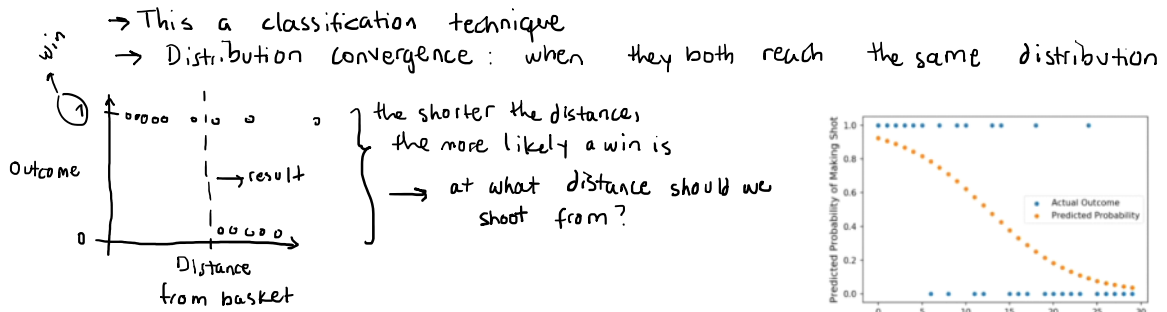


Week4

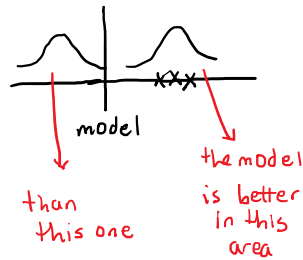
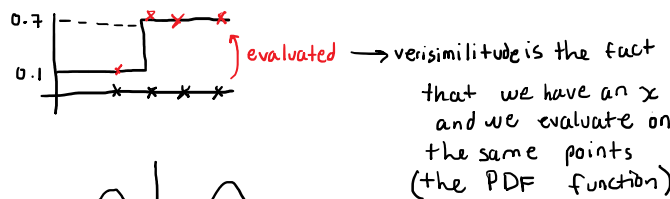
Saturday, February 26, 2022 12:32 PM

Logistic Regression



Verisimilitude (verosimilitud)

→ how similar my data is, called π



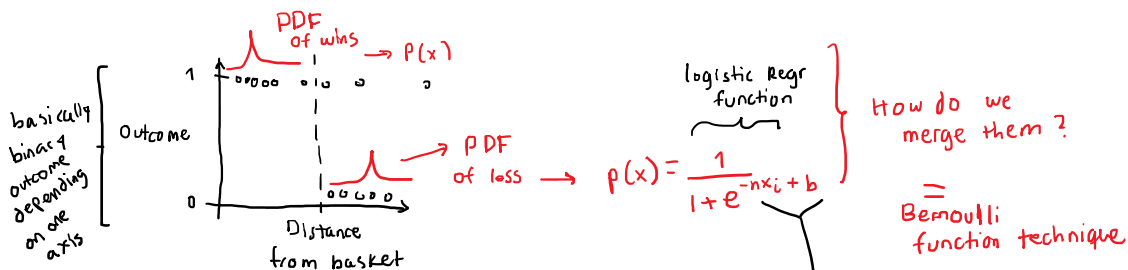
the more x points land on the distribution's density domain, the larger π we get, the better model

→ (probability density function)

look for the closest points to the density function's peak.

we want to maximize it

$$\pi = (0.1)(0.7)(0.7) = \boxed{0.049}$$



Bernoulli function: $f(x) = p(x_i)^{y_i} (1-p(x_i))^{1-y_i}$ (Bernoulli Trials)

filters what we consider when merging the PDF's for both win/loss

Thus x is verisimilitude

$$\max_{m,b} f(x) = \max_{m,b} \prod_{i=0}^N f(x_i)$$

we want to find the m, b parameters that optimize the function

$$= \max_{m,b} \prod_{i=1}^N p(x_i)^{y_i} (1-p(x_i))^{1-y_i}$$

After $\log()$,

$$= \max_{m,b} \log \left(\prod_{i=1}^N p(x_i)^{y_i} (1-p(x_i))^{1-y_i} \right)$$

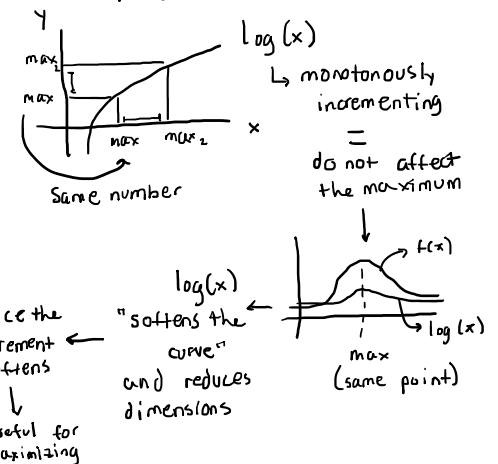
$$= \max_{m,b} \sum_{i=0}^{N-1} \log (p(x_i)^{y_i} (1-p(x_i))^{1-y_i})$$

Differentiate to maximize (and also equalize to 0)

$$\frac{d}{dm} \sum_{i=0}^{N-1} \log (p(x_i)^{y_i} (1-p(x_i))^{1-y_i}) = 0$$

$$\sum_{i=0}^{N-1} \dots$$

Why $\log(L(x_i, m, b))$?



$$\frac{\partial}{\partial m} \sum_{i=0}^{N-1} \log(p(x_i)^{y_i} (1-p(x_i))^{1-y_i}) = 0$$

useful for
maximizing

$$\frac{\partial}{\partial m} \sum_{i=0}^{N-1} \log(p(x_i)^{y_i} + \log(1-p(x_i))^{1-y_i}) = 0$$

$$\frac{\partial}{\partial m} \sum_{i=0}^{N-1} y_i \log(p(x_i)) + (1-y_i) \log(1-p(x_i)) = 0$$

$$\frac{\partial}{\partial m} \sum_{i=0}^{N-1} (y_i \log(p(x_i)) + \log(1-p(x_i)) - y_i \log(1-p(x_i))) = 0$$

Considering $p(x_i) = \frac{1}{1+e^{-mx_i+b}}$, with bias $b=0$

logistic
distribution

$$\frac{\partial}{\partial m} \sum \left[y_i \log\left(\frac{p(x_i)}{1-p(x_i)}\right) + \log(1-p(x_i)) \right] = 0$$

Evaluating the term $\frac{p(x_i)}{1-p(x_i)}$, considering $\alpha = mx_i + b$

$$\frac{p(x_i)}{1-p(x_i)} = \frac{\frac{1}{1+e^{-\alpha}}}{1 - \frac{1}{1+e^{-\alpha}}} = \frac{1}{1+e^{-\alpha}} \cdot \frac{1+e^{-\alpha}}{1+e^{-\alpha} - 1} = \frac{1+e^{-\alpha}}{e^{-\alpha}(1+e^{-\alpha})} = \frac{1}{e^{-\alpha}} = e^{\alpha}$$

$= e^{mx_i+b}$

Evaluating the term $1-p(x_i)$

$$1-p(x_i) = 1 - \frac{1}{1+e^{-\alpha}} = \frac{1+e^{-\alpha} - 1}{1+e^{-\alpha}} = \frac{e^{-\alpha}}{1+e^{-\alpha}}$$

Thus,

$$= \frac{\partial}{\partial m} \sum_{i=0}^{N-1} y_i \log(e^{mx_i+b}) + \sum_{i=0}^{N-1} \log\left(\frac{e^{-(mx_i+b)}}{1+e^{-(mx_i+b)}}\right) = 0$$

Now with $b=0$ (bias)

$$\frac{\partial}{\partial m} \left(-\sum_{i=0}^{N-1} mx_i y_i - \sum_{i=0}^{N-1} \log(1+e^{-mx_i}) \right) = \frac{\partial}{\partial m} \underbrace{\ell(x, m)}_{\max_m \ell(x, m)} = 0$$

Differentiating

$$= 2 \sum_{i=0}^{N-1} \left(x_i y_i + \frac{x_i (1+e^{-mx_i})}{1+e^{-mx_i}} \right)$$

$$= \sum_{i=0}^{N-1} x_i (y_i - p(x_i))$$

where the objective was to find parameters m and b that maximize the verisimilitude function

$$\max_{m,b} L(x_i, m, b) = \max_{m,b} \prod f(x_i)$$

this sum will output the m that maximizes the verisimilitude function

Types of Convergence

- Convergence in distribution, converges weakly or in Law (fractals)

- Convergence in mean (unbiased estimator)

- Convergence 'Almost-Sure':

o Almost surely, almost everywhere, strongly or with probability 1.

We thus said we try to find m, b parameters that maximize our function, that is, that separate the training data in almost-sure convergence. The problem is addressed with the Likelihood Method, which is defined for Bernoulli Trials as mentioned.

$$L(m, b) = \prod_{i=0}^{N-1} p(x_i)^{y_i} (1-p(x_i))^{1-y_i}$$

Thus $\max_{m,b} L(m, b)$

Our results for maximizing the Likelihood Function (verisimilitude)

were

$$\frac{\partial}{\partial m} \left(\sum_{i=0}^{N-1} x_i (y_i - p(x_i)) \right) = 0$$

Our results for maximizing the Likelihood Function (verisimilitude) were

$$\frac{\partial \mathcal{L}(m, b)}{\partial m} = \sum_{i=0}^{N-1} ((y_i - p(x_i)) x_i) = 0$$

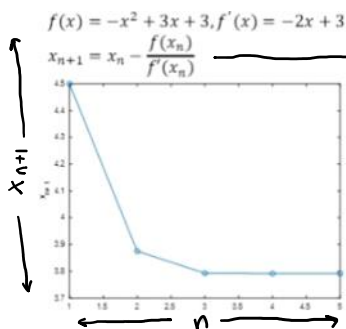
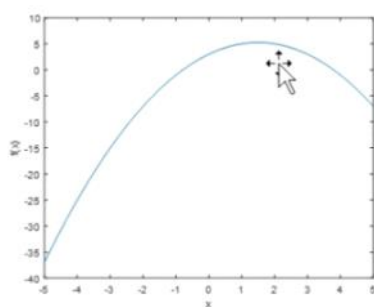
But unfortunately the last equation does not present a closed form, and thus we propose to use a Newton-Rapson method for numerical approximation

Basically we can't solve this for m since m is very inside of an exponential function, so we will use a num method for 'solving' for m .

$$m_{\text{iteration}+1} = m_{\text{iteration}} - \frac{\frac{\partial \mathcal{L}(m_{\text{iteration}})}{\partial m}}{\frac{\partial^2 \mathcal{L}(m_{\text{iteration}})}{\partial^2 m}}$$

Newton Rapson Example

- Given the $f(x)$ function, find its maximum



similar

So, our expression of Newton Rapson to the Likelihood Function

$$m_{\text{iteration}+1} = m_{\text{iteration}} - \frac{\sum_{i=0}^{N-1} ((y_i - p(x_i)) x_i)}{\sum_{i=0}^{N-1} (-x_i (1 - p(x_i)) - p(x_i))}$$

Computing
 $\frac{\partial \mathcal{L}}{\partial m}$
 $\frac{\partial^2 \mathcal{L}}{\partial^2 m}$

Example,

Given a set of training, find the parameter m and b such that separates into classes

$$X = \{-0.1, -0.5, -2, 2, 2.3, 4, 0, 5\}$$

$$Y = \{0, 0, 0, 1, 1, 1, 1\}$$

-The stop criteria between iterations is > 0.05

→ Code Available, resulting in:

$$b = 0 \text{ (bias)}$$

$$m = 1.2140$$

$$\text{error} = -0.043547$$

