

Classification $\not\equiv$ Regression

Probabilistic View of LINEAR Regression

Classification

Why not linear regression?

Logistic Regression

METHOD: Newton's METHOD

- ② converges quickly \rightarrow fast
- ② each step is expensive

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 $\theta \in \mathbb{R}^{d+1}$ st. $\theta = \underset{\theta}{\operatorname{argm}} \sum_{i=1}^n (y^{(i)} - h_{\theta}(x^{(i)}))^2$

$$h_{\theta}(x) = \theta^T x$$

Why?

Assume $y^{(i)} = \theta^T x^{(i)} + \varepsilon^{(i)}$

↑
plug-in

↖ random
error or noise term
= stuff we can't really

1. $E[\varepsilon^{(i)}] = 0$... it's unbiased

2. The errors independent $E[\varepsilon^{(i)} \varepsilon^{(j)}] = E[\varepsilon^{(i)}]E[\varepsilon^{(j)}]$ for $i \neq j$

How noisy $E[(\varepsilon^{(i)})^2] = \theta^2$

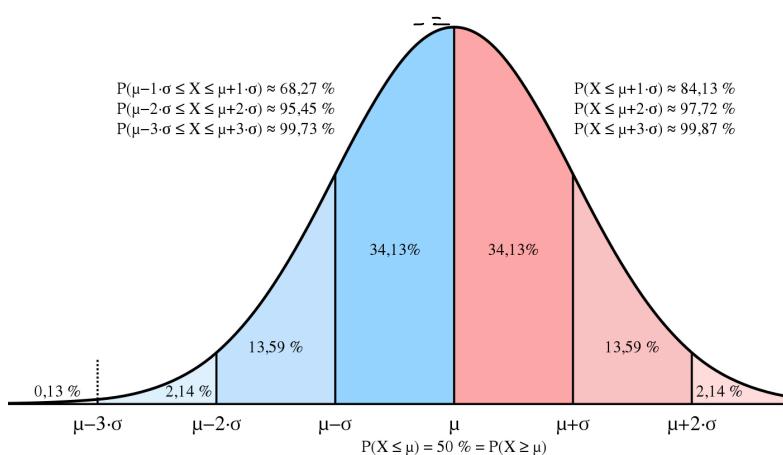
(unique of the above)

Gaussian or Normal Distribution

WRITE $\epsilon^{(i)} \sim N(\mu, \sigma^2)$

$P(z; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{(z-\mu)^2}{2\sigma^2}\right\}$

Not conditional
(density)



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

histograms

eventually converges to this

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

Picking $\theta \Rightarrow$ picks distribution

Likelihood among many distributions, "most likely"

$$\begin{aligned} L(\theta) &= P(y|X; \theta) = \prod_{i=1}^n P(y^{(i)} | x^{(i)}; \theta) \quad \text{iid} \\ &= \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{(y^{(i)} - \theta \cdot x^{(i)})^2}{2\sigma^2}\right\} \\ &\text{to find } L(\theta) \text{ maximize wrt } \theta \\ &\log L(\theta) = \sum_{i=1}^n \log \frac{1}{\sigma\sqrt{2\pi}} - \frac{(y^{(i)} - \theta \cdot x^{(i)})^2}{2\sigma^2} \\ &\log L(\theta) \text{ depends on } \theta \\ &x \text{ depend on data } (D) \end{aligned}$$

$$J(\theta) = \underset{\theta}{\operatorname{Max}} \ell(\theta) = \min_{\theta} \frac{1}{2} \sum_i (y^{(i)} - \theta \cdot x^{(i)})^2$$

Least squares?

Likelihoods Among many distributions, Pick most likely one

$$\mathcal{L}(\theta) =$$

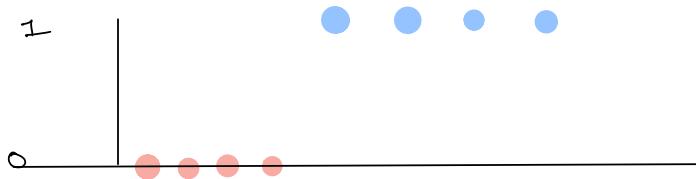
CLASSIFICATION

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

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

↑ positive class

↳ negative class



Same RECIPE AS linear Regression!

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

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

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



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

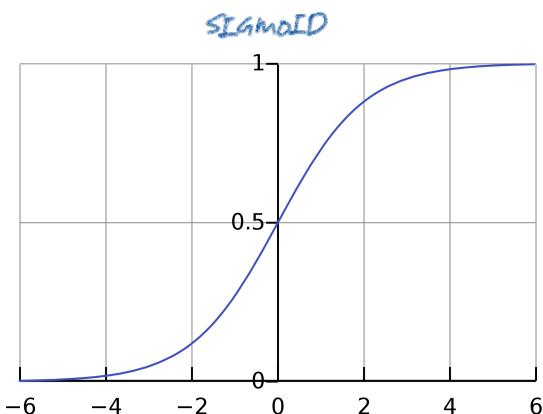
$$P(y=0|x; \theta) = 1 - h_{\theta}(x)$$

Likelihood (Probability)

$$\ell(\theta) = P(\vec{y}|\vec{x}; \theta)$$

$$\text{encoding } \sum_i = \prod_i P(y^{(i)} | x^{(i)}; \theta)$$

$$\sum_i = \prod_i h_{\theta}(x^{(i)})^{y^{(i)}} (1 - h_{\theta}(x^{(i)}))^{1-y^{(i)}}$$



① Smooth (nice transition)

②

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

↑ label

log likelihood

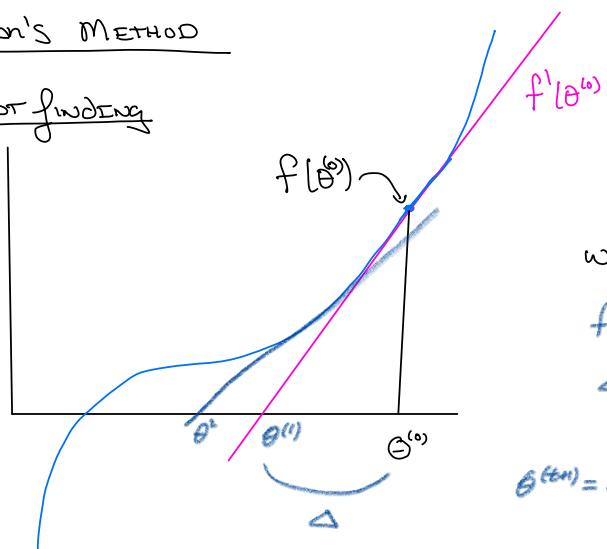
$$\text{SAME RECIPE: } \theta^{(t+1)} = \theta^{(t)} + \alpha \frac{\partial}{\partial \theta} J(\theta)$$

$$\frac{\partial}{\partial \theta_j} \ell(\theta) = \sum_{i=1}^n (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

Rule is very general

Newton's METHOD

Root finding



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

Do $f(\theta) = 0$

(Aside $\min \ell(\theta) \Rightarrow \ell(\theta) = 0$)

what is Δ ? $\theta^{(t+1)} = \theta^{(t)} - \Delta$

$$f(\theta) = f'(\theta^{(t)}) \cdot \Delta$$

$$\Delta = f'(\theta^{(t)})^{-1} \cdot f(\theta^{(t)})$$

$$\theta^{(t+1)} = \theta^{(t)} - \frac{f(\theta^{(t)})}{f'(\theta^{(t)})}$$

$$0.1 \rightarrow 0.01 \rightarrow 0.001$$

$$\theta^{(t+1)} = \theta^{(t)} - (H^T D \ell(\theta))^{-1} \in \mathbb{R}^d$$

$$\theta \in \mathbb{R}^{d+1} \quad \ell'(\theta) = f(\theta)$$

$$\downarrow \text{Hessian } \in \mathbb{R}^{(d+1) \times (d+1)}$$

$$\text{Hessian} = \frac{\partial^2}{\partial \theta_i \partial \theta_j} \ell(\theta)$$

\uparrow
Diagonal

TO find minimum,

Rough Comparison				
METHOD	Per iteration	Compute	Steps to Error ϵ^2	robustness against error
SGD	1 data point	$O(d)$	ϵ^{-2}	robust

Batch GD	N data points	$O(nd)$	ϵ^{-1}
Newton's method	N data points	$\Omega(nd^2)$	$\approx \log(\frac{1}{\epsilon}) \rightarrow \epsilon^{-2}$ robust
↴ mini batch (parallelism) batch size denominator (b)			