

# Point Estimation

CSE 446: Machine Learning  
Emily Fox  
University of Washington  
January 6, 2017

©2017 Emily Fox

Maximum likelihood estimation  
for a binomial distribution

## Your first consulting job

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



- You say: The probability is:  $\frac{3}{5}$
- He says: Why???
- You say: Because...

3

©2017 Emily Fox

CSF 446: Machine Learning

## Thumbtack – Binomial distribution

- $P(\text{Heads}) = \theta$ ,  $P(\text{Tails}) = 1 - \theta$  (param)
- Flips are i.i.d.:
  - Independent events
  - Identically distributed according to a binomial distribution

- Sequence D of  $\alpha_H$  heads (H) and  $\alpha_T$  tails (T)
- $P(D | \theta) = p(HH TT H | \theta) = \theta \cdot \theta \cdot (1 - \theta) \cdot (1 - \theta) \cdot \theta = \theta^3 (1 - \theta)^2$ 

cond. or given      generically... =  $\theta^{\alpha_H} (1 - \theta)^{\alpha_T}$       If I knew  $\theta$  I could compute the prob. of this seq.

4

©2017 Emily Fox

CSF 446: Machine Learning

## The learning task

- Want to learn a model of thumbtack flips from experience  $\leftarrow$  data *binomial dist.*
- Example 1: Maximum likelihood estimation**  
What value of  $\theta$  maximizes the **likelihood** of having seen the observed sequence (according to my model)?
- What is a **likelihood function**?

$$l(\theta) = p(D | \theta)$$

$\uparrow$  data       $\uparrow$  params      function of  $\theta$

Is this a prob.?  
Does it add to 1  
No

5

©2017 Emily Fox

CSF 446: Machine Learning

## Maximum likelihood estimation

- Data:** Observed set  $D$  of  $\alpha_H$  heads (H) and  $\alpha_T$  tails (T)  $\leftarrow$  "experience"
- Hypothesis:** Binomial distribution  $\leftarrow$  assumption of how the world works
- Learning  $\theta$  is an optimization problem
  - What's the objective function?

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

- MLE:** Choose  $\theta$  that maximizes the likelihood of observed data

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

*argument max(.) of fcn*

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

*$\arg \max_{\theta} f(\theta)$   
 $= \arg \max_{\theta} \ln f(x)$   
 $\ln$  is monotonic*

6

©2017 Emily Fox

CSF 446: Machine Learning

## Your first learning algorithm

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

$$= \arg \max_{\theta} \ln \theta^{\alpha_H} (1 - \theta)^{\alpha_T}$$

- Set derivative to zero:  $\frac{d}{d\theta} \ln P(D | \theta) = 0$

$$\frac{d}{d\theta} \ln \theta = \frac{1}{\theta}$$

$$\frac{d}{d\theta} \ln(1-\theta) = -\frac{1}{1-\theta}$$

$$\frac{d}{d\theta} \ln P(D | \theta) = \frac{d}{d\theta} [\alpha_H \ln \theta + \alpha_T \ln(1-\theta)]$$

$$= \frac{\alpha_H}{\theta} - \frac{\alpha_T}{(1-\theta)} = 0 \rightarrow \hat{\theta} = \frac{\alpha_H}{\alpha_H + \alpha_T} = \frac{3}{3+2} = \frac{3}{5}$$

7

©2017 Emily Fox

CSF 446: Machine Learning

## How many flips do I need?

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

- Billionaire says:** I flipped 3 heads and 2 tails.
- You say:**  $\hat{\theta} = 3/5$ , I can prove it!
- He says:** What if I flipped 30 heads and 20 tails?
- You say:** Same answer, I can prove it!  $\hat{\theta} = \frac{3}{5}$
- He says: What's better?**
- You say:** Humm... The more the merrier???
- He says:** Is this why I am paying you the big bucks???

8

©2017 Emily Fox

CSF 446: Machine Learning

## Simple bound (based on Hoeffding's Inequality)

- For  $N = \alpha_H + \alpha_T$  and  $\hat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T}$   
 $\uparrow$  total # datapoints
- Let  $\theta^*$  be the true parameter. For any  $\epsilon > 0$ :

$$P(|\hat{\theta}_{MLE} - \theta^*| \geq \epsilon) \leq 2e^{-2N\epsilon^2}$$

$\hat{\theta}_{MLE}$  is a bad est. of  $\theta^*$  (lose your job!)  
 error  
 small tolerance, say 0.1  
 # flips  
 goes down exp fast with  $N$ !  
 prob. of mistake

9

©2017 Emily Fox

CSF 446: Machine Learning

## PAC learning

- PAC:** Probably Approximate Correct
- Billionaire says:** I want to know the thumbtack parameter  $\theta$  within  $\epsilon = 0.1$ , with probability at least  $1 - \delta = 0.95$ . How many flips do I need?

$$P(|\hat{\theta}_{MLE} - \theta^*| \geq \epsilon) \leq 2e^{-2N\epsilon^2} \leq \delta$$

my tolerance to losing my job

$$\ln \delta \geq \ln 2 - 2N\epsilon^2$$

$$N \geq \frac{\ln 2 / \delta}{2\epsilon^2}$$

If  $\delta = 0.05$   $\epsilon = 0.1$   
 $\downarrow$   
 $N \geq 184.4$  flips  
 pretty loose bounds in practice

10

©2017 Emily Fox

CSF 446: Machine Learning

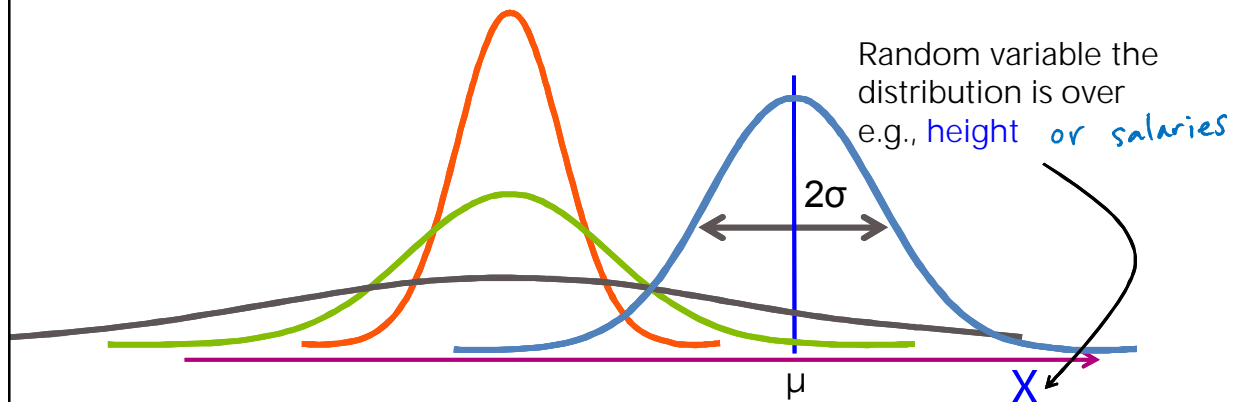
What about continuous-valued data?

What about continuous variables?

- **Billionaire says:** If I am measuring a continuous variable, what can you do for me? *salary of employees*
- **You say: Let me tell you about Gaussians...** *normal distributions*

# 1D Gaussians

Fully specified by **mean  $\mu$**  a **variance  $\sigma^2$**   
(or **standard deviation  $\sigma$** )

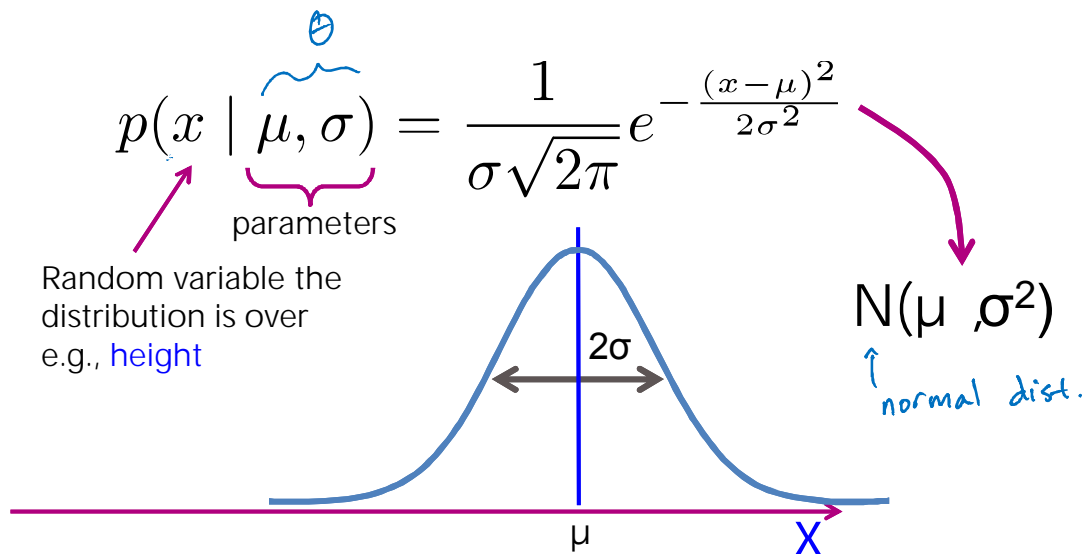


13

©2016 Emily Fox &amp; Carlos Guestrin

CSF 446: Machine Learning

## 1D Gaussian probability density function



14

©2016 Emily Fox &amp; Carlos Guestrin

CSF 446: Machine Learning

## Some properties of Gaussians

- Affine transformation (multiplying by scalar and adding a constant) *"distributed as"*
  - $X \sim N(\mu, \sigma^2)$
  - $Y = aX + b \rightarrow Y \sim N(a\mu + b, a^2\sigma^2)$  *same affine trans.  $E[aX+b] = aE[X] + b = a\mu + b$**known (deterministic) scalars*
- Sum of Gaussians
  - $X \sim N(\mu_X, \sigma_X^2)$
  - $Y \sim N(\mu_Y, \sigma_Y^2)$
  - $Z = X + Y \rightarrow Z \sim N(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2)$

15

©2017 Emily Fox

CSF 446: Machine Learning

## Learning a Gaussian

- Collect a bunch of data
  - Hopefully, i.i.d. samples
  - e.g., heights of students in class
- Learn parameters
  - Mean  $\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i$  *← why? MLE*
  - Variance  $\hat{\sigma}^2$

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

16

©2017 Emily Fox

CSF 446: Machine Learning



## MLE for Gaussian

- Prob. of i.i.d. samples  $D = \{x_1, \dots, x_N\}$ :

$$p(D | \mu, \sigma) = \left( \frac{1}{\sigma\sqrt{2\pi}} \right)^N \prod_{i=1}^N e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

dataset  $i$

- Log-likelihood of data:  $\hat{\mu}_{MLE}, \hat{\sigma}_{MLE} = \underset{\mu, \sigma}{\operatorname{argmax}} p(D | \mu, \sigma) = \underset{\mu, \sigma}{\operatorname{argmax}} \ln p(D | \mu, \sigma)$

$$\begin{aligned} \ln p(D | \mu, \sigma) &= \ln \left[ \left( \frac{1}{\sigma\sqrt{2\pi}} \right)^N \prod_{i=1}^N e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \right] \\ &= -N \ln \sigma\sqrt{2\pi} - \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2} \end{aligned}$$

17

©2017 Emily Fox

CSF 446: Machine Learning

## Your second learning algorithm: MLE for mean of a Gaussian

- What's MLE for the mean?

$$\frac{d}{d\mu} \ln p(D | \mu, \sigma) = \frac{d}{d\mu} \left[ -N \ln \sigma\sqrt{2\pi} - \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2} \right] = 0$$

no depend. on  $\mu$

$$\frac{d}{d\mu} \ln p(D | \mu, \sigma) = - \sum_{i=1}^N \frac{d}{d\mu} \frac{(x_i - \mu)^2}{2\sigma^2} = \sum_{i=1}^N \frac{x_i - \mu}{\sigma^2} = 0$$

$$N \hat{\mu} = \sum_{i=1}^N x_i$$

$$\hat{\mu}_{MLE} = \frac{1}{N} \sum_{i=1}^N x_i$$

MLE doesn't depend on choice of  $\sigma^2$

18

©2017 Emily Fox

CSF 446: Machine Learning

## MLE for variance

- Again, set derivative to zero:

$$\begin{aligned} \frac{d}{d\sigma} \ln p(D | \mu, \sigma) &= \frac{d}{d\sigma} \left[ -N \ln \sigma \sqrt{2\pi} - \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2} \right] \\ &= \frac{d}{d\sigma} \left[ -N \ln \sigma \sqrt{2\pi} \right] - \sum_{i=1}^N \frac{d}{d\sigma} \left[ \frac{(x_i - \mu)^2}{2\sigma^2} \right] = 0 \end{aligned}$$

Handwritten notes below the equation:

- $\Rightarrow -\frac{N}{\sigma} + \sum_{i=1}^N \frac{(x_i - \mu)^2}{\sigma^3} = 0$
- $\Rightarrow \hat{\sigma}_{MLE}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu}_{MLE})^2$
- use  $\mu = \hat{\mu}_{MLE}$  bc optimal choice of  $\mu$  doesn't depend on  $\sigma$

19

©2017 Emily Fox

CSF 446: Machine Learning

## Learning Gaussian parameters

- MLE:  $\hat{\mu}_{MLE} = \frac{1}{N} \sum_{i=1}^N x_i$
- $\hat{\sigma}_{MLE}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu}_{MLE})^2$
- FYI, MLE for the variance of a Gaussian is **biased**
  - Expected value of estimator is **not** true parameter!
  - Unbiased variance estimator:

$$\hat{\sigma}_{unbiased}^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \hat{\mu}_{MLE})^2$$

20

©2017 Emily Fox

CSF 446: Machine Learning

## Recap of concepts

## What you need to know...

- Learning is...
  - Collect some data
    - E.g., thumbtack flips
  - Choose a hypothesis class or model
    - E.g., binomial
  - Choose a loss function
    - E.g., data likelihood
  - Choose an optimization procedure
    - E.g., set derivative to zero to obtain MLE
  - Collect the big bucks
- Like everything in life, there is a lot more to learn...
  - Many more facets... Many more nuances...
  - The fun will continue...