

Classification \neq Regression

LINEAR regression \rightarrow Probabilistic Interpretation

Classification

Why not linear regression?

logistic Regression

METHOD: Newton's Method

RECALL LEAST SQUARES

Given $\{(x^{(i)}, y^{(i)}) \text{ for } i = 1 \dots n\}$

in which $x^{(i)} \in \mathbb{R}^{d+1}$ $y^{(i)} \in \mathbb{R}$

DO find $\Theta \in \mathbb{R}^{d+1}$ s.t. $\Theta = \underset{\Theta}{\operatorname{Argmin}} \sum_{i=1}^n (y^{(i)} - h_{\Theta}(x^{(i)}))^2$

where $h_{\Theta}(x) = \Theta^T x$

Why?

Assume $y^{(i)} = \Theta^T x^{(i)} + \epsilon^{(i)}$

\hookrightarrow error, unmodelled effects, random noise

Properties of $\epsilon^{(i)}$ we want (iid Gaussian)

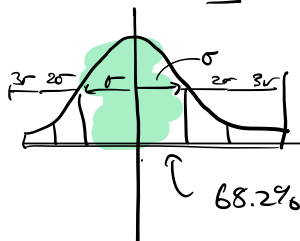
1. $\mathbb{E}[\epsilon^{(i)}] = 0$, IT'S UNBIASED

2. the errors are independent ($\mathbb{E}[\epsilon^{(i)} \epsilon^{(j)}] = \mathbb{E}[\epsilon^{(i)}] \mathbb{E}[\epsilon^{(j)}]$ if $i \neq j$)

How noisy? VARIANCE $\mathbb{E}[(\epsilon^{(i)})^2] = \sigma^2$

Turns out, unique distribution parameterized by this, the Gaussian \Rightarrow

$$\epsilon^{(i)} \sim N(0, \sigma^2) \quad \text{or} \quad P(\epsilon^{(i)}) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left\{-\frac{1}{2} \frac{(\epsilon^{(i)})^2}{\sigma^2}\right\}$$



within $1\sigma \approx 68.2\%$ of mass

$3\sigma \approx 99.7\%$ of mass.

$$\text{Therefore, } P(y^{(i)} | x^{(i)}; \Theta) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(y^{(i)} - \Theta \cdot x^{(i)})^2}{2\sigma^2}\right)$$

\hookrightarrow Parameterized by Θ

$$\text{or } y^{(i)} | x^{(i)}; \Theta \sim N(\Theta^T x^{(i)}, \sigma^2)$$

Picking $\Theta \Rightarrow$ Picks a distribution

Likelihoods Among many distributions, Pick "most likely" given all data

$$\begin{aligned} \mathcal{L}(\theta) &= p(y | x; \theta) \\ &= \prod_{i=1}^n p(y^{(i)} | x^{(i)}; \theta) \quad (\text{iid assumption}) \\ &= \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y^{(i)} - x^{(i)} \cdot \theta)^2}{2\sigma^2}\right) \end{aligned}$$

WE USE \log likelihood (convenient)

$$\begin{aligned} \ell(\theta) &= \log \mathcal{L}(\theta) \\ &= \sum_{i=1}^n \log\left(\frac{1}{\sigma\sqrt{2\pi}}\right) - \frac{(y^{(i)} - x^{(i)} \cdot \theta)^2}{2\sigma^2} \\ &= \sum_{i=1}^n \log\left(\frac{1}{\sigma\sqrt{2\pi}}\right) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y^{(i)} - x^{(i)} \cdot \theta) \end{aligned}$$

Does depend on θ

\hookrightarrow Doesn't Depend on θ

Thus, to find maximum likelihood, equivalently find

$$J(\theta) = \min_{\theta} \frac{1}{2} \sum_{i=1}^n (y^{(i)} - x^{(i)} \cdot \theta)^2$$

□

Classification

Given $\{(x^{(i)}, y^{(i)}) \text{ for } i=1 \dots n\}$

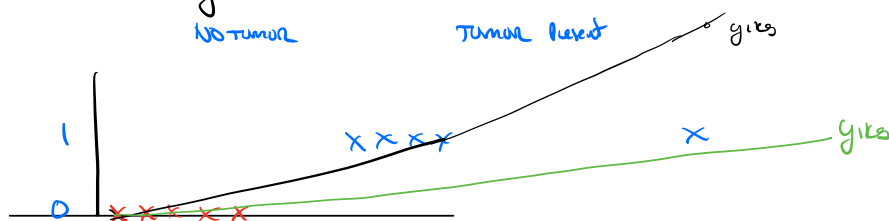
$$y^{(i)} \in \{0, 1\}$$

negative class

NOT TUMOR

Positive class

TUMOR Present

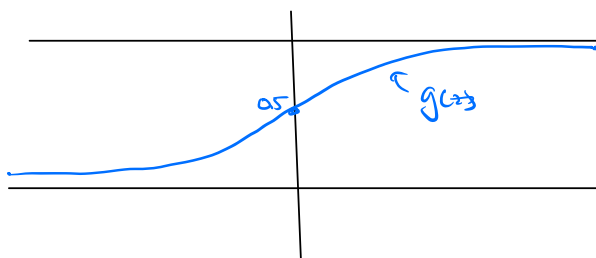


SAME RECIPE

WANT $h_{\theta}(x) \in [0, 1]$

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

$$g(z) = \frac{1}{1 + e^{-z}} \quad \text{"link function"}$$



Sigmoid or A logistic function.

$$P(y=1 | x; \theta) = h_{\theta}(x)$$

$$P(y=0 | x; \theta) = 1 - h_{\theta}(x)$$

$$\begin{aligned} \mathcal{L}(\theta) &= P(\vec{y} | x; \theta) = \prod_{i=1}^n P(y^{(i)} | x^{(i)}; \theta) \\ &= \prod_{i=1}^n h_{\theta}(x^{(i)})^{y^{(i)}} (1 - h_{\theta}(x^{(i)}))^{(1-y^{(i)})} \end{aligned} \quad \text{"ENCODER g"}$$

$$\ell(\theta) = \log \mathcal{L}(\theta) = \sum_{i=1}^n y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)}))$$

SAME RECIPE MAXIMIZE with gradient ascent

$$\theta_1 := \theta_1 + \alpha \frac{\partial \ell(\theta)}{\partial \theta_1} \quad (\text{Gradient Ascent})$$

of last week $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$ (least squares)

fun observation

$$\frac{\partial \ell(\theta)}{\partial \theta_j} = (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)} \quad \text{so ...}$$

$$\Theta := \Theta - (h_{\theta}(x) - y^{(i)}) x^{(i)}$$

We'll see later this rule is very general.

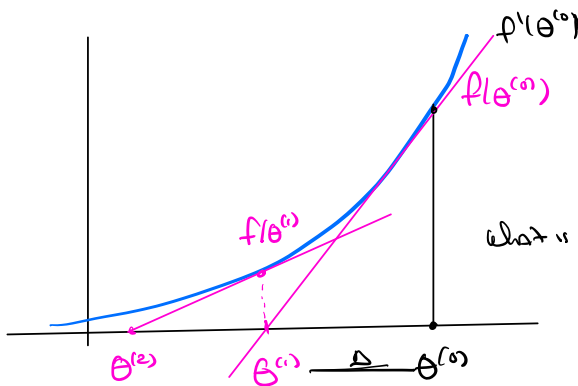
Newton's Method

Given $f: \mathbb{R}^D \rightarrow \mathbb{R}$

Do find x st. $f(x) = 0$.

ASIDE

$\max_{\theta} \ell(\theta)$ want $\ell'(\theta) = 0$
 \uparrow
 DERIVATIVE



$$\theta^{(n)} = \theta^{(n)} - \Delta$$

$$\text{What is } \Delta? \quad f(\theta^{(n)}) = f'(\theta^{(n)}) \cdot \Delta$$

$$\therefore \Delta = f'(\theta^{(n)})^{-1} f(\theta^{(n)})$$

$$\theta^{(n+1)} = \theta^{(n)} - \frac{f(\theta^{(n)})}{f'(\theta^{(n)})}$$

Converges fast! Quadratic $0.1 \rightarrow 0.01 \rightarrow 0.0001$ (Digits double!)

Generalizing to vector $\theta \in \mathbb{R}^{D+1}$ and $\ell'(\theta) = \nabla \ell(\theta)$ (i.e. minimization)

$$\theta^{(t+1)} = \theta^{(t)} - H^{-1} \nabla_{\theta} \ell(\theta)$$

\nwarrow Hessian $\in \mathbb{R}^{(D+1) \times (D+1)}$ $\searrow \in \mathbb{R}^{D+1}$

$$H_{ij} = \frac{\partial^2}{\partial \theta_i \partial \theta_j} \ell(\theta)$$

NB: NO STEPSIZE!

Rough Comparison

	<u>Per iteration</u>	<u>Compute</u>	<u>STEPS TO Error ϵ</u>
SGD	1 DATA point	$O(d)$	ϵ^{-2}
BATCH SGD	N DATA points	$O(nd)$	$\approx \epsilon^{-1}$
Newton	N DATA points	$\mathcal{O}(nd^2)$	$\log(\frac{1}{\epsilon})$

IN CLASSICAL STATS d is small 100 or so

AND Exact Answer MATTERS \Rightarrow Newton (LBFGS)

MODERN ML d is HUGE 175B $d^2 \Rightarrow \ddot{\cdot}$

~~4~~
⇒ SGD of TEN WORKHOUSE (Exact solution less so)