# Classification & Regression

Probabilistic View of Linear Regression

Classification

Why not linear Regression?

Logistic Regression

METHOD: Newton's METHOD

# Recall Least Squares

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

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

Do find $\theta \in \mathbb{R}^{d+1}$ s.t. $\theta = \arg\min_{\theta} \sum_{i=1}^{n} (y^{(i)} - h_\theta(x^{(i)}))^2$

where $h_\theta(x) = \theta^T x$

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

$\hookrightarrow$ Error or noise term.

## Properties

1. $\mathbb{E}[\epsilon^{(i)}] = 0$ — It's unbiased

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

How "Noisy" $\mathbb{E}[(\epsilon^{(i)})^2] = \sigma^2$
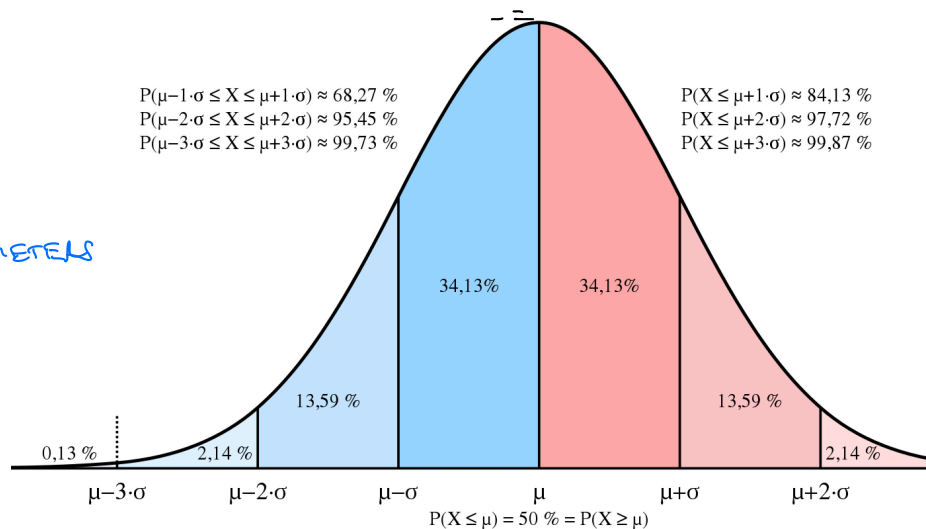
# Gaussian or Normal Distribution    (Unique of the Above)

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

$\mu = 0$ → Mean

→ Variance

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

→ Parameters

$$\exp(x) = e^x$$



P(μ−1·σ ≤ X ≤ μ+1·σ) ≈ 68,27 %
P(μ−2·σ ≤ X ≤ μ+2·σ) ≈ 95,45 %
P(μ−3·σ ≤ X ≤ μ+3·σ) ≈ 99,73 %

P(X ≤ μ+1·σ) ≈ 84,13 %
P(X ≤ μ+2·σ) ≈ 97,72 %
P(X ≤ μ+3·σ) ≈ 99,87 %

34,13%   34,13%

13,59 %        13,59 %

0,13 %   2,14 %                2,14 %

μ−3·σ   μ−2·σ   μ−σ   μ   μ+σ   μ+2·σ

P(X ≤ μ) = 50 % = P(X ≥ μ)

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

$= \epsilon^{(i)}$

→ Parameter

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

Picking $\theta$ $\Rightarrow$ Picks distribution

**Likelihood** among many distributions, "most likely"

$$\mathcal{L}(\theta) = P(y \mid X; \theta) = \prod_{i=1}^{n} P(y^{(i)} \mid x^{(i)}; \theta) \quad (iid)$$

$$= \prod_{i=1}^{n} \frac{1}{\sigma\sqrt{2\pi}} \cdot \exp\left\{ -\frac{(y^{(i)} - \theta \cdot x^{(i)})^2}{2\sigma^2} \right\}$$

log likelihood

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

→ Depends on Data & $\theta$.

$$J(\theta) = \underset{\theta}{\text{MAX }} \ell(\theta) = \underset{\theta}{\text{MIN }} \frac{1}{2} \sum_{i=1}^{n} (y^{(i)} - \theta \cdot x^{(i)})^2$$
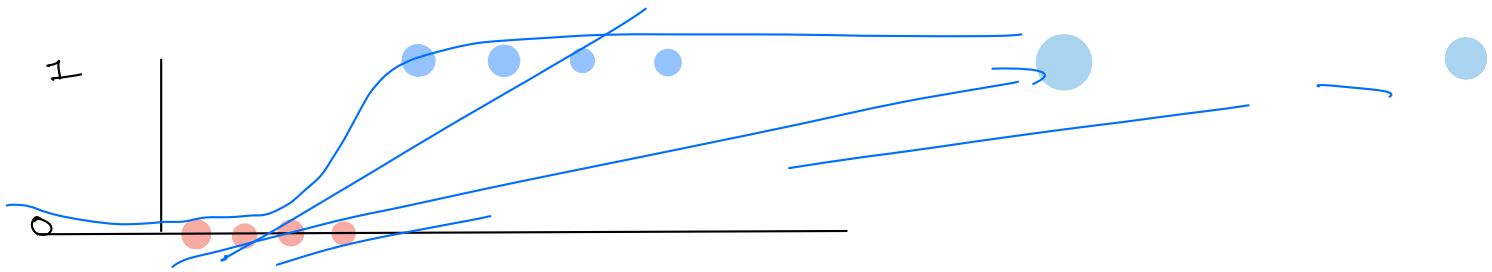
# Likelihoods

Among many <mark>distributions</mark>, Pick most likely one

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

# Classification

Given $(x^{(i)}, y^{(i)})$ for $i = 1 \ldots n$ $\qquad$ $y^{(i)} \in \{0, 1\}$

Positive class

Negative class

## SAME RECIPE AS linear Regression!
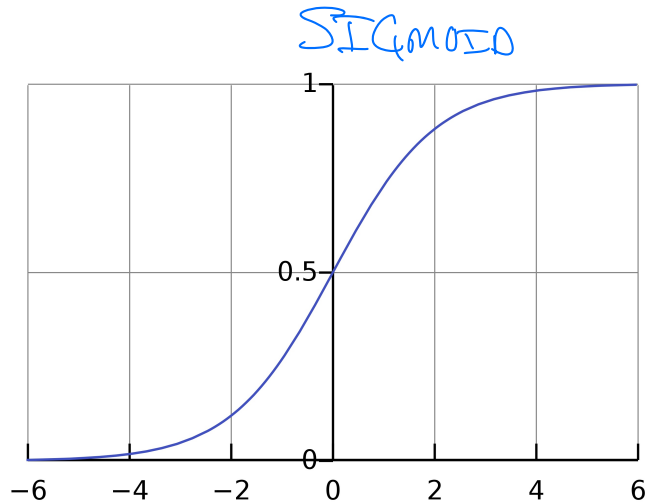
$$h_\theta(x) \in [0,1]$$

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

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

SIGMOID



$$P(y = 1 \mid x; \theta) = h_\theta(x)$$

$$P(y = 0 \mid x; \theta) = 1 - h_\theta(x)$$

↙ likelihood

$$\mathcal{L}(\theta) = P(\vec{y} \mid \vec{x}; \theta) = \prod_{i=1}^{n} P(y^{(i)} \mid x^{(i)}; \theta)$$

↙ log likeldo

$$= \prod_{i=1}^{n} h_\theta(x^{(i)})^{y^{(i)}} (1 - h_\theta(x^{(i)}))^{1 - y^{(i)}}$$

$$\ell(\theta) \log \mathcal{L}(\theta) = \sum_{i=1}^{n} y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)}))$$
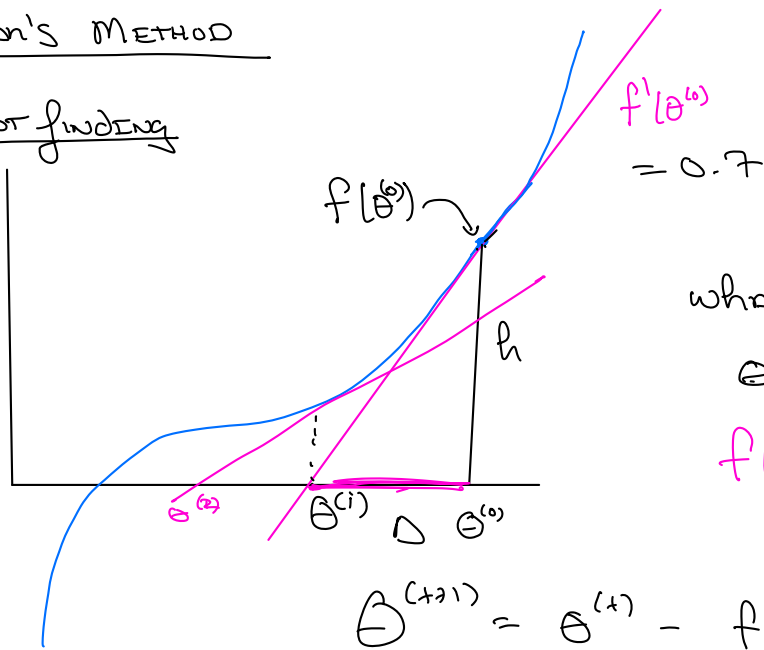
SAME RECIPE: $\quad \theta^{(t+1)} = \theta^{(t)} + \alpha \frac{\partial}{\partial \theta_i} J(\theta)$

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

RULE IS VERY general.

# Newton's Method

## Root finding



$f'(\theta^{(0)})$
$= 0.7$

$f(\theta^{(0)})$

$h$

$\theta^{(2)}$    $\theta^{(1)}$   $\Delta$  $\theta^{(0)}$

GIVEN $f: \mathbb{R}^d \to \mathbb{R}$

DO find $f(x) = 0$

( ASIDE $\min \ell(\theta) \Rightarrow \ell'(\theta=0)$ )

what is $\Delta$?

$$\theta^{(1)} = \theta^{(0)} - \Delta$$

$$f(\theta^{(0)}) = f'(\theta^{(0)}) \cdot \Delta$$

$$\Delta = f'(\theta^{(0)})^{-1} \cdot f(\theta^{(0)})$$

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

Error   $0.1 \longrightarrow$   $0.01$  $\rightarrow$  $0.0001$

Generalize & use for min.  $\theta \in \mathbb{R}^{d+1}$   $\ell'(\theta) = f(\theta)$

$$\theta^{(t+1)} = \theta^{(t)} - H^{-1} \nabla_\theta \ell(\theta) \quad \to \in \mathbb{R}^d$$

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

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

To find minimum,

# Rough Comparison

| METHOD | Per iteration | Compute | Steps to Error $\varepsilon^2$ |
|---|---|---|---|
| SGD | 1 Data Point | $\Theta(d)$ | $\varepsilon^{-2}$ |
| BATCH GD | N Data Points | $\Theta(nd)$ | $\approx \varepsilon^{-1}$ |
| Newton Method | N Data points | $\Omega(nd^2)$ | $\approx \log\left(\frac{1}{\varepsilon}\right)$ |

MINIBATCH

$I_{N}$ classical stats   $d$ was small

$n$ was moderate