

Naive Bayes

DRIPTA MJ

Department of Mathematics

RAMAKRISHNA MISSION VIVEKANANDA EDUCATIONAL AND RESEARCH INSTITUTE
BELUR MATH, INDIA

Machine Learning

DA 220

Sem 2, 2019-20

Introduction

- A Bayes rule based simple classification method.
- Naive because features are not usually conditionally independent.
- Popular in the field of natural language processing.
- Example: Classify an email as spam or not spam.

Example

- Want to classify an email as spam or not-spam.
- Outputs: $y \in \{c_1, c_2\}$, where c_1 indicates email is not-spam and c_2 spam.
- Training dataset: N emails (say 10000).
- Suppose we have a vocabulary comprising 5000 words.
- Dimension D of input \mathbf{x} = Size of the vocabulary (=5000).
- Each email is described in terms of the vocabulary. For example, the inputs of the n th email can be expressed as

$$\mathbf{x}^{(n)} = [x_1^{(n)}, x_2^{(n)}, \dots, x_D^{(n)}]^T$$

where

$$x_j^{(n)} = \begin{cases} 1 & \text{if the } j\text{th word of vocabulary exists in the } n\text{th email} \\ 0 & \text{otherwise} \end{cases}$$

Example

- Suppose vocabulary \mathcal{V} is of the form

$$\mathcal{V} = \left\{ \begin{array}{c} \text{abandon} \\ \text{ability} \\ \cdot \\ \cdot \\ \text{magic} \\ \cdot \\ \cdot \\ \text{zoo} \end{array} \right\}$$

then the feature vector

$$\mathbf{x} = \left[\begin{array}{c} 0 \\ 1 \\ \cdot \\ \cdot \\ 1 \\ \cdot \\ \cdot \\ 0 \end{array} \right]$$

indicates that the email contains the words **ability**, **magic**, but not **abandon** and **zoo**.

Bernoulli distribution

- Suppose x is a discrete random variable and $x \in \{0, 1\}$.
- Consider a parameter μ such that $0 \leq \mu \leq 1$.
- Let x take the value 1 with probability μ .
- Then x takes the value 0 with probability $1 - \mu$.
- The probability mass function of x can be written as

$$p(x) = \begin{cases} \mu & \text{if } x = 1 \\ 1 - \mu & \text{if } x = 0 \end{cases}$$

- The probability mass function $p(x)$ can be compactly written as

$$p(x) = \mu^x (1 - \mu)^{(1-x)}$$

- This is a Bernoulli distribution with parameter μ .

Class conditional density

- Consider a dataset comprising N data points with D features:
 - Inputs: $\mathbf{X} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$
 - Outputs: $\mathbf{y} = \{y^{(1)}, y^{(2)}, \dots, y^{(N)}\}$
- Suppose there are J output classes i.e. $y^{(n)} \in \{c_1, c_2, \dots, c_J\}$.
Note: In the slides $y^{(n)} \in \{1, 2, \dots, J\}$ is also used to indicate class.
- Features are assumed to conditionally independent, so the model is called “naive”.
- Class conditional densities as product of one-dimensional densities:

$$\begin{aligned} p(\mathbf{x}|y = c_j) &= p(x_1|y = c_j) p(x_2|y = c_j) \dots p(x_D|y = c_j) \\ &= \prod_{i=1}^D p(x_i|y = c_j) \end{aligned}$$

Class conditional distribution

- If we use a Bernoulli distribution to model $p(x_i|y = c_j)$ then

$$p(x_i|y = c_j) = \mu_{ij}^{x_i} (1 - \mu_{ij})^{(1-x_i)}$$

– Intuitively, μ_{ij} is the probability of the i th feature belonging to j th class.

- Therefore the class-conditional distribution can be written as

$$p(\mathbf{x}|y = c_j) = \prod_{i=1}^D \mu_{ij}^{x_i} (1 - \mu_{ij})^{(1-x_i)}$$

Likelihood function

$$\begin{aligned}p(\mathbf{y}|\mathbf{X}) &\propto p(\mathbf{y})p(\mathbf{X}|\mathbf{y}) = \prod_{n=1}^N p(y^{(n)})p(\mathbf{x}^{(n)}|y^{(n)}) \\&= \prod_{n=1}^N p(y^{(n)}) \prod_{i=1}^D p(x_i^{(n)}|y^{(n)}) \\&= \prod_{n=1}^N p(y^{(n)}) \prod_{i=1}^D \mu_{iy^{(n)}}^{x_i^{(n)}} (1 - \mu_{iy^{(n)}})^{(1-x_i^{(n)})} \\&= \prod_{n=1}^N \pi_{y^{(n)}} \prod_{i=1}^D \mu_{iy^{(n)}}^{x_i^{(n)}} (1 - \mu_{iy^{(n)}})^{(1-x_i^{(n)})} \\&= L(\boldsymbol{\theta})\end{aligned}$$

(where $\boldsymbol{\theta} = [\boldsymbol{\mu}, \boldsymbol{\pi}]$ are parameters of the model)

Log likelihood

$$\log L(\boldsymbol{\theta}) = \sum_{n=1}^N \left[\log \pi_{y^{(n)}} + \sum_{i=1}^D \log \mu_{iy^{(n)}}^{x_i^{(n)}} (1 - \mu_{iy^{(n)}})^{(1-x_i^{(n)})} \right]$$

$$= \sum_{n=1}^N \left[\log \pi_{y^{(n)}} + \sum_{i=1}^D \log \mu_{iy^{(n)}}^{x_i^{(n)}} + \sum_{i=1}^D \log (1 - \mu_{iy^{(n)}})^{(1-x_i^{(n)})} \right]$$

$$= \sum_{n=1}^N \left[\log \pi_{y^{(n)}} + \sum_{i=1}^D x_i^{(n)} \log \mu_{iy^{(n)}} + \sum_{i=1}^D (1 - x_i^{(n)}) \log (1 - \mu_{iy^{(n)}}) \right]$$

Maximum Likelihood Estimator

- Evaluate $\arg \max_{\theta} \log L(\theta)$ subject to $\sum_j \pi_j = 1$.
- Taking derivative with respect to μ_{ij}

$$\begin{aligned}\frac{\partial \log L(\theta)}{\partial \mu_{ij}} &= 0 \\ \Rightarrow \sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)} \left(\frac{x_i^{(n)}}{\mu_{ij}} - \frac{(1 - x_i^{(n)})}{(1 - \mu_{ij})} \right) &= 0 \\ \Rightarrow \sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)} \left[x_i^{(n)}(1 - \mu_{ij}) - (1 - x_i^{(n)})\mu_{ij} \right] &= 0 \\ \Rightarrow \sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)} \mu_{ij} &= \sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)} x_i^{(n)}\end{aligned}$$

$$\mu_{ij} = \frac{\sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)} x_i^{(n)}}{\sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)}}$$

Maximum Likelihood Estimator

- To derive π , consider the Lagrangian formulation:

$$\mathcal{L} = \log L(\boldsymbol{\theta}) + \lambda \left(\sum_j \pi_j - 1 \right)$$

- Taking derivative with respect to π_j

$$\frac{\partial \log L(\boldsymbol{\theta})}{\partial \pi_j} + \lambda \frac{\partial \sum_j \pi_j}{\partial \pi_j} = 0$$

$$\Rightarrow \lambda = - \sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)} \frac{1}{\pi_j}$$

$$\pi_j = - \frac{\sum_{i=1}^N \mathbb{1}_{(y^{(n)}=c_j)}}{\lambda}$$

Now

$$\sum_j \pi_j = 1 \quad \Rightarrow \quad - \frac{\sum_j \sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)}}{\lambda} = 1$$

$$\pi_j = \frac{\sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)}}{N}$$

Interpretation of parameters

$$\mu_{ij} = \frac{\sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)} x_i^{(n)}}{\sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)}}$$

- μ_{ij} is the fraction of the data points assigned class c_j in which the i th feature occurs.

$$\pi_j = \frac{\sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)}}{N}$$

- π_j is the fraction of the training dataset assigned to the j th output class.

Prediction

- Eventually want to predict class for a new (unobserved) example with features (say) \mathbf{x}^* .
- Compute the posterior probability using Baye's rule. For example, the probability that \mathbf{x}^* belongs to the j th class can be computed as

$$\begin{aligned} p(y^* = c_j | \mathbf{x}^*) &= \frac{p(\mathbf{x}^* | y^* = c_j) p(y^* = c_j)}{p(\mathbf{x}^*)} \\ &= \frac{\prod_{i=1}^D p(x_i^* | y^* = c_j) \pi_j}{\sum_{j=1}^J \prod_{i=1}^D p(x_i^* | y^* = c_j) \pi_j} \end{aligned}$$

- \mathbf{x}^* is assigned to the class which has the highest posterior probability.

Email – spam or not-spam?

- Goal: For a new email with features \mathbf{x}^* , predict the output class y^* .
 - Estimate μ_{i1} s and μ_{i2} s.
 - μ_{i1} is the fraction of the non-spam emails in which the i th word of the vocabulary occurs.
 - μ_{i2} is the fraction of the spam emails in which the i th word of the vocabulary occurs.
 - Estimate π_1 and π_2 .
 - π_1 is the fraction of emails which are not spam.
 - π_2 is the fraction of emails which are spam.
- Note: All the above estimations are made using the training dataset.
- Use the trained parameters to classify new emails.

GAUSSIAN MODEL

Gaussian model

- Suppose $p(\mathbf{x}|y = c_j)$ is the class likelihood and $p(y = c_j)$ is the prior of the j^{th} output class, then using Bayes rule we have

$$p(y = c_j|\mathbf{x}) = \frac{p(\mathbf{x}|y = c_j)p(y = c_j)}{\sum_{j=1}^M p(\mathbf{x}|y = c_j)p(y = c_j)} \\ \propto p(\mathbf{x}|y = c_j)p(y = c_j)$$

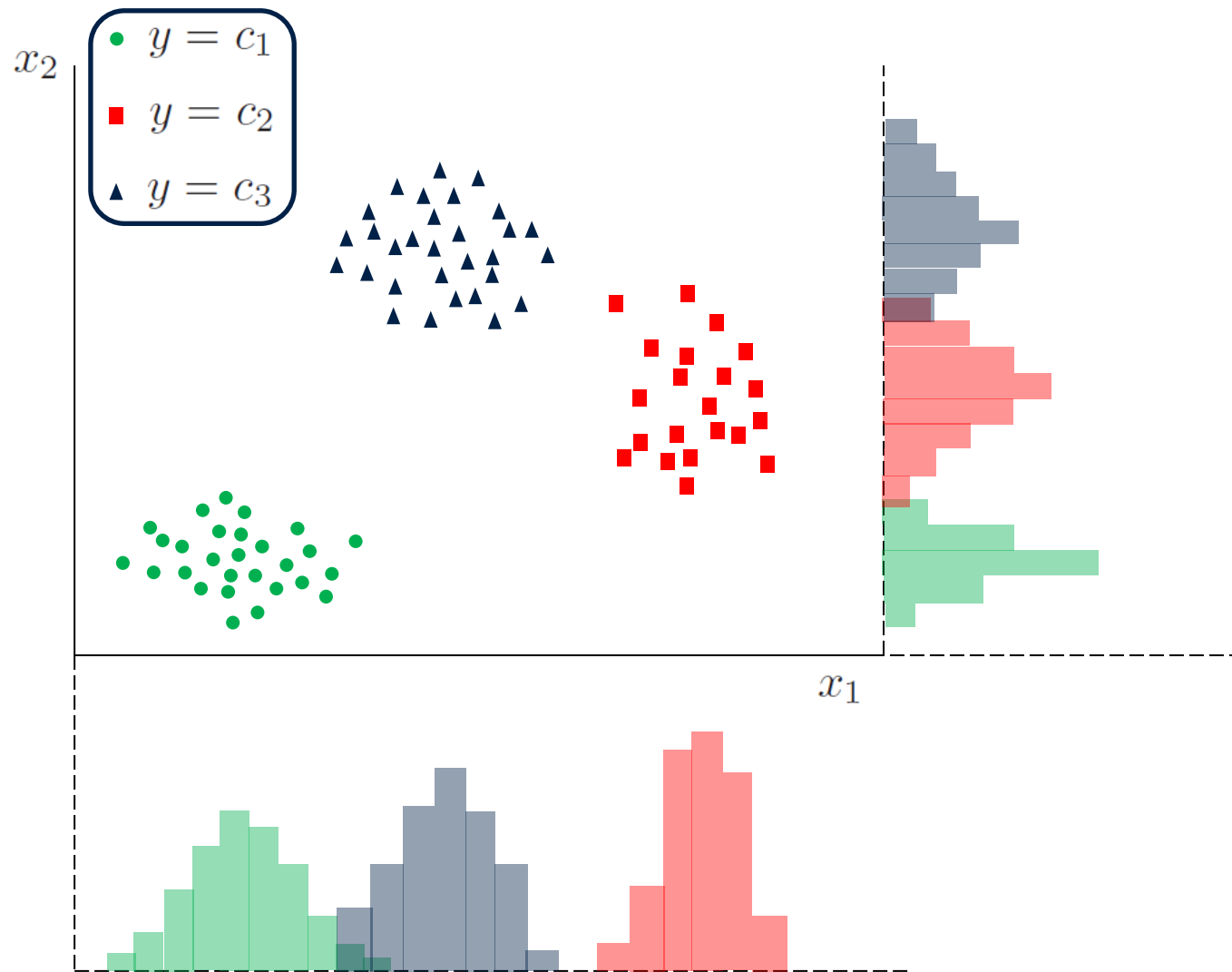
- Assuming conditional independence of the features given the class yields

$$p(\mathbf{x}|y = c_j) = \prod_{i=1}^D p(x_i|y = c_j) \\ = \prod_{i=1}^D \mathcal{N}(x_i|\mu_{ij}, \sigma_{ij}^2)$$

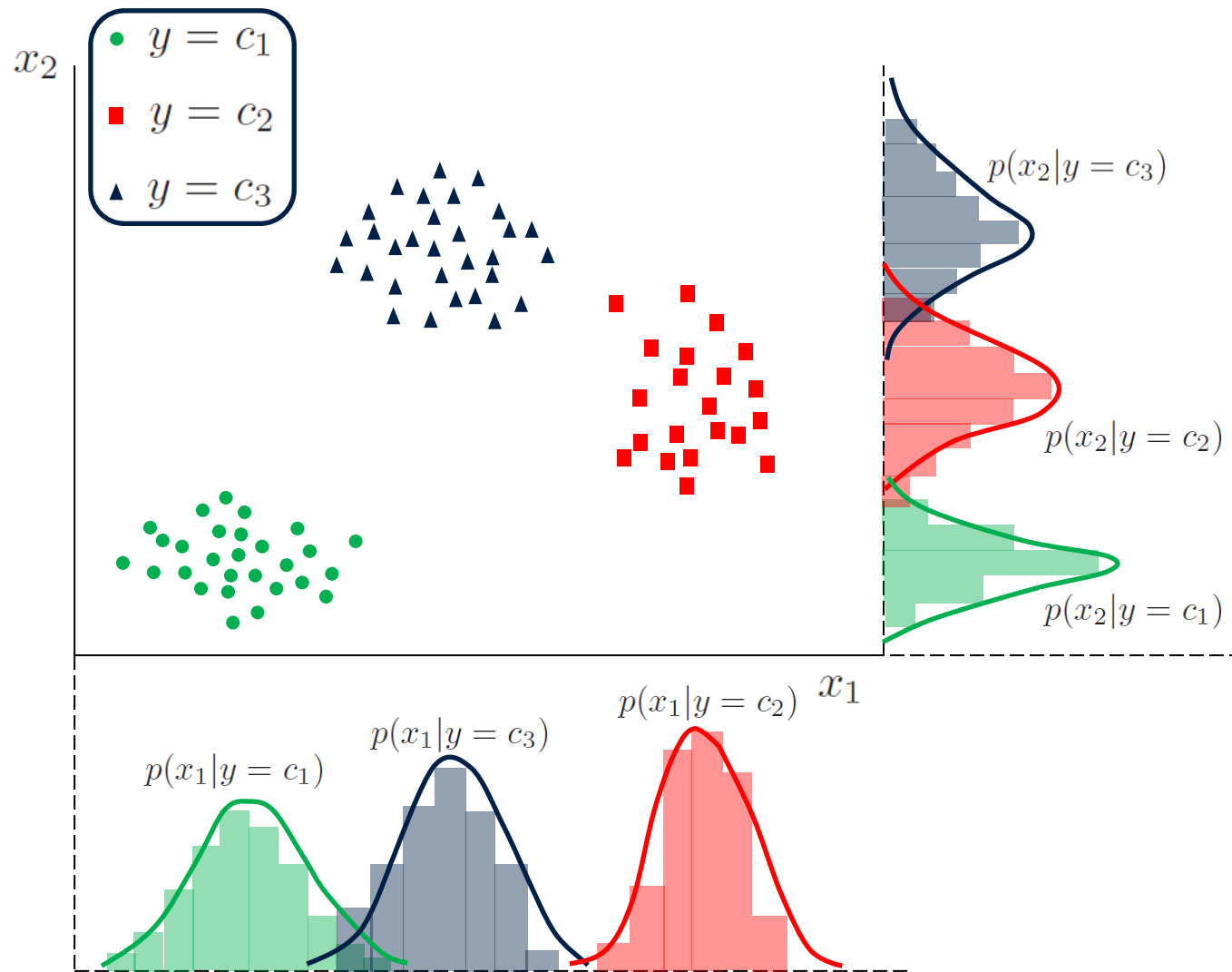
$$p(x_i|y = c_j) = \frac{1}{\sqrt{2\pi\sigma_{ij}^2}} \exp\left(-\frac{(x_i - \mu_{ij})^2}{2\sigma_{ij}^2}\right)$$

- μ_{ij} is the mean of the i th feature belonging to class c_j .
- σ_{ij} is the variance of the i th feature belonging to class c_j .
- The model parameters are $\theta = [\mu, \sigma, \pi]$.

Gaussian model



Gaussian model



Gaussian model – likelihood

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}) &\propto p(\mathbf{y})p(\mathbf{X}|\mathbf{y}) = \prod_{n=1}^N p(y^{(n)})p(\mathbf{x}^{(n)}|y^{(n)}) \\ &= \prod_{n=1}^N p(y^{(n)}) \prod_{i=1}^D p(x_i^{(n)}|y^{(n)}) \\ &= \prod_{n=1}^N p(y^{(n)}) \prod_{i=1}^D \frac{1}{\sqrt{2\pi}\sigma_{iy^{(n)}}} \exp \left[-\frac{(x_i - \mu_{iy^{(n)}})^2}{2\sigma_{iy^{(n)}}^2} \right] \\ &= \prod_{n=1}^N \pi_{y^{(n)}} \prod_{i=1}^D \frac{1}{\sqrt{2\pi}\sigma_{iy^{(n)}}} \exp \left[-\frac{(x_i - \mu_{iy^{(n)}})^2}{2\sigma_{iy^{(n)}}^2} \right] \\ &= L(\boldsymbol{\theta}) \\ &\quad (\text{where } \boldsymbol{\theta} = [\boldsymbol{\mu}, \boldsymbol{\sigma}, \boldsymbol{\pi}] \text{ are parameters of the model}) \end{aligned}$$

Gaussian model – log likelihood

$$L(\boldsymbol{\theta}) = \prod_{n=1}^N \pi_{y(n)} \prod_{i=1}^D \frac{1}{\sqrt{2\pi}\sigma_{iy(n)}} \exp \left[-\frac{(x_i - \mu_{iy(n)})^2}{2\sigma_{iy(n)}^2} \right]$$

- Taking logarithm on both sides

$$\begin{aligned} \log L(\boldsymbol{\theta}) &= \sum_{n=1}^N \left[\log \pi_{y(n)} + \sum_{i=1}^D \log \left(\frac{1}{\sqrt{2\pi}\sigma_{iy(n)}} \right) - \sum_{i=1}^D \frac{(x_i - \mu_{iy(n)})^2}{2\sigma_{iy(n)}^2} \right] \\ &= \sum_{n=1}^N \left[\log \pi_{y(n)} - \sum_{i=1}^D \log (\sqrt{2\pi}\sigma_{iy(n)}) - \sum_{i=1}^D \frac{(x_i - \mu_{iy(n)})^2}{2\sigma_{iy(n)}^2} \right] \end{aligned}$$

Gaussian model

- Evaluate $\arg \max_{\boldsymbol{\theta}} \log L(\boldsymbol{\theta})$ subject to $\sum_j \pi_j = 1$.
- Taking derivative with respect to μ_{ij}

$$\frac{\partial \log L(\boldsymbol{\theta})}{\partial \mu_{ij}} = 0$$

$$\Rightarrow \sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)} \left(\frac{(x_i^{(n)} - \mu_{ij})}{\sigma_{ij}^2} \right) = 0$$

$$\mu_{ij} = \frac{\sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)} x_i^{(n)}}{\sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)}}$$

Gaussian model

- Taking derivative with respect to σ_{ij}

$$\frac{\partial \log L(\boldsymbol{\theta})}{\partial \sigma_{ij}} = 0$$
$$\Rightarrow \sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)} \left[-\frac{\sqrt{2\pi}}{\sqrt{2\pi}\sigma_{ij}} - \frac{(x_i^{(n)} - \mu_{ij})^2}{2} \left(\frac{-2}{\sigma_{ij}^3} \right) \right] = 0$$
$$\Rightarrow \sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)} \left[\frac{(x_i^{(n)} - \mu_{ij})^2}{\sigma_{ij}^2} - 1 \right] = 0$$

$$\sigma_{ij}^2 = \frac{\sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)} (x_i^{(n)} - \mu_{ij})^2}{\sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)}}$$

- π_j 's can be estimated using a Lagrange formulation (as in case of Bernoulli model)

$$\pi_j = \frac{\sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)}}{N}$$

MULTINOMIAL MODEL

Shortcomings of Bernoulli model

- In the Bernoulli model, the feature vector (say \mathbf{x}) captures the presence or absence of features in the dataset.
- Such representation cannot capture the frequency of features in an example, i.e. multiple occurrence of a particular feature in the same example is not accounted.
- Multinomial model addressed this issue.
- First a look at multinomial coefficient....

Multinomial coefficient

- Number of distinct ways to permute a set of items (say a multiset of m elements).
 - Consider the word **KOLKATA**.
 - The word has 2 **K**s, 1 **O**, 1 **L**, 2 **A**s and 1 **T**.
 - The number of unique permutations of letters of the word:

$$\frac{7!}{2! 1! 1! 2! 1!}$$

Multinomial coefficient

- Suppose there are m items of D number of types.
 - Let us say there are z_i number of items of type i with $i = 1, 2, \dots, D$.
 - So $z_1 + z_2 + \dots + z_D = m$.

- The number of unique arrangements of the items is given by

$$\frac{m!}{z_1! z_2! \dots z_D!}$$

- This is the general form of the multinomial coefficient.
- Also, proportion of items of type i is

$$p_i = \frac{z_i}{m}$$

and therefore

$$\sum_{i=1}^D p_i = 1$$

Multinomial distribution

- Suppose m items are selected at random from a set comprising D types of items.
- Also suppose there are z_i number of items of kind i , such that

$$z_1 + z_2 + \dots + z_D = m$$

- Let the probability of the i^{th} kind of item be p_i , and so $\sum_{i=1}^D p_i = 1$.
- The probability distribution of $\mathbf{z} = [z_1, z_2, \dots, z_D]^T$ is given by

$$p(\mathbf{z}) = \frac{m!}{z_1! z_2! \dots z_D!} (p_1)^{z_1} (p_2)^{z_2} \dots (p_D)^{z_D}$$

- The product of the probabilities indicates the probability of having z_i number of items with feature i , with $i = 1, \dots, D$.
- The multinomial coefficient indicates the number of all possible combinations with features having counts $\mathbf{z} = [z_1, z_2, \dots, z_D]^T$.

Multinomial model – example

- Consider the problem of email classification: spam or not-spam.
- Suppose you are given a set of N emails which are already classified.
- The dataset can be represented as $\{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$.
 - The n^{th} input data point is $\mathbf{x}^{(n)} = [x_1^{(n)}, x_2^{(n)}, \dots, x_D^{(n)}]^T$.
 - $\{x_1, x_2, \dots, x_D\}$ are the features of the model.
 - The features correspond to a vocabulary (dictionary) of words.
 - $x_i^{(n)}$ indicates the number of times the i th word occurs in the n th email.
 - The outputs are either spam $y^{(n)} = 0$ or not-spam $y^{(n)} = 1$.

Multinomial model – likelihood function

$$p(\mathbf{y}|\mathbf{X}) \propto p(\mathbf{y})p(\mathbf{X}|\mathbf{y}) = \prod_{n=1}^N p(y^{(n)})p(\mathbf{x}^{(n)}|y^{(n)})$$

- The distribution $p(\mathbf{x}^{(n)}|y^{(n)})$ can be written as a multinomial distribution

$$p(\mathbf{x}^{(n)}|y^{(n)}) = \frac{m!}{x_1^{(n)}! x_2^{(n)}! \dots x_D^{(n)}!} (p(x_1|y^{(n)}))^{x_1^{(n)}} (p(x_2|y^{(n)}))^{x_2^{(n)}} \dots (p(x_D|y^{(n)}))^{x_D^{(n)}}$$

- Taking $p(x_i|y^{(n)}) = \mu_{iy^{(n)}}$, can write

$$p(\mathbf{x}^{(n)}|y^{(n)}) = \frac{m!}{x_1^{(n)}! x_2^{(n)}! \dots x_D^{(n)}!} \mu_{1y^{(n)}}^{x_1^{(n)}} \mu_{2y^{(n)}}^{x_2^{(n)}} \dots \mu_{Dy^{(n)}}^{x_D^{(n)}}$$

- Since the multinomial coefficient in front is a constant, can write

$$p(\mathbf{x}^n|y^{(n)}) \propto \prod_{i=1}^D \mu_{iy^{(n)}}^{x_i^{(n)}}$$

Multinomial model – likelihood function

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}) &\propto p(\mathbf{y})p(\mathbf{X}|\mathbf{y}) = \prod_{n=1}^N p(y^{(n)})p(\mathbf{x}^{(n)}|y^{(n)}) \propto \prod_{n=1}^N p(y^{(n)}) \prod_{i=1}^D \mu_{iy^{(n)}}^{x_i^{(n)}} \\ &= \prod_{n=1}^N \pi_{y^{(n)}} \prod_{i=1}^D \mu_{iy^{(n)}}^{x_i^{(n)}} \\ &= L(\boldsymbol{\theta}) \end{aligned}$$

- $\boldsymbol{\theta} = [\boldsymbol{\mu}, \boldsymbol{\pi}]$ are the parameters of the model.
- Log likelihood:

$$\begin{aligned} \log L(\boldsymbol{\theta}) &= \sum_{n=1}^N \left[\log \pi_{y^{(n)}} + \sum_{i=1}^D \log \mu_{iy^{(n)}}^{x_i^{(n)}} \right] \\ &= \sum_{n=1}^N \left[\log \pi_{y^{(n)}} + \sum_{i=1}^D x_i^{(n)} \log \mu_{iy^{(n)}} \right] \end{aligned}$$

Multinomial model – MLE

- Evaluate $\arg \max_{\theta} \log L(\theta)$

$$\text{subject to } \sum_{i=1}^D \mu_{ij} = 1, \quad \forall j = 1, 2, \dots, M$$

$$\sum_{j=1}^M \pi_j = 1$$

- Method of Lagrange multipliers:

$$\begin{aligned} \mathcal{L} = & \sum_{n=1}^N \left[\log \pi_{y(n)} + \sum_{i=1}^D x_i^{(n)} \log \mu_{iy(n)} \right] \\ & + \alpha_1 \left(\sum_{i=1}^D \mu_{i1} - 1 \right) + \alpha_2 \left(\sum_{i=1}^D \mu_{i2} - 1 \right) + \dots + \alpha_M \left(\sum_{i=1}^D \mu_{iM} - 1 \right) \\ & + \beta \left(\sum_{j=1}^M \pi_j - 1 \right) \end{aligned}$$

Model parameters

- Taking derivative with respect to μ_{ij} :

$$\frac{\partial \mathcal{L}}{\partial \mu_{ij}} = 0$$

$$\Rightarrow \sum_{n=1}^N \left[\mathbb{1}_{(y^{(n)}=c_j)} \left(\frac{x_i^{(n)}}{\mu_{ij}} \right) \right] + \alpha_j = 0$$

$$\Rightarrow \alpha_j = - \sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)} \frac{x_i^{(n)}}{\mu_{ij}}$$

$$\Rightarrow \mu_{ij} = - \sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)} \frac{x_i^{(n)}}{\alpha_j}$$

But α_j is unknown.

Model parameters

- Substituting μ_{ij} in

$$\sum_{i=1}^D \mu_{ij} = 1$$

we have

$$\Rightarrow - \sum_{i=1}^D \sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)} \frac{x_i^{(n)}}{\alpha_j} = 1$$

$$\Rightarrow - \sum_{i=1}^D \sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)} x_i^{(n)} = \alpha_j$$

$$\Rightarrow - \sum_{i=1}^D \sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)} x_i^{(n)} = - \sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)} \frac{x_i^{(n)}}{\mu_{ij}}$$

$$\mu_{ij} = \frac{\sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)} x_i^{(n)}}{\sum_{i=1}^D \sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)} x_i^{(n)}}$$

Model parameters

- π_j 's can be derived in a similar way by differentiating \mathcal{L} with respect to π_j , which gives

$$\pi_j = \frac{\sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)}}{N}$$

- Use the trained parameters to predict the output class of a new (unobserved) example.

- Posterior probability:

$$\begin{aligned} p(c_j | \mathbf{x}^*) &\propto p(c_j) \prod_{i=1}^D p(x_i^* | c_j) \\ &\propto \pi_j \prod_{i=1}^D \mu_{ij}^{x_i^*} \end{aligned}$$

- **Problem:** If any $p(x_i^* | c_j) = 0$, then the posterior probability becomes 0.
- Example:
 - Consider the problem of classifying documents between two classes: **Science** and **Sports**.
 - Suppose the word “laser” is there in the vocabulary, but not present in the training dataset.
 - If the word “laser” occurs in a test document, then the posterior probabilities of both classes become equal to 0.
- A word not present in the training dataset of a particular document class does not imply that the word cannot occur in any document belonging to that class.

Laplace smoothing

- Derived earlier:

$$\mu_{ij} = \frac{\sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)} x_i^{(n)}}{\sum_{i=1}^D \sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)} x_i^{(n)}}$$

- Laplace Smoothing:** Add 1 to the numerator and D (size of the vocabulary) to the denominators.

$$\mu_{ij} = \frac{\sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)} x_i^{(n)} + 1}{\sum_{i=1}^D \sum_{n=1}^N \mathbb{1}_{(y^{(n)}=c_j)} x_i^{(n)} + D}$$

- $\sum_{i=1}^D \mu_{ij} = 1$ still holds.

Term frequency (tf)

- Consider two documents on the same topic – Doc-1 and Doc-2.
 - Doc-1 is much longer than Doc-2.
 - Therefore the average count values in Doc-1 will be much larger than in Doc-2.
 - For such a case occurrence count may not be a good choice.
- Solution: Divide the number of occurrences of each word in a document by the total number of words in the document.
- The new features are known as term frequencies (tf).

Inverse document frequency (idf) and tf-idf

- Some words in the vocabulary can occur frequently in different types of documents.
- These words are less informative as they don't provide any extra information to that enable distinction between document classes.
- Solution: Downscale the weights assigned to such words using inverse document frequency (idf).
 - Let the document frequency df_j be defined as the number of documents that contain the j th word of the vocabulary.
 - Let N be the total number of documents in the collection.
 - The inverse document frequency idf_j of the j th word is defined as

$$idf_j = \log \frac{N}{df_j}$$

- idf is usually taken as a refinement on top of tf.
- The combination known as tf-idf is computed as $tf-idf = tf \times idf$.

N-gram

- N -gram is a sequence of N words.
 - Bigram (2-gram) is a two-word sequence.
 - Trigram (3-gram) is a three-word sequence.

- Example: “We are visiting Belur Math today”

- Bigram:

```
[('We', 'are'), ('are', 'visiting'), ('visiting', 'Belur'),  
 ('Belur', 'Math'), ('Math', 'today')]
```

- Trigram:

```
[('We', 'are', 'visiting'), ('are', 'visiting', 'Belur'),  
 ('visiting', 'Belur', 'Math'), ('Belur', 'Math', 'today')]
```

NLTK CountVectorizer

- CountVectorizer can build a dictionary of features.
 - Enables transformation of documents to feature vectors.
- Using CountVectorizer text document collection can be converted into a matrix of token counts.
 - Can also specify an a-priori dictionary or some specific form of analyzer for feature selection.
 - Otherwise the number of features will be equal to the vocabulary size of the dataset.