

# Bayesian Methods (Lab)

David Rosenberg

New York University

October 29, 2016

# Coin Flipping

- **Parameter space**  $\theta \in \Theta = [0, 1]$ :

$$\mathbb{P}(\text{Heads} \mid \theta) = \theta.$$

- **Data**  $\mathcal{D} = \{H, H, T, T, T, T, T, H, \dots, T\}$ 
  - $n_h$ : number of heads
  - $n_t$ : number of tails

- **Likelihood model** (Bernoulli Distribution):

$$p(\mathcal{D} \mid \theta) = \theta^{n_h} (1 - \theta)^{n_t}$$

- (probability of getting the flips in the order they were received)

# Coin Flipping: Beta Prior

- **Prior:**

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

$$p(\theta) \propto \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

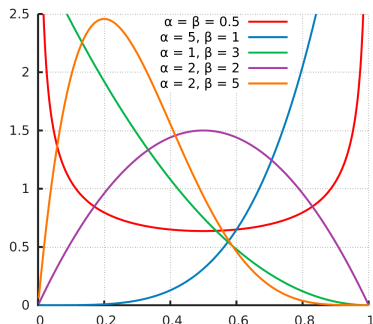


Figure by Horas based on the work of Krishnavedala (Own work) [Public domain], via Wikimedia Commons  
[http://commons.wikimedia.org/wiki/File:Beta\\_distribution\\_pdf.svg](http://commons.wikimedia.org/wiki/File:Beta_distribution_pdf.svg).

# Coin Flipping: Beta Prior

- **Prior:**

$$\begin{aligned}\theta &\sim \text{Beta}(h, t) \\ p(\theta) &\propto \theta^{h-1} (1-\theta)^{t-1}\end{aligned}$$

- **Mean of Beta distribution:**

$$\mathbb{E}\theta = \frac{h}{h+t}$$

# Coin Flipping: Posterior

- **Prior:**

$$\begin{aligned}\theta &\sim \text{Beta}(h, t) \\ p(\theta) &\propto \theta^{h-1} (1-\theta)^{t-1}\end{aligned}$$

- **Likelihood model:**

$$p(\mathcal{D} | \theta) = \theta^{n_h} (1-\theta)^{n_t}$$

- **Posterior density:**

$$\begin{aligned}p(\theta | \mathcal{D}) &\propto p(\theta)p(\mathcal{D} | \theta) \\ &\propto \theta^{h-1} (1-\theta)^{t-1} \times \theta^{n_h} (1-\theta)^{n_t} \\ &= \theta^{h-1+n_h} (1-\theta)^{t-1+n_t}\end{aligned}$$

# Posterior is Beta

- **Prior:**

$$\begin{aligned}\theta &\sim \text{Beta}(h, t) \\ p(\theta) &\propto \theta^{h-1} (1-\theta)^{t-1}\end{aligned}$$

- **Posterior density:**

$$p(\theta | \mathcal{D}) \propto \theta^{h-1+n_h} (1-\theta)^{t-1+n_t}$$

- **Posterior is in the beta family:**

$$\theta | \mathcal{D} \sim \text{Beta}(h + n_h, t + n_t)$$

- **Interpretation:**

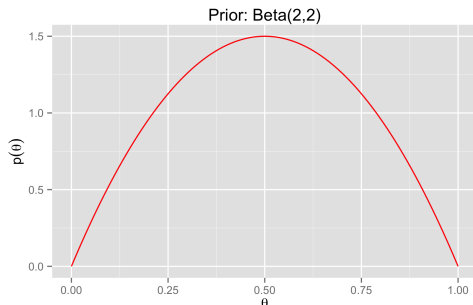
- Prior initializes our counts with  $h$  heads and  $t$  tails.
- Posterior increments counts by observed  $n_h$  and  $n_t$ .

# Example: Coin Flipping

- Suppose we have a coin, possibly biased

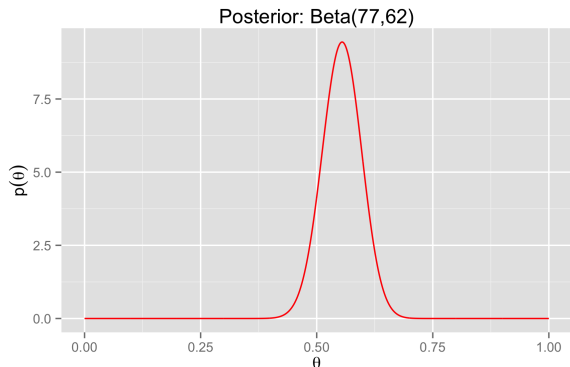
$$\mathbb{P}(\text{Heads} \mid \theta) = \theta.$$

- **Parameter space**  $\theta \in \Theta = [0, 1]$ .
- **Prior distribution:**  $\theta \sim \text{Beta}(2, 2)$ .



## Example: Coin Flipping

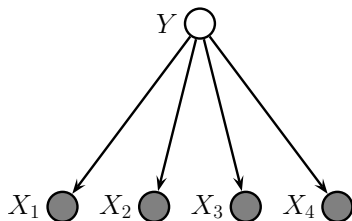
- Next, we gather some data  $\mathcal{D} = \{H, H, T, T, T, T, T, H, \dots, T\}$ :
- Heads: 75      Tails: 60
  - $\hat{\theta}_{\text{MLE}} = \frac{75}{75+60} \approx 0.556$
- **Posterior distribution:**  $\theta \mid \mathcal{D} \sim \text{Beta}(77, 62)$ :





# Naive Bayes: A Generative Model for Classification

- $\mathcal{X} = \left\{ (X_1, X_2, X_3, X_4) \in \{0, 1\}^4 \right\}$       $\mathcal{Y} = \{0, 1\}$  be a class label.
- Consider the Bayesian network depicted below:



- BN structure implies joint distribution factors as:

$$p(x_1, x_2, x_3, x_4, y) = p(y)p(x_1 | y)p(x_2 | y)p(x_3 | y)p(x_4 | y)$$

- Features  $X_1, \dots, X_4$  are independent given the class label  $Y$ .

# Example: Message Classification

- $\mathcal{X} = \{\text{Message Text}\}$
- $\mathcal{Y} = \{\text{BUSINESS, PERSONAL}\}$
- **Training Data**
  - BUSINESS
    - "Lunch meeting?"
    - "Expenses submitted EOM."
    - "LOL"
  - PERSONAL
    - "Meet for lunch? EOM"
    - "LOL"

# Bag of Words Representation (Bernoulli Version)

- Represent a message by the set of words it contains:
  - ignores word order
  - ignores word count (some bag of words models keep the count)
  - typically ignores punctuation and capitalization
- Generate vocabulary from training data:

$W = \{\text{eom}, \text{expenses}, \text{for}, \text{lol}, \text{lunch}, \text{meet}, \text{meeting}, \text{submitted}, \text{UNKNOWN}\}$

- Add in an UNKNOWN value, in case we encounter new words in deployment.
- Message  $M$  is represented by binary vector of length  $|W| = 9$ .

# Bag of Words Representation (Bernoulli Version)

- Input: “Lunch? EOM”  $\implies M = \{\text{lunch, eom}\}$ :
- Vector representation:  $x = (x_1 \dots, x_{|W|})$

Word ( $w$ )	$x_w$
lunch	1
meeting	0
expenses	0
submitted	0
eom	1
meet	0
for	0
lol	0
UNKNOWN	0

# Bernoulli Naive Bayes Model

- Joint probability of message  $x = (x_1, \dots, x_{|W|})$  and class  $y$  is

$$p(x, y) = p(y) \prod_{i=1}^{|W|} p(x_i | y),$$

where each  $x_i \in \{0, 1\}$ , and  $y \in \{B, P\}$ .

- We need to estimate:

$$\mathbb{P}(Y = B)$$

$$\mathbb{P}(Y = P)$$

$$\mathbb{P}(X_w = 1 | Y = B) \forall w \in W$$

$$\mathbb{P}(X_w = 1 | Y = P) \forall w \in W$$

# Bernoulli Naive Bayes: Parameter Estimation

- Using relative frequencies in training, we have:

$$\hat{p}(Y = B) = 3/5 \quad \hat{p}(Y = P) = 2/5$$

and

Word ( $w$ )	$\hat{p}(X_w = 1   B)$	$\hat{p}(X_w = 1   P)$
lunch	1/3	1/2
meeting	1/3	0
expenses	1/3	0
submitted	1/3	0
eom	1/3	1/2
meet	0	1/2
for	0	1/2
lol	1/3	1/2
UNKNOWN	0	0

## Naive Bayes Prediction for "Lunch? EOM"

Word ( $w$ )	$x_w$	$\hat{p}(X_w = 1   B)$	$\hat{p}(x_w   B)$	$\hat{p}(X_w = 1   P)$	$\hat{p}(x_w   P)$
lunch	1	1/3	<b>1/3</b>	1/2	<b>1/2</b>
meeting	0	1/3	<b>2/3</b>	0	<b>1</b>
expenses	0	1/3	<b>2/3</b>	0	<b>1</b>
submitted	0	1/3	<b>2/3</b>	0	<b>1</b>
eom	1	1/3	<b>1/3</b>	1/2	<b>1/2</b>
meet	0	0	<b>1</b>	1/2	<b>1/2</b>
for	0	0	<b>1</b>	1/2	<b>1/2</b>
lol	0	1/3	<b>2/3</b>	1/2	<b>1/2</b>
UNKNOWN	0	0	<b>1</b>	0	<b>1</b>

$$p(M | B) = \frac{1}{3} \cdot \frac{2}{3} \cdot \frac{2}{3} \cdot \frac{2}{3} \cdot \frac{1}{3} \cdot 1 \cdot 1 \cdot \frac{2}{3} \cdot 1 = \frac{16}{243} \approx .07$$

$$p(M | P) = \frac{1}{2} \cdot 1 \cdot 1 \cdot 1 \cdot \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} \cdot 1 = \frac{1}{32} = .03$$

# Naive Bayes Prediction for “Lunch? EOM”

- Input: “Lunch? EOM”  $\implies M = \{\text{lunch, eom}\}$
- Message probability, conditional on message type:

$$p(M | B) = \frac{1}{3} \cdot \frac{2}{3} \cdot \frac{2}{3} \cdot \frac{2}{3} \cdot \frac{1}{3} \cdot 1 \cdot 1 \cdot \frac{2}{3} \cdot 1 = \frac{16}{243} \approx .07$$

$$p(M | P) = \frac{1}{2} \cdot 1 \cdot 1 \cdot 1 \cdot \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} \cdot 1 = \frac{1}{32} = .03$$

- What does it mean that  $p(M | P) = .03$ ?
  - 3% of personal messages have same bag of words as  $M$ .



# Naive Bayes Prediction

- Input: “Lunch? EOM”  $\implies M = \{\text{lunch, eom}\}$
- Output:

$$\begin{aligned}
 p(\text{BUSINESS} \mid M) &\propto p(B)p(M \mid B) \\
 &= \frac{3}{5} \cdot \frac{16}{243} = \frac{16}{405} \\
 p(\text{PERSONAL} \mid M) &\propto p(P)p(M \mid P) \\
 &= \frac{2}{5} \cdot \frac{1}{32} = \frac{1}{90}
 \end{aligned}$$

- Now just renormalize:

$$\begin{aligned}
 p(\text{BUSINESS} \mid M) &= \frac{16}{405} / \left( \frac{1}{90} + \frac{16}{405} \right) \approx 0.78 \\
 p(\text{PERSONAL} \mid M) &= \frac{1}{90} / \left( \frac{1}{90} + \frac{16}{405} \right) \approx 0.22
 \end{aligned}$$

# Naive Bayes Prediction: Issue With Zeros

- Input:  $M = \text{"Meeting?"}$
- Output:

$$p(\text{BUSINESS} \mid M) \propto \frac{1}{3}$$

$$p(\text{PERSONAL} \mid M) \propto 0$$

- Renormalizing:

$$p(\text{BUSINESS} \mid M) = 1$$

$$p(\text{PERSONAL} \mid M) = 0$$

- This is bad:
  - Never want to predict probability 0 if something is possible.
- Worse: Zero counts common for small sample sizes and rare features.

# Laplace Smoothing

- **Laplace Smoothing** is a traditional fix to the 0 count issue.
- Idea is to add 1 to every empirical count:

$$\hat{p}(\text{lunch} \mid \text{PERSONAL}) = \frac{1 + \sum 1(\text{lunch and PERSONAL})}{1 + \sum 1(\text{PERSONAL})}$$

- The added 1 is called a **pseudocount**.
- Like assuming every outcome that can occur was observed at least once.
- Seems to solve the problem – but is there a more principled approach?

# Bayesian Naive Bayes

- Be **Bayesian** and put a beta prior on each parameter.
- **Option 1:** Use posterior mean as point estimate for each parameter, then continue as before.
  - Laplace smoothing is a special case, in which priors are all  $\text{Beta}(1, 1)$ .
- **Option 2: Go full Bayesian.**
  - No parameter estimates. Base everything on posterior  $\theta \mid \mathcal{D}$ .
- Predict with the predictive distribution:

$$y \mid x, \mathcal{D}$$

- Recall, this is integrating out the parameter  $\theta$  w.r.t. the posterior distribution.