

Dirichlet-Multinomial and Naive Bayes

Instructor: Alan Ritter

Many Slides from Pedro Domingos

Last time: Beta-Binomial

- Binary random variable: bent coin

Data Likelihood:

$$P(x_1, x_2, \dots, x_n | \theta_H) = \theta_H^{\#H} (1 - \theta_H)^{\#T}$$

Last time: Beta-Binomial

- Binary random variable: bent coin

Data Likelihood:

$$P(x_1, x_2, \dots, x_n | \theta_H) = \theta_H^{\#H} (1 - \theta_H)^{\#T}$$



Last time: Beta-Binomial

- Binary random variable: bent coin

Data Likelihood:

$$P(x_1, x_2, \dots, x_n | \theta_H) = \theta_H^{\#H} (1 - \theta_H)^{\#T}$$

Prior (Beta distribution):

$$P(\theta_H | \alpha, \beta) = \frac{1}{B(\alpha, \beta)} \theta_H^{\alpha-1} (1 - \theta_H)^{\beta-1}$$



Last time: Beta-Binomial

- Binary random variable: bent coin

Data Likelihood:

$$P(x_1, x_2, \dots, x_n | \theta_H) = \theta_H^{\#H} (1 - \theta_H)^{\#T}$$

Prior (Beta distribution):

$$P(\theta_H | \alpha, \beta) = \frac{1}{B(\alpha, \beta)} \theta_H^{\alpha-1} (1 - \theta_H)^{\beta-1}$$

Posterior:

$$P(\theta_H | \alpha, \beta, x_1, \dots, x_n) = \frac{1}{B(\alpha + \#H, \beta + \#T)} \theta_H^{\#H + \alpha - 1} (1 - \theta_H)^{\#T + \beta - 1}$$



Last time: Beta-Binomial

- Binary random variable: bent coin

Maximum Likelihood:

$$\theta^{ML} = \frac{\#H}{\#T + \#H}$$

Last time: Beta-Binomial

- Binary random variable: bent coin

Maximum Likelihood:

$$\theta^{ML} = \frac{\#H}{\#T + \#H}$$



Last time: Beta-Binomial

- Binary random variable: bent coin

Maximum Likelihood:

$$\theta^{ML} = \frac{\#H}{\#T + \#H}$$



Maximum a Posteriori:

$$\theta^{MAP} = \frac{\#H + \alpha - 1}{\#T + \#H + \alpha + \beta - 2}$$

K-Sided Dice



- Weighted
 - (Generalization of Bent Coin)
- Assume an observed sequence of rolls:

1123213213

$$\theta_1$$

$$\theta_2$$

$$\theta_3$$

K-Sided Dice



- Weighted
 - (Generalization of Bent Coin)
- Assume an observed sequence of rolls:

1123213213

$$\theta_1 \quad \theta_2 \quad \theta_3$$

$$P(x; \theta) = \theta_x$$

Likelihood

$$P(1123213213|\theta) = \theta_1 \times \theta_1 \times \theta_2 \times \dots \times \theta_3$$

$$= \theta_1^4 \times \theta_2^3 \times \theta_3^3$$

Likelihood In General

- N Dice Rolls, K possible outcomes:

$$P(D|\theta) = \prod_{k=1}^K \theta_k^{N_k}$$

Likelihood In General

- N Dice Rolls, K possible outcomes:

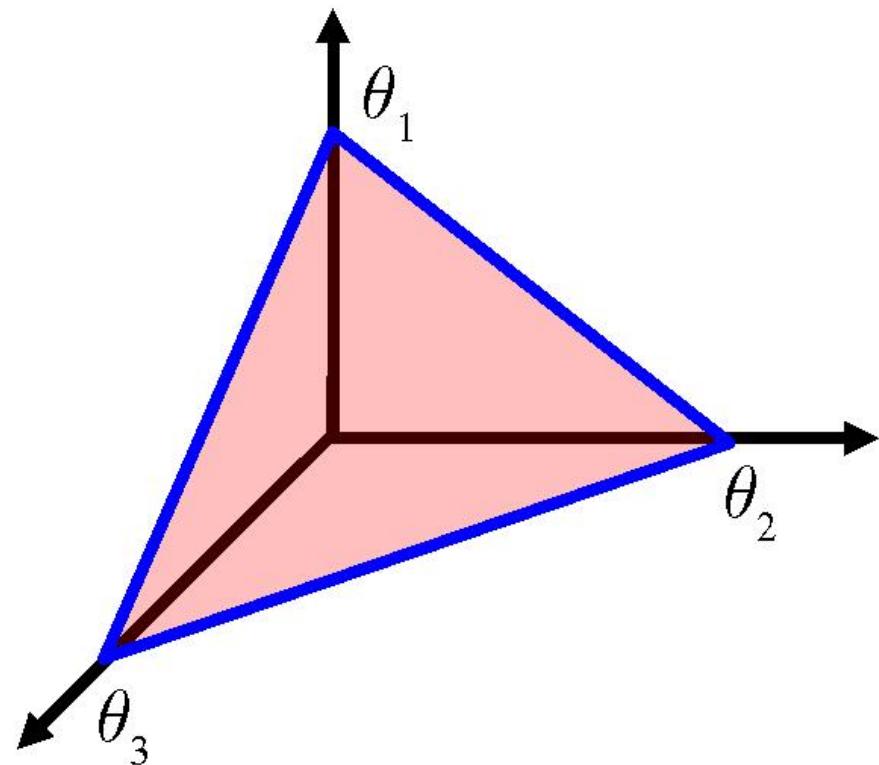
$$P(D|\theta) = \prod_{k=1}^K \theta_k^{N_k}$$

- Likelihood is a multivariable function

$$= f(\theta_1, \theta_2, \dots, \theta_K)$$

3D Probability Simplex

- 3 parameters
- Constraint that parameters sum to 1

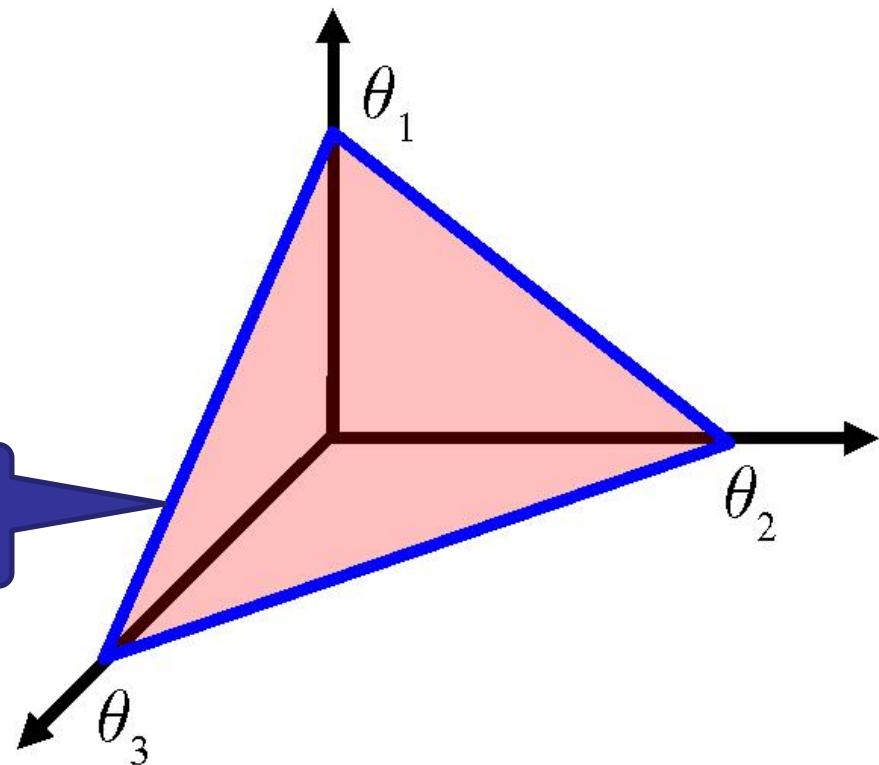


$$S_K = \{\theta : 0 \leq \theta_k \leq 1, \sum_{k=1}^K \theta_k = 1\}$$

3D Probability Simplex

- 3 parameters
- Constraint that parameters sum to 1

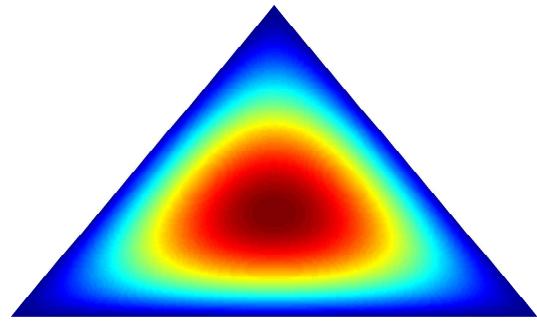
We want a probability distribution over this



$$S_K = \{\theta : 0 \leq \theta_k \leq 1, \sum_{k=1}^K \theta_k = 1\}$$

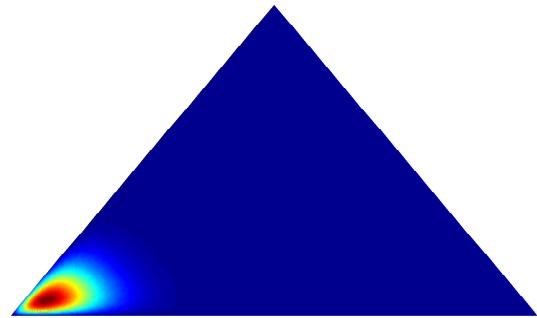
Dirichlet distribution

- Multivariate generalization of Beta distribution
- Conjugate prior to multinomial

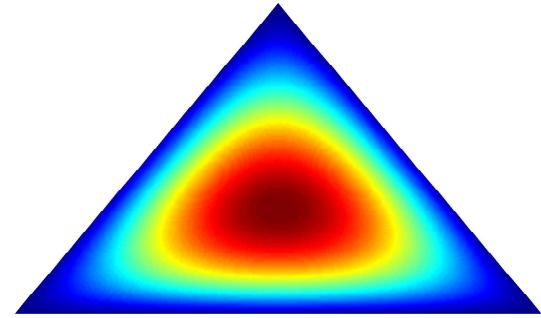


$$\text{Dir}(\theta|\alpha) = \frac{1}{B(\alpha)} \prod_{k=1}^K \theta_k^{\alpha_k - 1} \mathbb{1}(\theta \in S_K)$$

Dirichlet distribution



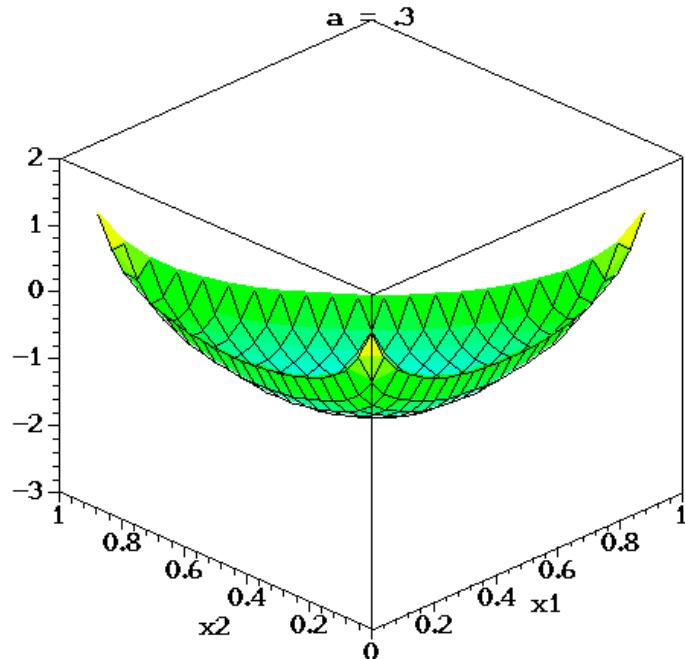
$\alpha = <20, 2, 2>$



$\alpha = <2, 2, 2>$

$$\text{Dir}(\theta|\alpha) = \frac{1}{B(\alpha)} \prod_{k=1}^K \theta_k^{\alpha_k - 1} \mathbb{1}(\theta \in S_K)$$

(log) Dirichlet distribution



$\alpha = <0.3, 0.3, 0.3>$ to $<2.0, 2.0, 2.0>$

Posterior



$$P(\theta|D) \propto P(D|\theta)P(\theta)$$

Posterior

$$P(\theta|D) \propto P(D|\theta)P(\theta)$$



Posterior



$$P(\theta|D) \propto P(D|\theta)P(\theta)$$

$$\propto \prod_{k=1}^K \theta_k^{N_k} \theta_k^{\alpha_k - 1} = \prod_{k=1}^K \theta_k^{N_k + \alpha_k - 1}$$

Posterior



$$P(\theta|D) \propto P(D|\theta)P(\theta)$$

$$\begin{aligned} & \propto \prod_{k=1}^K \theta_k^{N_k} \theta_k^{\alpha_k - 1} = \prod_{k=1}^K \theta_k^{N_k + \alpha_k - 1} \\ & = \text{Dir}(\theta | \alpha_1 + N_1, \dots, \alpha_K + N_K) \end{aligned}$$

Posterior



$$P(\theta|D) \propto P(D|\theta)P(\theta)$$

$$\begin{aligned} & \propto \prod_{k=1}^K \theta_k^{N_k} \theta_k^{\alpha_k - 1} = \prod_{k=1}^K \theta_k^{N_k + \alpha_k - 1} \\ & = \text{Dir}(\theta | \alpha_1 + N_1, \dots, \alpha_K + N_K) \end{aligned}$$

Dirichlet is
Conjugate to
Multinomial

MAP Point Estimate

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

MAP Point Estimate

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

$$= \frac{N_k + \alpha_k - 1}{\sum_{k=1}^K N_k + \sum_{k=1}^K \alpha_k - K}$$

Naive Bayes Classifier

Naive Bayes Classifier

Assume target function $f : X \rightarrow V$, where each instance x described by attributes $\langle a_1, a_2 \dots a_n \rangle$.

Most probable value of $f(x)$ is:

$$\begin{aligned} v_{MAP} &= \operatorname{argmax}_{v_j \in V} P(v_j | a_1, a_2 \dots a_n) \\ v_{MAP} &= \operatorname{argmax}_{v_j \in V} \frac{P(a_1, a_2 \dots a_n | v_j) P(v_j)}{P(a_1, a_2 \dots a_n)} \\ &= \operatorname{argmax}_{v_j \in V} P(a_1, a_2 \dots a_n | v_j) P(v_j) \end{aligned}$$

Naive Bayes assumption:

$$P(a_1, a_2 \dots a_n | v_j) = \prod_i P(a_i | v_j)$$

which gives

Naive Bayes classifier:

$$v_{NB} = \operatorname{argmax}_{v_j \in V} P(v_j) \prod_i P(a_i | v_j)$$

Naive Bayes Algorithm

Naive_Bayes_Learn(*examples*)

For each target value v_j

- $\hat{P}(v_j) \leftarrow$ estimate $P(v_j)$
- For each attribute value a_i of each attribute a
 $\hat{P}(a_i|v_j) \leftarrow$ estimate $P(a_i|v_j)$

Classify_New_Instance(x)

$$v_{NB} = \operatorname{argmax}_{v_j \in V} \hat{P}(v_j) \prod_{a_i \in x} \hat{P}(a_i|v_j)$$

Naive Bayes: Example

Consider *PlayTennis* again, and new instance

$$\langle Outlk = sun, Temp = cool, Humid = high, Wind = strong \rangle$$

Want to compute:

$$v_{NB} = \operatorname{argmax}_{v_j \in V} P(v_j) \prod_i P(a_i | v_j)$$

$$P(y) P(sun|y) P(cool|y) P(high|y) P(strong|y) = .005$$

$$P(n) P(sun|n) P(cool|n) P(high|n) P(strong|n) = .021$$

$$\rightarrow v_{NB} = n$$

Naive Bayes: Subtleties (1)

Conditional independence assumption is often violated

$$P(a_1, a_2 \dots a_n | v_j) = \prod_i P(a_i | v_j)$$

... but it works surprisingly well anyway. Don't need estimated posteriors $\hat{P}(v_j | x)$ to be correct; need only that

$$\operatorname{argmax}_{v_j \in V} \hat{P}(v_j) \prod_i \hat{P}(a_i | v_j) = \operatorname{argmax}_{v_j \in V} P(v_j) P(a_1 \dots, a_n | v_j)$$

Naive Bayes posteriors often unrealistically close to 1 or 0

Naive Bayes: Subtleties (2)

What if none of the training instances with target value v_j have attribute value a_i ? Then

$$\hat{P}(a_i|v_j) = 0, \text{ and } \dots$$

$$\hat{P}(v_j) \prod_i \hat{P}(a_i|v_j) = 0$$

MAP Point Estimate

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

MAP Point Estimate

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

$$= \frac{N_k + \alpha_k - 1}{\sum_{k=1}^K N_k + \sum_{k=1}^K \alpha_k - K}$$

Naïve Bayes with Log Probabilities

$$\begin{aligned} c_{MAP} &= \operatorname{argmax}_c P(c|x_1, \dots, x_n) \\ &= \operatorname{argmax}_c P(c) \prod_{i=1}^n P(x_i|c) \\ &= \operatorname{argmax}_c \log \left(P(c) \prod_{i=1}^n P(x_i|c) \right) \\ &= \operatorname{argmax}_c \log P(c) + \sum_{i=1}^n \log P(x_i|c) \end{aligned}$$

Naïve Bayes with Log Probabilities

$$c_{MAP} = \operatorname{argmax}_c \log P(c) + \sum_{i=1}^n \log P(x_i|c)$$

Naïve Bayes with Log Probabilities

$$c_{MAP} = \operatorname{argmax}_c \log P(c) + \sum_{i=1}^n \log P(x_i|c)$$

- Q: Why don't we have to worry about floating point underflow anymore?

Learning to Classify Text

Why?

- Learn which news articles are of interest
- Learn to classify web pages by topic

Naive Bayes is among most effective algorithms

What attributes shall we use to represent text documents?

Learning to Classify Text

Target concept $\text{Interesting?} : \text{Document} \rightarrow \{+, -\}$

1. Represent each document by vector of words:
one attribute per word position in document
2. Learning: Use training examples to estimate
 - $P(+)$
 - $P(-)$
 - $P(\text{doc}|+)$
 - $P(\text{doc}|-)$

Naive Bayes conditional independence assumption

$$P(doc|v_j) = \prod_{i=1}^{length(doc)} P(a_i = w_k | v_j)$$

where $P(a_i = w_k | v_j)$ is probability that word in position i is w_k , given v_j

One more assumption:

$$P(a_i = w_k | v_j) = P(a_m = w_k | v_j), \forall i, m$$

`LEARN_NAIVE_BAYES_TEXT(Examples, V)`

1. *Collect all words & tokens that occur in Examples*
 - $\text{Vocabulary} \leftarrow$ all distinct words & tokens in *Examples*
2. *Compute all probabilities $P(v_j)$ and $P(w_k|v_j)$*
 - For each target value v_j in *V* do
 - $\text{docs}_j \leftarrow \text{Examples}$ for which the target value is v_j
 - $P(v_j) \leftarrow \frac{|\text{docs}_j|}{|\text{Examples}|}$
 - $\text{Text}_j \leftarrow$ concatenate all members of docs_j
 - $n \leftarrow$ total number of words in Text_j (counting duplicate words multiple times)
 - for each word w_k in *Vocabulary*
 - * $n_k \leftarrow$ number of times word w_k occurs in Text_j
 - * $P(w_k|v_j) \leftarrow \frac{n_k+1}{n+|\text{Vocabulary}|}$

`CLASSIFY_NAIVE_BAYES_TEXT(Doc)`

- $positions \leftarrow$ all word positions in *Doc* that contain tokens found in *Vocabulary*
- Return v_{NB} , where

$$v_{NB} = \operatorname{argmax}_{v_j \in V} P(v_j) \prod_{i \in positions} P(a_i | v_j)$$

Example: 20 Newsgroups

Given 1000 training documents from each group

Learn to classify new documents according to which newsgroup it came from

comp.graphics	misc.forsale
comp.os.ms-windows.misc	rec.autos
comp.sys.ibm.pc.hardware	rec.motorcycles
comp.sys.mac.hardware	rec.sport.baseball
comp.windows.x	rec.sport.hockey
alt.atheism	sci.space
soc.religion.christian	sci.crypt
talk.religion.misc	sci.electronics
talk.politics.mideast	sci.med
talk.politics.misc	talk.politics.guns

Naive Bayes: 89% classification accuracy

Article from rec.sport.hockey

Path: cantaloupe.srv.cs.cmu.edu!...

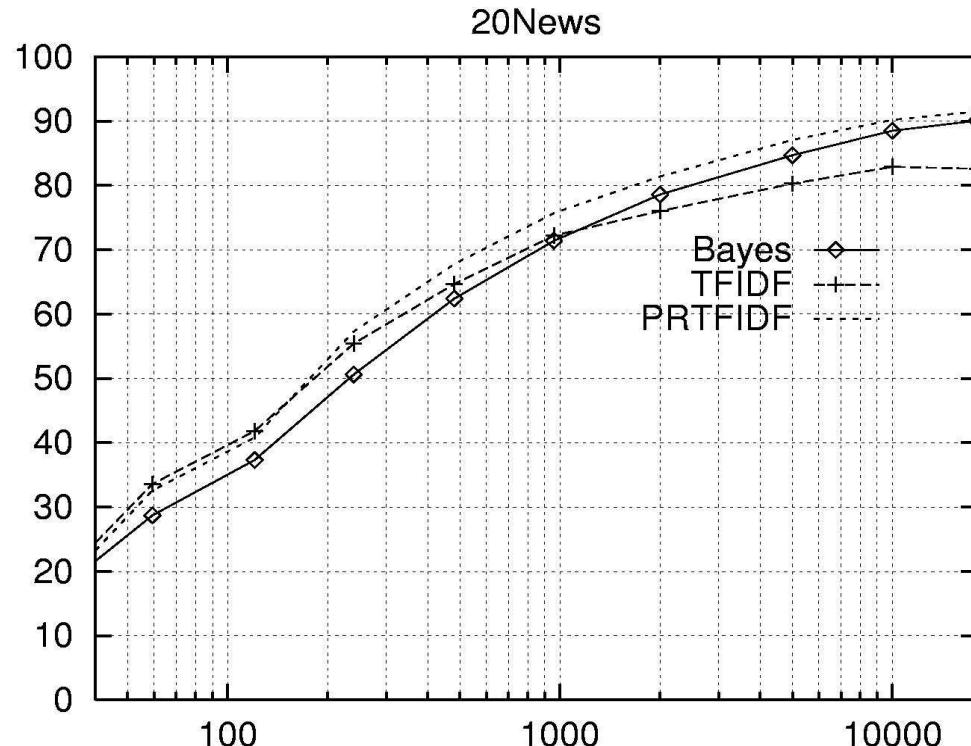
From: xxx@yyy.zzz.edu (John Doe)

Subject: Re: This year's biggest and worst (opinion)

Date: 5 Apr 93 09:53:39 GMT

I can only comment on the Kings, but the most obvious candidate for pleasant surprise is Alex Zhitnik. He came highly touted as a defensive defenseman, but he's clearly much more than that. Great skater and hard shot (though wish he were more accurate). In fact, he pretty much allowed the Kings to trade away that huge defensive liability Paul Coffey. Kelly Hrudey is only the biggest disappointment if you thought he was any good to begin with. But, at best, he's only a mediocre goaltender. A better choice would be Tomas Sandstrom, though not through any fault of his own, but because some thugs in Toronto decided

Learning Curve for 20 Newsgroups



Accuracy vs. Training set size (1/3 withheld for test)