

```
31 def __init__(self, *args, **kwargs):
32     self.file = None
33     self.fingerprints = {}
34     self.logduplicates = True
35     self.debug = debug
36     self.logger = logging.getLogger(__name__)
37     if path:
38         self.file = open(os.path.join(path, 'fingerprint'), 'w')
39         self.file.seek(0)
40         self.fingerprints.update(json.load(self.file))
41
42     @classmethod
43     def from_settings(cls, settings):
44         debug = settings.getboolean('debug')
45         return cls(job_dir=settings['job_dir'],
46                    request_size=settings['request_size'],
47                    fp_size=settings['fp_size'],
48                    fp=fp,
49                    logduplicates=settings['logduplicates'],
50                    loglevel=settings['loglevel'],
51                    max_fingerprints=settings['max_fingerprints'])
```

GAUSSIAN DISCRIMINANT ANALYSIS (GDA)

Nancy, Ornella, Avotra, Allassan, Fenosoa

April 22, 2022

Overview

1. Introduction
2. Gaussian Discriminant Analysis
3. Sentimental Analysis using GDA
4. Conclusion

Introduction

Introduction

The main aim of this project is to implement Gaussian Discriminative Analysis. Machine Learning models can be classified as discriminative or generative models.

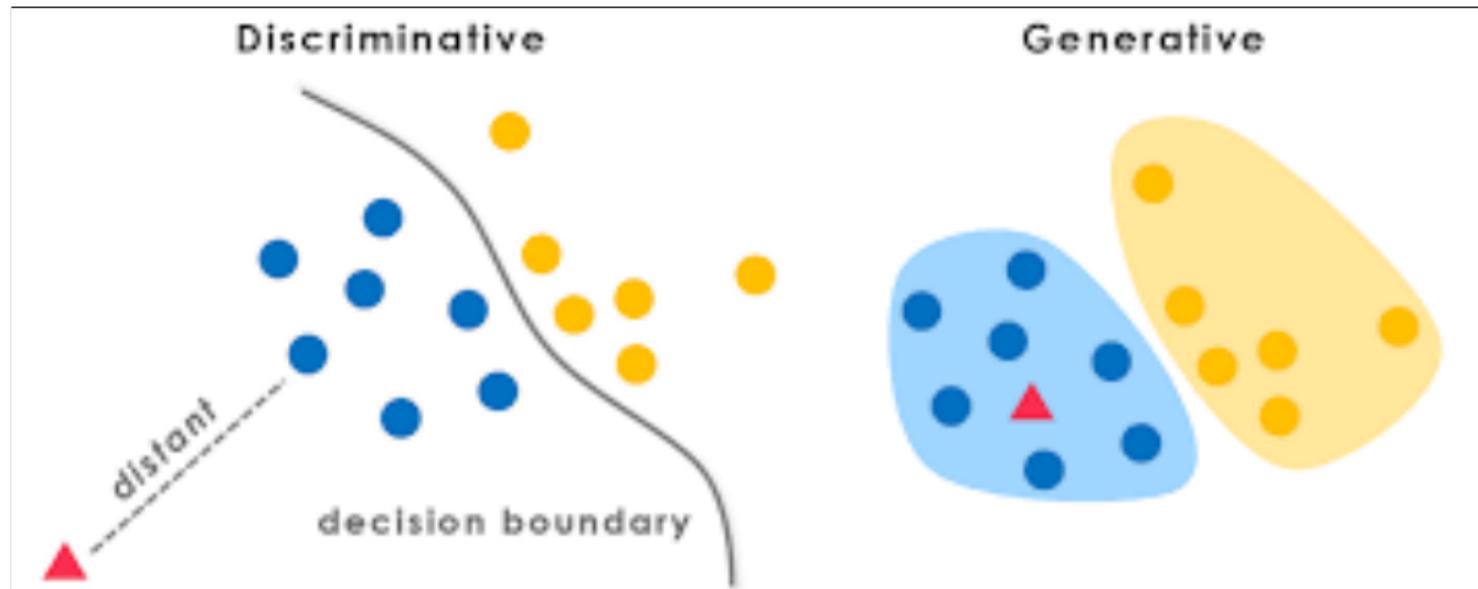
- ▶ Discriminative Models

We model probability distribution $\mathbb{P}(\mathbf{y}|\mathbf{x})$ with the training data.

- ▶ Generative Models

We This is a part of unsupervised learning task in

Machine Learning Models



Discriminative models Vs Generative models

Gaussian Discriminant Analysis

Gaussian Discriminant Analysis

GDA is a generative learning algorithm.

- ▶ Input features are continuous random variables
- ▶ Assume $\mathbb{P}(\mathbf{x}|\mathbf{y})$ and $\mathbb{P}(\mathbf{y})$ are distributed according to multivariate normal distribution and Bernoulli distribution respectively.

Gaussian Discriminant Analysis

Assumptions

$$\mathbb{P}[\mathbf{y}_i; \phi] = \phi^{\mathbf{y}_i} (1 - \phi)^{1 - \mathbf{y}_i}, \quad \mathbf{y}_i \in \{0, 1\}. \quad (1)$$

$$\mathbb{P}[\mathbf{x}_i | \mathbf{y}_i = 0; \boldsymbol{\mu}_0, \boldsymbol{\Sigma}] = \frac{1}{(2\pi)^{n/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left(-\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_0) \right) \quad (2)$$

$$\mathbb{P}[\mathbf{x}_i | \mathbf{y}_i = 1; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}] = \frac{1}{(2\pi)^{n/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left(-\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_1) \right) \quad (3)$$

Gaussian Discriminant Analysis

From Bayes rule , we can compute $\mathbb{P}(\mathbf{y}|\mathbf{x})$.

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

$$\operatorname{argmax}_{\mathbf{y}} \mathbb{P}(\mathbf{y}|\mathbf{x}) = \operatorname{argmax}_{\mathbf{y}} \frac{\mathbb{P}(\mathbf{x}|\mathbf{y}) \times \mathbb{P}(\mathbf{y})}{\mathbb{P}(\mathbf{x})}$$

$$\operatorname{argmax}_{\mathbf{y}} \mathbb{P}(\mathbf{y}|\mathbf{x}) \propto \operatorname{argmax}_{\mathbf{y}} \mathbb{P}(\mathbf{x}|\mathbf{y}) \times \mathbb{P}(\mathbf{y})$$

The likelihood function L is given by

$$L(\mathbf{y}|\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}, \phi) = \mathbb{P}(\mathbf{x}|\mathbf{y}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}) \mathbb{P}(\mathbf{y}; \phi) \text{ where } k = \mathbb{1}_{\{\mathbf{y}_i=1\}}$$

$$\phi = \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{\mathbf{y}_i=1\}}$$

$$\boldsymbol{\mu_0} = \frac{\sum_{i=1}^n \mathbb{1}_{\{\mathbf{y}_i=0\}} \mathbf{x}_i}{\mathbb{1}_{\{\mathbf{y}_i=0\}}}$$

$$\boldsymbol{\mu_1} = \frac{\sum_{i=1}^n \mathbb{1}_{\{\mathbf{y}_i=1\}} \mathbf{x}_i}{\mathbb{1}_{\{\mathbf{y}_i=1\}}}$$

$$\Sigma = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu_k}) (\mathbf{x}_i - \boldsymbol{\mu_k})^T \text{ where } k = \mathbb{1}_{\{\mathbf{y}_i=1\}}$$

Sentimental Analysis using GDA

Sentimental Analysis using GDA

Preprocessing

Steps

- ▶ lower case
- ▶ remove stopwords
- ▶ remove punctuations
- ▶ remove apostrophe
- ▶ tokenize
- ▶ lemmatization
- ▶ convert numbers

TF -IDF

- ▶ Term- Frequency

$$TF = \frac{\text{Nbr of repetitions of word in a doc}}{\text{\# of words in a doc}}$$

- ▶ Inverse Document Frequency

$$IDF = \log \left[\frac{\text{\# Number of docs}}{\text{Nbr of docs containing the word}} \right]$$

Sentimental Analysis using GDA

The estimated parameters from the training dataset

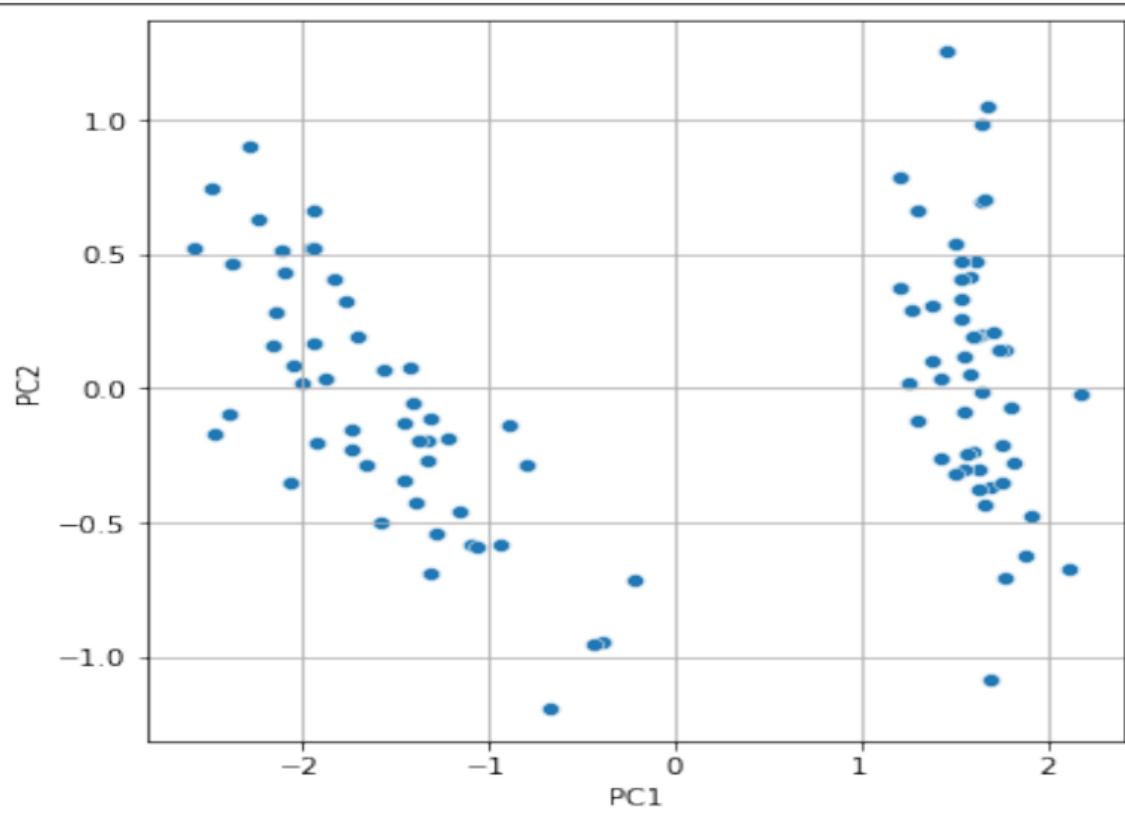
$$\phi^* = 0.75$$

$$\mu_0^* = (-0.0262559, -0.01121153)^T$$

$$\mu_1^* = (0.00595905, 0.0021744)^T$$

$$\Sigma^* = \begin{pmatrix} 4.62948277 \cdot 10^{-4} & -9.65664189 \cdot 10^{-6} \\ -9.65664189 \cdot 10^{-6} & 1.49865708 \cdot 10^{-3} \end{pmatrix}$$

The accuracy after prediction is 50%.



Visualization of the Iris Dataset using PCA

Conclusion

In conclusion,

- ▶ GDA model does not need any hyperparameter.
- ▶ GDA models works perfectly with small dataset.
- ▶ Time complexity is expensive



THANK
YOU