

# Learning From Data

## Lecture 4: Generative Learning Algorithms

Yang Li [yangli@sz.tsinghua.edu.cn](mailto:yangli@sz.tsinghua.edu.cn)

October 9, 2021

# Introduction

# Today's Lecture

## Supervised Learning (Part III)

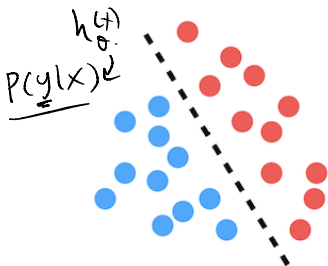
- ▶ Discriminative & Generative Models
- ▶ Gaussian Discriminant Analysis

## Discriminative & Generative Models

# Two Learning Approaches

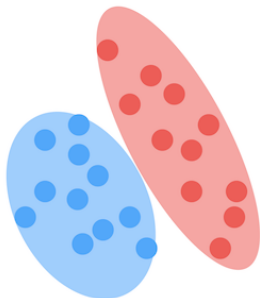
Classify input data  $x$  into two classes  $y \in \{0, 1\}$

Discriminative



Discriminate between  
classes of data points

Generative



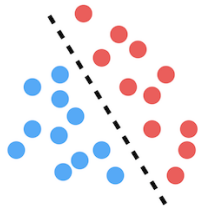
Model the underlying distribu-  
tion of the data

$P(x, y)$   
joint distribution

$$x \xrightarrow{h_0(x)} y$$

## Discriminative Learning Algorithms

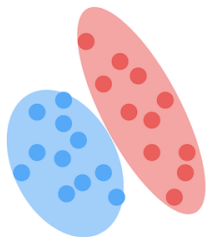
A class of learning algorithms that try to learn the **conditional probability**  $p(y|x)$  directly or learn mappings directly from  $\mathcal{X}$  to  $\mathcal{Y}$ .



- ▶ e.g. linear regression, logistic regression, k-Nearest Neighbors ...

## Generative Learning Algorithms

A class of learning algorithms that model the joint probability  $p(x, y) = p(x|y)p(y)$



- ▶ Equivalently, generative algorithms model  $p(x|y)$  and  $p(y)$
- ▶  $p(y)$  is called the class prior
- ▶ Learned models are transformed to  $p(y|x)$  later to classify data using Bayes' rule

## Bayes Rule

The posterior distribution on  $y$  given  $x$ :

$$p(y|x) = \frac{p(x|y)p(y)}{p(x)} \leftarrow p(x|y=0)p(0) + p(x|y=1)p(1)$$

$y \in \{0, 1\}$

## Bayes Rule

The posterior distribution on  $y$  given  $x$ : <sup>evidence</sup>

$$p(y|x) = \frac{p(x|y)p(y)}{p(x)}$$

Make predictions in a generative model:

$$\begin{aligned} \underset{y}{\operatorname{argmax}} p(y|x) &= \underset{\textcircled{y}}{\operatorname{argmax}} \left( \frac{p(x|y)p(y)}{p(x)} \right) \text{ doesn't depend on } y \\ &= \underset{\textcircled{y}}{\operatorname{argmax}} p(x|y)p(y) \end{aligned}$$

No need to calculate  $p(x)$ .



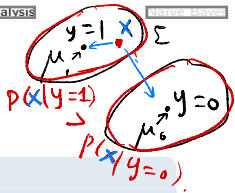
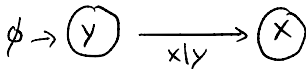
# Generative Models

Generative classification algorithms:

- ▶ Continuous input: Gaussian Discriminant Analysis (GDA)
- ▶ Discrete input: Naïve Bayes

## Gaussian Discriminant Analysis

# Gaussian Discriminant Analysis: Overview



## Goal

$$\begin{cases} N(\mu_0, \Sigma) & \text{if } y=0 \\ N(\mu_1, \Sigma) & \text{if } y=1 \end{cases}$$

Binary classification with input in  $\mathcal{X} = \mathbb{R}^n$  and label in  $\mathcal{Y} = \{0, 1\}$

## Main steps

1. Select a *data generating distribution*.

$$y \sim \text{Bernoulli}(\phi) \quad p(y)$$

$$x|y=0 \sim N(\mu_0, \Sigma), x|y=1 \sim N(\mu_1, \Sigma)$$

2. Estimate model parameters  $\phi$ ,  $\mu_0$ ,  $\mu_1$  and  $\Sigma$  from training data.  $p(x, y)$
3. For any new sample  $x'$ , predict its label by computing  $p(y|x = x'; \phi, \mu_0, \mu_1, \Sigma)$

# Multivariate Normal Distribution

Multivariate normal (or multivariate Gaussian) distribution  $N(\mu, \Sigma)$

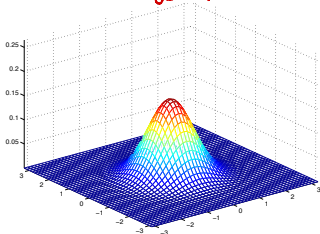
- ▶  $\mu \in \mathbb{R}^n$  is the mean vector,
- ▶  $\Sigma \in \mathbb{R}^{n \times n}$  is the covariance matrix.  $\Sigma$  is symmetric and SPD.

$\overline{n \times 1}$   $\overline{n \times n}$

Density function:

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

↑  
determinant of  $\Sigma$



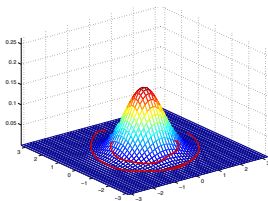
# Multivariate Normal Distribution

Let  $X \in \mathbb{R}^n$  be a random vector. If  $X \sim N(\mu, \Sigma)$ ,

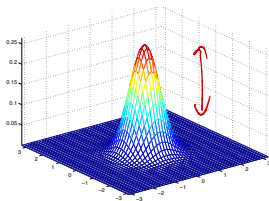
$$\mathbb{E}[X] = \int_{\mathbf{x}} p(\mathbf{x}; \mu, \Sigma) d\mathbf{x} = \mu$$

$$\text{Cov}(X) = \mathbb{E} [(X - \mathbb{E}[X])(X - \mathbb{E}[X])^T] = \Sigma$$

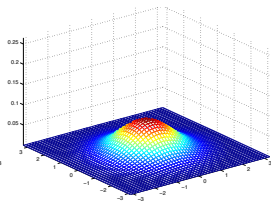
# Gaussian Discriminative Analysis



$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



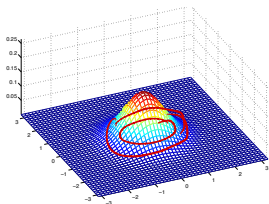
$$\Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 0.6 \end{bmatrix}$$



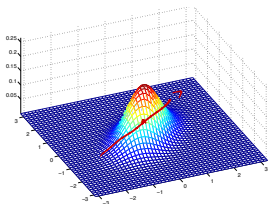
$$\Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

Diagonal entries of  $\Sigma$  controls the “spread” of the distribution

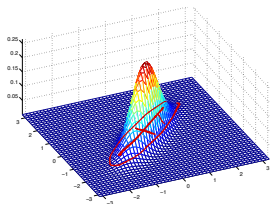
# Gaussian Discriminative Analysis



$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



$$\Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$



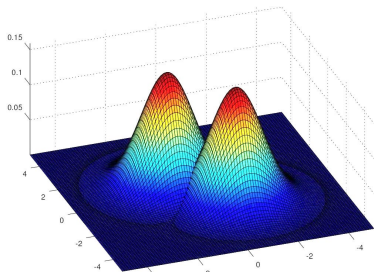
$$\Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$

The distribution is no longer oriented along the axes when off-diagonal entries of  $\Sigma$  are non-zero.

# Gaussian Discriminant Analysis (GDA) Model

Given parameters  $\phi, \mu_0, \mu_1, \Sigma$ ,

$$\begin{aligned} y &\sim \text{Bernoulli}(\phi) \\ x|y=0 &\sim \mathcal{N}(\underline{\mu}_0, \Sigma) \\ x|y=1 &\sim \mathcal{N}(\underline{\mu}_1, \Sigma) \end{aligned}$$



$$p(x|y=y_i) \sim \mathcal{N}(\underline{\mu}_{y_i}, \Sigma)$$

Probability density functions:

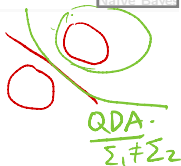
$$\begin{aligned} p(y) &= \underline{\phi^y (1 - \phi)^{1-y}} \\ p(x|y=0) &= \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\underline{\mu}_0)^T \Sigma^{-1} (x-\underline{\mu}_0)} \\ p(x|y=1) &= \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\underline{\mu}_1)^T \Sigma^{-1} (x-\underline{\mu}_1)} \end{aligned}$$



Log likelihood of the data:

$$p(y|x)$$

$$\text{GDA: } \Sigma = \Sigma_1 = \Sigma_2$$



$$l(\phi, \mu_0, \mu_1, \Sigma) = \log \prod_{i=1}^m p(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma)$$

$$= \log \prod_{i=1}^m p(x^{(i)}|y^{(i)}; \mu_0, \mu_1, \Sigma) p(y^{(i)}; \phi)$$

class prior

$$= \sum_{i=1}^m \log p(x^{(i)}|y^{(i)}; \mu_0, \mu_1, \Sigma) + \sum_{i=1}^m \log p(y^{(i)}; \phi)$$

①

$$p(x|y; \mu_0, \mu_1, \Sigma) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(x-\mu_y)^T \Sigma^{-1} (x-\mu_y)}$$

$$\text{② } p(y; \phi) = \phi^y (1-\phi)^{1-y}$$

$$L = \sum_{i=1}^m \log \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} + (-\frac{1}{2}(x^{(i)}-\mu_{y^{(i)}})^T \Sigma^{-1} (x^{(i)}-\mu_{y^{(i)}})) + \sum_{i=1}^m y^{(i)} \log \phi + (1-y^{(i)}) \log(1-\phi)$$

Next step:

$$\frac{\partial L}{\partial \phi} = 0, \text{ solve for } \phi$$

$$\frac{\partial L}{\partial \mu_0} = 0 \quad \frac{\partial L}{\partial \mu_1} = 0 \quad \frac{\partial L}{\partial \Sigma} = 0, \text{ solve for } \mu_0, \mu_1, \Sigma$$

$$L = \sum_{i=1}^m \log \frac{1}{(2\pi)^{\frac{1}{2}} |\Sigma|^{\frac{1}{2}}} + \underbrace{\left(-\frac{1}{2} (x^i - \mu y^i)^T \Sigma^{-1} (x^i - \mu y^i)\right)}_{\text{}} \\ + \sum_{i=1}^m \underbrace{y^i \log \phi + (1-y^i) \log(1-\phi)}_{\text{}}$$

$$\frac{\partial L}{\partial \phi} = \sum_{i=1}^m y^{(i)} \frac{\partial}{\partial \phi} \log \phi + (1-y^{(i)}) \frac{\partial}{\partial \phi} \log(1-\phi) \\ = \left( \sum_{i=1}^m \frac{y^{(i)}}{\phi} \right) + \frac{(1-y^{(i)})}{1-\phi} (-1) \\ = \frac{1}{\phi} \sum_{i=1}^m y^{(i)} - \frac{1}{1-\phi} \sum_{i=1}^m (1-y^{(i)})$$

$$c := \left[ \begin{array}{l} \# \text{ of } 1\text{'s in} \\ \text{training data} \end{array} \right] = \sum_{i=1}^m 1\{y^{(i)} = 1\}$$

$$= \frac{c}{\phi} - \frac{m-c}{1-\phi} = 0$$

$$c(1-\phi) = \phi(m-c)$$

$$c - c\phi - \phi(m-c) = 0$$

$$c - \phi(c+m-c) = 0$$

$$\phi = \frac{c}{m} = \frac{1}{m} \sum_{i=1}^m (y^{(i)} = 1)$$

Facts

$$\textcircled{1} \nabla_x (x^T A x)$$

$$= Ax + A^T x$$

A is symmetric

$$= 2Ax$$

$$\frac{\partial L}{\partial \mu_0} = \sum_{i=1}^m \nabla_{\mu_0} \left[ -\frac{1}{2} (x^{(i)} - \mu_0)^T \Sigma^{-1} (x^{(i)} - \mu_0) \right]$$

$$= \sum_{\substack{i=1 \\ y^{(i)}=0}}^m -\frac{1}{2} \cdot 2 \Sigma^{-1} (x^{(i)} - \mu_0) (-1)$$

$$= \sum_{\substack{i=1 \\ y^{(i)}=0}}^m \Sigma^{-1} (x^i - \mu_0) = 0 \Rightarrow \sum_{\substack{i=1 \\ y^{(i)}=0}}^m (x^i - \mu_0) = 0$$

$$\sum_{i=1}^m 1\{y^i=0\} x^i - \sum_{i=1}^m 1\{y^i=0\} \mu_0 = 0$$

$$\mu_0 = \frac{\sum_{i=1}^m 1\{y^i=0\} x^i}{\sum_{i=1}^m 1\{y^i=0\}} \quad \left. \begin{array}{l} \text{mean of} \\ x \text{ in} \\ \text{class } 0. \end{array} \right\}$$

$$\frac{\partial L}{\partial \mu_1} = 0 \Rightarrow \mu_1 = \frac{\sum_{i=1}^m 1\{y^i=1\} x^i}{\sum_{i=1}^m 1\{y^i=1\}}$$

$$L = \sum_{i=1}^m \log \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} + (-\frac{1}{2} (x^{(i)} - \mu_{y^{(i)}})^T \Sigma^{-1} (x^{(i)} - \mu_{y^{(i)}}))$$

$$+ \sum_{i=1}^m y_i \log \phi + (1 - y_i) \log(1 - \phi)$$

$$= \sum_{i=1}^m \log \frac{1}{(2\pi)^{\frac{n}{2}}} + \sum_{i=1}^m \log |\Sigma|^{-\frac{1}{2}} - \frac{1}{2} (x^{(i)} - \mu_{y^{(i)}})^T \Sigma^{-1} (x^{(i)} - \mu_{y^{(i)}}) + \sum_{i=1}^m y_i \log \phi + (1 - y_i) \log(1 - \phi)$$

$$\textcircled{2} \nabla_A |A| = |A| (A^{-1})^T$$

$$\textcircled{3} \nabla_A x^T A^{-1} y$$

$$= -A^{-T} x y^T A^{-T}$$

$$\nabla_{\Sigma} L(\phi, \mu, \Sigma) = \nabla_{\Sigma} \sum_{i=1}^m \log |\Sigma|^{-\frac{1}{2}} - \frac{1}{2} (x^{(i)} - \mu_{y_i})^T \Sigma^{-1} (x^{(i)} - \mu_{y_i})$$

$$= \nabla_{\Sigma} \sum_{i=1}^m -\frac{1}{2} \log |\Sigma| - \frac{1}{2} (x^{(i)} - \mu_{y_i})^T \Sigma^{-1} (x^{(i)} - \mu_{y_i})$$

$$= \sum_{i=1}^m -\frac{1}{2} \frac{1}{|\Sigma|} \nabla_{\Sigma} |\Sigma| - \frac{1}{2} (-\Sigma^{-T}) (x^{(i)} - \mu_{y_i}) (x^{(i)} - \mu_{y_i})^T \Sigma^{-T}$$

$|\Sigma| (\Sigma^{-1})^T$  - symmetric

$$= \sum_{i=1}^m \left( -\frac{1}{2} (\Sigma^{-1}) \right) + \left( \frac{1}{2} \Sigma^{-1} (x^{(i)} - \mu_{y_i}) (x^{(i)} - \mu_{y_i})^T \Sigma^{-1} \right) = 0$$

$$\sum_{i=1}^m (1 - \Sigma^{-1} (x^{(i)} - \mu_{y_i}) (x^{(i)} - \mu_{y_i})^T) = 0$$

$$m - \Sigma^{-1} \sum_{i=1}^m (x^{(i)} - \mu_{y_i}) (x^{(i)} - \mu_{y_i})^T = 0$$

$$\underline{\Sigma} = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu_{y_i}) (x^{(i)} - \mu_{y_i})^T$$

Log likelihood of the data:

$$\begin{aligned}l(\phi, \mu_0, \mu_1, \Sigma) &= \log \prod_{i=1}^m p(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma) \\ &= \log \prod_{i=1}^m p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \Sigma) p(y^{(i)}; \phi)\end{aligned}$$

Log likelihood of the data:

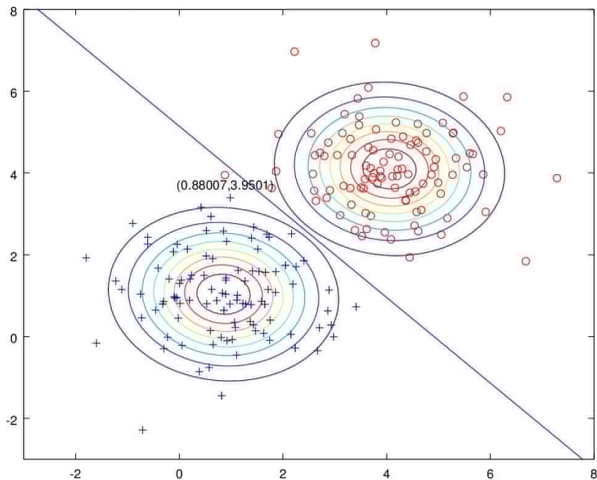
$$\begin{aligned}
 l(\phi, \mu_0, \mu_1, \Sigma) &= \log \prod_{i=1}^m p(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma) \\
 &= \log \prod_{i=1}^m p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \Sigma) p(y^{(i)}; \phi)
 \end{aligned}$$

Maximum likelihood estimate of the parameters:

$$\left. \begin{aligned}
 \phi &= \frac{1}{m} \sum_{i=1}^m \mathbf{1}\{y^{(i)} = 1\} \\
 \mu_b &= \frac{\sum_{i=1}^m \mathbf{1}\{y^{(i)} = b\} x^{(i)}}{\sum_{i=1}^m \mathbf{1}\{y^{(i)} = b\}} \text{ for } b = 0, 1 \\
 \Sigma &= \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T
 \end{aligned} \right\}$$

# Maximum likelihood estimation of GDA

GDA finds a linear decision boundary at which  
 $p(y = 1|x) = p(y = 0|x) = 0.5$



# GDA and Logistic Regression

## Proposition

$p(y = 1|x; \phi, \mu_0, \mu_1, \Sigma)$  can be written in the form:

$$p(y = 1|x; \phi, \Sigma, \mu_0, \mu_1) = \frac{1}{\underbrace{1 + e^{-\theta^T x}}}$$

# GDA and Logistic Regression

## Proposition

$p(y = 1|x; \phi, \mu_0, \mu_1, \Sigma)$  can be written in the form:

$$p(y = 1|x; \phi, \Sigma, \mu_0, \mu_1) = \frac{1}{1 + e^{-\theta^T x}}$$

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} \left[ \frac{1}{2} (\mu_0^T \Sigma^{-1} \mu_0 - \mu_1^T \Sigma^{-1} \mu_1) - \log \frac{1-\phi}{\phi} \right], \quad x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \\ 1 \end{bmatrix}$$

$$p(y=1|x; H) = \frac{p(x|y=1; H) P(y=1; H)}{P(x; H)} = \frac{p(x|y=1; H) P(y=1; H)}{p(x|y=1; H) P(y=1; H) + p(x|y=0; H) P(y=0; H)}$$

$$= \frac{1}{1 + \exp\left( - \underbrace{(\mu_1 - \mu_0)^T \Sigma^{-1} x}_{\theta_0} - \underbrace{\left( \frac{1}{2} (\mu_0^T \Sigma^{-1} \mu_0 - \mu_1^T \Sigma^{-1} \mu_1) - \log \frac{1-\phi}{\phi} \right)}_{\theta_1 \in \mathbb{R}} \right)}$$



# GDA and Logistic Regression

## Proposition

$p(y = 1|x; \phi, \mu_0, \mu_1, \Sigma)$  can be written in the form:

$$p(y = 1|x; \phi, \Sigma, \mu_0, \mu_1) = \frac{1}{1 + e^{-\theta^T x}}$$

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} \Sigma^{-1}(\mu_1 - \mu_0) \\ \frac{1}{2}(\mu_0^T \Sigma^{-1} \mu_0 - \mu_1^T \Sigma^{-1} \mu_1) - \log \frac{1-\phi}{\phi} \end{bmatrix}, x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \\ 1 \end{bmatrix}$$

Similarly,

$$p(y = 0|x; \phi, \Sigma, \mu_0, \mu_1) = \frac{1}{1 + e^{\theta^T x}}$$

If  $\underline{p(x|y)} \sim \underline{\mathcal{N}(\mu, \Sigma)}$ ,  $\underline{p(y|x)}$  is a logistic function.

# GDA and Logistic Regression

$$\Sigma_1 = \Sigma_2$$

## GDA

- ▶ Maximizes the joint likelihood  $\prod_{i=1}^m p(x^{(i)}, y^{(i)})$
- ▶ Modeling assumptions:  $x|y=b \sim \mathcal{N}(\mu_b, \Sigma)$ ,  $y \sim \text{Bernoulli}(\phi)$   $P(y)$
- ▶ When modeling assumptions are correct, GDA is asymptotically efficient and data efficient

## Logistic Regression

- ▶ Maximizes the conditional likelihood  $\prod_{i=1}^m p(y^{(i)}|x^{(i)})$
- ▶ Modeling assumptions:  $p(y|x)$  is a logistic function; no restriction on  $p(x)$
- ▶ More robust and less sensitive to incorrect modeling assumptions.

## Naïve Bayes

# Naïve Bayes: Motivating Example

A simple generative learning algorithm for discrete input variables

## Example: Spam filter (document classification)

Classify email messages  $x$  to spam ( $y = 1$ ) and non-spam ( $y = 0$ ) classes.

Hello [REDACTED]

We need to confirm your info...

(1) FINAL MESSAGE: Payout Verification - \$3000 PAYOUT is ready to be addressed in your Name and we want to be sure it gets to the right place. Click below to start the confirmation process. The sooner you act, the sooner it can be in your hands!

[Raging Bull Casino](#)

A sample spam email

## Example: Spam Filter

### Binary text features

Given a dictionary of size  $n$ , represent a message composed of dictionary words as  $x \in \{0, 1\}^n$ :

$$x_i = \begin{cases} 1 & i\text{-th dictionary word is in message} \\ 0 & \text{otherwise} \end{cases}$$

$$x = \begin{matrix} & \text{A} & \\ \begin{matrix} 0 \\ 0 \\ \vdots \\ \textcircled{1} \\ \vdots \\ \textcircled{1} \\ \vdots \\ 0 \end{matrix} & \begin{matrix} \underline{a} \\ \underline{aardvark} \\ \vdots \\ \underline{casino} \\ \vdots \\ \underline{payout} \\ \vdots \\ \underline{zyzzyva} \\ \text{z.} \end{matrix} & \left. \vphantom{\begin{matrix} 0 \\ 0 \\ \vdots \\ \textcircled{1} \\ \vdots \\ \textcircled{1} \\ \vdots \\ 0 \end{matrix}} \right\} n \end{matrix}$$

# Naïve Bayes Model

$$\underline{P(y)} \quad \underline{P(x|y)}$$

Probability of observing email  $x_1, \dots, x_n$  given spam class  $y$  :

$$\underline{P(x_1, \dots, x_n|y)} = \underline{P(x_1|y)P(x_2|y, x_1), \dots, P(x_n|y, x_1, \dots, x_{n-1})}$$

## Naïve Bayes (NB) assumption

$x_i$ 's are conditionally independent given  $y$ :

$$\underline{P(x_i|y, x_1, \dots, x_{i-1})} = \underline{P(x_i|y)}$$

$$\begin{aligned} \uparrow \\ \text{message} \quad P(x|y) &= \underline{P(x_1, \dots, x_n|y)} = \underline{P(x_1|y)P(x_2|y) \dots P(x_n|y)} = \underbrace{\prod_{i=1}^n P(x_i|y)}_{\substack{\text{size of} \\ \text{dictionary}}} \end{aligned}$$

# Naïve Bayes Parameters

## Multi-variate Bernoulli event model

$x|y$  generated from  $n$  independent Bernoulli trials

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

- ▶  $y \sim \text{Bernoulli}(\phi_y)$  : assume email class (spam vs no-spam) is randomly generated with prior  $p(y) = \phi_y^y (1 - \phi_y)^{1-y}$
- ▶  $x_i|y=b \sim \text{Bernoulli}(\phi_{i|y=b})$ ,  $b = 1, 2$  : given  $y = b$ , each word  $x_i$  is included in the message independently with  $p(x_i = 1|y=b) = \phi_{i|y=b}$ . i.e.  $(1, 0)$

$$p(x_i=1|y=1) = \phi_{i|y=1}$$

$$p(x_i=0|y=0) = \phi_{i|y=0}$$

$$p(x_i|y=b) = \phi_{i|y=b}^{x_i} (1 - \phi_{i|y=b})^{1-x_i}$$

Model parameters:

- ▶  $\phi_y \in [0, 1]$
- ▶  $\phi_{i|y=1}, \phi_{i|y=0}$  for  $i = 1, \dots, n$   $2n$ .

## Naïve Bayes Parameter Learning

Likelihood of training data  $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$ : i.i.d

$$L(\theta) = \prod_{i=1}^m p(x^{(i)}, y^{(i)}) = \prod_{i=1}^m p(x^{(i)} | y^{(i)}) p(y^{(i)})$$

$\log L(\dots)$

Maximum likelihood estimation of parameters:

$$\phi_y = \frac{1}{m} \sum_{i=1}^m \mathbf{1}\{y^{(i)} = 1\} \quad \text{\% of spam emails}$$

$$y^i = 1$$

$$\phi_{j|y=b} = \frac{\sum_{i=1}^m \mathbf{1}\{x_j^{(i)} = 1, y^{(i)} = b\}}{\sum_{i=1}^m \mathbf{1}\{y^{(i)} = b\}} \quad \text{for } b = 1, 0$$

$$\frac{\# \text{ of messages of both spam and } x_j = 1}{\# \text{ spam messages}}$$

*\%* of spam(non-spam) emails containing  $j$ th dictionary word

if  $y^i = 0$ .

$$\phi_{j|y=0} = \frac{\# \text{ of non spam messages containing } x_j}{\# \text{ of non spam messages}}$$

$j$ th word appear in the  $i$ th message



# Naïve Bayes Prediction

 $[x_1, \dots, x_n]$ 

Given new example with feature  $\underline{x}$ , compute the posterior probability

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

known  $\phi_y^*$ ,  $\phi_i|y=b$ .

$$= \frac{p(x|y = 1)p(y = 1)}{p(x|y = 1)p(y = 1) + p(x|y = 0)p(y = 0)}$$

By the N.B assumption,

$$= \frac{\prod_{i=1}^n p(x_i|y = 1)p(y = 1)}{\prod_{i=1}^n p(x_i|y = 1)p(y = 1) + \prod_{i=1}^n p(x_i|y = 0)p(y = 0)}$$

Choose label  $y = 1$  (spam) if  $\underline{p(y = 1|x)} > T$  where  $\underline{T} \in [0, 1]$  is a threshold .. e.g.  $T = 0.5$

$\underline{T}$  tradeoff between wrongly blocked non-spam (FPs) vs. wrongly blocked spams (FNs).  
 false negatives. false positives

# Laplace smoothing

Issue with Naïve Bayes prediction:

- Suppose word  $x_j$  hasn't been seen in the training data,

$$\phi_{j|y=1} =$$

For all  $i=1, \dots, m$ ,  $1\{x_j=1\} = 0$

$$\phi_{j|y=b} = \frac{\sum_{i=1}^m 1\{y^i=b, x_j=1\}}{\sum_{i=1}^m 1\{y^i=b\}} \rightarrow 0 = 0$$

$$\left. \begin{array}{l} \phi_{j|y=1} \\ \phi_{j|y=0} \end{array} \right\} = 0$$

Given  $x$ ,

$$P(y=1|x) = \frac{\prod_{l=1}^n \underbrace{P(x_l|y=1)}_0 \underbrace{P(y=1)}_0}{\prod_{l=1}^n \underbrace{P(x_l|y=1)}_0 P(y=1) + \prod_{l=1}^n \underbrace{P(x_l|y=0)}_0 P(y=0)} \quad \text{if } l=j \quad = \frac{0}{0}$$

# Laplace smoothing

Issue with Naïve Bayes prediction:

- ▶ Suppose word  $x_j$  hasn't been seen in the training data,  
 $\phi_{j|y=1} = \phi_{j|y=0} = 0$
- ▶ Can not compute class posterior  $p(y = 1|x) = \frac{0}{0}$ .

# Laplace smoothing

Issue with Naïve Bayes prediction:

- ▶ Suppose word  $x_j$  hasn't been seen in the training data,  
 $\phi_{j|y=1} = \phi_{j|y=0} = 0$
- ▶ Can not compute class posterior  $p(y = 1|x) = \frac{0}{0}$ .

$$\phi_j = \frac{1}{n} \sum_{i=1}^m \mathbb{1}\{z^i = j\}$$

$y \in \{0, 1\}$   $k=2$   
 $\times | y$

## Laplace smoothing

Let  $z \in \{1, \dots, k\}$  be a multinomial random variable. Given  $m$  independent observations  $z^{(1)} \dots z^{(m)}$ , maximum likelihood estimation of  $\phi_j = p(z = j)$  with **Laplace smoothing** is

$$\phi_j = \frac{\left( \sum_{i=1}^m \mathbb{1}\{z^{(i)} = j\} \right) + 1}{m + k}$$

$\uparrow$   $\epsilon$

$$\sum_{j=1}^k \phi_j = 1.$$

- ▶  $\phi_j \neq 0$  for all  $j$
- ▶  $\sum_{j=1}^k \phi_j = 1$

# Naïve Bayes with Laplace smoothing

Apply Laplace smoothing to  $\phi_{j|y=b}$  for  $b \in \{0, 1\}$

$$\phi_{j|y=b} = \frac{\sum_{i=1}^m \mathbf{1}\{x_j^{(i)} = 1, y^{(i)} = b\} + 1}{\sum_{i=1}^m \mathbf{1}\{y^i = b\} + 2}$$

In practice we don't apply Laplace smoothing to  $\phi_y = p(y = 1)$ , which is greater than 0.

# Naïve Bayes Summary

## Naïve Bayes (NB) assumption

$x_i$ 's are conditionally independent given  $y$ :

$$\underline{p(x_1, \dots, x_n | y)} = \prod_{i=1}^n \underline{p(x_i | y)}$$

Different event models

- ▶ **Multi-variate Bernoulli model:** represent document of vocab size  $n$  as  $n$  independent Bernoulli trials
- ▶ **Multinomial event model:** represent document of  $N$  words as  $x = \{x_1, \dots, x_n\}$  where  $x_i \in \{1, \dots, K\}$ . (not covered)

# Homework

- ▶ Programming Assignment 1 late submission.
- ▶ TA sessions