# Fundamentals of Machine Learning

## LINEAR MODELS

Amit K Roy-Chowdhury
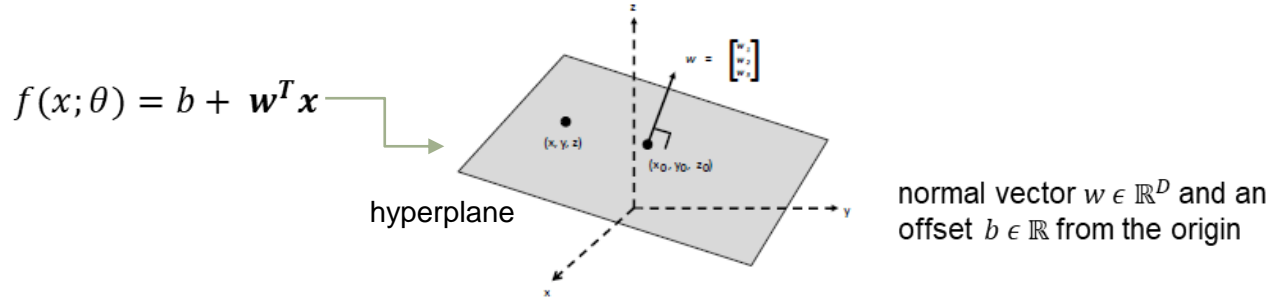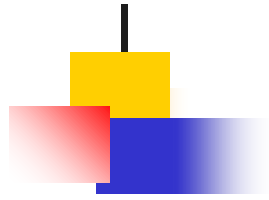
# Linear Classifier

The prediction can be written as

$$f(x) = \mathbb{I}\left(p(y=1|x) > p(y=0|x)\right) = \mathbb{I}\left(\log \frac{p(y=1|x)}{p(y=0|x)} > 0\right) = \mathbb{I}\left(a > 0\right)$$

$$f(x;\theta) = b + \boldsymbol{w}^T \boldsymbol{x}$$

hyperplane

normal vector $w \in \mathbb{R}^D$ and an offset $b \in \mathbb{R}$ from the origin

This linear hyperplane separate 3d space into half $\square$ decision boundary

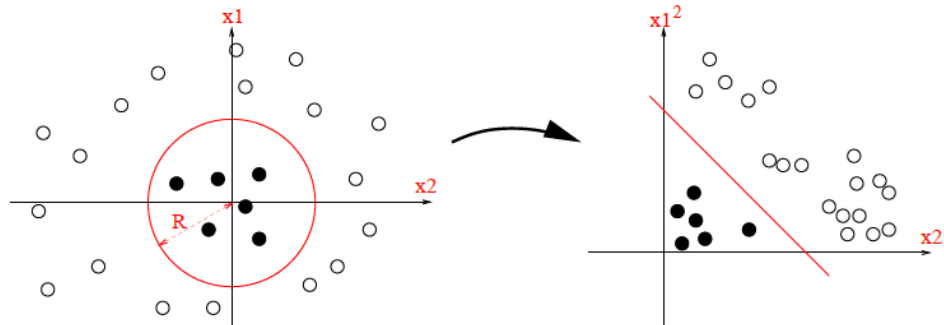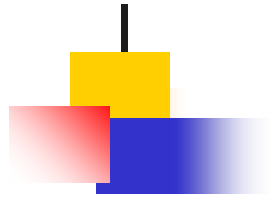# Non Linear Classifier

Transform input features in suitable way

$$\phi(x_1, x_2) = [1, x_1^2, x_2^2]$$

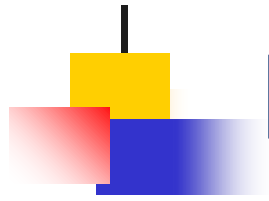$$w = [-R^2, 1, 1]. \text{ Then } w^\mathsf{T}\phi(x) = x_1^2 + x_2^2 - R^2$$

Decision boundary (where f(x) = 0) defines a circle with radius R

# Outline

- Logistic Regression
- Linear Regression
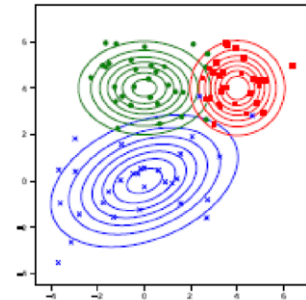- **Linear Discriminant Analysis**
- Naïve Bayes

# Linear Discriminant Analysis

class conditional

prior

$$p(y = c|\boldsymbol{x}; \boldsymbol{\theta}) = \frac{p(\boldsymbol{x}|y = c; \boldsymbol{\theta})p(y = c; \boldsymbol{\theta})}{\sum_{c'} p(\boldsymbol{x}|y = c'; \boldsymbol{\theta})p(y = c'; \boldsymbol{\theta})}$$

posterior

Linear Discriminant Analysis: $\quad \log p(y = c|\boldsymbol{x}; \boldsymbol{\theta}) = \boldsymbol{w}^{\mathsf{T}}\boldsymbol{x} + \text{const}$

Gaussian Discriminant Analysis: $\quad p(\boldsymbol{x}|y = c, \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)$

$\implies p(y = c|\boldsymbol{x}, \boldsymbol{\theta}) \propto \pi_c \mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) \quad$ where $\pi_c = p(y = c)$

# Linear Discriminant Analysis

class conditional

prior

$$p(y = c|\boldsymbol{x}; \boldsymbol{\theta}) = \frac{p(\boldsymbol{x}|y = c; \boldsymbol{\theta})p(y = c; \boldsymbol{\theta})}{\sum_{c'} p(\boldsymbol{x}|y = c'; \boldsymbol{\theta})p(y = c'; \boldsymbol{\theta})}$$

posterior

Gaussian Discriminant Analysis:
$$p(\boldsymbol{x}|y = c, \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)$$

$$\Longrightarrow \quad p(y = c|\boldsymbol{x}, \boldsymbol{\theta}) \propto \pi_c \mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) \quad \text{where } \pi_c = p(y = c)$$

Discriminant Function:
$$\log p(y = c|\boldsymbol{x}, \boldsymbol{\theta}) = \log \pi_c - \frac{1}{2}\log|2\pi\boldsymbol{\Sigma}_c| - \frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu}_c)^\mathsf{T}\boldsymbol{\Sigma}_c^{-1}(\boldsymbol{x} - \boldsymbol{\mu}_c) + \text{const}$$

# Linear Discriminant Analysis

class conditional

prior

$$p(y = c | \boldsymbol{x}; \boldsymbol{\theta}) = \frac{p(\boldsymbol{x} | y = c; \boldsymbol{\theta}) p(y = c; \boldsymbol{\theta})}{\sum_{c'} p(\boldsymbol{x} | y = c'; \boldsymbol{\theta}) p(y = c'; \boldsymbol{\theta})}$$

posterior

**Gaussian Discriminant Analysis:** $\quad p(\boldsymbol{x} | y = c, \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{x} | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)$

$$\implies p(y = c | \boldsymbol{x}, \boldsymbol{\theta}) \propto \pi_c \mathcal{N}(\boldsymbol{x} | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) \quad \text{where } \pi_c = p(y = c)$$

**Discriminant Function:** $\quad \log p(y = c | \boldsymbol{x}, \boldsymbol{\theta}) = \log \pi_c - \frac{1}{2} \log |2\pi \boldsymbol{\Sigma}_c| - \frac{1}{2} (\boldsymbol{x} - \boldsymbol{\mu}_c)^{\mathsf{T}} \boldsymbol{\Sigma}_c^{-1} (\boldsymbol{x} - \boldsymbol{\mu}_c) + \text{const}$

$\boldsymbol{\Sigma}_c = \boldsymbol{\Sigma} \quad \log p(y = c | \boldsymbol{x}, \boldsymbol{\theta}) = \log \pi_c - \frac{1}{2} (\boldsymbol{x} - \boldsymbol{\mu}_c)^{\mathsf{T}} \boldsymbol{\Sigma}^{-1} (\boldsymbol{x} - \boldsymbol{\mu}_c) + \text{const}$

$$= \underbrace{\log \pi_c - \frac{1}{2} \boldsymbol{\mu}_c^{\mathsf{T}} \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c}_{\gamma_c} + \boldsymbol{x}^{\mathsf{T}} \underbrace{\boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c}_{\beta_c} + \underbrace{\text{const} - \frac{1}{2} \boldsymbol{x}^{\mathsf{T}} \boldsymbol{\Sigma}^{-1} \boldsymbol{x}}_{\kappa}$$

**LInear Discriminant Analysis:** $\quad = \gamma_c + \boldsymbol{x}^{\mathsf{T}} \boldsymbol{\beta}_c + \kappa$

# Interpretation of LDA

## Uniform prior over classes

$$\hat{y}(x) = \underset{c}{\arg\max} \log p(y = c | x, \theta) = \underset{c}{\arg\min}(x - \mu_c)^\mathsf{T} \Sigma^{-1}(x - \mu_c)$$

nearest centroid classifier or nearest class mean classifier

# Fisher's LDA

Reduce feature dimensionality (PCA!) and classify.

find the matrix $\mathbf{W}$ such that the low-dimensional data can be classified as well as possible

$x \in \mathbb{R}^D$

$z \in \mathbb{R}^K$

$$z_n = \mathbf{W}x_n \qquad m_c = \frac{1}{N_c}\sum_{n:y_n=c} z_n \qquad m = \frac{1}{N}\sum_{c=1}^{C} N_c m_c$$

data points ⟶ mean for class c ⟶ overall mean

$$\tilde{\mathbf{S}}_W = \sum_{c=1}^{C}\sum_{n:y_n=c}(z_n - m_c)(z_n - m_c)^{\mathsf{T}}$$

scatter matrices

$$\tilde{\mathbf{S}}_B = \sum_{c=1}^{C} N_c(m_c - m)(m_c - m)^{\mathsf{T}}$$

maximize objective function

$$J(\mathbf{W}) = \frac{|\tilde{\mathbf{S}}_B|}{|\tilde{\mathbf{S}}_W|} = \frac{|\mathbf{W}^{\mathsf{T}}\mathbf{S}_B\mathbf{W}|}{|\mathbf{W}^{\mathsf{T}}\mathbf{S}_W\mathbf{W}|}$$

Leads to a generalized eigenvalue problem – advanced reading

UC**R** | Bourns College of Engineering

# FLDA – 2 classes

$$\mathbf{S}_B = (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^\mathsf{T}$$

$$J(\boldsymbol{w}) = \frac{\boldsymbol{w}^\mathsf{T}\mathbf{S}_B\boldsymbol{w}}{\boldsymbol{w}^\mathsf{T}\mathbf{S}_W\boldsymbol{w}}$$

$$\mathbf{S}_W = \sum_{n:y_n=1}(\boldsymbol{x}_n - \boldsymbol{\mu}_1)(\boldsymbol{x}_n - \boldsymbol{\mu}_1)^\mathsf{T} + \sum_{n:y_n=2}(\boldsymbol{x}_n - \boldsymbol{\mu}_2)(\boldsymbol{x}_n - \boldsymbol{\mu}_2)^\mathsf{T}$$

$$\boldsymbol{\mu}_1 = \frac{1}{N_1}\sum_{n:y_n=1}\boldsymbol{x}_n, \quad \boldsymbol{\mu}_2 = \frac{1}{N_2}\sum_{n:y_n=2}\boldsymbol{x}_n$$

Take derivative wrt $\boldsymbol{w}$

$$\mathbf{S}_B\boldsymbol{w} = \lambda\mathbf{S}_W\boldsymbol{w} \quad \text{where} \quad \lambda = \frac{\boldsymbol{w}^\mathsf{T}\mathbf{S}_B\boldsymbol{w}}{\boldsymbol{w}^\mathsf{T}\mathbf{S}_W\boldsymbol{w}}$$

Generalized eigenvalue problem, becomes regular eigenvalue problem if $\mathbf{S}_W^{-1}\mathbf{S}_B\boldsymbol{w} = \lambda\boldsymbol{w}$

**Interpretation:**

Let $m_k = \boldsymbol{w}^\mathsf{T}\boldsymbol{\mu}_k$ be the projection of each mean onto the line $\boldsymbol{w}$.

$$\mathbf{S}_B\boldsymbol{w} = (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^\mathsf{T}\boldsymbol{w} = (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)(m_2 - m_1)$$

$$\lambda\,\boldsymbol{w} = \mathbf{S}_W^{-1}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)(m_2 - m_1)$$

$$\boldsymbol{w} \propto \mathbf{S}_W^{-1}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) \qquad \boldsymbol{w} \text{ is proportional to the vector that joins the class means}$$

# Outline

- Logistic Regression
- Linear Regression
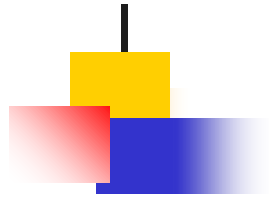- Linear Discriminant Analysis
- Naïve Bayes

# Naïve Bayes

$$P(y|\mathbf{x}) = \frac{P(\mathbf{x}|y)P(y)}{P(\mathbf{x})}$$

Naïve Bayes Assumption:

$$P(\mathbf{x}|y) = \prod_{\alpha=1}^{d} P(x_\alpha|y), \text{ where } x_\alpha = [\mathbf{x}]_\alpha \text{ is the value for feature } \alpha$$

$$
\begin{aligned}
h(\mathbf{x}) &= \operatorname*{argmax}_{y} P(y|\mathbf{x}) \\
&= \operatorname*{argmax}_{y} \frac{P(\mathbf{x}|y)P(y)}{P(\mathbf{x})} \\
&= \operatorname*{argmax}_{y} P(\mathbf{x}|y)P(y) && (P(\mathbf{x}) \text{ does not depend on } y) \\
&= \operatorname*{argmax}_{y} \prod_{\alpha=1}^{d} P(x_\alpha|y)P(y) && (\text{by the naive Bayes assumption}) \\
&= \operatorname*{argmax}_{y} \sum_{\alpha=1}^{d} \log(P(x_\alpha|y)) + \log(P(y)) && (\text{as log is a monotonic function})
\end{aligned}
$$

# Simple Example

- Given N1 emails which are spam and N2 not spam; $p(S) = N1/(N1+ N2)$; $p(NS) = N2/N1+N2$

- Consider some words that occur in each category with some frequency: $p(w1|S)$, $p(w2|S)$, $P(w3|S)$,......;

$$p(w1|NS), p(w2|NS),.....$$

- Say you observe $W = \{w1, w5, w7\}$ in an email. Is it spam or not spam?

- $p(S|W) \propto p(S) \times p(W|S)$ – now use conditional independence

- $P(NS|W) \propto p(NS) \times p(W|NS)$ – now use conditional independence

- What happens if one word, say w7, never occurred in the training data for NS?

- What is it we needed to know to calculate this?
    - Number of times each word occurred in each class as a fraction of all the words in that class.