# L1
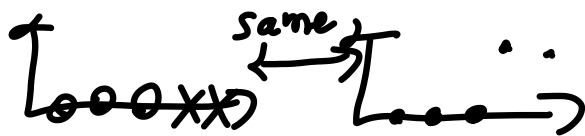
ML : machine learn with task, performance, experience

1. Supervised learning : $(x, y) \in E$, $x \to y$

   ① regression : $y \to$ continuous
     logistic regression : seperate good/bad
   ② classification : $y \to$ discrete



   CV : digitize $\to$ learn neural
              $\to$ execute

ML decisions many : strategic important
     never work first-time : debug important
                            strategy

2. Deep learning:

3. unsupervised learning: label outputs X
→ find interesting clustering

eg. cocktail party problem

4. reinforce learning: do stuff → reward

eg. dog training

# L2

Linear regression

Train Set T

↓

feed algorithm

↓ function

ed.
size $\longrightarrow$ [ h ] $\longrightarrow$ house price

hypothesis

How to represent h ?

$$h(x) = \sum_{j=0}^{n} \theta_j x_j \quad, \quad x_0 = 1$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} \quad, \quad x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

$\theta$ : parameter

m : #training examples

X : input / features

y : output / target variable

(x, y) : train examples

$(x^{(i)}, y^{(i)})$ = ith example

n : # features

Goal : ? classify performance?

Linear regression: $J(\theta) = \min\left(\frac{m}{2} \sum_{i=1}^{m} (h(x) - y)^2\right)$

$\downarrow$

Gradient descent : 1. start with $\theta$

2. change $\theta$ to $\theta'$

? no local minima?

$\theta_j := \theta_j - \boxed{\alpha} \frac{\partial}{\partial \theta_j} J(\theta)$
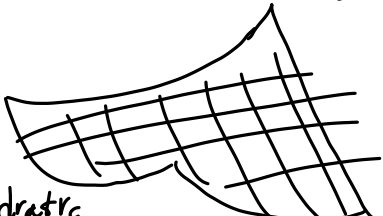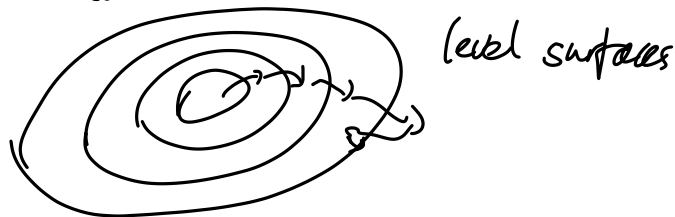
$\downarrow$ learning rate

$$\frac{\partial}{\partial \theta_{ij}} (J(\theta)) = (h_\theta(x) - y) x_j$$

Repeat until converge

$$\theta_j := \theta_j - \alpha \frac{\partial (J(\theta))}{\partial \theta_j}$$
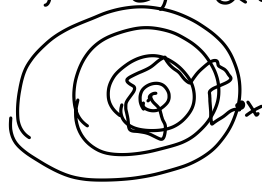
$J(\theta)$ :

↓

because quadratic

so no local minima

a bowl

level surfaces

$\alpha$ value:  ↓ → slow

↑ → overshoot → slow

? make sure one → one ? how?

no need I think

Batch gradient descent :  disadvantage : m↑ → slow

Alternative Stachastic gradient descent

Repeat

For i to m

$$\theta_j := \theta_j - \alpha (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Never converge!

switch stocastic to batch?

no. instead ↓↓α

$$\nabla_\theta J(\theta) = \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \vdots \\ \frac{\partial J}{\partial \theta_2} \end{bmatrix}$$

$$\nabla_A f(\theta) = \begin{bmatrix} \frac{\partial f}{\partial A_{1,1}} & \frac{\partial f}{\partial A_{1,2}} \\ \vdots & \ddots \end{bmatrix} \text{ if } f : \mathbb{R}^{m \times n} \to \mathbb{R}$$

$$\nabla_\theta J(\theta) = \vec{0} \longrightarrow \text{global min}$$

$$f(A) = tr \, AB, \quad \nabla_A f(A) = B^T$$

$$\nabla_A tr \, AA^T C = CA + C^T A$$

$$\vec{X} = \begin{bmatrix} - & x^{(1)T} & - \\ - & x^{(2)T} & - \\ & \vdots & \end{bmatrix}$$

similar for $\vec{y}$

$$J(\theta) = \frac{1}{2}(X\theta - y)^T (X\theta - y)$$

$$\nabla_\theta J(\theta) = X^T X \theta - X^T y \stackrel{set}{=} \vec{0}$$

$$\theta = \underbrace{\boxed{(X^T X)}^{-1} X^T y}_{\text{must exist}}$$

# L3

locally weighted regression $\Rightarrow$ fit non-linear

parametric / non-parametric learning algor
-ithm

$\downarrow$

fixedset $\theta(\theta_i)$

$\downarrow$

Size $\theta$ (set $\theta_i$)
increases

Locally weighted — regression

$$J(\theta) = \sum_{i=1}^{m} w^i (y^{(i)} - \theta^T x^{(i)})^2$$

$$w^i = e^{-(x^i - x)^2 / 2\sigma^2} \rightarrow \text{bandwidth}$$
$$\searrow \text{decide width}$$

if $|x^{(i)} - x| \rightarrow 0$, $w^i \rightarrow 1$

if $|x^{(i)} - x| \rightarrow \infty$, $w^{(i)} \rightarrow 0$

only these matter

# Probabilistic interpretation

Why LS? (Assume) $\textcircled{1} y^{(i)} = \theta^T x^{(i)} + \varepsilon^{(i)}$

$\uparrow$ error

$\theta \varepsilon^{(i)} \sim N(0, \sigma^2)$

Assume $\textcircled{2}$ IID $\to$ independent & identically distributed

$\Rightarrow p(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma} e^{\frac{-(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}}$

$\llcorner$ parametrization, $\theta$ not random variable

$\to$ likelyhood

$L(\theta) = p(\vec{y} | \vec{x}; \theta)$

$= \prod_{i=1}^{m} p(y^{(i)} | x^{(i)}; \theta)$

likelyhood $\to$ fixed $\vec{x}$, vary $\theta \to$ likelyhood of $\theta$

prob $\to$ fixed $\theta$, vary $\vec{x} \to$ prob of data

log likelyhood

$\ell(\theta) = \log L(\theta)$

$= m \log \frac{1}{\sqrt{2\pi}\sigma} + \sum_{i=1}^{M} - \frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}$

maximize $\ell(\theta) \to$ minimize

$\Rightarrow$ maximal likelyhood to LS

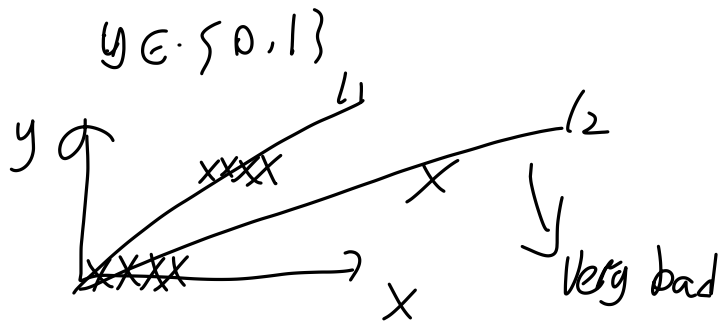Non IID → not bother to have better computation

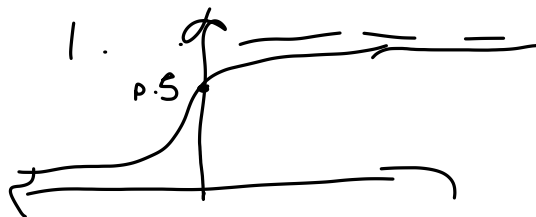Maximal likelyhood estimation

classification problem

$y \in \{0, 1\}$



SO linear regression bad

Logistic regression

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

"sigmoid" or "logistic function"



$$P(y \mid x; \theta) = h^y (1-h)^{1-y}$$
$$\text{for } y \in \{0, 1\}$$

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

To maximize likelyhood : use Batch

all for GLM

$$\theta_j = \theta_j + \frac{\partial}{\partial \theta_j} \ell(\theta) \cdot \alpha$$

$$\theta_j := \theta_j + \alpha \sum_{i=1}^{m} (y^{(i)} - h_\theta(x^{(i)})) x^{(i)}_j$$

There is no local minima/maxima
for the log function, so for logistic function

Sadly no normal equation

# Newton's method

$$\ell'(\theta) = 0 \implies \text{maxima}$$

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

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

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

"Quadratic error"
↓

#sigmfract digits doubles on
one iteration

When $\theta$ is a vector

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

actually
↗

$H$ : Hessian matrix → $\nabla_\theta^2$

$H^{-1}$ equivalent to $\frac{1}{H}$

$$H_{ij} : \frac{\partial \ell^2}{\partial \theta_i \partial \theta_j}$$

Newton pro: fast   cons: matrix large
then bad.

Perceptron algorithm

$$g(z) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$$

$$h_\theta(x) = g(\theta^T x) \quad \theta_j := \theta_j + \alpha \left( y^{(i)} - h_\theta(x)^{(i)} \right) x_j^{(i)}$$

$0$ : right $\quad y^{(i)} = 1$

$\pm 1$ : if wrong $\quad y^{(i)} = 0$