

Lecture 04 Logistic Regression

Readings: Murphy: 10-10.2.4, 10.3-10.3.3.
Abu-Mostafa: 3.3

Classification

Learn $f: X \rightarrow Y$ $X \in \mathbb{R}^d$ features; $Y = \{1, 2, \dots, k\}$ target classes.

"0-1 Loss Function": $l(f(x), y) = 1_{\{f(x) \neq y\}}$

Expected Loss:

$$\begin{aligned}
 E_{X,Y} [1_{\{f(x) \neq y\}}] & \stackrel{\text{law of total expectation (tower property)}}{=} E_X [E_{Y|X} [1_{\{f(x) \neq y\}} | X=x]] \\
 & \stackrel{\text{Expectation of indicator is probability of event.}}{=} E_X [P(f(x) \neq y | X=x)] \\
 & = 1 - P(f(x) = y | X=x) \\
 & = E_X [1 - P(f(x) = y | X=x)]
 \end{aligned}$$

To minimize expected loss, we must maximize:

$$f^*(x) = \underset{y \in Y}{\operatorname{argmax}} P(Y=y | X=x), \quad \forall x$$

Bayes optimal classifier

Bayes optimal classifier.

(2)

$$f^*(x) = \underset{y}{\operatorname{argmax}} P(Y=y | X=x)$$

In practice, we do not know it.

We only get n iid samples $\{(x_i, y_i)\}_{i=1}^n$

- Suppose X is **discrete** so that $x \in \{1, 2, \dots, m\}$
What is a natural estimator for $P(Y=y | X=x)$?

$$\hat{f}(x) = \underset{y}{\operatorname{argmax}} \frac{\sum_{i=1}^n \mathbb{1}\{x_i=x, y_i=y\}}{\sum_{i=1}^n \mathbb{1}\{x_i=x\}}$$

- If X is **continuous**, we need a **model** to explain observations

Maximum Likelihood Estimation (MLE) for Classification

$$f^*(x) = \underset{y}{\operatorname{argmax}} P(Y=y | X=x)$$

$$\hat{f}(x) = \underset{y}{\operatorname{argmax}} P_{\omega}(Y=y | X=x)$$

General MLE problem:

(1) Parameterize $P_{\omega}(Y=y | X=x)$ as a function of ω

(2) Learn ω on iid training data $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$

$$\begin{aligned} \hat{\omega}_{\text{MLE}} &= \underset{\omega}{\operatorname{argmax}} \prod_{i=1}^n P_{\omega}(y_i | x_i) \\ &= \underset{\omega}{\operatorname{argmax}} \sum_{i=1}^n \log(P_{\omega}(y_i | x_i)) \end{aligned}$$

maximum likelihood of training data!

Modeling conditional probabilities

- Recall linear regression:

- $P_w(Y=y|X=x) = \frac{1}{\sqrt{2\pi}\sigma} \cdot e^{-\frac{(y-\omega^T x)^2}{2\sigma^2}}$
- Prediction: $E[Y|X=x] = \omega^T x$

- in Logistic regression, we use a specialized model for binary classification:

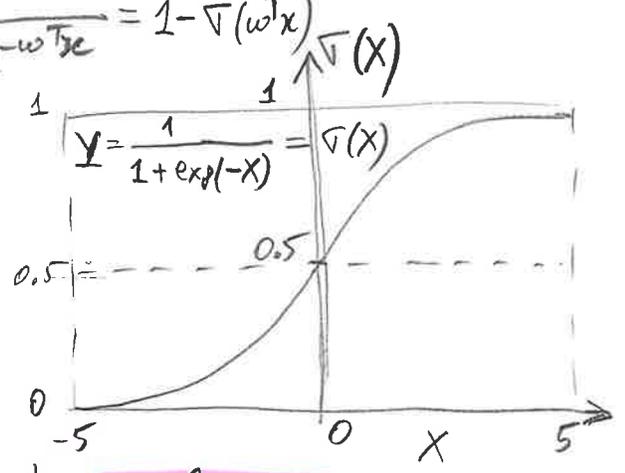
Posterior for class "1"

$$P_w(Y=1|X=x) = \frac{1}{1 + e^{-\omega^T x}} = \sigma(\omega^T x)$$

pass $\omega^T x$ through sigmoid function σ .

$$P_w(Y=0|X=x) = \frac{e^{-\omega^T x}}{1 + e^{-\omega^T x}} = 1 - \sigma(\omega^T x)$$

- Prediction: just round sigmoid($\omega^T x$)



Training data $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ iid, $x_i \in \mathbb{R}^d$, $y_i \in \{-1, +1\}$ for notation convenience!!!

$$\hat{\omega}_{MLE} = \underset{\omega}{\text{argmax}} \sum_{i=1}^n \log P_w(y_i|x_i)$$

we fact that:
 $1 - \sigma(s) = \sigma(-s)$
 $P_w(y_i|x_i) = \sigma(y_i \cdot \omega^T x_i)$
 can be ± 1

$$\hat{\omega}_{MLE} = \underset{\omega}{\text{argmin}} \sum_{i=1}^n \log(1 + e^{-y_i \omega^T x_i})$$

To compute, there is no closed form. But, it is a smooth convex problem!

=> Compute via gradient descent!
(GD)

(See Abu-Mostafa 3.3)

Softmax Classification

Binary: 2 classes $y \in \{-1, 1\}$

Multi-class: k classes $y \in \{c_1, c_2, \dots, c_k\}$

Conditional Probabilities:

$$P(y=1|x) = \frac{1}{1 + e^{-\omega^T x}} = \frac{e^{\omega^T x}}{1 + e^{\omega^T x}}$$

$$P(y=-1|x) = \frac{1}{1 + e^{\omega^T x}}$$

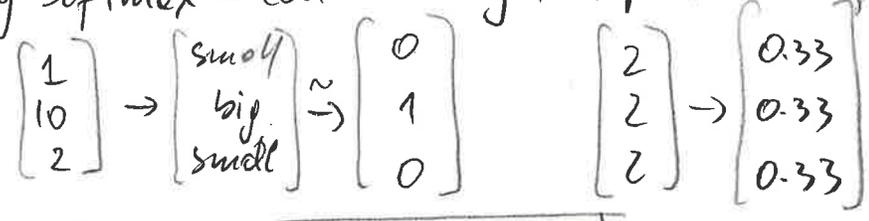
$$P(y=c_j|x) = \frac{e^{\omega_j^T x}}{\sum_{j'} e^{\omega_{j'}^T x}}$$

MLE:

$$\arg \max_{\omega} \sum_{i=1}^n \log \left(\frac{1}{1 + e^{-y_i \omega^T x_i}} \right)$$

$$\arg \max_{\omega} \sum_{i=1}^n \sum_{j=1}^k 1\{y_i = c_j\} \cdot \log \left(\frac{e^{\omega_j^T x_i}}{\sum_{j'=1}^k e^{\omega_{j'}^T x_i}} \right)$$

Why Softmax - can be easily interpreted as probabilities



Regression and Classification

ML paradigm: (define predictor $f_{\omega}(x)$ and Loss/Cost $l(f_{\omega}(x), y)$)

{Then, optimize $\hat{\omega} = \arg \min_{\omega} \sum_{i=1}^n l(f_{\omega}(x_i), y_i)$ }

Regression - squared error loss: $l(f_{\omega}(x), y) = (y - f_{\omega}(x))^2$

Logistic loss: $l(f_{\omega}(x), y) = \log [1 + \exp(-y f_{\omega}(x))]$