

Announcements



Assignments

- HW6
 - Due Mon, 11/2, 11:59 pm

Midterm 2

- Mon, 11/9, during lecture
- See Piazza for details
- Forms for conflicts / tech issues due Fri, 10/30

Fireside Chat about the CMU ML PhD Program

- Fri, 10/30, 8:00 pm
- See Piazza for details, including form to show interest


Plan

Last Time

- PAC Criteria and Learning Theorems
- Bias-Variance trade-off as we change $|\mathcal{H}|$ or N
- Started VC dimension for infinite $|\mathcal{H}|$


$$N \geq \frac{\log |\mathcal{H}|}{\epsilon}$$

Today

- VC dimensions
- Learning theory and regularization 
- MLE
 - MLE for linear regression
- MAP (Maximum a posteriori) estimation
 - MAP for linear regression

Wrap up Learning Theory

Learning theory slides

An abstract graphic on the left side of the slide, featuring a sphere-like shape composed of a dense grid of intersecting red, green, and blue lines. The lines are curved and follow the contours of the sphere, creating a complex, woven pattern. The sphere is set against a dark gray background.

Introduction to Machine Learning

MLE & MAP

Instructor: Pat Virtue

Reminder MLE

$Y=1$ Heads
 $[1, 0, 1, 1]$

Trick coin

| | | | | |
|------------------|----------------------------|----------------------------|----------------------------|------------------|
| All T | $\frac{1}{3}$ H | Fair | $\frac{2}{3}$ H | All H |
| $\phi^{(A)} = 0$ | $\phi^{(B)} = \frac{1}{3}$ | $\phi^{(C)} = \frac{1}{2}$ | $\phi^{(D)} = \frac{2}{3}$ | $\phi^{(E)} = 1$ |

$$p(y^{(1)} \dots y^{(4)} | \phi^A) = \prod p(y^{(i)} | \phi^A) \quad \begin{array}{l} \leftarrow \phi \text{ heads} \\ \leftarrow (1-\phi) \text{ tail} \end{array}$$

$$= \phi^A \cdot (1-\phi^A) \cdot \phi^A \cdot \phi^A \dots$$

$$= 0 \cdot 1 \cdot 0 \cdot 0 = 0$$

$$\hat{\phi}_{MLE} = \operatorname{argmax}_{\phi} \prod_i^N p(y^{(i)} | \phi)$$

Previous Piazza Poll

$$\hat{\phi}_{MLE} = \operatorname{argmax}_{\phi} \prod_i^N p(y^{(i)} | \phi)$$

We model the outcome of a single mysterious weighted-coin flip as a Bernoulli random variable:

$$Y \sim \text{Bern}(\phi)$$
$$p(y | \phi) = \begin{cases} \phi, & y = 1 \text{ (Heads)} \\ 1 - \phi, & y = 0 \text{ (Tails)} \end{cases}$$

Given the ordered sequence of coin flip outcomes:

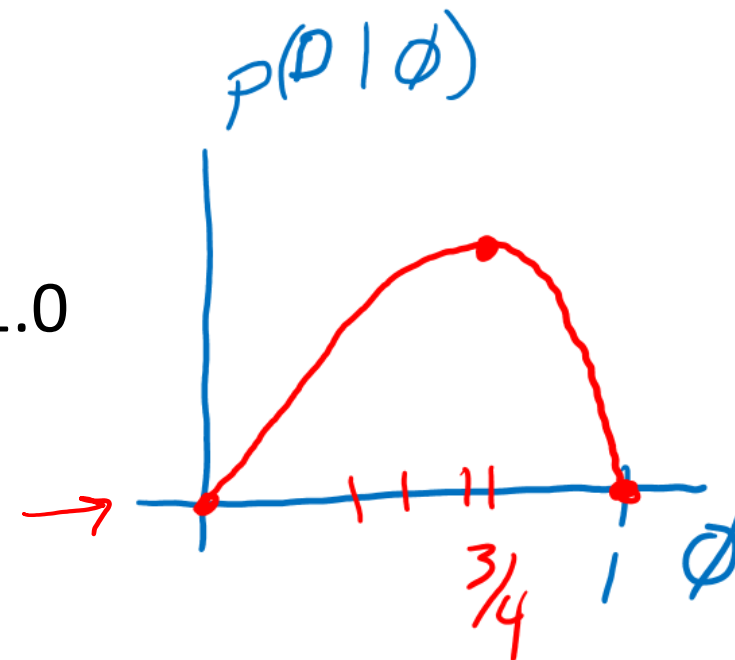
[1, 0, 1, 1] ←

What is the estimate of parameter $\hat{\phi}$?

A. 0.0 B. 1/8 C. 1/4 D. 1/2 E. 3/4 F. 3/8 G. 1.0

Why?

$$p(\mathcal{D} | \phi) = \phi^3 (1 - \phi)^1$$



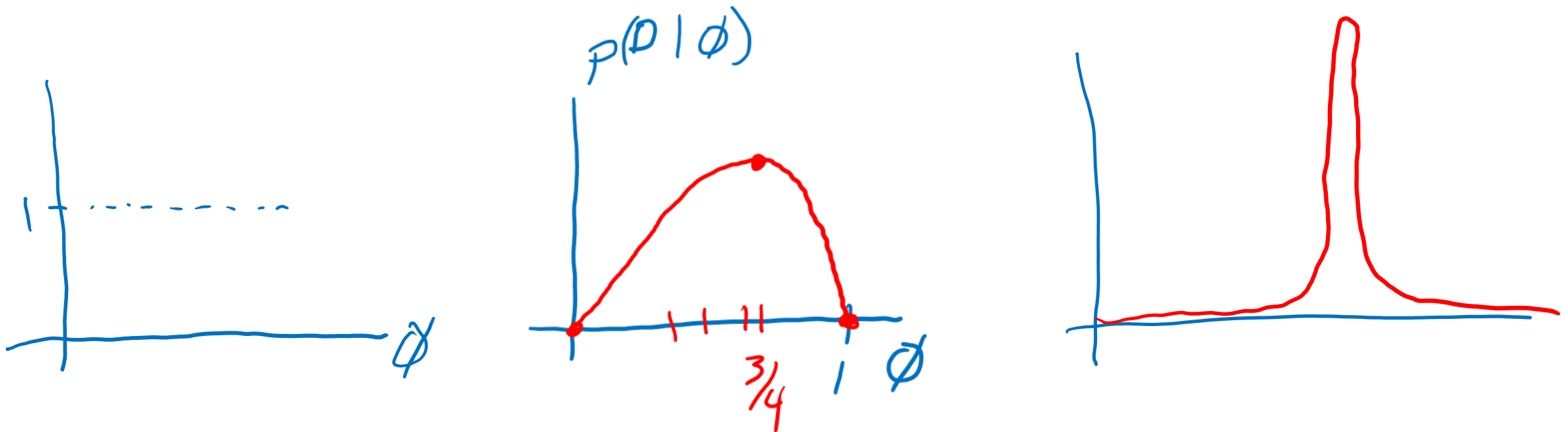
MLE as Data Increases

Given the ordered sequence of coin flip outcomes:

$[1, 0, 1, 1]$

$$p(\mathcal{D} \mid \phi) = \prod_i^N p(y^{(i)} \mid \phi) = \phi^{N_{y=1}} (1 - \phi)^{N_{y=0}}$$

What happens as we flip more coins?



MLE for Gaussian

Gaussian distribution:

$$Y \sim \mathcal{N}(\mu, \sigma^2)$$

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

What is the log likelihood for three i.i.d. samples, given parameters μ, σ^2 ?

$$\mathcal{D} = \{\underline{y^{(1)}} = 65, \underline{y^{(2)}} = 95, \underline{y^{(3)}} = 85\}$$

$$\rightarrow L(\mu, \sigma^2) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y^{(i)}-\mu)^2}{2\sigma^2}}$$

$$\ell(\mu, \sigma^2) = \sum_{i=1}^N -\log \sqrt{2\pi\sigma^2} - \frac{(y^{(i)} - \mu)^2}{2\sigma^2}$$

$$\hat{\theta}_{MLE} = \operatorname{argmax}_{\theta} \prod_i^N p(y^{(i)} \mid \theta) \quad \checkmark$$

$$\hat{\theta}_{MLE} = \operatorname{argmax}_{\theta} \sum_i^N \log p(y^{(i)} \mid \theta) \quad \checkmark$$

Recipe for Estimation

MLE

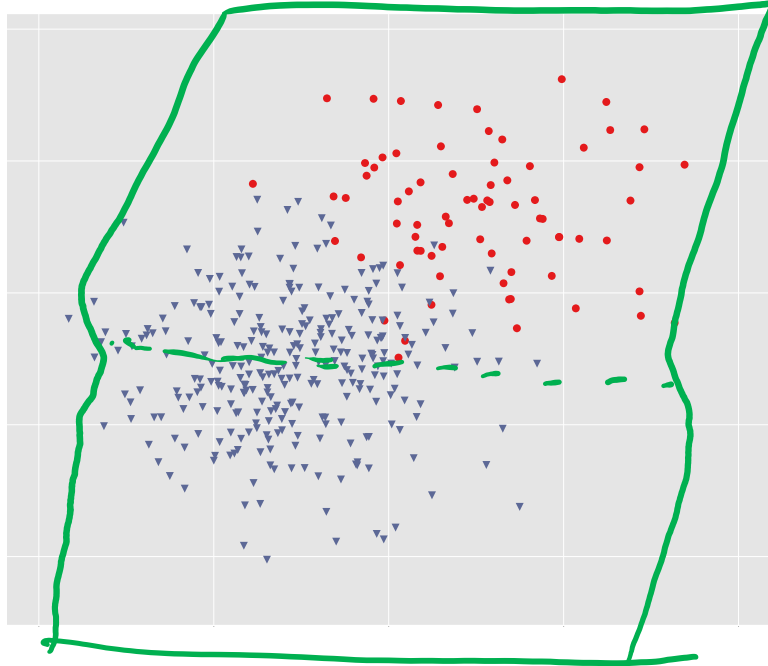
1. Formulate the likelihood, $p(\mathcal{D} \mid \theta)$
2. Set objective $J(\theta)$ equal to negative log of likelihood
$$J(\theta) = -\log p(\mathcal{D} \mid \theta)$$
3. Compute derivative of objective, $\partial J / \partial \theta$
4. Find $\hat{\theta}$, either
 - a. Set derivate equal to zero and solve for θ
 - b. Use (stochastic) gradient descent to step towards better θ

M(C)LE for Logistic Regression

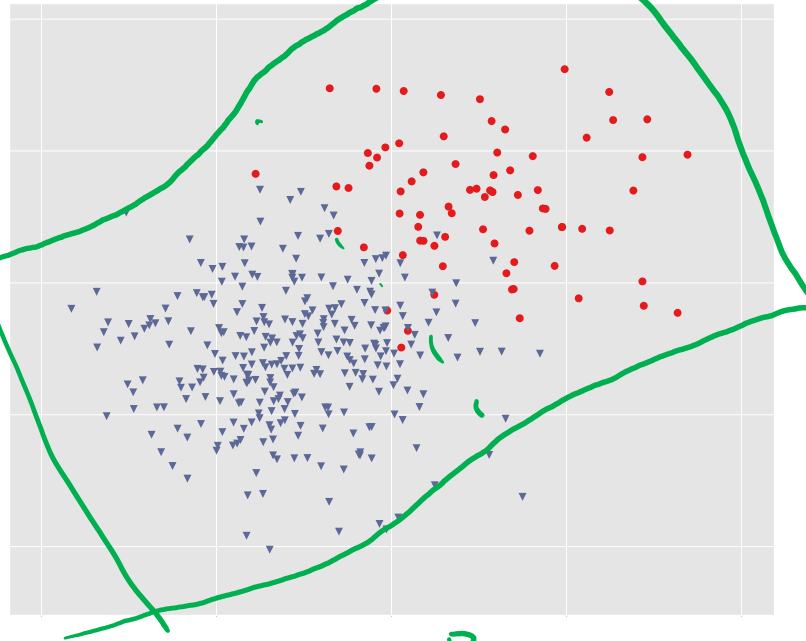
$$p(y|x, \theta)$$

Learn to predict if a patient has cancer ($Y = 1$) or not ($Y = 0$) given the input of just one test result, X_A .

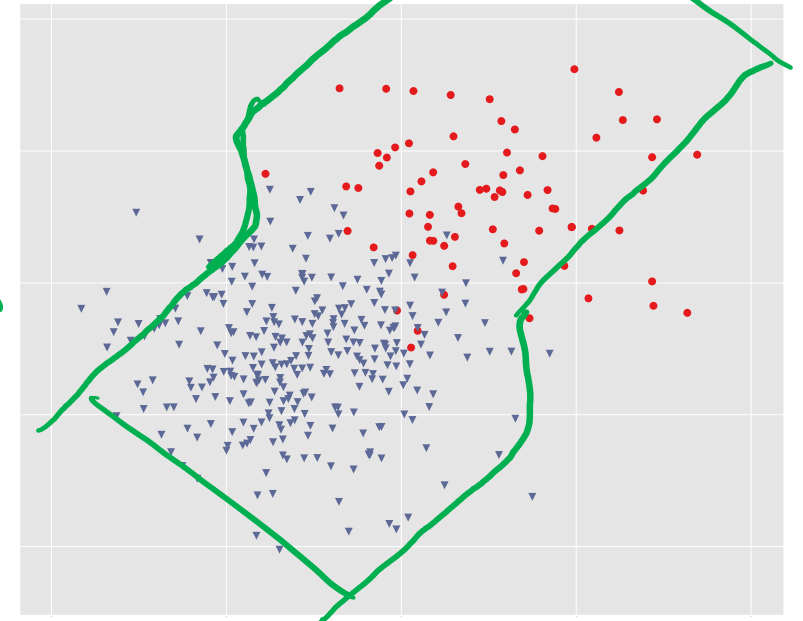
$$\hat{\theta}_{MLE} = \operatorname{argmax}_{\theta} \prod_i^N \frac{1}{1 + e^{-\theta^T x^{(i)}}}^{\mathbb{I}(y^{(i)=1})} \left(1 - \frac{1}{1 + e^{-\theta^T x^{(i)}}} \right)^{\mathbb{I}(y^{(i)=0})}$$



θ^A



θ^B

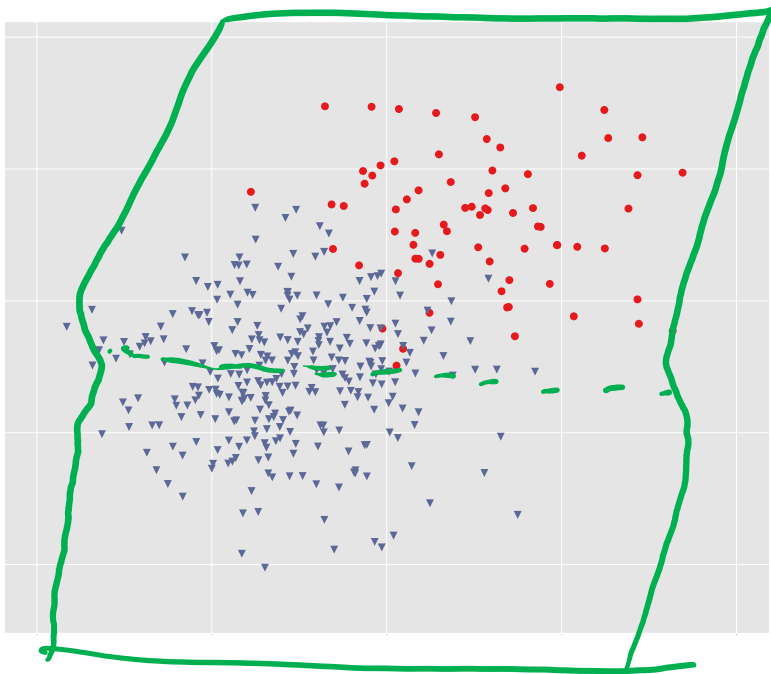


θ^C

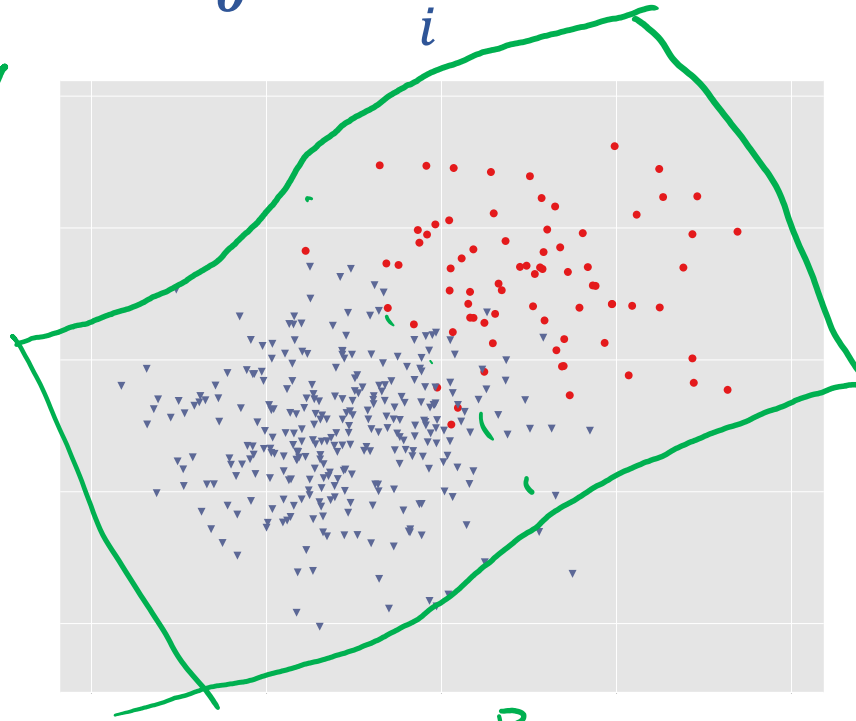
M(C)LE for Logistic Regression

Learn to predict if a patient has cancer ($Y = 1$) or not ($Y = 0$) given the input of just one test result, X_A .

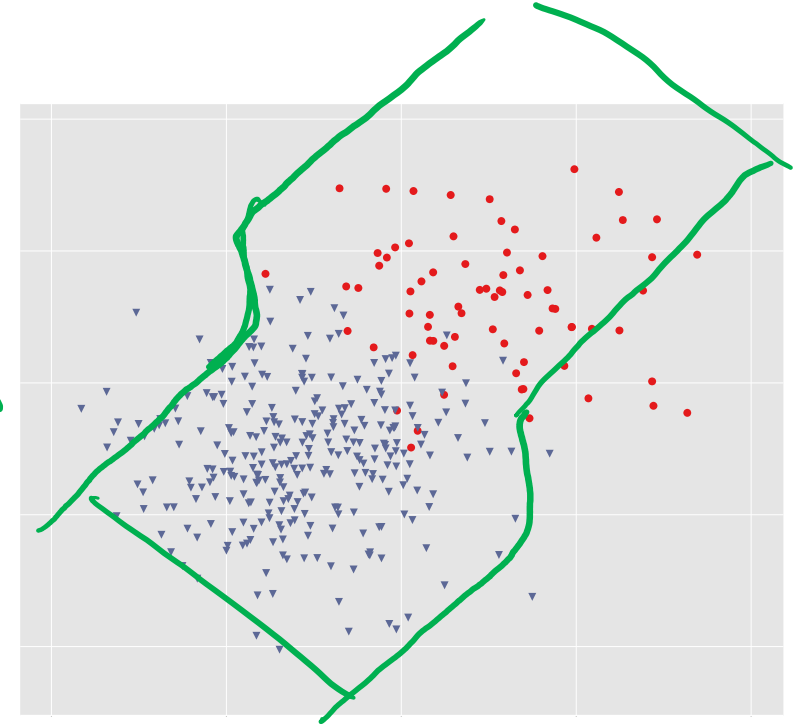
$$\hat{\theta}_{MLE} = \operatorname{argmax}_{\theta} \prod_i^N p(y^{(i)} | x^{(i)}, \theta)$$



θ^A



θ^B

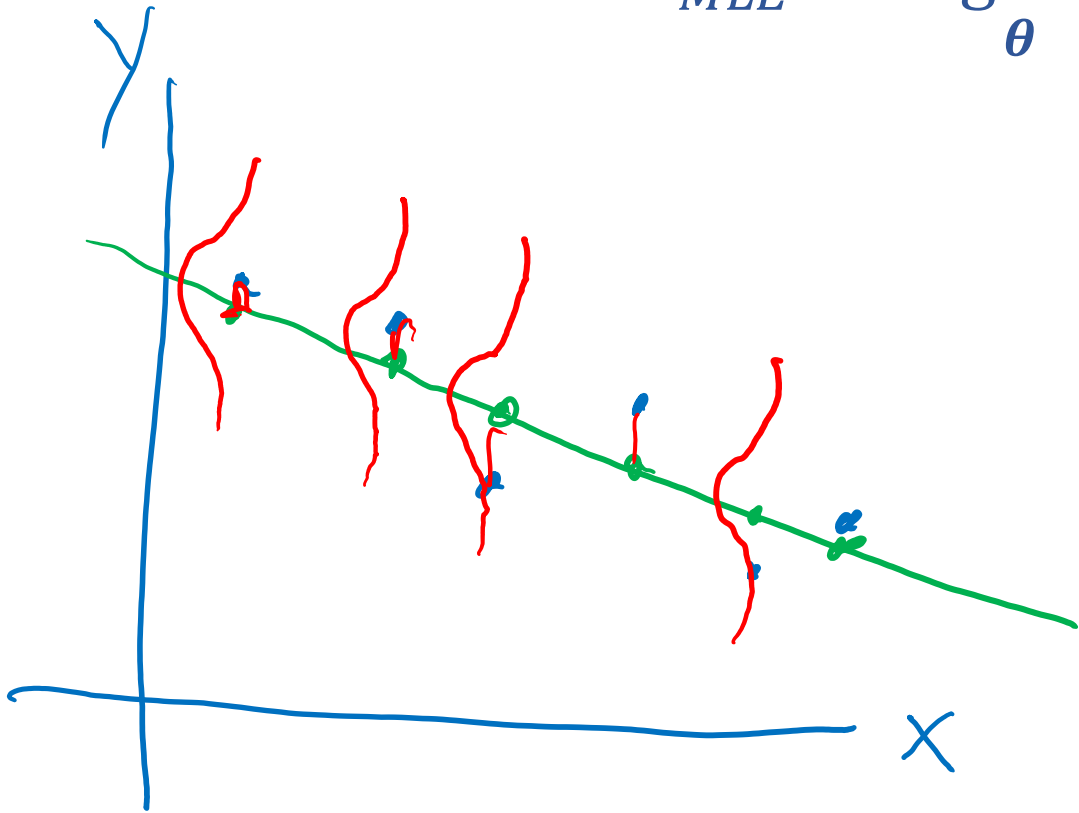


θ^C

M(C)LE for Linear Regression

Probabilistic interpretation of linear regression

$$\hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} \prod_i^N p(y^{(i)} | x^{(i)}, \theta)$$



$$y = \underline{w^T x + b} + \underline{\epsilon}$$
$$\epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$y \sim \underline{\hspace{10em}}$$

$$p(y^{(i)} | x^{(i)}, w, b) =$$

FROM MLE TO MAP

Product Rule

Construct the joint by multiplying the conditional by the appropriate marginal

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

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

Also works when something is given everywhere

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

$$P(A, B | C, D, E) = P(A | B, C, D, E)P(B | C, D, E)$$

Coin Flipping Example

Trick coin: Suppose I know how many coins are in each container in the store. How can I use this information both before and after flipping coins?

| | | | | | |
|------------------|--------------------|--------------------|--------------------|------------------|------------------|
| $8/20$ | $2/20$ | $4/20$ | $3/20$ | $3/20$ | $p(\phi_A) = .4$ |
| All T | $1/3$ H | Fair | $2/3$ H | All H | |
| $\phi^{(A)} = 0$ | $\phi^{(B)} = 1/3$ | $\phi^{(C)} = 1/2$ | $\phi^{(D)} = 2/3$ | $\phi^{(E)} = 1$ | |

$$p(D|\phi_A)$$

$$\rightarrow p(\theta_A|D)$$

Likelihood, Prior, and Posterior


Likelihood: $p(\mathcal{D} \mid \theta)$

Joint: $p(\mathcal{D}, \theta) = p(\mathcal{D} \mid \theta) p(\theta)$


Prior: $p(\theta)$

Posterior: $p(\theta \mid \mathcal{D})$

Relating these with Bayes rule

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


X

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


MLE and MAP

Likelihood: $p(\mathcal{D} \mid \theta)$

Joint: $p(\mathcal{D}, \theta)$

Prior: $p(\theta)$

Posterior: $p(\theta \mid \mathcal{D})$

$$\underline{p(\theta \mid \mathcal{D})} \propto \underbrace{p(\mathcal{D} \mid \theta)p(\theta)}$$

MLE: $\hat{\theta}_{MLE} = \operatorname{argmax}_{\theta} p(\mathcal{D} \mid \theta)$

MAP: $\hat{\theta}_{MAP} = \operatorname{argmax}_{\theta} p(\mathcal{D} \mid \theta)p(\theta)$

$p(\mathcal{D}, \theta)$

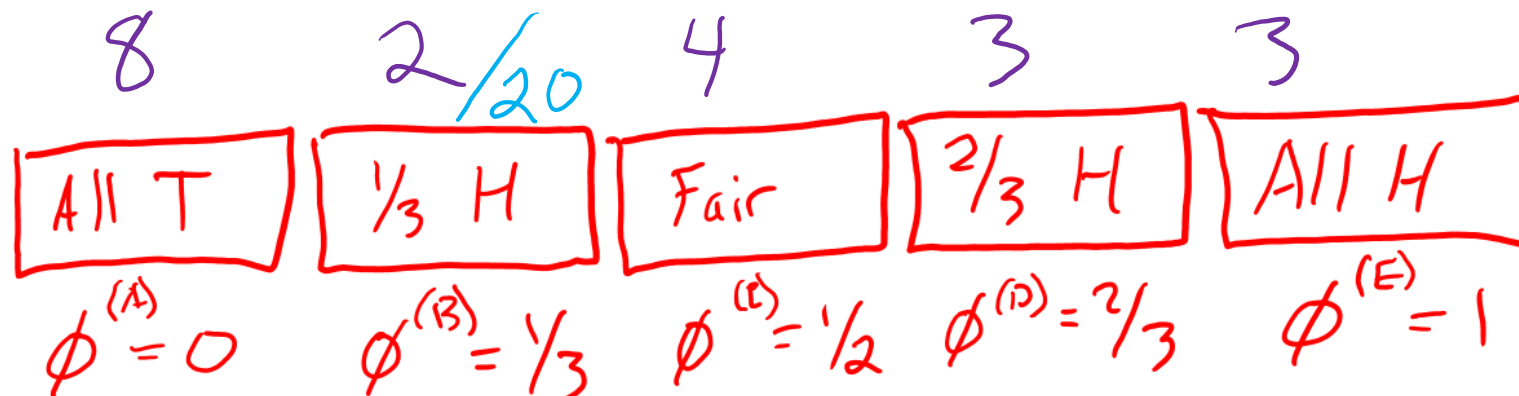
$p(\theta \mid \mathcal{D})$

Maximum *a posteriori* estimation

Coin Flipping Example

$$\hat{\theta}_{MAP} = \operatorname{argmax}_{\theta} \underbrace{\prod_{i=1}^N p(y^{(i)} | \theta)}_{\text{likelihood}} \underbrace{p(\theta)}_{\text{prior}}$$

Trick coin: Suppose I know how many coins are in each container in the store. How can I use this information both before and after flipping coins?



1 1 0 1

$$\phi_B \phi_B (1 - \phi_B) \phi_B \cdot \frac{2}{20}$$

Piazza Poll 1:

$$\underline{p(\theta \mid \mathcal{D})} \propto \underline{p(\mathcal{D} \mid \theta)} \underline{p(\theta)} \quad \underline{\text{posterior}} \propto \underline{\text{likelihood}} \cdot \underline{\text{prior}}$$

$$p(\theta \mid \mathcal{D}) \propto \underbrace{\prod p(\mathcal{D}^{(n)} \mid \theta)} p(\theta)$$

As the number of data points increases, which of the following are true?

Select ALL that apply

- A. The MAP estimate approaches the MLE estimate
- B. The **posterior** distribution approaches the **prior** distribution
- C. The **likelihood** distribution approaches the **prior** distribution
- D. The **posterior** distribution approaches the **likelihood** distribution
- E. The **likelihood** has a lower impact on the **posterior**
- F. The **prior** has a lower impact on the **posterior**

Piazza Poll 1:

$$p(\theta \mid \mathcal{D}) \propto p(\mathcal{D} \mid \theta) p(\theta) \quad \text{posterior} \propto \text{likelihood} \cdot \text{prior}$$

$$p(\theta \mid \mathcal{D}) \propto \prod \underbrace{p(\mathcal{D}^{(n)} \mid \theta)} \underbrace{p(\theta)}$$

As the number of data points increases, which of the following are true?

Select ALL that apply

- A. The MAP estimate approaches the MLE estimate
- B. The posterior distribution approaches the prior distribution
- C. The likelihood distribution approaches the prior distribution
- D. The posterior distribution approaches the likelihood distribution ←
- E. The likelihood has a lower impact on the posterior
- F. The prior has a lower impact on the posterior ←

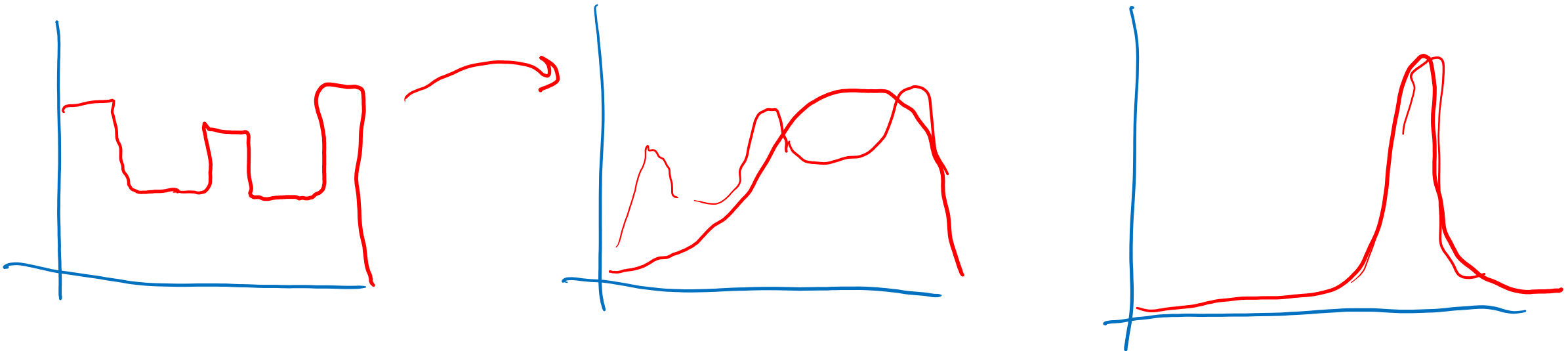
MAP as Data Increases

Given the ordered sequence of coin flip outcomes:

$$\mathcal{D} = [1, 0, 1, 1]$$

$$p(\mathcal{D} \mid \phi) p(\phi) = \prod_i^N p(y^{(i)} \mid \phi) p(\phi) = \phi^{N_{y=1}} (1 - \phi)^{N_{y=0}} p(\phi)$$

What happens as we flip more coins?



Recipe for Estimation

MLE

1. Formulate the likelihood, $p(\mathcal{D} \mid \theta)$
2. Set objective $J(\theta)$ equal to negative log of likelihood
$$J(\theta) = -\log p(\mathcal{D} \mid \theta)$$
3. Compute derivative of objective, $\partial J / \partial \theta$
4. Find $\hat{\theta}$, either
 - a. Set derivate equal to zero and solve for θ
 - b. Use (stochastic) gradient descent to step towards better θ

Recipe for Estimation

MAP

1. Formulate the likelihood **times the prior**, $p(\mathcal{D} \mid \theta) \underline{p(\theta)}$
2. Set objective $J(\theta)$ equal to negative log of likelihood **times the prior**
$$J(\theta) = -\log[p(\mathcal{D} \mid \theta) \underline{p(\theta)}]$$
3. Compute derivative of objective, $\partial J / \partial \theta$
4. Find $\hat{\theta}$, either
 - a. Set derivate equal to zero and solve for θ
 - b. Use (stochastic) gradient descent to step towards better θ

M(C)LE for Linear Regression

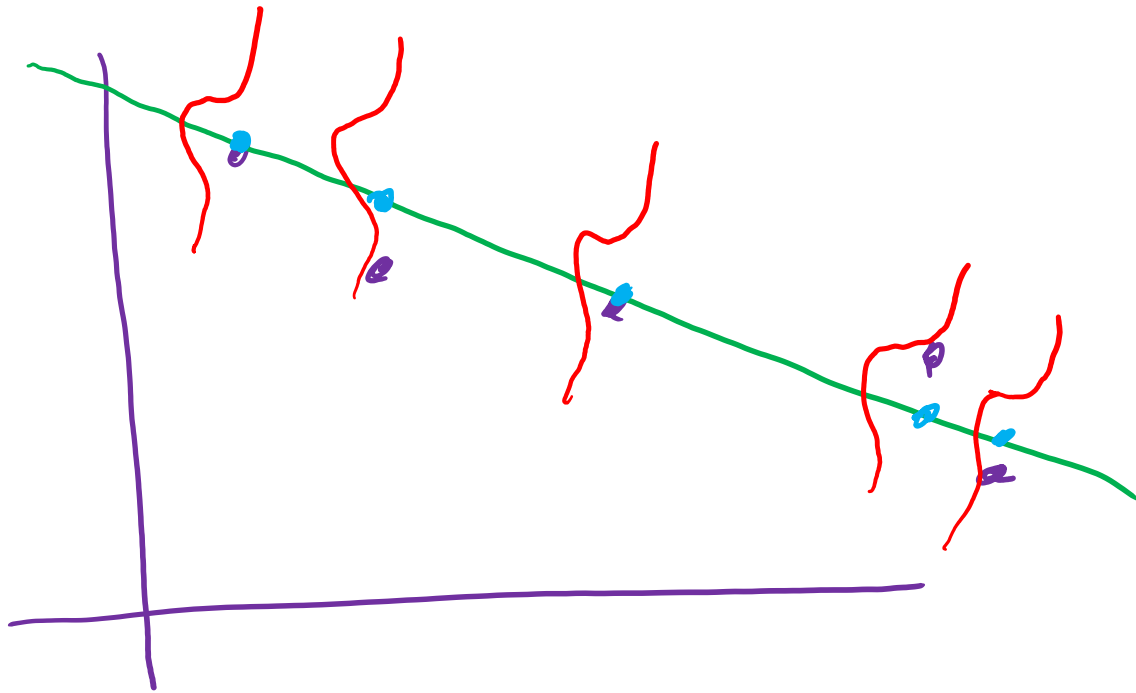
Probabilistic interpretation of linear regression

$$\hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} \prod_i^N p(y^{(i)} | x^{(i)}, \theta)$$

$$\begin{cases} y = \underline{w}^T x + \epsilon \\ \epsilon \sim \mathcal{N}(0, \tau) \end{cases}$$

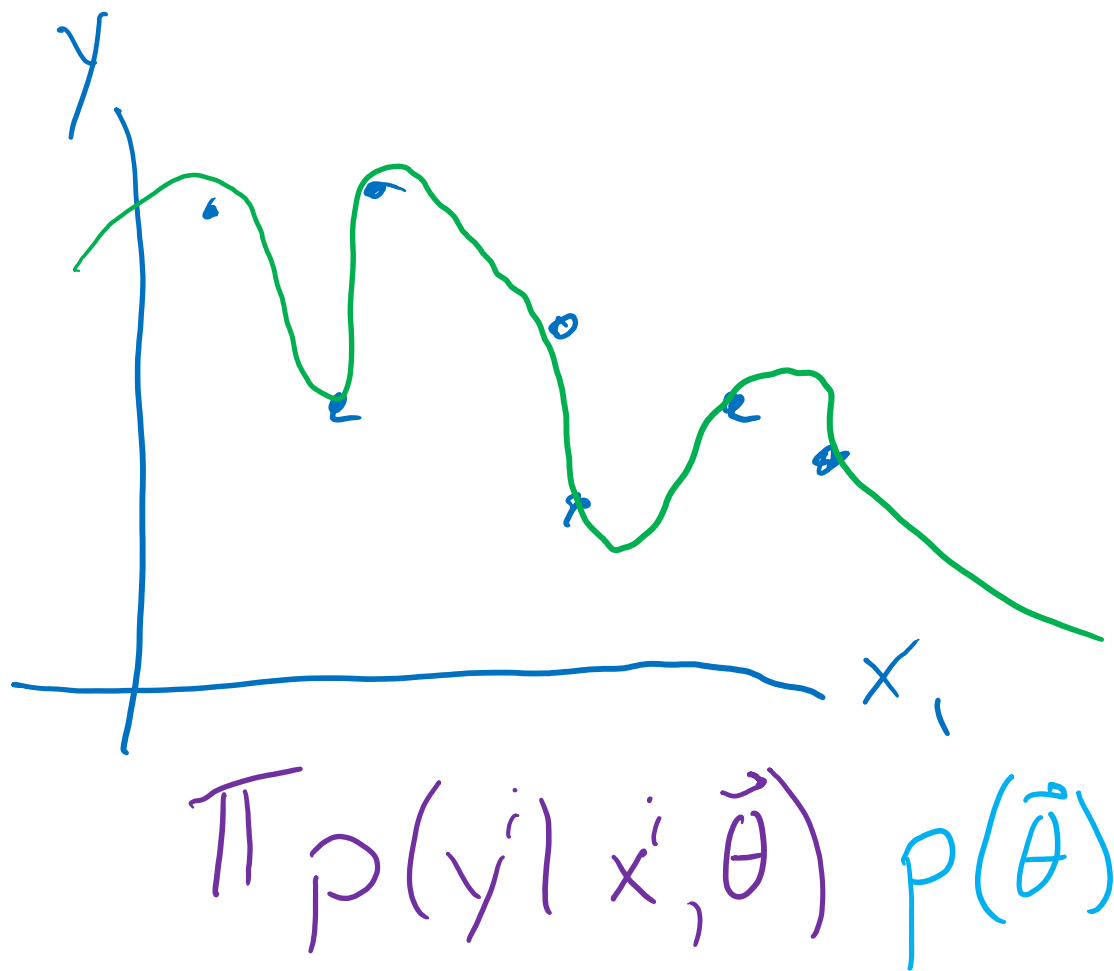
$$y \sim \underbrace{\hspace{10em}}$$

$$p(y | x, \vec{w}) = \underbrace{\hspace{10em}}$$



MAP for Linear Regression

What assumptions are we making about our parameters?



$$\begin{bmatrix} 1 \\ x_i \end{bmatrix} \xrightarrow{\phi}$$

$$\begin{bmatrix} 1 \\ x_{i,1} \\ x_{i,2} \\ x_{i,3} \\ x_{i,4} \end{bmatrix} \quad \begin{matrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \end{matrix}$$

$$\theta_i \sim \mathcal{N}(0, \tau^2)$$

$$\vec{w}, b \quad \theta \quad p(\theta_i) \rightarrow$$

