

# Lecture 5: 2/1/17

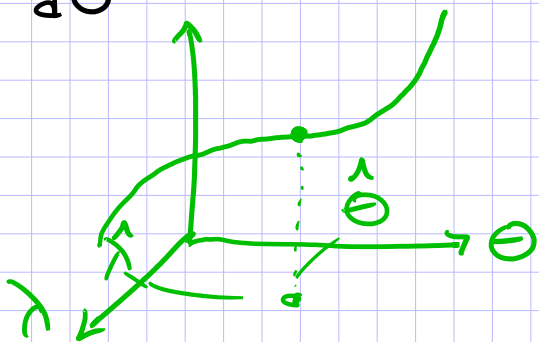
## Aside: Method of Lagrange Multipliers

① Given problem:  $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} f(\theta)$

s.t.  $g(\theta) = c$

② Introduce Lagrangian:  $L(\theta, \lambda) = f(\theta) - \lambda(g(\theta) - c)$

③ Solve  $\frac{dL(\theta, \lambda)}{d\theta} = 0$  and  $\frac{dL(\theta, \lambda)}{d\lambda} = 0$



## Back to MLE of Cat.

★ Apply the Method Lagrange Mult.

$$L(\vec{\theta}, \lambda) = l(\theta) - \lambda \left( \sum_{k=1}^K \theta_k - 1 \right)$$

# times  
that  $k$  appears  
in  $D$

$$= \sum_{k=1}^K N_k \log \theta_k - \lambda \left( \sum_{k=1}^K \theta_k - 1 \right)$$

$$\frac{dL(\vec{\theta}, \lambda)}{d\theta_j} = \frac{N_j}{\theta_j} - \lambda = 0 \Rightarrow \theta_j = \frac{N_j}{\lambda}$$

$$\frac{dL(\vec{\theta}, \lambda)}{d\lambda} = \sum_{k=1}^K \theta_k - 1 = 0 \Rightarrow \sum_{k=1}^K \theta_k = 1$$

$$\forall j \quad \theta_j = \frac{N_j}{\sum_{k=1}^K N_k}$$

# Apply MLE to Data

Ex #1:  $D = \{x^{(1)}, \dots, x^{(n)}\} = [1, 1, 1, 1, 1, 1, 2, 2, 2, 2, 3]$

$X^{(i)} \sim \text{Categorical}(\vec{\theta}) \quad K \triangleq |\vec{\theta}| = 3$

$$\vec{\theta}_{MLE} = \begin{bmatrix} \theta_{1,MLE} \\ \theta_{2,MLE} \\ \theta_{3,MLE} \end{bmatrix} = ? = \begin{bmatrix} 6/11 \\ 4/11 \\ 1/11 \end{bmatrix}$$

Ex #2: Same dataset  $D$  as above.  
But  $K = 6$

$$\vec{\theta}_{MLE} = \begin{bmatrix} \theta_{1,MLE} \\ \theta_{2,MLE} \\ \vdots \\ \theta_{6,MLE} \end{bmatrix} = ? = \begin{bmatrix} 6/11 \\ 4/11 \\ 1/11 \\ 0/11 \\ 0/11 \\ 0/11 \end{bmatrix}$$

Suppose  $D$   
is from six-sided  
die.

Is MLE  
still a good  
estimate?

# Learning from Data (Bayesian)

## MAP Estimation:

- "Maximum a posteriori"  $\iff$  MAP
- Bayesian statistics
- key idea: inject world knowledge into estimation problem via a prior over parameters
- MAP estimate is the mode of the posterior

Def: posterior

$$p(\theta | D) = \frac{\overbrace{p(D|\theta)}^{\text{likelihood}} \overbrace{p(\theta)}^{\text{prior}}}{\underbrace{p(D)}_{\text{posterior}}}$$

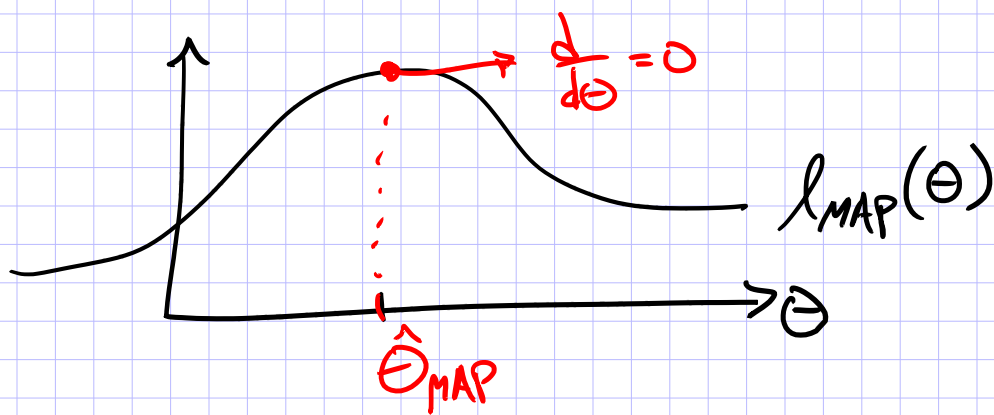
Bayes' Rule

## Opt. for MAP Est.

$$\begin{aligned}\hat{\theta}_{\text{MAP}} &= \operatorname{argmax}_{\theta \in \Theta} p(\theta | D) \\ &= \operatorname{argmax}_{\theta \in \Theta} \log p(\theta | D) \\ &= \operatorname{argmax}_{\theta \in \Theta} \log \left[ \frac{p(D|\theta) p(\theta)}{p(D)} \right] \\ &= \operatorname{argmax}_{\theta \in \Theta} \log [p(D|\theta) p(\theta)] \\ &= \operatorname{argmax}_{\theta \in \Theta} l_{\text{MAP}}(\theta)\end{aligned}$$

$p(D)$  is constant wrt.  $\theta$   
why convenient?  
 $p(D) = \int_{\Theta} p(D|\theta) p(\theta) d\theta$

where  $l_{\text{MAP}}(\theta) \triangleq$



# Ex: MAP Estimate of a Bernoulli

① Write the likelihood

$$X^{(i)} \sim \text{Bernoulli}(\phi) \quad \text{where } 0 \leq \phi \leq 1$$

$$p(X=x) = \begin{cases} \phi & \text{if } x=1 \\ 1-\phi & \text{if } x=0 \end{cases}$$

$$= (\phi)^x (1-\phi)^{(1-x)}$$

notation  
trick

$$\Rightarrow p(D|\phi) = \prod_{i=1}^N (\phi)^{X^{(i)}} (1-\phi)^{(1-X^{(i)})} = \phi^{N_1} (1-\phi)^{N_0}$$

$$\text{where } N_1 = \# \text{ 1s}$$

$$N_0 = \# \text{ 0s}$$

② Write prior

$$p(\phi) = \begin{cases} ? & \text{if } 0 \leq \phi \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

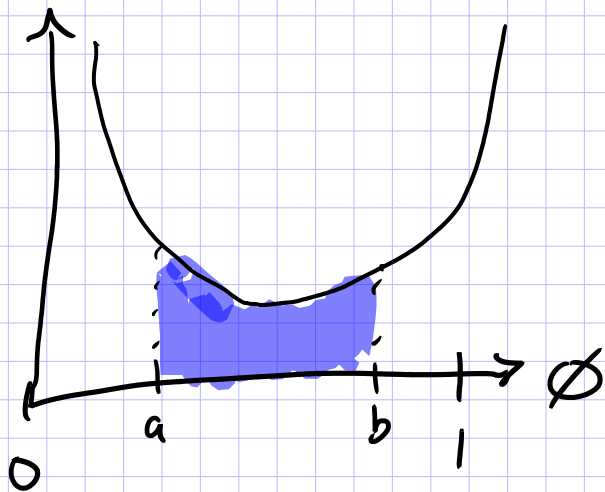
$$p(\phi) \sim \text{Beta}(\alpha, \beta)$$

probability density function:

$$f(\phi|\alpha, \beta) = \frac{1}{B(\alpha, \beta)}$$

$$\phi^{(\alpha-1)} (1-\phi)^{(\beta-1)}$$

$$p(a \leq \phi \leq b) = \int_a^b f(\phi|\alpha, \beta) d\phi$$



③ Write Joint (likelihood times prior)

$$p(D|\theta)p(\theta) = [\phi^{N_1} (1-\phi)^{N_0}] [\phi^{(\alpha-1)} (1-\phi)^{(\beta-1)}] \cdot \frac{1}{B(\alpha, \beta)}$$

$$= \underbrace{\phi^{(N_1 + \alpha - 1)} (1-\phi)^{(N_0 + \beta - 1)}}_{\text{identical to likelihood where } (N_1 + \alpha - 1) \text{ "heads" and } (N_0 + \beta - 1) \text{ "tails"}} \cdot \frac{1}{B(\alpha, \beta)}$$

$$D = \left[ \underbrace{000\dots 0}_{N_0} \underbrace{00\dots 0}_{\beta-1} \underbrace{111\dots 1}_{N_1} \underbrace{11\dots 1}_{\alpha-1} \right]$$

④ Solve the optimization

$$\phi_{\text{MAP}} = \underset{\phi \in \mathbb{R}}{\text{argmax}} \ell_{\text{MAP}}(\phi) \quad \text{s.t. } 0 \leq \phi \leq 1$$

Q: Do we need Lagrange Mult? No!

a) Take the derivative

$$\ell_{\text{MAP}}(\phi) = N_\alpha \log \phi + N_\beta \log(1-\phi)$$

$$\text{where } N_\alpha \triangleq N_1 + \alpha - 1$$

$$N_\beta \triangleq N_0 + \beta - 1$$

$$\frac{d \ell_{\text{MAP}}(\phi)}{d \phi} = \frac{N_\alpha}{\phi} + \frac{N_\beta}{\phi-1}$$

b) set to zero and solve

$$\frac{d \ell_{\text{MAP}}(\phi)}{d \phi} = 0 \Rightarrow N_\alpha(\phi-1) + N_\beta(\phi) = 0$$

$$\Rightarrow (N_\alpha + N_\beta)\phi = N_\alpha$$

$$\Rightarrow \phi = \frac{N_\alpha}{N_\alpha + N_\beta}$$

$$= \frac{N_1 + \alpha - 1}{N_1 + \alpha - 1 + N_0 + \beta - 1}$$

Intuition: proportional to # observed and hallucinated heads

# Probabilistic View of Classification (Generative)

①  $\exists$  some unknown dist.  $p^*$  parameterized by  $\theta^*$  that generates labeled iid examples

$$(\vec{x}^{(i)}, y^{(i)}) \sim p^*(\vec{x}, y | \theta^*), \forall i$$

$$y^{(i)} \sim p(y | \theta^*)$$

$$\vec{x}^{(i)} \sim p(\vec{x} | y, \theta^*)$$

② Learning Algo.

$$\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmax}} \log \prod_{i=1}^N p(\vec{x}^{(i)}, y^{(i)} | \theta) \quad \underbrace{p(\theta)}_{\text{optimal}}$$

③ Convert to a decision rule:  
output most likely label  $\hat{y}$  for input  $\vec{x}$

$$\hat{y} = h(\vec{x}) \triangleq \underset{y \in \{0,1\}}{\operatorname{argmax}} p(y | \vec{x}, \hat{\theta})$$

$$= \underset{y \in \{0,1\}}{\operatorname{argmax}} p(\vec{x} | y, \hat{\theta}) p(y | \hat{\theta})$$

Bayes' Rule

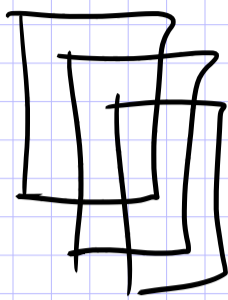
# Natural Data

$y$  (label)

$E_{con} = 1$

$O_{un} = 0$

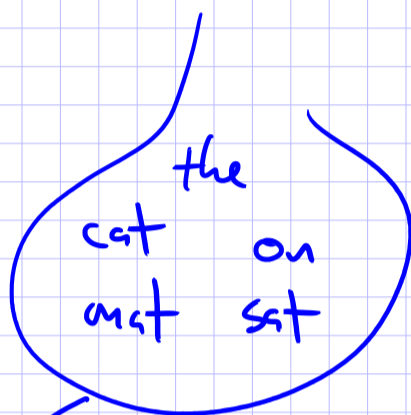
$\vec{x}$  (words)



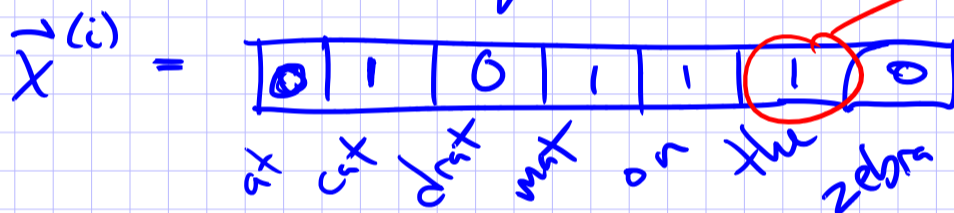
Conversion #1  
i<sup>th</sup> Doc.



i<sup>th</sup> bag-of-words



Conversion #2



not counts  
just indicators

## Generate Synthetic Docs

$y \in \{0, 1\}$  where  $E_{con} = 1$   $O_{un} = 0$

$\vec{x} \in \{0, 1\}^M$  where  $M = \#$  of words in vocabulary

$\vec{x}^T = [x_1, \dots, x_M]$

Generative Story:

$\theta \in \mathbb{R}$

$$\theta = \begin{bmatrix} \theta_{10} & \theta_{11} \\ \theta_{20} & \theta_{21} \\ \vdots & \vdots \\ \theta_{M0} & \theta_{M1} \end{bmatrix}$$

$$y \sim \text{Bernoulli}(\theta) = p(y|\theta)$$

$$x_1 \sim \text{Bernoulli}(\theta_{1,y}) = p(x_1|y, \theta)$$

$$x_2 \sim \text{Bern.}(\theta_{2,y}) = p(x_2|y, \theta)$$

$$\vdots$$

$$x_M \sim \text{Bern}(\theta_{M,y}) = p(x_M|y, \theta)$$

$$p(x_1, x_2, \dots, x_M, y|\theta, \theta) = p(y|\theta) p(x_1|y, \theta) \dots p(x_M|y, \theta)$$

Because each coin flip (for  $x_n$ ) is conditionally independent given  $y$

## Naive Bayes Model

$$p(\vec{x}, y | \phi, \theta) = p(y | \phi) \prod_{m=1}^M p(x_m | y, \theta)$$