

Approximating Class Posteriors with a Model

Function whose Parameters are Estimated via MLE
for the 2-class case

Given dataset $D = \{(x_1, l_1), \dots, (x_N, l_N)\}$ ^{assume iid}

and model $h(x; \theta)$ to approximate

$P(L=1|x)$, and $1 - h(x; \theta)$ to approximate
 $P(L=0|x)$, use MLE to optimize θ with).

$$\hat{\theta}_{ML} = \arg\max_{\theta} p(D|\theta) = \arg\max_{\theta} \ln p(D|\theta)$$

$$= \arg\max_{\theta} \ln [p(x_1, l_1, \dots, x_N, l_N | \theta)]$$

$$= \arg\max_{\theta} \ln \left[\prod_{n=1}^N p(x_n, l_n | \theta) \right] \quad \leftarrow \text{Assume iid samples}$$

$$= \arg\max_{\theta} \sum_{n=1}^N \ln p(x_n, l_n | \theta) \quad \leftarrow \text{Bayes Rule}$$

$$= \arg\max_{\theta} \sum_{n=1}^N \ln [p(l_n | x_n, \theta) p(x_n | \theta)]$$

$$= \arg\max_{\theta} \sum_{n=1}^N \ln p(l_n | x_n, \theta) + \sum_{n=1}^N \ln p(x_n | \theta) \quad \leftarrow \text{Assume } x_n \perp \theta$$

Assuming $\{x_1, \dots, x_N\}$ are independent from θ makes the second term $\sum_{n=1}^N \ln p(x_n | \theta) = \sum_{n=1}^N \ln p(x_n)$ and this term does not depend on θ anymore, so it is an additive constant. Therefore,

$$\hat{\theta}_{ML} = \arg \max_{\theta} \sum_{n=1}^N \ln p(l_n | x_n, \theta)$$

Each $l_n \in \{0, 1\}$, takes one of 2 values so there are two cases to consider

$$\text{if } l_n = 1, p(l_n | x_n, \theta) \approx h(x_n; \theta)$$

$$\text{if } l_n = 0, p(l_n | x_n, \theta) \approx 1 - h(x_n; \theta)$$

$$\text{Let } \ln p(l_n | x_n, \theta) = \cancel{\ln p(l_n | x_n, \theta)} = l_n \ln h(x_n; \theta) + (1 - l_n) \ln(1 - h(x_n; \theta))$$

Substituting this:

$$\hat{\theta}_{ML} = \arg \max_{\theta} \sum_{n=1}^N \left[\underbrace{l_n \ln h(x_n; \theta)}_{\text{contributes for } x_n \text{ if } l_n = 1} + \underbrace{(1 - l_n) \ln(1 - h(x_n; \theta))}_{\text{contributes for } x_n \text{ if } l_n = 0} \right]$$

Multiplying with $-1/N$:

$$\hat{\theta}_{ML} = \underset{\theta}{\operatorname{argmin}} \quad -\frac{1}{N} \sum_{n=1}^N \left[\ell_n \ln h(x_n; \theta) + (1 - \ell_n) \ln (1 - h(x_n; \theta)) \right]$$

Now let's consider logistic-generalized-linear models:

$$h(x; \theta) = \frac{1}{1 + e^{-w^T b(x)}}$$

where $b(x) = \begin{bmatrix} b_0(x) \\ b_1(x) \\ \vdots \\ b_M(x) \end{bmatrix}$, typically with $b_0(x) = 1$

Logistic-linear model: $b(x) = \begin{bmatrix} 1 \\ x \end{bmatrix}$ $w = \begin{bmatrix} w_0 \\ w_x \end{bmatrix}$

e.g. $x \in \mathbb{R}^3 \Rightarrow b(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$ $w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \end{bmatrix}$

Logistic-quadratic-polynomial model: $b(x) = \begin{bmatrix} 1 \\ x \\ q(x) \end{bmatrix}$

e.g. $x \in \mathbb{R}^2 \Rightarrow b(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1^2 \\ x_1 x_2 \\ x_2^2 \end{bmatrix}$

all monic
quadratic
polynomials \rightarrow

Choose $\dim w = \dim b(x)$