

# Text as Data

Justin Grimmer

Professor  
Department of Political Science  
University of Chicago

May 8th, 2019

# Naive Bayes and General Problem Setup

Suppose we have document  $i$ , ( $i = 1, \dots, N$ ) with  $J$  features

# Naive Bayes and General Problem Setup

Suppose we have document  $i$ , ( $i = 1, \dots, N$ ) with  $J$  features

$$\mathbf{x}_i = (x_{1i}, x_{2i}, \dots, x_{Ji})$$

# Naive Bayes and General Problem Setup

Suppose we have document  $i$ , ( $i = 1, \dots, N$ ) with  $J$  features

$$\mathbf{x}_i = (x_{1i}, x_{2i}, \dots, x_{Ji})$$

Set of  $K$  categories. Category  $k$  ( $k = 1, \dots, K$ )

$$\{C_1, C_2, \dots, C_K\}$$

# Naive Bayes and General Problem Setup

Suppose we have document  $i$ , ( $i = 1, \dots, N$ ) with  $J$  features

$$\mathbf{x}_i = (x_{1i}, x_{2i}, \dots, x_{Ji})$$

Set of  $K$  categories. Category  $k$  ( $k = 1, \dots, K$ )

$$\{C_1, C_2, \dots, C_K\}$$

Subset of labeled documents  $\mathbf{Y} = (Y_1, Y_2, \dots, Y_{N_{\text{train}}})$  where

$$Y_i \in \{C_1, C_2, \dots, C_K\}.$$

# Naive Bayes and General Problem Setup

Suppose we have document  $i$ , ( $i = 1, \dots, N$ ) with  $J$  features

$$\mathbf{x}_i = (x_{1i}, x_{2i}, \dots, x_{Ji})$$

Set of  $K$  categories. Category  $k$  ( $k = 1, \dots, K$ )

$$\{C_1, C_2, \dots, C_K\}$$

Subset of labeled documents  $\mathbf{Y} = (Y_1, Y_2, \dots, Y_{N_{\text{train}}})$  where

$$Y_i \in \{C_1, C_2, \dots, C_K\}.$$

**Goal:** classify every document into **one** category.

# Naive Bayes and General Problem Setup

Suppose we have document  $i$ , ( $i = 1, \dots, N$ ) with  $J$  features

$$\mathbf{x}_i = (x_{1i}, x_{2i}, \dots, x_{Ji})$$

Set of  $K$  categories. Category  $k$  ( $k = 1, \dots, K$ )

$$\{C_1, C_2, \dots, C_K\}$$

Subset of labeled documents  $\mathbf{Y} = (Y_1, Y_2, \dots, Y_{N_{\text{train}}})$  where

$$Y_i \in \{C_1, C_2, \dots, C_K\}.$$

**Goal:** classify every document into **one** category.

Learn a function that maps from space of (possible) documents to categories

# Naive Bayes and General Problem Setup

Suppose we have document  $i$ , ( $i = 1, \dots, N$ ) with  $J$  features

$$\mathbf{x}_i = (x_{1i}, x_{2i}, \dots, x_{Ji})$$

Set of  $K$  categories. Category  $k$  ( $k = 1, \dots, K$ )

$$\{C_1, C_2, \dots, C_K\}$$

Subset of labeled documents  $\mathbf{Y} = (Y_1, Y_2, \dots, Y_{N_{\text{train}}})$  where

$$Y_i \in \{C_1, C_2, \dots, C_K\}.$$

**Goal:** classify every document into **one** category.

Learn a function that maps from space of (possible) documents to categories

To do this: use hand coded observations to estimate (train) regression model



# Naive Bayes and General Problem Setup

Suppose we have document  $i$ , ( $i = 1, \dots, N$ ) with  $J$  features

$$\mathbf{x}_i = (x_{1i}, x_{2i}, \dots, x_{Ji})$$

Set of  $K$  categories. Category  $k$  ( $k = 1, \dots, K$ )

$$\{C_1, C_2, \dots, C_K\}$$

Subset of labeled documents  $\mathbf{Y} = (Y_1, Y_2, \dots, Y_{N_{\text{train}}})$  where

$$Y_i \in \{C_1, C_2, \dots, C_K\}.$$

**Goal:** classify every document into **one** category.

Learn a function that maps from space of (possible) documents to categories

To do this: use hand coded observations to estimate (train) regression model

Apply model to test data, classify those observations

# Naive Bayes and General Problem Setup (Jurafsky Inspired Slide)

Goal: For each document  $\mathbf{x}_i$ , we want to infer most likely **category**

(0.1)

# Naive Bayes and General Problem Setup (Jurafsky Inspired Slide)

Goal: For each document  $\mathbf{x}_i$ , we want to infer most likely **category**

$$C_{\text{Max}} = \arg \max_k p(C_k | \mathbf{x}_i)$$

(0.1)

# Naive Bayes and General Problem Setup (Jurafsky Inspired Slide)

Goal: For each document  $\mathbf{x}_i$ , we want to infer most likely **category**

$$C_{\text{Max}} = \arg \max_k p(C_k | \mathbf{x}_i)$$

We're going to use Bayes' rule to estimate  $p(C_k | \mathbf{x}_i)$ .

(0.1)

# Naive Bayes and General Problem Setup (Jurafsky Inspired Slide)

Goal: For each document  $\mathbf{x}_i$ , we want to infer most likely **category**

$$C_{\text{Max}} = \arg \max_k p(C_k | \mathbf{x}_i)$$

We're going to use Bayes' rule to estimate  $p(C_k | \mathbf{x}_i)$ .

$$p(C_k | \mathbf{x}_i) = \frac{p(C_k, \mathbf{x}_i)}{p(\mathbf{x}_i)}$$

(0.1)

# Naive Bayes and General Problem Setup (Jurafsky Inspired Slide)

Goal: For each document  $\mathbf{x}_i$ , we want to infer most likely **category**

$$C_{\text{Max}} = \arg \max_k p(C_k | \mathbf{x}_i)$$

We're going to use Bayes' rule to estimate  $p(C_k | \mathbf{x}_i)$ .

$$\begin{aligned} p(C_k | \mathbf{x}_i) &= \frac{p(C_k, \mathbf{x}_i)}{p(\mathbf{x}_i)} \\ &= \frac{p(C_k)p(\mathbf{x}_i | C_k)}{p(\mathbf{x}_i)} \end{aligned} \tag{0.1}$$

# Naive Bayes and General Problem Setup (Jurafsky Inspired Slide)

Goal: For each document  $\mathbf{x}_i$ , we want to infer most likely **category**

$$C_{\text{Max}} = \arg \max_k p(C_k | \mathbf{x}_i)$$

We're going to use Bayes' rule to estimate  $p(C_k | \mathbf{x}_i)$ .

$$\begin{aligned} p(C_k | \mathbf{x}_i) &= \frac{p(C_k, \mathbf{x}_i)}{p(\mathbf{x}_i)} \\ &\quad \text{Proportion in } C_k \\ &\quad \underbrace{p(C_k)} \quad \underbrace{p(\mathbf{x}_i | C_k)} \\ &= \frac{\text{Language model}}{p(\mathbf{x}_i)} \end{aligned}$$

# Naive Bayes and Optimization (Jurafsky Inspired Slide)



# Naive Bayes and Optimization (Jurafsky Inspired Slide)

$$C_{\text{Max}} = \arg \max_k p(C_k | \mathbf{x}_i)$$

# Naive Bayes and Optimization (Jurafsky Inspired Slide)

$$C_{\text{Max}} = \arg \max_k p(C_k | \mathbf{x}_i)$$

$$C_{\text{Max}} = \arg \max_k \frac{p(C_k)p(\mathbf{x}_i | C_k)}{p(\mathbf{x}_i)}$$

# Naive Bayes and Optimization (Jurafsky Inspired Slide)

$$C_{\text{Max}} = \arg \max_k p(C_k | \mathbf{x}_i)$$

$$C_{\text{Max}} = \arg \max_k \frac{p(C_k)p(\mathbf{x}_i | C_k)}{p(\mathbf{x}_i)}$$

$$C_{\text{Max}} = \arg \max_k p(C_k)p(\mathbf{x}_i | C_k)$$

# Naive Bayes and Optimization (Jurafsky Inspired Slide)

$$C_{\text{Max}} = \arg \max_k p(C_k | \mathbf{x}_i)$$

$$C_{\text{Max}} = \arg \max_k \frac{p(C_k)p(\mathbf{x}_i | C_k)}{p(\mathbf{x}_i)}$$

$$C_{\text{Max}} = \arg \max_k p(C_k)p(\mathbf{x}_i | C_k)$$

Two probabilities to estimate:

# Naive Bayes and Optimization (Jurafsky Inspired Slide)

$$C_{\text{Max}} = \arg \max_k p(C_k | \mathbf{x}_i)$$

$$C_{\text{Max}} = \arg \max_k \frac{p(C_k)p(\mathbf{x}_i | C_k)}{p(\mathbf{x}_i)}$$

$$C_{\text{Max}} = \arg \max_k p(C_k)p(\mathbf{x}_i | C_k)$$

Two probabilities to estimate:

$$p(C_k) = \frac{\text{No. Documents in } k}{\text{No. Documents}} \text{ (training set)}$$

# Naive Bayes and Optimization (Jurafsky Inspired Slide)

$$C_{\text{Max}} = \arg \max_k p(C_k | \mathbf{x}_i)$$

$$C_{\text{Max}} = \arg \max_k \frac{p(C_k)p(\mathbf{x}_i | C_k)}{p(\mathbf{x}_i)}$$

$$C_{\text{Max}} = \arg \max_k p(C_k)p(\mathbf{x}_i | C_k)$$

Two probabilities to estimate:

$$p(C_k) = \frac{\text{No. Documents in } k}{\text{No. Documents}} \text{ (training set)}$$

$p(\mathbf{x}_i | C_k)$  **complicated** without assumptions

# Naive Bayes and Optimization (Jurafsky Inspired Slide)

$$C_{\text{Max}} = \arg \max_k p(C_k | \mathbf{x}_i)$$

$$C_{\text{Max}} = \arg \max_k \frac{p(C_k)p(\mathbf{x}_i | C_k)}{p(\mathbf{x}_i)}$$

$$C_{\text{Max}} = \arg \max_k p(C_k)p(\mathbf{x}_i | C_k)$$

Two probabilities to estimate:

$$p(C_k) = \frac{\text{No. Documents in } k}{\text{No. Documents}} \text{ (training set)}$$

$$p(\mathbf{x}_i | C_k) \text{ \textcolor{red}{complicated} without assumptions}$$

- Imagine each  $x_{ij}$  just binary indicator. Then  $2^J$  possible  $\mathbf{x}_i$  documents

# Naive Bayes and Optimization (Jurafsky Inspired Slide)

$$C_{\text{Max}} = \arg \max_k p(C_k | \mathbf{x}_i)$$

$$C_{\text{Max}} = \arg \max_k \frac{p(C_k)p(\mathbf{x}_i | C_k)}{p(\mathbf{x}_i)}$$

$$C_{\text{Max}} = \arg \max_k p(C_k)p(\mathbf{x}_i | C_k)$$

Two probabilities to estimate:

$$p(C_k) = \frac{\text{No. Documents in } k}{\text{No. Documents}} \text{ (training set)}$$

$p(\mathbf{x}_i | C_k)$  **complicated** without assumptions

- Imagine each  $x_{ij}$  just binary indicator. Then  $2^J$  possible  $\mathbf{x}_i$  documents
- Simplify: assume each feature is independent



# Naive Bayes and Optimization (Jurafsky Inspired Slide)

$$C_{\text{Max}} = \arg \max_k p(C_k | \mathbf{x}_i)$$

$$C_{\text{Max}} = \arg \max_k \frac{p(C_k)p(\mathbf{x}_i | C_k)}{p(\mathbf{x}_i)}$$

$$C_{\text{Max}} = \arg \max_k p(C_k)p(\mathbf{x}_i | C_k)$$

Two probabilities to estimate:

$$p(C_k) = \frac{\text{No. Documents in } k}{\text{No. Documents}} \text{ (training set)}$$

$p(\mathbf{x}_i | C_k)$  **complicated** without assumptions

- Imagine each  $x_{ij}$  just binary indicator. Then  $2^J$  possible  $\mathbf{x}_i$  documents
- Simplify: assume each feature is independent

$$p(\mathbf{x}_i | C_k) = \prod_{j=1}^J p(x_{ij} | C_k)$$

# Naive Bayes and Optimization (Jurafsky Inspired Slide)

Two components to estimation:

- $p(C_k) = \frac{\text{No. Documents in } k}{\text{No. Documents}}$  (training set)
- $p(\mathbf{x}_i | C_k) = \prod_{j=1}^J p(x_{ij} | C_k)$

# Naive Bayes and Optimization (Jurafsky Inspired Slide)

Two components to estimation:

- $p(C_k) = \frac{\text{No. Documents in } k}{\text{No. Documents}}$  (training set)
- $p(\mathbf{x}_i|C_k) = \prod_{j=1}^J p(x_{ij}|C_k)$

Maximum likelihood estimation (training set):

# Naive Bayes and Optimization (Jurafsky Inspired Slide)

Two components to estimation:

- $p(C_k) = \frac{\text{No. Documents in } k}{\text{No. Documents}}$  (training set)
- $p(\mathbf{x}_i|C_k) = \prod_{j=1}^J p(x_{ij}|C_k)$

Maximum likelihood estimation (training set):

$$p(x_{im} = z|C_k) = \frac{\text{No}(\text{ Docs}_{ij} = z \text{ and } C = C_k )}{\text{No}(C = C_k)}$$

# Naive Bayes and Optimization (Jurafsky Inspired Slide)

Two components to estimation:

- $p(C_k) = \frac{\text{No. Documents in } k}{\text{No. Documents}}$  (training set)
- $p(\mathbf{x}_i|C_k) = \prod_{j=1}^J p(x_{ij}|C_k)$

Maximum likelihood estimation (training set):

$$p(x_{im} = z|C_k) = \frac{\text{No}(\text{ Docs}_{ij} = z \text{ and } C = C_k )}{\text{No}(C = C_k)}$$

**Problem:** What if  $\text{No}(\text{ Docs}_{ij} = z \text{ and } C = C_k ) = 0$  ?

# Naive Bayes and Optimization (Jurafsky Inspired Slide)

Two components to estimation:

- $p(C_k) = \frac{\text{No. Documents in } k}{\text{No. Documents}}$  (training set)
- $p(\mathbf{x}_i | C_k) = \prod_{j=1}^J p(x_{ij} | C_k)$

Maximum likelihood estimation (training set):

$$p(x_{im} = z | C_k) = \frac{\text{No}(\text{ Docs}_{ij} = z \text{ and } C = C_k )}{\text{No}(C = C_k)}$$

**Problem:** What if  $\text{No}(\text{ Docs}_{ij} = z \text{ and } C = C_k ) = 0$  ?

$$\prod_{j=1}^J p(x_{ij} | C_k) = 0$$

# Naive Bayes and General Problem Setup (Jurafsky Inspired Slide)

# Naive Bayes and General Problem Setup (Jurafsky Inspired Slide)

Solution: smoothing (Bayesian estimation)



# Naive Bayes and General Problem Setup (Jurafsky Inspired Slide)

Solution: smoothing (Bayesian estimation)

$$p(x_{ij} = z | C_k) = \frac{\text{No}(\text{ Docs}_{ij} = z \text{ and } C = C_k ) + 1}{\text{No}(C = C_k) + k}$$

# Naive Bayes and General Problem Setup (Jurafsky Inspired Slide)

Solution: smoothing (Bayesian estimation)

$$p(x_{ij} = z | C_k) = \frac{\text{No}(\text{Docs}_{ij} = z \text{ and } C = C_k) + 1}{\text{No}(C = C_k) + k}$$

Algorithm steps:

# Naive Bayes and General Problem Setup (Jurafsky Inspired Slide)

Solution: smoothing (Bayesian estimation)

$$p(x_{ij} = z | C_k) = \frac{\text{No}(\text{ Docs}_{ij} = z \text{ and } C = C_k ) + 1}{\text{No}(C = C_k) + k}$$

Algorithm steps:

- 1) Learn  $\hat{p}(C)$  and  $\hat{p}(\mathbf{x}_i | C_k)$  on **training data**

# Naive Bayes and General Problem Setup (Jurafsky Inspired Slide)

Solution: smoothing (Bayesian estimation)

$$p(x_{ij} = z | C_k) = \frac{\text{No}(\text{Docs}_{ij} = z \text{ and } C = C_k) + 1}{\text{No}(C = C_k) + k}$$

Algorithm steps:

- 1) Learn  $\hat{p}(C)$  and  $\hat{p}(\mathbf{x}_i | C_k)$  on **training data**
- 2) Use this to identify most likely  $C_k$  for each document  $i$  in **test set**

# Naive Bayes and General Problem Setup (Jurafsky Inspired Slide)

Solution: smoothing (Bayesian estimation)

$$p(x_{ij} = z | C_k) = \frac{\text{No}(\text{Docs}_{ij} = z \text{ and } C = C_k) + 1}{\text{No}(C = C_k) + k}$$

Algorithm steps:

- 1) Learn  $\hat{p}(C)$  and  $\hat{p}(\mathbf{x}_i | C_k)$  on **training data**
- 2) Use this to identify most likely  $C_k$  for each document  $i$  in **test set**

$$C_i = \arg \max_k \hat{p}(C_k) \hat{p}(\mathbf{x}_i | C_k)$$

# Naive Bayes and General Problem Setup (Jurafsky Inspired Slide)

Solution: smoothing (Bayesian estimation)

$$p(x_{ij} = z | C_k) = \frac{\text{No}(\text{ Docs}_{ij} = z \text{ and } C = C_k ) + 1}{\text{No}(C = C_k) + k}$$

Algorithm steps:

- 1) Learn  $\hat{p}(C)$  and  $\hat{p}(\mathbf{x}_i | C_k)$  on **training data**
- 2) Use this to identify most likely  $C_k$  for each document  $i$  in **test set**

$$C_i = \arg \max_k \hat{p}(C_k) \hat{p}(\mathbf{x}_i | C_k)$$

Simple intuition about Naive Bayes:

# Naive Bayes and General Problem Setup (Jurafsky Inspired Slide)

Solution: smoothing (Bayesian estimation)

$$p(x_{ij} = z | C_k) = \frac{\text{No}(\text{ Docs}_{ij} = z \text{ and } C = C_k ) + 1}{\text{No}(C = C_k) + k}$$

Algorithm steps:

- 1) Learn  $\hat{p}(C)$  and  $\hat{p}(\mathbf{x}_i | C_k)$  on **training data**
- 2) Use this to identify most likely  $C_k$  for each document  $i$  in **test set**

$$C_i = \arg \max_k \hat{p}(C_k) \hat{p}(\mathbf{x}_i | C_k)$$

Simple intuition about Naive Bayes:

- Learn what documents in class  $j$  look like

# Naive Bayes and General Problem Setup (Jurafsky Inspired Slide)

Solution: smoothing (Bayesian estimation)

$$p(x_{ij} = z | C_k) = \frac{\text{No}(\text{ Docs}_{ij} = z \text{ and } C = C_k ) + 1}{\text{No}(C = C_k) + k}$$

Algorithm steps:

- 1) Learn  $\hat{p}(C)$  and  $\hat{p}(\mathbf{x}_i | C_k)$  on **training data**
- 2) Use this to identify most likely  $C_k$  for each document  $i$  in **test set**

$$C_i = \arg \max_k \hat{p}(C_k) \hat{p}(\mathbf{x}_i | C_k)$$

Simple intuition about Naive Bayes:

- Learn what documents in class  $j$  look like
- Find class  $k$  that document  $i$  is most similar to



# Naive Bayes and Unigram Language Models

Assume the following data generating process (should look familiar)

$$\boldsymbol{\pi} \sim \text{Dirichlet}(\boldsymbol{\alpha})$$

$$\boldsymbol{\theta} \sim \text{Dirichlet}(\boldsymbol{\lambda})$$

$$\boldsymbol{\tau}_i \sim \text{Multinomial}(1, \boldsymbol{\pi})$$

$$\mathbf{x}_i | \tau_{ik} = 1, \boldsymbol{\theta} \sim \text{Multinomial}(n_i, \boldsymbol{\theta}_k)$$

# Naive Bayes and Unigram Language Models

Assume the following data generating process (should look familiar)

$$\boldsymbol{\pi} \sim \text{Dirichlet}(\boldsymbol{\alpha})$$

$$\boldsymbol{\theta} \sim \text{Dirichlet}(\boldsymbol{\lambda})$$

$$\boldsymbol{\tau}_i \sim \text{Multinomial}(1, \boldsymbol{\pi})$$

$$\mathbf{x}_i | \tau_{ik} = 1, \boldsymbol{\theta} \sim \text{Multinomial}(n_i, \boldsymbol{\theta}_k)$$

If we randomly sample documents  $N_{\text{train}}$  and label them ( $\mathbf{Y}$ ), then we can estimate

# Naive Bayes and Unigram Language Models

Assume the following data generating process (should look familiar)

$$\begin{aligned}\boldsymbol{\pi} &\sim \text{Dirichlet}(\boldsymbol{\alpha}) \\ \boldsymbol{\theta} &\sim \text{Dirichlet}(\boldsymbol{\lambda}) \\ \boldsymbol{\tau}_i &\sim \text{Multinomial}(1, \boldsymbol{\pi}) \\ \mathbf{x}_i | \tau_{ik} = 1, \boldsymbol{\theta} &\sim \text{Multinomial}(n_i, \boldsymbol{\theta}_k)\end{aligned}$$

If we randomly sample documents  $N_{\text{train}}$  and label them ( $\mathbf{Y}$ ), then we can estimate

$$\hat{\pi}_k = \frac{\sum_{i=1}^N I(Y_i = k) + \alpha_k}{N_{\text{train}} + \sum_{k=1}^K \alpha_k}$$

# Naive Bayes and Unigram Language Models

Assume the following data generating process (should look familiar)

$$\begin{aligned}\boldsymbol{\pi} &\sim \text{Dirichlet}(\boldsymbol{\alpha}) \\ \boldsymbol{\theta} &\sim \text{Dirichlet}(\boldsymbol{\lambda}) \\ \boldsymbol{\tau}_i &\sim \text{Multinomial}(1, \boldsymbol{\pi}) \\ \mathbf{x}_i | \tau_{ik} = 1, \boldsymbol{\theta} &\sim \text{Multinomial}(n_i, \boldsymbol{\theta}_k)\end{aligned}$$

If we randomly sample documents  $N_{\text{train}}$  and label them ( $\mathbf{Y}$ ), then we can estimate

$$\begin{aligned}\hat{\pi}_k &= \frac{\sum_{i=1}^N I(Y_i = k) + \alpha_k}{N_{\text{train}} + \sum_{k=1}^K \alpha_k} \\ \hat{\theta}_{jk} &= \frac{\sum_{i=1}^N I(Y_i = k) x_{ij} + \lambda_j}{\sum_{j=1}^J \sum_{i=1}^N I(Y_i = k) x_{ij} + \sum_{j=1}^J \lambda_j}\end{aligned}$$

# Naive Bayes and Unigram Language Models

The probability a new document has  $\tau_{ik} = 1$  is then

# Naive Bayes and Unigram Language Models

The probability a new document has  $\tau_{ik} = 1$  is then

$$p(\tau_{ik} = 1 | \mathbf{x}_i, \hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\theta}}) \propto p(\tau_{ik} = 1) p(\mathbf{x}_i | \boldsymbol{\theta}, \tau_{ik} = 1)$$

# Naive Bayes and Unigram Language Models

The probability a new document has  $\tau_{ik} = 1$  is then

$$\begin{aligned} p(\tau_{ik} = 1 | \mathbf{x}_i, \hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\theta}}) &\propto p(\tau_{ik} = 1) p(\mathbf{x}_i | \boldsymbol{\theta}, \tau_{ik} = 1) \\ &\propto \hat{\pi}_k \prod_{j=1}^J (\hat{\theta}_{jk})^{x_{ij}} \end{aligned}$$

# Naive Bayes and Unigram Language Models

The probability a new document has  $\tau_{ik} = 1$  is then

$$\begin{aligned} p(\tau_{ik} = 1 | \mathbf{x}_i, \hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\theta}}) &\propto p(\tau_{ik} = 1) p(\mathbf{x}_i | \boldsymbol{\theta}, \tau_{ik} = 1) \\ &\propto \hat{\pi}_k \prod_{j=1}^J \left( \hat{\theta}_{jk} \right)^{x_{ij}} \\ &\propto \underbrace{\hat{\pi}_k}_{p(C_k)} \underbrace{\prod_{j=1}^J \left( \hat{\theta}_{jk} \right)^{x_{ij}}}_{\text{Unigram model}} \end{aligned}$$



# Some R Code

```
library(e1071)
dep<- c(labels, rep(NA, no.testSet))
dep<- as.factor(dep)
out<- naiveBayes(dep~., as.data.frame(tdm))
predicts<- predict(out, as.data.frame(tdm[-training.set,]))
```

# ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Naive Bayes, LASSO, ...: focused on individual document classification.

# ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Naive Bayes, LASSO, ...: focused on individual document classification.  
But what if we're focused on **proportions only**?

# ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Naive Bayes, LASSO, ...: focused on individual document classification.

But what if we're focused on **proportions only**?

Hopkins and King (2010): method for characterizing distribution of classes

# ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Naive Bayes, LASSO, ...: focused on individual document classification.

But what if we're focused on **proportions only**?

Hopkins and King (2010): method for characterizing distribution of classes

**Can be much more accurate than individual classifiers**, requires fewer assumptions (**do not need random sample of documents** ) .

# ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Naive Bayes, LASSO, ...: focused on individual document classification.

But what if we're focused on **proportions only**?

Hopkins and King (2010): method for characterizing distribution of classes

**Can be much more accurate than individual classifiers**, requires fewer assumptions (**do not need random sample of documents** ) .

- King and Lu (2008): derive method for characterizing causes of deaths for verbal autopsies

# ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Naive Bayes, LASSO, ...: focused on individual document classification.

But what if we're focused on **proportions only**?

Hopkins and King (2010): method for characterizing distribution of classes

**Can be much more accurate than individual classifiers**, requires fewer assumptions (**do not need random sample of documents** ) .

- King and Lu (2008): derive method for characterizing causes of deaths for verbal autopsies
- Hopkins and King (2010): extend the method to text documents

# ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Naive Bayes, LASSO, ...: focused on individual document classification.

But what if we're focused on **proportions only**?

Hopkins and King (2010): method for characterizing distribution of classes

**Can be much more accurate than individual classifiers**, requires fewer assumptions (**do not need random sample of documents** ) .

- King and Lu (2008): derive method for characterizing causes of deaths for verbal autopsies
- Hopkins and King (2010): extend the method to text documents

Basic intuition:



# ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Naive Bayes, LASSO, ...: focused on individual document classification.

But what if we're focused on **proportions only**?

Hopkins and King (2010): method for characterizing distribution of classes

**Can be much more accurate than individual classifiers**, requires fewer assumptions (**do not need random sample of documents** ) .

- King and Lu (2008): derive method for characterizing causes of deaths for verbal autopsies
- Hopkins and King (2010): extend the method to text documents

Basic intuition:

- Examine joint distribution of characteristics (without making Naive Bayes like assumption)
- Focus on distributions (only) makes this analysis possible

# ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Measure **only** presence/absence of each term [ $J \times 1$  vector]

# ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Measure **only** presence/absence of each term  $[(J \times 1)$  vector]

$$\mathbf{x}_i = (1, 0, 0, 1, \dots, 0)$$

# ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Measure **only** presence/absence of each term  $[(J \times 1)$  vector]

$$\mathbf{x}_i = (1, 0, 0, 1, \dots, 0)$$

What are the possible realizations of  $\mathbf{x}_i$ ?

# ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Measure **only** presence/absence of each term [ $(J \times 1)$  vector]

$$\mathbf{x}_i = (1, 0, 0, 1, \dots, 0)$$

What are the possible realizations of  $\mathbf{x}_i$ ?

- $2^J$  possible vectors

# ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Measure **only** presence/absence of each term  $[(J \times 1) \text{ vector}]$

$$\mathbf{x}_i = (1, 0, 0, 1, \dots, 0)$$

What are the possible realizations of  $\mathbf{x}_i$ ?

- $2^J$  possible vectors

Define:

# ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Measure **only** presence/absence of each term  $[(J \times 1) \text{ vector}]$

$$\mathbf{x}_i = (1, 0, 0, 1, \dots, 0)$$

What are the possible realizations of  $\mathbf{x}_i$ ?

-  $2^J$  possible vectors

Define:

$$P(\mathbf{x}) = \text{probability of observing } \mathbf{x}$$

# ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Measure **only** presence/absence of each term [ $(J \times 1)$  vector]

$$\mathbf{x}_i = (1, 0, 0, 1, \dots, 0)$$

What are the possible realizations of  $\mathbf{x}_i$ ?

-  $2^J$  possible vectors

Define:

$P(\mathbf{x})$  = probability of observing  $\mathbf{x}$

$P(\mathbf{x}|C_j)$  = Probability of observing  $\mathbf{x}$  conditional on category  $C_j$



# ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Measure **only** presence/absence of each term [ $(J \times 1)$  vector]

$$\mathbf{x}_i = (1, 0, 0, 1, \dots, 0)$$

What are the possible realizations of  $\mathbf{x}_i$ ?

-  $2^J$  possible vectors

Define:

$P(\mathbf{x})$  = probability of observing  $\mathbf{x}$

$P(\mathbf{x}|C_j)$  = Probability of observing  $\mathbf{x}$  conditional on category  $C_j$

$P(\mathbf{X}|C)$  = Matrix collecting vectors

# ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Measure **only** presence/absence of each term [ $(J \times 1)$  vector]

$$\mathbf{x}_i = (1, 0, 0, 1, \dots, 0)$$

What are the possible realizations of  $\mathbf{x}_i$ ?

-  $2^J$  possible vectors

Define:

$P(\mathbf{x})$  = probability of observing  $\mathbf{x}$

$P(\mathbf{x}|C_j)$  = Probability of observing  $\mathbf{x}$  conditional on category  $C_j$

$P(\mathbf{X}|C)$  = Matrix collecting vectors

$P(C)$  =  $P(C_1, C_2, \dots, C_K)$  target quantity of interest

# ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

Measure **only** presence/absence of each term [ $(J \times 1)$  vector]

$$\mathbf{x}_i = (1, 0, 0, 1, \dots, 0)$$

What are the possible realizations of  $\mathbf{x}_i$ ?

-  $2^J$  possible vectors

Define:

$P(\mathbf{x})$  = probability of observing  $\mathbf{x}$

$P(\mathbf{x}|C_j)$  = Probability of observing  $\mathbf{x}$  conditional on category  $C_j$

$P(\mathbf{X}|C)$  = Matrix collecting vectors

$P(C)$  =  $P(C_1, C_2, \dots, C_K)$  target quantity of interest

# ReadMe: Optimization for a Different Goal (Hopkins and King 2010)

$$\underbrace{P(\mathbf{x})}_{2^J \times 1} = \underbrace{P(\mathbf{x}|C)}_{2^J \times K} \underbrace{P(C)}_{K \times 1}$$

Matrix algebra problem to solve, for  $P(C)$

Like Naive Bayes, requires two pieces to estimate

Complication  $2^J \gg$  no. documents

**Kernel Smoothing Methods** (without a formal model)

- $P(\mathbf{x})$  = estimate directly from test set
- $P(\mathbf{x}|C)$  = estimate from training set
  - Key assumption:  $P(\mathbf{x}|C)$  in training set is equivalent to  $P(\mathbf{x}|C)$  in test set
- If true, can perform biased sampling of documents, worry less about drift...

# Algorithm Summarized

- Estimate  $\hat{p}(\mathbf{x})$  from test set
- Estimate  $\hat{p}(\mathbf{x}|C)$  from training set
- Use  $\hat{p}(\mathbf{x})$  and  $\hat{p}(\mathbf{x}|C)$  to solve for  $p(C)$

# Assessing Model Performance

Not classifying individual documents  $\rightarrow$  different standards

Mean Square Error :

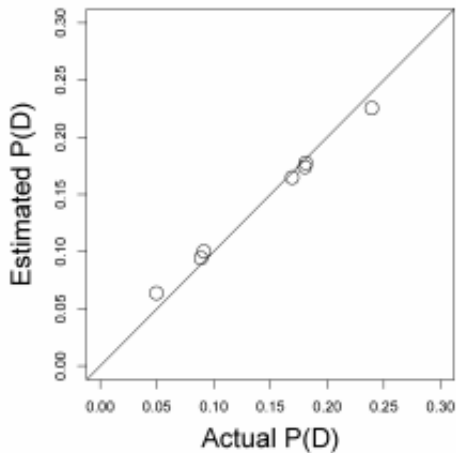
$$E[(\hat{\theta} - \theta)^2] = \text{var}(\hat{\theta}) + \text{Bias}(\hat{\theta}, \theta)^2$$

Suppose we have true proportions  $P(C)^{\text{true}}$ . Then, we'll estimate **Root Mean Square Error**

$$\text{RMSE} = \sqrt{\frac{\sum_{j=1}^J (P(C_j)^{\text{true}} - P(C_j))^2}{J}}$$

$$\text{Mean Abs. Prediction Error} = \left| \frac{\sum_{j=1}^J (P(C_j)^{\text{true}} - P(C_j))}{J} \right|$$

**Visualize:** plot true and estimated proportions



**TABLE 1 Performance of Our Nonparametric Approach and Four Support Vector Machine Analyses**

	Percent of Blog Posts Correctly Classified			
	In-Sample Fit	In-Sample Cross-Validation	Out-of-Sample Prediction	Mean Absolute Proportion Error
Nonparametric	—	—	—	1.2
Linear	67.6	55.2	49.3	7.7
Radial	67.6	54.2	49.1	7.7
Polynomial	99.7	48.9	47.8	5.3
Sigmoid	15.6	15.6	18.2	23.2

*Notes:* Each row is the optimal choice over numerous individual runs given a specific kernel. Leaving aside the sigmoid kernel, individual classification performance in the first three columns does not correlate with mean absolute error in the document category proportions in the last column.



# Using the House Press Release Data

Method	RMSE	APSE
ReadMe	0.036	0.056
NaiveBayes	0.096	0.14
SVM	0.052	0.084

# Code to Run in R

Control file:

filename	truth	trainingset
20July2009LEWIS53.txt	4	1
26July2006LEWIS249.txt	2	0

```
tdm<- undergrad(control=control, fullfreq=F)
process<- preprocess(tdm)
output<- undergrad(process)
output$est.CSMF ## proportion in each category
output$true.CSMF ## if labeled for validation set (but not
used in training set)
```