

DS 5220

Supervised Machine Learning and Learning Theory

Alina Oprea
Associate Professor, CCIS
Northeastern University

September 23 2019

Logistics

- HW 1 is due today, Sept. 23
- Midterm
 - Monday, Oct. 28
- Final exam
 - In the last class, Wed. Dec. 4
- Project report due during the final exam week
- Project milestone
 - Form teams of 2-3 people
 - Submit project proposal on Oct. 16
 - Project pitch on Oct. 21 in class

Outline

- Linear classifier
 - LDA on one dimension
 - Linear discriminant functions
 - Multi-variate LDA
- Bias-variance tradeoff
 - Derivation for linear regression

LDA

- Classify to one of k classes (multi-class)
- LDA uses Bayes Theorem to estimate it

$$- P[Y = k|X = x] = \frac{P[X = x|Y = k]P[Y=k]}{P[X=x]}$$

– Let $\pi_k = P[Y = k]$ be the prior probability of class k and $f_k(x) = P[X = x|Y = k]$

- **Generative model**
 - Given X and Y , learns the joint probability $P(X, Y)$
- **Discriminative model**
 - Given X and Y , learns a decision function for classification

LDA

$$\Pr(Y = k|X = x) = \frac{\pi_k f_k(x)}{\sum_{l=1}^K \pi_l f_l(x)}.$$

Assume $f_k(x)$ is Gaussian!
Unidimensional case (d=1)

$$f_k(x) = \frac{1}{\sqrt{2\pi}\sigma_k} \exp\left(-\frac{1}{2\sigma_k^2}(x - \mu_k)^2\right)$$

$$p_k(x) = \frac{\pi_k \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x - \mu_k)^2\right)}{\sum_{l=1}^K \pi_l \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x - \mu_l)^2\right)}.$$

Assumption: $\sigma_1 = \dots \sigma_k = \sigma$

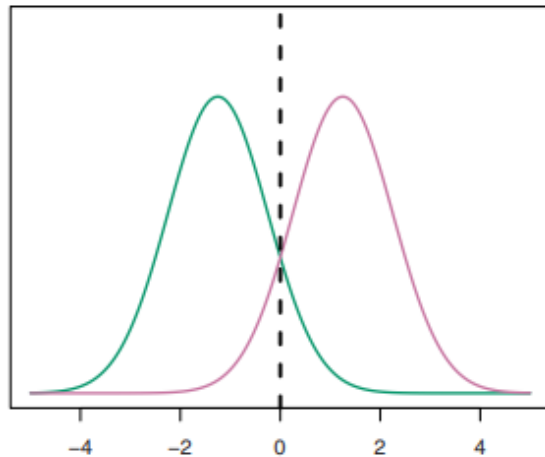
LDA decision boundary

Pick class k to maximize

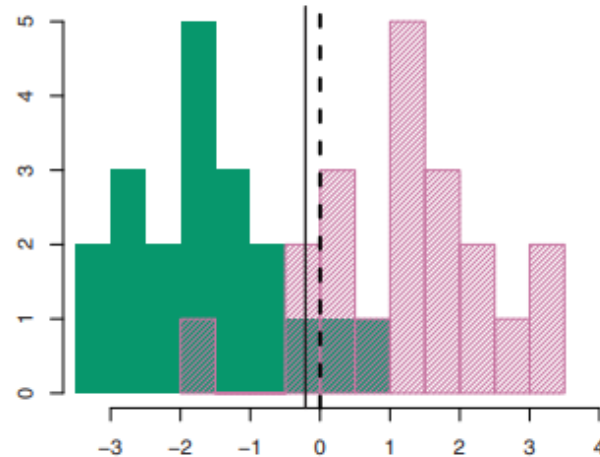
$$\delta_k(x) = x \cdot \frac{\mu_k}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2} + \log(\pi_k)$$

Example: $k = 2, \pi_1 = \pi_2$

Classify as class 1 if $x > \frac{\mu_1 + \mu_2}{2}$



True decision boundary



Estimated decision boundary

LDA in practice

Given training data $(x_i, y_i), i = 1, \dots, n, y_i \in \{1, \dots, K\}$

1. Estimate mean and variance

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{i:y_i=k} x_i$$
$$\hat{\sigma}^2 = \frac{1}{n - K} \sum_{k=1}^K \sum_{i:y_i=k} (x_i - \hat{\mu}_k)^2$$

2. Estimate prior

$$\hat{\pi}_k = n_k / n.$$

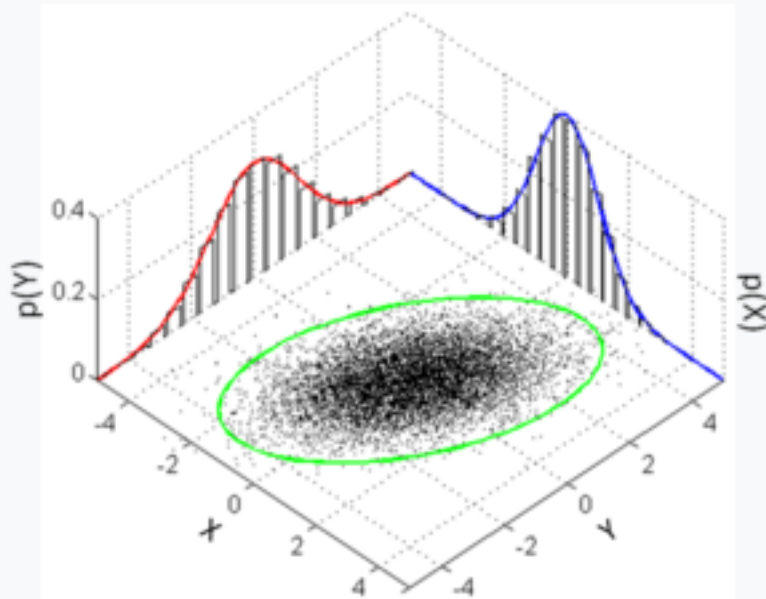
Given testing point x , predict k that maximizes:

$$\hat{\delta}_k(x) = x \cdot \frac{\hat{\mu}_k}{\hat{\sigma}^2} - \frac{\hat{\mu}_k^2}{2\hat{\sigma}^2} + \log(\hat{\pi}_k)$$

Multi-Variate Normal

Multivariate normal

Probability density function



Many sample points from a multivariate normal distribution with $\boldsymbol{\mu} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ and $\boldsymbol{\Sigma} = \begin{bmatrix} 1 & 3/5 \\ 3/5 & 2 \end{bmatrix}$, shown along with the 3-sigma ellipse, the two marginal distributions, and the two 1-d histograms.

$$\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}),$$

with k -dimensional mean vector

$$\boldsymbol{\mu} = \mathbb{E}[\mathbf{X}] = [\mathbb{E}[X_1], \mathbb{E}[X_2], \dots, \mathbb{E}[X_k]]^T,$$

and $k \times k$ covariance matrix

$$\Sigma_{i,j} =: \mathbb{E}[(X_i - \mu_i)(X_j - \mu_j)] = \text{Cov}[X_i, X_j]$$

$$\boldsymbol{\Sigma} =: \mathbb{E}[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T] = [\text{Cov}[X_i, X_j]; 1 \leq i, j \leq k].$$

$$f_k(\mathbf{x}) = \frac{1}{(2\pi)^{p/2} |\boldsymbol{\Sigma}_k|^{1/2}} e^{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)}.$$

Multi-variate LDA

$$\Pr(Y = k|X = x) = \frac{\pi_k f_k(x)}{\sum_{l=1}^K \pi_l f_l(x)}.$$

Assume $\Sigma_k = \Sigma$

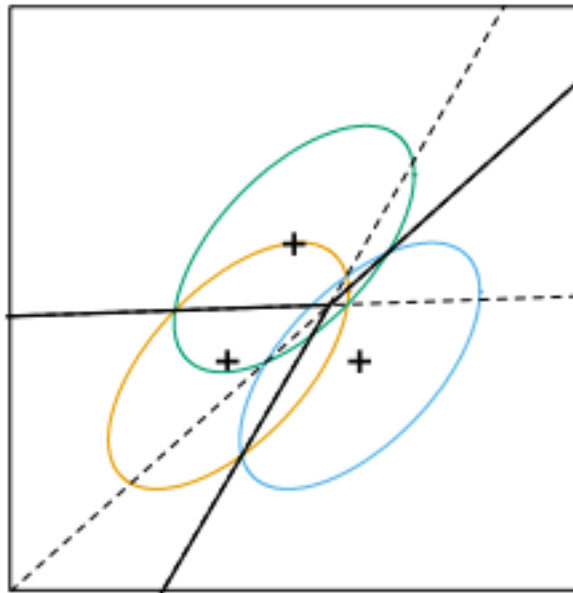
$$\begin{aligned} \log \frac{\Pr(Y = k|X = x)}{\Pr(Y = l|X = x)} &= \log \frac{f_k(x)}{f_l(x)} + \log \frac{\pi_k}{\pi_l} \\ &= \log \frac{\pi_k}{\pi_l} - \frac{1}{2}(\mu_k + \mu_l)^T \Sigma^{-1}(\mu_k - \mu_l) \\ &\quad + x^T \Sigma^{-1}(\mu_k - \mu_l), \end{aligned}$$

Linear decision boundary between classes k and l

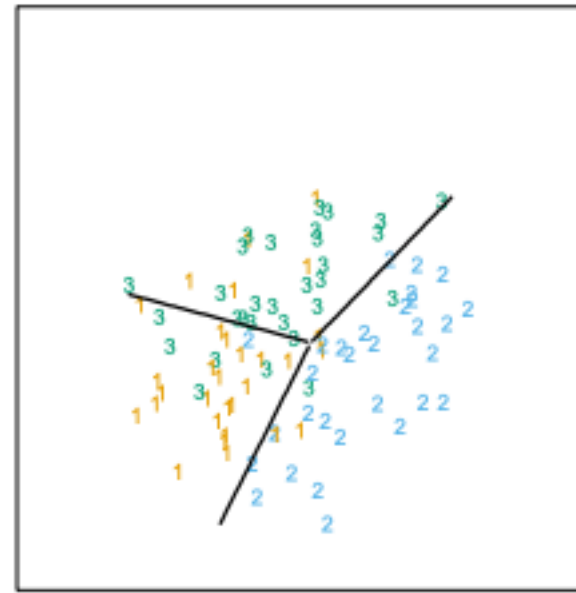
Linear discriminant functions $\delta_k(x) = x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k$

Given x , classify to class k : $\operatorname{argmax}_k \delta_k(x)$

Example 3 classes



3 Normal distributions
with same co-variance,
but different means



LDA decision boundary

Multi-variate LDA

Given training data $(x_i, y_i), i = 1, \dots, n, y_i \in \{1, \dots, K\}$

1. Estimate mean and variance

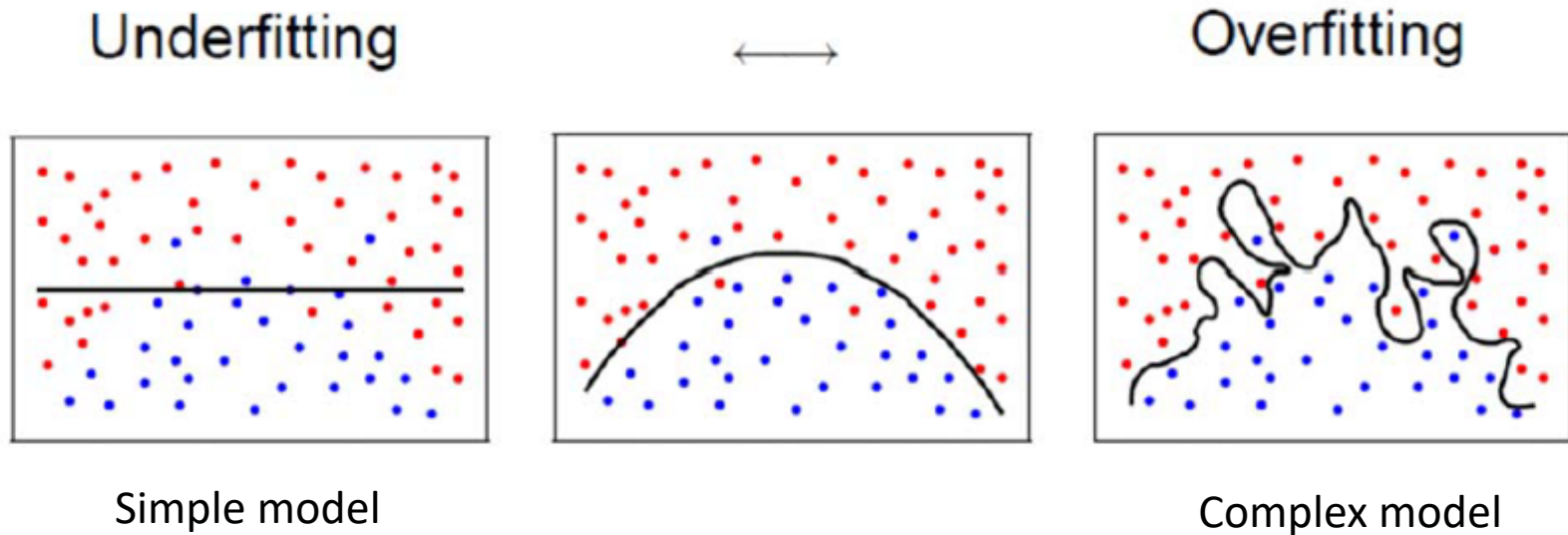
- $\hat{\pi}_k = N_k/N$, where N_k is the number of class- k observations;
- $\hat{\mu}_k = \sum_{g_i=k} x_i / N_k$;
- $\hat{\Sigma} = \sum_{k=1}^K \sum_{g_i=k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T / (N - K)$.

2. Estimate prior

Given testing point x , predict k that maximizes:

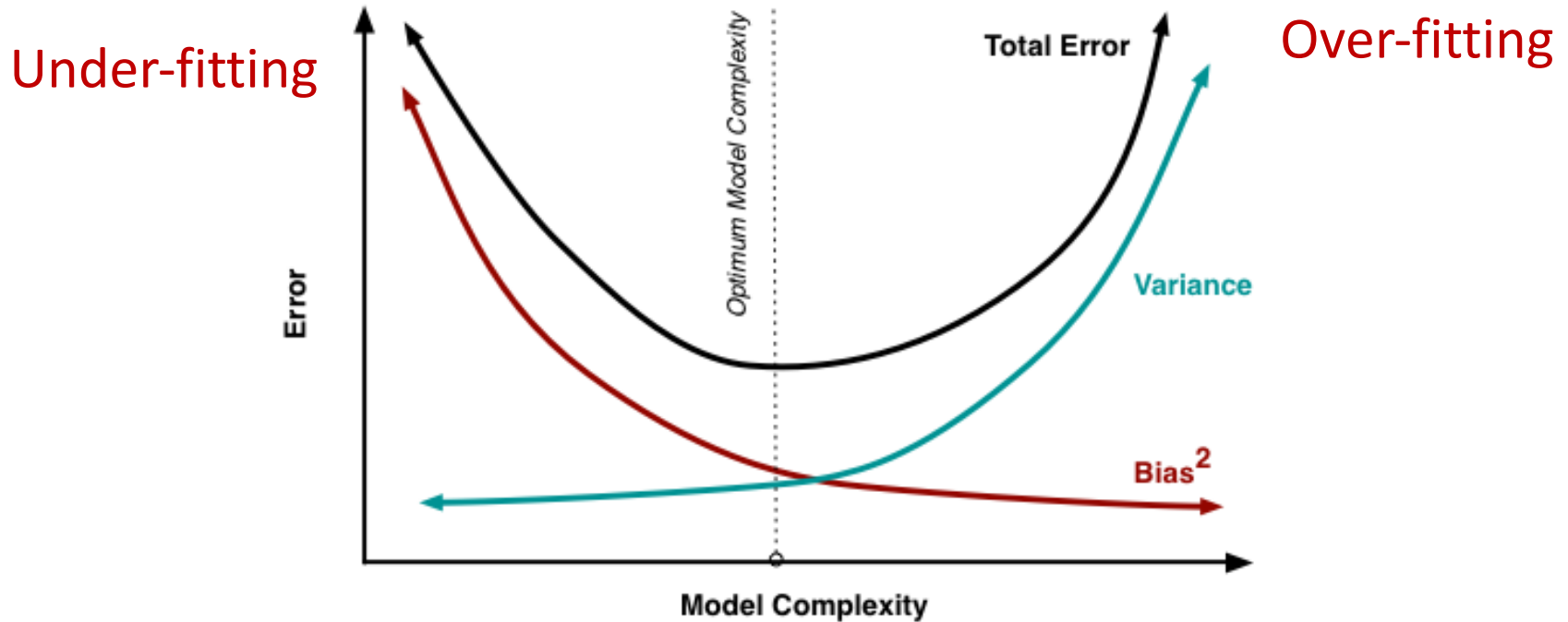
$$\delta_k(x) = x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k$$

Generalization in ML



- Goal is to generalize well on new testing data
- Risk of overfitting to training data
 - MSE close to 0, but performs poorly on test data

Bias-Variance Tradeoff



- Bias = Difference between estimated and true models
 - Variance = Model difference on different training sets
- MSE is proportional to Bias + Variance**

Bias-Variance Decomposition

- Assume that $y = f(\mathbf{x}) + \epsilon$
 - Noise ϵ is sampled from a normal distribution with 0 mean and variance σ^2 : $\epsilon \sim N(0, \sigma^2)$
 - Noise lower-bounds the performance we can achieve
- Recall the following objective function:

$$J(\theta) = \sum_{i=1}^n (h_{\theta}(x_i) - y_i)^2$$

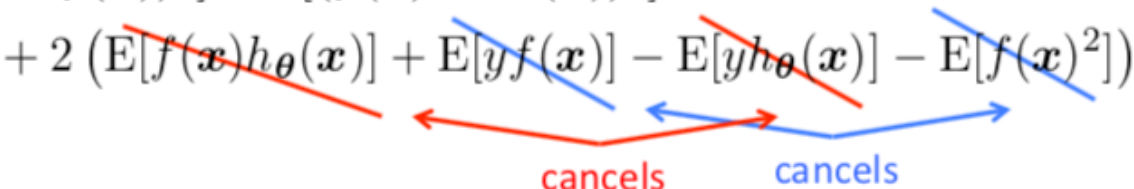
- We can re-write this as the expected value of the squared error: $E(y - h_{\theta}(\mathbf{x}))^2$

$f(\mathbf{x})$: True model (not necessarily linear)

$h_{\theta}(\mathbf{x})$: Learned model (linear)

Bias-Variance Decomposition

$$\begin{aligned} E[(y - h_{\theta}(\mathbf{x}))^2] &= E[(y - f(\mathbf{x}) + f(\mathbf{x}) - h_{\theta}(\mathbf{x}))^2] \\ &= E[(y - f(\mathbf{x}))^2] + E[(f(\mathbf{x}) - h_{\theta}(\mathbf{x}))^2] \\ &\quad + 2 E[(f(\mathbf{x}) - h_{\theta}(\mathbf{x}))(y - f(\mathbf{x}))] \\ &= E[(y - f(\mathbf{x}))^2] + E[(f(\mathbf{x}) - h_{\theta}(\mathbf{x}))^2] \\ &\quad + 2 (E[f(\mathbf{x})h_{\theta}(\mathbf{x})] + E[yf(\mathbf{x})] - E[yh_{\theta}(\mathbf{x})] - E[f(\mathbf{x})^2]) \end{aligned}$$




Therefore,

$$\begin{aligned} E[(y - h_{\theta}(\mathbf{x}))^2] &= E[(y - f(\mathbf{x}))^2] + E[(f(\mathbf{x}) - h_{\theta}(\mathbf{x}))^2] \\ &= E[\epsilon^2] + E[(f(\mathbf{x}) - h_{\theta}(\mathbf{x}))^2] \end{aligned}$$

Aside:

Definition of Variance

$$\text{var}(z) = E[(z - E[z])^2]$$

 This is actually $\text{var}(\epsilon)$, since mean is 0

Bias-Variance Decomposition

$$\begin{aligned}
 E[(y - h_{\theta}(\mathbf{x}))^2] &= \text{var}(\epsilon) + E[(f(\mathbf{x}) - h_{\theta}(\mathbf{x}))^2] \\
 &= \text{var}(\epsilon) + E[(f(\mathbf{x}) - E[h_{\theta}(\mathbf{x})] + E[h_{\theta}(\mathbf{x})] - h_{\theta}(\mathbf{x}))^2] \\
 &= \text{var}(\epsilon) + E[(f(\mathbf{x}) - E[h_{\theta}(\mathbf{x})])^2] + E[(E[h_{\theta}(\mathbf{x})] - h_{\theta}(\mathbf{x}))^2] \\
 &\quad + 2E[(E[h_{\theta}(\mathbf{x})] - h_{\theta}(\mathbf{x}))(f(\mathbf{x}) - E[h_{\theta}(\mathbf{x})])] \\
 &= \text{var}(\epsilon) + E[(f(\mathbf{x}) - E[h_{\theta}(\mathbf{x})])^2] + E[(E[h_{\theta}(\mathbf{x})] - h_{\theta}(\mathbf{x}))^2] \\
 &\quad + 2(E[f(\mathbf{x})E[h_{\theta}(\mathbf{x})]] - E[E[h_{\theta}(\mathbf{x})]^2] - E[f(\mathbf{x})h_{\theta}(\mathbf{x})] + E[h_{\theta}(\mathbf{x})E[h_{\theta}(\mathbf{x})]])
 \end{aligned}$$

Therefore,

$$E[(y - h_{\theta}(\mathbf{x}))^2] = \underbrace{\text{var}(\epsilon)}_{\text{noise}} + \underbrace{E[(f(\mathbf{x}) - E[h_{\theta}(\mathbf{x})])^2]}_{\text{bias}} + \underbrace{E[(E[h_{\theta}(\mathbf{x})] - h_{\theta}(\mathbf{x}))^2]}_{\text{variance}}$$

$$E[(y - h_{\theta}(\mathbf{x}))^2] = \text{bias}(h_{\theta}(\mathbf{x}))^2 + \text{var}(h_{\theta}(\mathbf{x})) + \sigma^2$$

Review

- LDA
 - Example of generative model
 - Use Bayes Theorem to estimate the probability that label is from each class k
 - Assumes normal distribution of features in each class
- Bias-Variance tradeoff for linear regression
 - Decompose expectation of testing error as a sum of squared bias, variance, and noise term
 - Shows that bias and variance need to be simultaneously minimized

Acknowledgements

- Slides made using resources from:
 - Andrew Ng
 - Eric Eaton
 - David Sontag
- Thanks!