# Lecture 6: Classification with Mixtures of Gaussians CS480/680 Intro to Machine Learning

2023-1-26
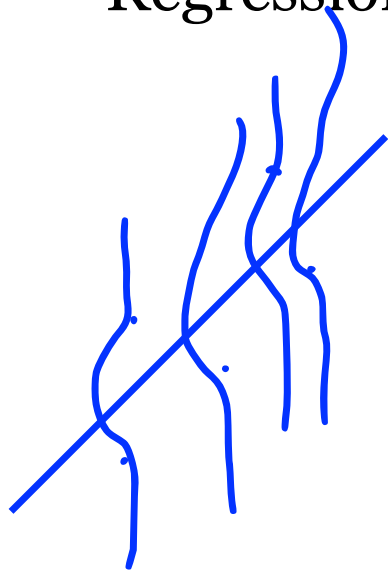
Pascal Poupart
David R. Cheriton School of Computer Science

UNIVERSITY OF
WATERLOO

# Linear Models

- Probabilistic Generative Models

Regression

Classification

UNIVERSITY OF
WATERLOO

# Probabilistic Generative Model

- $\Pr(C)$: prior probability of class $C$

- $\Pr(\boldsymbol{x}|C)$: class conditional distribution of $\boldsymbol{x}$

- Classification: compute posterior $\Pr(C|\boldsymbol{x})$ according to Bayes' theorem

$$\Pr(C|\boldsymbol{x}) = \frac{\Pr(\boldsymbol{x}|C)\Pr(C)}{\sum_C \Pr(\boldsymbol{x}|C)\Pr(C)}$$

$$= k\Pr(\boldsymbol{x}|C)\Pr(C)$$

*normalization constant*

UNIVERSITY OF
WATERLOO

# Assumptions

- In classification, the number of classes is finite, so a natural prior $\Pr(C)$ is the multinomial

$$\Pr(C = c_k) = \pi_k$$

- When $\boldsymbol{x} \in \Re^d$, then it is often OK to assume that $\Pr(\boldsymbol{x}|C)$ is Gaussian.

- Furthermore, assume that the same covariance matrix $\boldsymbol{\Sigma}$ is used for each class.

$$\Pr(\boldsymbol{x}|c_k) \propto e^{-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu_k})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu_k})}$$

UNIVERSITY OF
WATERLOO

# Posterior Distribution

$$\Pr(c_k | \boldsymbol{x}) = \frac{\pi_k e^{-\frac{1}{2}(x-\boldsymbol{\mu_k})^T \Sigma^{-1}(x-\boldsymbol{\mu_k})}}{\sum_k \pi_k e^{-\frac{1}{2}(x-\boldsymbol{\mu_k})^T \Sigma^{-1}(x-\boldsymbol{\mu_k})}}$$

$$= \frac{\pi_k e^{-\frac{1}{2}\left(\cancel{x^T \Sigma^{-1}x} - 2\mu_k^T \Sigma^{-1}x + \mu_k^T \Sigma^{-1}\mu_k\right)}}{\sum_k \pi_k e^{-\frac{1}{2}\left(\cancel{x^T \Sigma^{-1}x} - 2\mu_k^T \Sigma^{-1}x + \mu_k^T \Sigma^{-1}u_k\right)}}$$

Consider two classes $c_k$ and $c_j$

$$= \frac{1}{1 + \dfrac{\pi_j e^{\mu_j^T \Sigma^{-1}x - \frac{1}{2}\mu_j^T \Sigma^{-1}\mu_j}}{\pi_k e^{\mu_k^T \Sigma^{-1}x - \frac{1}{2}\mu_k^T \Sigma^{-1}\mu_k}}}$$

UNIVERSITY OF
WATERLOO

# Posterior Distribution

$$= \frac{1}{1+e^{-\left(\mu_k^T - \mu_j^T\right)\Sigma^{-1}x + \frac{1}{2}\mu_k^T\Sigma^{-1}\mu_k - \frac{1}{2}\mu_j^T\Sigma^{-1}\mu_j - \ln\frac{\pi_k}{\pi_j}}}$$

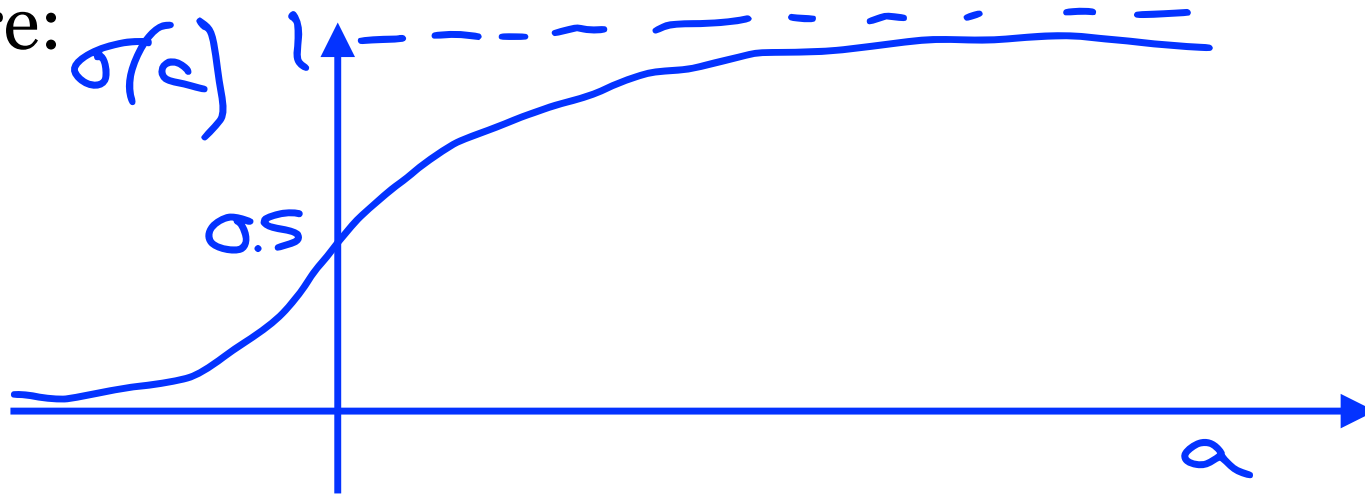$$= \frac{1}{1+e^{-(w^T x + w_0)}}$$

where $\boldsymbol{w} = \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_k - \boldsymbol{\mu}_j)$

and $w_0 = -\frac{1}{2}\boldsymbol{\mu}_k^T\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_k + \frac{1}{2}\boldsymbol{\mu}_j^T\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_j + \ln\frac{\pi_k}{\pi_j}$
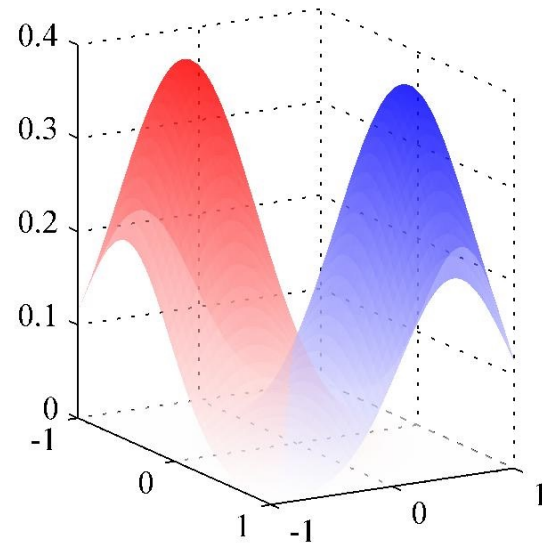
UNIVERSITY OF
WATERLOO

# Logistic Sigmoid

- Let $\sigma(a) = \dfrac{1}{1+e^{-a}}$
  $$\longrightarrow \text{Logistic sigmoid}$$

- Then $\Pr(c_k|\boldsymbol{x}) = \sigma(\boldsymbol{w}^T\boldsymbol{x} + w_0)$

- Picture:

# Logistic Sigmoid

class conditionals

posterior

UNIVERSITY OF
WATERLOO

# Prediction

$$best\ class = argmax_k \Pr(c_k | \boldsymbol{x})$$

$$= \begin{cases} c_1 & \sigma(\boldsymbol{w}^T \boldsymbol{x} + w_0) \geq 0.5 \\ c_2 & \text{otherwise} \end{cases}$$

Class boundary: $\sigma(\boldsymbol{w}_k^T \overline{\boldsymbol{x}}) = 0.5$

$$\implies \frac{1}{1+e^{-(\boldsymbol{w}_k^T \overline{\boldsymbol{x}})}} = 0.5$$

$$\implies \boldsymbol{w}_k^T \overline{\boldsymbol{x}} = 0$$

$$\therefore \textbf{linear separator}$$

# Multi-class Problems

- Consider Gaussian conditional distributions with identical $\Sigma$

$$\Pr(c_k|\boldsymbol{x}) = \frac{\Pr(c_k)\Pr(\boldsymbol{x}|c_k)}{\sum_j \Pr(c_j)\Pr(\boldsymbol{x}|c_j)}$$

$$= \frac{\pi_k e^{-\frac{1}{2}(x-\mu_k)^T \Sigma^{-1}(x-\mu_k)}}{\sum_j \pi_j e^{-\frac{1}{2}(x-\mu_j)^T \Sigma^{-1}(x-\mu_j)}}$$

$$= \frac{\pi_k e^{-\frac{1}{2}(-2\mu_k^T \Sigma^{-1}x + \mu_k^T \Sigma^{-1}\mu_k)}}{\sum_j \pi_j e^{-\frac{1}{2}(-2\mu_j^T \Sigma^{-1}x + \mu_j^T \Sigma^{-1}\mu_j)}}$$

$$= \frac{e^{\mu_k^T \Sigma^{-1}x - \frac{1}{2}\mu_k^T \Sigma^{-1}u_k + \ln \pi_k}}{\sum_j e^{\mu_j^T \Sigma^{-1}x - \frac{1}{2}\mu_j^T \Sigma^{-1}\mu_j + \ln \pi_j}} = \frac{e^{w_k^T \bar{x}}}{\sum_j e^{w_j^T \bar{x}}} \implies \text{softmax}$$

where $\boldsymbol{w}_k^T = (-\frac{1}{2}\boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_k + \ln \pi_k, \ \boldsymbol{\mu}_k^T \boldsymbol{\Sigma}^{-1})$

# Softmax

- When there are several classes, the posterior is a **softmax** (sigmoid generalization)

- Softmax distribution: $\Pr(c_k|\boldsymbol{x}) = \dfrac{e^{f_k(\boldsymbol{x})}}{\sum_j e^{f_j(\boldsymbol{x})}}$

- Argmax distribution:

$$\Pr(c_k|\boldsymbol{x}) = \begin{cases} 1 & \text{if } k = argmax_j \, f_j(x) \\ 0 & \text{otherwise} \end{cases}$$

$$= \lim_{base\to\infty} \frac{base^{f_k(x)}}{\sum_j base^{f_j(x)}}$$

$$\approx \frac{e^{f_k(x)}}{\sum_j e^{f_j(x)}} \quad \text{(softmax approximation)}$$

UNIVERSITY OF
**WATERLOO**

# Softmax

class conditionals

posterior

# Parameter Estimation

- Where do $\Pr(c_k)$ and $\Pr(\boldsymbol{x}|c_k)$ come from?

- Parameters: $\pi, \boldsymbol{\mu_1}, \boldsymbol{\mu_2}, \boldsymbol{\Sigma}$

    $$\Pr(c_1) = \pi, \qquad \Pr(\boldsymbol{x}|c_1) = k_{\boldsymbol{\Sigma}}\, e^{-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu_1})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu_1})}$$

    $$\Pr(c_2) = 1 - \pi, \quad \Pr(\boldsymbol{x}|c_2) = k_{\boldsymbol{\Sigma}}\, e^{-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu_2})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu_2})}$$
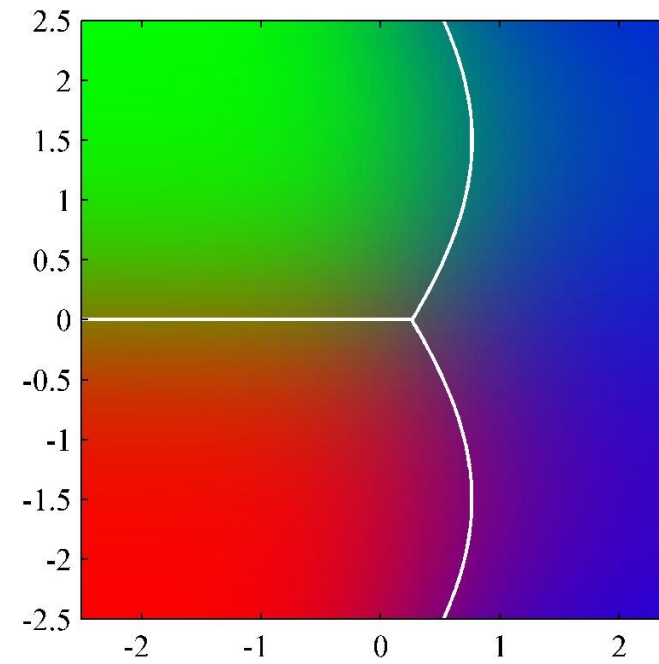
    where $k_{\boldsymbol{\Sigma}}$ is the normalization constant that depends on $\boldsymbol{\Sigma}$

- Estimate parameters by
    - **Maximum likelihood**
    - Maximum a posteriori
    - Bayesian learning

UNIVERSITY OF
**WATERLOO**

# Maximum Likelihood Solution

- Likelihood:

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

$$L(\mathbf{X}, \mathbf{y}) = \Pr(\mathbf{X}, \mathbf{y} | \pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \prod_n [\pi N(\boldsymbol{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})]^{y_n} [(1-\pi) N(\boldsymbol{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})]^{1-y_n}$$

- ML hypothesis:

$$< \pi^*, \boldsymbol{\mu}_1^*, \boldsymbol{\mu}_2^*, \boldsymbol{\Sigma}^* > = argmax_{\pi, \boldsymbol{\mu}_2, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}} \sum_n y_n \left[ \ln \pi + \ln k_{\boldsymbol{\Sigma}} - \frac{1}{2} (\boldsymbol{x}_n - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{x}_n - \boldsymbol{\mu}_1) \right]$$

$$+ (1 - y_n) \left[ \ln(1 - \pi) + \ln k_{\boldsymbol{\Sigma}} - \frac{1}{2} (\boldsymbol{x}_n - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{x}_n - \boldsymbol{\mu}_2) \right]$$

UNIVERSITY OF
WATERLOO

# Maximum Likelihood Solution

- Set derivative to 0

$$0 = \frac{\partial \ln L(\boldsymbol{X,y})}{\partial \pi}$$

$$\implies 0 = \sum_n y_n \left[\frac{1}{\pi}\right] + (1 - y_n)\left[-\frac{1}{1-\pi}\right]$$

$$\implies 0 = \sum_n y_n(1 - \pi) + (1 - y_n)(-\pi)$$

$$\implies \sum_n y_n = \pi\left(\sum_n y_n + \sum_n(1 - y_n)\right)$$

$$\implies \sum_n y_n = \pi N \text{ (where } N \text{ is the \# of training points)}$$

$$\therefore \frac{\sum_n y_n}{N} = \pi$$

UNIVERSITY OF
WATERLOO

# Maximum Likelihood Solution

$$0 = \partial \ln L(\boldsymbol{X}, \boldsymbol{y}) / \partial \boldsymbol{\mu}_1$$

$$\Longrightarrow 0 = \sum_n y_n [-\boldsymbol{\Sigma}^{-1}(\boldsymbol{x}_n - \boldsymbol{\mu}_1)]$$

$$\Longrightarrow \sum_n y_n \boldsymbol{x}_n = \sum_n y_n \boldsymbol{\mu}_1$$

$$\Longrightarrow \sum_n y_n \boldsymbol{x}_n = N_1 \boldsymbol{\mu}_1$$

$$\therefore \frac{\sum_n y_n \boldsymbol{x}_n}{N_1} = \boldsymbol{\mu}_1 \quad \text{Similarly: } \frac{\sum_n (1-y_n)\boldsymbol{x}_n}{N_2} = \boldsymbol{\mu}_2$$

where $N_1$ is the # of data points in class 1

$N_2$ is the # of data points in class 2

UNIVERSITY OF
WATERLOO

# Maximum Likelihood

$$\frac{\partial \ln L(\boldsymbol{X}, \boldsymbol{y})}{\partial \Sigma} = 0$$

$$\Longrightarrow \cdots$$

$$\Longrightarrow \boxed{\Sigma = \frac{N_1}{N} \boldsymbol{S}_1 + \frac{N_2}{N} \boldsymbol{S}_2}$$

where $\boldsymbol{S}_1 = \frac{1}{N_1} \sum_{n \in c_1} (\boldsymbol{x}_n - \boldsymbol{\mu}_1)(\boldsymbol{x}_n - \boldsymbol{\mu}_1)^T$

$$\boldsymbol{S}_2 = \frac{1}{N_2} \sum_{n \in c_2} (\boldsymbol{x}_n - \boldsymbol{\mu}_2)(\boldsymbol{x}_n - \boldsymbol{\mu}_2)^T$$

($\boldsymbol{S}_k$ is the empirical covariance matrix of class $k$)

UNIVERSITY OF
WATERLOO