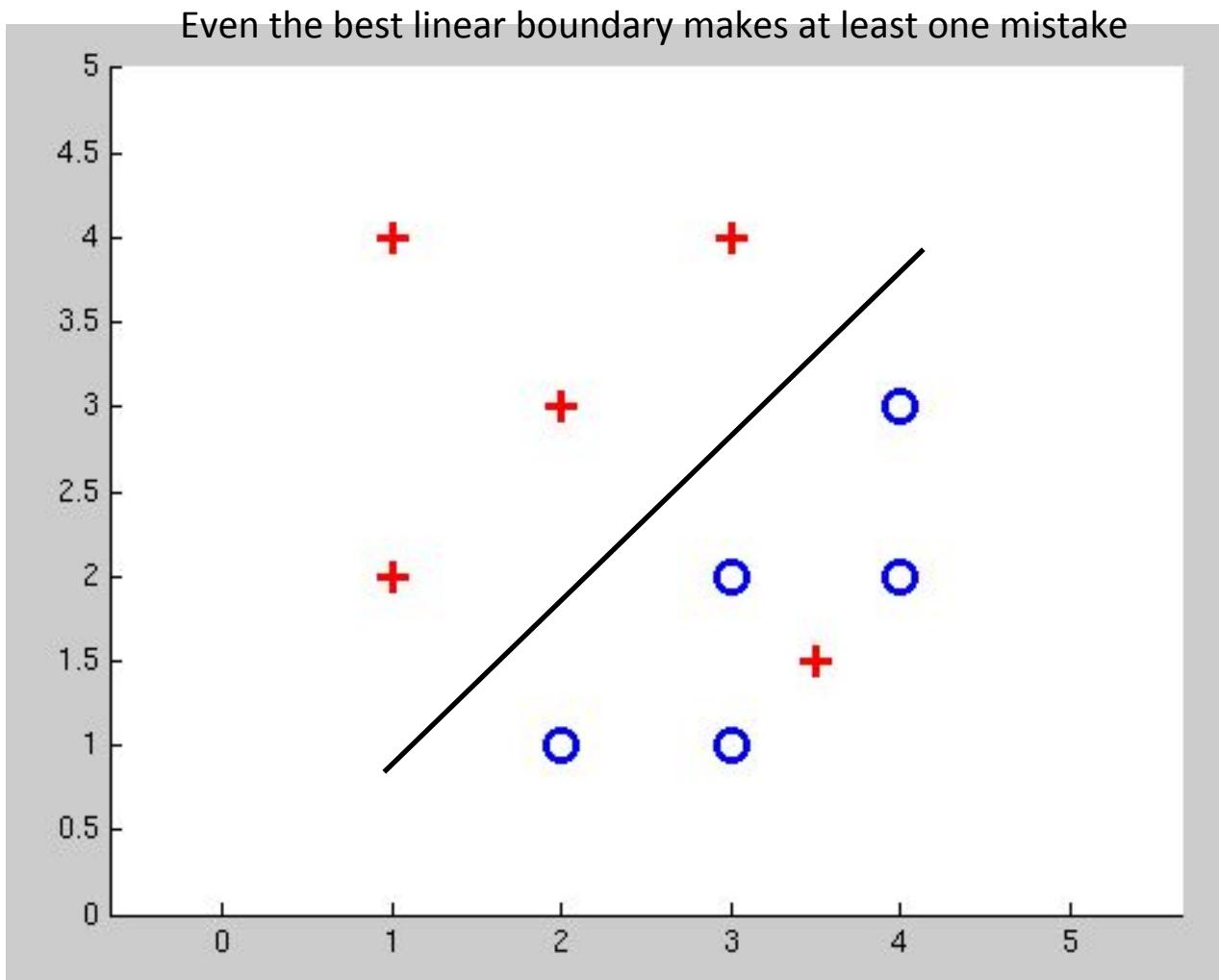


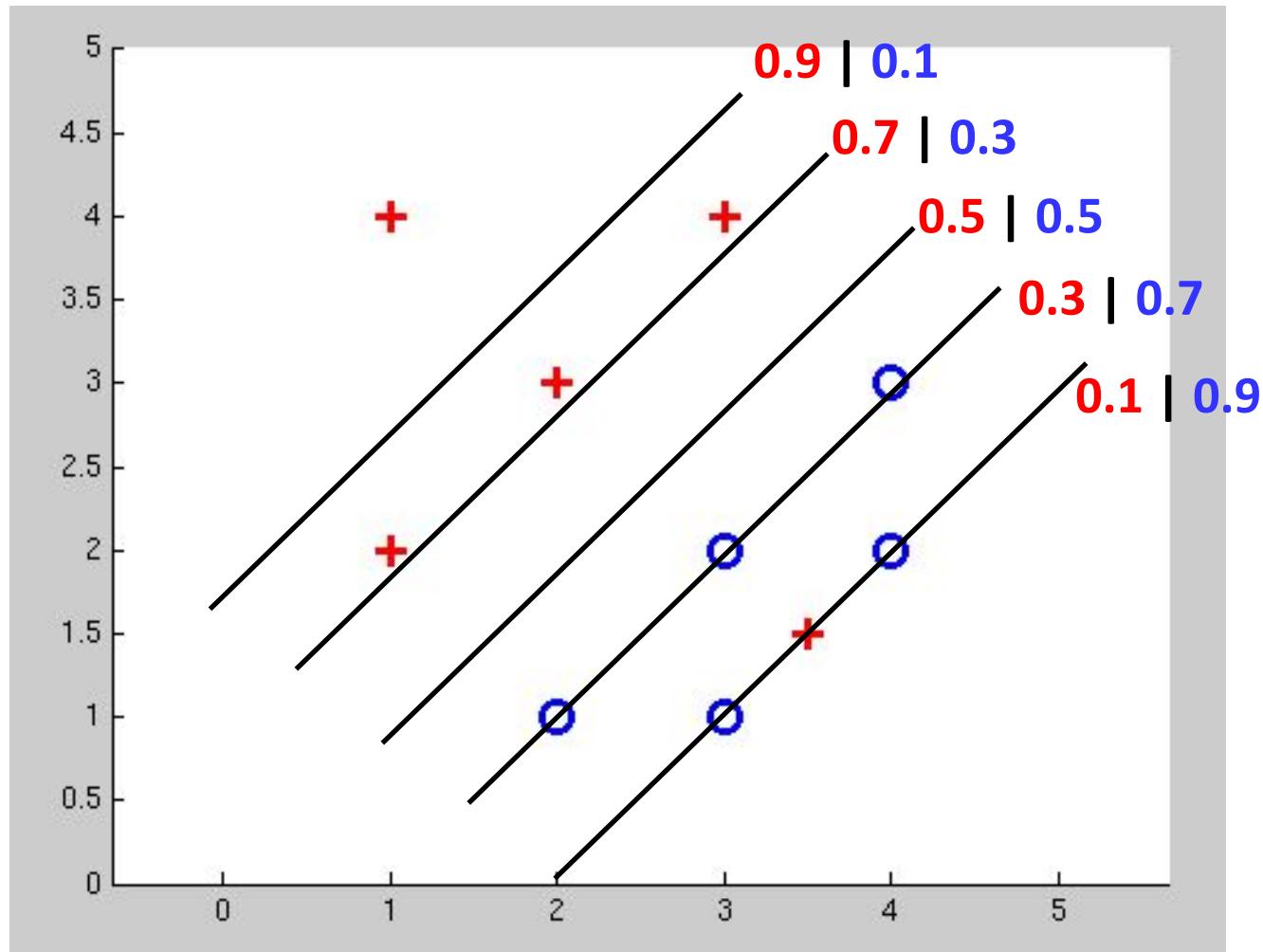
## Improving the Perceptron

## Non-Separable Case: Deterministic Decision



## Non-Separable Case: Probabilistic Decision

$$[\underline{0} \quad \overline{1}] \xleftarrow{\quad} \\ w \cdot f(x) \rightarrow [-\infty, +\infty]$$



$$w \cdot f(x) \Rightarrow + + \text{ve}$$

$\hookrightarrow 1$

## How to get probabilistic decisions?

- Perceptron scoring:  $z = w \cdot f(x)$
- If  $z = w \cdot f(x)$  very positive → want probability going to 1
- If  $z = w \cdot f(x)$  very negative → want probability going to 0

$$w \cdot f(x) \Rightarrow - - \text{ve}$$

$\hookrightarrow 0$

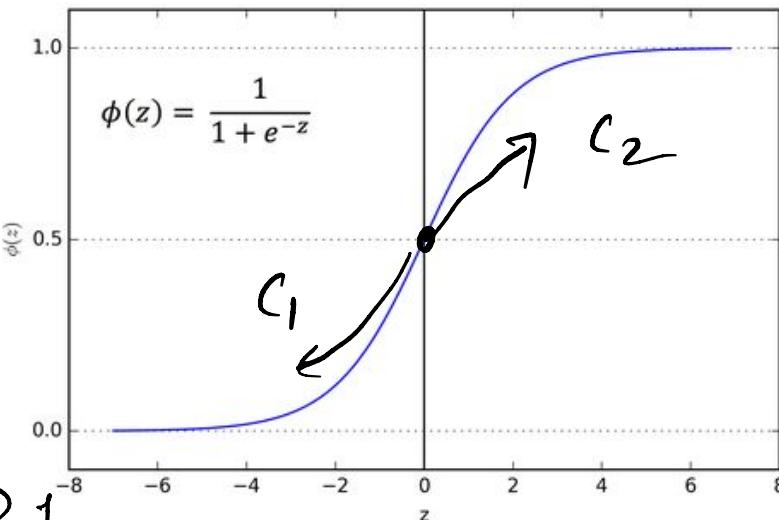
## How to get probabilistic decisions?

- Perceptron scoring:  $z = w \cdot f(x)$
- If  $z = w \cdot f(x)$  very positive → want probability going to 1
- If  $z = w \cdot f(x)$  very negative → want probability going to 0
- Sigmoid function

$$\phi(z) = \frac{1}{1 + e^{-z}}$$

$\infty \rightarrow 1$

$-\infty \rightarrow 0$



Best w?

loss funct' / cost funct'

- Maximum likelihood estimation:

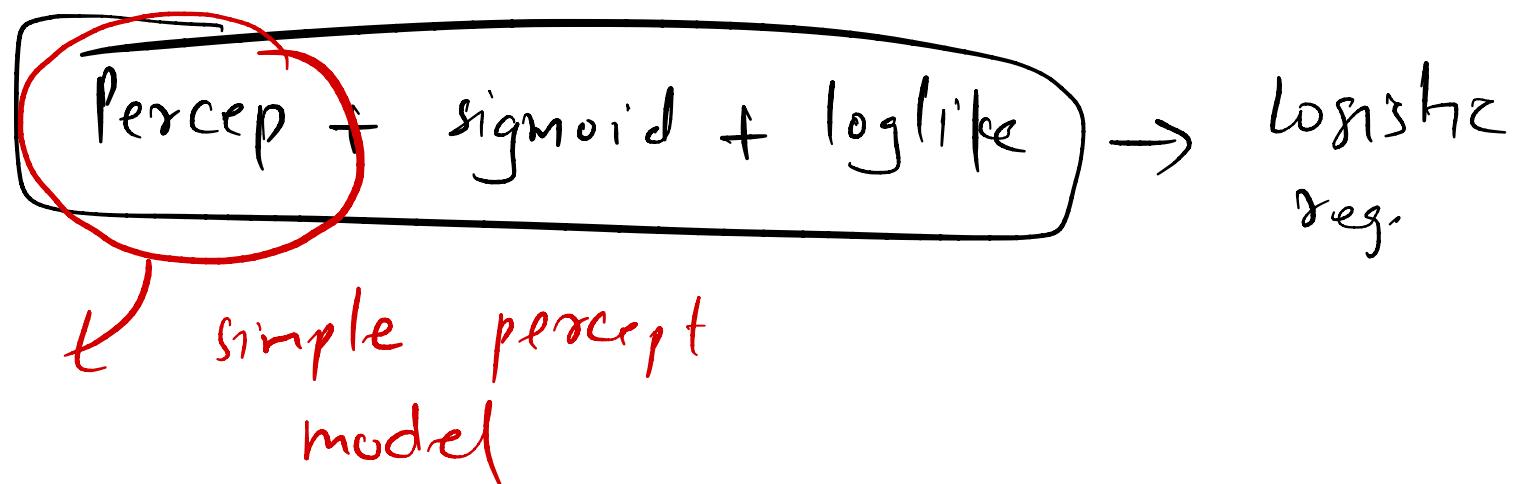
$$\max_w ll(w) = \max_w \left[ \sum_i \log P(y^{(i)} | x^{(i)}; w) \right]$$

with:

$$P(y^{(i)} = +1 | x^{(i)}; w) = \frac{1}{1 + e^{-w \cdot f(x^{(i)})}} \Rightarrow p_1$$

$$P(y^{(i)} = -1 | x^{(i)}; w) = 1 - \frac{1}{1 + e^{-w \cdot f(x^{(i)})}} = 1 - p_1$$

= Logistic Regression



$10^0, 1000$

(1.) init  $\omega$

(2)  $z = f(x) \cdot \omega$

(3)

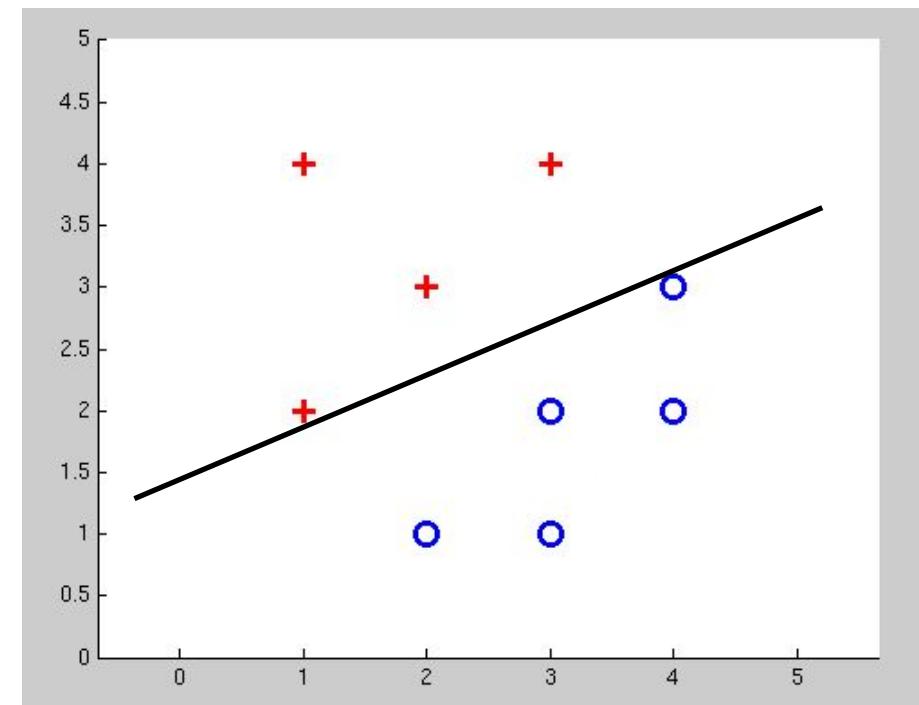
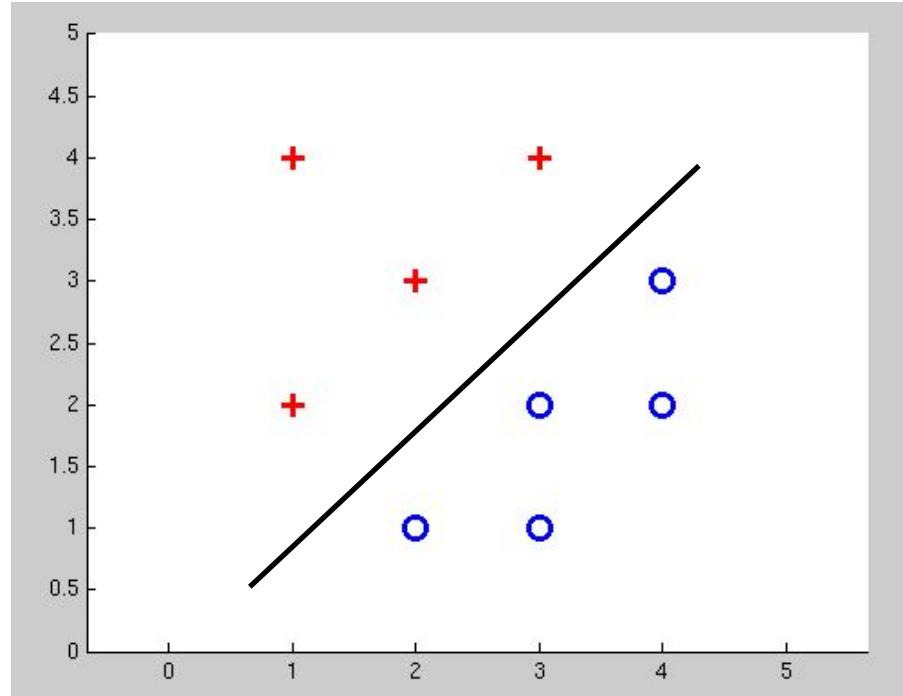
$$\frac{1}{1 + e^{-z}}$$

$$y = 1 \Rightarrow p_1$$

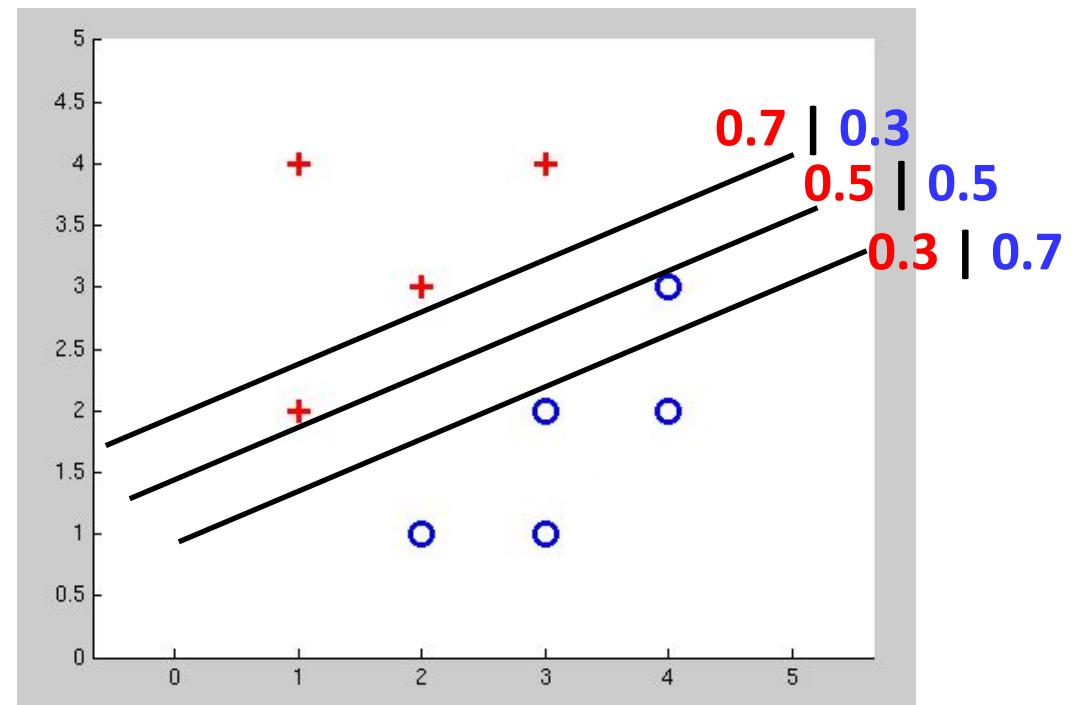
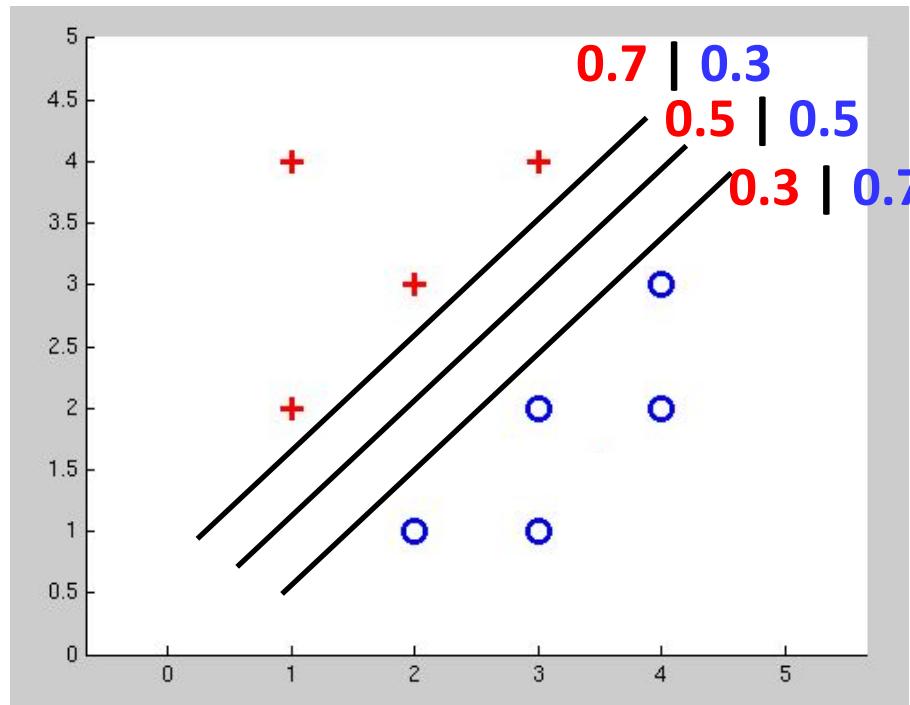
$$y = 0 \Rightarrow 1 - p_1$$

$$\text{score} = \sum p$$

## Separable Case: Deterministic Decision – Many Options



## Separable Case: Probabilistic Decision – Clear Preference



$$\omega_{SIP} \quad \omega_{poli} \quad \omega_{Sports}$$

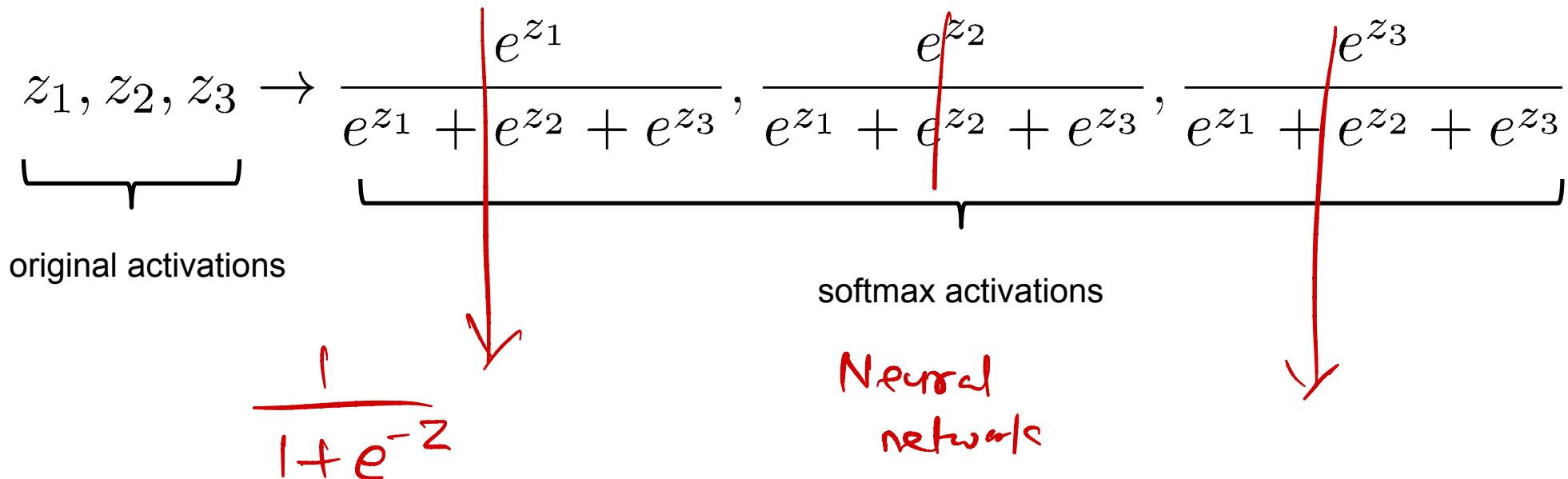
## Multiclass Logistic Regression

- Recall Perceptron:

- A weight vector for each class:  $w_y$
- Score (activation) of a class  $y$ :  $w_y \cdot f(x)$

- Prediction highest score wins  $y = \arg \max_y w_y \cdot f(x)$

- How to make the scores into probabilities?



Best w?

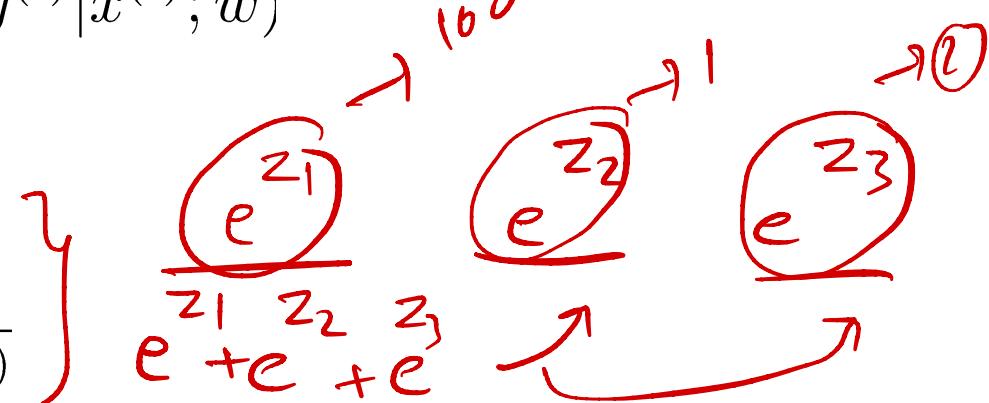
$$\frac{P_1}{1+e^{-z}} \quad 1 - \frac{1}{1+e^{-z}}$$

- Maximum likelihood estimation:

$$\max_w ll(w) = \max_w \sum_i \log P(y^{(i)}|x^{(i)}; w)$$

with:

$$P(y^{(i)}|x^{(i)}; w) = \left[ \frac{e^{w_y \cdot f(x^{(i)})}}{\sum_y e^{w_y \cdot f(x^{(i)})}} \right] y$$



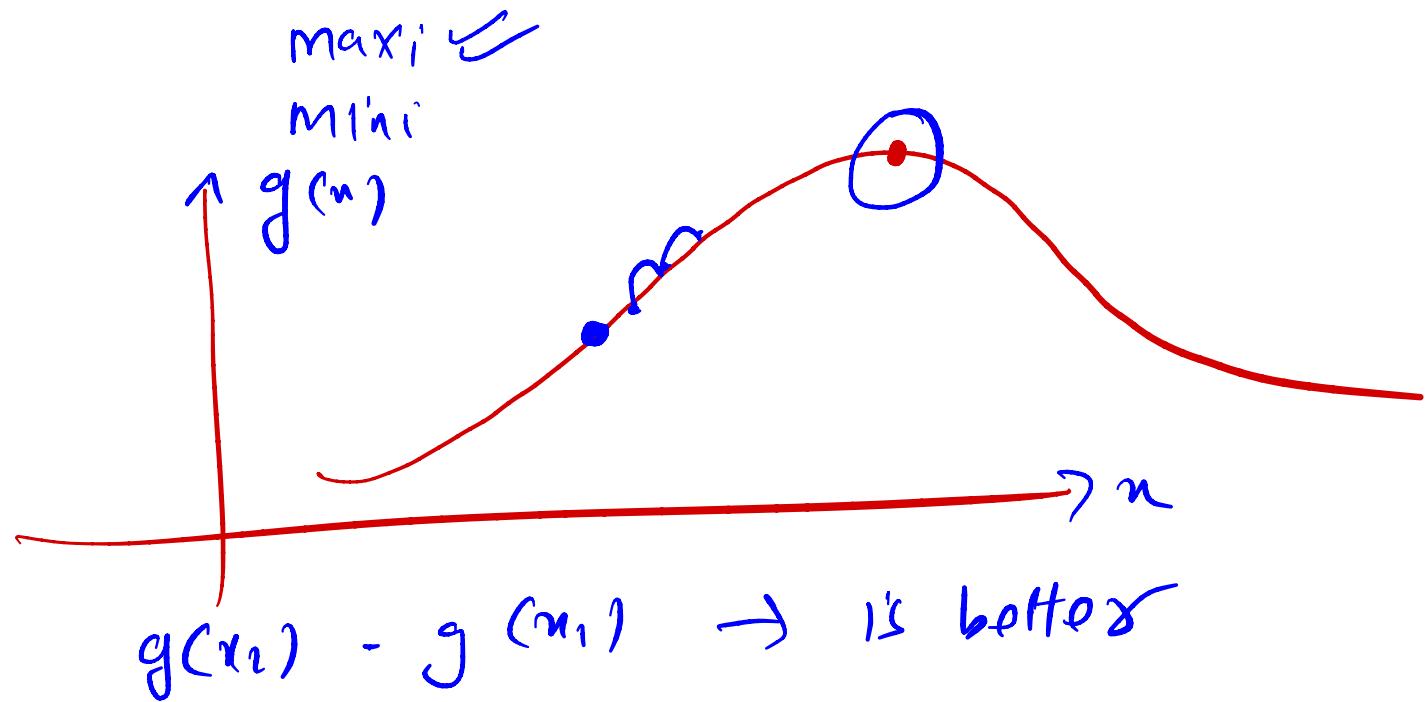
= Multi-Class Logistic Regression

Choosing of  $\omega$  ?.

loop  $\rightarrow$  ite  $\omega$   
 $\omega \rightarrow 0.1 - 100$

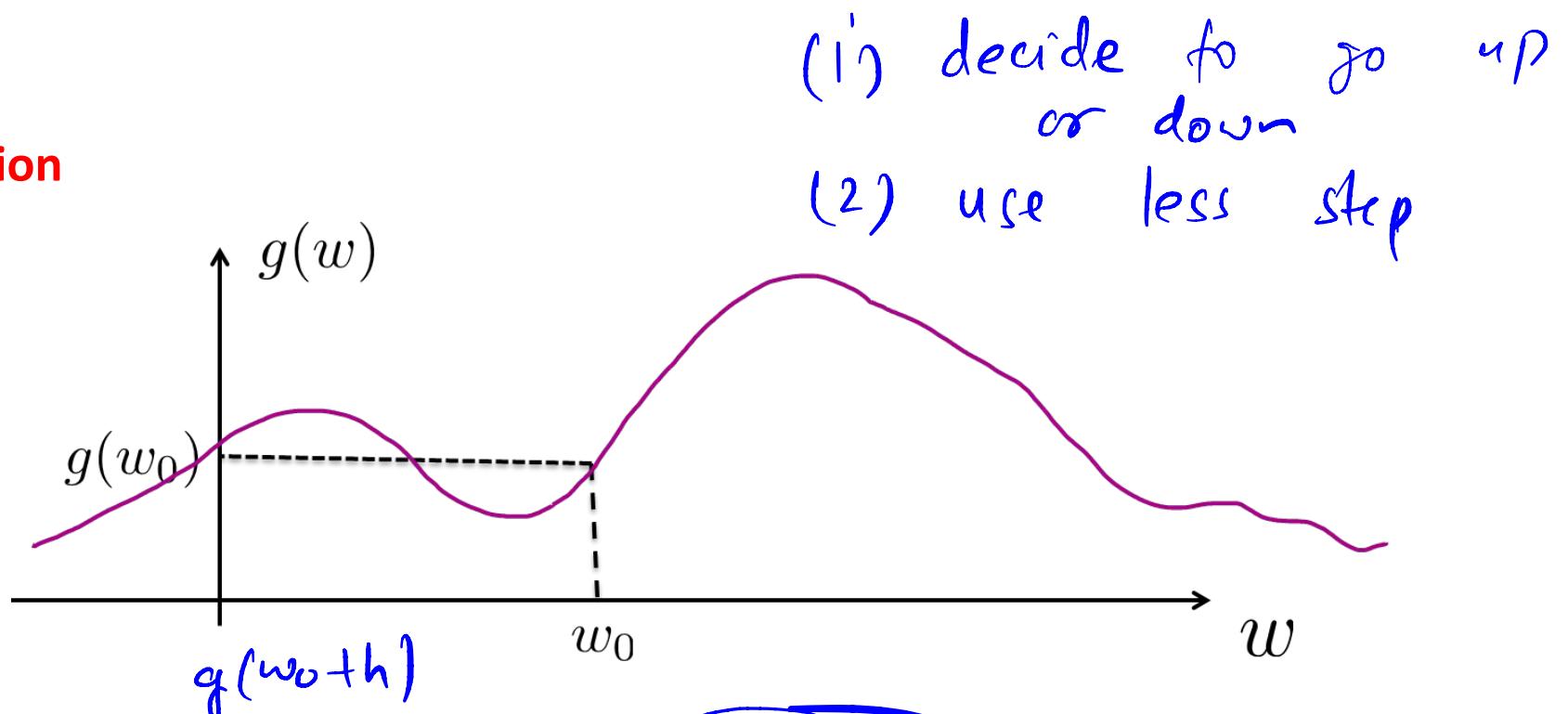
$\omega \rightarrow -5000$  to  $5000$   
 $0.01$

## Hill Climbing

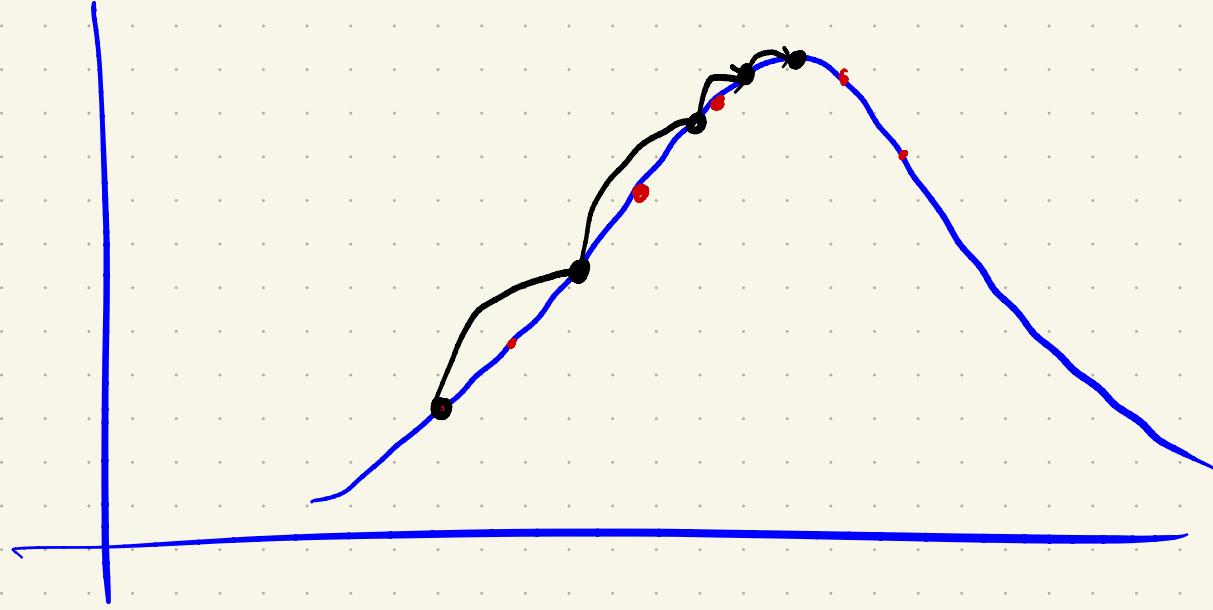


gradient

## 1-D Optimization

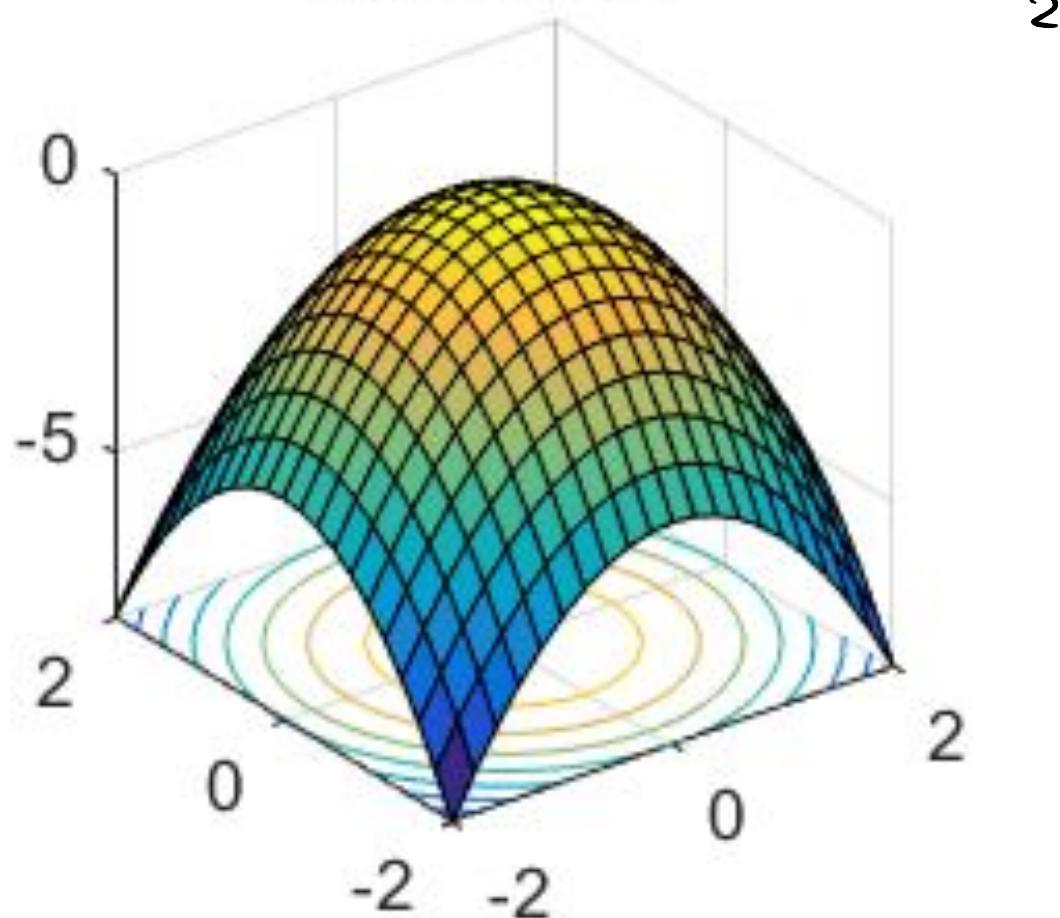


- Could evaluate  $g(w_0 + h)$ 
  - Then step in best direction
- Or, evaluate derivative:  $\frac{\partial g(w_0)}{\partial w} = \lim_{h \rightarrow 0} \frac{g(w_0 + h) - g(w_0 - h)}{2h}$ 
  - Tells which direction to step into



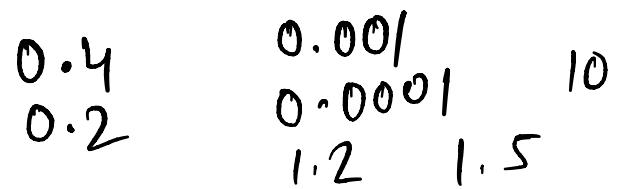
## 2-D Optimization

Rishita



$n=1$

$20 \times 20 \Rightarrow \underline{\underline{400}}$



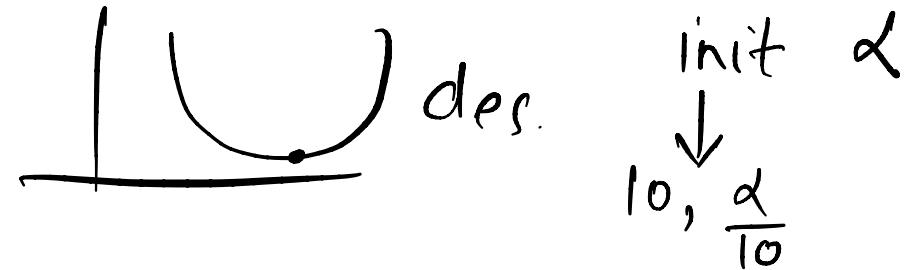
## Gradient Ascent

- Perform update in uphill direction for each coordinate
- The steeper the slope (i.e. the higher the derivative) the bigger the step for that coordinate
- E.g., consider:
  - Updates:

$$g(w_1, w_2) \xrightarrow{f(x_1)} f(x_2)$$

$$w_1 \xleftarrow{\text{new}} w_1 \xrightarrow{\text{old}} + \alpha * \frac{\partial g}{\partial w_1}(w_1, w_2)$$

$$w_2 \xleftarrow{\text{new}} w_2 \xrightarrow{\text{old}} + \alpha * \frac{\partial g}{\partial w_2}(w_1, w_2)$$



$\alpha \rightarrow$  learning rate  
 $\eta \rightarrow$

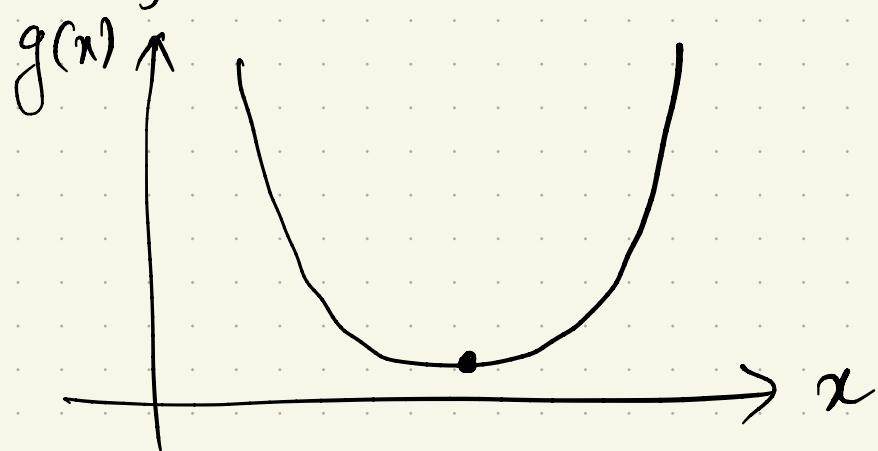
- Updates in vector notation:

$$w \leftarrow w + \alpha * \nabla_w g(w)$$

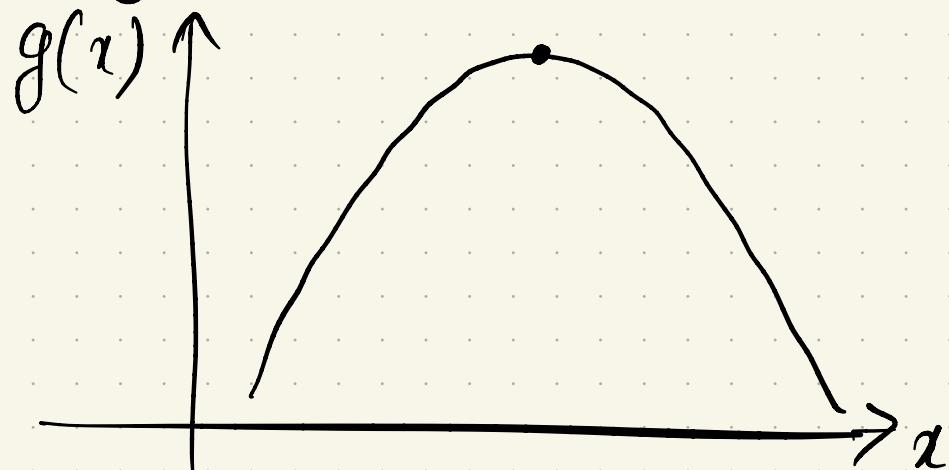
with:  $\nabla_w g(w) = \begin{bmatrix} \frac{\partial g}{\partial w_1}(w) \\ \frac{\partial g}{\partial w_2}(w) \end{bmatrix}$  = gradient

For **gradient descent**, the sign changes to  $-$  in weight update equation.

gradient descent



gradient ascent



## Gradient Ascent

- Idea:
  - Start somewhere
  - Repeat: Take a step in the gradient direction

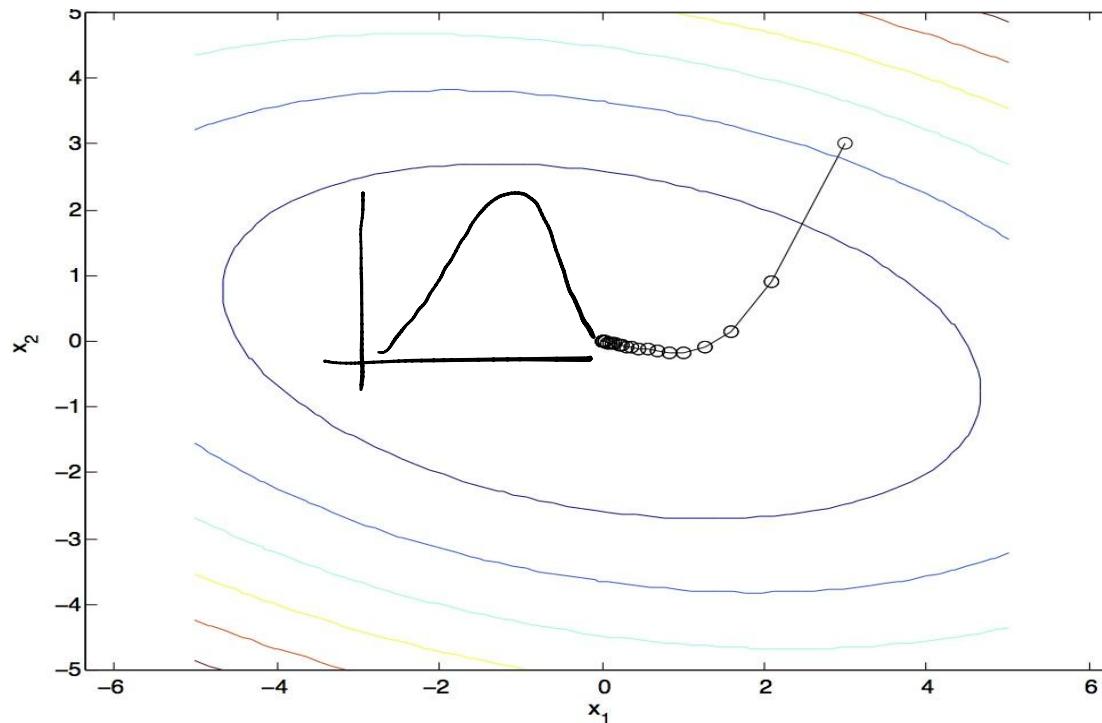
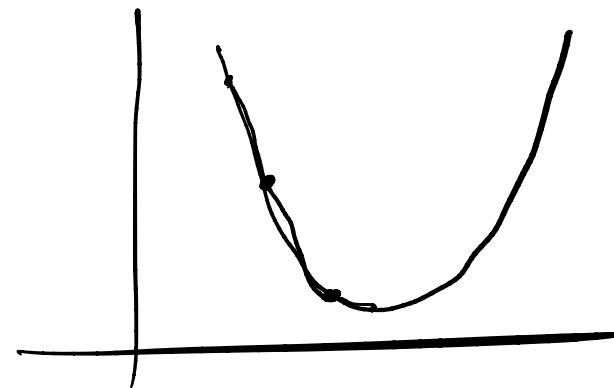


Figure source: Mathworks

## Gradient in n dimensions

n features

$$Z = x_0 w_0 + x_1 w_1 + x_2 w_2 + \dots + x_n w_n$$

$$\nabla g = \begin{bmatrix} \frac{\partial g}{\partial w_1} \\ \frac{\partial g}{\partial w_2} \\ \dots \\ \frac{\partial g}{\partial w_n} \end{bmatrix}$$

$$\frac{\partial g}{\partial w_1} \rightarrow$$

$$\frac{\partial g}{\partial w_2} \rightarrow$$

.

$$\frac{\partial g}{\partial w_n} \rightarrow$$

$$w \leftarrow w + \alpha \cdot \frac{\partial g}{\partial w}$$

## Optimization Procedure: Gradient Ascent

batch gradient  
stochastic gradient  
mini-batch gradient

- init  $w$
- for iter = 1, 2, ...

$$w \leftarrow w + \alpha * \nabla g(w)$$

- $\alpha$ : learning rate --- tweaking parameter that needs to be chosen carefully
- How? Try multiple choices
  - Crude rule of thumb: update changes  $w$  about 0.1 – 1 %

$$\min_{\tau} - \text{batch} = 10$$

$$(\tau = 10)$$

samp = 100 → each → sto.  
 → once → batch  
 → 100 times → mini batch

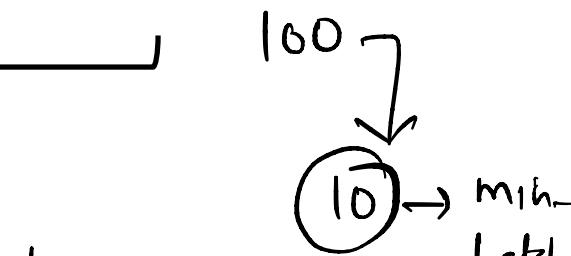
### Batch Gradient Ascent on the Log Likelihood Objective

$$\frac{5000}{100} \rightarrow 500 \times 500 = 250 \text{ min}$$

$$\frac{10}{10} = \underline{\underline{y_0}} = \underline{\underline{y_1}}$$

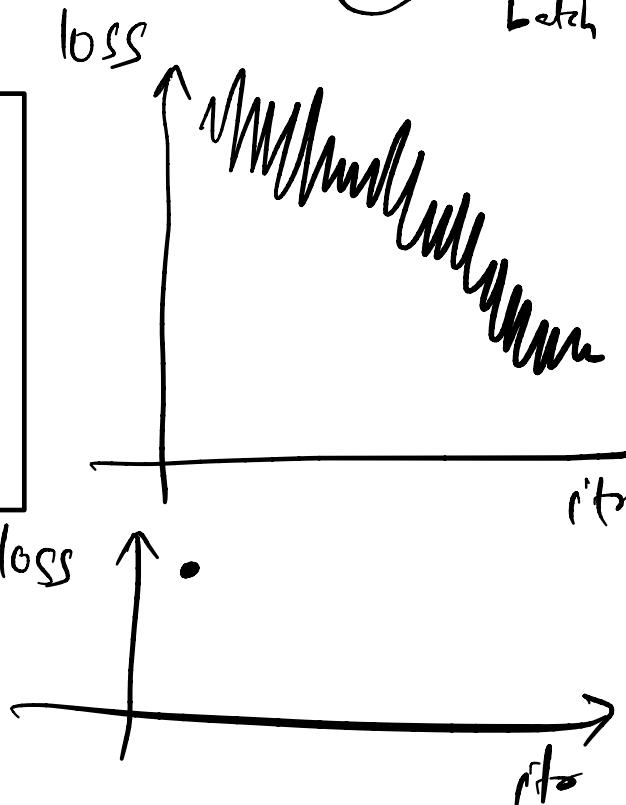
$$\max_w ll(w) = \max_w \sum_i \log P(y^{(i)} | x^{(i)}; w)$$

$$g(w)$$



- init  $w$
- for iter = 1, 2, ...

$$w \leftarrow w + \alpha * \sum_i \nabla \log P(y^{(i)} | x^{(i)}; w)$$



## Stochastic Gradient Ascent on the Log Likelihood Objective

$$\max_w \text{ll}(w) = \max_w \sum_i \log P(y^{(i)} | x^{(i)}; w)$$

**Observation:** once gradient on one training example has been computed, might as well incorporate before computing next one

- init  $w$
- for iter = 1, 2, ...
  - pick random  $j$

$$w \leftarrow w + \alpha * \nabla \log P(y^{(j)} | x^{(j)}; w)$$

neural networks sk learn  
(MLP)

parameter.  $\rightarrow \omega, b$   
hyperpara.  $\rightarrow \alpha, \beta, \gamma$ , min-batch,  
epi.

## Mini-Batch Gradient Ascent on the Log Likelihood Objective

$$\max_w ll(w) = \max_w \sum_i \log P(y^{(i)} | x^{(i)}; w)$$

**Observation:** gradient over small set of training examples (=mini-batch)  
can be computed in parallel, might as well do that instead of a single one

- init  $w$
- for iter = 1, 2, ...
  - pick random subset of training examples  $J$

$$w \leftarrow w + \alpha * \sum_{j \in J} \nabla \log P(y^{(j)} | x^{(j)}; w)$$

## Mini-Batch Gradient Ascent on the Log Likelihood Objective

$$\max_w \text{ll}(w) = \max_w \sum_i \log P(y^{(i)} | x^{(i)}; w)$$

**Observation:** gradient over small set of training examples (=mini-batch) can be computed in parallel, might as well do that instead of a single one

- init  $w$
- for iter = 1, 2, ...
  - pick random subset of training examples  $J$

$$w \leftarrow w + \alpha * \sum_{j \in J} \nabla \log P(y^{(j)} | x^{(j)}; w)$$

$$x_0 = 1 \quad x \rightarrow x_1, x_2$$

{ logistic regression } (gradient descent)

$$z = x_0 w_0 + x_1 w_1 + x_2 w_2$$

$$A = \frac{1}{1 + e^{-z}}$$

$$\text{loss} = - \left[ y \log A + (1-y) \log (1-A) \right] \Rightarrow \text{binary loss}$$

$$\frac{\partial \text{loss}}{\partial w} = \frac{\partial \text{loss}}{\partial A} \cdot \frac{\partial A}{\partial z} \cdot \frac{\partial z}{\partial w}$$

$$w_1 \leftarrow w_1 - \alpha ( )$$

$$w_2 \leftarrow w_2 - \alpha ( )$$

$$w_0 \leftarrow w_0 - \alpha ( )$$

$$\begin{array}{c} \frac{\partial \text{loss}}{\partial w_1} \\ \frac{\partial \text{loss}}{\partial w_2} \\ \frac{\partial \text{loss}}{\partial w_0} \end{array}$$