

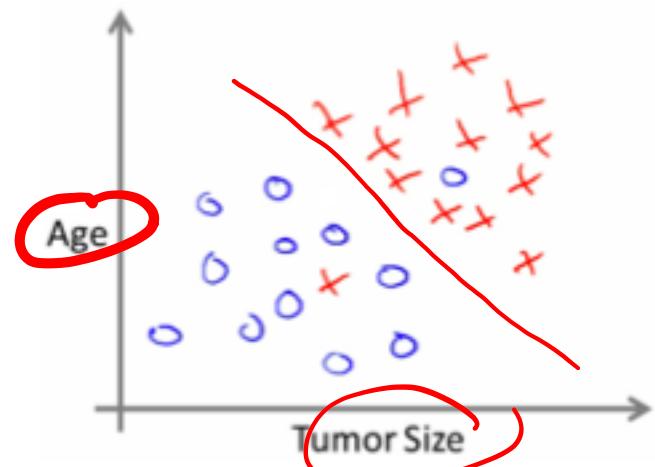
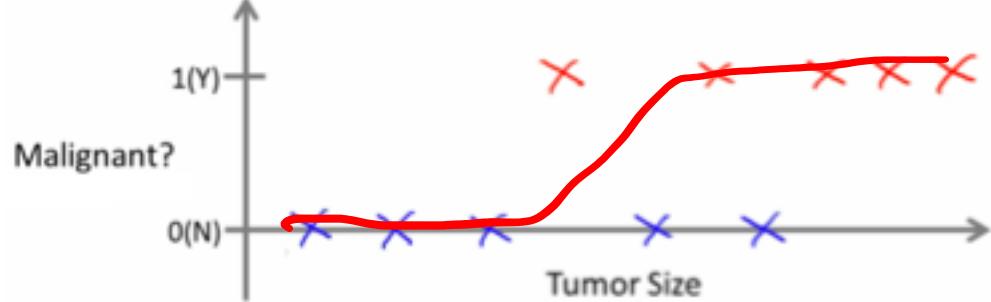
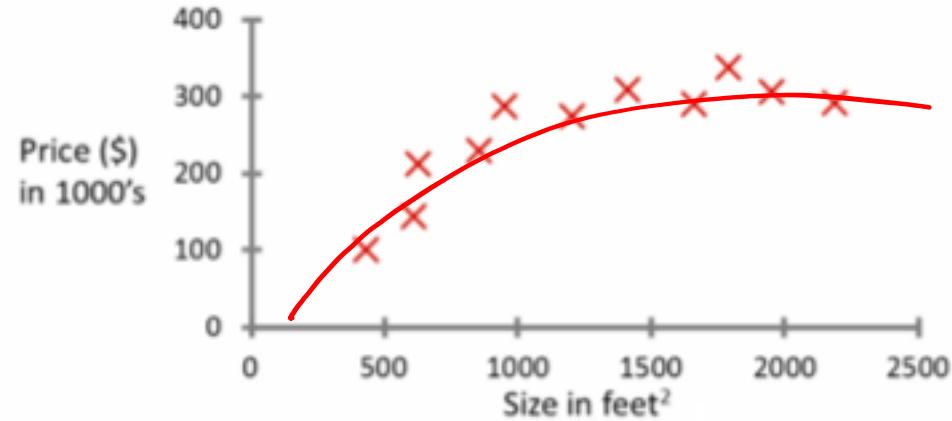
深度学习

第三讲

王胤

Recap: Regression vs. Classification

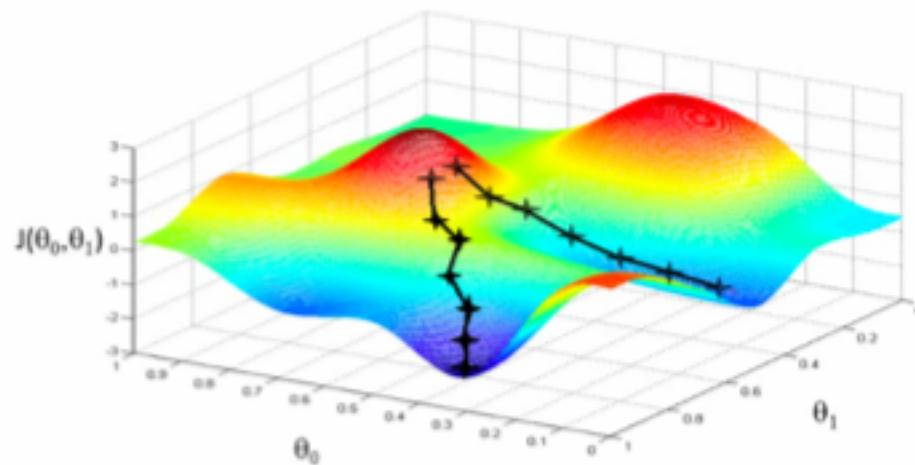
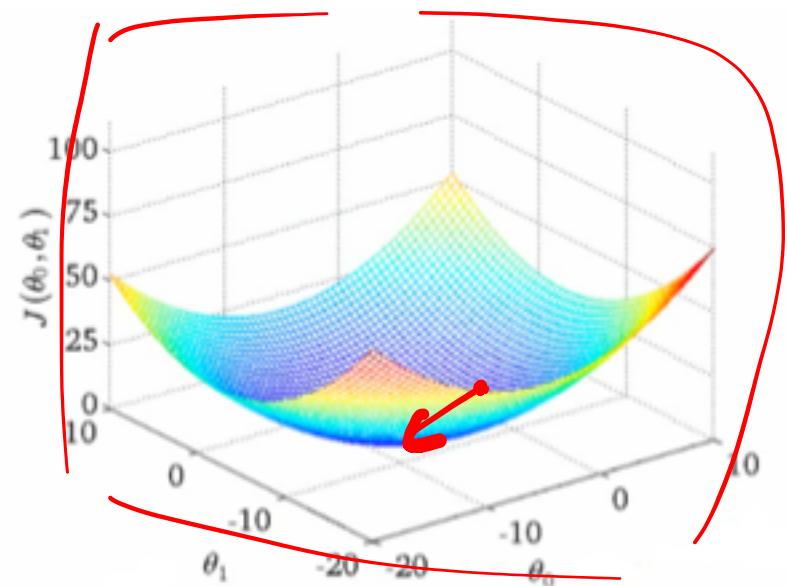
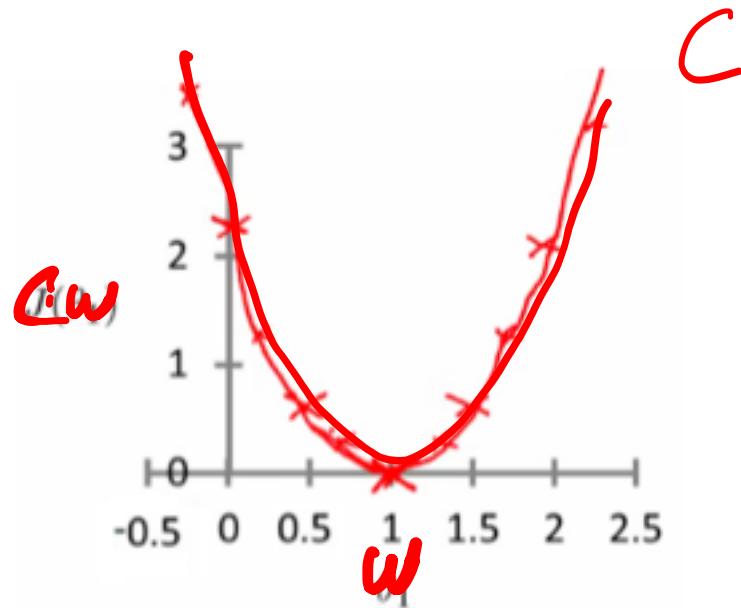
Housing price prediction.



Recap: Linear Regression

- 训练数据是 m 个 (x^i, t^i) 数据对， x^i 是 n 纬向量
- $y_{\omega}(x) = \omega^T x = \omega_0 + \omega_1 x_1 + \dots + \omega_n x_n$
- $C(\omega) = \frac{1}{2m} \sum_i (y_i - t_i)^2$
- 目标？ $\min C(\omega)$

Recap: 1D case



Recap: Gradient Descent

- $\omega_j := \omega_j - \alpha \nabla \omega$

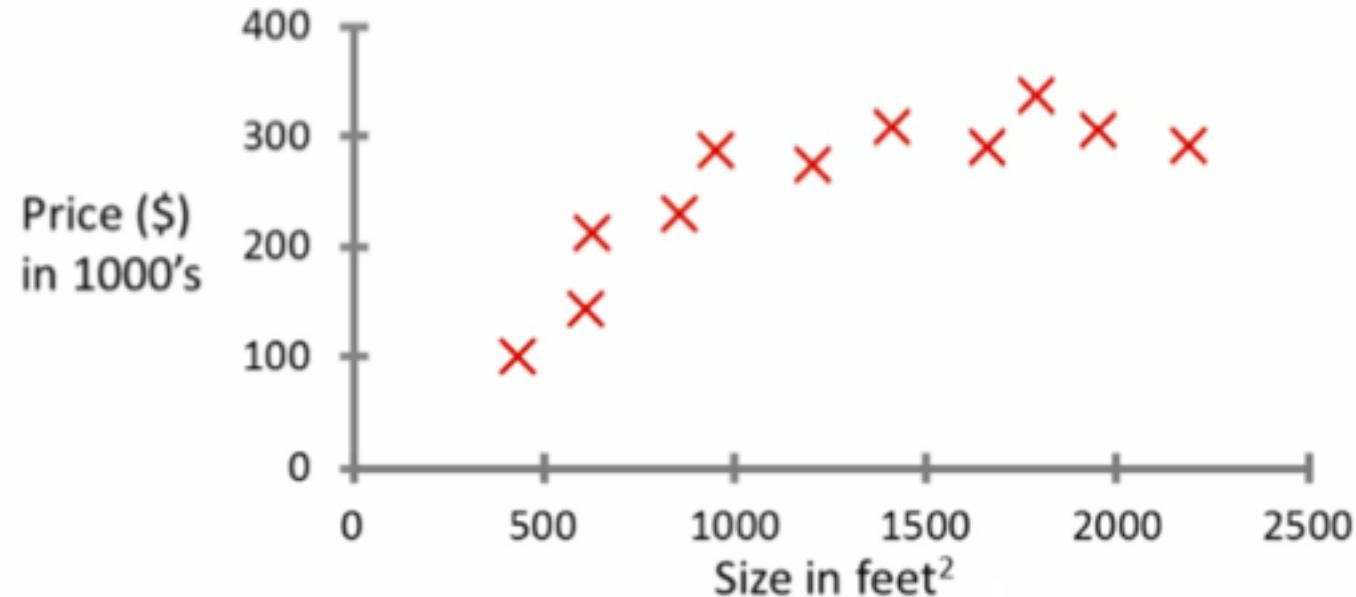
$$\omega_j = \omega_j - \alpha \sum_i (y_i - t_i) \cdot X_{ij}$$

Diagram illustrating the gradient descent update rule:

- The equation shows the update step: $\omega_j = \omega_j - \alpha \sum_i (y_i - t_i) \cdot X_{ij}$.
- The term $\sum_i (y_i - t_i)$ is highlighted with a blue oval.
- The term $\cdot X_{ij}$ is highlighted with a blue oval.
- The variables are labeled:
 - w_1, w_2, w_n above the weights $\omega_1, \omega_2, \omega_n$.
 - t_1, t_2, t_m below the targets y_1, y_2, y_m .
 - f_1, f_2, f_n below the function values y_1, y_2, y_m .
 - x below the feature vector X .
 - $y_i - t_i$ to the right of the error term.
- Below the diagram, examples are listed:
 - ex 1
 - ex 2
 - ex 3
 - ⋮
 - ex m

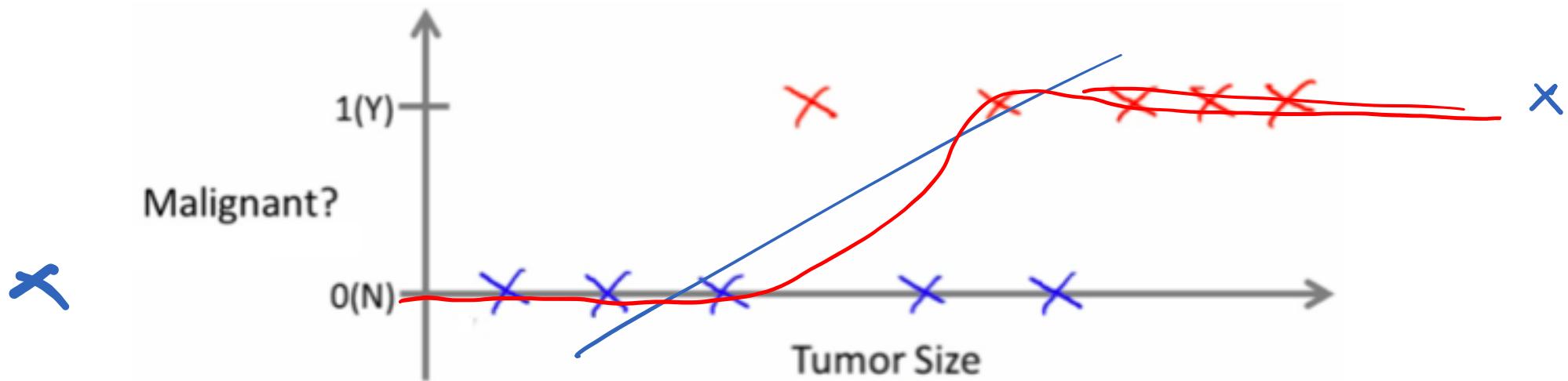
Recap: Polynomial Regression

Housing price prediction.



$$x_1 \quad x_1^2 \quad x_1^3 \quad \dots$$

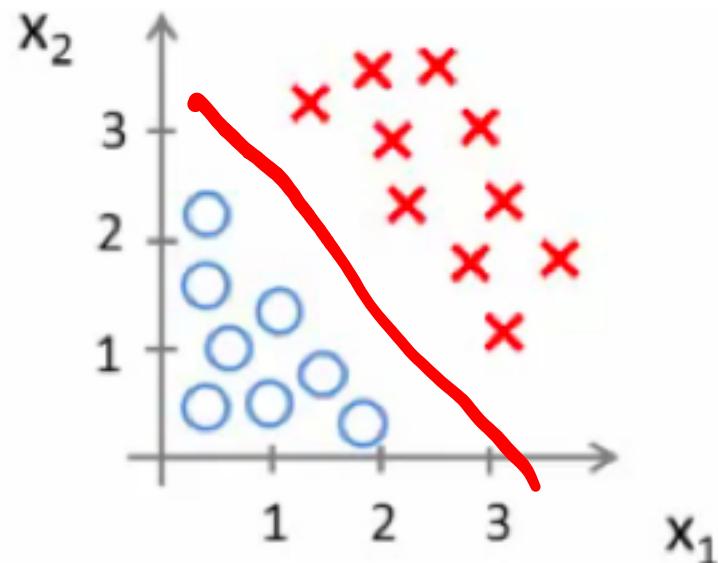
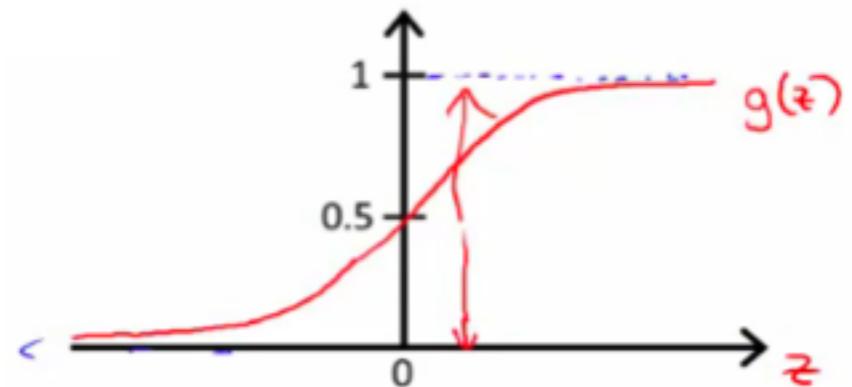
Recap: Logistic Regression



Recap: Logistic Function

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

$$y_{\omega}(x) = \frac{1}{1 + e^{-\omega^T x}}$$



Recap: Logistic Regression

- 训练数据是 m 个 (x^i, y^i) 数据对， x^i 是 n 纬向量

$$y_{\omega}(x) = \frac{1}{1 + e^{-\omega^T x}}$$

$$\begin{aligned} C(\omega) &= -\frac{1}{m} \sum_i \left(t_i \log y_{\omega}(x_i) + (1-t_i) \log (1-y_{\omega}(x_i)) \right) \\ &= -\frac{1}{m} \sum_i t_i \log y_i \end{aligned}$$

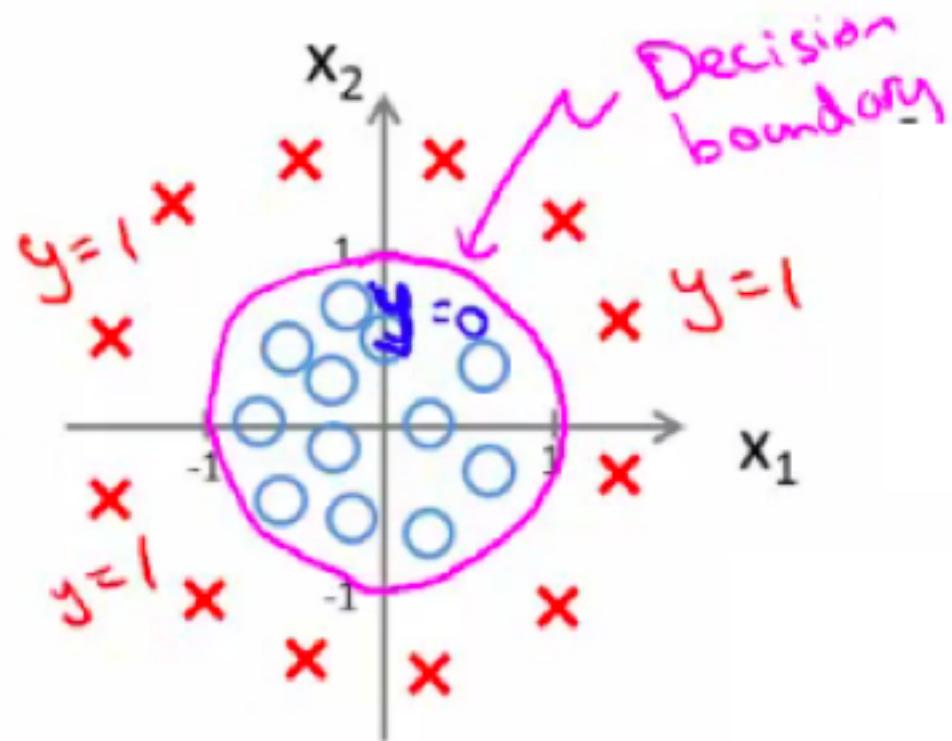
目标？

$$\min_{\omega} C(\omega)$$

Recap: Gradient Descent

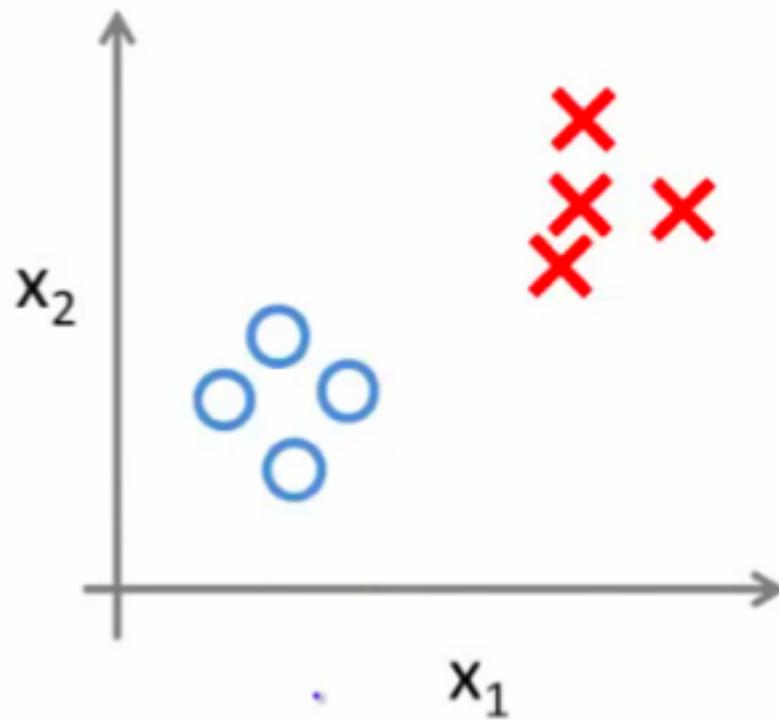
- $w_j := w_j - \alpha \Delta w_j$
 $= w_j - \alpha \sum_i \underset{T}{\overbrace{(y_i - t_i)}} x_{ij}$

Recap: Non-Linear

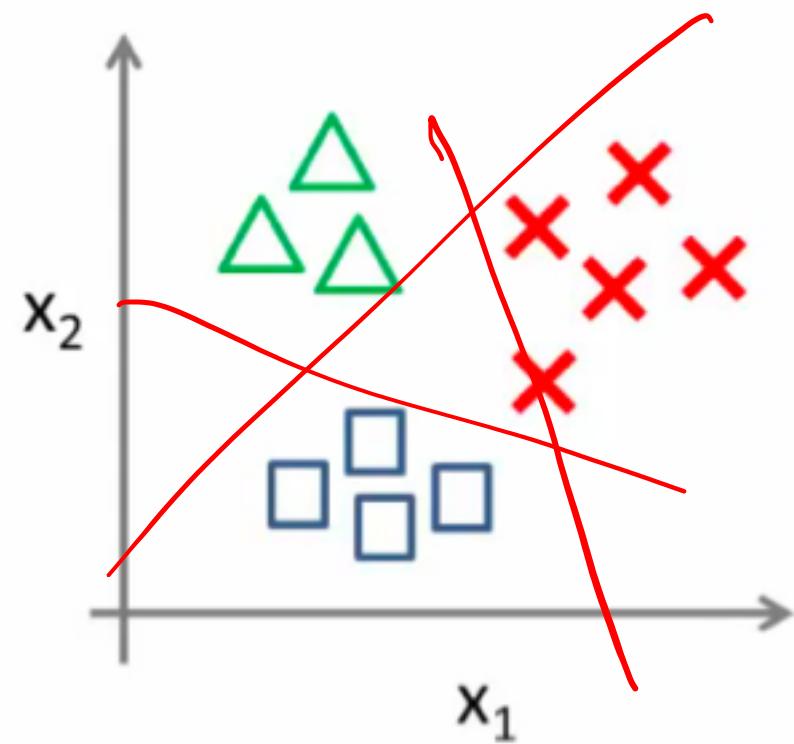


Recap: Multiclass Classification

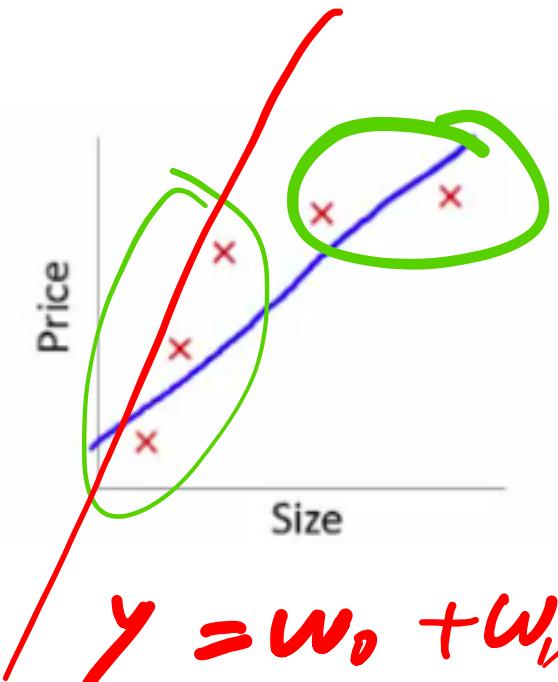
Binary classification:



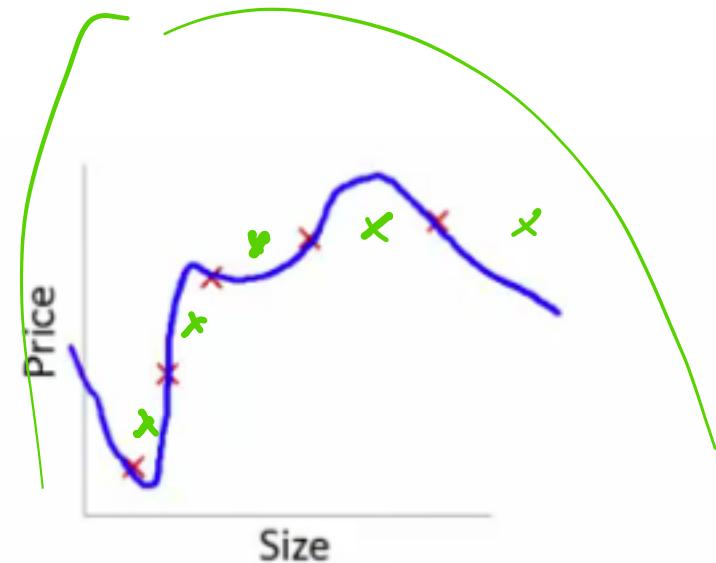
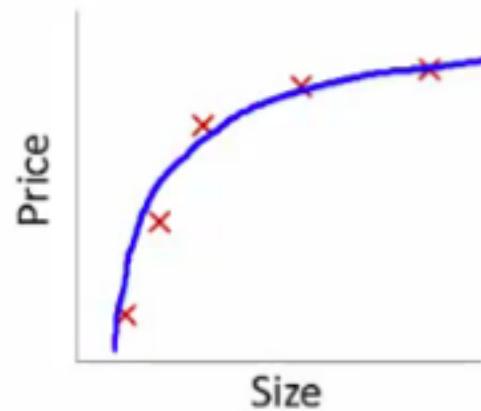
Multi-class classification:



Overfitting: Linear Regression

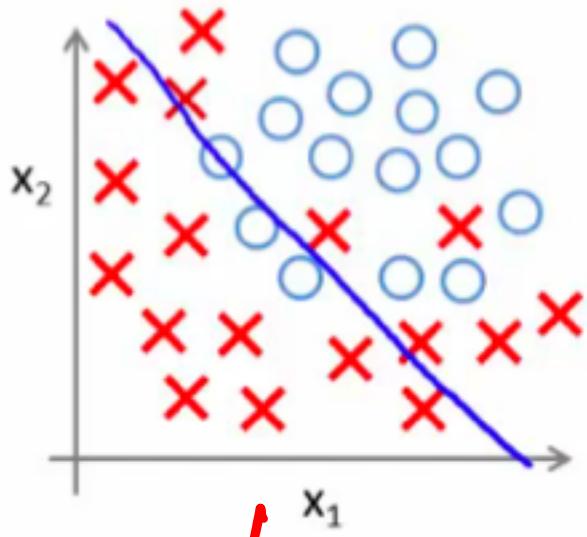


s-fold CV

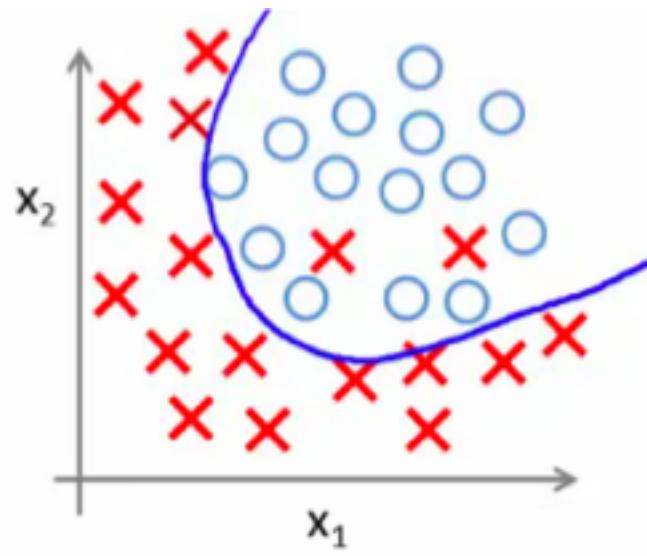


$y = \dots - w_3 x^3 \dots - w_n x^n$

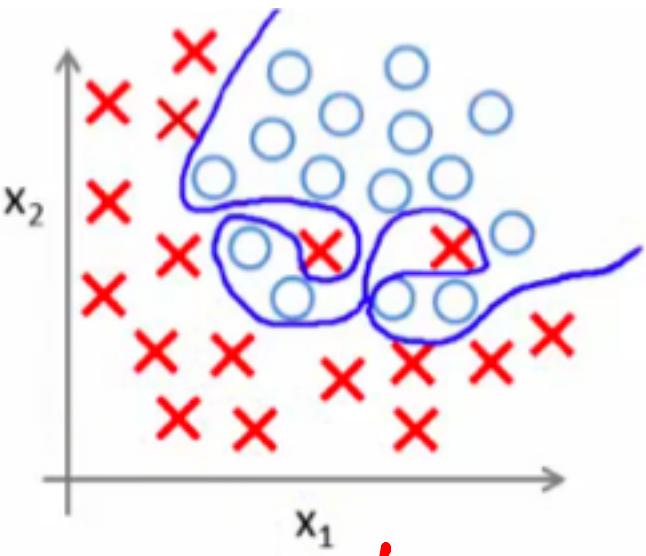
Overfitting: Logistic Regression



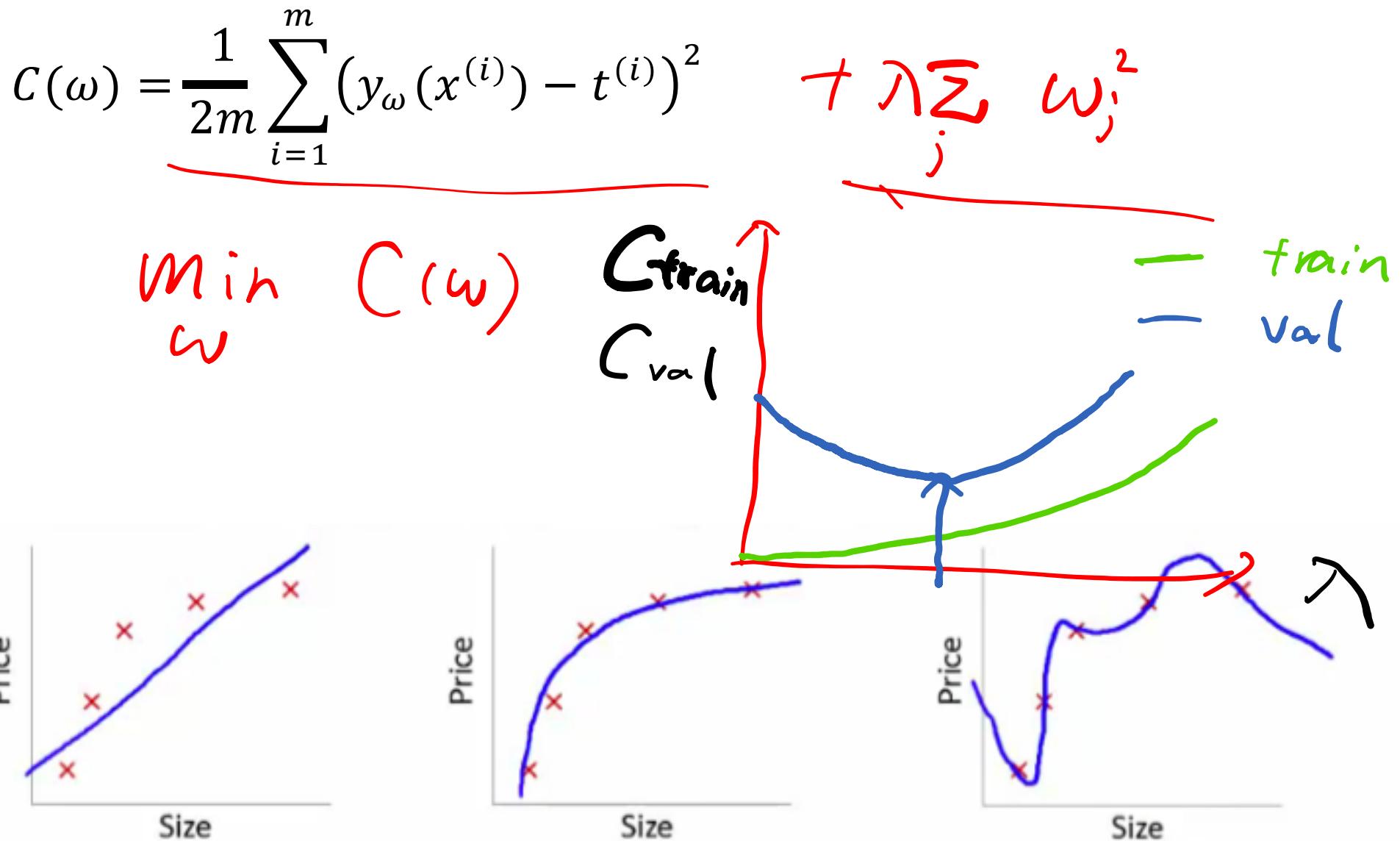
$$f + e^{-(w_0 + w_1 x_1 + w_2 x_2)}$$



$$\begin{aligned} & \downarrow \\ & w_0 + w_1 x_1 + w_2 x_2 + \\ & w_3 x_1^2 + w_4 x_1 x_2 + w_5 x_2^2 \end{aligned}$$



Regularization



Regularized Linear Regression

$$C(\omega) = \frac{1}{2m} \left[\sum_{i=1}^m (y_\omega(x^{(i)}) - t^{(i)})^2 + \lambda \sum_{j=1}^n \omega_j^2 \right]$$

$$\omega_j := \underline{\omega_j} - \alpha \Delta \underline{\omega_j}$$

$$\underline{\omega_j} := \omega_j \left(1 - \alpha \frac{\lambda}{m} \right) - \alpha \frac{1}{m} \sum_{i=1}^m (y_\omega(x^{(i)}) - t^{(i)}) x_j^{(i)}$$

- Is going to be a number less than 1 usually
- Usually learning rate is small and m is large
 - So the term is often around 0.99 to 0.95

Regularized Logistic Regression

$$C(\omega) = -\frac{1}{m} \sum_{i=1}^m t^{(i)} \log(y_\omega(x^{(i)})) + (1 - t^{(i)}) \log(1 - y_\omega(x^{(i)}))$$

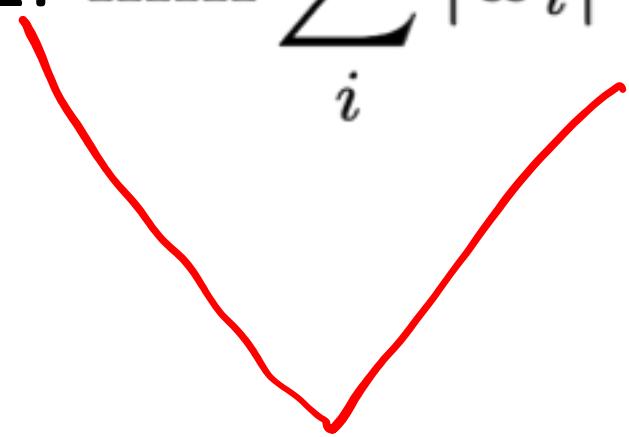
$+ \lambda \sum_j w_j^2$

$$\underline{\omega_j := \omega_j \left(1 - \alpha \frac{\lambda}{m}\right) - \alpha \frac{1}{m} \sum_{i=1}^m (y_\omega(x^{(i)}) - t^{(i)}) x_j^{(i)}}$$

Regularization

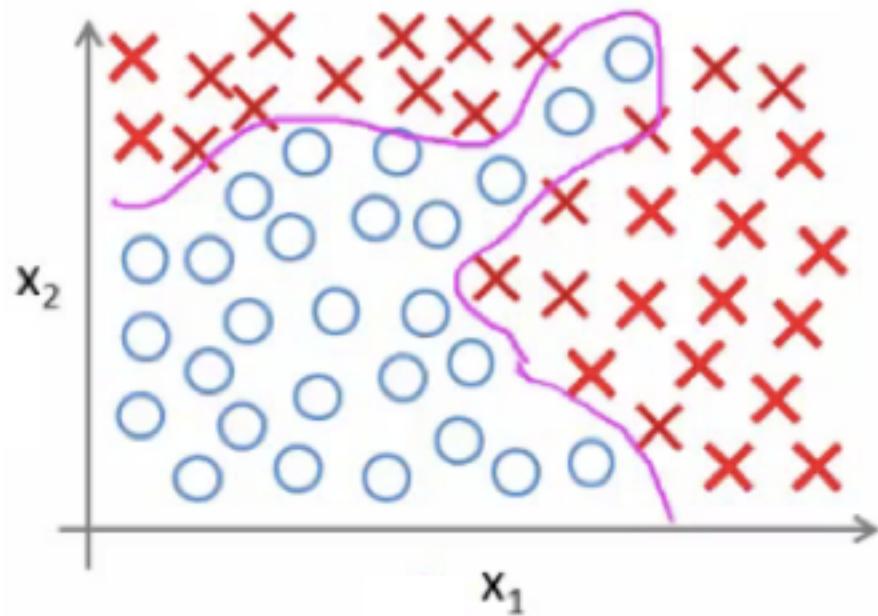
L1 vs. L2

$$\text{L1: } \min \sum_i |w_i|$$



$$\text{L2: } \min \sum_i w_i^2$$

Curse of Dimensionality



$x_1^L \quad x_2^L$
 $x_1 x_2$
 $x_1^3 \quad x_1^2 x_2 \quad - \cdot$
⋮
⋮
⋮

Computer Vision: classification vs. regression

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

Is there a cat in the image? Yes or No

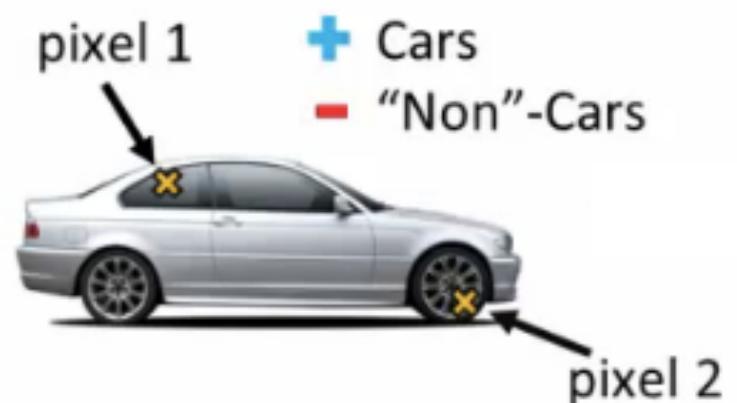
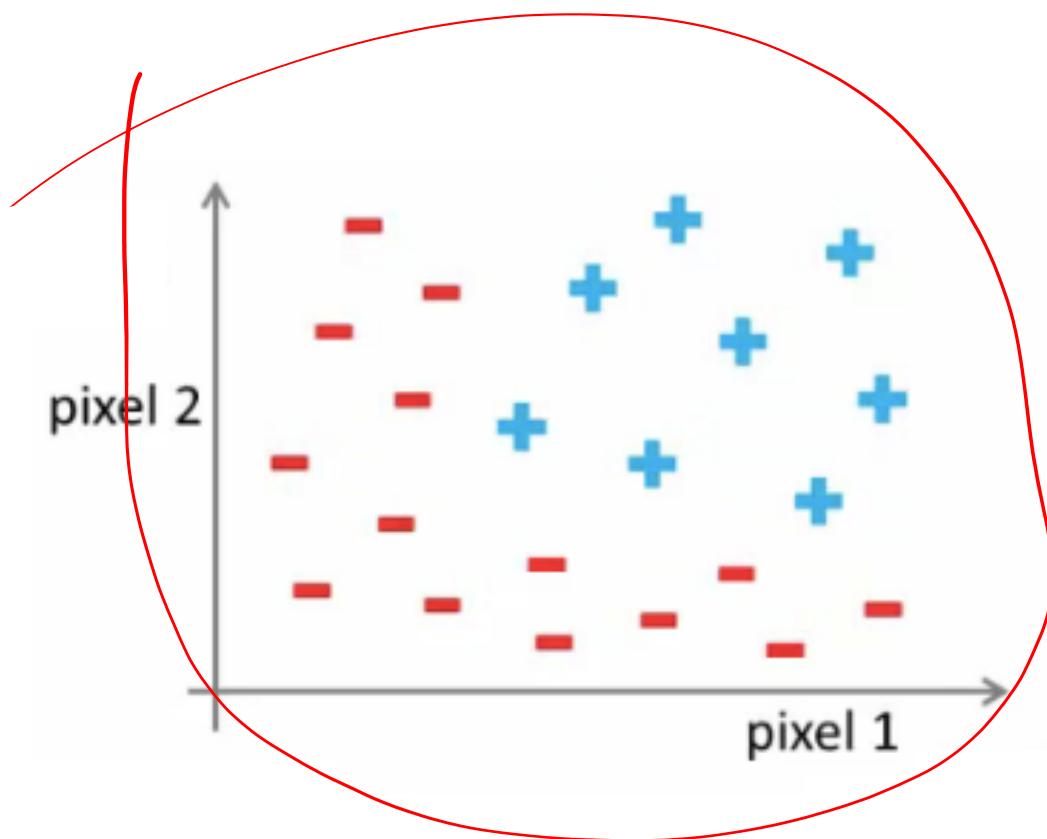


$y \in \mathbb{R}$ regression

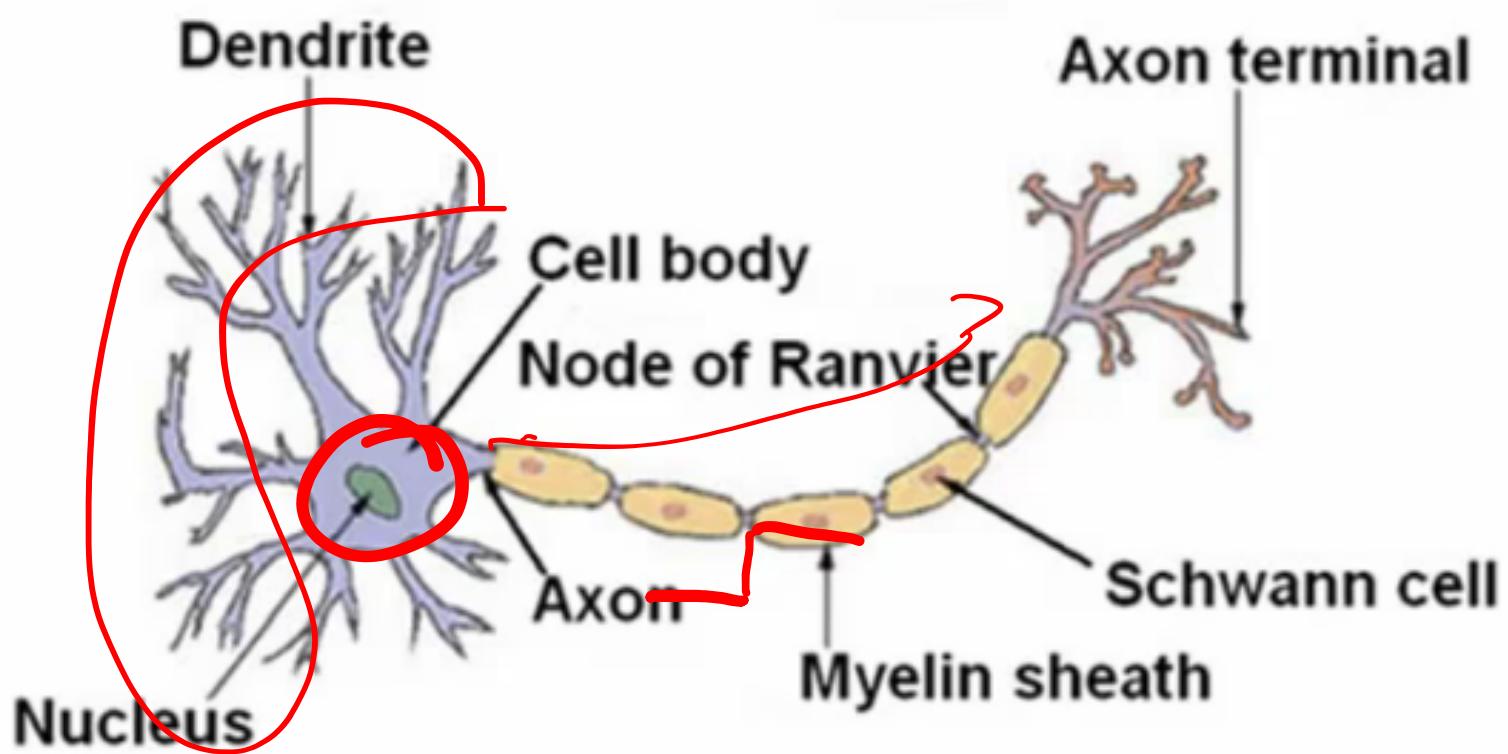
How cute is the cat (scale 1-10)? (9.4)



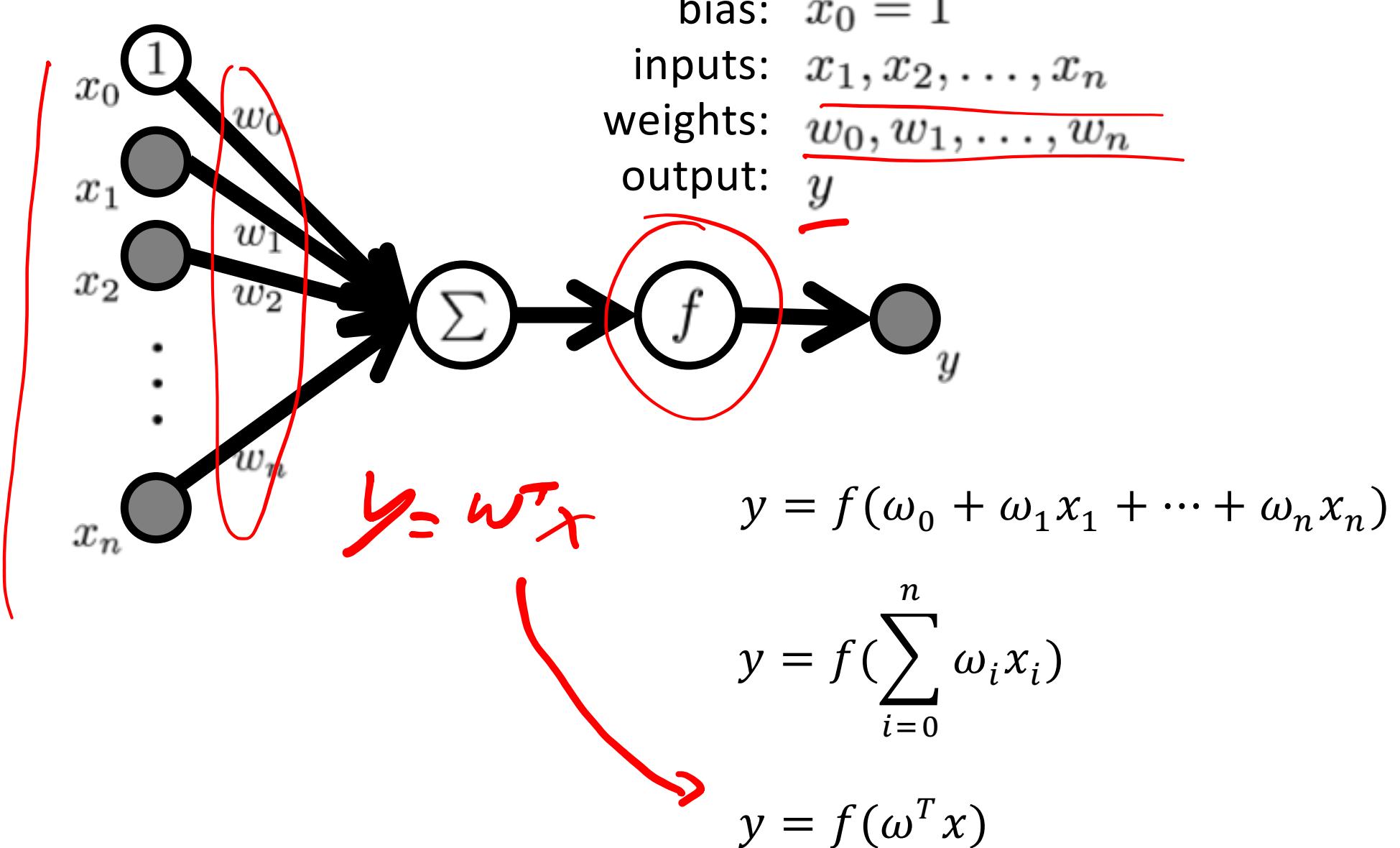
Problems where n is large: Computer Vision



Neuron



A Single Neuron



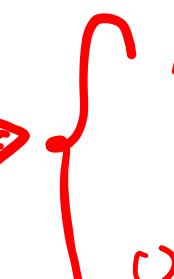
Activation function: f

$$y = f\left(\sum_{i=0}^n w_i x_i\right)$$

linear: $f(x) = x$



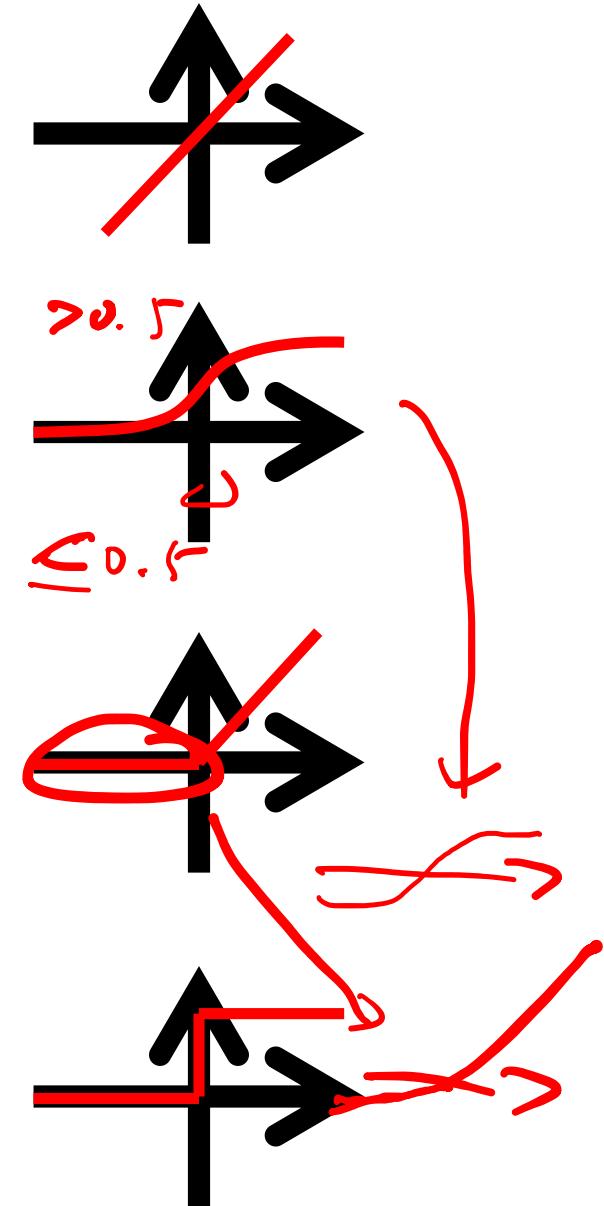
sigmoid: $f(x) = \frac{1}{1+e^{-x}}$



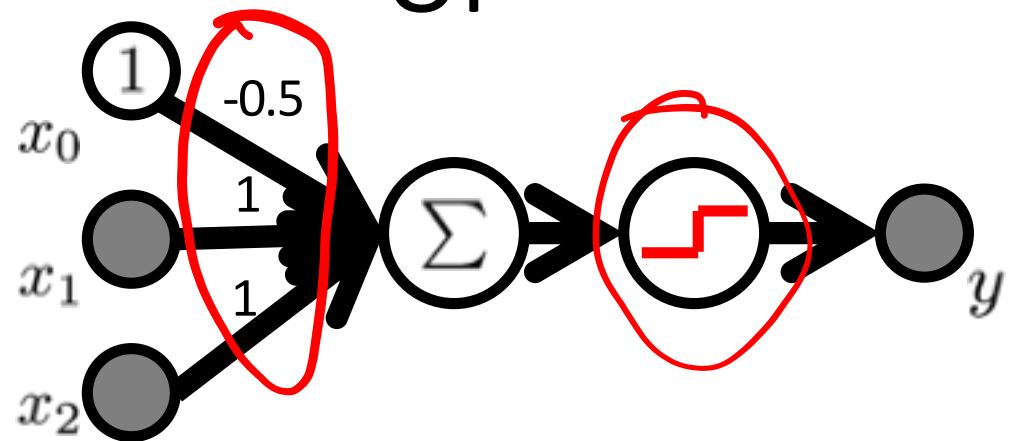
ReLU: $f(x) = \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases}$

Rectified Linear Unit

unit step: $f(x) = \begin{cases} 1 & x > 0 \\ 0 & x \leq 0 \end{cases}$

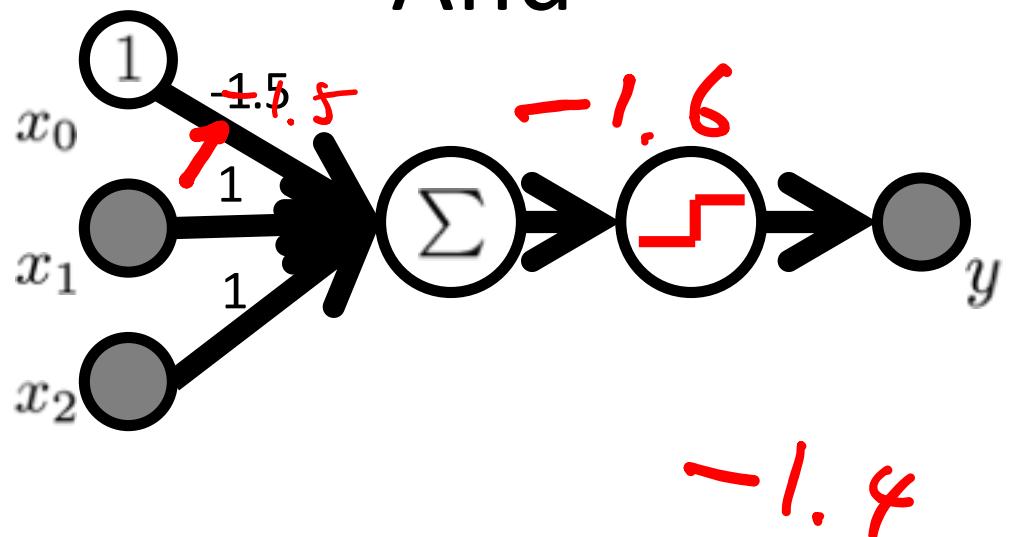


Or



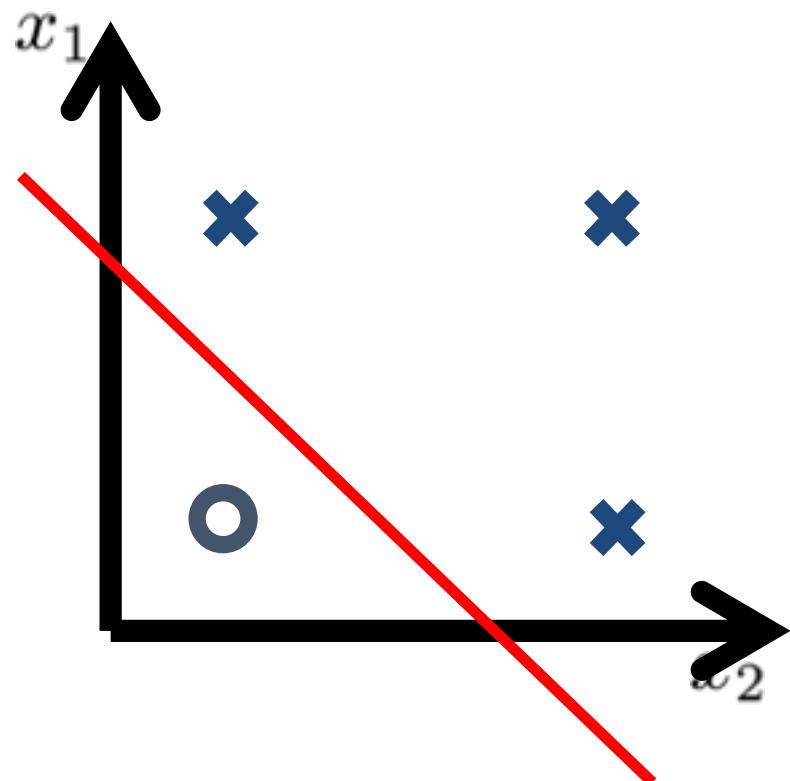
x_1	x_2	Σ	y
0	0	-0.5	0
0	1	0.5	1
1	0	0.5	1
1	1	1.5	1

And



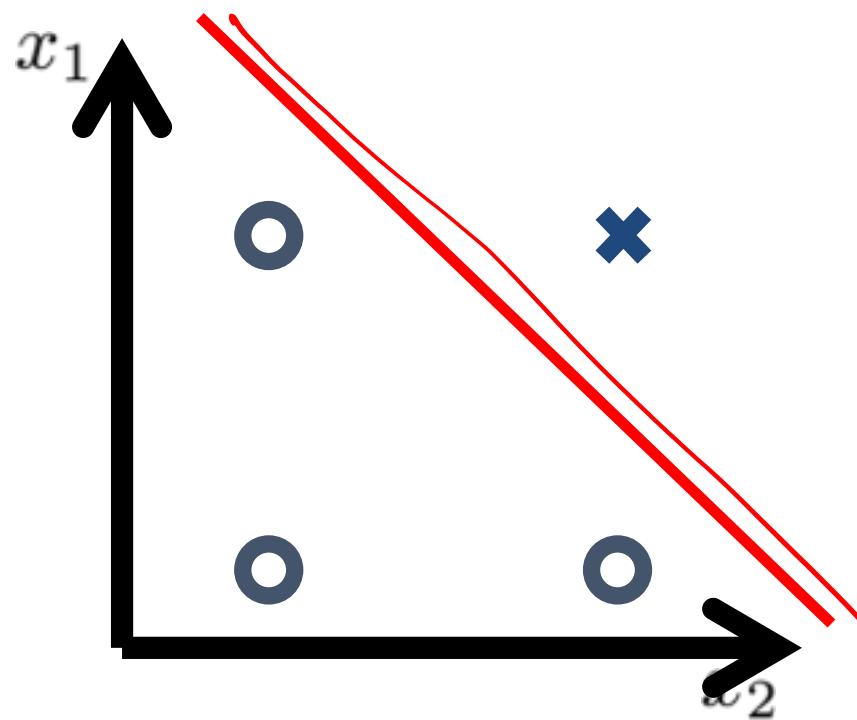
x_1	x_2	Σ	y
0	0	-1.5	0
0	1	-0.5	0
1	0	-0.5	0
1	1	0.5	1

Or



$$y = x_2 + x_1 - 0.5$$

And



$$y = x_2 + x_1 - 1.5$$

Objective function

How do we measure our error?

training example: $\{\underline{\mathbf{x}}_j, \underline{t}_j\}$

target output: \underline{t}_j

network parameters: $\underline{\omega}$

$C(\underline{\omega})$

$$\min_{\underline{\omega}} \sum_j (y_j - t_j)^2$$

Also called the error or loss function.

$$\frac{\partial E}{\partial w_i} = \frac{\partial}{\partial w_i} \frac{1}{2} \sum_j (t_j - y_j)^2$$

$$= (t_j - y_j) \sum_j \frac{\partial}{\partial w_i} (t_j - y_j)$$

$$= (t_j - y_j) \sum_j \frac{\partial}{\partial w_i} f(\sum_k w_k x_k)$$

$$= - (t_j - y_j) \sum_j f'(.) x_i$$

$$\frac{\partial E}{\partial w_i} = \frac{\partial}{\partial w_i} \frac{1}{2} \sum_j (t_j - y_j)^2$$

$$= \frac{1}{2} \sum_j \frac{\partial}{\partial w_i} (t_j - y_j)^2$$

$$= \frac{1}{2} \sum_j 2(t_j - y_j) \frac{\partial}{\partial w_i} (t_j - y_j)$$

$$= \sum_j (t_j - y_j) (-\frac{\partial}{\partial w_i} y_j)$$

$$= \sum_j (t_j - y_j) (-f'(net_j) x_{ij})$$

$$= - \sum_j (t_j - y_j) f'(net_j) x_{ij}$$

Function derivatives

Derivatives

linear: $f(x) = x$

$$f'(x) = 1$$

sigmoid: $f(x) = \frac{1}{1+e^{-x}}$

$$f'(x) = f(x)(1 - f(x))$$

ReLU: $f(x) = \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases}$

$$f'(x) = \begin{cases} 1 & x \geq 0 \\ 0 & x < 0 \end{cases}$$

unit step: $f(x) = \begin{cases} 1 & x > 0 \\ 0 & x \leq 0 \end{cases}$

$$f'(x) = 0 \quad ?$$

Weight update

new weight



old weight



increment



$$w_i \leftarrow w_i + \Delta w_i$$

step size



derivative of activation function



$$\Delta w_i = \alpha(t_j - y_j) f'(net_j) \underline{x_{ij}}$$



target



output

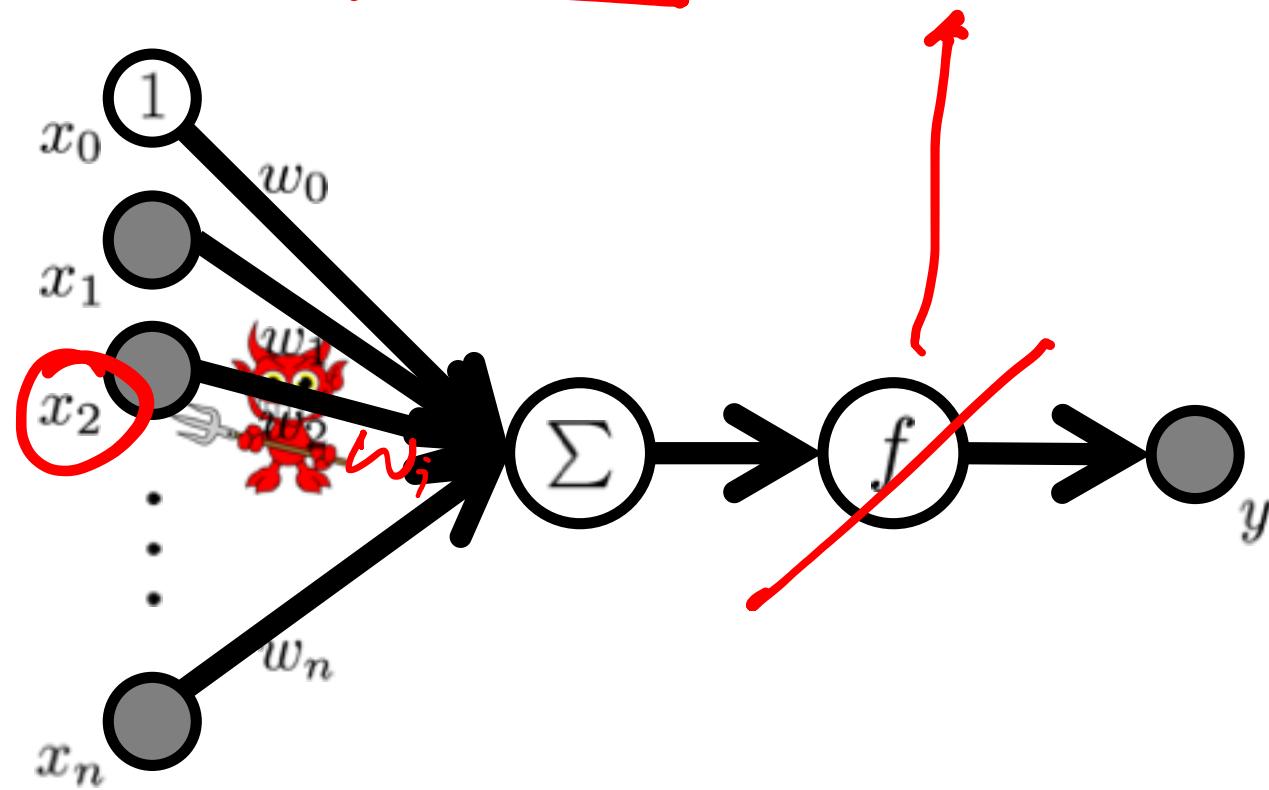


input

$$net_j = \sum_{i=0}^n w_i x_{ij}$$

Reducing the error

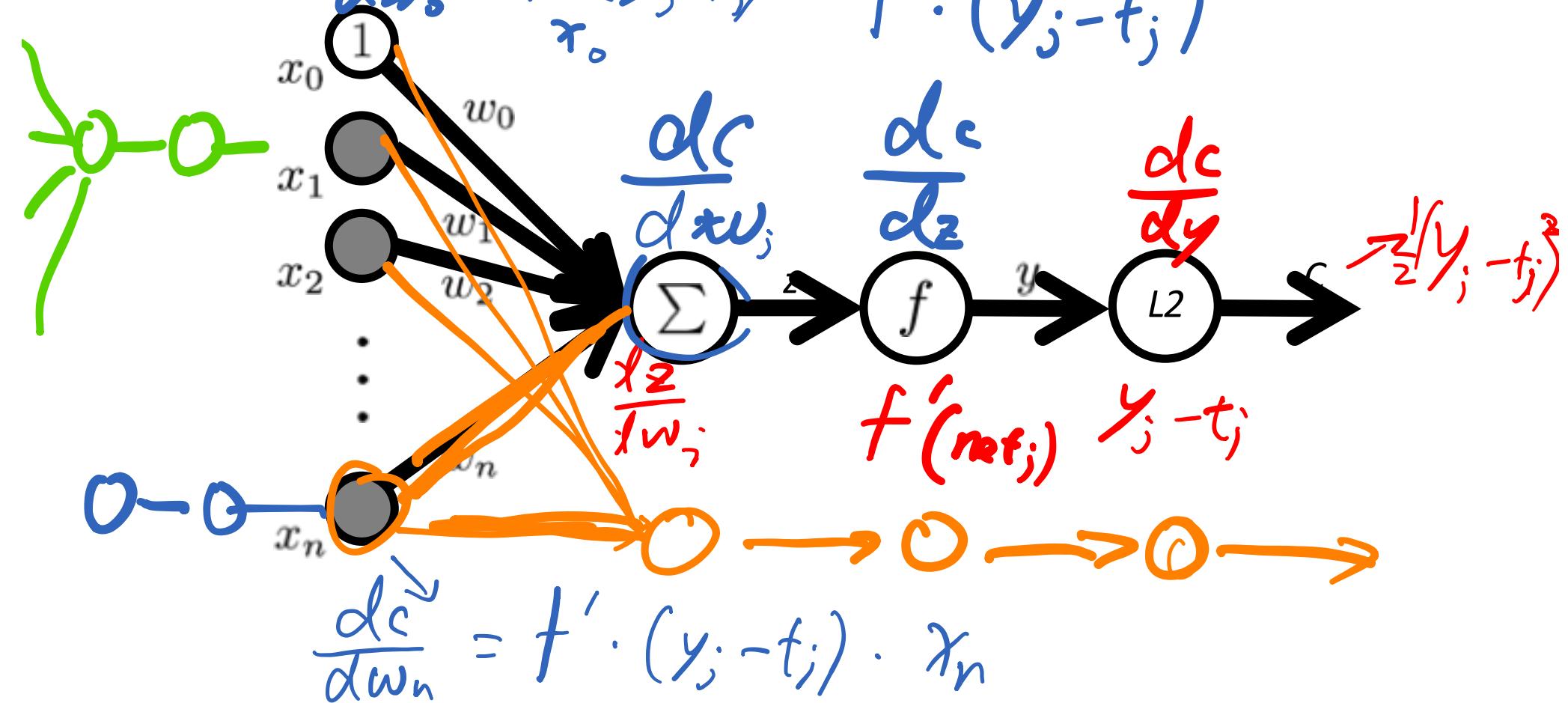
$$\Delta w_i = \alpha(t_j - y_j) f'(net_j) \underline{x_{ij}}$$



Chain Rule

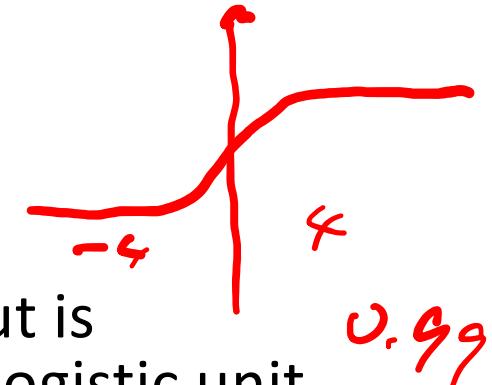
$$\Delta w_i = \alpha(t_j - y_j) f'(net_j) x_{ij}$$

$$\frac{\partial c}{\partial w_0} = f' \cdot (y_j - t_j) \cdot f' \cdot (y_j - t_j)$$



Problems with logistic activation and squared error

- The squared error measure has some drawbacks:
 - If the desired output is 1 and the actual output is 0.00000001 there is almost no gradient for a logistic unit to fix up the error.
 - If we are trying to assign probabilities to mutually exclusive class labels, we know that the outputs should sum to 1, but we are depriving the network of this knowledge.
- Is there a different cost function that works better?
 - Yes: Force the outputs to represent a probability distribution across discrete alternatives.

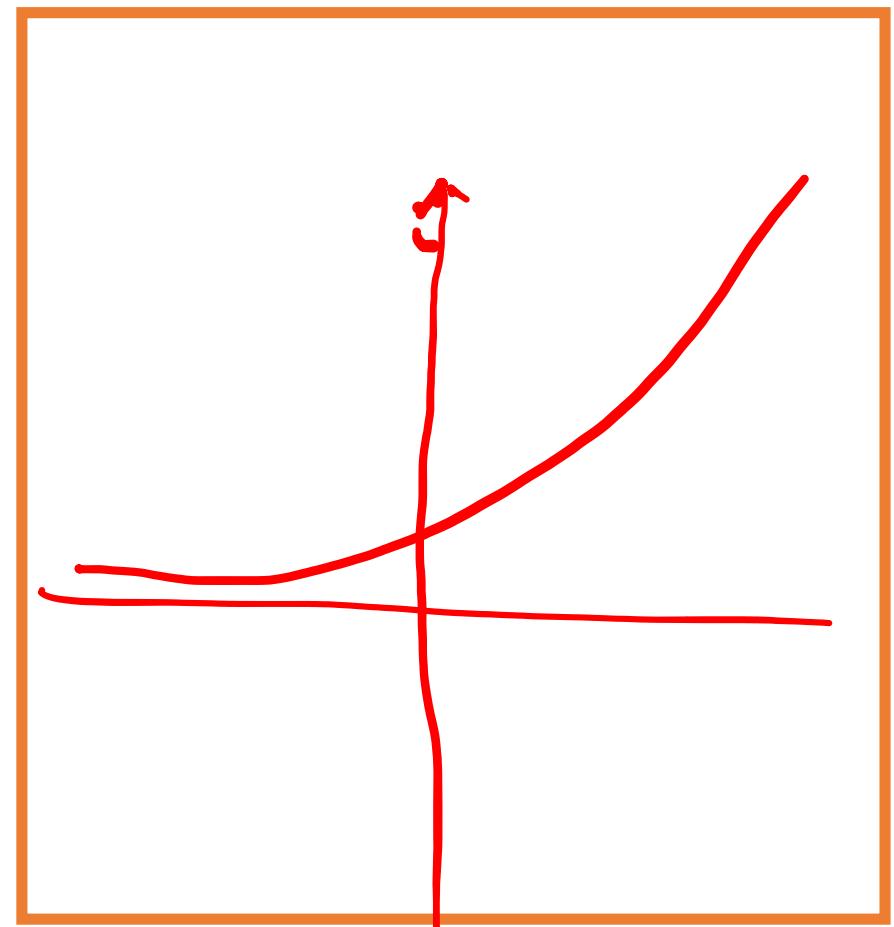


Softmax

Used for most classification problems:

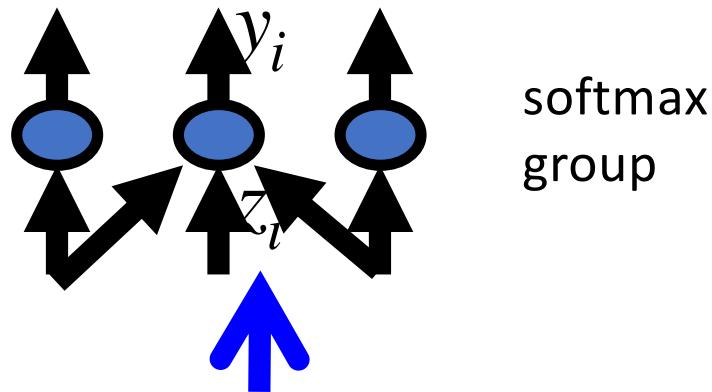
$$f(\underline{\text{net}_k}) = \frac{e^{\text{net}_k}}{\sum_l e^{\text{net}_l}}$$

$$\sum_{k=1}^c y_k = 1$$



Softmax

The output units in a softmax group use a non-local non-linearity:



$$y_i = \frac{e^{z_i}}{\sum_{j \in group} e^{z_j}}$$

$$\frac{\partial y_i}{\partial z_i} = y_i (1 - y_i)$$

$$z_i = \sum w_i x$$

Cross-entropy: the right cost function to use with softmax

- The right cost function is the negative log probability of the right answer.
- C has a very big gradient when the target value is 1 and the output is almost zero.
 - A value of 0.000001 is much better than 0.000000001
 - The steepness of dC/dy exactly balances the flatness of dy/dz

$$C = - \sum_j t_j \log y_j$$


target value

$$\frac{\partial C}{\partial z_i} = \sum_j \frac{\partial C}{\partial y_j} \frac{\partial y_j}{\partial z_i} = y_i - t_i$$

$$\frac{\partial}{\partial a^{(L+1)}(\mathbf{x})_c} - \log f(\mathbf{x})_y$$

$$= \frac{-1}{f(\mathbf{x})_y} \frac{\partial}{\partial a^{(L+1)}(\mathbf{x})_c} f(\mathbf{x})_y$$

$$= \frac{-1}{f(\mathbf{x})_y} \frac{\partial}{\partial a^{(L+1)}(\mathbf{x})_c} \text{softmax}(\mathbf{a}^{(L+1)}(\mathbf{x}))_y$$

$$= \frac{-1}{f(\mathbf{x})_y} \frac{\partial}{\partial a^{(L+1)}(\mathbf{x})_c} \frac{\exp(a^{(L+1)}(\mathbf{x})_y)}{\sum_{c'} \exp(a^{(L+1)}(\mathbf{x})_{c'})}$$

$$= \frac{-1}{f(\mathbf{x})_y} \left(\frac{\frac{\partial}{\partial a^{(L+1)}(\mathbf{x})_c} \exp(a^{(L+1)}(\mathbf{x})_y)}{\sum_{c'} \exp(a^{(L+1)}(\mathbf{x})_{c'})} - \frac{\exp(a^{(L+1)}(\mathbf{x})_y) \left(\frac{\partial}{\partial a^{(L+1)}(\mathbf{x})_c} \sum_{c'} \exp(a^{(L+1)}(\mathbf{x})_{c'}) \right)}{\left(\sum_{c'} \exp(a^{(L+1)}(\mathbf{x})_{c'}) \right)^2} \right)$$

$$= \frac{-1}{f(\mathbf{x})_y} \left(\frac{1_{(y=c)} \exp(a^{(L+1)}(\mathbf{x})_y)}{\sum_{c'} \exp(a^{(L+1)}(\mathbf{x})_{c'})} - \frac{\exp(a^{(L+1)}(\mathbf{x})_y)}{\sum_{c'} \exp(a^{(L+1)}(\mathbf{x})_{c'})} \frac{\exp(a^{(L+1)}(\mathbf{x})_c)}{\sum_{c'} \exp(a^{(L+1)}(\mathbf{x})_{c'})} \right)$$

$$= \frac{-1}{f(\mathbf{x})_y} \left(1_{(y=c)} \text{softmax}(\mathbf{a}^{(L+1)}(\mathbf{x}))_y - \text{softmax}(\mathbf{a}^{(L+1)}(\mathbf{x}))_y \text{softmax}(\mathbf{a}^{(L+1)}(\mathbf{x}))_c \right)$$

$$= \frac{-1}{f(\mathbf{x})_y} (1_{(y=c)} f(\mathbf{x})_y - f(\mathbf{x})_y f(\mathbf{x})_c)$$

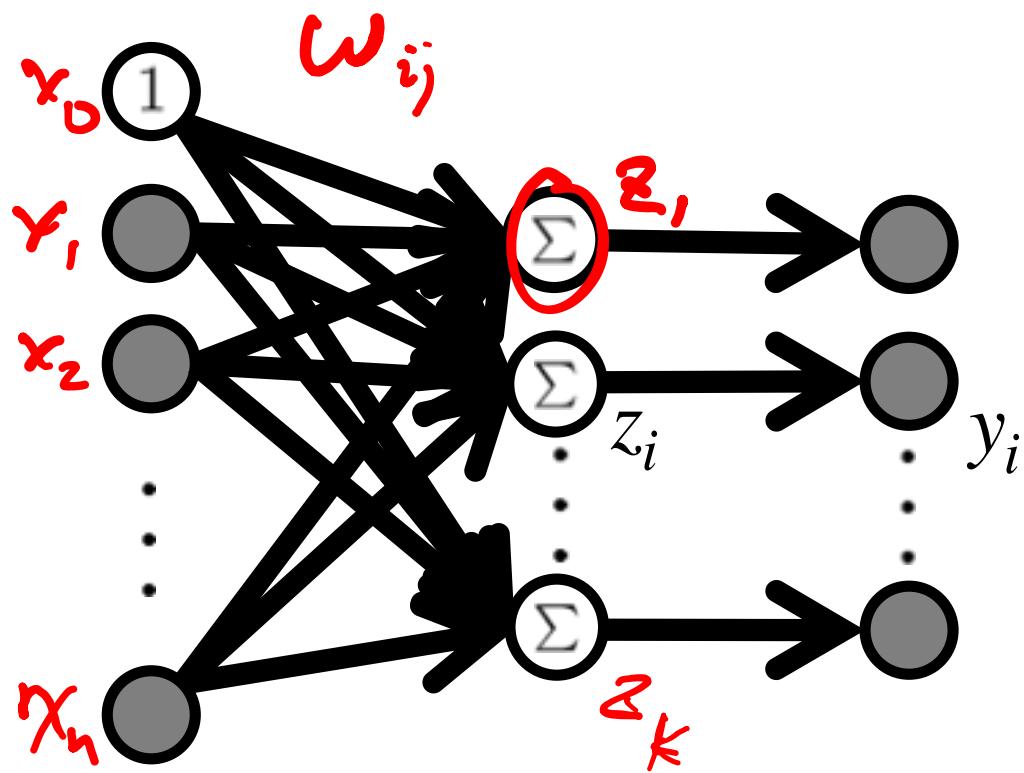
$$= - (1_{(y=c)} - f(\mathbf{x})_c)$$

$\leftarrow t_i - y_i$

$$\frac{\partial \frac{g(x)}{h(x)}}{\partial x} = \frac{\partial g(x)}{\partial x} \frac{1}{h(x)} - \frac{g(x)}{h(x)^2} \frac{\partial h(x)}{\partial x}$$

(from Hugo Larochelle)

Softmax



$$y_i = \frac{e^{z_i}}{\sum_{j \in group} e^{z_j}}$$

$$\frac{\partial y_i}{\partial z_i} = y_i (1 - y_i)$$

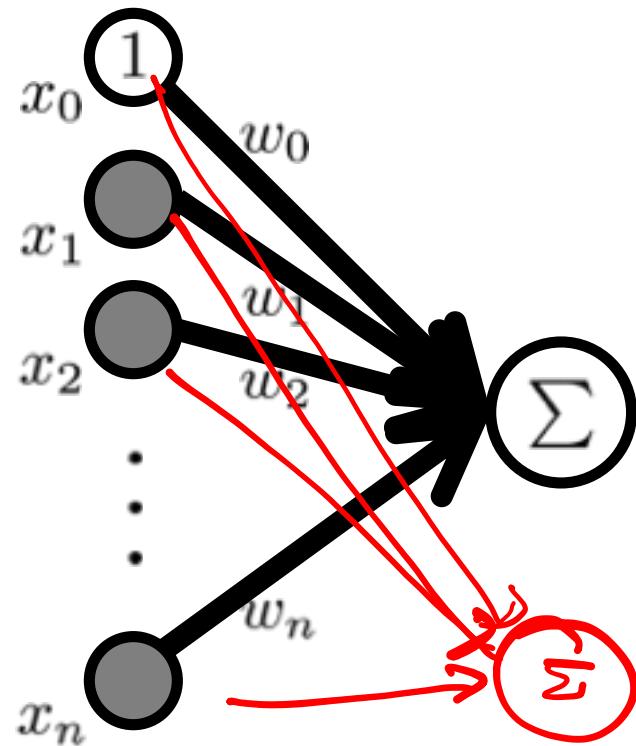
$$C = - \sum_j t_j \log y_j$$

▼ target value

$$\Delta w_i = \alpha(t_j - y_j) f'(net_j) x_{ij}$$

$$\frac{\partial C}{\partial z_i} = \sum_j \frac{\partial C}{\partial y_j} \frac{\partial y_j}{\partial z_i} = y_i - t_i$$

Logistic Regression (special case)



$$y_i = \frac{e^{w_i^T x}}{1 + e^{-(w_0 - w_i)^T x}}$$
$$y_1 = \frac{e^{w_1^T x}}{e^{w_0^T x} + e^{w_1^T x}}$$
$$y_2 = \frac{e^{w_2^T x}}{e^{w_0^T x} + e^{w_2^T x}}$$

