

CSCI567 Machine Learning (Spring 2021)

Sirisha Rambhatla

University of Southern California

Feb 10, 2021

1 / 19

Logistics

Outline

- 1 Logistics
- 2 Review of last lecture
- 3 Linear Discriminant Analysis and Quadratic Discriminant Analysis
- 4 Relationship between Logistic Regression and LDA

3 / 19

Outline

- 1 Logistics
- 2 Review of last lecture
- 3 Linear Discriminant Analysis and Quadratic Discriminant Analysis
- 4 Relationship between Logistic Regression and LDA

2 / 19

Logistics

Logistics

- HW 2 was assigned. Solutions for HW 1 will be delayed, stay tuned!
- Please form the groups by Friday, let us know if cannot find a group.

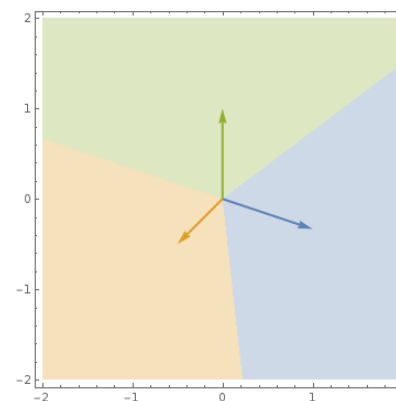
4 / 19

Outline

- 1 Logistics
- 2 Review of last lecture
- 3 Linear Discriminant Analysis and Quadratic Discriminant Analysis
- 4 Relationship between Logistic Regression and LDA

5 / 19

Linear models: from binary to multiclass



$$\begin{aligned} \mathbf{w}_1 &= (1, -\frac{1}{3}) \\ \mathbf{w}_2 &= (-\frac{1}{2}, -\frac{1}{2}) \\ \mathbf{w}_3 &= (0, 1) \end{aligned}$$

- Blue class:
 $\{\mathbf{x} : 1 = \operatorname{argmax}_k \mathbf{w}_k^T \mathbf{x}\}$
- Orange class:
 $\{\mathbf{x} : 2 = \operatorname{argmax}_k \mathbf{w}_k^T \mathbf{x}\}$
- Green class:
 $\{\mathbf{x} : 3 = \operatorname{argmax}_k \mathbf{w}_k^T \mathbf{x}\}$

$$\mathcal{F} = \left\{ f(\mathbf{x}) = \operatorname{argmax}_{k \in [C]} \mathbf{w}_k^T \mathbf{x} \mid \mathbf{w}_1, \dots, \mathbf{w}_C \in \mathbb{R}^D \right\}$$

6 / 19

Softmax + MLE = minimizing cross-entropy loss

Maximize probability of see labels y_1, \dots, y_N given $\mathbf{x}_1, \dots, \mathbf{x}_N$

$$P(\mathbf{W}) = \prod_{n=1}^N \mathbb{P}(y_n \mid \mathbf{x}_n; \mathbf{W}) = \prod_{n=1}^N \frac{e^{\mathbf{w}_{y_n}^T \mathbf{x}_n}}{\sum_{k \in [C]} e^{\mathbf{w}_k^T \mathbf{x}_n}}$$

By taking **negative log**, this is equivalent to minimizing

$$F(\mathbf{W}) = \sum_{n=1}^N \ln \left(\frac{\sum_{k \in [C]} e^{\mathbf{w}_k^T \mathbf{x}_n}}{e^{\mathbf{w}_{y_n}^T \mathbf{x}_n}} \right) = \sum_{n=1}^N \ln \left(1 + \sum_{k \neq y_n} e^{(\mathbf{w}_k - \mathbf{w}_{y_n})^T \mathbf{x}_n} \right)$$

This is the **multiclass logistic loss**, a.k.a **cross-entropy loss**.

7 / 19

Comparisons of multiclass-to-binary reductions

In big O notation,

Reduction	#training points	test time	Idea
OvA	CN	C	is class k or not?
OvO	CN	C^2	is class k or class k' ?
ECOC	LN	L	is bit b on or off?
Tree	$(\log_2 C)N$	$\log_2 C$	belong to which half of the label set?

8 / 19

Outline

- 1 Logistics
- 2 Review of last lecture
- 3 Linear Discriminant Analysis and Quadratic Discriminant Analysis
- 4 Relationship between Logistic Regression and LDA

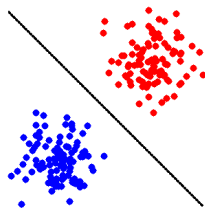
9 / 19

What do we know?

Ok, so we know that by **Bayes theorem** for a C class classification task

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

Let's consider a Binary Classification task, $C = \{0, 1\}$. *What is the decision boundary?*



$$\mathcal{P}(y = 0|X = \mathbf{x}) = \mathcal{P}(y = 1|X = \mathbf{x})$$

11 / 19

Revisiting Bayes optimal classifier

Tells us what to predict for \mathbf{x} , *knowing* $\mathcal{P}(y|\mathbf{x})$

Bayes optimal classifier: $f^*(\mathbf{x}) = \operatorname{argmax}_{c \in [C]} \mathcal{P}(c|\mathbf{x})$.

But the main issue was that in practice we don't know what $\mathcal{P}(y|\mathbf{x})$ is!

10 / 19

The main bottleneck is *not knowing* $\mathcal{P}(X = \mathbf{x}|y = c)$

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

LDA makes **two simplifying assumptions**:

- Let $\mathcal{P}(X = \mathbf{x}|y = c) \sim \mathcal{N}(\boldsymbol{\mu}_c, \Sigma_c)$, and
- Let all class covariances be the same i.e. $\Sigma_c = \Sigma$ for all $c \in [C]$

If so, the *decision boundary* (for binary classification) is given by

$$\begin{aligned} \mathcal{P}(y = 0|X = \mathbf{x}) &= \mathcal{P}(y = 1|X = \mathbf{x}) \\ \frac{\mathcal{P}(X = \mathbf{x}|y = 0)\mathcal{P}(y = 0)}{\mathcal{P}(X = \mathbf{x})} &= \frac{\mathcal{P}(X = \mathbf{x}|y = 1)\mathcal{P}(y = 1)}{\mathcal{P}(X = \mathbf{x})} \end{aligned}$$

12 / 19

The main bottleneck is *not knowing* $\mathcal{P}(X = \mathbf{x}|y = c)$

$$\mathcal{P}(y = c|X = \mathbf{x}) = \frac{\mathcal{P}(X = \mathbf{x}|y = c)\mathcal{P}(y = c)}{\mathcal{P}(X = \mathbf{x})}.$$

LDA makes **two simplifying assumptions**:

- Let $\mathcal{P}(X = \mathbf{x}|y = c) \sim \mathcal{N}(\boldsymbol{\mu}_c, \Sigma_c)$, and
- Let all class covariances be the same i.e. $\Sigma_c = \Sigma$ for all $c \in [C]$

If so, the *decision boundary* (for binary classification) is given by

$$\mathcal{P}(y = 0|X = \mathbf{x}) = \mathcal{P}(y = 1|X = \mathbf{x})$$

$$\mathcal{P}(X = \mathbf{x}|y = 0)\mathcal{P}(y = 0) = \mathcal{P}(X = \mathbf{x}|y = 1)\mathcal{P}(y = 1)$$

Now, $\mathcal{P}(X = \mathbf{x}|y = 0) \sim \mathcal{N}(\boldsymbol{\mu}_0, \Sigma_0)$ and $\mathcal{P}(X = \mathbf{x}|y = 1) \sim \mathcal{N}(\boldsymbol{\mu}_1, \Sigma_1)$.

For $\boldsymbol{\mu}_c \in \mathbb{R}^d$ and $\Sigma_c^{-1} \in \mathbb{R}^{d \times d}$, we have

$$\mathcal{P}(X = \mathbf{x}|y = c) = \frac{1}{(2\pi)^{d/2}|\Sigma_c|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_c)^\top \Sigma_c^{-1}(\mathbf{x} - \boldsymbol{\mu}_c)\right)$$

The *decision boundary* is given by

$$\mathcal{P}(X = \mathbf{x}|y = 0)\mathcal{P}(y = 0) = \mathcal{P}(X = \mathbf{x}|y = 1)\mathcal{P}(y = 1).$$

Substituting for $\mathcal{P}(X = \mathbf{x}|y = c)$, and taking $\log(\cdot)$ and simplifying¹

$$\log\left(\frac{\mathcal{P}(y = 0)}{\mathcal{P}(y = 1)}\right) - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_0)^\top \Sigma_0^{-1}(\mathbf{x} - \boldsymbol{\mu}_0) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^\top \Sigma_1^{-1}(\mathbf{x} - \boldsymbol{\mu}_1)$$

For simplicity of this exposition we have ignored the $1/(2\pi)^{d/2}|\Sigma_c|^{1/2}$ terms since these will just add to the constants.

Let's simplify this

$$\log \frac{\mathcal{P}(y=0)}{\mathcal{P}(y=1)} - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_0)^\top \Sigma_0^{-1}(\mathbf{x} - \boldsymbol{\mu}_0) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^\top \Sigma_1^{-1}(\mathbf{x} - \boldsymbol{\mu}_1)$$

$$\begin{aligned} \log \frac{\mathcal{P}(y = 0)}{\mathcal{P}(y = 1)} - \frac{1}{2}\mathbf{x}^\top \Sigma_0^{-1}\mathbf{x} + \boldsymbol{\mu}_0^\top \Sigma_0^{-1}\mathbf{x} - \frac{1}{2}\boldsymbol{\mu}_0^\top \Sigma_0^{-1}\boldsymbol{\mu}_0 = \\ -\frac{1}{2}\mathbf{x}^\top \Sigma_1^{-1}\mathbf{x} + \boldsymbol{\mu}_1^\top \Sigma_1^{-1}\mathbf{x} - \frac{1}{2}\boldsymbol{\mu}_1^\top \Sigma_1^{-1}\boldsymbol{\mu}_1 \end{aligned}$$

Setting $\Sigma_0 = \Sigma_1 = \Sigma$,

$$\mathbf{w}^\top \mathbf{x} + w_0 = 0$$

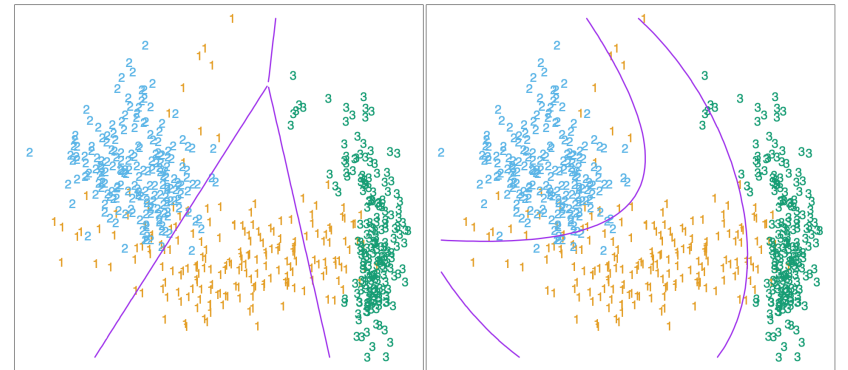
Here,

$$\mathbf{w}^\top = \boldsymbol{\mu}_0^\top \Sigma^{-1} - \boldsymbol{\mu}_1^\top \Sigma^{-1}$$

and

$$w_0 = \log \frac{\mathcal{P}(y = 0)}{\mathcal{P}(y = 1)} - \frac{1}{2}\boldsymbol{\mu}_0^\top \Sigma^{-1}\boldsymbol{\mu}_0 + \frac{1}{2}\boldsymbol{\mu}_1^\top \Sigma^{-1}\boldsymbol{\mu}_1$$

What do the decision boundaries look like?



The decision boundaries are a quadratic when Σ 's are not the same, this is known as *Quadratic Discriminant Analysis*!

Outline

- 1 Logistics
- 2 Review of last lecture
- 3 Linear Discriminant Analysis and Quadratic Discriminant Analysis
- 4 Relationship between Logistic Regression and LDA

17 / 19

For **Logistic Regression** we assumed:

$$\mathcal{P}(y = 1|X = \mathbf{x}) = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}}$$

This can be written as follows, in terms of *log-odds*

$$\log \frac{\mathcal{P}(y = 1|X = \mathbf{x})}{\mathcal{P}(y = 0|X = \mathbf{x})} = \mathbf{w}^\top \mathbf{x}$$

The log-odds can be modeled as a linear function of \mathbf{x} .

18 / 19

For **LDA** we can write the *log-odds* as

$$\begin{aligned} \log \frac{\mathcal{P}(y = 1|X = \mathbf{x})}{\mathcal{P}(y = 0|X = \mathbf{x})} &= \log \frac{\mathcal{P}(X = \mathbf{x}|y = 1)\mathcal{P}(y = 1)}{\mathcal{P}(X = \mathbf{x}|y = 0)\mathcal{P}(y = 0)} \\ &= \log \frac{\mathcal{P}(y = 1)}{\mathcal{P}(y = 0)} + (\boldsymbol{\mu}_1^\top \Sigma^{-1} - \boldsymbol{\mu}_0^\top \Sigma^{-1})\mathbf{x} - \frac{1}{2}\boldsymbol{\mu}_1^\top \Sigma^{-1}\boldsymbol{\mu}_1 + \frac{1}{2}\boldsymbol{\mu}_0^\top \Sigma^{-1}\boldsymbol{\mu}_0 \\ &= \mathbf{w}^\top \mathbf{x} \end{aligned}$$

LDA satisfies the assumptions of the Logistics Regression model!

LDA imposes additional assumptions on the data, i.e., it **assumes that the class conditional densities are Gaussian**.

19 / 19