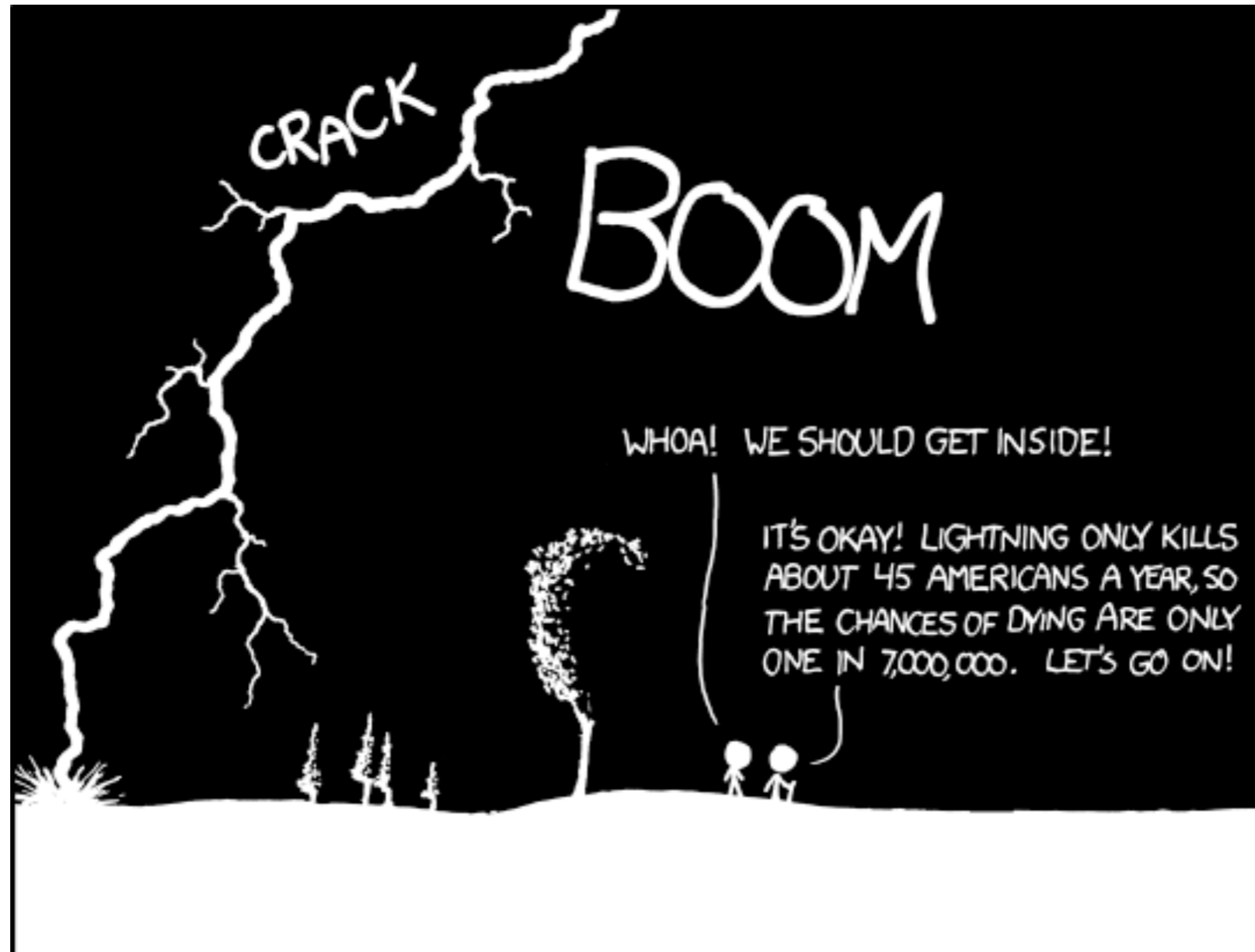


Parameter estimation

Conditional risk



THE ANNUAL DEATH RATE AMONG PEOPLE WHO KNOW THAT STATISTIC IS ONE IN SIX.

Formalizing the problem

- Specify random variables we care about
 - e.g., Commute Time
 - e.g., Heights of buildings in a city
- We might then pick a particular distribution over these random variables
 - e.g., Say we think our variable is Gaussian
- Now want a way to use the data to inform the model choice
 - e.g., Let data tell us the parameters for that Gaussian

Parameter estimation

- Assume that we are given some model class, M ,
 - e.g., Gaussian with parameters μ and σ $\mathcal{N}(\mu, \sigma^2)$
 - selection of model from the class corresponds to selecting μ , σ
- Now want to select “best” model; how do we define best?
 - Generally assume data comes from that model class; might want to find model that best explains the data (or most likely given the data)
 - Might want most likely model, with preference for “important” samples
 - Might want most likely model, that also matches expert prior info
 - Might want most likely model, that is the simplest (least parameters)
- These additional requirements are usually in place to enable better generalization to unseen data

How can we identify the parameters for a distribution?

- All we get is access to samples x_1, \dots, x_n from the true underlying distribution $p(x)$
- We want to find parameters θ so that p_θ approximates the true distribution p well
 - e.g., θ is the mean and variance for a Gaussian
 - e.g., θ is the parameters to a Gamma distribution (commute times)
- Any ideas?

What about in the prediction setting? That feels easier...

- Here, the goal is to learn a function f for some inputs \mathbf{x} to make predictions about targets y
- How do we identify the “best” f in our set of functions F ?

$$\mathcal{F} = \{f : \mathbb{R}^d \rightarrow \mathbb{R} \mid f(\mathbf{x}) = \mathbf{x}^\top \mathbf{w} \text{ for some } \mathbf{w} \in \mathbb{R}^d\}$$

- Any ideas?

- Maybe $\min_{f \in \mathcal{F}} \sum_{i=1}^n (f(\mathbf{x}_i) - y_i)^2$

- What does that entail? What if y_i is always 0 or 1?

Using the language of probability

- It can be hard to guess an objective function that tells you what is the best model in your class
- A natural objective is: find the model that is the most likely given the data

$$\arg \max_{\theta} p(\theta | \mathcal{D})$$

- Reflects the inherent uncertainty in identifying the model, since many models are possible given only samples

Some notation

$$\sum_{i=1}^n x_i = x_1 + x_2 + \dots + x_n$$

$$\prod_{i=1}^n x_i = x_1 x_2 \dots x_n$$

$$c : \mathbb{R}^d \rightarrow \mathbb{R}$$

$$\mathbf{w}^* = \arg \max_{\mathbf{w} \in \mathbb{R}^d} c(\mathbf{w})$$

$$c(\mathbf{w}^*) = \max_{\mathbf{w} \in \mathbb{R}^d} c(\mathbf{w})$$

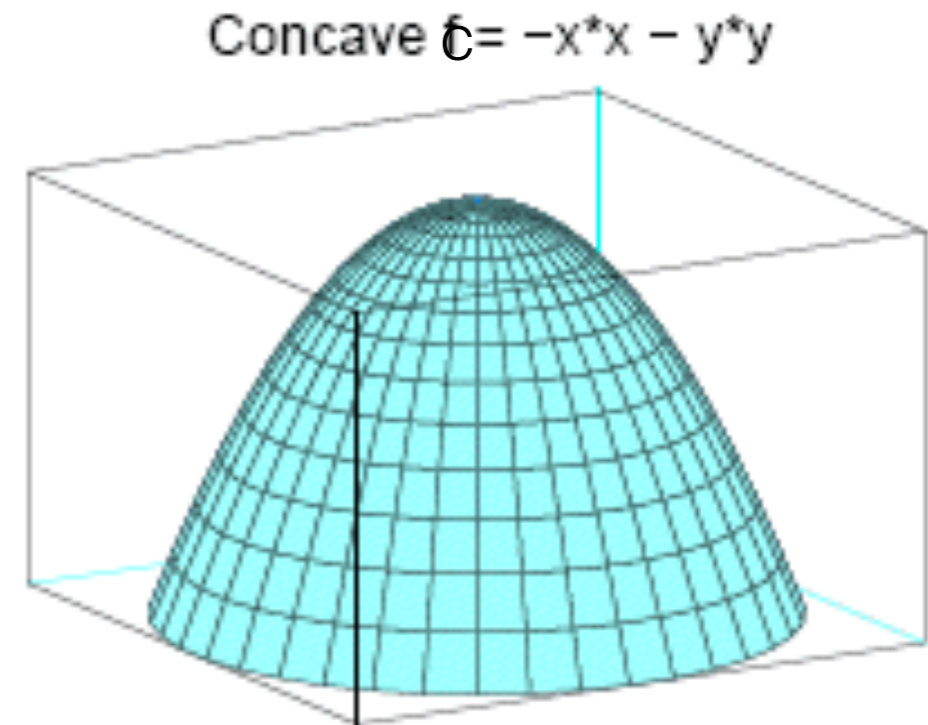
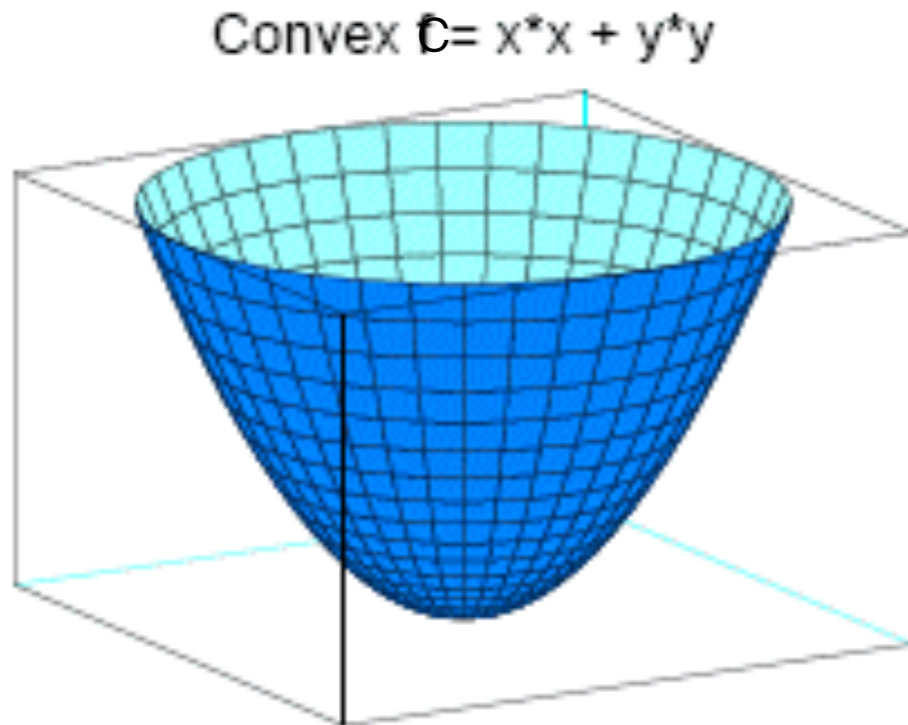
\mathcal{F} is a set of models

$$\text{e.g., } \mathcal{F} = \{\mathcal{N}(\mu, \sigma) \mid (\mu, \sigma) \in \mathbb{R}^2, \sigma > 0\}$$

$$\text{e.g., } \mathcal{F} = \{\mathbf{w} \in \mathbb{R}^d \mid f(\mathbf{x}) = \mathbf{x}^\top \mathbf{w}\}$$

Definition of optimization

- We select some (error) function c we care about



- Maximizing c means we are finding largest point
- Minimizing c means we are finding smallest point

Maximum a posteriori (MAP) estimation

$$f_{\text{MAP}} = \arg \max_{f \in \mathcal{F}} p(f | \mathcal{D})$$

- Want the f that is most likely, given the data
- $p(f | \mathcal{D})$ is the **posterior distribution** of the model given data
 - e.g., \mathcal{F} could be the space of Gaussian distributions, the model is f and $f(x)$ returns probability/density of a point x
 - e.g., we could assume x is Gaussian distributed with variance = 1, and so \mathcal{F} could be the reals, and the model f is the mean

Question: What is the function we are optimizing and what are the parameters we are learning?

MAP

$$f_{\text{MAP}} = \arg \max_{f \in \mathcal{F}} p(f | \mathcal{D})$$

- $p(f | \mathcal{D})$ is the **posterior distribution** of the model given data
- e.g., we could assume x is Gaussian distributed with variance = 1, and so \mathcal{F} could be the reals, and the model f is the mean

Question: What is the function we are optimizing and what are the parameters we are learning?

$$c(f) = p(\text{mean is } f | \mathcal{D})$$

$$\max_{f \in \mathbb{R}} c(f)$$

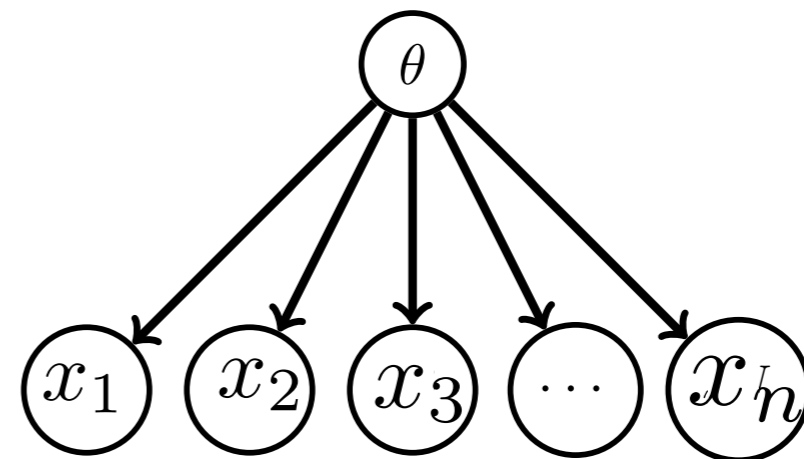
Maximum a posteriori (MAP)

$$f_{\text{MAP}} = \arg \max_{f \in \mathcal{F}} p(f | \mathcal{D})$$

- $p(f | \mathcal{D})$ is the **posterior distribution** of the model given data
- In discrete spaces: $p(f | \mathcal{D})$ is a PMF
 - the MAP estimate is exactly the most probable model
 - e.g., bias of coin is 0.1, 0.5, or 0.7, $p(f = 0.1 | \mathcal{D})$, ...
- In continuous spaces: $p(f | \mathcal{D})$ is a PDF
 - the MAP estimate is the model with the largest value of the posterior density function
 - e.g., bias of a coin is in $[0, 1]$
- But what is $p(f | \mathcal{D})$? Do we pick it? If so, how?

Example: Posterior for discrete distributions

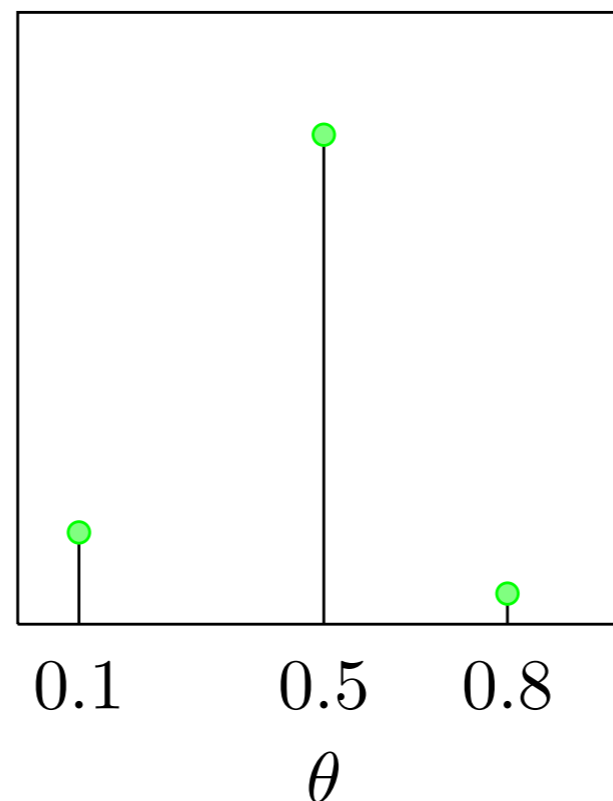
- Imagine you are flipping a biased coin; the model parameter is the bias of the coin, **theta**
- You get a dataset $D = \{x_1, \dots, x_n\}$ of coin flips, where $x_i = 1$ if it was heads, and $x_i = 0$ if it was tails



Example: Prior before seeing data

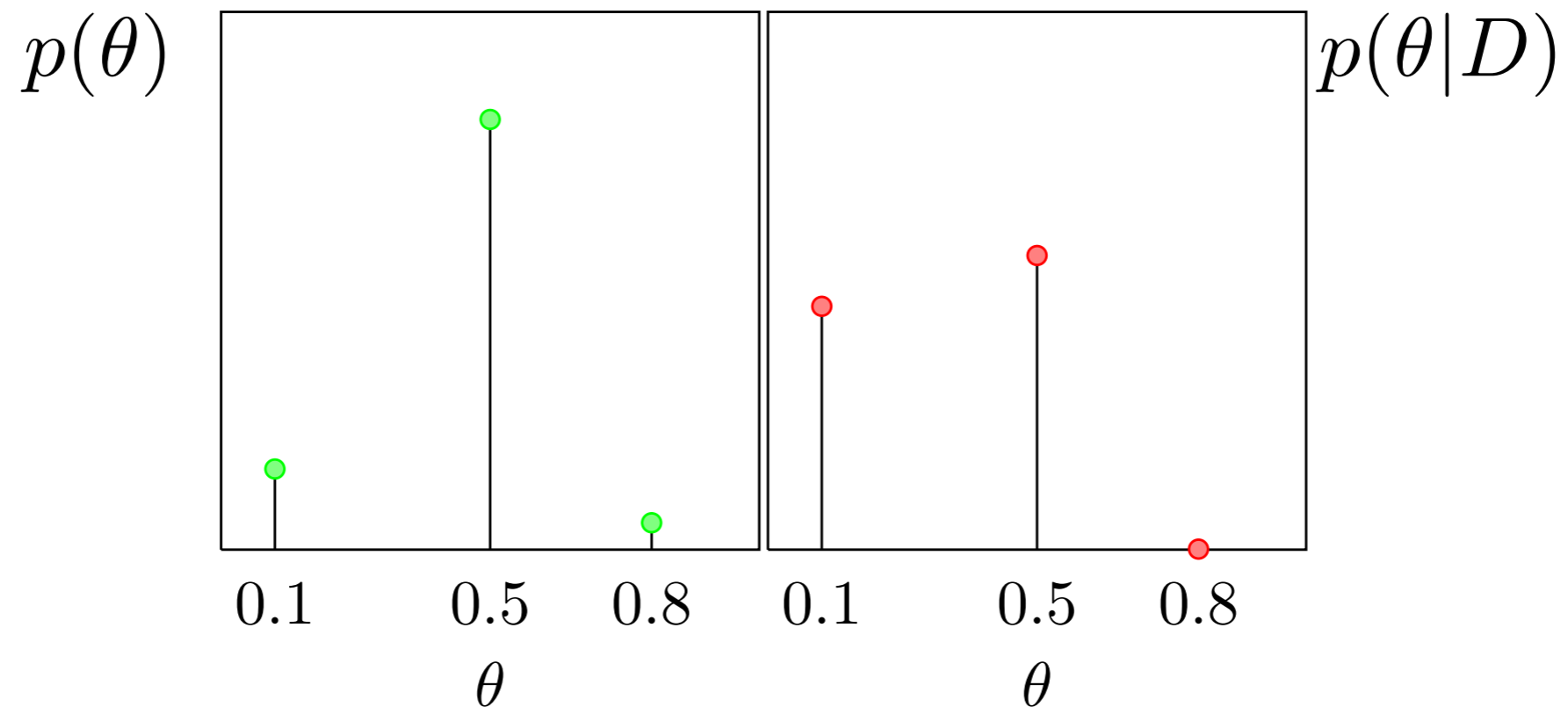
To avoid complexities resulting from continuous variables, we'll consider a discrete θ with only three possible states, $\theta \in \{0.1, 0.5, 0.8\}$. Specifically, we assume

$$p(\theta = 0.1) = 0.15, \quad p(\theta = 0.5) = 0.8, \quad p(\theta = 0.8) = 0.05$$



Example: Posterior after data

For an experiment with $N_H = 2$, $N_T = 8$, the posterior distribution is



But how do we get this posterior?

MAP calculation

- Start by applying Bayes rule

$$p(f|\mathcal{D}) = \frac{p(\mathcal{D}|f)p(f)}{p(\mathcal{D})}$$

- $p(\mathcal{D} | f)$ is the **likelihood** of the data, under the model
- $p(f)$ is the **prior** of the model
- $p(\mathcal{D})$ is the marginal distribution of the data
 - we will often be able to ignore this term

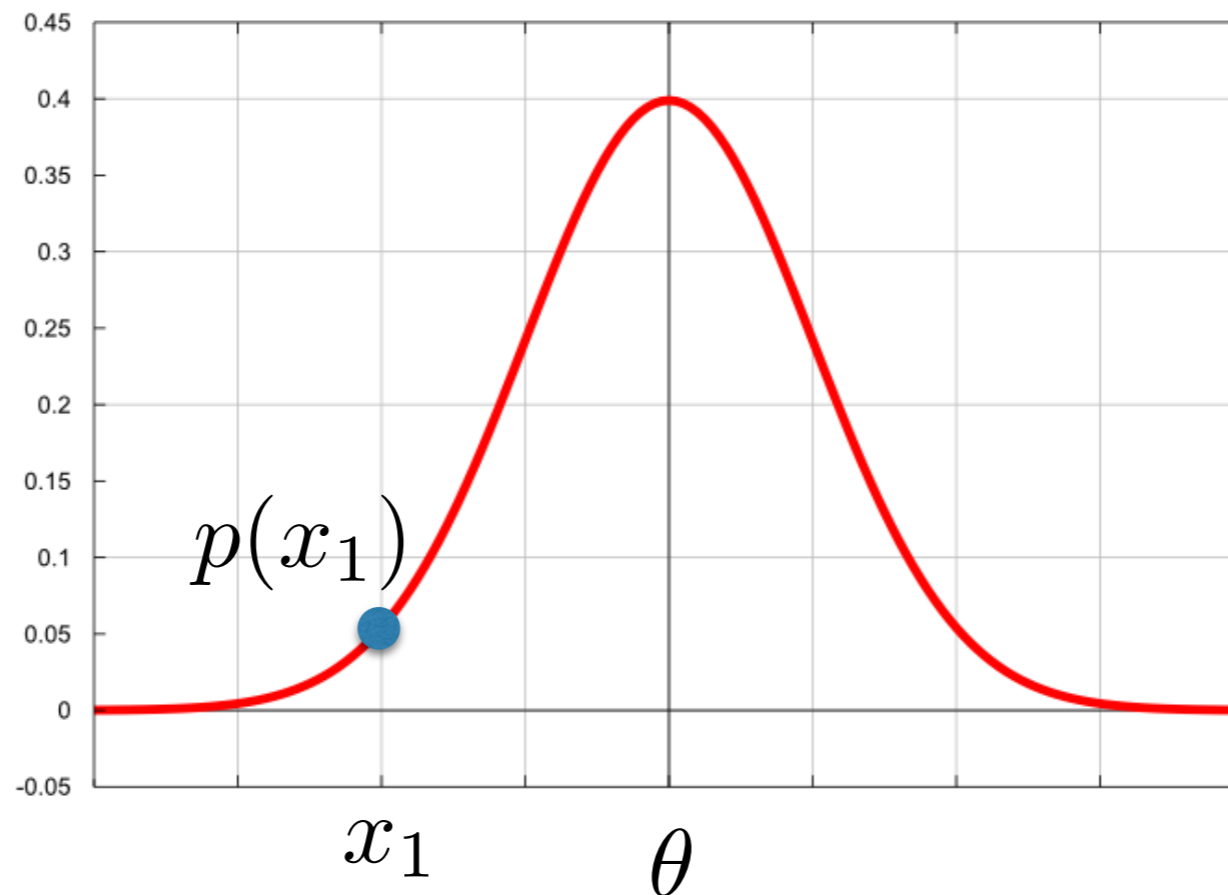
Why is this conversion important?

- Do not always have a known form for $p(f \mid D)$
- We usually have chosen (known) forms for $p(D \mid f)$ and $p(f)$
- **Example:** Let $D = \{x_1\}$ (one sample). Then, if model class is the set of Gaussians: $p(D \mid f) = p(x_1 \mid \mu, \sigma)$
 - $p(f \mid D)$ is not obvious, since specified our model class for $p(x \mid f)$
 - What is $p(f)$ in this case? We may put some prior “preferences” on μ and σ , e.g., normal distribution around μ , specifying that really large magnitude values in μ are unlikely
 - Specifying and using $p(f)$ is related to regularization and Bayesian parameter estimation, which will will discuss more later

Example of $p(D | \theta)$

$$= \mathcal{N}(x_1 | \mu = \theta, \sigma^2 = 1)$$
$$= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x_1 - \theta)^2\right)$$

Probabilities for fixed theta



September 12, 2019

- Thought questions and Assignment coming up
- Clarified biased coin example on eClass
- Any questions?

Recap

$$f_{\text{MAP}} = \arg \max_{f \in \mathcal{F}} p(f | \mathcal{D})$$

- We talked about finding the MAP solution
- This means we are thinking of the parameters as a random variable
- First step: understand $p(\mathcal{D} | \theta)$
- Second step: understand the role of the prior

What is the prior?

- The prior gives a chance to incorporate expert knowledge
 - The knowledge or information **prior** to seeing any data
- Example: maybe ahead of time you know what you think are reasonable values for the heights of people
- The prior could be a Gaussian distribution around $\theta =$ height, where $p(\theta)$ is a Gaussian with mean = 170 (cm) and variance = 60
 - What does the variance mean here, for the prior?

$p(D \mid \theta)$ with n samples

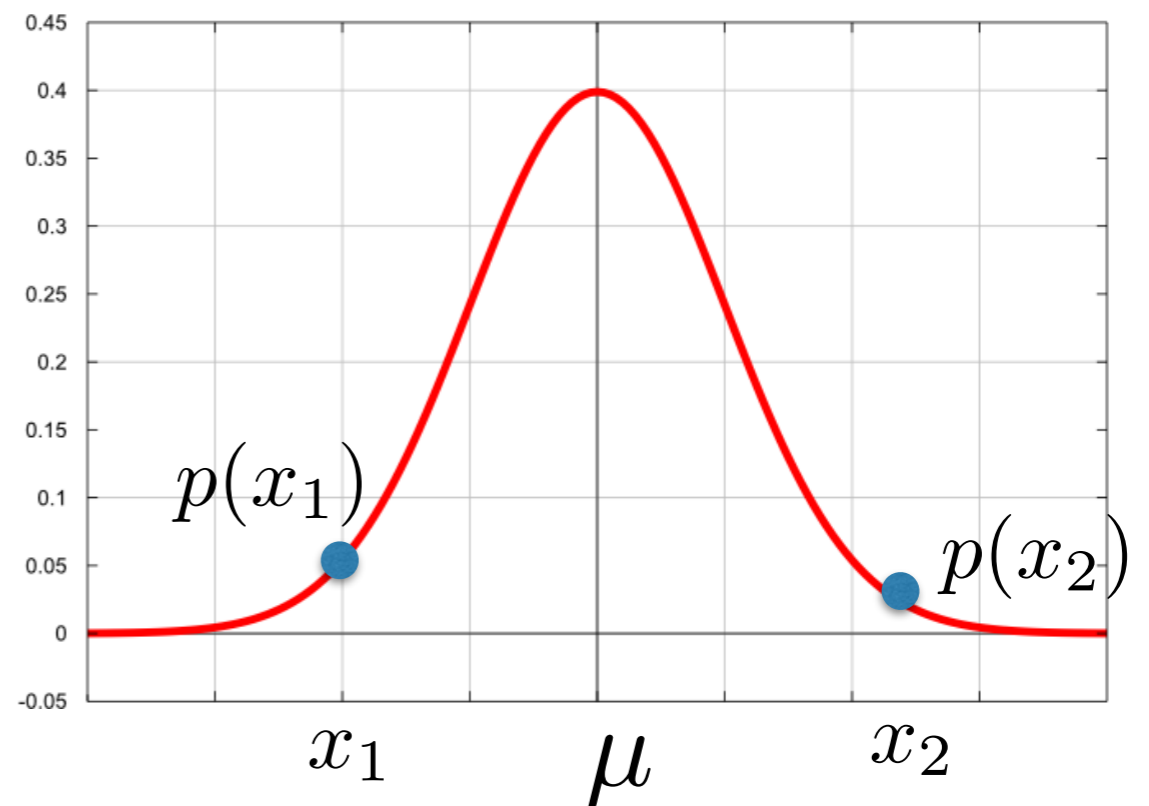
- Imagine we get iid samples x_1, \dots, x_n
- We could choose a Gaussian distribution for $P(x_i \mid \theta)$, with $\theta = \{\mu, \sigma\}$
- Recall iid = independent and identically distributed

What does it mean to say X_1 and X_2 are independent?

- Independent samples X_1, \dots, X_n
- e.g., $P(X_1, X_2 \mid \theta) = P(X_1 \mid \theta) P(X_2 \mid \theta)$
- X_1 and X_2 are random variables, from sampling the distribution twice
- Because of randomness, X_1 could have been many things
 - it represents the random variable that is the first sample

With n samples

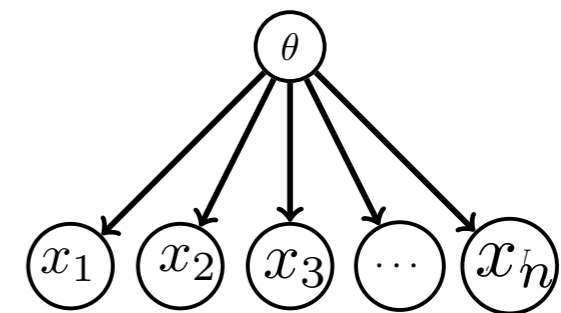
- For iid samples x_1, \dots, x_n , we could choose a Gaussian distribution for $P(x_i | \theta)$, with $\theta = \{\mu, \sigma\}$
- $P(x_1, \dots, x_n | \theta) = P(x_1 | \theta) \dots P(x_n | \theta)$
- Example: $D = \{x_1, x_2\}$



$$\begin{aligned} p(\mathcal{D} | \theta = (\mu, \sigma^2)) &= p(\{x_1, x_2\} | \theta = (\mu, \sigma^2)) \\ &= p(x_1 | \theta = (\mu, \sigma^2)) p(x_2 | \theta = (\mu, \sigma^2)) \end{aligned}$$

Example: MAP for discrete distributions

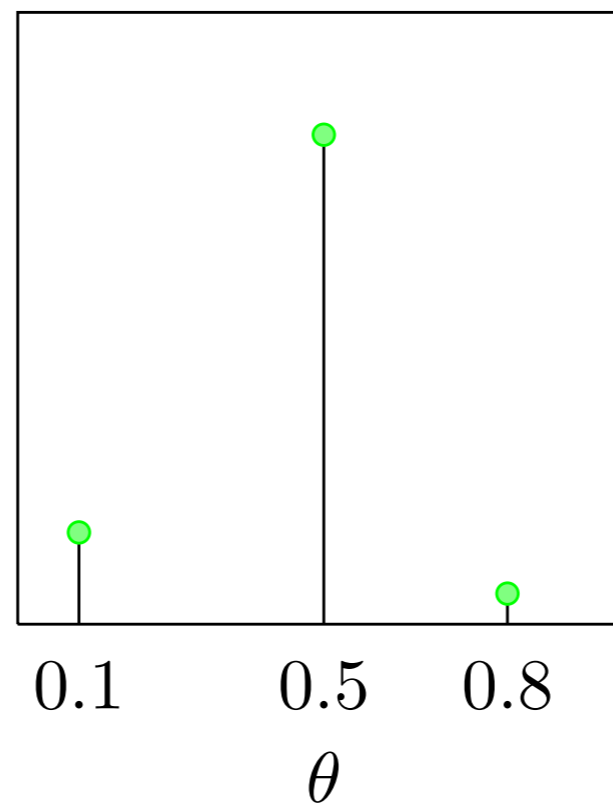
- Imagine you are flipping a biased coin; the model parameter is the bias of the coin, θ
- You get a dataset $D = \{x_1, \dots, x_n\}$ of coin flips, where $x_i = 1$ if it was heads, and $x_i = 0$ if it was tails
- How do we specify $p(\theta)$?
- What is the MAP estimate?



Example: MAP for discrete distributions

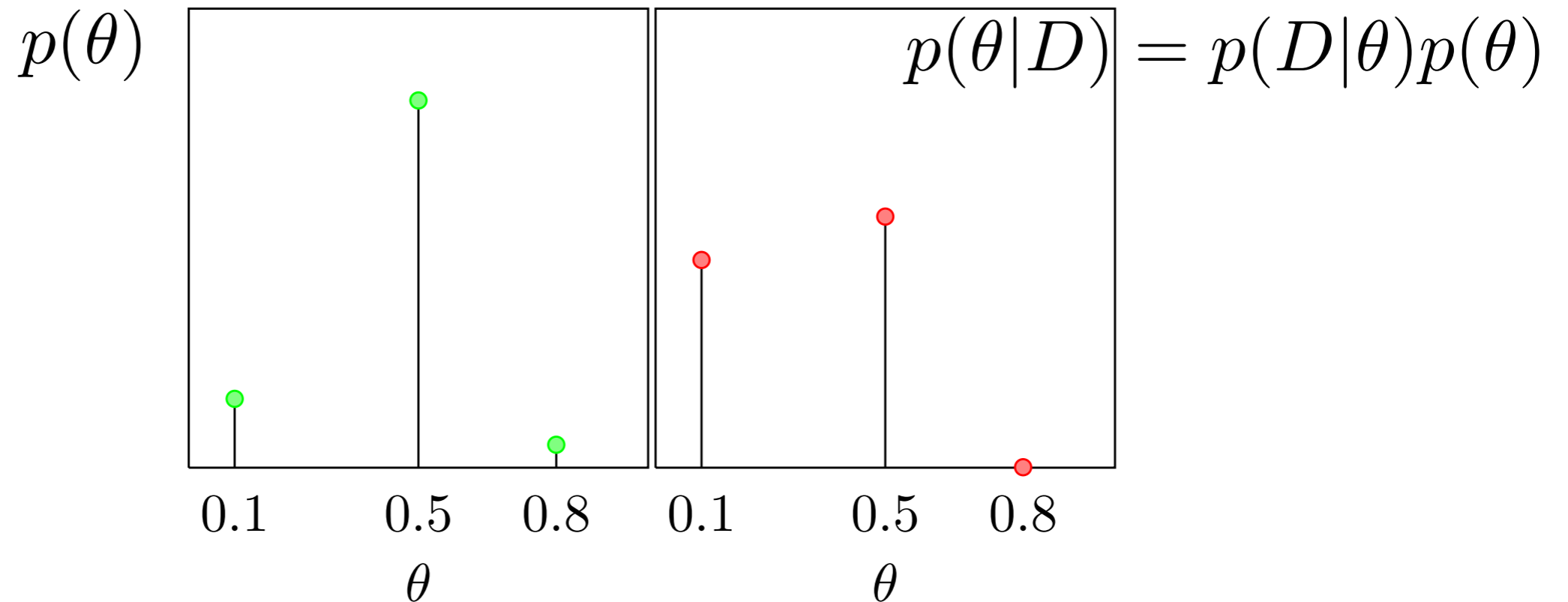
We still need to fully specify the prior $p(\theta)$. To avoid complexities resulting from continuous variables, we'll consider a discrete θ with only three possible states, $\theta \in \{0.1, 0.5, 0.8\}$. Specifically, we assume

$$p(\theta = 0.1) = 0.15, \quad p(\theta = 0.5) = 0.8, \quad p(\theta = 0.8) = 0.05$$



Example: MAP for discrete distributions

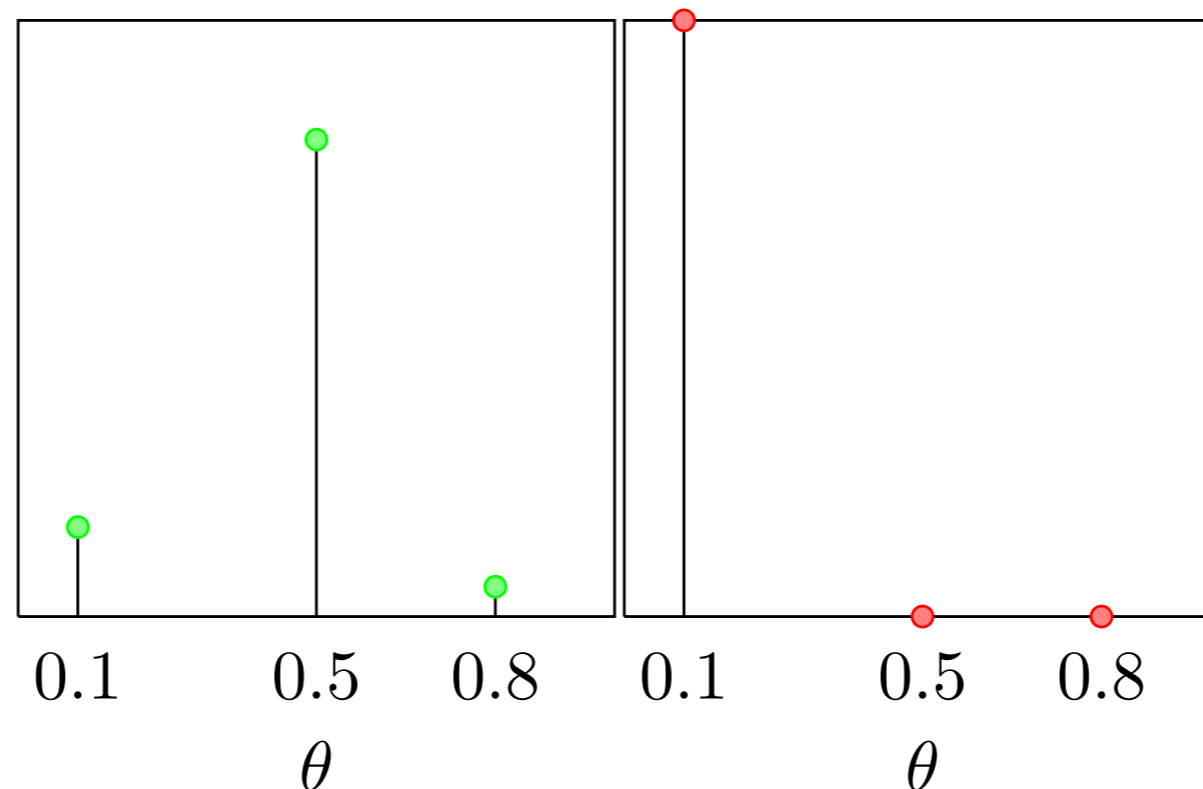
For an experiment with $N_H = 2$, $N_T = 8$, the posterior distribution is



If we were asked to choose a single *a posteriori* most likely value for θ , it would be $\theta = 0.5$, although our confidence in this is low since the posterior belief that $\theta = 0.1$ is also appreciable. This result is intuitive since, even though we observed more Tails than Heads, our prior belief was that it was more likely the coin is fair.

Example: MAP for discrete distributions

Repeating the above with $N_H = 20$, $N_T = 80$, the posterior changes to



so that the posterior belief in $\theta = 0.1$ dominates. There are so many more tails than heads that this is unlikely to occur from a fair coin. Even though we *a priori* thought that the coin was fair, *a posteriori* we have enough evidence to change our minds.

Reformulating MAP

The constant is irrelevant, as it is the same for all \mathcal{D}

$$\begin{aligned}\theta_{\text{MAP}} &= \arg \max_{\theta} \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})} \\ &= \arg \max_{\theta} p(\mathcal{D}|\theta)p(\theta)\end{aligned}$$

Will often write:

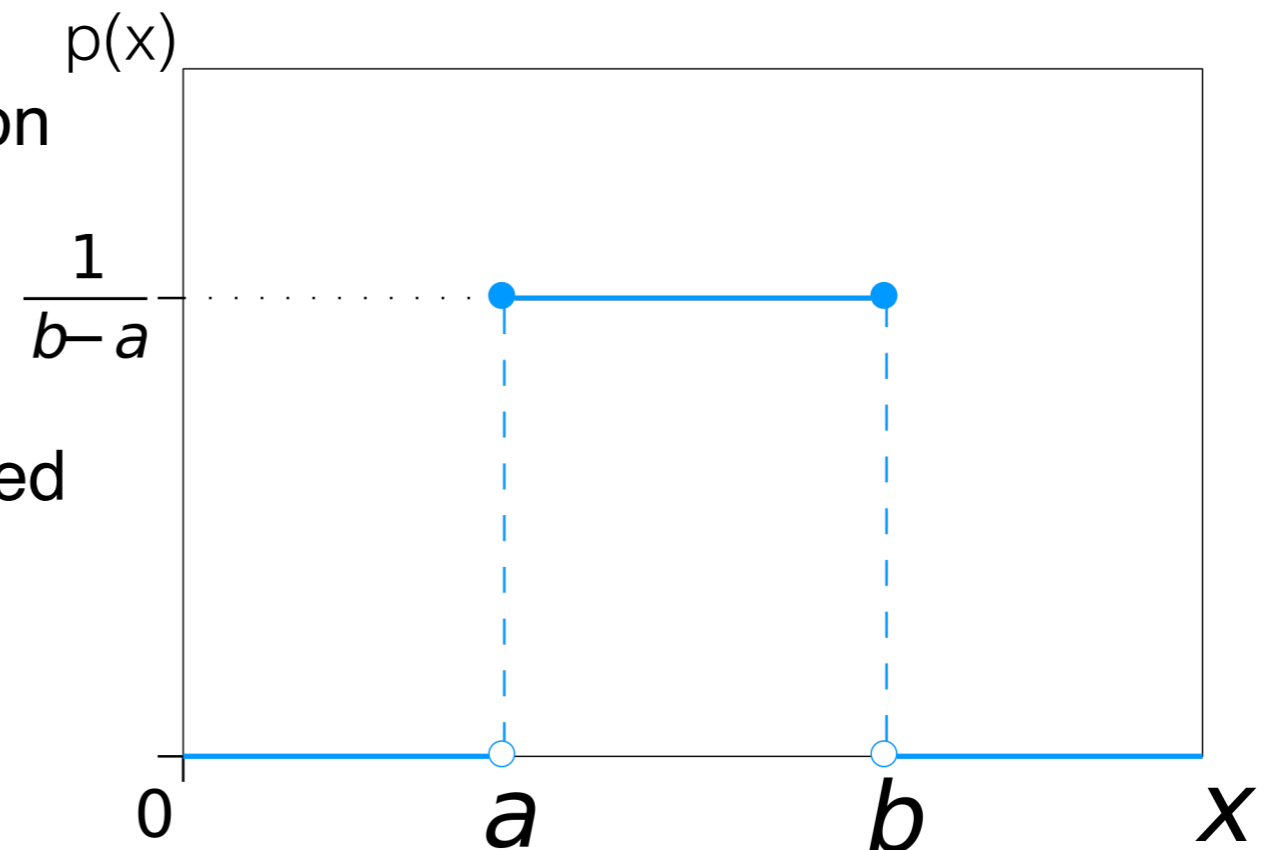
$$\begin{aligned}p(\theta|\mathcal{D}) &= \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})} \\ &\propto p(\mathcal{D}|\theta)p(\theta)\end{aligned}$$

Maximum likelihood

$$\theta_{\text{ML}} = \arg \max_{\theta} p(\mathcal{D}|\theta)$$

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})}$$
$$\propto p(\mathcal{D}|\theta)p(\theta)$$

- In some situations, may not have a reason to prefer one model over another (i.e., no prior knowledge or preferences)
- Can loosely think of maximum likelihood as instance of MAP, with uniform prior $p(\theta) = u$ for some constant u
 - If domain is infinite (example, the set of reals), the uniform distribution is not defined!
 - but the interpretation is still similar
 - in practice, typically have a bounded space in mind for the model class

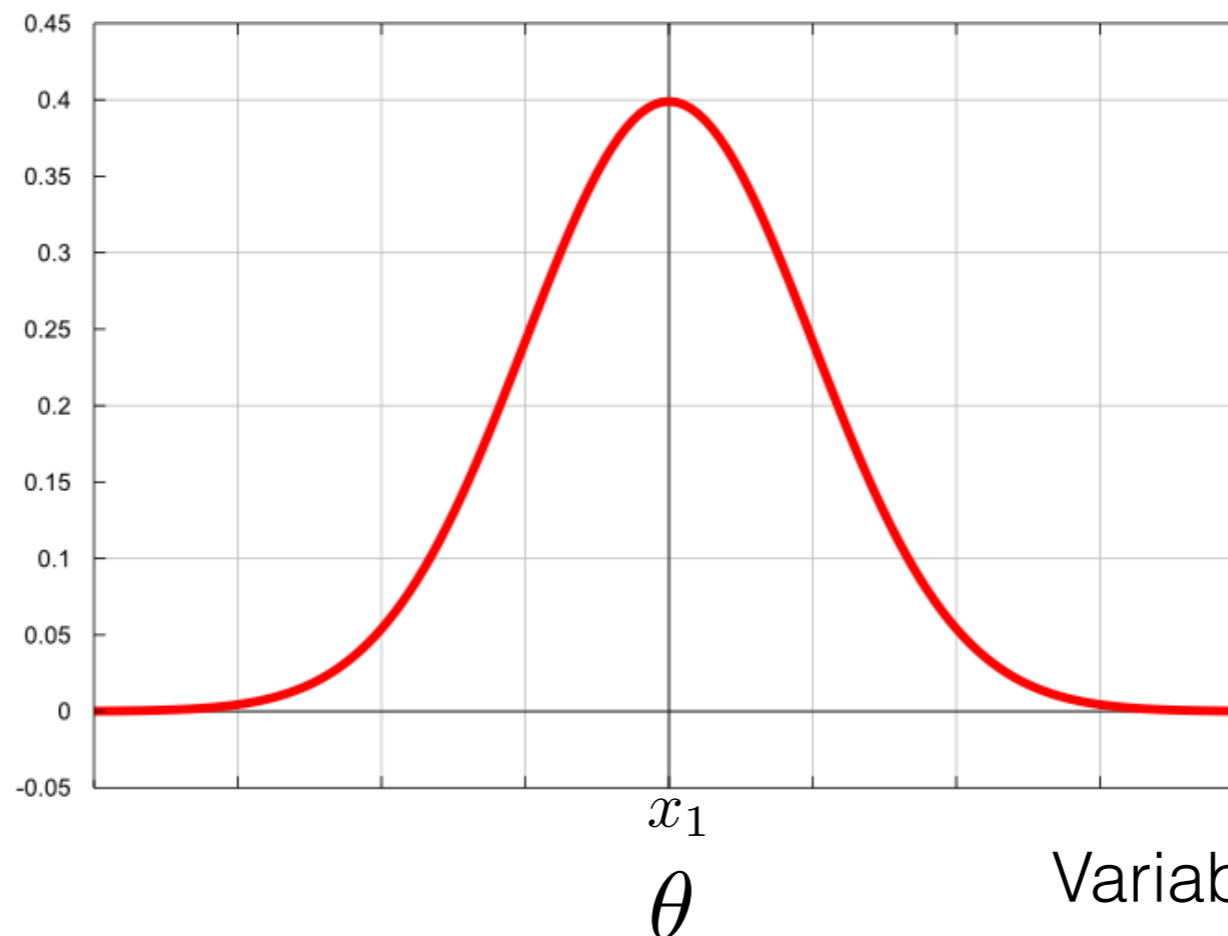


ML example

e.g., $\mathcal{F} = \mathbb{R}$, θ is the mean of a Gaussian, fixed $\sigma = 1$

$$\begin{aligned}c(\theta) &= p(\mathcal{D}|\theta) \\ &= \mathcal{N}(x_1 | \mu = \theta, \sigma^2 = 1) \\ &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x_1 - \theta)^2\right)\end{aligned}$$

$$c(\theta) = p(\mathcal{D}|\theta)$$



Maximizing the log-likelihood

- We want to maximize the **likelihood**, but often instead maximize the **log-likelihood**

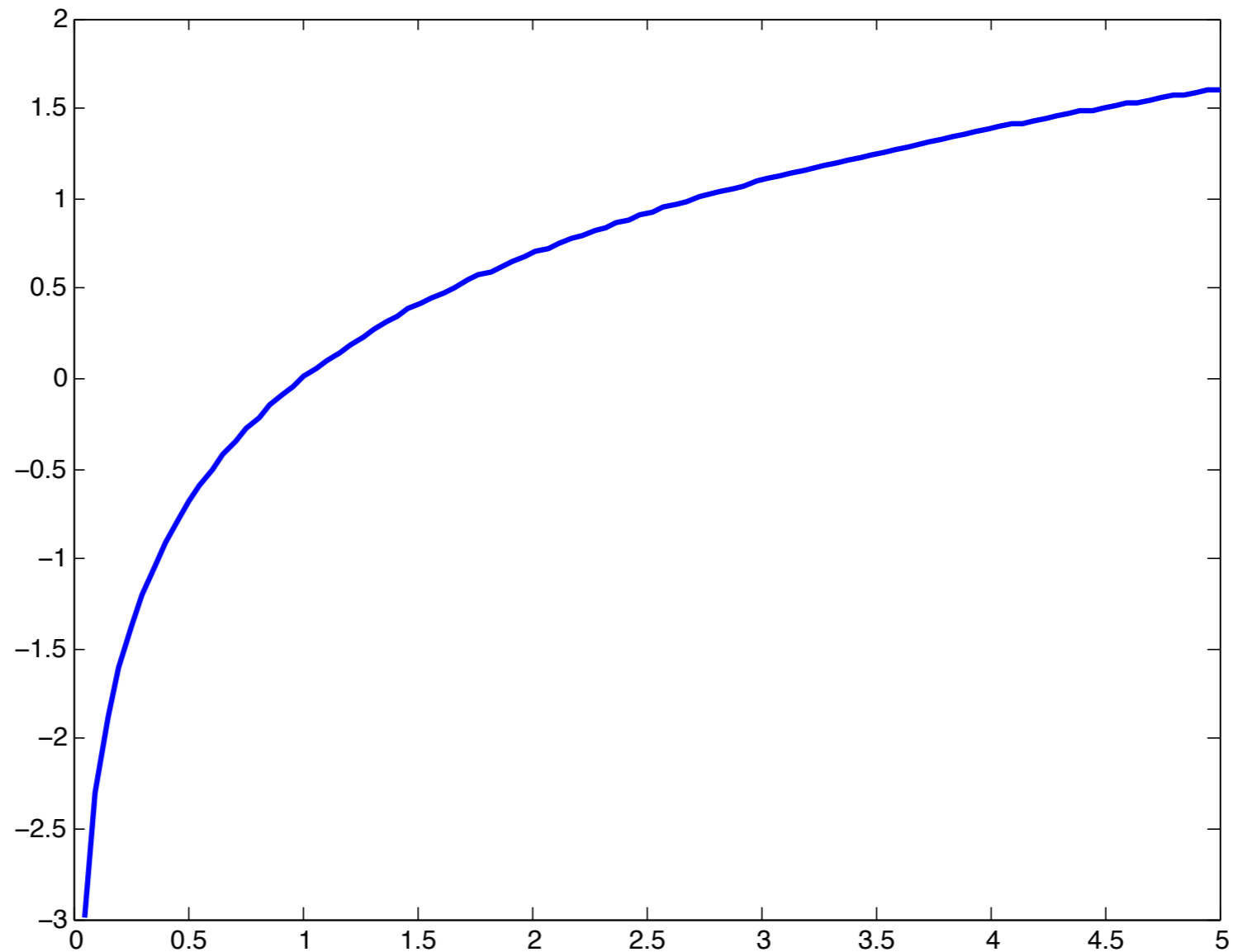
$$\arg \max_{\theta \in \mathcal{F}} p(\mathcal{D}|\theta) = \arg \max_{\theta \in \mathcal{F}} \log p(\mathcal{D}|\theta)$$

- Why? Or maybe first, is this equivalent?
 - The Why is that it makes the optimization much simpler, when we have more than one sample

Why can we shift by log?

$$c(\theta_1) > c(\theta_2) \iff \log c(\theta_1) > \log c(\theta_2) \quad \text{for } c(\theta) > 0$$

Monotone
increasing



Likelihood values always > 0

Maximizing the log-likelihood

e.g., $\mathcal{F} = \mathbb{R}$, θ is the mean of a Gaussian, fixed $\sigma = 1$

$$\log(ab) = \log a + \log b$$

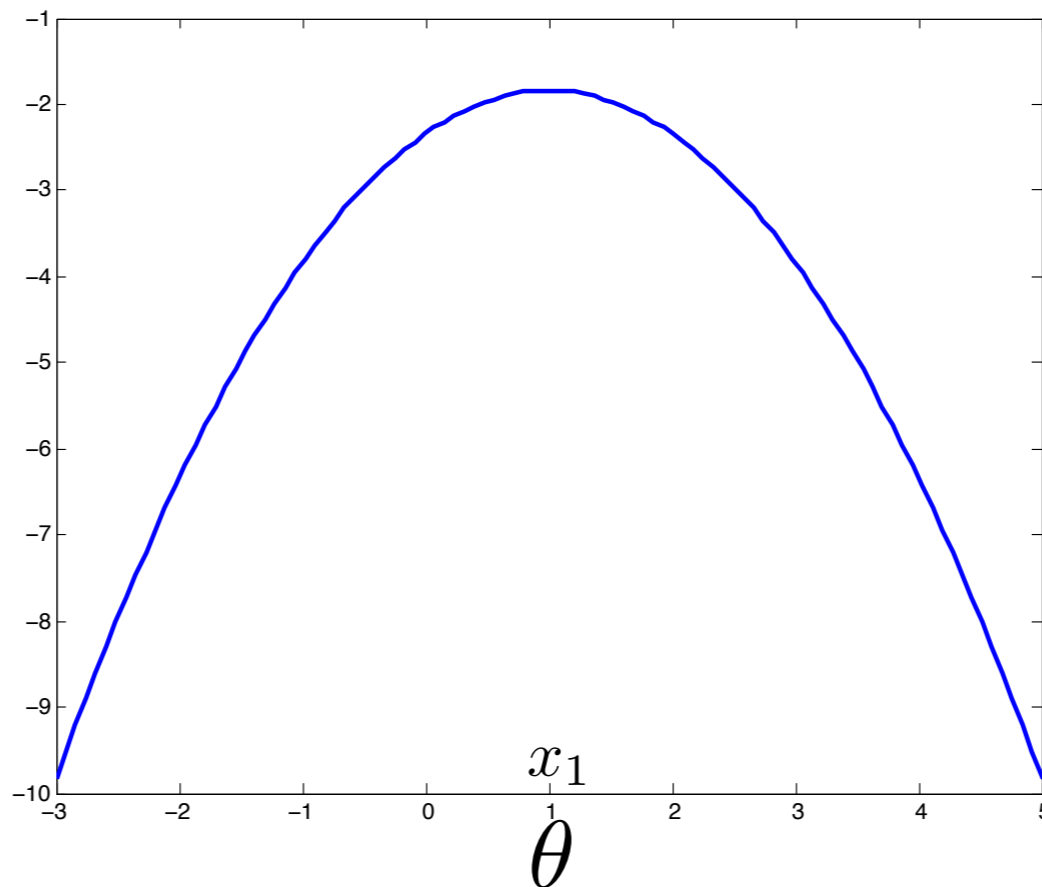
$$\log(a^c) = c \log a$$

$$c(\theta) = \log p(\mathcal{D}|\theta)$$

$$= \log \left(\frac{1}{2\pi} \exp \left(-\frac{1}{2} (x_1 - \theta)^2 \right) \right)$$

$$= -\log(2\pi) - \frac{1}{2} (x_1 - \theta)^2$$

$$c(\theta) = \log p(\mathcal{D}|\theta)$$



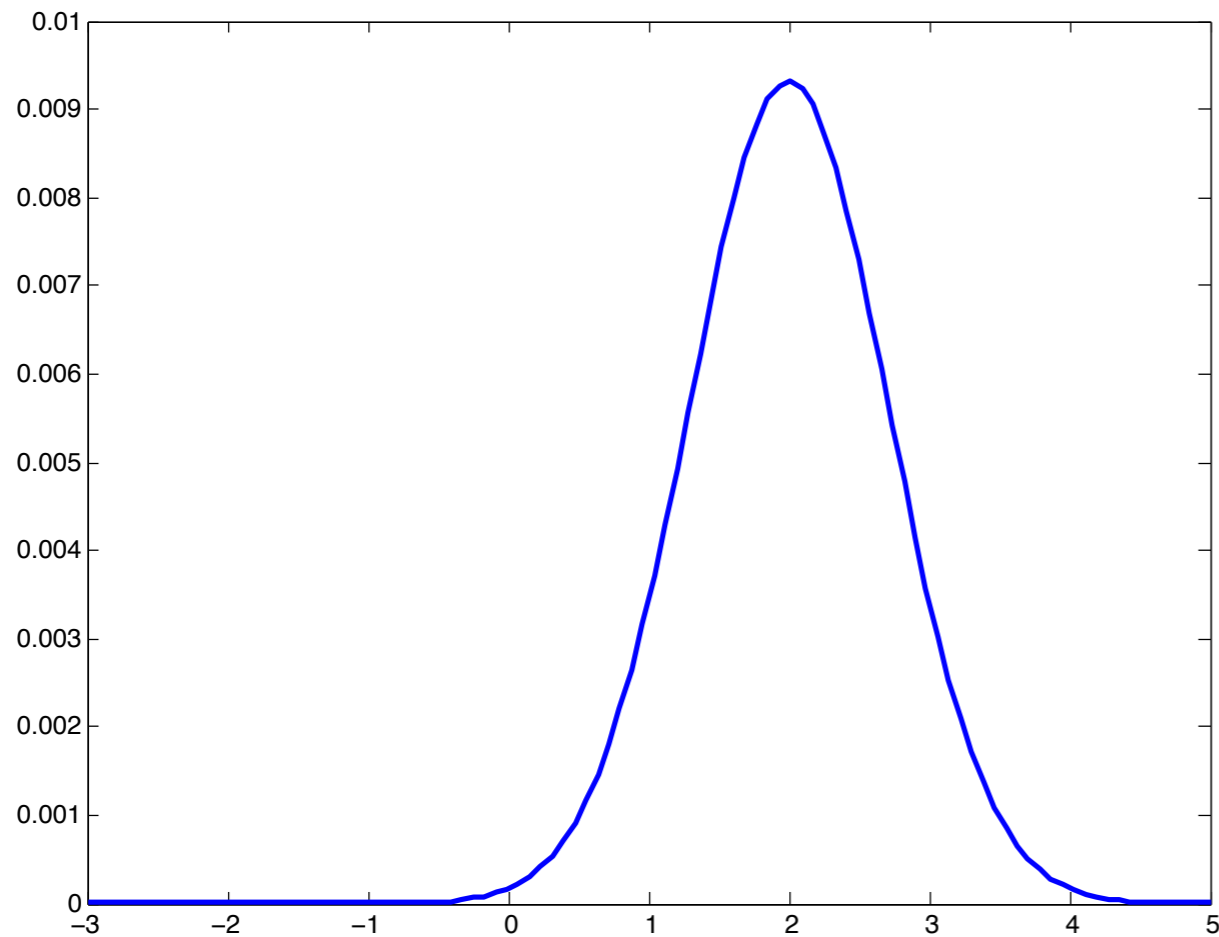
This conversion is even more important when we have more than one sample

- **Example:** Let $D = \{x_1, x_2\}$ (two samples).
- If x_1 and x_2 are independent samples from same distribution (same model), then $P(x_1, x_2 | \theta) = P(x_1 | \theta) P(x_2 | \theta)$

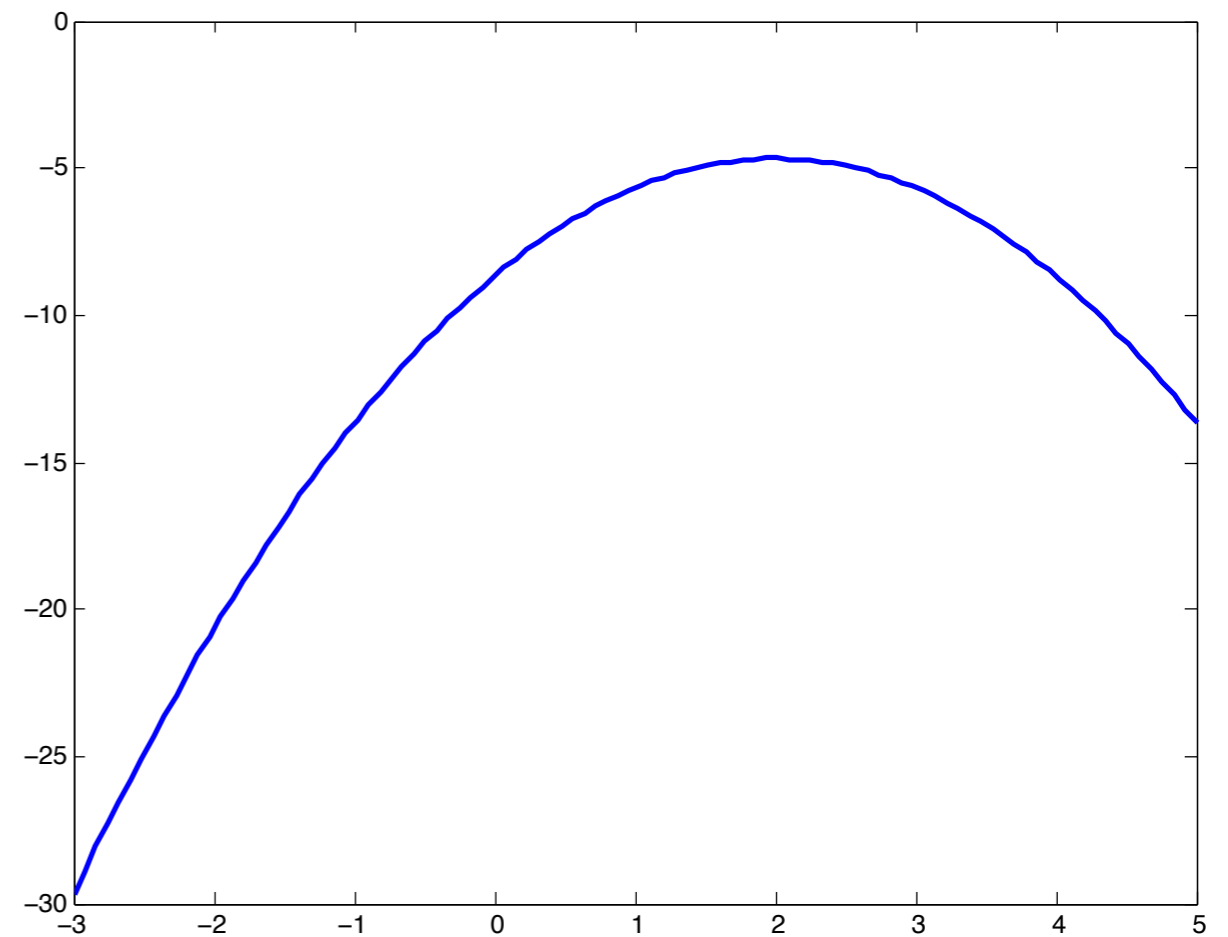
$$p(x_1 | \theta) p(x_2 | \theta) = \frac{1}{2\pi} \exp\left(-\frac{1}{2}(x_1 - \theta)^2\right) \times \frac{1}{2\pi} \exp\left(-\frac{1}{2}(x_2 - \theta)^2\right)$$

$$\begin{aligned} \log(p(x_1 | \theta) p(x_2 | \theta)) &= \log p(x_1 | \theta) + \log p(x_2 | \theta) \\ &= -2 \log(2\pi) - \frac{1}{2}(x_1 - \theta)^2 - \frac{1}{2}(x_2 - \theta)^2 \end{aligned}$$

Visualizing the objective



$$p(x_1|\theta)p(x_2|\theta)$$



$$\log(p(x_1|\theta)p(x_2|\theta))$$

Exercise: about the marginal over the data $p(D)$

- Why do we avoid computing $p(D)$?
- How do we compute $p(D)$? We have the tools, since we know about marginals

How do we compute $p(D)$?

- If we have $p(D, f)$, can we obtain $p(D)$?
 - Marginalization

$$p_{X_i}(x_i) = \sum_{x_1} \cdots \sum_{x_{i-1}} \sum_{x_{i+1}} \cdots \sum_{x_k} p_{\mathbf{X}}(x_1, \dots, x_k)$$

$$p_{X_i}(x_i) = \int_{x_1} \cdots \int_{x_{i-1}} \int_{x_{i+1}} \cdots \int_{x_k} p_{\mathbf{X}}(\mathbf{x}) dx_1 \cdots dx_{i-1} dx_{i+1} \cdots dx_k$$

- If we have $p(D|f)$ and $p(f)$, do we have $p(D, f)$?

Data marginal

- Using the formula of total probability

$$p(\mathcal{D}) = \begin{cases} \sum_{f \in \mathcal{F}} p(\mathcal{D}|f)p(f) & f : \text{discrete} \\ \int_{\mathcal{F}} p(\mathcal{D}|f)p(f)df & f : \text{continuous} \end{cases}$$

- Fully expressible in terms of likelihood and prior

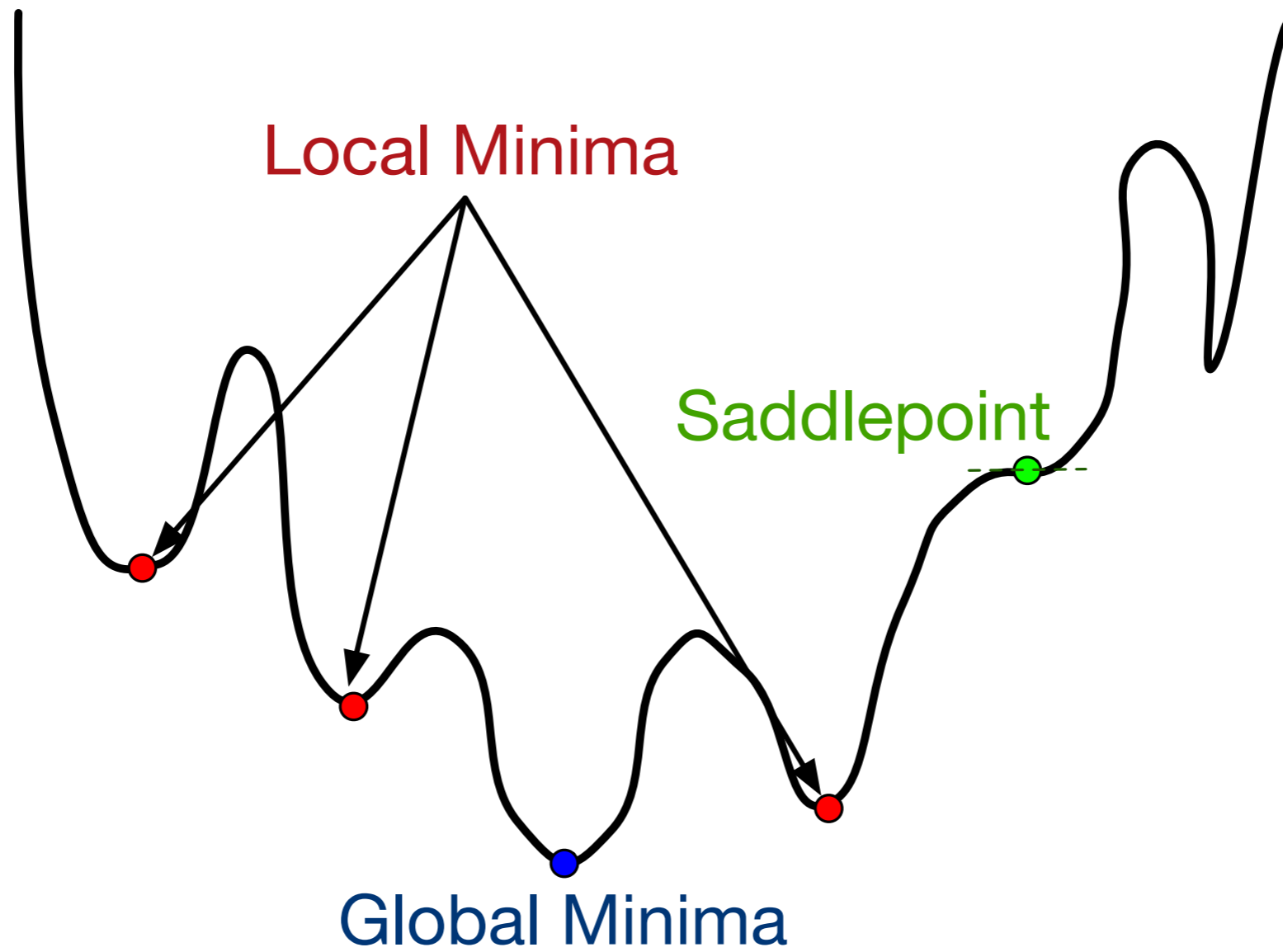
How do we solve the MAP and ML maximization problems?

- Naive strategy:
 - 1. Guess 100 solutions θ
 - 2. Pick the one with the largest value
- Can we do something better?

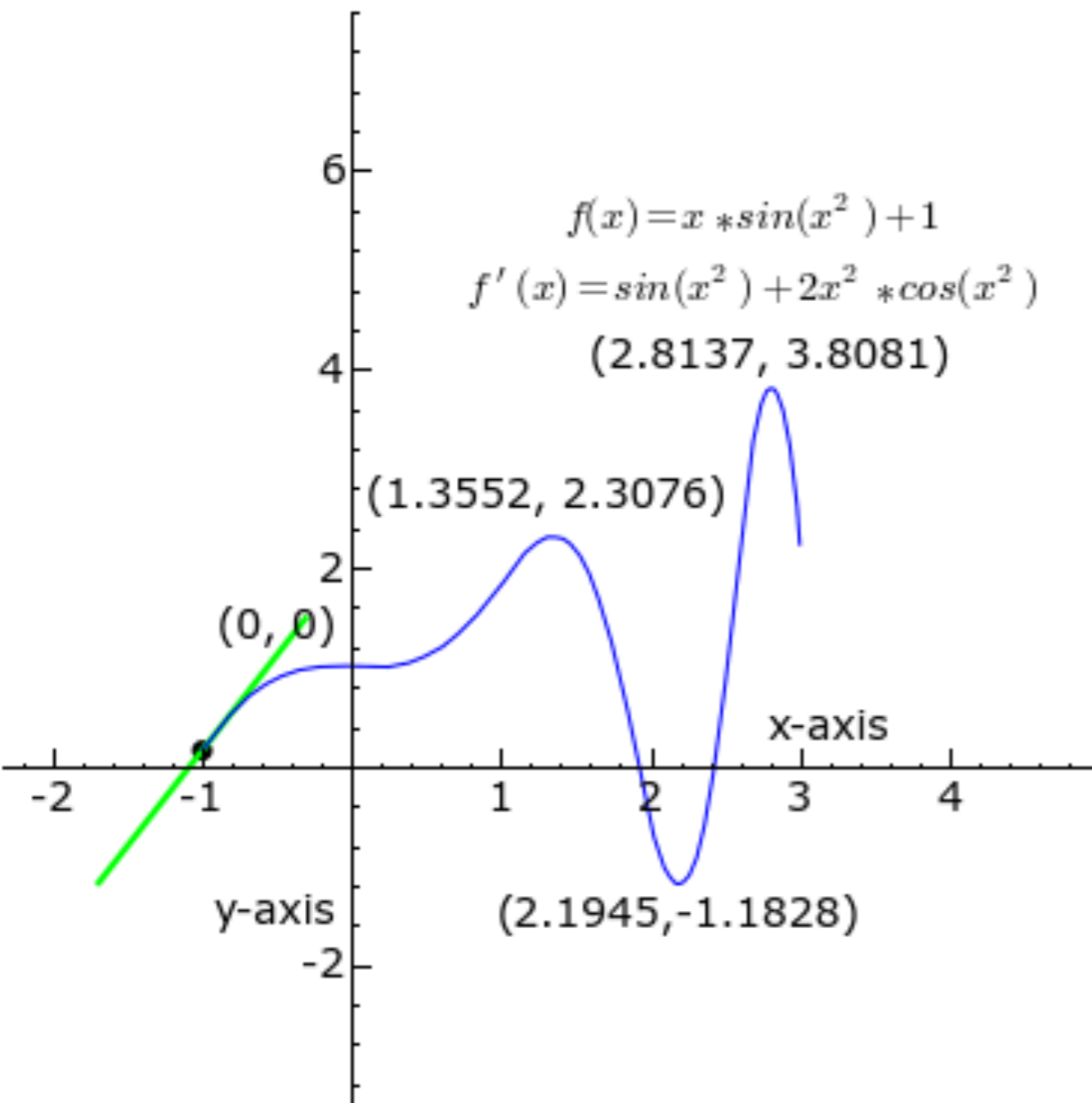
Crash course in optimization

- Goal: find maximal (or minimal) points of functions
- Generally assume functions are smooth, use gradient descent
- Derivative: direction of ascent from a scalar point $\frac{d}{dx} c(x)$
- Gradient: direction of ascent from a vector point $\nabla c(\mathbf{x})$

Function surface



Single-variate calculus



GIF from Wikipedia: Tangent

For a function f defined on a scalar x , the derivative is

$$\frac{df}{dx}(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

At any point, x , $\frac{df}{dx}(x)$ gives the slope of the tangent to the function at $f(x)$

Why don't constants matter?

$$\max_x c(x)$$

$$\frac{d}{dx} c(x) = 0$$

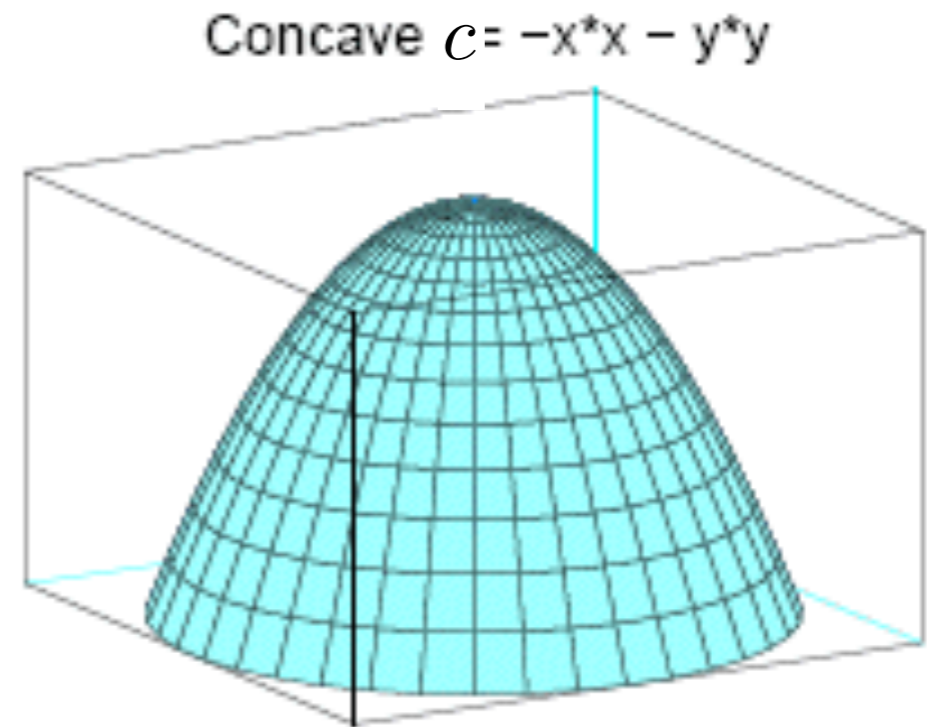
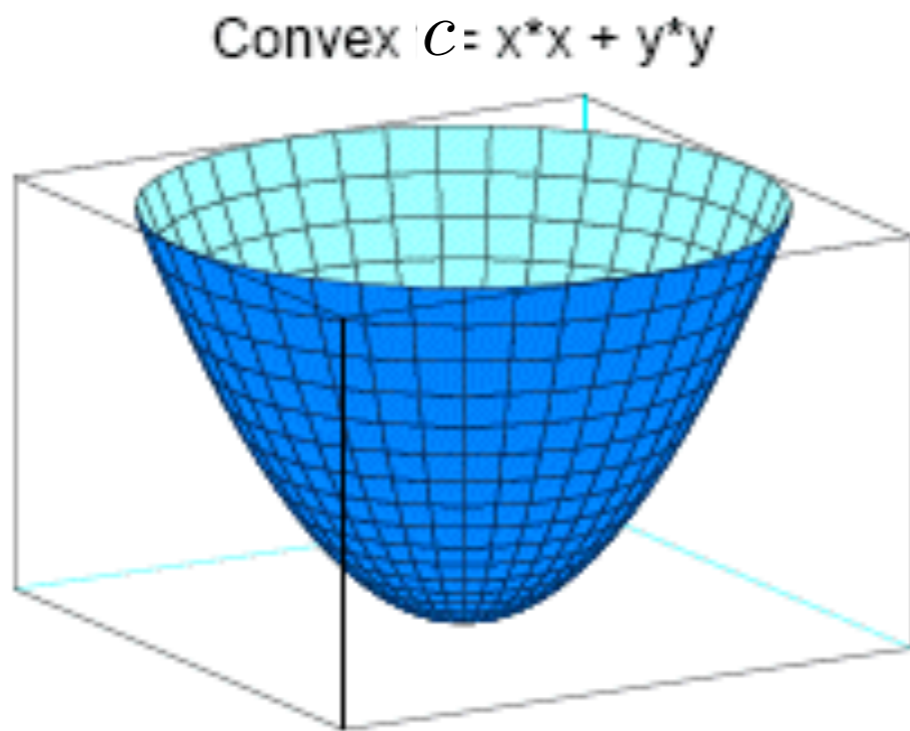
$$\max_x u c(x), \quad u > 0$$

$$\frac{d}{dx} u c(x) = u \frac{d}{dx} c(x) = 0$$

Both have derivative zero under same condition
regardless of $u > 0$

Can either minimize or maximize

$$\arg \min_{\theta} c(\theta) = \arg \max_{\theta} -c(\theta)$$



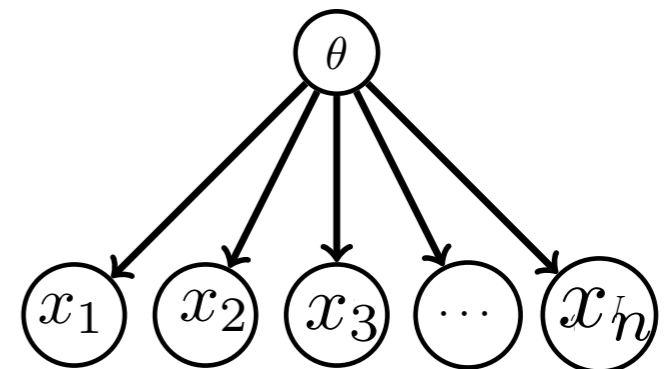
$$c(t\mathbf{w}_1 + (1-t)\mathbf{w}_2) \leq tc(\mathbf{w}_1) + (1-t)c(\mathbf{w}_2)$$

Example: maximum likelihood for discrete distributions

- Imagine you are flipping a biased coin; the model parameter is the bias of the coin, theta
- You get a dataset $D = \{x_1, \dots, x_n\}$ of coin flips, where $x_i = 1$ if it was heads, and $x_i = 0$ if it was tails
- What is $p(D | \theta)$?

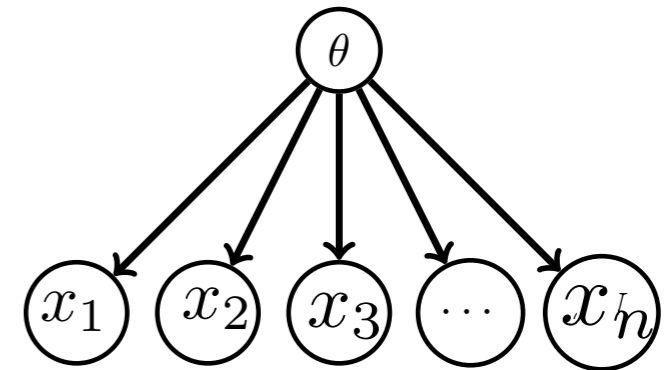
$$\begin{aligned} p(\mathcal{D}|\theta) &= p(x_1, \dots, x_n | \theta) \\ &= \prod_{i=1}^n p(x_i | \theta) \end{aligned}$$

$$p(x|\theta) = \theta^x (1 - \theta)^{1-x}$$



Example: maximum likelihood for discrete distributions

- How do we estimate theta?
- Counting:
 - count the number of heads N_h
 - count the number of tails N_t
 - normalize: $\theta = N_h / (N_h + N_t)$
- What if you actually try to maximize the likelihood?
 - i.e., solve $\operatorname{argmax} p(\mathcal{D} | \theta)$

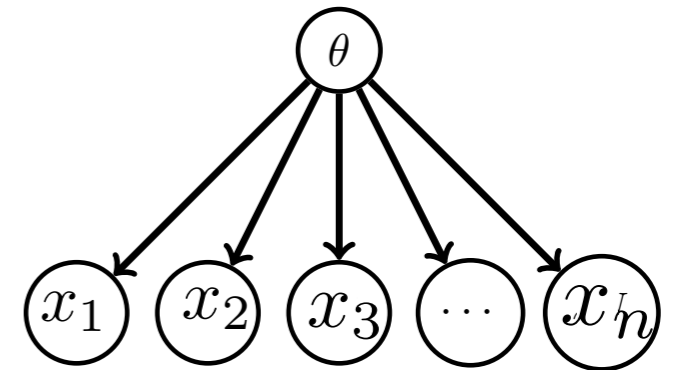


$$\begin{aligned} p(\mathcal{D} | \theta) &= p(x_1, \dots, x_n | \theta) \\ &= \prod_{i=1}^n p(x_i | \theta) \end{aligned}$$

$$p(x | \theta) = \theta^x (1 - \theta)^{1-x}$$

Example: maximum likelihood for discrete distributions

- What if you actually try to maximize the likelihood to get theta?
 - i.e., solve $\operatorname{argmax} p(\mathcal{D} | \theta)$



$$\max_{\theta} \prod_{i=1}^n p(x_i | \theta) = \max_{\theta} c(\theta)$$

$$c(\theta) = \prod_{i=1}^n p(x_i | \theta)$$

$$\operatorname{arg} \max_{\theta} c(\theta) = \operatorname{arg} \max_{\theta} \log c(\theta)$$

$$p(\mathcal{D} | \theta) = p(x_1, \dots, x_n | \theta)$$

$$= \prod_{i=1}^n p(x_i | \theta)$$

$$p(x | \theta) = \theta^x (1 - \theta)^{1-x}$$

Example: maximum likelihood for discrete distributions

$$\arg \max_{\theta} c(\theta) = \arg \max_{\theta} \log c(\theta)$$

$$\log(ab) = \log a + \log b$$

$$\log(a^c) = c \log a$$

$$\log c(\theta) = \log \prod_{i=1}^n p(x_i|\theta)$$

$$= \sum_{i=1}^n \log p(x_i|\theta)$$

$$p(x|\theta) = \theta^x (1 - \theta)^{1-x}$$

$$\begin{aligned} \log p(x|\theta) &= \log(\theta^x) + \log((1 - \theta)^{1-x}) \\ &= x \log(\theta) + (1 - x) \log(1 - \theta) \end{aligned}$$

Example: maximum likelihood for discrete distributions

$$\begin{aligned}\sum_{i=1}^n \log p(x_i|\theta) &= \sum_{i=1}^n x_i \log(\theta) + \sum_{i=1}^n (1 - x_i) \log(1 - \theta) \\ &= \log(\theta) \left(\sum_{i=1}^n x_i \right) + \log(1 - \theta) \left(\sum_{i=1}^n (1 - x_i) \right)\end{aligned}$$

$$\bar{x} = \sum_{i=1}^n x_i$$

$$\frac{d}{d\theta} = \frac{1}{\theta} \bar{x} - \frac{1}{1 - \theta} (n - \bar{x}) = 0$$

Example: maximum likelihood for discrete distributions

$$\bar{x} = \sum_{i=1}^n x_i$$

$$\frac{d}{d\theta} = \frac{1}{\theta} \bar{x} - \frac{1}{1-\theta} (n - \bar{x}) = 0$$

$$\implies \frac{\bar{x}}{\theta} = \frac{n - \bar{x}}{1 - \theta}$$

$$\implies (1 - \theta)\bar{x} = \theta(n - \bar{x})$$

$$\implies \bar{x} - \theta\bar{x} = \theta n - \theta\bar{x}$$

$$\implies \theta = \frac{\bar{x}}{n}$$

We solved for a stationary point

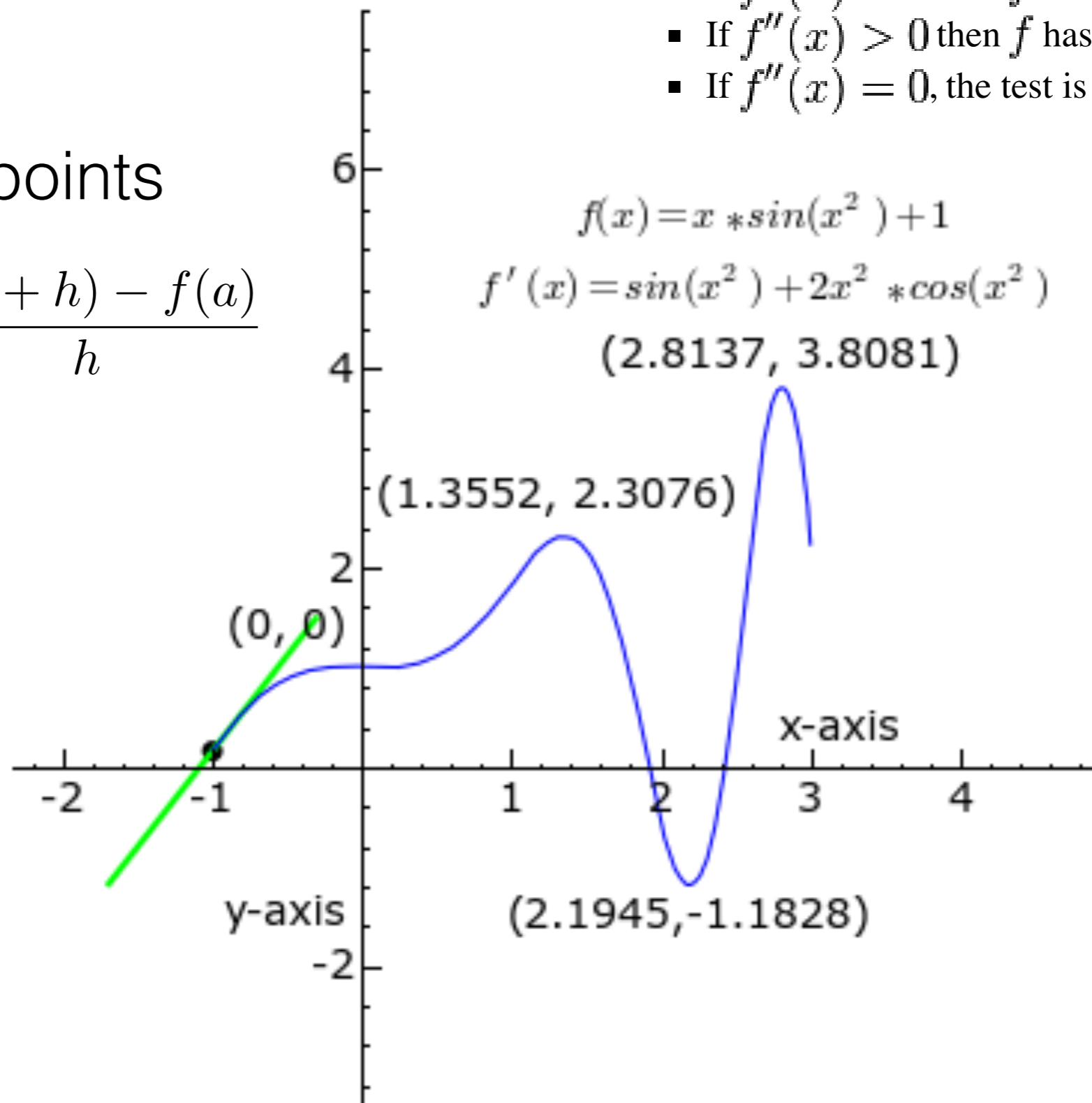
- But how do we know that this is the optimal solution to the optimization problem we specified?

Univariate optimization

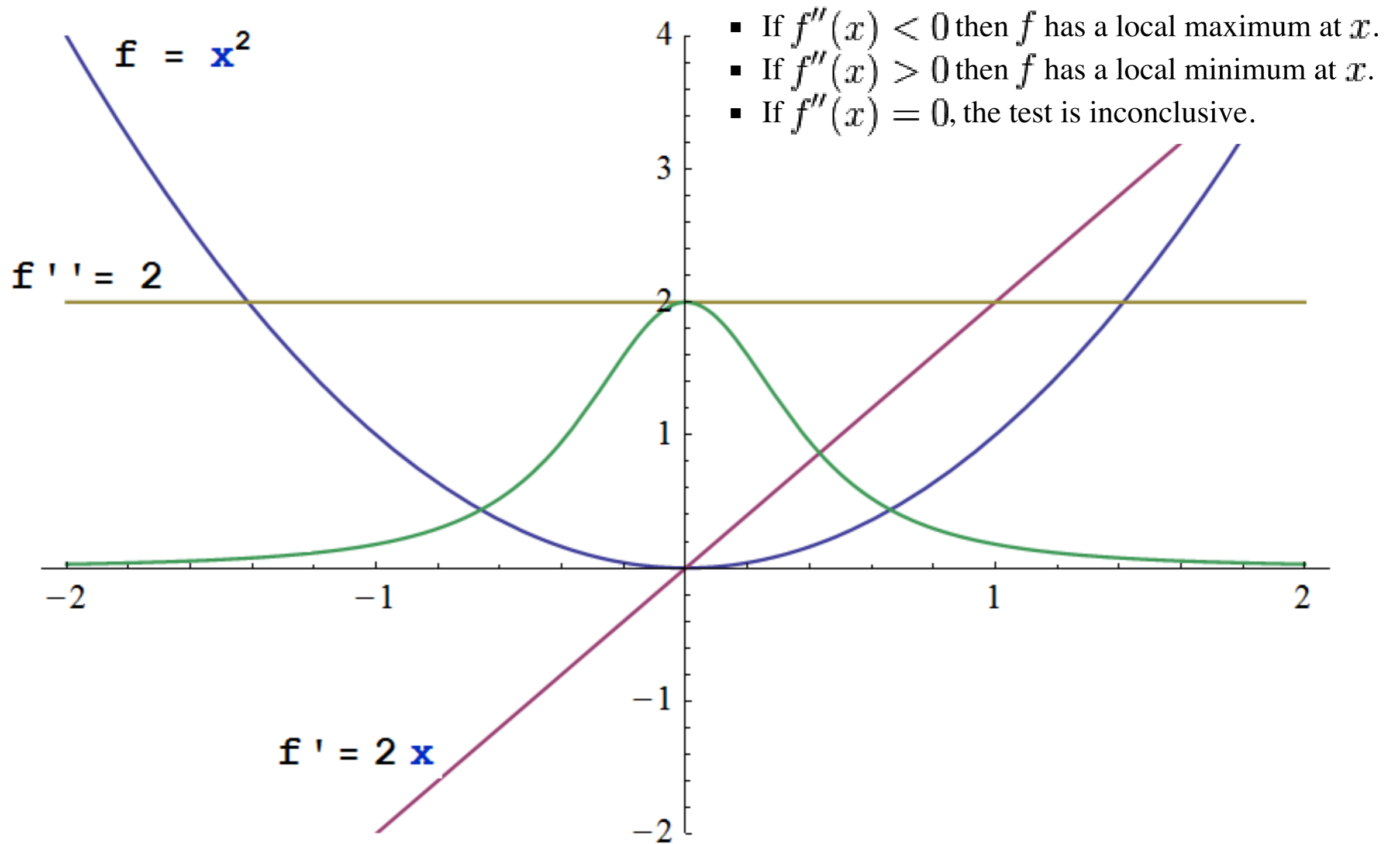
- Minima
- Maxima
- Saddle points

$$f'(a) = \lim_{h \rightarrow 0} \frac{f(a+h) - f(a)}{h}$$

- If $f''(x) < 0$ then f has a local maximum at x .
- If $f''(x) > 0$ then f has a local minimum at x .
- If $f''(x) = 0$, the test is inconclusive.



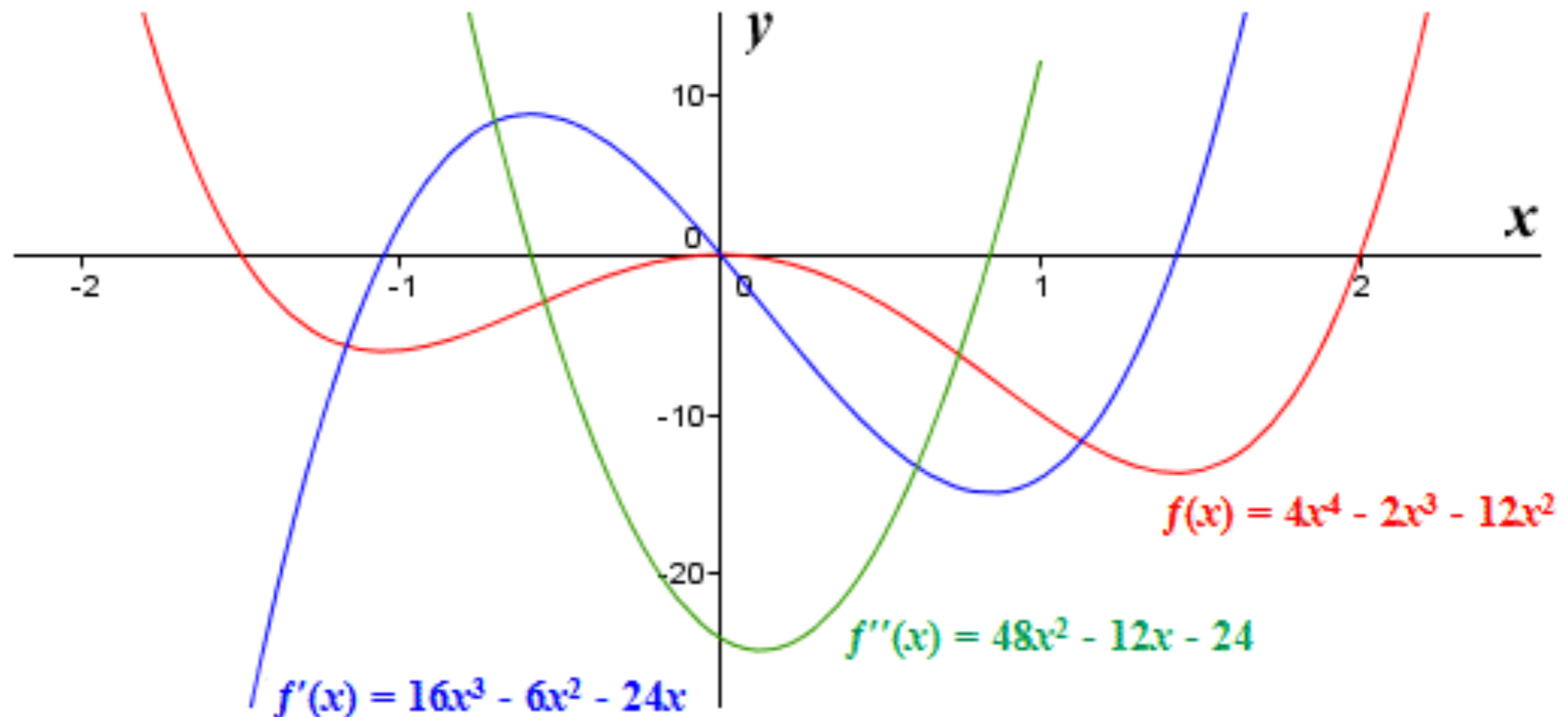
Second derivative test



Ignore the green line

Second derivative test

- If $f''(x) < 0$ then f has a local maximum at x .
- If $f''(x) > 0$ then f has a local minimum at x .
- If $f''(x) = 0$, the test is inconclusive.



Now on to some careful examples of MAP!

- Whiteboard for Examples 8, 9, 10, 11
- More fun with derivatives and finding the minimum of a function
- Next:
 - introduction to prediction problems for ML

