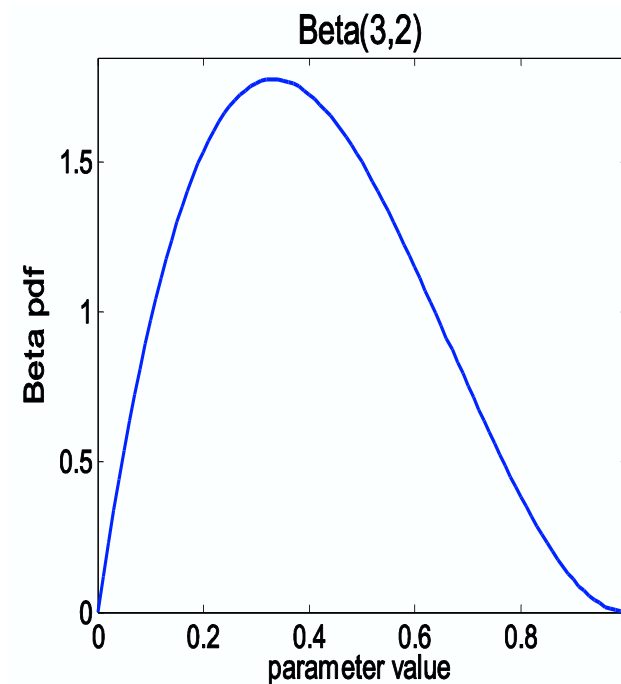
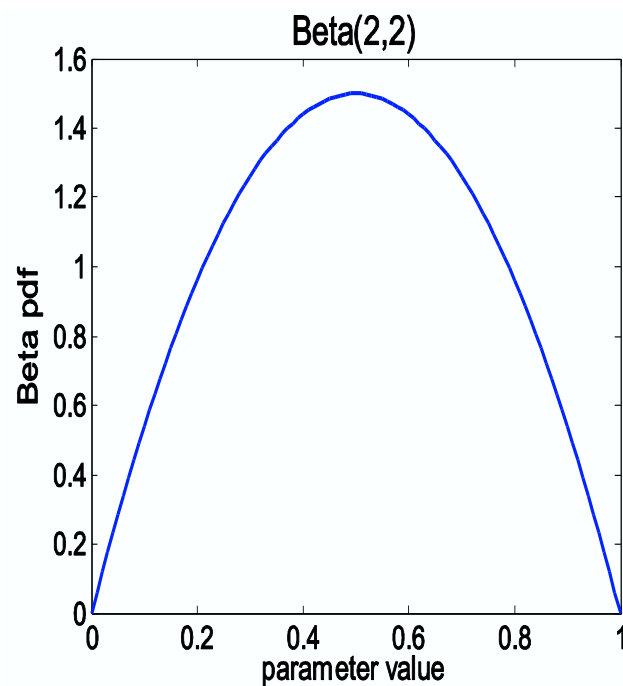
More concentrated as values of β_H, β_T increase



Beta conjugate prior

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

$$P(\theta|D) \sim \text{Beta}(\beta_H + \alpha_H, \beta_T + \alpha_T)$$



As $n = \alpha_H + \alpha_T$
increases

As we get more samples, effect of prior is “washed out”

Conjugate Prior

- $P(\theta)$ and $P(\theta|D)$ have the same form

Eg. 2 Dice roll problem (6 outcomes instead of 2)

Likelihood is $\sim \text{Multinomial}(\theta = \{\theta_1, \theta_2, \dots, \theta_k\})$

$$P(\mathcal{D} | \theta) = \theta_1^{\alpha_1} \theta_2^{\alpha_2} \dots \theta_k^{\alpha_k}$$

If prior is Dirichlet distribution,

$$P(\theta) = \frac{\prod_{i=1}^k \theta_i^{\beta_i-1}}{B(\beta_1, \dots, \beta_k)} \sim \text{Dirichlet}(\beta_1, \dots, \beta_k)$$

Then posterior is Dirichlet distribution

$$P(\theta|D) \sim \text{Dirichlet}(\beta_1 + \alpha_1, \dots, \beta_k + \alpha_k)$$



For Multinomial, conjugate prior is Dirichlet distribution.

Posterior Distribution

- The approach seen so far is what is known as a **Bayesian** approach
- Prior information encoded as a **distribution** over possible values of parameter
- Using the Bayes rule, you get an updated **posterior** distribution over parameters, which you provide with flourish to the Billionaire
- But the billionaire is not impressed
 - Distribution? I just asked for one number: is it $3/5$, $1/2$, what is it?
 - How do we go from a distribution over parameters, to a single estimate of the true parameters?

Maximum A Posteriori Estimation

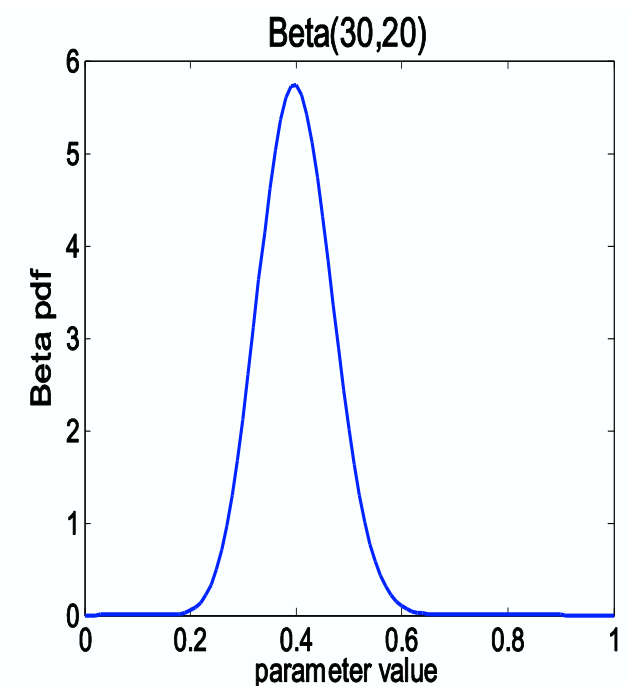
Choose θ that maximizes a posterior probability

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

MAP estimate of probability of head:

$$P(\theta|D) \sim \text{Beta}(\beta_H + \alpha_H, \beta_T + \alpha_T)$$

$$\hat{\theta}_{MAP} = \frac{\alpha_H + \beta_H - 1}{\alpha_H + \beta_H + \alpha_T + \beta_T - 2}$$

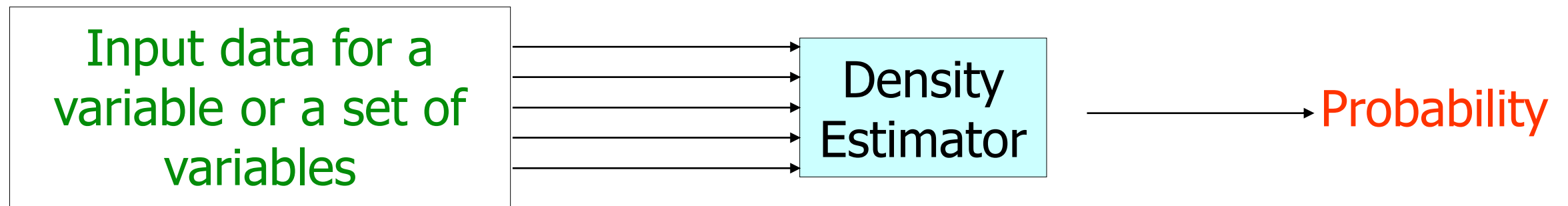


Mode of Beta
distribution

Density estimation: Learning

Density Estimation

- A Density Estimator learns a mapping from a set of attributes to a Probability



Density estimation

- Estimate the distribution (or conditional distribution) of a random variable
- Types of variables:
 - Binary
coin flip, alarm
 - Discrete
dice, car model year
 - Continuous
height, weight, temp.,

When do we need to estimate densities?

- Density estimators are critical ingredients in several of the ML algorithms we will discuss
- In some cases these are combined with other inference types for more involved algorithms (i.e. EM) while in others they are part of a more general process (learning in BNs and HMMs)

Density estimation

- Binary and discrete variables:

Easy: Just count!

- Continuous variables:

Harder (but just a bit): Fit a model

Learning a density estimator for discrete variables

$$\hat{P}(x_i = u) = \frac{\text{\# records in which } x_i = u}{\text{total number of records}}$$

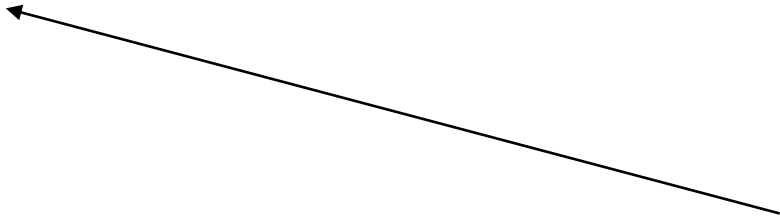
A trivial learning algorithm!

But why is this true?

Maximum Likelihood Principle

We can define the likelihood of the data given the model as follows:

$$\hat{P}(\text{dataset} \mid M) = \hat{P}(x_1 \wedge x_2 \dots \wedge x_n \mid M) = \prod_{k=1}^n \hat{P}(x_k \mid M)$$



M is our model (usually a collection of parameters)

For example M is

- The probability of 'head' for a coin flip
- The probabilities of observing 1,2,3,4 and 5 for a dice
- etc.

Maximum Likelihood Principle

$$\hat{P}(\text{dataset} \mid M) = \hat{P}(x_1 \wedge x_2 \dots \wedge x_n \mid M) = \prod_{k=1}^n \hat{P}(x_k \mid M)$$

- Our goal is to determine the values for the parameters in M
- We can do this by maximizing the probability of generating the observed samples
- For example, let Θ be the probabilities for a coin flip
- Then

$$L(x_1, \dots, x_n \mid \Theta) = p(x_1 \mid \Theta) \dots p(x_n \mid \Theta)$$

- The observations (different flips) are assumed to be independent
- For such a coin flip with $P(H)=q$ the best assignment for Θ_h is

$$\operatorname{argmax}_q = \#H/\#\text{samples}$$

- Why?

Maximum Likelihood Principle: Binary variables

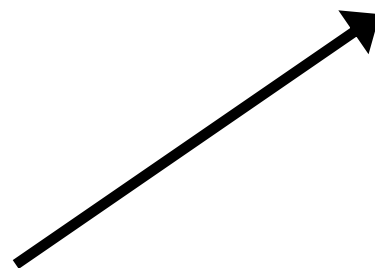
- For a binary random variable A with $P(A=1)=q$
 $\operatorname{argmax}_q = \#1/\#\text{samples}$

- Why?

Data likelihood: $P(D | M) = q^{n_1} (1 - q)^{n_2}$

We would like to find: $\operatorname{arg max}_q q^{n_1} (1 - q)^{n_2}$

Omitting terms that do
not depend on q



Maximum Likelihood Principle

Data likelihood: $P(D | M) = q^{n_1} (1 - q)^{n_2}$

We would like to find: $\arg \max_q q^{n_1} (1 - q)^{n_2}$

$$\frac{\partial}{\partial q} q^{n_1} (1 - q)^{n_2} = n_1 q^{n_1-1} (1 - q)^{n_2} - q^{n_1} n_2 (1 - q)^{n_2-1}$$

$$\frac{\partial}{\partial q} = 0 \Rightarrow$$

$$n_1 q^{n_1-1} (1 - q)^{n_2} - q^{n_1} n_2 (1 - q)^{n_2-1} = 0 \Rightarrow$$

$$q^{n_1-1} (1 - q)^{n_2-1} (n_1 (1 - q) - q n_2) = 0 \Rightarrow$$

$$n_1 (1 - q) - q n_2 = 0 \Rightarrow$$

$$n_1 = n_1 q + n_2 q \Rightarrow$$

$$q = \frac{n_1}{n_1 + n_2}$$

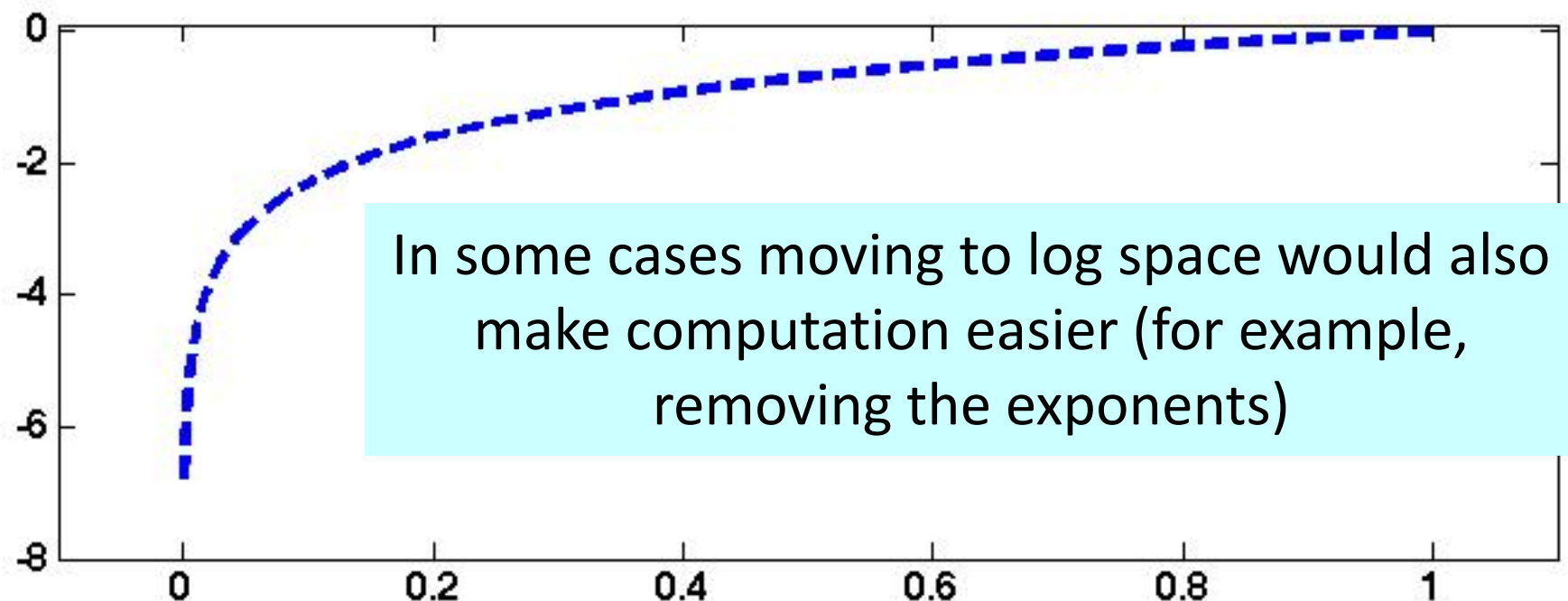
Log Probabilities

When working with products, probabilities of entire datasets often get too small. A possible solution is to use the log of probabilities, often termed 'log likelihood'

$$\log \hat{P}(\text{dataset} \mid M) = \log \prod_{k=1}^n \hat{P}(x_k \mid M) = \sum_{k=1}^n \log \hat{P}(x_k \mid M)$$

Maximizing this likelihood function is the same as maximizing $P(\text{dataset} \mid M)$

Log values
between 0 and 1



Density estimation

- Binary and discrete variables:

Easy: Just count!

- Continuous variables:

Harder (but just a bit): Fit a model

But what if we
only have very few
samples?



Maximum Likelihood Principle

- We can fit statistical models by maximizing the probability of generating the observed samples:

$$L(x_1, \dots, x_n \mid \Theta) = p(x_1 \mid \Theta) \dots p(x_n \mid \Theta)$$

(the samples are assumed to be independent)

- In the Gaussian case we simply set the mean and the variance to the sample mean and the sample variance:

$$\overline{\mu} = \frac{1}{n} \sum_{i=1}^n x_i \qquad \overline{\sigma^2} = \frac{1}{n} \sum_{i=1}^n (x_i - \overline{\mu})^2$$

Why?

MLE vs. MAP

- Maximum Likelihood estimation (MLE)

Choose value that maximizes the probability of observed data

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

- Maximum *a posteriori* (MAP) estimation

Choose value that is most probable given observed data and prior belief

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