

# Applied Text Mining in Python

## *Naïve Bayes Classifier*

# Case study: Classifying text search queries


- Suppose you are interested in classifying search queries in three classes: Entertainment, Computer Science, Zoology
- Most common class of the three is Entertainment.

# Case study: Classifying text search queries

- Suppose the query is “Python”
  - Python, the snake (Zoology)
  - Python, the programming language (Computer Science)
  - Python, as in Monty Python (Entertainment)
- Most common class, given “Python”, is Zoology.





# Case study: Classifying text search queries

- Suppose the query is “Python download”
- Most probable class, given “Python download”, is Computer Science. 

# Probabilistic Model





- Update the likelihood of the class given new information
- **Prior Probability:**  $\Pr(y = \text{Entertainment})$ ,  $\Pr(y = \text{CS})$ ,  $\Pr(y = \text{Zoology})$  
- **Posterior probability:**  $\Pr(y = \text{Entertainment} | x = \text{“Python”})$  

# Bayes' Rule

- Posterior probability = 
$$\frac{\text{Prior probability} \times \text{Likelihood}}{\text{Evidence}}$$
- $$\Pr(y | X) = \frac{\Pr(y) \times \Pr(X | y)}{\Pr(X)}$$

# Naïve Bayes Classification

- $\Pr(y=\text{CS} | \text{"Python"}) = \frac{\Pr(y=\text{CS}) \times \Pr(\text{"Python"} | y=\text{CS})}{\Pr(\text{"Python"})}$  
- $\Pr(y=\text{Zoology} | \text{"Python"}) = \frac{\Pr(y=\text{Zoology}) \times \Pr(\text{"Python"} | y=\text{Zoology})}{\Pr(\text{"Python"})}$
- $\Pr(y=\text{CS} | \text{"Python"}) > \Pr(y=\text{Zoology} | \text{"Python"}) \Rightarrow$  

# Naïve Bayes Classification

- $\Pr(y | X) = \frac{\Pr(y) \times \Pr(X | y)}{\Pr(X)}$

$$y^* = \underset{y}{\operatorname{argmax}} \Pr(y | X) = \underset{y}{\operatorname{argmax}} \Pr(y) \times \Pr(X | y) \quad \square$$

- **Naïve assumption:** Given the class label, features are assumed to be independent of each other

$$y^* = \underset{y}{\operatorname{argmax}} \Pr(y | X) = \underset{y}{\operatorname{argmax}} \Pr(y) \times \prod_{i=1}^n \Pr(x_i | y) \quad \square$$



# Naïve Bayes Classifier

$$y^* = \underset{y}{\operatorname{argmax}} \operatorname{Pr}(y \mid X) = \underset{y}{\operatorname{argmax}} \operatorname{Pr}(y) \times \prod_{i=1}^n \operatorname{Pr}(x_i \mid y)$$

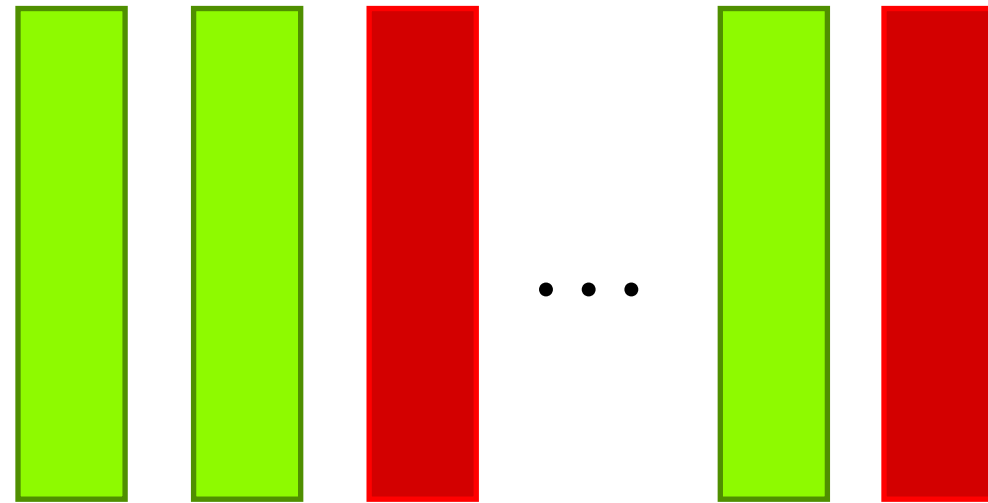

Query: “Python download”

$$y^* = \underset{y}{\operatorname{argmax}} \operatorname{Pr}(y) \times \operatorname{Pr}(\text{“Python”} \mid y) \times \operatorname{Pr}(\text{“download”} \mid y)$$

# Naïve Bayes: What are the parameters?

- Prior probabilities:  $\Pr(y)$  for all  $y$  in  $Y$
- Likelihood:  $\Pr(x_i | y)$  for all features  $x_i$  and labels  $y$  in  $Y$
- If there are 3 classes ( $|Y| = 3$ ) and 100 features in  $X$ , how many parameters does naïve Bayes models have?

# Naïve Bayes: Learning parameters

- **Prior probabilities:**  $\Pr(y)$  for all  $y$  in  $Y$
- Remember training data?
- Count the number of instances in each class
- If there are  $N$  instances in all, and  $n$  out of those are labeled as class  $y$   $\Rightarrow \Pr$  

# Naïve Bayes: Learning parameters

- **Likelihood:**  $\Pr(\mathbf{x}_i \mid \mathbf{y})$  for all features  $\mathbf{x}_i$  and labels  $\mathbf{y}$  in  $\mathbf{Y}$
- Count how many times feature  $\mathbf{x}_i$  appears in instances labeled as class  $\mathbf{y}$
- If there are  $p$  instances of class  $\mathbf{y}$ , and  $\mathbf{x}_i$  appears in  $k$  of those,  $\Pr(\mathbf{x}_i \mid \mathbf{y}) = k / p$

# Naïve Bayes: Smoothing

- What happens if  $\Pr(\mathbf{x}_i \mid \mathbf{y}) = 0$ ?
- Feature  $\mathbf{x}_i$  never occurs in documents labeled  $\mathbf{y}$
- But then, the posterior probability  $\Pr(\mathbf{y} \mid \mathbf{x}_i)$  will be 0!!
- Instead, smooth the parameters
- Laplace smoothing or Additive smoothing: Add a dummy count
- $\Pr(\mathbf{x}_i \mid \mathbf{y}) = (k+1) / (p+n)$ ; where  $n$  is number of features

# Take Home Concepts

- Naïve Bayes is a probabilistic model
- Naïve, because it assumes features are independent of each other, given the class label
- For text classification problems, naïve Bayes models typically provide very strong baselines
- Simple model, easy to learn parameters