

Foundations of Machine Learning (ECE 5984)

- Generative Learning Algorithms -

Eunbyung Park

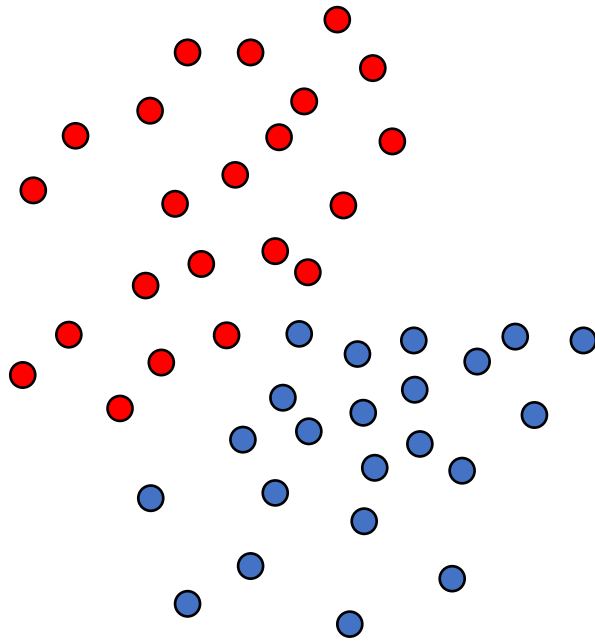
Assistant Professor

School of Electronic and Electrical Engineering

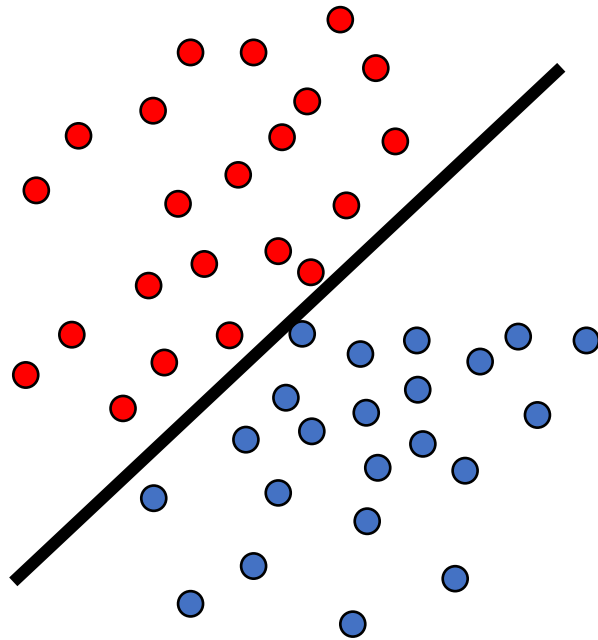
[Eunbyung Park \(silverbottlep.github.io\)](https://silverbottlep.github.io)

Generative vs Discriminative

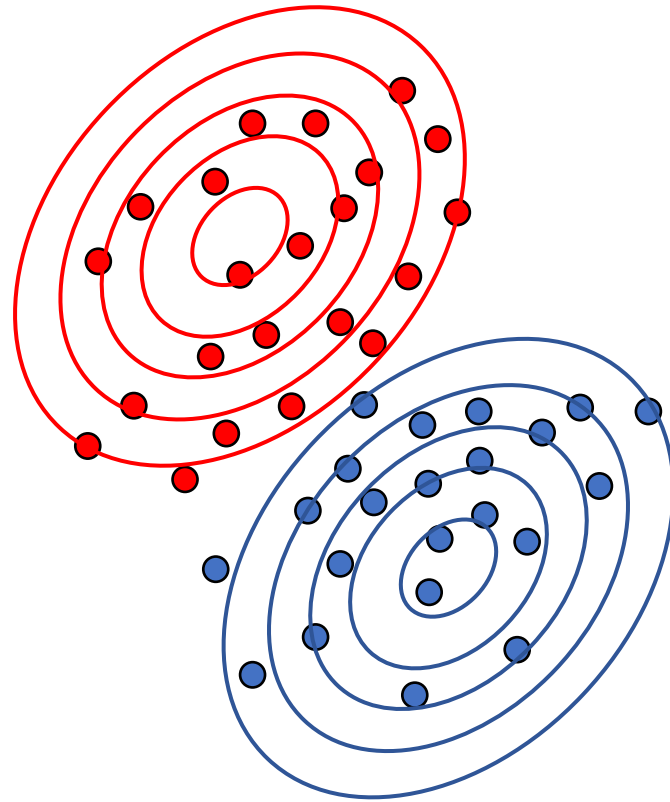
Discriminative Models



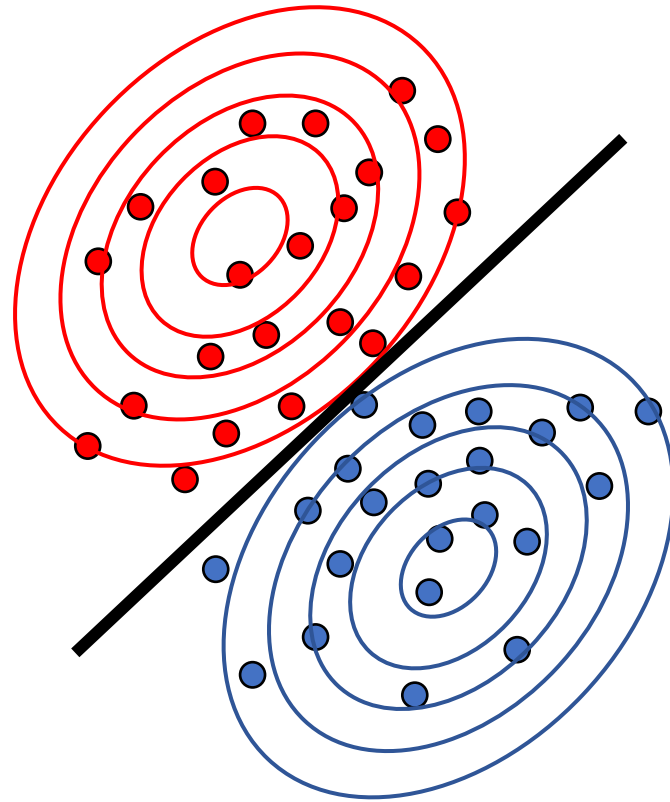
Discriminative Models



Generative Models



Generative Models



Discriminative vs. Generative

Discriminative Approach

$$p(y|x)$$

$$\operatorname{argmax}_y p(y|x)$$

Generative Approach

$$p(x|y), p(y)$$

$$\begin{aligned}\operatorname{argmax}_y p(y|x) &= \operatorname{argmax}_y \frac{p(x|y)p(y)}{p(x)} \\ &= \operatorname{argmax}_y p(x|y)p(y)\end{aligned}$$

Gaussian Discriminant Analysis

Gaussian Discriminant Analysis

- Classification where input feature x are continuous variables

$$D = \{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}, \quad x^{(i)} \in \mathbb{R}^d, \quad y^{(i)} \in \{0, 1\}$$

Gaussian Discriminant Analysis

- Classification where input feature x are continuous variables

$$D = \{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}, \quad x^{(i)} \in \mathbb{R}^d, \quad y^{(i)} \in \{0, 1\}$$

$$\mu_0, \mu_1 \in \mathbb{R}^d \quad \Sigma \in \mathbb{R}^{d \times d}$$

$$p(y) = \text{Bern}(\phi)$$

$$p(x|y = 0) = N(\mu_0, \Sigma)$$

$$p(x|y = 1) = N(\mu_1, \Sigma)$$

Shared Covariance

Gaussian Discriminant Analysis

- Classification where input feature x are continuous variables

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

$$p(x|y = 0) = \frac{1}{(2\pi)^{d/2}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(x - \mu_0)^\top \Sigma^{-1}(x - \mu_0)\right)$$

$$p(x|y = 1) = \frac{1}{(2\pi)^{d/2}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(x - \mu_1)^\top \Sigma^{-1}(x - \mu_1)\right)$$

MLE

- Given data D , we want to maximize likelihood!

$$\operatorname{argmax} \log p(D|\phi, \mu_0, \mu_1, \Sigma) =$$

MLE

- Given data D , we want to maximize likelihood!

$$\begin{aligned}\operatorname{argmax} \log p(D; \phi, \mu_0, \mu_1, \Sigma,) &= \log \prod_{i=1}^N p(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma,) \\ &= \log \prod_{i=1}^N p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \Sigma) p(y^{(i)}; \phi) \\ &= \sum_{i=1}^N \log p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \Sigma) + \sum_{i=1}^N \log p(y^{(i)}; \phi)\end{aligned}$$

MLE

- Given data D , we want to maximize likelihood!

$$\frac{\partial \log p}{\partial \phi} =$$

MLE

- Given data D, we want to maximize likelihood!

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

$$= \frac{1}{\phi} \sum_{i=1}^N y^{(i)} - \frac{1}{1 - \phi} \sum_{i=1}^N (1 - y^{(i)}) = 0$$

$$\frac{1}{\phi} \sum_{i=1}^N y^{(i)} = \frac{1}{1 - \phi} \sum_{i=1}^N (1 - y^{(i)})$$

$$\sum_{i=1}^N y^{(i)} = N\phi$$

$$\phi^* = \frac{1}{N} \sum_{i=1}^N y^{(i)}$$

MLE

- Given data D, we want to maximize likelihood!

$$\frac{\partial \log p}{\partial \mu_1} =$$

MLE

- Given data D, we want to maximize likelihood!

$$\begin{aligned}\frac{\partial \log p}{\partial \mu_1} &= \frac{\partial}{\partial \mu_1} \sum_{i=1}^N \log p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \Sigma) \\ &= \frac{\partial}{\partial \mu_1} \sum_{i \in \{j | y^{(j)} = 1\}} \log p(x^{(i)} | y^{(i)}; \mu_1, \Sigma)\end{aligned}$$

$$\mu_1^* = \frac{\sum_{i=1}^N 1\{y^{(i)} = 1\} x^{(i)}}{\sum_{i=1}^N 1\{y^{(i)} = 1\}}$$

$$\mu_0^* = \frac{\sum_{i=1}^N 1\{y^{(i)} = 0\} x^{(i)}}{\sum_{i=1}^N 1\{y^{(i)} = 0\}}$$

MLE

- Given data D , we want to maximize likelihood!

$$\frac{\partial \log p}{\partial \Sigma} =$$

MLE

- Given data D, we want to maximize likelihood!

$$\frac{\partial \log p}{\partial \Sigma} = \frac{\partial}{\partial \Sigma} \sum_{i=1}^N \log p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \Sigma)$$

$$\Sigma^* = \frac{1}{N} \sum_{i=1}^N \left(x^{(i)} - \mu_{y^{(i)}} \right) \left(x^{(i)} - \mu_{y^{(i)}} \right)^{\top}$$

Testing

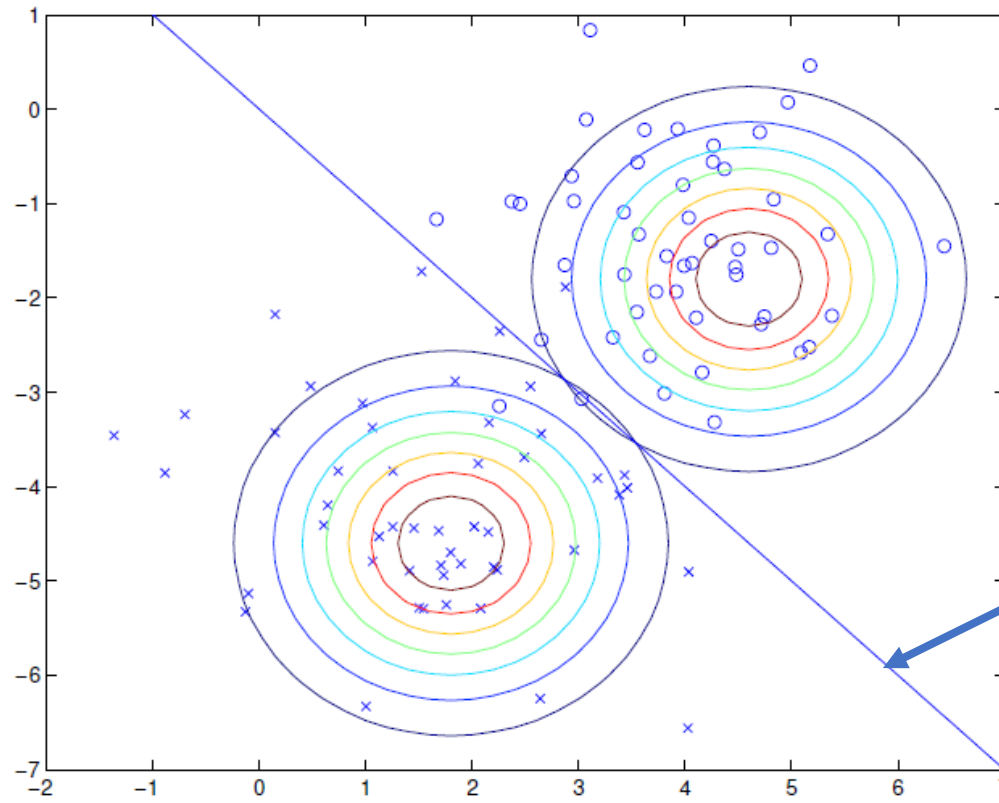
- Given a new data x^{new}

$$\operatorname{argmax}_y p(y|x^{new}) = \operatorname{argmax}_y \frac{p(x^{new}|y)p(y)}{p(x)} = \operatorname{argmax}_y p(x^{new}|y)p(y)$$

Compute $p(x^{new}|y = 0), p(y = 0), p(x^{new}|y = 1), p(y = 1)$

Linear Decision Boundary

- Linear Discriminant Analysis



$$p(y = 1|x) = 0.5 \\ = p(y = 0|x)$$

Linear Decision Boundary

- Linear Discriminant Analysis

$$\log p(y = 1|x) = \log p(y = 0|x)$$

Linear Decision Boundary

- Linear Discriminant Analysis

$$\log p(y = 1|x) = \log p(y = 0|x)$$

$$(x - \mu_0)^\top \Sigma^{-1} (x - \mu_0) = (x - \mu_1)^\top \Sigma^{-1} (x - \mu_1) + \text{const}$$

$$x^\top \Sigma^{-1} x - 2\mu_0^\top \Sigma^{-1} x + \mu_0^\top \mu_0 = x^\top \Sigma^{-1} x - 2\mu_1^\top \Sigma^{-1} x + \mu_1^\top \mu_1 + \text{const}$$

$$-2(\mu_0^\top \Sigma^{-1} + \mu_1^\top \Sigma^{-1})x + \mu_0^\top \mu_0 - \mu_1^\top \mu_1 + \text{const} = 0$$

Quadratic Decision Boundary

- Quadratic Discriminant Analysis, no shared covariance!

$$\log p(y = 1|x) = \log p(y = 0|x)$$

Quadratic Decision Boundary

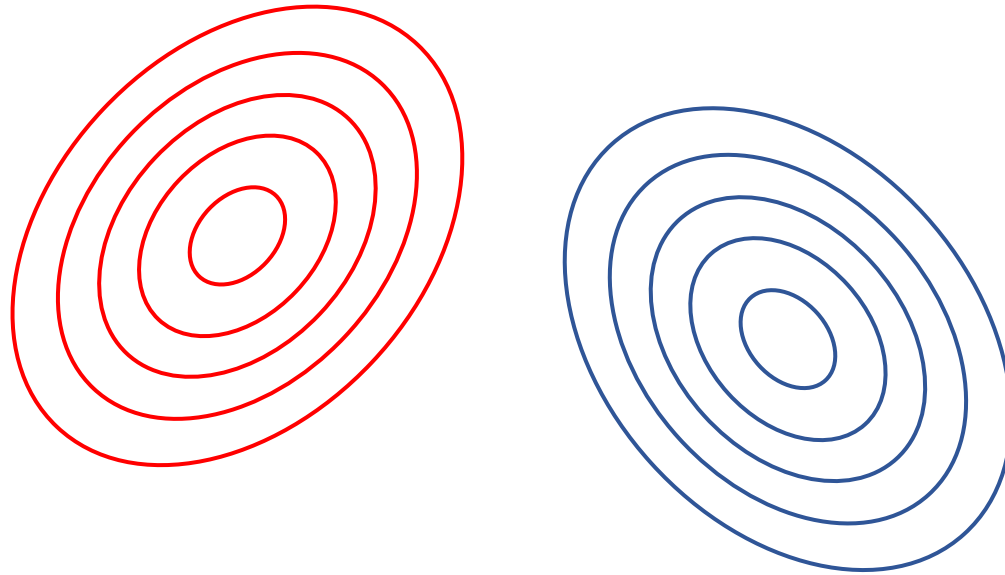
- Quadratic Discriminant Analysis, no shared covariance!

$$\log p(y = 1|x) = \log p(y = 0|x)$$

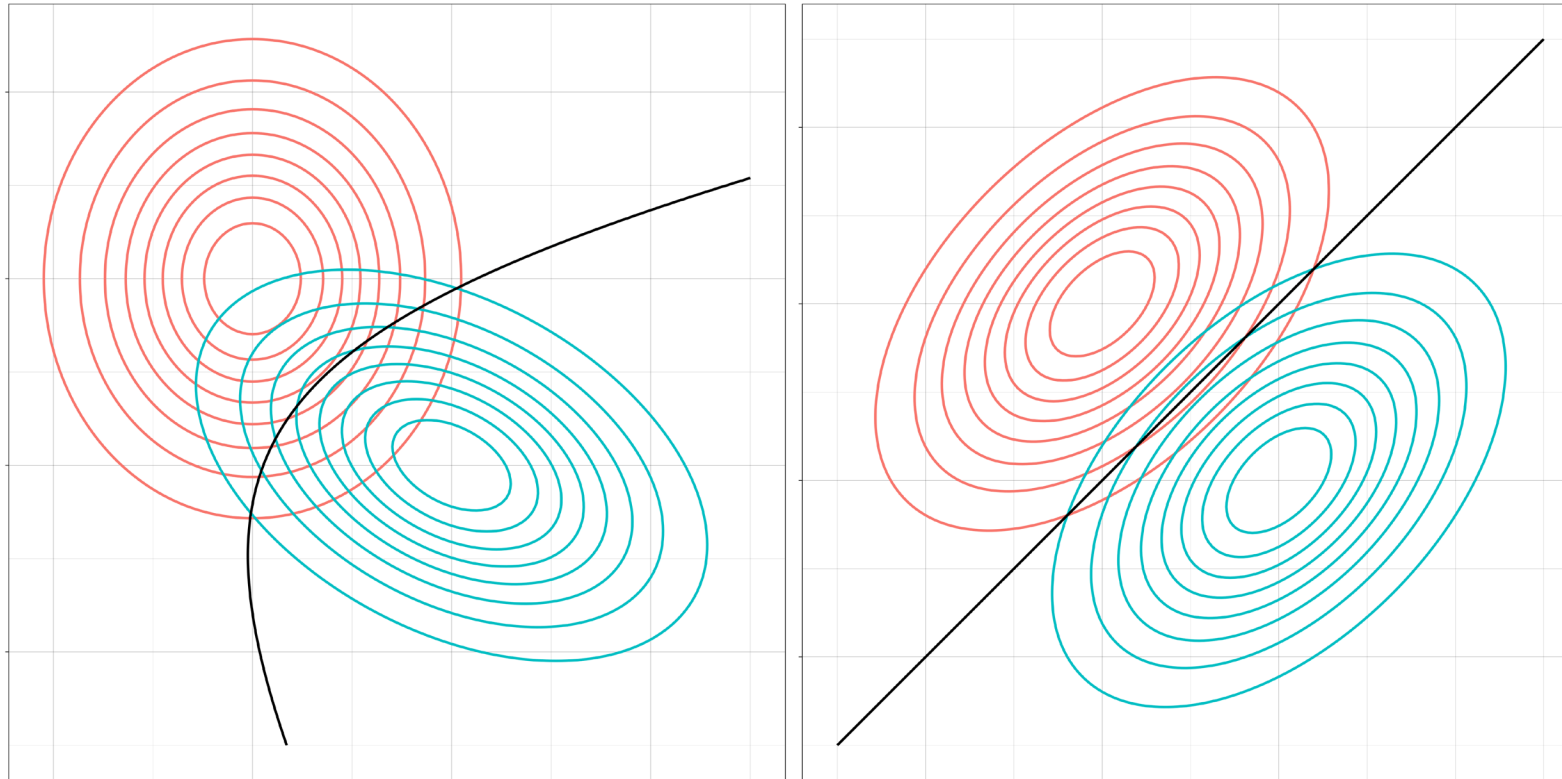
$$\log \frac{1}{|\Sigma_0|^{\frac{1}{2}}} + (x - \mu_0)^\top \Sigma_0^{-1} (x - \mu_0) = \log \frac{1}{|\Sigma_1|^{\frac{1}{2}}} + (x - \mu_1)^\top \Sigma_1^{-1} (x - \mu_1) + \text{const}$$

Quadratic Decision Boundary

- Quadratic Discriminant Analysis, no shared covariance!

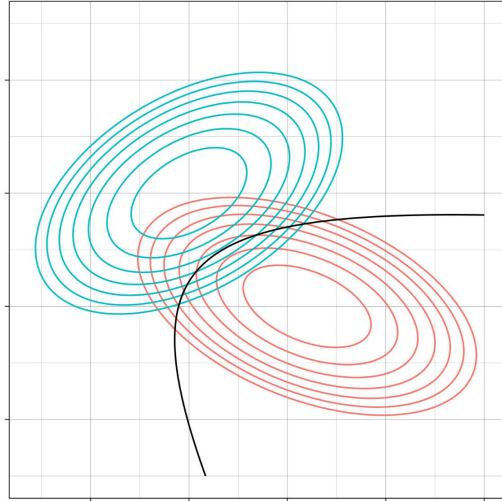


LDA vs QDA

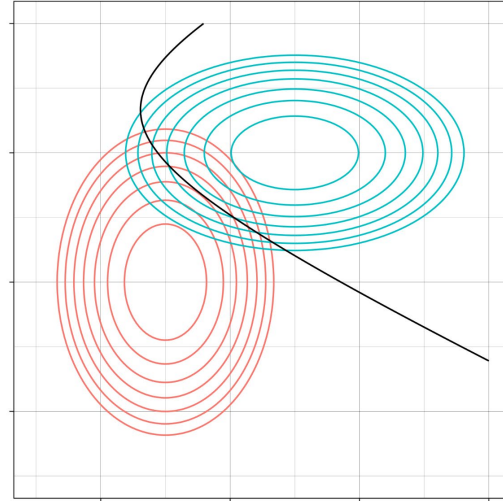


LDA vs QDA

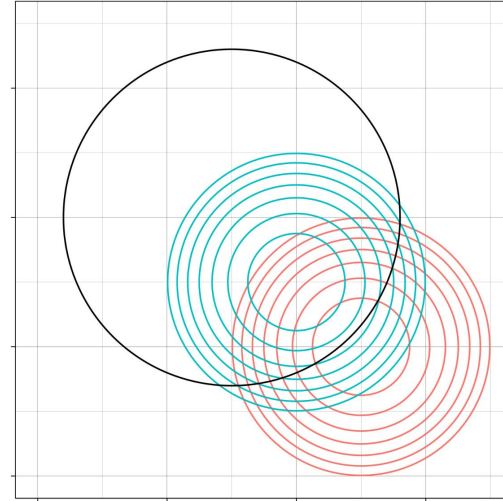
QDA



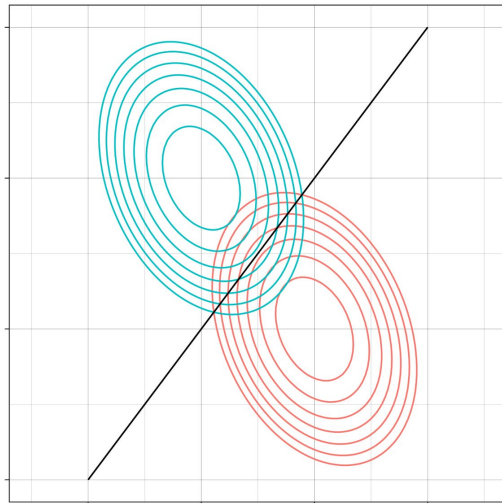
Diagonal QDA



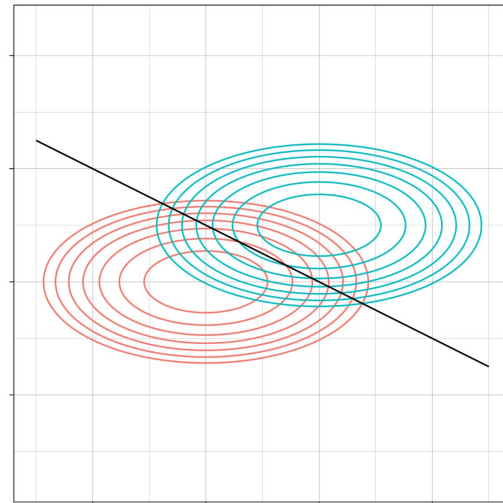
Spherical QDA



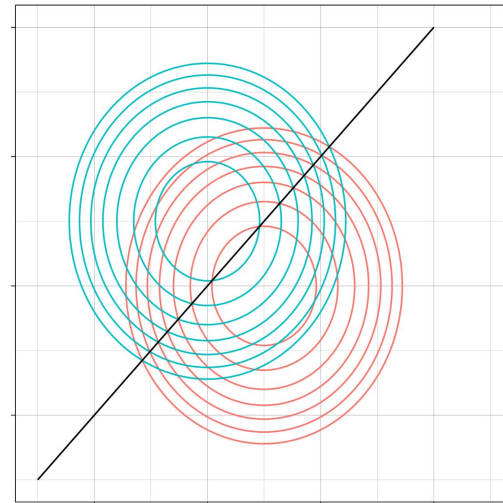
LDA



Diagonal LDA



Spherical LDA



GDA vs Logistic Regression

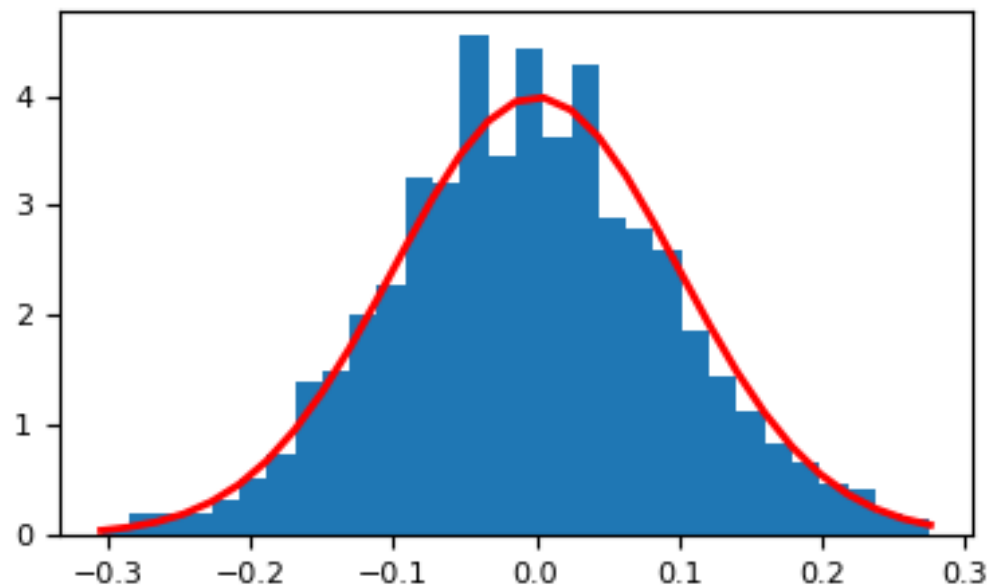
- GDA makes stronger assumption: $p(x|y)$ is a gaussian
- If the assumption is true, then GDA is “asymptotically efficient”
 - The best possible model when $N \rightarrow \infty$
- When data is not Gaussian, logistic regression beats GDA when N is large
- GDA is usually better than logistic regression when N is small
- GDA is a generative model, so we can sample!

How To Generate Samples?

How to Generate Samples?

`numpy.random.normal`

`random.normal(loc=0.0, scale=1.0, size=None)`



Random Number Generators

- Sampling from a uniform distribution over $[0,1)$

`numpy.random.rand`

`random.rand($d0, d1, \dots, dn$)`

Pseudo Random Number Generators

- Linear Congruential Generators

$$X_{n+1} = (aX_n + c) \bmod d$$

$$\begin{aligned} \text{seed} &= X_0 = 1 \\ a &= 5, c = 3, d = 9 \end{aligned}$$

$$0 < d$$

$$0 < a < d$$

$$0 \leq c < d$$

$$X_0 = 1$$

$$X_1 = (5 \cdot 1 + 3) \bmod 9 = 8$$

$$X_2 = (5 \cdot 8 + 3) \bmod 9 = 7$$

$$X_3 = (5 \cdot 7 + 3) \bmod 9 = 2$$

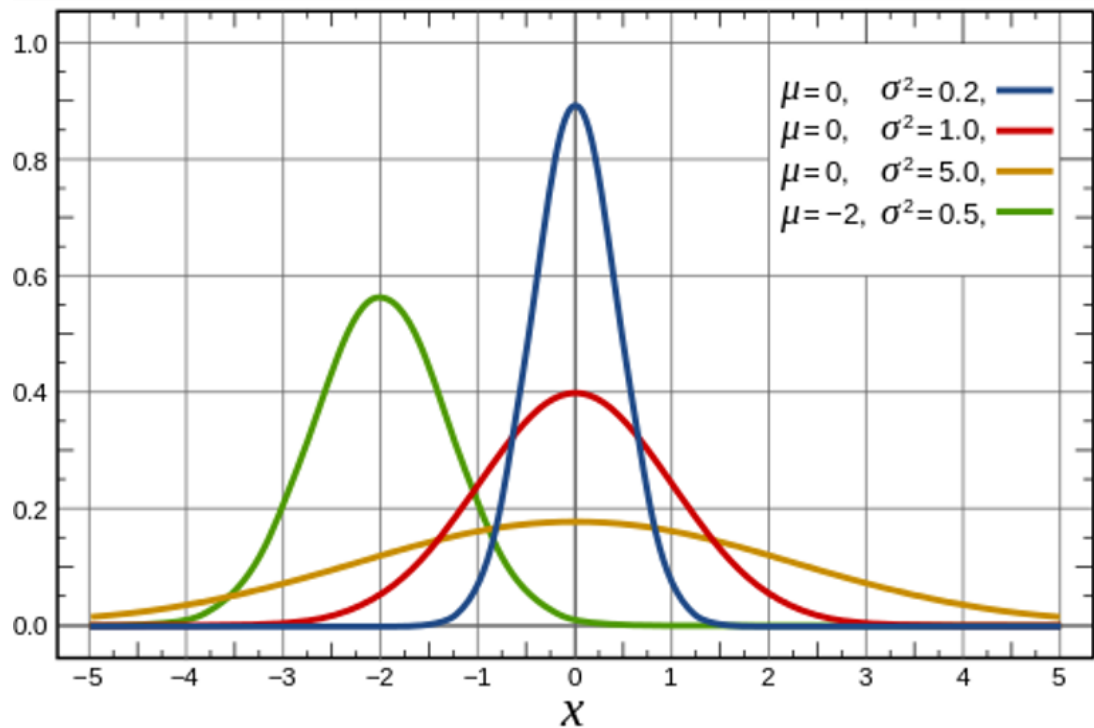
$$X_4 = (5 \cdot 2 + 3) \bmod 9 = 4$$

...

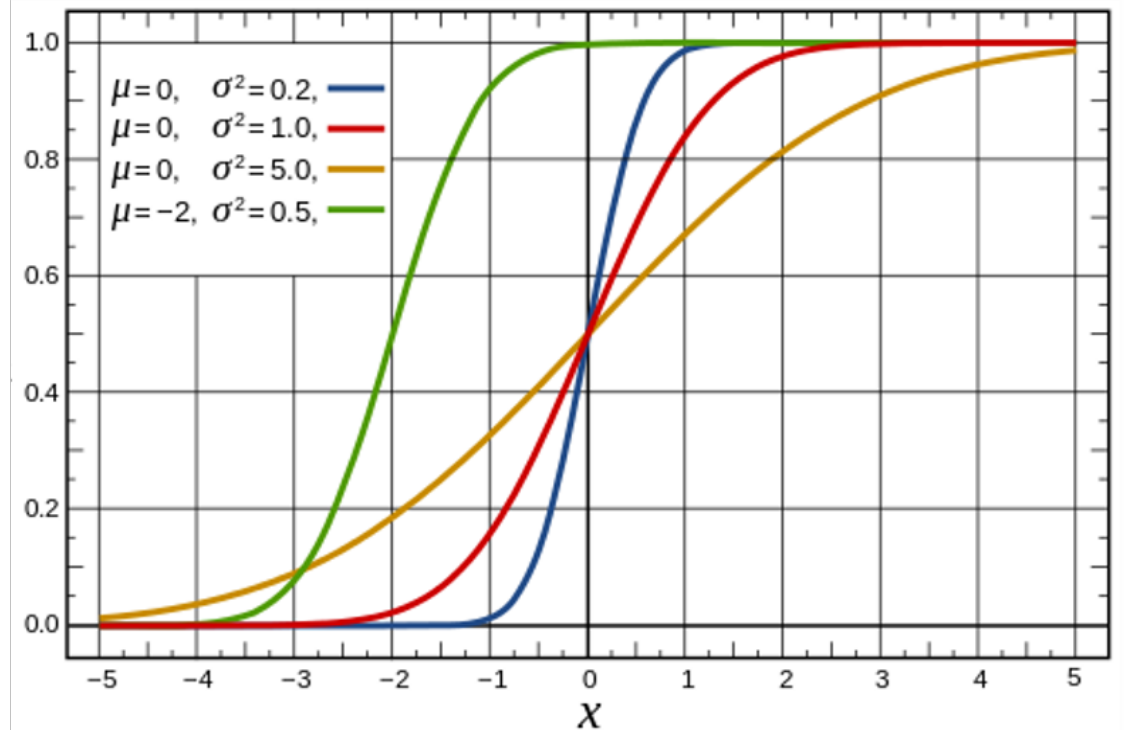
Sampling From Gaussian Distribution

- Inverse transform sampling

Probability Density Function



Cumulative Density Function



Sampling From Gaussian Distribution

- Standard Normal Distribution $N(0, 1)$

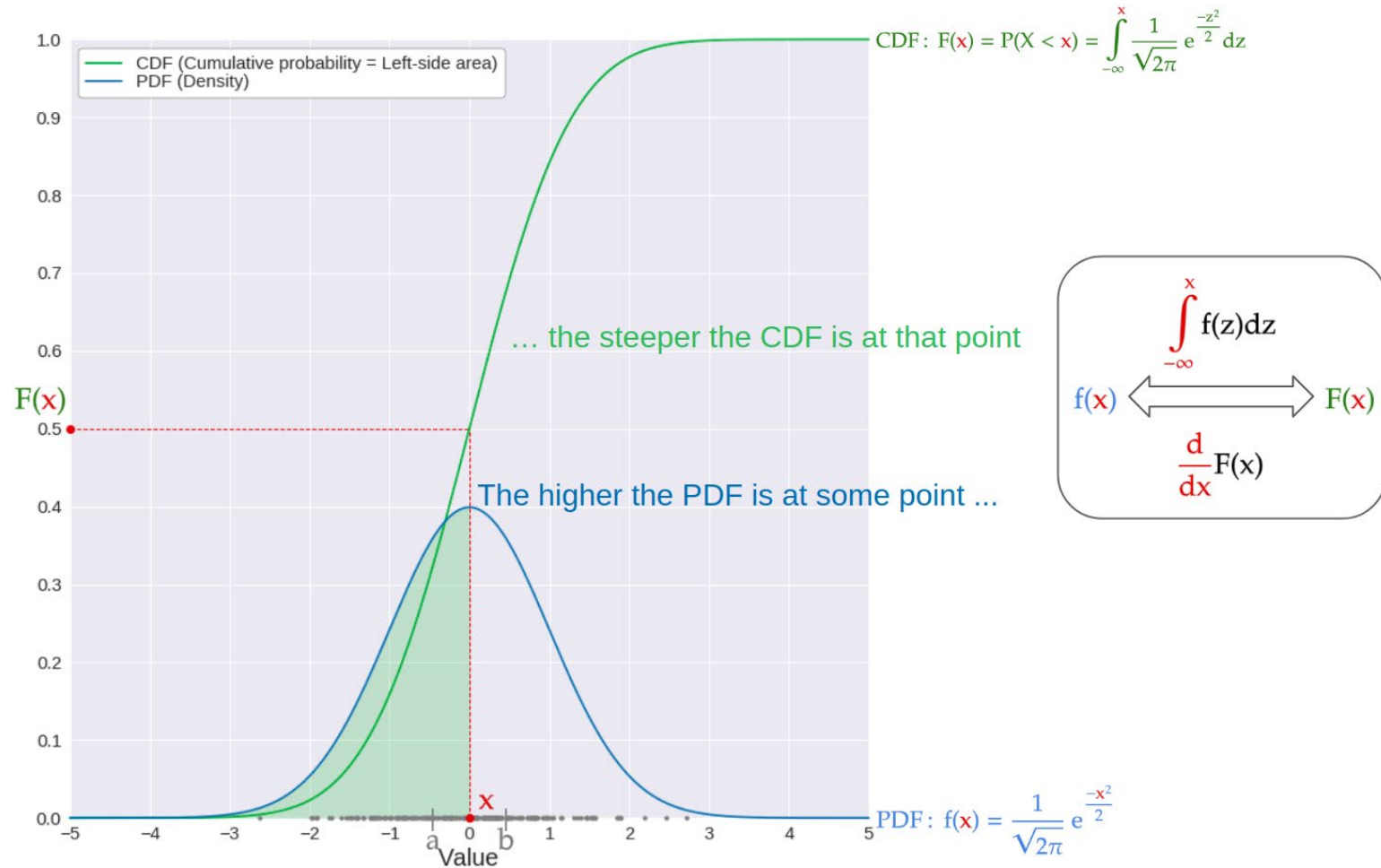
Probability Density Function

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

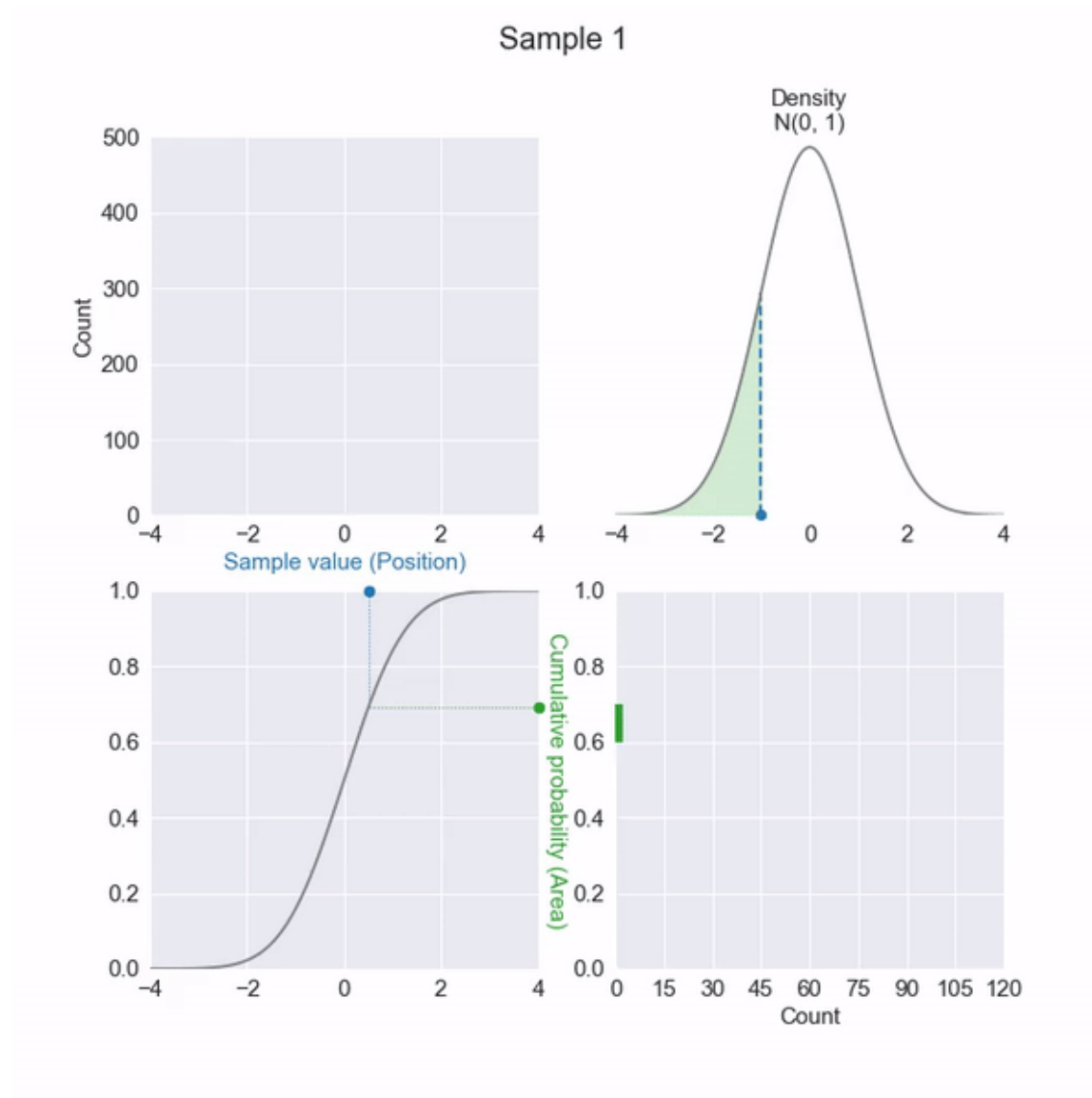
Cumulative Density Function

$$F(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{u^2}{2}} du$$

Inverse transform sampling



Inverse transform sampling



Naïve Bayes

(Discrete Input Features)

Example: Spam Classification

$$\begin{aligned}\operatorname{argmax}_y p(y|x) &= \operatorname{argmax}_y \frac{p(x|y)p(y)}{p(x)} \\ &= \operatorname{argmax}_y p(x|y)p(y)\end{aligned}$$

Example: Spam Classification

- Each vocabulary is one feature dimension
- We encode each email as a feature vector $x \in \{0,1\}^{|V|}$
 - One-hot encoding
- $x_j = 1$, iff the vocabulary x_j appears in the email

$$x = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \\ \dots \\ 0 \end{bmatrix} \in \{0,1\}^{|V|}$$

a
Skku
University
Buy
He
She
...

$$\operatorname{argmax}_y p(\textcolor{green}{x}|\textcolor{red}{y})p(y)$$

y: spam or not

x: input

Example: Spam Classification

- We want to model the probability of any word x_j appearing in an email given the email is spam or not

Issues

- What if $|V|$ (the number of vocabulary) is large?
- Example: $|V| = 3$
- 2^3 possible outcome, 2^3 classification

$$p(x|y)$$

(categorical distribution)

$$\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix} \quad \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \quad \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

Example: Spam Classification

- Naïve Bayes
- $p(x_j|y = k)$ is a Bernoulli distribution

$$p(x|y = k) = \prod_{j=1}^{|V|} p(x_j|y = k)$$

- Not a right assumption in practice
 - If $y=1$ (spam), then, knowledge of 'buy' presence does not affect on your beliefs about other words, e.g., 'price'

Example: Spam Classification

- Both $p(x_j|y = 1)$ and $p(y = 1)$ are Bernoulli distributions

$$p(x|y = k) = \prod_{j=1}^{|V|} p(x_j|y = k) \quad \phi \in [0,1]^{|V| \times 2} \quad \theta \in [0,1]$$

$$p(x_j = 1|y = 1) = \phi_{j1}$$

$$p(y = 1) = \theta$$

$$p(x_j = 0|y = 1) = 1 - \phi_{j1}$$

$$p(y = 0) = 1 - \theta$$

$$p(x_j = 1|y = 0) = \phi_{j0}$$

$$p(x_j = 0|y = 0) = 1 - \phi_{j0}$$

MLE

$$\log L(\phi, \theta) = \log \prod_{i=1}^N p(x^{(i)}, y^{(i)}; \phi, \theta)$$

MLE

$$\begin{aligned}\log L(\phi, \theta) &= \log \prod_{i=1}^N p(x^{(i)}, y^{(i)}; \phi, \theta) \\&= \log \prod_{i=1}^N p(x^{(i)} | y^{(i)}; \phi) p(y^{(i)}; \theta) \\&= \log \prod_{i=1}^N p(y^{(i)}; \theta) \prod_{j=1}^{|V|} p(x_j^{(i)} | y^{(i)}; \phi) \\&= \sum_{i=1}^N \log p(y^{(i)}; \theta) + \sum_{i=1}^N \sum_{j=1}^{|V|} \log p(x_j^{(i)} | y^{(i)}; \phi)\end{aligned}$$

MLE

$$\begin{aligned}\frac{\partial L}{\partial \phi_{l1}} &= \frac{\partial}{\partial \phi_{l1}} \sum_{i=1}^N \log p(y^{(i)}; \theta) + \sum_{i=1}^N \sum_{j=1}^{|V|} \log p(x_j^{(i)} | y^{(i)}; \phi) \\&= \frac{\partial}{\partial \phi_{l1}} \sum_{i \in \{k | y^{(k)}=1\}} \sum_{j=1}^{|V|} \log p(x_j^{(i)} | y^{(i)}; \phi) \\&= \frac{\partial}{\partial \phi_{l1}} \sum_{i \in \{k | y^{(k)}=1\}} \sum_{j=1}^{|V|} \log \left(\phi_{j1}^{x_j^{(i)}} (1 - \phi_{j1})^{(1-x_j^{(i)})} \right) = \frac{\partial}{\partial \phi_{\textcolor{red}{l}1}} \sum_{i \in \{k | y^{(k)}=1\}} \log \left(\phi_{\textcolor{red}{l}1}^{x_{\textcolor{red}{l}}^{(i)}} (1 - \phi_{l1})^{(1-x_{\textcolor{red}{l}}^{(i)})} \right) \\&= \frac{\partial}{\partial \phi_{l1}} \sum_{i \in \{k | y^{(k)}=1\}} x_l^{(i)} \log \phi_{l1} + (1 - x_l^{(i)}) \log(1 - \phi_{l1}) \\&= \sum_{i \in \{k | y^{(k)}=1\}} \frac{x_l^{(i)}}{\phi_{l1}} - \frac{1 - x_l^{(i)}}{(1 - \phi_{l1})} = 0\end{aligned}$$

MLE

$$\sum_{i \in \{k|y^{(k)}=1\}} \frac{x_l^{(i)}}{\phi_{l1}} = \sum_{i \in \{k|y^{(k)}=1\}} \frac{1 - x_l^{(i)}}{(1 - \phi_{l1})}$$

$$\frac{1}{\phi_{l1}} \sum_{i \in \{k|y^{(k)}=1\}} x_l^{(i)} = \frac{1}{(1 - \phi_{l1})} \sum_{i \in \{k|y^{(k)}=1\}} 1 - x_l^{(i)}$$

$$(1 - \phi_{l1}) \sum_{i \in \{k|y^{(k)}=1\}} x_l^{(i)} = \phi_{l1} \sum_{i \in \{k|y^{(k)}=1\}} 1 - x_l^{(i)}$$

$$\sum_{i \in \{k|y^{(k)}=1\}} x_l^{(i)} = \phi_{l1} \sum_{i \in \{k|y^{(k)}=1\}} 1$$

$$\sum_{i \in \{k|y^{(k)}=1\}} x_l^{(i)} = \phi_{l1} \sum_{i=1}^N 1\{y^{(i)} = 1\}$$

$$\phi_{l1}^* = \frac{\sum_{i=1}^N 1\{y^{(i)} = 1\} x_l^{(i)}}{\sum_{i=1}^N 1\{y^{(i)} = 1\}} = \frac{\sum_{i=1}^N 1\{y^{(i)} = 1\} x_l^{(i)}}{N_1}$$

MLE

$$\sum_{i \in \{k|y^{(k)}=1\}} \frac{x_l^{(i)}}{\phi_{l1}} = \sum_{i \in \{k|y^{(k)}=1\}} \frac{1 - x_l^{(i)}}{(1 - \phi_{l1})}$$

$$\frac{1}{\phi_{l1}} \sum_{i \in \{k|y^{(k)}=1\}} x_l^{(i)} = \frac{1}{(1 - \phi_{l1})} \sum_{i \in \{k|y^{(k)}=1\}} 1 - x_l^{(i)}$$

$$(1 - \phi_{l1}) \sum_{i \in \{k|y^{(k)}=1\}} x_l^{(i)} = \phi_{l1} \sum_{i \in \{k|y^{(k)}=1\}} 1 - x_l^{(i)}$$

$$\sum_{i \in \{k|y^{(k)}=1\}} x_l^{(i)} = \phi_{l1} \sum_{i \in \{k|y^{(k)}=1\}} 1$$

$$\sum_{i \in \{k|y^{(k)}=1\}} x_l^{(i)} = \phi_{l1} \sum_{i=1}^N 1\{y^{(i)} = 1\}$$

$$\phi_{l1}^* = \frac{\sum_{i=1}^N 1\{y^{(i)} = 1\} x_l^{(i)}}{\sum_{i=1}^N 1\{y^{(i)} = 1\}} = \frac{\sum_{i=1}^N 1\{y^{(i)} = 1\} x_l^{(i)}}{N_1}$$

$$\phi_{l0}^* = \frac{\sum_{i=1}^N 1\{y^{(i)} = 0\} x_l^{(i)}}{\sum_{i=1}^N 1\{y^{(i)} = 0\}} = \frac{\sum_{i=1}^N 1\{y^{(i)} = 0\} x_l^{(i)}}{N_0}$$

MLE

$$\sum_{i \in \{k|y^{(k)}=1\}} \frac{x_l^{(i)}}{\phi_{l1}} = \sum_{i \in \{k|y^{(k)}=1\}} \frac{1 - x_l^{(i)}}{(1 - \phi_{l1})}$$

$$\frac{1}{\phi_{l1}} \sum_{i \in \{k|y^{(k)}=1\}} x_l^{(i)} = \frac{1}{(1 - \phi_{l1})} \sum_{i \in \{k|y^{(k)}=1\}} 1 - x_l^{(i)}$$

$$(1 - \phi_{l1}) \sum_{i \in \{k|y^{(k)}=1\}} x_l^{(i)} = \phi_{l1} \sum_{i \in \{k|y^{(k)}=1\}} 1 - x_l^{(i)}$$

$$\sum_{i \in \{k|y^{(k)}=1\}} x_l^{(i)} = \phi_{l1} \sum_{i \in \{k|y^{(k)}=1\}} 1$$

$$\sum_{i \in \{k|y^{(k)}=1\}} x_l^{(i)} = \phi_{l1} \sum_{i=1}^N 1\{y^{(i)} = 1\}$$

$$\phi_{l1}^* = \frac{\sum_{i=1}^N 1\{y^{(i)} = 1\} x_l^{(i)}}{\sum_{i=1}^N 1\{y^{(i)} = 1\}} = \frac{\sum_{i=1}^N 1\{y^{(i)} = 1\} x_l^{(i)}}{N_1}$$

$$\phi_{l0}^* = \frac{\sum_{i=1}^N 1\{y^{(i)} = 0\} x_l^{(i)}}{\sum_{i=1}^N 1\{y^{(i)} = 0\}} = \frac{\sum_{i=1}^N 1\{y^{(i)} = 0\} x_l^{(i)}}{N_0}$$

$$\phi_{lk}^* = \frac{\sum_{i=1}^N 1\{y^{(i)} = k\} x_l^{(i)}}{\sum_{i=1}^N 1\{y^{(i)} = k\}} = \frac{\sum_{i=1}^N 1\{y^{(i)} = k\} x_l^{(i)}}{N_k}$$

MLE

$$\begin{aligned}\frac{\partial L}{\partial \theta} \log L(\phi, \theta) &= \frac{\partial L}{\partial \theta} \sum_{i=1}^N \log p(y^{(i)}; \theta) + \sum_{i=1}^N \sum_{j=1}^{|V|} \log p(x_j^{(i)} | y^{(i)}; \phi) \\ &= \frac{\partial L}{\partial \theta} \sum_{i=1}^N \log p(y^{(i)}; \theta)\end{aligned}$$

$$\theta^* = \frac{\sum_{i=1}^N 1\{y^{(i)} = 1\}}{N} = \frac{N_1}{N}$$

$$\theta_k^* = \frac{\sum_{i=1}^N 1\{y^{(i)} = k\}}{N} = \frac{N_k}{N}$$

Testing

$$\operatorname{argmax}_k p(y = k) \prod_{j=1}^d p(x_j^{new} | y = k)$$

Testing

$$\operatorname{argmax}_k p(y = k) \prod_{j=1}^d p(x_j^{new} | y = k)$$

$$\operatorname{argmax}_k \theta_k^* \prod_{j=1}^d \phi_{jk}^* x_j^{new} (1 - \phi_{jk}^*)^{1-x_j^{new}}$$

Laplace Smoothing

- What if we have not seen a word “skku” before?
- Then,

$$p(x_{30}|y = 1; \phi) = \phi_{30,1} = 0$$

$$\begin{aligned} p(y = 1|x) &= \frac{p(y = 1) \prod_{j=1}^d p(x_j|y = 1)}{p(x)} \\ &= \frac{p(y = 1) \prod_{j=1}^d p(x_j|y = 1)}{p(y = 0) \prod_{j=1}^d p(x_j|y = 0) + p(y = 1) \prod_{j=1}^d p(x_j|y = 1)} = \frac{0}{0} \end{aligned}$$

Laplace Smoothing

- Statistically, it is a bad idea to say probability is 'zero' just because you haven't seen it!
- So, add '1' to numerator, K to denominator

$$\theta_k^* = \frac{\sum_{i=1}^N 1\{y^{(i)} = k\}}{N} \longrightarrow \theta_k^* = \frac{1 + \sum_{i=1}^N 1\{y^{(i)} = k\}}{K + N}$$

$$\phi_{lk}^* = \frac{\sum_{i=1}^N 1\{y^{(i)} = 1\}x_l^{(i)}}{N_k} \longrightarrow \phi_{lk}^* = \frac{1 + \sum_{i=1}^N 1\{y^{(i)} = 1\}x_l^{(i)}}{K + N_k}$$