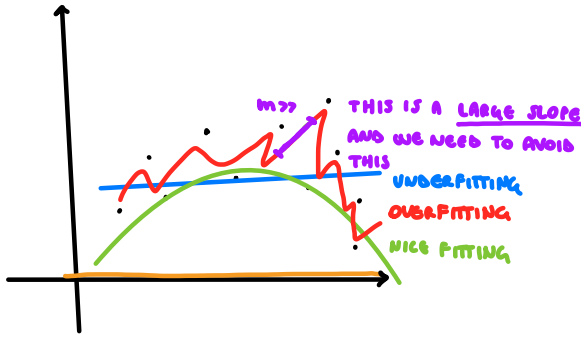# REGULARIZATION

IN PARAMETRIC MODELS, $\theta$ IS EXTREMELY DATA DEPENDENT AND THIS CAN BE A PROBLEM
THANKS TO REGULARIZATION, WE START FROM AN UNDERFITTING PREDICTOR $h_\theta$ AND OBTAIN SOMETHING OPTIMAL



m?? THIS IS A LARGE SLOPE
AND WE NEED TO AVOID
THIS

UNDERFITTING
OVERFITTING
NICE FITTING

RECALL THE LOSS FUNCTION

$$J(\theta) = \underbrace{\sum_{i=1}^{n} \left( y^{(i)} - \theta^T x^{(i)} \right)^2}_{1} + \underbrace{\lambda \; REG}_{2}$$

1. ACTUAL SQUARED LOSS, $\nabla_\theta (1)$ AND $\underset{\hat{\theta}}{\arg\min}$ MAKES THE MODEL OVERFITTING

IN A CERTAIN SENSE, WE NEED TO MAKE $J(\theta)$ <u>WORSE</u> BY A TERM REGULARIZATION

$$REG = \lambda \underbrace{\|\theta\|_2^2}_{\substack{L_2 \; NORM \\ \sum_i \theta_i^2}} , \; \lambda \underbrace{\|\theta\|}_{\substack{L_1 \; NORM \\ \sum_i |\theta_i|}}$$

$\|\theta\| \gg \longrightarrow COST \gg$

# BAYESIAN INTERPRETATION OF REGULARIZATION AND M.A.P.

2 APPROACHES

- FREQUENTISTIC : MAX LIKELIHOOD, COMPLETELY DATA DEPENDENT

$$\hat{\ell}(\theta) = \log \prod_{i=1}^{n} P(y^{(i)} | x^{(i)}, \theta) \rightarrow \hat{\theta} = \arg\max_{\theta} \ell(\theta) \rightarrow \text{PREDICTION}$$

  (IN LOG. REG AND LIN. REG IF THE MEAN OF THE DISTRIBUTION IS 0, THEN MIN(COST) = MAX(LIKEL.)

- BAYESIAN : PRIOR $P(\theta)$ → DATA OBSERVATION → POSTERIOR $P(\theta | y, x)$ → BAYES

$$P(\theta | y, x) = \frac{P(\theta) P(y | x, \theta)}{\int P(\theta) P(y | x, \theta)} \longrightarrow \text{POSTERIOR PREDICTIVE DISTRIBUTION}$$

M.A.P. ESTIMATION, NEW APPROACH! (MAXIMUM A POSTERIOR EST.)
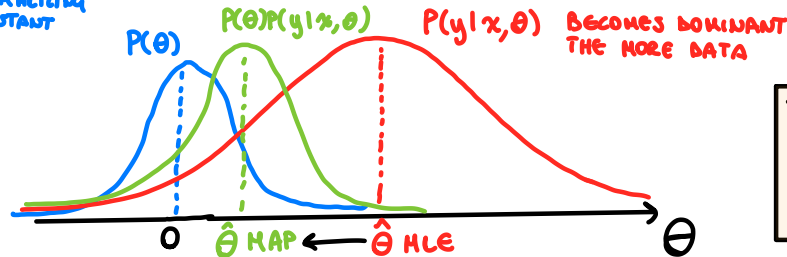IT'S A COMPROMISE BETWEEN THE 2 PREVIOUSLY MENTIONED APPROACHES
CONSISTS OF CALCULATING $\theta$ POINT BY ==MAXIMIZING THE POSTERIOR DISTRIBUTION==    $\left( \sum \log P(y^i | x^i, \theta) \right)$

$$P(\theta | x, y) = \frac{P(\theta) P(y | x, \theta)}{K} \longrightarrow \arg\max_{\hat{\theta}} P(\theta | x, y) = \arg\max \; \overset{\text{PRIOR}}{P(\theta)} \overset{\text{LIKELYHOOD}}{P(y | x, \theta)} = \arg\max \; \log P(\theta) + \log P(y | x, \theta)$$

$K \rightarrow$ NORMALIZING CONSTANT

BY NORMALIZING BY THE PRIOR
IS AS IF WE WERE PULLING
$\hat{\theta}_{MLE}$ TOWARDS $\emptyset$

$P(\theta)$

$P(\theta) P(y | x, \theta)$

$P(y | x, \theta)$ BECOMES DOMINANT THE MORE DATA



0     $\hat{\theta}$ MAP ← $\hat{\theta}$ MLE     $\theta$

THE CHOICE OF $P(\theta)$ DETERMINES THE REGULARIZATION ( ex $L_1, L_2, \ldots$ )

# NEURAL NETWORKS AND DEEP LEARNING

UNTIL NOW, WE COULD INTRODUCE NON-LINEARITY IN OUR MODELS THANKS TO <u>FEATURE MAPS</u>

$$\bar{y} = h_\theta(x) = \Theta^T \underset{\underset{\text{NON-LINEAR!}}{\uparrow}}{\phi(x)} \quad \left( \phi : \mathbb{R}^d \mapsto \mathbb{R}^p \right)$$
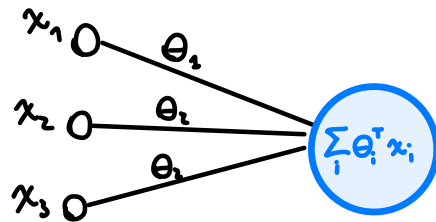
AND WE CAN SEE THIS AS ADDING "WIDTH" TO THE MODEL. BUT ANOTHER WAY TO INTRODUCE NON-LINEARITY IS <u>DEEP LEARNING</u>, BY INTRODUCING <u>DEPTH</u> TO THE MODEL
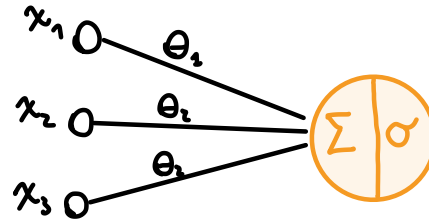
<u>*ex*</u>

$$x \in \mathbb{R}^3 \longrightarrow \begin{cases} x_1 : \text{SIZE} \\ x_2 : n \text{ BEDS} \\ x_3 : \text{ZIPCODE} \end{cases}$$

LINEAR REGRESSION

LOGISTIC REGRESSION



$x_1$ — $\theta_1$ ⟶ $\sum_i \theta_i^T x_i$

$x_2$ — $\theta_2$

$x_3$ — $\theta_3$

$x_1$ — $\theta_1$ ⟶ $\Sigma \mid \sigma$

$x_2$ — $\theta_2$

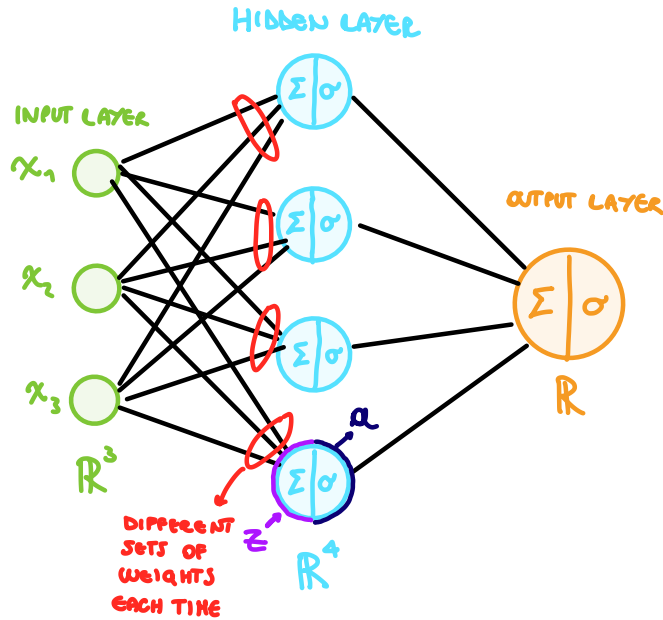$x_3$ — $\theta_3$

◯ AND ◯ ARE CALLED NEURONS!

# ex NEURAL NETWORK

INPUT LAYER TAKES INPUT $x \in \mathbb{R}^d$

HIDDEN LAYER (S) PERFORMS LINEAR
COMBINATION, APPLIES NON-LINEARITY

$$z \rightarrow \bigcirc \rightarrow a$$

USUALLY PERFORM THE SAME OPER.
BUT WITH DIFFERENT PARAMETERS

OUTPUT LAYER OUTPUTS $\hat{y}$

HIDDEN LAYER

INPUT LAYER

$x_1$

$x_2$

$x_3$

$\mathbb{R}^3$

$\Sigma | \sigma$

$\Sigma | \sigma$

$\Sigma | \sigma$

$\Sigma | \sigma$

$a$

$z$

$\mathbb{R}^4$

DIFFERENT
SETS OF
WEIGHTS
EACH TIME

OUTPUT LAYER

$\Sigma | \sigma$

$\mathbb{R}$

WE CAN WORK WITH THE ARCHITECTURE BY PLAYING AROUND :
- — N OF NEURONS
- — OPERATIONS
- — NON-LINEARITY
- — OTHER CHOICES

- INPUTS: $x \in \mathbb{R}^d$

- $a_i^{[l]}$, $l$: LAYER, $i = 1, ..., n$ → THING COMING OUT A HD NEURON, TRANSFORMED SCALAR (VECTOR
  FOR INSTANCE, WE COULD DEFINE INPUTS $x$ SUCH AS $x = a^{[0]}$ ($x_1 = a_1^{[0]}$, $x_2 = a_2^{[0]}$, ...)

- $z^{[l]}$, $l$: LAYER → THING GOING INSIDE A LAYER, IS A LINEAR COMBINATION (SCALAR) DEFINED AS FOLLOWS

$$z^{[l]} = \sum_i a_i^{[l-1]} \theta_i + b_i$$

$$z^{[l]} = W^{[l]} a^{[l-1]} + \vec{b}$$

WEIGHT MATRIX        BIAS



N OF NEURONS ON PREV. LAYER

$$z = \begin{bmatrix} \\ \\ \end{bmatrix} \times \begin{bmatrix} \\ \\ \end{bmatrix} \Rightarrow \begin{bmatrix} \\ \\ \end{bmatrix} + \begin{bmatrix} \\ \\ \end{bmatrix}$$

N OF NEURONS CURR. LAYER    NN PREV. LAYER    CURR    CURR

$b$

$$W^{l \; ? \times d}, \; b^{l \; ?}$$

THEN WHAT HAPPENS IS

$$z^{[2]} \longrightarrow \sigma \text{ (NON-LINEARITY)} \longrightarrow a_i^{[2]}$$

THE MOST COMMON NON-LINEARITY FUNCTION ARE:

| | | |
|---|---|---|
| SIGMOID | $\sigma(z) = \dfrac{1}{1 + e^{-z}}$ |  |
| TANH | $\tanh(z) = \dfrac{e^z - e^{-z}}{e^z + e^{-z}}$ |  |
| ReLU | $ReLU(z) = \max\{z, 0\}$ |  |

SO, TO RECAP

$$a^{[L]} = \sigma\left(z^{[L]}\right)$$

AND, FINALLY

$$\hat{y} = a^{[L]}$$