

CSE 575: Statistical Machine Learning (Spring 2021)

Instructor: Nupur Thakur

Supervised Learning



Table of contents

1. Project Part - 1
2. Density Estimation using MLE
3. Generative v/s Discriminative Models

Project Part 1 - Density Estimation and Supervised Learning

- Due on: 17 Feb, 11:59 PM MST.
- Binary classification using Fashion-MNIST dataset.
- Two classes - T-shirt (label 0) and Trouser (label 1)
- 28x28 grayscale images
- Training set: "Tshirt": 6000; "Trouser": 6000.
- Testing set: "Tshirt": 6000; "Trouser": 1000.



Project Part 1 - Tasks

- Calculate two features of each image - average and standard deviation.
- Estimate parameters for 2D Gaussian distribution for each class.
- Implement Naive Bayes to perform classification
- Train Logistic regression using gradient ascent to perform classification

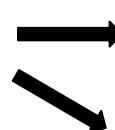
Project Part 1 - Deliverables

- Code:
 - Acceptable file types are .py/.m or .zip.
 - Well-commented code. Be sure to read through the directions carefully to ensure you have included all necessary parts in your code.
- Report:
 - Acceptable File types: .pdf
 - Length: 2-5 A4 pages
 - Content: Refer Canvas

Project Part 1 - Languages/Software

- What can you use?
 - Python/MATLAB to code.
 - Libraries or in-built functions to manipulate the data i.e. NumPy, SciPy libraries etc.
- What cannot be used?
 - Packages like scikit-learn which have ready-to-use algorithms.

What is supervised learning?

- Data - <sample, label> pairs
- Objective - Learn from the given data such that label can be predicted for a new data sample.
- Based on labels -  Regression
Classification

Density Estimation

- Estimating underlying probability density function, based on the training data
- Parametric: each class of images (the feature vectors) may be modeled by a density function $p_\theta(x)$ with parameter θ .
- Non-parametric: makes no assumption about the distribution of data or modeling the data without any parameters

Maximum Likelihood Estimation

- Given some training data and assuming a parametric model $p(x|\theta)$; what specific θ will fit/explain the data best?
- To consider all the samples denoted by $D=\{x_1, x_2, \dots, x_n\}$, assume that all the samples are i.i.d - independent and identically distributed.
- So, data likelihood represented by $L(\theta)$ is -

$$L(\theta) = P(D|\theta) = \prod_i P(x_i|\theta)$$

Maximum Likelihood Estimation

- ① Likelihood function
- ② take derivative
- ③ Set to 0.

- Maximum Likelihood Estimation (MLE): Finding the parameter that maximizes the likelihood function

$$\hat{\theta} = \operatorname{argmax}_{\theta} p(x|\theta) \rightarrow L1$$

$$\begin{aligned}\hat{\theta} &= \operatorname{argmax}_{\theta} \prod_{i=1}^n p(x_i|\theta) \\ &\stackrel{?}{=} \operatorname{argmax}_{\theta} \log \prod_{i=1}^n p(x_i|\theta) = \operatorname{argmax}_{\theta} \sum_{i=1}^n \log p(x_i|\theta)\end{aligned}$$

MLE - Example 1

$$\theta = \{\mu, \sigma^2\}$$

Given n i.i.d. samples $\{x_i\}$ from the 1-D normal distribution $N(\mu, \sigma^2)$, find the MLE for μ and σ^2 .

① Likelihood func

$$P(D|\theta) = \prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}$$

$$P(D|\theta) = \prod_{i=1}^n P(x_i|\theta)$$

$$\log P(D|\theta) = \underbrace{\log \prod_{i=1}^n}_{\text{log } P(x_i|\theta)} + \underbrace{\sum_{i=1}^n \log \left(\frac{1}{\sigma \sqrt{2\pi}} \right)}_{-\frac{n}{2} \log \sigma^2} - \frac{n}{2} \sum_{i=1}^n \frac{(x_i-\mu)^2}{\sigma^2}$$

$$(\mu, \sigma^2)$$

MLE - Example 1

② Take derivative w.r.t μ

$$\begin{aligned}\frac{\partial}{\partial \mu} \log P(D|\theta) &= 0 - \sum_{i=1}^n \frac{\partial}{\partial \mu} \frac{f(x_i - \mu)}{\sigma^2} \\ &= \sum_{i=1}^n \frac{(x_i - \mu)}{\sigma^2}\end{aligned}$$

③ Equate derivative to 0.

$$\begin{aligned}\sum_{i=1}^n (x_i - \mu) &= 0 \\ \sum_{i=1}^n x_i - n\mu &= 0 \\ n\bar{x} - n\mu &= 0 \\ \bar{x} - \mu &= 0\end{aligned}$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} \Rightarrow \text{Sample mean}$$

$$\sum_{i=1}^n x_i - \frac{\sum_{i=1}^n x_i}{n} = 0$$

Generative vs Discriminative Models

- Generative models -
 - Learn $P(y)$ and $P(x|y)$.
 - Ex: Bayesian classifier, Naive Bayes.
- Discriminative models -
 - Directly learn $P(y|x)$
 - Ex: Logistic Regression

Questions?