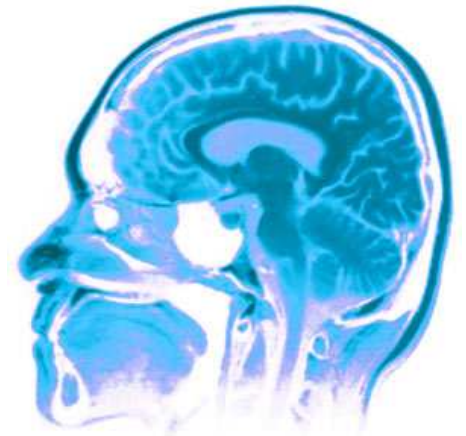




CPS C340



Dirichlet and Categorical
variables:
Naïve Bayes classifier



Nando de Freitas
November, 2012
University of British Columbia

Outline of the lecture

This lecture introduces the Dirichlet and categorical distributions, as well as the Naïve Bayes classifier. The goal is for you to:

- Learn categorical distributions.
- Derive the Dirichlet posterior from the Dirichlet prior and categorical likelihood.
- Understand how a classifier for text is set up.
- Understand the Naïve Bayes classifier for text classification.

Revision: Beta-Bernoulli

Suppose $X_i \sim \text{Ber}(\theta)$, so $X_i \in \{0, 1\}$. We know that the likelihood has the form

$$p(\mathcal{D}|\theta) = \theta^{N_1} (1 - \theta)^{N_0}$$

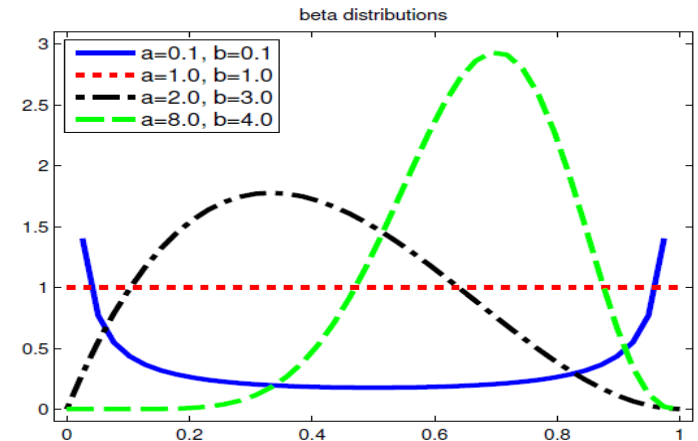
Handwritten note: $X_{1:N} = \text{data}$

where we have $N_1 = \sum_{i=1}^N \mathbb{I}(x_i = 1)$ heads and $N_0 = \sum_{i=1}^N \mathbb{I}(x_i = 0)$ tails.

The **beta prior** has pdf:

$$0 \leq \theta \leq 1$$

$$\text{Beta}(\theta|\alpha_1, \alpha_2) = \frac{1}{B(\alpha_1, \alpha_2)} \theta^{\alpha_1-1} (1 - \theta)^{\alpha_2-1}$$



Revision: Beta-Bernoulli

If we multiply the Bernoulli likelihood by the beta prior we get

$$\begin{aligned} p(\theta|\mathcal{D}) &\propto p(\mathcal{D}|\theta)p(\theta) \\ &\propto [\theta^{N_1}(1-\theta)^{N_2}] [\theta^{\alpha_1-1}(1-\theta)^{\alpha_2-1}] \\ &= \theta^{N_1+\alpha_1-1}(1-\theta)^{N_2+\alpha_2-1} \\ &\propto \text{Beta}(\theta|N_1 + \alpha_1, N_2 + \alpha_2) \end{aligned}$$

We see that the posterior has the same functional form (beta) as the prior (beta), since it is conjugate.

Categorical distribution

The multivariate version of the Bernoulli distribution is the **Categorical distribution** (an instance of the **multinomial distribution**).

We are given n data points, $\mathbf{x}_{1:n} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$. Each point \mathbf{x}_i indicates one of K values. For example if $K=3$, then the possible vectors are (100) , (010) and (001) .

$$\theta_k = 1 - \sum_{j=1}^{k-1} \theta_j$$

$$\mathbf{x}_i = \begin{cases} 1 & 100 \\ 2 & 010 \\ 3 & 001 \end{cases}$$

The likelihood of the data is then:

$$\begin{aligned} P(x_i=1|\theta) &= \theta_1 \\ P(x_i=2|\theta) &= \theta_2 \\ &\vdots \end{aligned}$$

$$p(\mathbf{x}_i | \theta) = \text{Cat}(\mathbf{x}_i | \theta) = \prod_{j=1}^K \theta_j^{\mathbb{I}(x_{ij}=1)}$$

$$p(\mathbf{x}_{1:n} | \theta) = \prod_{i=1}^n \prod_{j=1}^K \theta_j^{\mathbb{I}(x_{ij}=1)}$$

$x_i = 001$ $x_{i1}=0$ $x_{i3}=1$

$$P(x_i=3|\theta) = \theta_1^0 \theta_2^0 \theta_3^1$$

$$= \theta_3$$

$$\theta_1 + \theta_2 + \theta_3 = 1$$

Dirichlet distribution

The conjugate prior is the **Dirichlet distribution** which is the natural generalization of the beta distribution to multiple dimensions.

The pdf is defined as follows:

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

Beta

$$\propto \frac{\theta_1^{\alpha_1 - 1} (1 - \theta_1)^{\alpha_2 - 1}}{\theta_1^{\alpha_1 - 1} \theta_2^{\alpha_2 - 2}}$$

defined on the **probability simplex**, i.e., the set of vectors such that $0 \leq \theta_k \leq 1$ and $\sum_{k=1}^K \theta_k = 1$.

In addition, $B(\alpha_1, \dots, \alpha_K)$ is the natural generalization of the beta function to K variables:

$$B(\boldsymbol{\alpha}) := \frac{\prod_{i=1}^K \Gamma(\alpha_i)}{\Gamma(\alpha_0)}$$

where $\alpha_0 := \sum_{k=1}^K \alpha_k$.

Beta Dir Die

$$\theta_1^{\alpha_1 - 1} \theta_2^{\alpha_2 - 1} \theta_3^{\alpha_3 - 1} \theta_4^{\alpha_4 - 1} \theta_5^{\alpha_5 - 1}$$

$$(1 - \theta_1 - \theta_2 - \theta_3 - \theta_4 - \theta_5)^{\alpha_6 - 1}$$

Dirichlet-categorical model

$$P(\theta | x_{1:n}) \propto P(x_{1:n} | \theta) P(\theta)$$

$$\propto \prod_{i=1}^n \prod_{j=1}^k \theta_j^{\mathbb{I}(x_{ij}=1)} \prod_{j=1}^k \theta_j^{\alpha_j - 1}$$

$$N_j = \sum_{i=1}^n \mathbb{I}(x_{ij}=1)$$

For the die

N_5 is the # times
you saw a 5.

$$= \prod_{j=1}^k \theta_j^{N_j} \prod_{j=1}^k \theta_j^{\alpha_j - 1}$$

$$= \prod_{j=1}^k \theta_j^{(N_j + \alpha_j) - 1}$$

$$\alpha'_j = N_j + \alpha_j$$

Posterior is Dirichlet \square

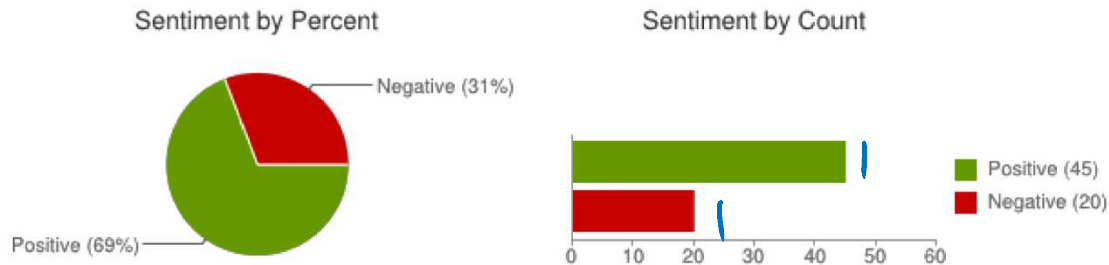
Text classification example

Sentiment140

Tweet 381 Like 173 85

obama English Search Save this search

Sentiment analysis for obama



Tweets about: obama

lillian1984: RT @BarackObama: President **Obama**: "You know that I know what real change looks like because you've seen me fight for it."

Posted 22 seconds ago

a_girl_irl: Romney endorsed by cool H'wood celebs: Kid Rock, Chuck Woolery, hot chick from Clueless, that's literally it, everyone else likes **Obama**

Posted 22 seconds ago

AmericanWoman8: RT @RBReich: If **Obama** wins, will radical right see it as a repudiation and become more reasonable, or as a provocation and grow even more extreme?

Posted 22 seconds ago

y is used to indicate C classes. E.g., the classes could be *positive*, *negative* and *neutral*. That is, $C=3$.

The input x in this example is a vector of d zeros with ones indicating which words occur in the tweet.

$x_i = \begin{bmatrix} 1 & 0 & 1 & 0 & 1 \end{bmatrix}$ tokens $d=5$
 -cut rat sat hi omg dictionary
 Nobody: omg the cut sat!

Naïve Bayes classifier $Y \in \{1, 2, \dots, C\}$

We are interested in the posterior distribution of y given the model parameters θ and π and the inputs (features) x :

$$P(y_i/x_i, \theta, \pi) = p(y_i/\pi) p(x_i/y_i, \theta) / p(x_i/\theta, \pi)$$

$$= \frac{P(y_i/\pi) P(x_i/y_i, \theta)}{\sum_{y_i=c} P(y_i=c/\pi) P(x_i/y_i=c, \theta)}$$

$$P(y_i/\pi) = \prod_{c=1}^C \pi_c \mathbb{I}(y_i=c)$$

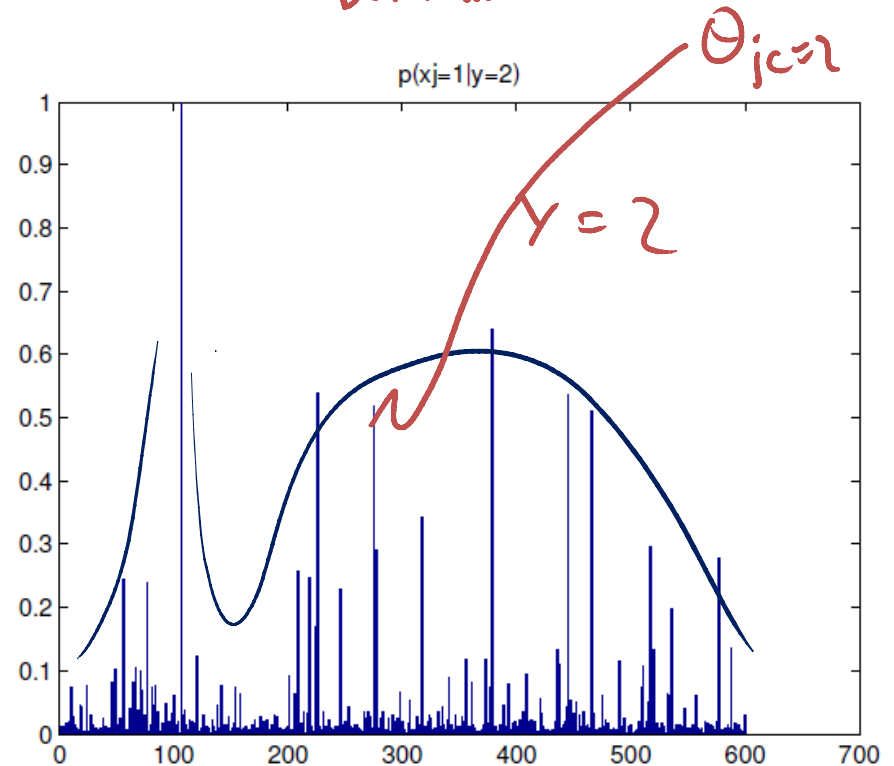
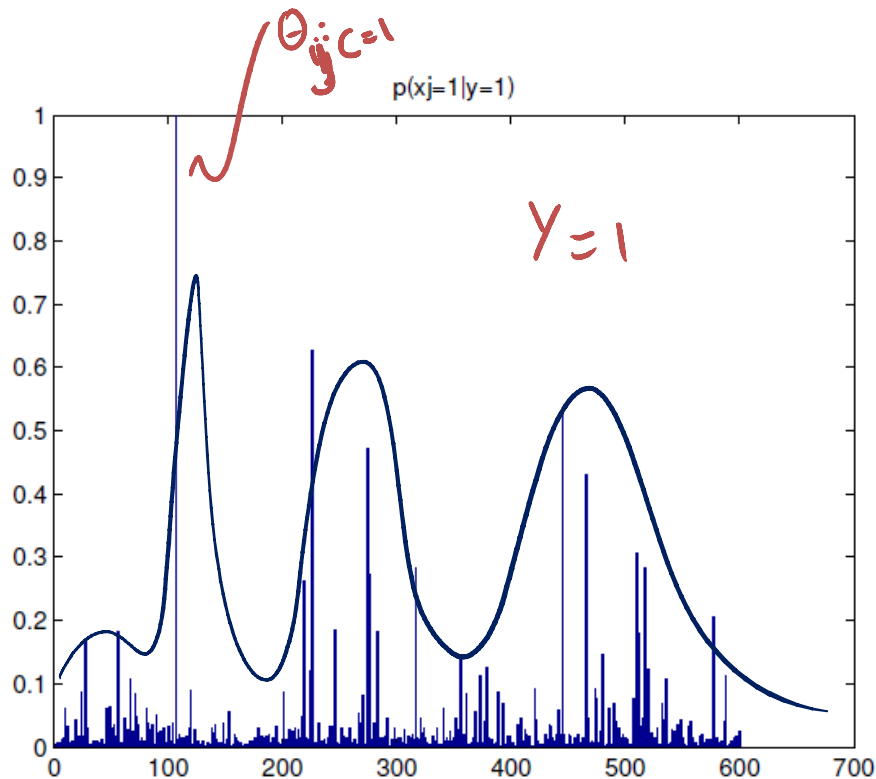
$$P(y_i=c/\pi) = \pi_c$$

$C=2$ $c=1$ positive
 $c=2$ negative

Naïve Bayes classifier

Assume the features are conditionally independent given the class label.
That is,

$$\underbrace{P(x_i | \theta, y_i = c)}_{\text{ith tweet}} = \prod_{j=1}^d \underbrace{P(x_{ij} | \theta, y_i = c)}_{\text{Bernoulli}}$$



Naïve Bayes classifier with binary features \mathbf{x}

$$P(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}, \boldsymbol{\pi}) \propto p(\mathbf{y} | \boldsymbol{\pi}) p(\mathbf{x} | \mathbf{y}, \boldsymbol{\theta})$$

$$X_i = [1 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 1]$$

$j=1 \rightarrow j=d=8$

$$p(\mathbf{y} | \boldsymbol{\pi}, \boldsymbol{\theta}, \mathbf{x}) \propto \prod_{i=1}^n \prod_{c=1}^C \left(\pi_c^{\mathbb{I}_c(y_i)} \prod_{j=1}^d \theta_{jc}^{\mathbb{I}_c(y_i) \mathbb{I}_1(x_{ij})} (1 - \theta_{jc})^{\mathbb{I}_c(y_i) \mathbb{I}_0(x_{ij})} \right)$$

$\gamma_i = \begin{cases} 1 & c=1 \\ 2 & c=2 \end{cases}$

$$P(\mathbf{y} | \boldsymbol{\pi}) = \prod_{i=1}^n \prod_{c=1}^C \pi_c^{\mathbb{I}_c(y_i)}$$

$$\longleftrightarrow P(y_i = c | \boldsymbol{\pi}) = \pi_c$$

$$P(x_{ij} | y_i = c, \boldsymbol{\theta}) = \theta_{jc}^{\mathbb{I}_c(y_i) \mathbb{I}_1(x_{ij})} (1 - \theta_{jc})^{\mathbb{I}_c(y_i) \mathbb{I}_0(x_{ij})}$$

$$\theta_{jc} = P(x_{ij} = 1 | y_i = c, \boldsymbol{\theta})$$

MLE for Naïve Bayes classifier with binary features \mathbf{x}

$$P(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}, \boldsymbol{\pi}) \propto p(\mathbf{y} | \boldsymbol{\pi}) p(\mathbf{x} | \mathbf{y}, \boldsymbol{\theta})$$

Note: $\mathbb{I}_c(y_i)$
is the same as $\mathbb{I}(y_i=c)$

$$p(\mathbf{y} | \boldsymbol{\pi}, \boldsymbol{\theta}, \mathbf{x}) \propto \prod_{i=1}^n \prod_{c=1}^C \left(\pi_c^{\mathbb{I}_c(y_i)} \prod_{j=1}^d \theta_{jc}^{\mathbb{I}_c(y_i) \mathbb{I}_1(x_{ij})} (1 - \theta_{jc})^{\mathbb{I}_c(y_i) \mathbb{I}_0(x_{ij})} \right)$$

n is the number of data points

$N_c = \sum_{i=1}^n \mathbb{I}(y_i=c)$ is the number of data of class c .

$N_{jc} = \sum_{i=1}^n \mathbb{I}_c(y_i) \mathbb{I}_1(x_{ij})$ #times $x_{ij}=1$ given we are in class c .

$$\hat{\pi}_c = N_c / n$$

$$\hat{\theta}_{jc} = N_{jc} / N_c$$

Predicting the class of new data

Given a new data point (say tweet)
 x^* , the class prediction is:

$$P(\underline{y=c} | \underline{x^*}, D) \propto \hat{\pi}_c \prod_{j=1}^d \hat{\theta}_{jc}^{\mathbb{I}(x_j^*=1)} (1 - \hat{\theta}_{jc})^{\mathbb{I}(x_j^*=0)}$$

training data

Do this for all classes and then normalize
so that

$$\sum_{c'} P(Y=c' | x^*, D) = 1$$

Naïve Bayes classifier with binary features

```
1  $N_c = 0, N_{jc} = 0;$ 
2 for  $i = 1 : \eta$  do
3    $c = y_i$  // Class label of  $i$ 'th example
4    $N_c := N_c + 1;$ 
5   for  $j = 1 : d$  do
6     if  $x_{ij} = 1$  then
7        $N_{jc} := N_{jc} + 1$ 
8  $\hat{\pi}_c = \frac{N_c}{N}, \hat{\theta}_{jc} = \frac{N_{jc}}{N_c}$ 
```

Log-sum-exp trick

$$\log p(y = c | \mathbf{x}) = b_c - \log \left[\sum_{c'=1}^C e^{b_{c'}} \right]$$
$$b_c := \log p(\mathbf{x} | y = c) + \log p(y = c)$$

$$\log \left[\sum_{c'} e^{b_{c'}} \right] = \log \sum_{c'} p(y = c', \mathbf{x}) = \log p(\mathbf{x}) \quad \text{log sum exp function}$$

$$\log \sum_c e^{b_c} = \log \left[\left(\sum_c e^{b_c - B} \right) e^B \right] = \left[\log \left(\sum_c e^{b_c - B} \right) \right] + B \quad \text{where } B = \max_c b_c.$$

For example,

$$\log(\underline{e^{-120}} + \underline{e^{-121}}) = \log(\underline{e^{-120}}(e^0 + e^{-1})) = \log(e^0 + e^{-1}) - \underline{120}$$

NBC prediction with log-sum-exp trick

```
1 for  $i = 1 : \tilde{n}$  do
2   for  $c = 1 : C$  do
3      $L_{ic} = \log \hat{\pi}_c;$ 
4     for  $j = 1 : \mathbf{d}$  do
5       if  $x_{ij} = 1$  then  $L_{ic} := L_{ic} + \log \hat{\theta}_{jc}$  else  $L_{ic} := L_{ic} + \log(1 - \hat{\theta}_{jc})$ 
6      $p_{ic} = \exp(L_{ic} - \text{logsumexp}(L_{i,:}));$ 
7      $\hat{y}_i = \text{argmax}_c p_{ic};$ 
```

$$2^3 2^4 = 2^{3+4}$$

MLE

$$\prod_i \prod_c \prod_c^{\mathbb{I}_c(y_i)} = \prod_c \prod_c^{\sum_i \mathbb{I}_c(y_i)}$$

$$P(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}, \boldsymbol{\pi}) \propto p(\mathbf{y} | \boldsymbol{\pi}) p(\mathbf{x} | \mathbf{y}, \boldsymbol{\theta})$$

$$p(\mathbf{y} | \boldsymbol{\pi}, \boldsymbol{\theta}, \mathbf{x}) \propto \prod_{i=1}^n \prod_{c=1}^C \left(\pi_c^{\mathbb{I}_c(y_i)} \prod_{j=1}^d \theta_{jc}^{\mathbb{I}_c(y_i) \mathbb{I}_1(x_{ij})} (1 - \theta_{jc})^{\mathbb{I}_c(y_i) \mathbb{I}_0(x_{ij})} \right)$$

$$= \prod_{c=1}^C \pi_c^{N_c} \prod_{j=1}^d \theta_{jc}^{N_{jc}} (1 - \theta_{jc})^{N_c - N_{jc}}$$

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\pi}) = \underbrace{\left(\sum_{c=1}^C N_c \log \pi_c \right)} + \underbrace{\left(\sum_{c=1}^C \sum_{j=1}^d N_{jc} \log \theta_{jc} + (N_c - N_{jc}) \log (1 - \theta_{jc}) \right)}$$

MLE for π

$$l(\pi, \lambda) = \left(\sum_{c=1}^C N_c \log \pi_c \right) + \lambda \left[1 - \sum_{c=1}^C \pi_c \right]$$

$$\frac{\partial l(\pi, \lambda)}{\partial \lambda} = 0 + \left[1 - \sum_{c=1}^C \pi_c \right] \stackrel{=0}{\rightarrow} \boxed{\sum_{c=1}^C \pi_c = 1} \text{ (1)}$$

$$\frac{\partial l(\pi, \lambda)}{\partial \pi_c} = N_c \frac{1}{\pi_c} + \lambda (-1) \stackrel{=0}{\rightarrow} \frac{N_c}{\pi_c} = \lambda$$

$$\sum_{c=1}^C N_c = \sum_{c=1}^C \pi_c \lambda$$
$$n = \lambda \left[\sum_{c=1}^C \pi_c \right] = \lambda \cdot 1$$

$$\lambda = n$$
$$N_c = \lambda \pi_c$$
$$\therefore \pi_c = N_c / n$$

MLE for π

MLE for θ

$$l(\theta) = \sum_{j=1}^d \sum_{c=1}^C N_{jc} \log \theta_{jc} + (N_c - N_{jc}) \log (1 - \theta_{jc})$$

$$\begin{aligned} \frac{\partial l(\theta)}{\partial \theta_{jc}} &= N_{jc} \frac{1}{\theta_{jc}} + (N_c - N_{jc}) \frac{-1}{1 - \theta_{jc}} \\ &= \left[N_{jc} (1 - \theta_{jc}) - (N_{jc} - N_c) \theta_{jc} \right] / \theta_{jc} (1 - \theta_{jc}) \end{aligned}$$

Equating to zero:

$$\begin{aligned} N_{jc} - \cancel{N_{jc}} \theta_{jc} + \cancel{N_{jc}} \theta_{jc} - N_c \theta_{jc} &= 0 \\ \hat{\theta}_{jc} &= N_{jc} / N_c \end{aligned}$$

Bayesian analysis

Likelihood:

$$p(\mathbf{y} | \boldsymbol{\pi}, \boldsymbol{\theta}, \mathbf{X}) \propto \prod_{i=1}^n \prod_{c=1}^C \left(\pi_c^{\mathbb{I}_c(y_i)} \prod_{j=1}^d \theta_{jc}^{\mathbb{I}_c(y_i) \mathbb{I}_1(x_{ij})} (1 - \theta_{jc})^{\mathbb{I}_c(y_i) \mathbb{I}_0(x_{ij})} \right)$$
$$= \prod_{c=1}^C \left(\pi_c^{N_c} \prod_{j=1}^d \theta_{jc}^{N_{jc}} (1 - \theta_{jc})^{N_c - N_{jc}} \right)$$

Prior:

$$P(\boldsymbol{\pi}) \propto \prod_{c=1}^C \pi_c^{\alpha_c - 1} \quad \text{Dirichlet}$$

$$P(\theta_{jc}) = \theta_{jc}^{\beta_1 - 1} (1 - \theta_{jc})^{\beta_2 - 1} \quad \text{dxC beta priors}$$

$$P(\boldsymbol{\pi}, \boldsymbol{\theta}) = \prod_{c=1}^C \left(\pi_c^{\alpha_c - 1} \prod_{j=1}^d \theta_{jc}^{\beta_1 - 1} (1 - \theta_{jc})^{\beta_2 - 1} \right)$$

Bayesian analysis

$$P(\theta, \pi | x, y) \propto \prod_{c=1}^C \pi_c^{N_c + \alpha_c - 1} \prod_{j=1}^d \theta_{jc}^{\beta_1 + N_{jc} - 1} (1 - \theta_{jc})^{\beta_2 + N_c - N_{jc} - 1}$$

↑
Dirichlet !

$$\text{Post. mean } (\pi) = \frac{N_c + \alpha_c}{n + \sum_{c=1}^C \alpha_c}$$

$$\text{post. mean } (\theta_{jc}) = \frac{N_{jc} + \beta_1}{N_c + \beta_1 + \beta_2}$$

Next lecture

In the next lecture, we learn about another very popular classifier: logistic regression. This classifier will be a building block for neural networks.