

Learning from Data

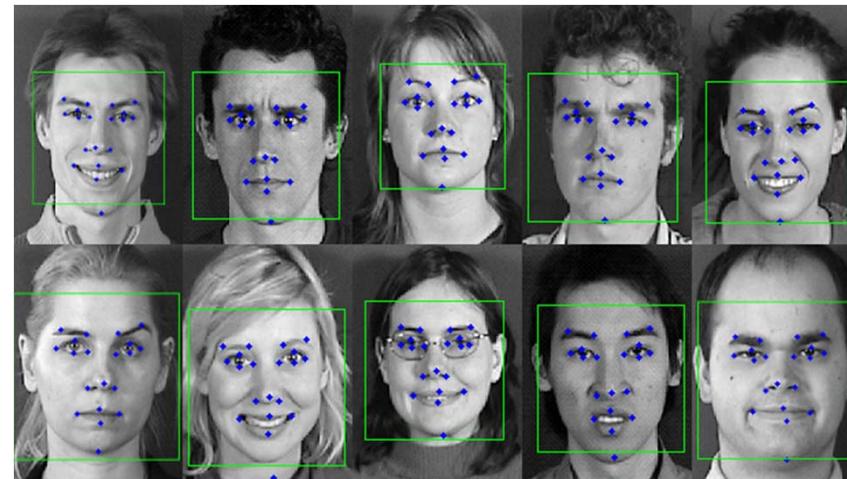
6. Classification: Naïve Bayes

© Prakash Ishwar

Spring 2017

Classification

- Supervised (predictive) learning: given examples with labels, predict labels for all unseen examples
 - Classification:
 - label = category,
 - $y \in \mathcal{Y} = \{1, \dots, m\}$, m = number of classes
 - $\ell(\mathbf{x}, y, h) = 1(h(\mathbf{x}) \neq y)$, Risk = $P(Y \neq h(\mathbf{X})) = P(\text{Error})$



\mathbf{x} = facial geometry features
 y = gender label

Naïve Bayes

- Generative learning
- d features (components) in each feature vector:

$$\mathbf{x} = (x_1, \dots, x_i, \dots, x_d)^\top$$

- **Naïve Bayes Assumption:** all features are **conditionally independent** given the class label

$$\forall y, p(\mathbf{x}|y, \theta) = p(x_1|y, \theta_1) \cdots p(x_d|y, \theta_d) = \prod_{i=1}^d p(x_i|y, \theta_i)$$

$$\theta = (\theta_0, \theta_1, \dots, \theta_d)^\top, \quad \theta_0 = (p(y=1), \dots, p(y=m))^\top$$

- Resulting MPE decision rule: Naïve Bayes classifier
- Why Naïve? We do not expect all features to be independent even conditional on class label

Naïve Bayes

Why Naïve Bayes?

- Low Complexity

$\begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix} \xrightarrow{\text{Binary feature } (0/1)}$

$P(x_1, \dots, x_d | y)$ $\therefore 2^d$ distinct feature vectors

$\begin{matrix} \text{e.g. } d=2 \\ (0)(0)(1)(1) \end{matrix}$ (4) 

$\begin{matrix} \text{e.g. } d=3 \\ (0)(0)(1)(1)(1)(1) \end{matrix}$ (8) 

$2^d - 1$ parameters to specify the joint pmf

of params needed is growing exponentially fast

- Although assumption is false, the resulting classifier works quite well in practice.
- It has only on the order of $O(md)$ parameters to learn → relatively immune to overfitting
 - Example: Suppose that all features are binary, say taking only values +1 or -1, then:
 - In NB model, we need to specify only 1 scalar parameter per feature per class, namely $P(X_i = +1|y), y = 1, \dots, m$ or md scalar parameters in total
 - But specifying the **joint** pmf of all d binary features requires $(2^d - 1)$ scalar parameters per class or $m(2^d - 1)$ scalar parameters in total

Naïve Bayes

- Class-conditional likelihood functions:

$$p(x_i|y, \theta_i), i = 1, \dots, d$$

- may be either completely parametric or completely non-parametric or some components may be modeled parametrically and others non-parametrically
- some or all feature components may be discrete, continuous, or mixed
- Gaussian Naïve Bayes: $\forall i, y, p(x_i|y, \theta_i) \sim \mathcal{N}(\mu_{iy}, \sigma_{iy}^2)$
- Categorical Naïve Bayes: all features are discrete and take only a finite number of possible values

Naïve Bayes

- Notation summary:

y = class label $\in \{1, \dots, m\}$

j = sample index $\in \{1, \dots, n\}$

i = feature (component) index $\in \{1, \dots, d\}$

θ_i = parameters for class-conditional pdfs/pdfs of feature i across all classes

$\theta_i = \{\theta_{iy}, y = 1, \dots, m\}$, where

θ_{iy} = parameters for class-conditional pdf/pmf of feature i in class y , and
 $p(x_i|y, \theta_i) = p(x_i|y, \theta_{iy})$ for all i, y

$\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ (training set)

$$\begin{array}{c}
 \text{feature index} \rightarrow \\
 \begin{matrix}
 i = 1 & \mathbf{x}_1 = \begin{pmatrix} x_{11} \\ \vdots \\ x_{i1} \\ \vdots \\ x_{d1} \end{pmatrix}, & \dots & \mathbf{x}_j = \begin{pmatrix} x_{1j} \\ \vdots \\ x_{ij} \\ \vdots \\ x_{dj} \end{pmatrix}, & \dots & \mathbf{x}_n = \begin{pmatrix} x_{1n} \\ \vdots \\ x_{in} \\ \vdots \\ x_{dn} \end{pmatrix} \\
 j = 1, & \dots & j, & \dots & j = n
 \end{matrix}
 \end{array}$$

sample index \longrightarrow

Naïve Bayes

- Training set: $\mathcal{D} = \{(\mathbf{x}_j, y_j), j = 1, \dots, n\}$
- Training set for feature i :
$$\mathcal{D}_i = \{(x_{ij}, y_j), j = 1, \dots, n\}$$
- Training set for feature i and class y :
$$\mathcal{D}_{iy} = \{(x_{ij}, y_j), j : y_j = y\}$$
- **Key Result:** for each feature i and each class y , the ML estimate of θ_{iy} based on \mathcal{D} is the same as the ML estimate of θ_{iy} based on \mathcal{D}_{iy}
- **Intuition:** features independent in every class \Rightarrow parameters of feature distribution can be estimated independently for each feature in each class

$$abb \alpha \ll \alpha^3 b^3 c^2 = \frac{\pi}{\alpha b c} \sqrt[3]{abc}$$

Naïve Bayes

Proof:

$$p(\mathcal{D}|\theta) = \prod_{j=1}^n p(\mathbf{x}_j^{\mathcal{D}}, y_j | \theta)$$

$$\theta()^n \ln p() \ln np()$$

$$= \prod_{j=1}^n \left[p(y_j) \prod_{i=1}^d p(x_{ij} | y_j, \theta_i) \right]$$

Naïve Bayes assumption:

$$p(\mathbf{x}_j | y_j) = \prod_{i=1}^d p(x_{ij} | y_j, \theta_i)$$

$$= \left(\prod_{j=1}^n p(y_j) \right) \prod_{i=1}^d \prod_{j=1}^n p(x_{ij} | y_j, \theta_i)$$

$$= \left(\prod_{y=1}^m p(y)^{n_y} \right) \prod_{i=1}^d \prod_{y=1}^m \prod_{j:y_j=y} p(x_{ij} | y, \theta_{iy})$$

$$n_y = \# \text{ class } y \text{ samples} \\ = \sum_{j=1}^n \mathbb{1}(y_j = y)$$

$$\Rightarrow \ln p(\mathcal{D}|\theta) = \sum_{y=1}^m n_y \ln p(y) + \sum_{i=1}^d \sum_{y=1}^m \sum_{j:y_j=y} \ln p(x_{ij} | y, \theta_{iy})$$

$$\Rightarrow \frac{1}{n} \ln p(\mathcal{D}|\theta) = \sum_{y=1}^m \left(\frac{n_y}{n} \right) \ln \theta_{0y} + \sum_{i=1}^d \sum_{y=1}^m \left(\frac{n_y}{n} \right) \left[\frac{1}{n_y} \sum_{j:y_j=y} \ln p(x_{ij} | y, \theta_{iy}) \right]$$

maximize

can ignore other features

Naïve Bayes

Thus,

$$\hat{\theta}_{ML}(\mathcal{D}) = \arg \max_{\theta} \frac{1}{n} \ln p(\mathcal{D}|\theta),$$

$$\theta = (\theta_0, \{\theta_{iy}, i = 1, \dots, d, y = 1, \dots, m\})$$

$$\Rightarrow \hat{\theta}_{0,ML}(\mathcal{D}) = \arg \max_{\theta_0} \left[\sum_{y=1}^m \left(\frac{n_y}{n} \right) \ln \theta_{0y} \right],$$

$$\begin{aligned} \forall i, y, \hat{\theta}_{iy,ML}(\mathcal{D}) &= \arg \max_{\theta_{iy}} \left[\frac{1}{n_y} \sum_{j:y_j=y} \ln p(x_{ij}|y, \theta_{iy}) \right] \\ &= \hat{\theta}_{iy,ML}(\mathcal{D}_{iy}) \end{aligned}$$

$$p(x_{ij}|y, \theta_{iy}) \stackrel{?}{=} \frac{(\# \text{ words in } \mathcal{D} == i)_y (\text{counts of } i)_y}{\text{Number of total words in class } Y}$$

Categorical Naïve Bayes

- All features are categorical:

– Without loss of generality (w.l.o.g.): $\forall i, x_i \in \{1, \dots, W\}$

- Notation summary:

$$y = \text{class label} \in \{1, \dots, m\}$$

$$j = \text{sample index} \in \{1, \dots, n\}$$

$$\stackrel{?}{=} i = \text{feature (component) index} \in \{1, \dots, d\}$$

$$w = \text{feature value} \in \{1, \dots, W\}$$

- Model:

$$\forall y, j, i, w, P(X_{ij} = w | Y_j = y) = \beta_{w,y,i}, \text{ where}$$

$$\forall y, i, w, \beta_{w,y,i} \geq 0, \text{ and } \forall y, i, \sum_{w=1}^W \beta_{w,y,i} = 1$$

- Parameters:

$$\theta = \underbrace{\{p(y)\}}_{\theta_0 y}, \underbrace{\{\beta_{w,y,i}, w = 1, \dots, W\}}_{\theta_{iy} = \{\beta_{w,y,i}, w = 1, \dots, W\}}, y = 1, \dots, m, i = 1, \dots, d\}$$

- Total number of scalar parameters = $m + mdW$

Categorical Naïve Bayes

- Bayes classifier for 0-1 loss (MPE rule = MAP rule):

$$h_{\text{MPE}}(\mathbf{x}) = \arg \max_{y=1, \dots, m} p(y)p(\mathbf{x}|y, \theta)$$

$$= \arg \max_{y=1, \dots, m} p(y) \prod_{i=1}^d \beta_{x_i, y, i} \quad \stackrel{\text{def}}{=} \sum_{i=1}^d \ln \beta_{x_i, y, i}$$

- ML (frequentist) estimate of θ : Let

- n_y = number of class y training examples

- $n_{w,y,i}$ = number of class y training examples in which the i -th feature value = w

- then, $\sum_{y=1}^m n_y = n$ and for all i , $\sum_{w=1}^W n_{w,y,i} = n_y$

- Also, let $n_{w,y} = \sum_{i=1}^d n_{w,y,i}$ = number of times w occurs across all d components in class y examples
then, $\sum_{w=1}^W n_{w,y} = dn_y$

Categorical Naïve Bayes

- ML (frequentist) estimate of θ :

$$\hat{\theta}_{ML}(\mathcal{D}) = \arg \max_{\theta} \frac{1}{n} \ln p(\mathcal{D}|\theta),$$

$$\theta = (\theta_0, \{\theta_{iy}, i = 1, \dots, d, y = 1, \dots, m\})$$

$$\hat{\theta}_{0,ML}(\mathcal{D}) = \arg \max_{\theta_0} \left[\sum_{y=1}^m \left(\frac{n_y}{n} \right) \ln \theta_{0y} \right],$$

$$\forall i, y, \quad \hat{\theta}_{iy,ML}(\mathcal{D}) = \arg \max_{\theta_{iy}} \left[\frac{1}{n_y} \sum_{j:y_j=y} \ln p(x_{ij}|y, \theta_{iy}) \right]$$

$$= \arg \max_{\theta_{iy}} \left[\frac{1}{n_y} \sum_{j:y_j=y} \ln \beta_{x_{ij}, y, i} \right]$$

$$= \arg \max_{\theta_{iy}} \left[\sum_{w=1}^W \left(\frac{n_{w,y,i}}{n_y} \right) \ln \beta_{w,y,i} \right]$$

Categorical Naïve Bayes

- ML (frequentist) estimate of θ solution:

$$\forall y, \hat{\theta}_{0y,ML} = \hat{p}(y) = \frac{n_y}{n}$$

$$\forall w, y, i, \hat{\beta}_{w,y,i} = \frac{n_{w,y,i}}{n_y}$$

- If for all w, y, i , $\beta_{w,y,i} = \beta_{w,y}$ then,

$$\forall y, \hat{\theta}_{0y,ML} = \hat{p}(y) = \frac{n_y}{n}$$

$$\begin{aligned}\forall w, y, i, \hat{\beta}_{w,y,i} &= \frac{\sum_{i=1}^d n_{w,y,i}}{\sum_{i=1}^d n_y} \\ &= \frac{n_{w,y}}{d \cdot n_y}\end{aligned}$$

Categorical Naïve Bayes

- Solution based on following result
- **Result:** Let p_1, \dots, p_L , and q_1, \dots, q_L , denote two pmfs over L items, i.e., they are non-negative and sum to one. Then,

$$\arg \max_{q_1, \dots, q_L} \sum_{l=1}^L p_l \ln q_l = \{p_l, l = 1, \dots, L\}$$

- *Proof:* $\forall t > 0, \ln t \leq t - 1$ with equality, if, and only if (iff), $t = 1$. Replacing t with q_l/p_l , multiplying by p_l , and summing over all l , we get

$$\sum_{l=1}^L p_l \ln \left(\frac{q_l}{p_l} \right) \leq \sum_{l=1}^L p_l \left(\frac{q_l}{p_l} - 1 \right) = \sum_{l=1}^L \left(q_l - p_l \right) = 1 - 1 = 0.$$

Alternatively, this follows from the fact that the KL-divergence $D(p||q)$ is non-negative and is zero if, and only if, the two pmfs p, q are identical

Document Classification (2 classes; m=2)

Vocabulary size = $\omega = 3$ e.g. $\{a, b, c\}$

Training Set			
class 1 examples		class 2 examples	
x_1	x_2	c	c
a	b	c	b
b	c	a	c
c	a	b	c
a	b	c	c
c	c	c	a
doc 1	doc 2	$L=5$	$L=3$
length = 7	length = 5	$L=4$	

likelihoods $p(x|y)$

$$\hat{\beta}_{a,1}^1 = \frac{2+1}{7+5} \leftarrow \begin{matrix} \# \text{ a's} \\ \leftarrow \text{total words} \end{matrix} = \frac{3}{12}$$

$$\hat{\beta}_{a,2}^1 = \frac{2}{12}$$

$$\hat{\beta}_{b,1}^1 = \frac{2+2}{12} = \frac{4}{12}$$

$$\hat{\beta}_{b,2}^1 = \frac{2}{12}$$

$$\hat{\beta}_{c,1}^1 = \frac{3+2}{12} = \frac{5}{12}$$

$$\hat{\beta}_{c,2}^1 = \frac{8}{12}$$

Priors

$$\hat{p}(y=1) = \frac{2}{2+3} = \frac{2}{5}, \quad \hat{p}(y=2) = \frac{3}{2+3} = \frac{3}{5}$$



Text document

NB assumption

b	$p(b)$	1 ₁	.	1 ₂
b	$p(b)$	1 ₁	:	1 ₂
a	$p(a)$	1 ₁	:	1 ₂
c	:			
c	:			
a	:			
b	:			
a				
c				
y	d			
c				
c				
		\bullet	\bullet	
		$p(\text{class 1})$	$p(\text{class 2})$	
				\downarrow
				$\hat{p}(y=1) \cdot \hat{p}(\underline{x} y=1)$
				$= \frac{3}{5} \cdot \hat{\beta}_{a,1}^3 \cdot \hat{\beta}_{b,1}^4 \cdot \hat{\beta}_{c,1}^1$
				$= \frac{3}{5} \cdot (\hat{\beta}_{a,1})^3 \cdot (\hat{\beta}_{b,1})^4 \cdot (\hat{\beta}_{c,1})^1$
				\vdots

possibly take log probabilities
avoid underflow

$$\text{argmax}_y \frac{1}{n} \sum \ln p(\underline{x}, y)$$

$$\text{where } p(\underline{x}|y) = \prod_{i=1}^{n+1} p(x_i|y)$$

so don't multiply by 0 if d not in training

OR can do a Bayesian regularization

Laplace Smoothing
(default counts)

$$\hat{p}(y=1) \cdot \hat{p}(\underline{x} | y=1)$$

$$= \frac{3}{5} \cdot \hat{\beta}_{a,1}^3 \cdot \hat{\beta}_{b,1}^4 \cdot \hat{\beta}_{c,1}^1$$

$$= \frac{3}{5} \cdot (\hat{\beta}_{a,1})^3 \cdot (\hat{\beta}_{b,1})^4 \cdot (\hat{\beta}_{c,1})^1$$

$$\underbrace{\dots}_{\sim \propto p^1, \max p(x|y)} \cdot \underbrace{\hat{p}(y)}_{\max p(y)}$$

then compare with

$$\hat{p}(y=2) \cdot \dots$$

$$= \frac{3}{5} \cdot (\hat{\beta}_{a,2})^3 \cdot (\hat{\beta}_{b,2})^4 \cdot (\hat{\beta}_{c,2})^1$$

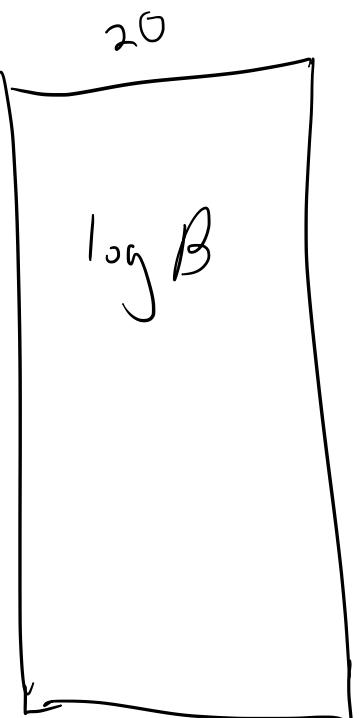
$$p(x_i | y_j, \gamma) \cdot \beta_{x_i, y_j}$$

$$p(\underline{x}|y) = \prod_{i=1}^n \beta_{x_i, y_i}$$

$$\log B(x_{-test}(\delta_1, 2), \cdot)$$

$\sim 61K$
words

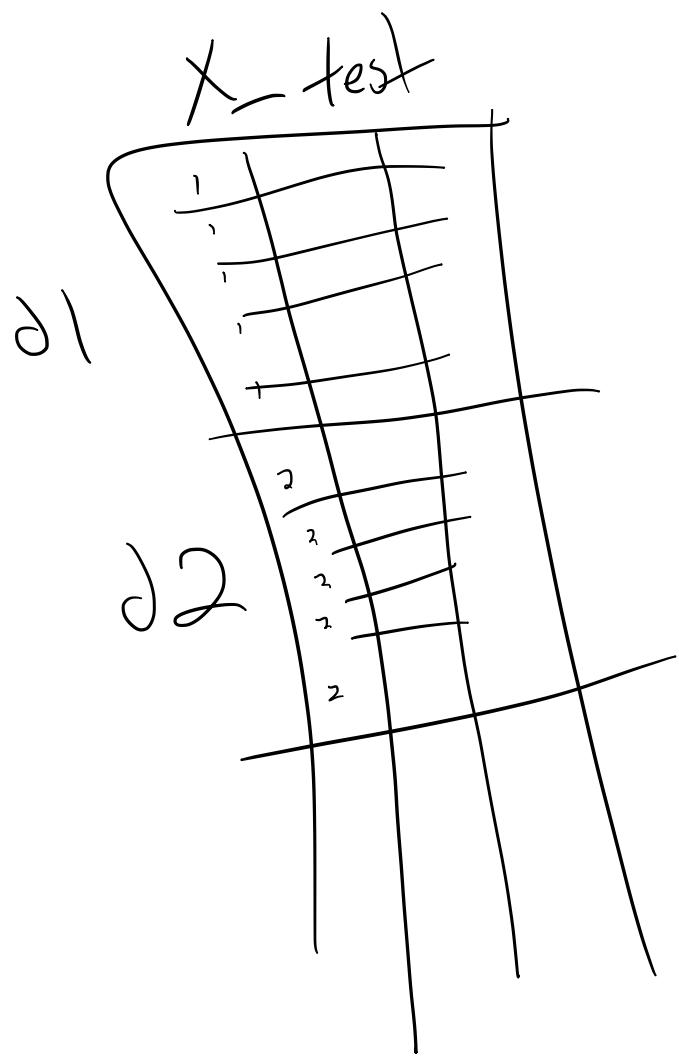
$\log B$



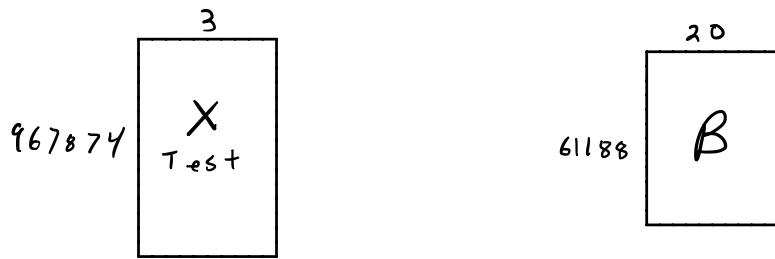
$\bullet \log(\text{likelihoods})$

$$x_{\text{-test}}(\delta_1, 2)$$

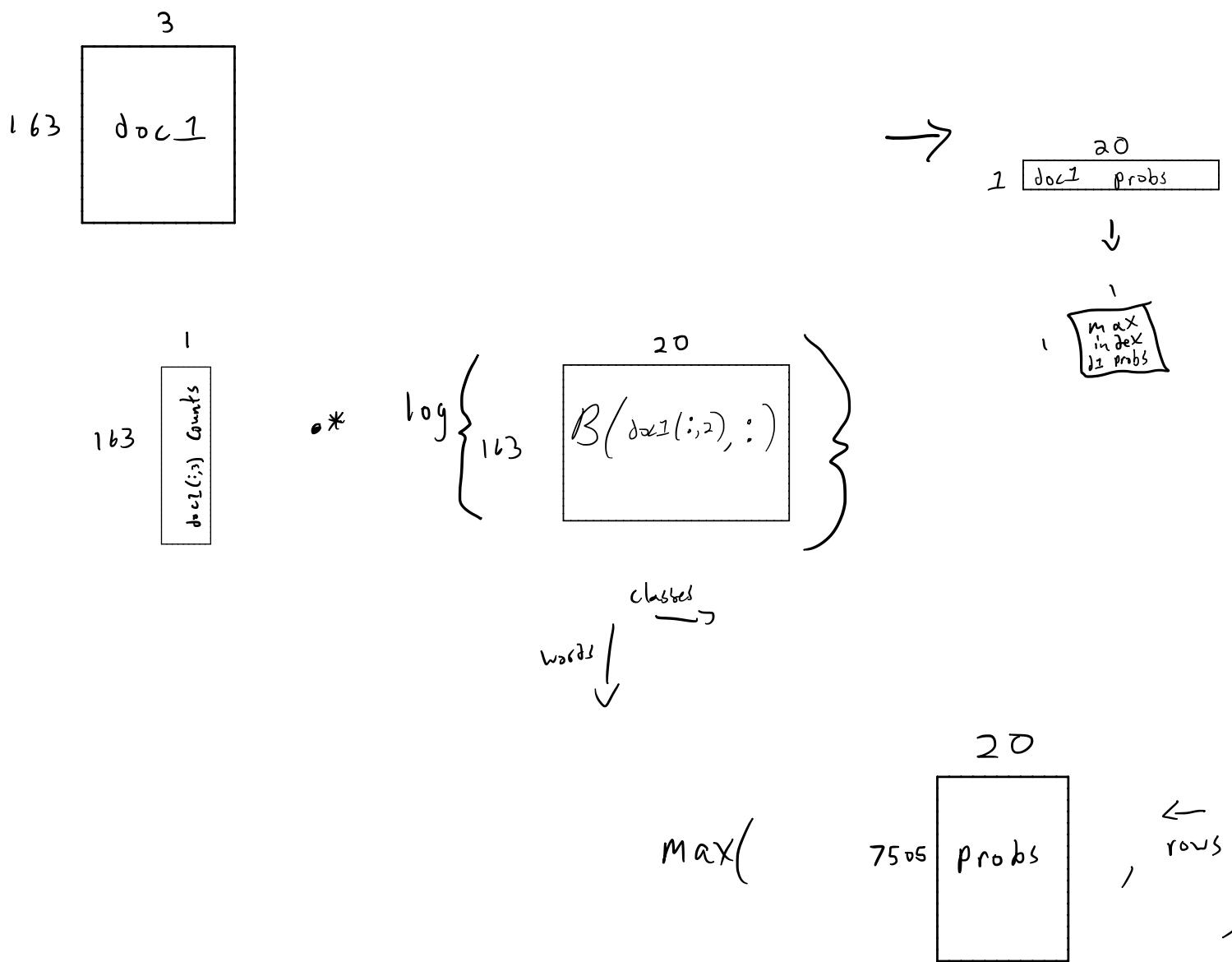
$\log(\text{likelihood})$



$$n_{DocsTest} = 7505$$



$\text{doc1}(:, 3)$: multiply



Overfitting problem

- If $W \gg n$, it is quite likely that there is a value w_0 which never occurs in the training set, but occurs in a test example \mathbf{x}_{test}
- For such a w_0 , $n_{w_0,y,i} = \beta_{w_0,y,i} = 0$ for all y, i , and $p(\mathbf{x}_{\text{test}}|y, \hat{\theta}_{\text{ML}}) = 0, \forall y$, and the **decision reduces to random guessing**. This ignores information from values that were seen in both training and test sets
- **Solution 1:** remove words that were not seen during training and proceed as before. Better than random guessing, but still ignores information in new words
- **Solution 2:** Regularize estimation of β by incorporating prior beliefs via a pdf $\pi(\beta)$.

Bayesian Naïve Bayes with Dirichlet prior

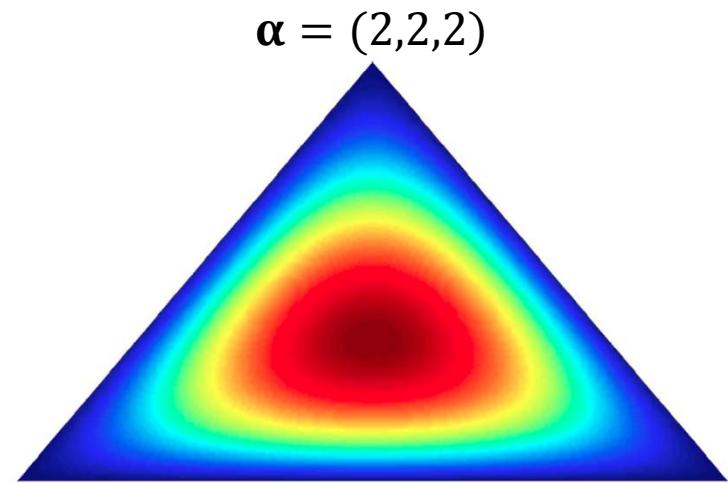
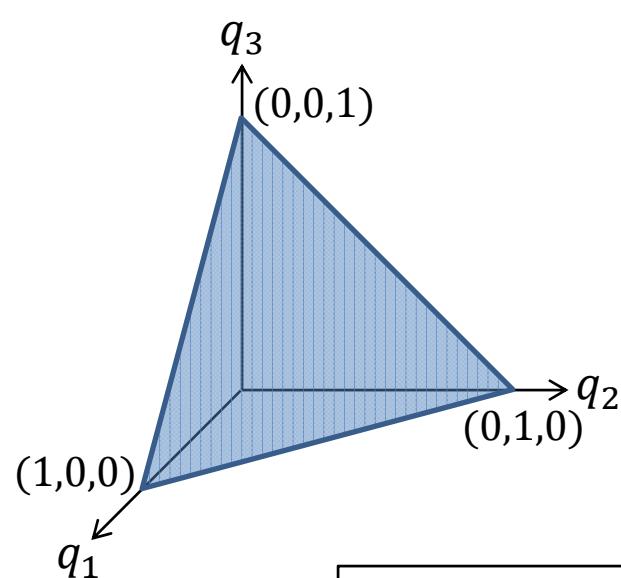
- **Dirichlet distribution** $\text{Dir}(\mathbf{q}|\boldsymbol{\alpha})$: is a family of continuous multivariate probability distributions parameterized by a W -dimensional vector $\boldsymbol{\alpha}$ of **positive reals** called the **concentration parameters**.
- It is a pdf over all pmfs over W values, i.e., a pdf over the $(W-1)$ -dimensional probability simplex:

$$\mathcal{S}_W := \left\{ \mathbf{q} : 0 \leq q_i \leq 1, i = 1, \dots, W, \sum_{i=1}^W q_i = 1 \right\}$$

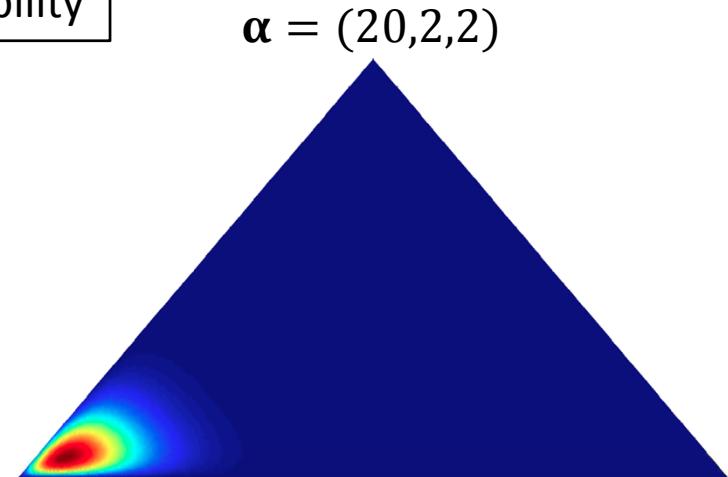
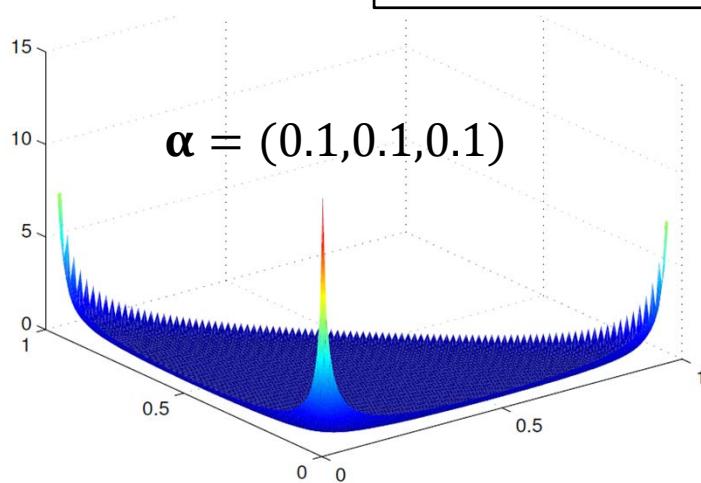
$$\text{Dir}(\mathbf{q}|\boldsymbol{\alpha}) := \frac{1(\mathbf{q} \in \mathcal{S}_W)}{Z(\boldsymbol{\alpha})} \prod_{i=1}^W q_i^{\alpha_i - 1}$$

$$Z(\boldsymbol{\alpha}) = \text{normalization constant} = \frac{1}{\Gamma\left(\sum_{i=1}^W \alpha_i\right)} \prod_{i=1}^W \Gamma(\alpha_i)$$

Dirichlet distribution



Red (hotter) \rightarrow larger probability
Blue (cooler) \rightarrow lower probability



Dirichlet distribution

- $\alpha_0 := \sum_{i=1}^W \alpha_i$ controls how peaked the distribution is (larger \Rightarrow more peaked)
- $\text{Dir}(1,1,1)$ is uniform over the probability simplex
- $\text{Dir}(2,2,2)$ is a broad distribution centered at $(1/3,1/3,1/3)$
- $\text{Dir}(20,20,20)$ is a narrow distribution centered at $(1/3,1/3,1/3)$
- If $\alpha_i < 0$ for all i , we get “spikes” at the corners of the probability simplex

Bayesian Naïve Bayes with Dirichlet prior

- Dirichlet Prior for β :

$$\pi(\beta) = \prod_{y=1}^m \prod_{i=1}^d \frac{1}{Z(\alpha_{yi})} \prod_{w=1}^W (\beta_{w,y,i})^{\alpha_{w,y,i}-1}$$

- MAP (Bayesian) estimate of θ :

$$\begin{aligned}\hat{\theta}_{\text{MAP}} &= \arg \max_{\theta} p(\theta | \mathcal{D}) \\ &= \arg \max_{\theta} p(\mathcal{D} | \theta) \pi(\beta) \\ &= \arg \max_{\theta} \left[\prod_{y=1}^m (p(y))^{n_y} \right] \cdot \left[\prod_{w=1}^W \prod_{y=1}^m \prod_{i=1}^d (\beta_{w,y,i})^{n_{w,y,i} + \alpha_{w,y,i} - 1} \right]\end{aligned}$$

Bayesian Naïve Bayes with Dirichlet prior

- MAP (Bayesian) estimate of θ solution:

$$\forall y, \hat{p}(y) = \frac{n_y}{n}$$

$$\forall w, y, i, \hat{\beta}_{w,y,i} = \frac{n_{w,y,i} + \alpha_{w,y,i} - 1}{n_y + \sum_{w=1}^W (\alpha_{w,y,i} - 1)}$$

- If for all w, y, i , $\beta_{w,y,i} = \beta_{w,y}$ then taking $\alpha_{w,y,i} = \alpha_{w,y}$ for all w, y, i , we get

$$\forall y, \hat{p}(y) = \frac{n_y}{n}$$

$$\forall w, y, i, \hat{\beta}_{w,y,i} = \frac{n_{w,y} + \alpha_{w,y} - 1}{dn_y + \sum_{w=1}^W (\alpha_{w,y} - 1)}$$

Remarks

- alphas can be interpreted as “prior” counts and the MAP solution as updating these prior counts with empirical counts from the likelihood
- If all alphas are equal to 2, then the prior counts are equal to **one**. This is referred to as **add-one smoothing** of the ML estimate or **Laplace’s rule of succession**.
- Can also incorporate a separate Dirichlet prior for $p(y)$ in a similar way
- The ML estimates of θ are asymptotically consistent

Example: Binary classification

$Y = (\text{lifestyle}) \quad \text{Peaceful / Stressed} \quad \text{Binary } (0, 1)$

Feature Vector

Values

$\underline{x} = \begin{pmatrix} 1 & \text{Income} \\ 2 & \text{charitable giving} \\ 3 & \text{Avg # hrs sleep per day} \end{pmatrix}$

- \rightarrow Low / middle / high ... Ternary, say Categorical over $(0, 1, 2)$
- \rightarrow No / Yes ... Binary, say Bernoulli over $(0, 1)$
- \rightarrow 2 - 10 ... Real values, say Gaussian (scalar)

$d = 3$

Training data

$y = 0 \quad (\text{peaceful})$

(7)

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35		
1	2	2	2	1	0	0	1	1	1	0	2	2	2	2	1	0	0	0	0	0	0	0	1	1	1	0	0	0	0	0	0	0	0	0	0	0	
2	1	1	1	1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
3	9	8	6.5	5	4	3.5	7	4	5	7	4	5	7	4	5	7	4	5	7	4	5	7	4	5	7	4	5	7	4	5	7	4	5	7	4	5	7

$y = 1 \quad (\text{stressed})$

(7)

Naive Bayes

$$\theta_{0y} = p(y) \quad \hat{p}_{ML}(0) = \frac{7}{14} \xleftarrow{n=7} = \hat{p}_{ML}(1)$$

sample mean

First feature

$$\theta_{1y} = \text{pmf over } (0, 1, 2)$$

for class '0' q_0, q_1, q_2

for class '1' r_0, r_1, r_2

$$y=0 \quad \hat{q}_{0ML} = ? \quad \frac{2}{7} \quad \hat{q}_{1ML} = ? \quad \frac{2}{7} \quad \hat{q}_{2ML} = ? \quad \frac{3}{7}$$

$$y=1 \quad \hat{r}_{0ML} = ? \quad \frac{2}{7} \quad \hat{r}_{1ML} = ? \quad \frac{2}{7} \quad \hat{r}_{2ML} = ? \quad \frac{3}{7}$$

Second feature

$$\theta_{2y} = \text{pmf over } (0, 1)$$

(charity)

$$\hat{\theta}_{2,0_{ML}} = \left(\frac{2}{7}, \frac{5}{7} \right)$$

(Class 0)
poor

$$\hat{\theta}_{2,1_{ML}} = \left(\frac{6}{7}, \frac{1}{7} \right)$$

(Class 1)
rich

Feature 3

$$\theta_{3y} = \text{for } y=0, \theta_{30} = (\mu_0, \sigma_0^2)$$

(slope)

$$\hat{\mu}_{0_{ML}} = \frac{9+8+6.5+\dots}{7}$$

estimated params.

$$\text{for } y=1, \theta_{31} = (\mu_1, \sigma_1^2)$$

$$\hat{\sigma}_{0_{ML}}^2 = \dots$$

$$\hat{\mu}_{1_{ML}} = \dots$$

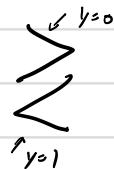
$$\hat{\sigma}_{1_{ML}}^2 = \dots$$

Test individual

$$\underline{x}_{\text{test}} = \begin{pmatrix} 1 \\ 1 \\ 3 \end{pmatrix}$$

$$y = ?$$

<u>Compare</u>	
$y=0$ $(p(y))$ $\frac{1}{2} \cdot \frac{2}{7} \cdot \frac{5}{7}$ $\cdot N(\hat{\mu}_{0_{ML}}, \hat{\sigma}_{0_{ML}}^2)$ <small>evaluated at 3</small>	$y=1$ $\frac{1}{2} \cdot \frac{3}{7} \cdot \frac{1}{7}$ $\cdot N(\hat{\mu}_{1_{ML}}, \hat{\sigma}_{1_{ML}}^2) (3)$



Why Laplace smoothing?

iPad 11:59 PM stats.stackexchange.com 90%

I was reading over Naive Bayes Classification today. I read, under the heading of **Parameter Estimation with add 1 smoothing**:

Let c refer to a class (such as Positive or Negative), and let w refer to a token or word.

The maximum likelihood estimator for $P(w|c)$ is

$$\frac{\text{count}(w, c)}{\text{count}(c)} = \frac{\text{counts w in class } c}{\text{counts of words in class } c}.$$

This estimation of $P(w|c)$ could be problematic since it would give us probability 0 for documents with unknown words. A common way of solving this problem is to use Laplace smoothing.

Let V be the set of words in the training set, add a new element UNK (for unknown) to the set of words.

Define

$$P(w|c) = \frac{\text{count}(w, c) + 1}{\text{count}(c) + |V| + 1},$$

where V refers to the vocabulary (the words in the training set).

In particular, any unknown word will have probability

$$\frac{1}{\text{count}(c) + |V| + 1}.$$

My question is this: why do we bother with this Laplacian smoothing at all? If these unknown words that we encounter in the testing set have a probability that is obviously almost zero, ie, $\frac{1}{\text{count}(c) + |V| + 1}$, what is the point of including them in the model? Why not just disregard and delete them?

share improve this question edited Oct 26 '16 at 13:03 asked Jul 22 '14 at 4:29 Mark Amery 103 340 2 6 11 Matt O'Brien 340 2 6 14

2 If you don't then any statement you encounter containing a previously unseen word will have $p = 0$. This means that an impossible event has come to pass. Which means your model was an incredibly bad fit. Also in a proper Bayesian model this could never happen, as the unknown word probability would have a numerator given by the prior (possibly not 1). So I don't know why this requires the fancy name 'Laplace smoothing'. – conjectures Jan 29 '16 at 9:55

add a comment

7 Answers active oldest votes

You always need this 'fail-safe' probability.

To see why consider the worst case where none of the words in the training sample appear in the test sentence. In this case, under your model we would conclude that the sentence is impossible but it clearly exists creating a contradiction.

Another extreme example is the test sentence "Alex met Steve." where "met" appears several times in the training sample but "Alex" and "Steve" don't. Your model would conclude this statement is very likely which is not true.

share improve this answer answered Jul 22 '14 at 5:21 Sid 1,532 4 14

I hate to sound like a complete moron, but would you mind elaborating? How does removing "Alex" and "Steve" change the likelihood of the statement occurring? – Matt O'Brien Jul 22 '14 at 6:21

2 If we assume independence of the words $P(\text{Alex})P(\text{Steve})P(\text{met}) < P(\text{met})$ – Sid Jul 22 '14 at 8:48

we could build a vocabulary when training the model on the training data set, so why not just remove all new words not occur in vocabulary when make predictions on test data set? – loganecolss Sep 28 '16 at 7:16

Get personalized job matches now
Get started stackoverflow jobs

Linked

- 0 Regarding probabilities for naiveBayes algo
- 0 Do information entropy probabilities have to sum to one?
- 1 Laplace smoothing and naive bayes

Related

- 2 Naive Bayes non-Dictionary Term in Test Document
- 1 Word probabilities in a Naive Bayes filter
- 10 In Kneser-Ney smoothing, how are unseen words handled?
- 1 Alternative Smoothing techniques for Naive Bayes?
- 1 Is the Laplace/Lidstone smoothing parameter (talking about Multinomial/Bernoulli Naive Bayes) related to the particular structure of the dataset?
- 2 Text categorization using Naive Bayes: Why isn't this working?
- 0 Naive Bayes and smoothing
- 1 bayesian classification unknown domain
- 2 Why does training naive Bayes on a data set in which all the features are repeated increase the confidence of the naive Bayes probability estimates?
- 0 Naive Bayes with Laplace Smoothing Probabilities Not Adding Up

Hot Network Questions

- How to choose between a na- and a no-adjective form of a word that takes both?
- Is the poetic device in "silence was golden" best described as metaphor or synesthesia?
- "Hat" operation not working
- If I define the word risk as "chances of losing" what would be a word for "chances of winning/gaining"?
- Word for someone who loves/wishes to live in the forest
- What is this insect in North Carolina?
- Can I use a credit card in an ATM to withdraw cash?
- C++ Coin flip simulator and data collector
- How do I center a 2x2 plate on a 3x3 plate?
- MCU - Minimum list of external components
- What do I need to legally and unproblematically take someone else's child on a trip to West Virginia?
- A player wants to always recruit NPCs into the party. How do I handle this?