

CSE 575: Statistical Machine Learning (Spring 2021)

Instructor: Nupur Thakur

Supervised Learning



Table of contents

- 1. MLE solution for variance**
- 2. Naive Bayes**
- 3. Logistic Regression**

MLE - Example 1 (Continued)

$$\sigma \rightarrow \sum \begin{cases} \sigma^2 \\ \sigma^{-2} \end{cases}$$

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 .

$$-\frac{n}{2} \log(2\pi\sigma^2) =$$

$$\log P(D|\mu, \sigma) = -n \log(\sigma\sqrt{2\pi}) - \sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2}$$

$$\frac{\partial}{\partial \sigma} \log P(D|\mu, \sigma) = -\frac{n}{\sigma} + \frac{\sum_{i=1}^n (x_i - \mu)^2}{\sigma^3}$$

Equate it to 0,

$$-\frac{n}{\sigma} + \frac{\sum_{i=1}^n (x_i - \mu)^2}{\sigma^3} = 0$$

MLE - Example 1 (Continued)

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

$$\left\{ \hat{\sigma}_{\text{MLE}}^2 = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n} \right\} \text{Biased}$$

$$\text{Bias} - E(\hat{\sigma}_{\text{MLE}}^2) - \sigma^2 = 0 \quad \text{Unbiased}$$

$$P(X|Y)$$

Naive Bayes

$$P(X_1, X_2, X_3 | Y)$$

2^3 2^3

$$Y=0$$

$$Y=1$$

- The "naive" conditional independence assumption: each feature is (conditionally) independent of every other feature, given the label, i.e.,
 $p(x_i | \{x_j \text{ for any } j \neq i\}, y) = p(x_i | y)$
- How does this simplify?
 - Consider the previous example again: d-dimensional binary features, and y is also binary.

Bayesian \rightarrow No. of probabilities $P(X_1, X_2, X_3 | Y) = 2^{3+1}$

\nexists d-dimensional $X \rightarrow 2^{d+1} \doteq 4d$

$$P(X_1, X_2, X_3 | Y) = P(X_1 | Y) \cdot P(X_2 | Y) \cdot P(X_3 | Y)$$
$$Y=0 \rightarrow P(X_1=0 | Y=0) \cdot P(X_2=0 | Y=0) \cdot P(X_3=0 | Y=0) = 6 \cdot 6 = 36$$
$$Y=1 \rightarrow P(X_1=1 | Y=1) \cdot P(X_2=1 | Y=1) \cdot P(X_3=1 | Y=1) = 6 \cdot 6 = 36$$

Naive Bayes

- The predicted label is given by -

$$\hat{y} = \operatorname{argmax}_y P(y) \prod_{i=1}^d p(x_i|y)$$

$P(y|x) = P(x|y) \cdot P(y)$

$P(x|y)$ is a probability distribution over x for a fixed y .

- Model parameters- x_1, x_2, x_3, y

$$P(x_1=0|y=0)$$

$$P(x_1=0|y=1)$$

$$P(x_2=0|y=0)$$

$$P(x_2=1|y=1)$$

$$P(x_3=0|y=0)$$

$$P(x_3=1|y=1)$$

$$P(x=0)$$

$y = \text{binary}$
 $x_i = \text{two values}$

$$(2d+1)$$

$$(k-1)$$

Naive Bayes - Example

X

Y

Sky	Temp	Humid	Wind	Water	Forecast	EnjoySpt
Sunny	Warm	Normal	Strong	Warm	Same	Yes
Sunny	Warm	High	Strong	Warm	Same	Yes
Rainy	Cold	High	Strong	Warm	Change	No
Sunny	Warm	High	Strong	Cool	Change	Yes

$$P(Y = \text{Yes}) = \frac{3}{4} \quad P(x_1 = \text{Sunny} | Y = \text{Yes}) = \frac{3}{3}$$

$(2d+1) = 13$ independent parameters

$$\underline{\underline{P(X|Y)}}$$

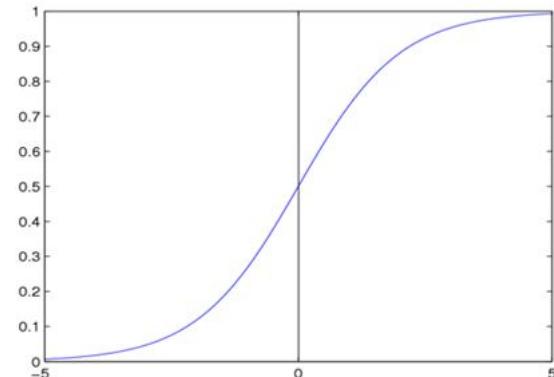
Logistic Regression

$$\underline{\underline{w^T x \geq 0}}$$

- Training set: n labelled samples $\langle \underline{\underline{x(i)}}, \underline{\underline{y(i)}} \rangle$
- Use the logistic function for modeling $\overbrace{P(y|x)}$, considering only the case of $y \in \{0,1\}$

$$\underline{\underline{P(y=0|x) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^d w_i x_i)}}} = 1 - \sigma(\underline{\underline{w^T x}})$$

$$\underline{\underline{P(y=1|x) = \frac{\exp(w_0 + \sum_{i=1}^d w_i x_i)}{1 + \exp(w_0 + \sum_{i=1}^d w_i x_i)}}} = \sigma(\underline{\underline{w^T x}})$$



$$\text{Sigmoid} \quad \underline{\underline{\sigma(t) = \frac{1}{1+e^{-t}} = \frac{e^t}{1+e^t}}}$$

Logistic Regression

$w \rightarrow$ model parameter

- Model parameters -

$$P(Y|X) = \left[1 - \sigma(w^T x) \right]^{1-y} \left[\sigma(w^T x) \right]^y$$
$$P(y^{(1)}, y^{(2)}, \dots, y^{(n)} | x^{(1)}, x^{(2)}, \dots, x^{(n)}) = \prod_{i=1}^n \left[1 - \sigma(w^T x^{(i)}) \right]^{1-y^{(i)}} \left[\sigma(w^T x^{(i)}) \right]^{y^{(i)}}$$

Argmax

$$\hat{L}(w) = \sum_{i=1}^m [1-y^{(i)}] \log [1 - \sigma(w^T x^{(i)})] + y^{(i)} \log [\sigma(w^T x^{(i)})]$$

Gradient Ascent

Update Equation -



$$w_{t+1} = w_t + \alpha \nabla_w L(w)$$

Updated weights Current weight Learning rate Gradient of the loss function

The equation $w_{t+1} = w_t + \alpha \nabla_w L(w)$ is displayed. Above the equation, arrows point from the terms to their definitions: 'Updated weights' points to w_{t+1} , 'Current weight' points to w_t , 'Learning rate' points to α , and 'Gradient of the loss function' points to $\nabla_w L(w)$.

Descent -

Gradient Ascent

Questions?