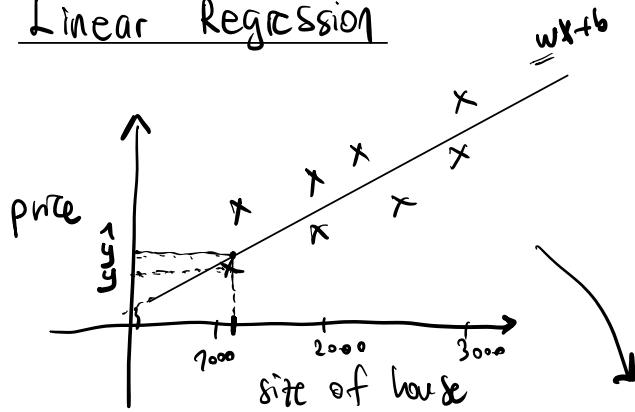


## Linear Regression



input var / feature  
 $x = (x, y)$  → output var / target.

$i^{th}$  training ex.:  $(x^{(i)}, y^{(i)})$

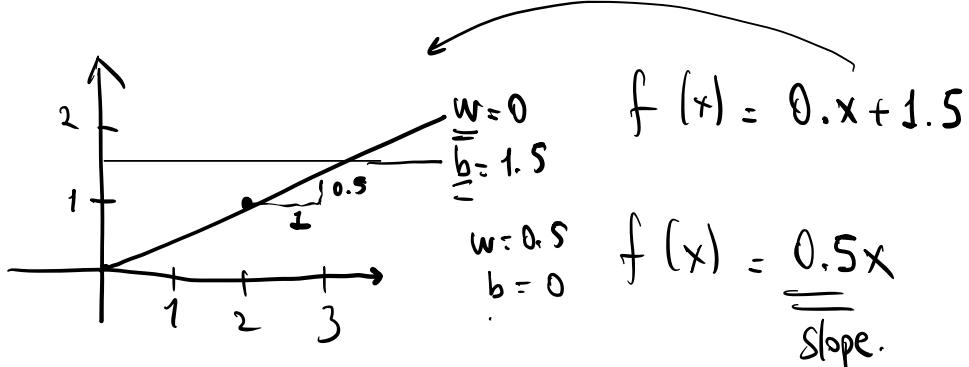
site of house	price
2010	600
:	:
:	:
:	:

$f_{w,b} = wx + b$

model parameters (coefficients) weights

linear regression with one variable.

univariate linear regression.  
one variable.



find  $w, b$ :

$\hat{y}^{(i)}$  is close to  $y^{(i)}$  for all  $(x^{(i)}, y^{(i)})$

prediction

ground truth

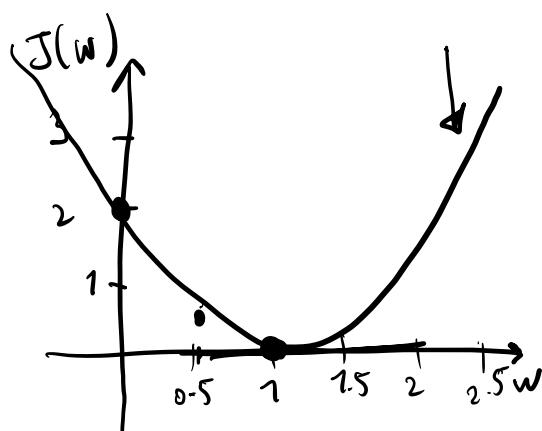
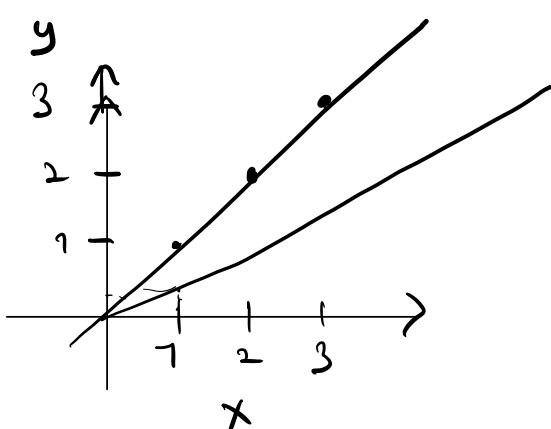
Cost function: Squared error cost function.

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2 \quad \# \text{ of training ex.}$$

$$= \frac{1}{2m} \sum_{i=1}^m \left( f_{w,b}(x^{(i)}) - y^{(i)} \right)^2$$

Goal:  $\boxed{\underset{w, b}{\text{minimize}} J(w, b)}$   $\rightarrow$   $\boxed{\underset{w}{\text{min}} J(w)}$

$$f_w(x) = w x \quad \underline{b = 0} \quad 0.5x$$



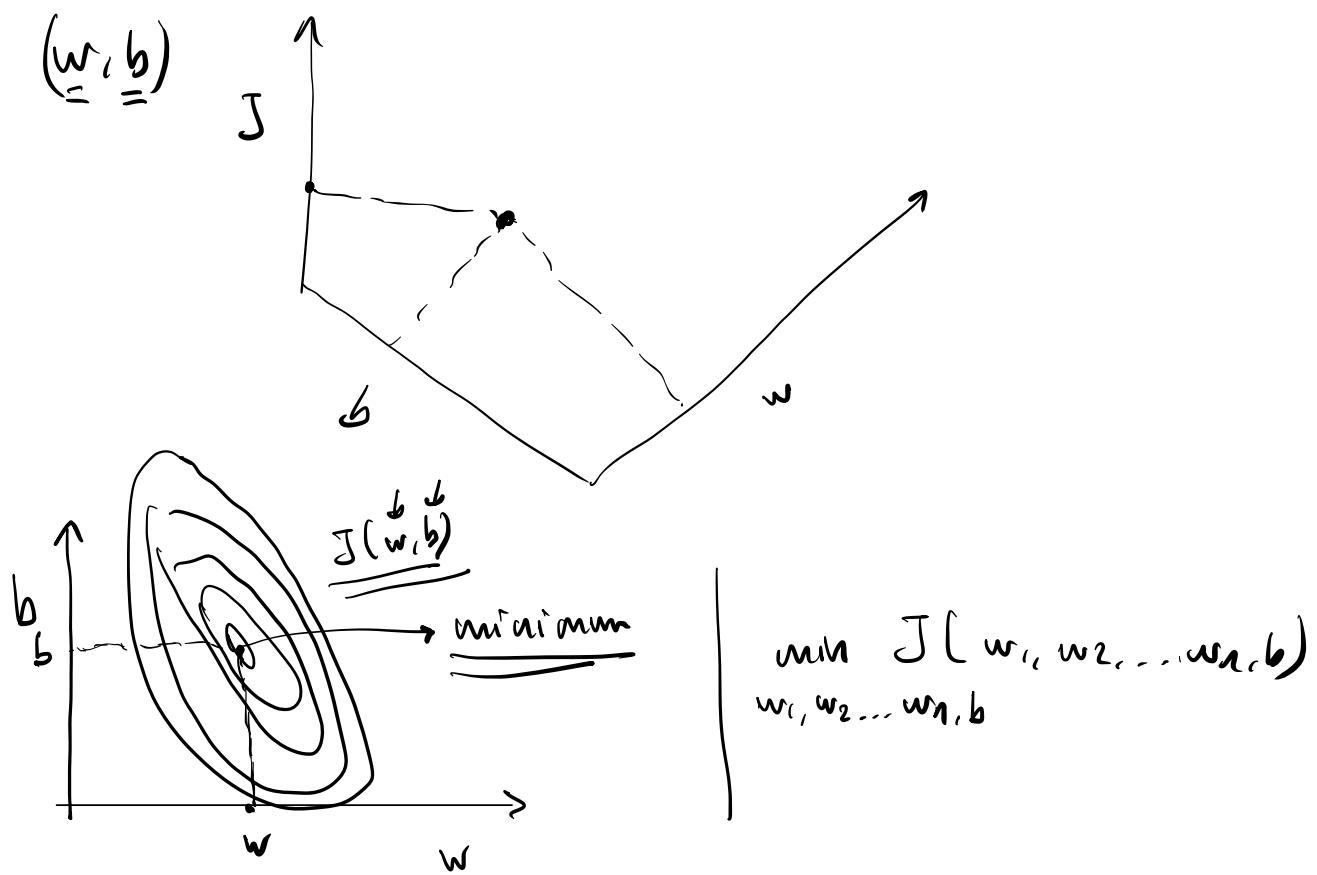
$$w = 0.5$$

$$J = \frac{1}{m} \left[ (0.5-1)^2 + (1-2)^2 + (1.5-3)^2 \right] \approx 0.58$$

$$w=0 \Rightarrow J \approx 2.3 \quad w=1$$

→ how to choose  $w$ ?

→ choose  $w$  to minimize  $J(w)$



## Gradient Descent Algorithm.

Outline:

Start with some  $w, b$  (ex:  $w=0.01, b=0.01$ )  $(0-1)$

keep changing  $w, b$  to reduce  $J(w, b)$

until we settle at or near a minimum

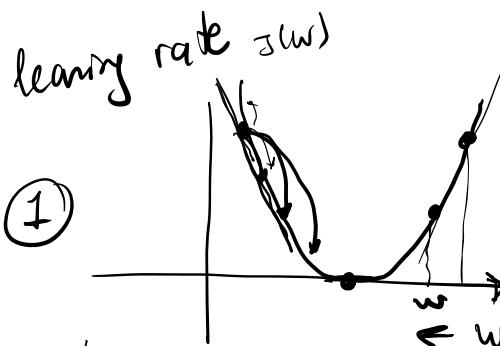
### Algorithm:

repeat until convergence {

$$w = w - \alpha \frac{\partial J(w, b)}{\partial w} \quad (1)$$

$$b = b - \alpha \frac{\partial J(w, b)}{\partial b} \quad (2)$$

} simultaneously update  $w, b$



$$\text{temp\_}w = w - \alpha \frac{\partial J(w, b)}{\partial w}$$

$$\text{temp\_}b = b - \alpha \frac{\partial J(w, b)}{\partial b}$$

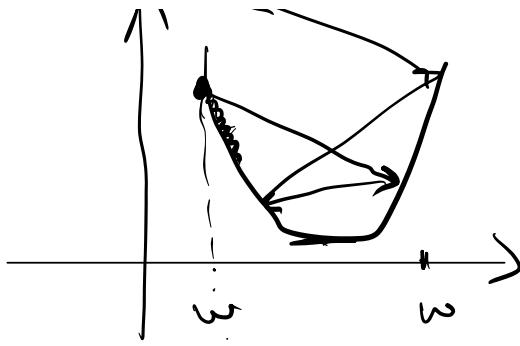
$$w = \text{temp\_}w$$

$$b = \text{temp\_}b$$

fix size  $\alpha = \underline{0.1} \rightarrow \underline{\underline{10^{-3}}}$   
 $\wedge \quad \curvearrowleft$

if  $\alpha$  is too small

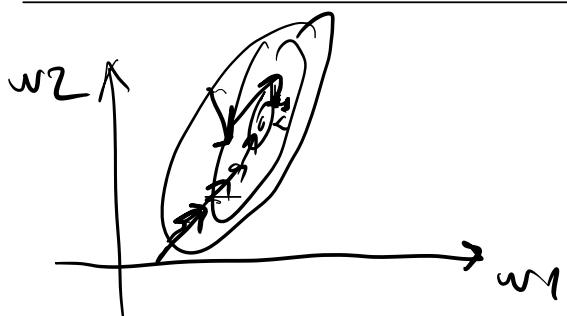
$$\alpha = 0.000001$$



if  $\alpha$  is too large

$\alpha = 100$ , overshoot, never reach minimum

### Adam optimization.



$$w_1 = w_1 - \alpha_1 \frac{\nabla J(\vec{w}, b)}{\| \nabla w_1 \|}$$

$$w_{10} = w_2 = (\alpha_{10})$$

if  $w_j$  keeps moving in same direction,

increase  $\alpha_j$ .

If  $w_j$  keeps oscillating, reduce  $\alpha_j$ .

$$\frac{\partial J(w, b)}{\partial w} = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

$\downarrow$

$wx^{(i)} + b$

$$= \frac{1}{2m} \sum_{i=1}^m (wx^{(i)} + b - y^{(i)}) \cancel{2x^{(i)}}$$

$$\frac{\partial J(w, b)}{\partial w} = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x^{(i)}$$

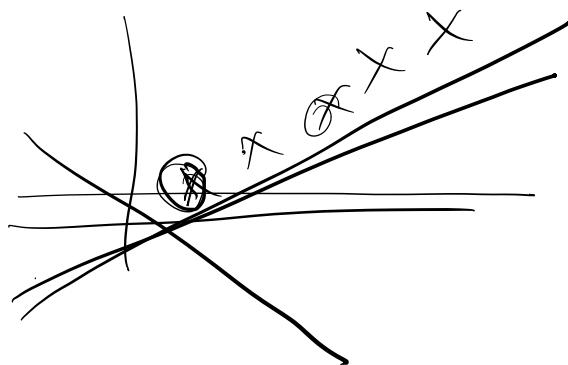
$$\frac{\partial J(w, b)}{\partial b} = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$$

Batch gradient

→ each step of gradient descent.  
we use all the training ex.

→ Stochastic Gradient Descent  $\otimes$

↓  
1 sample



→ Mini-batch.

1 calculation  $\downarrow$

1M data  $\rightarrow$  batch gradient  $\rightarrow$  1 iteration

→ stochastic  $\rightarrow$  1M in 1 iteration

→ mini-batch (200)  $\rightarrow$  10000 in 1 iteration

Multiple features

size in feet	# of bedrooms	# of floors	Age in years	price
-----------------	------------------	----------------	-----------------	-------

$$(x_1) \quad | \quad (x_2) \quad | \quad (x_3) \quad | \quad (x_4) \quad | \quad \dots$$

$$f_{w,b}(x) = wx + b.$$

$$f_{w,b} = w_1 x_1 + w_2 x_2 - \dots + b.$$

$$\vec{w} = [w_1 \dots w_n] \quad \vec{x} = [x_1 \dots x_n]$$

$\vec{w}$  is a vector  
 $b$  is a number

$$f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$$

↓  
dot product.

multiple linear regression.

## Vectorization

$$f_{\vec{w}, b}(\vec{x}) = \sum_{j=1}^n w_j x_j + b.$$

$$\vec{w} = [w_1 \ w_2 \ w_3]$$

$b$  is a number

$$\vec{x} = [x_1 \ x_2 \ x_3]$$

$$\left\{ \begin{array}{l} \vec{w} = [2, 2, 3], \quad \vec{x} = [2, 3, 4] \\ b = 10 \end{array} \right. , \quad \text{total} = 0$$

for  $j$  in range  $(0, n)$ :

sequential  
↓

$$\left\{ \begin{array}{l} \text{total} = \text{total} + w[j] * x[j] \\ \text{total} = \text{total} + b \end{array} \right.$$

Vectorization.

$$f_{w,b}(\vec{x}) = \vec{w}_j \vec{x} + b.$$

$$f = n_p. f_{w,b}(w, x) + b$$

→ Shorter  
much faster  
use parallel comp.  
CPU  
GPU.

6.4 for multiple LR

repeat

$$w_1 = w_1 - \alpha \frac{\partial J(w, b)}{\partial w_1}$$

$$w_{10} = w_{10} - \alpha \frac{\partial J(w, b)}{\partial w_{10}}$$

} Simultaneously update

Normal Equation. if you  $\sim < 10^6$

solve  $w, b$  without iterations.

$$w = (x^T x)^{-1} (x^T y)$$

Feature Scaling.

① max.

$$x_1 \rightarrow \frac{300 \leq x_1 \leq 2000}{2000 \quad 2000}$$

$$x_1, \text{scaled} = \frac{x_1}{2000}$$

② mean normalization.  $= [ \quad \quad \quad ]$

$$x_1 = \frac{x_1 - M}{\max - \min} \stackrel{M=800}{=} \frac{300 \leq x_1 \leq 2000}{2000 - 300} \stackrel{0.18 \leq x_1 \leq 0.8}{\approx}$$

③ Z-score

$$Z_1 = \frac{x_1 - M_1}{\sigma_1} \quad -0.67 \leq Z_1 \leq 3.1$$

~~Ex~~

$$\text{price} = w_1 x_1 + w_2 x_2 + b$$

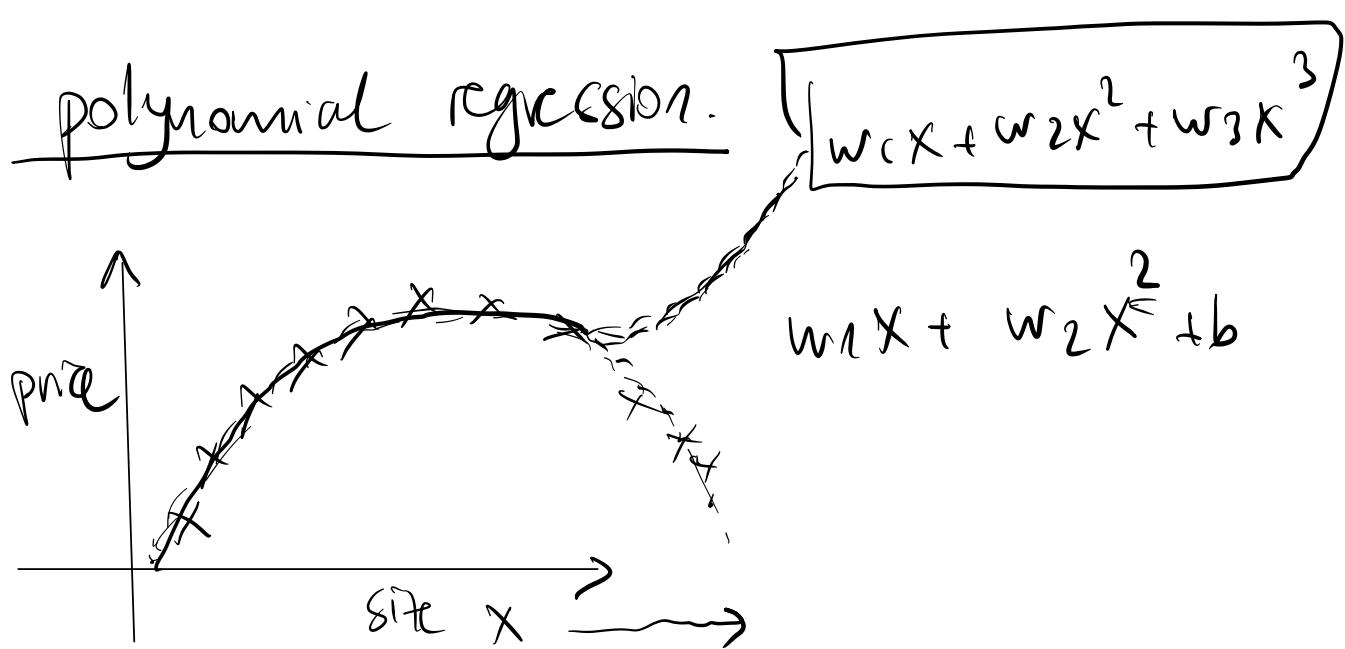
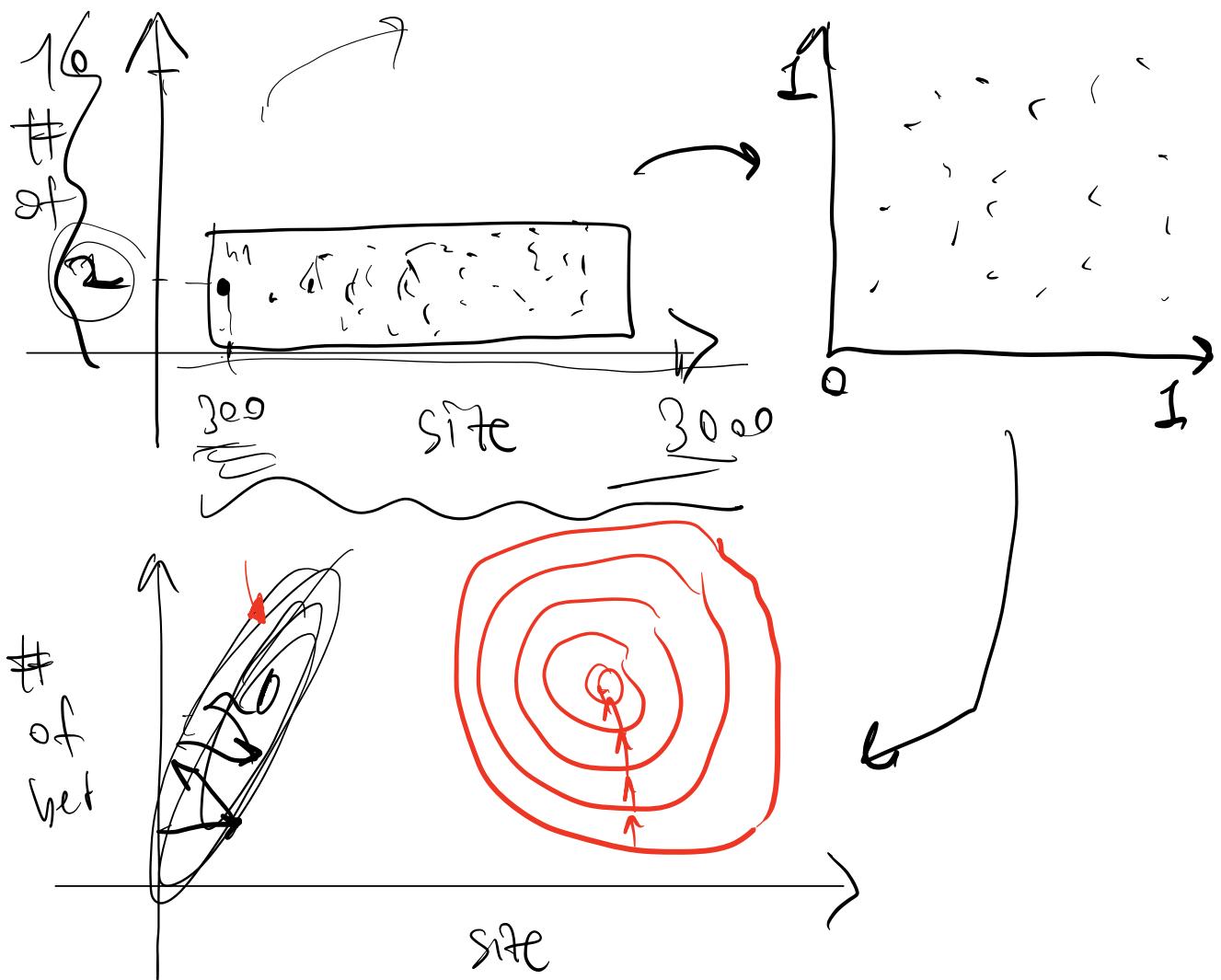
↓                            ↓  
size                              # bed rooms.

House 1:  $x_1 = 2000$ ,  $x_2 = 5$  price = \$50K

$w_1 = 0.1$	$w_2 = 50$	$b = 50$
-------------	------------	----------

$$\text{price} = \underbrace{0.1 \cdot 2000}_{200K} + \underbrace{50 \times 5}_{250K} + \underbrace{50}_{50K}$$

$$= 300K \$$$



# feature engineering.

$$w_1x_1 + w_2x_2 + b$$

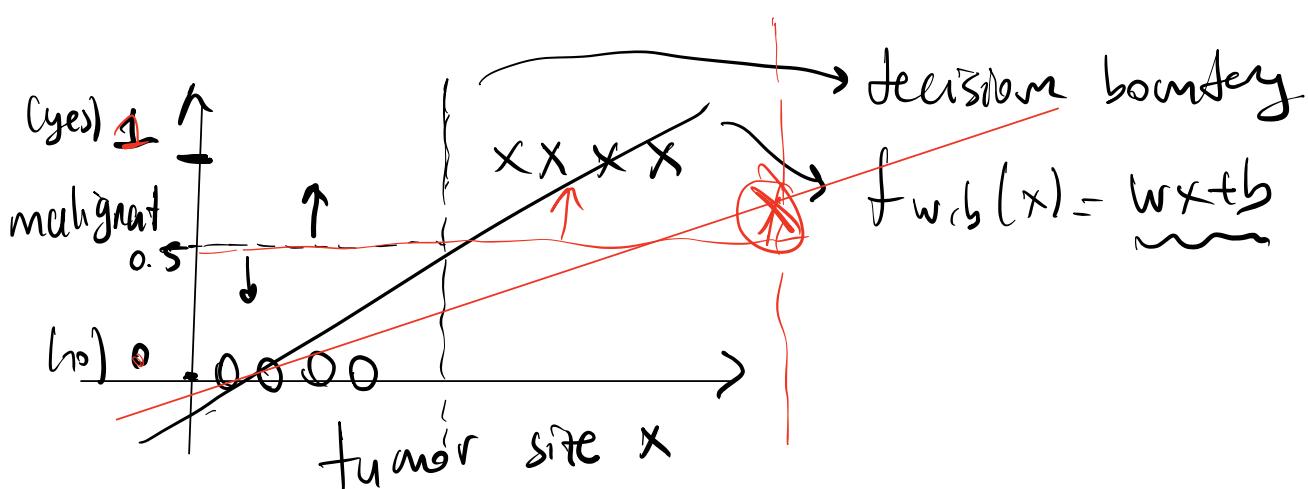
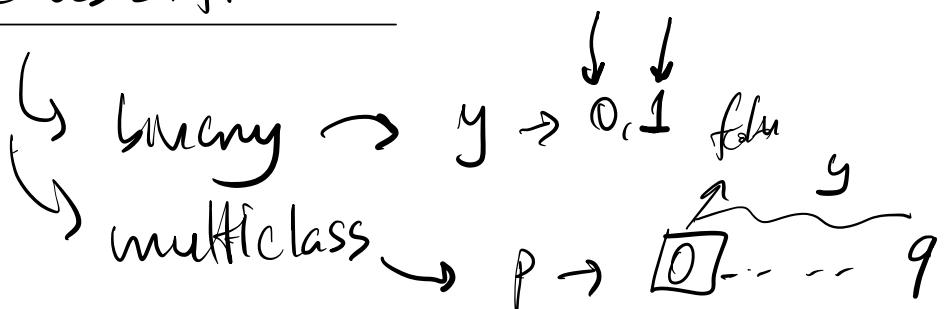
$\downarrow$        $\downarrow$   
frontage      depth

$$x_3 = \underbrace{\text{frontage} \times \text{depth}}_{\text{area}}$$



$$f_{w,b}(x) = w_1x_1 + w_2x_2 + w_3x_3$$

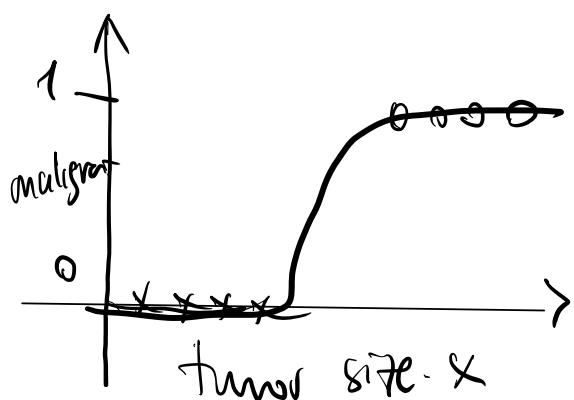
## Classification.



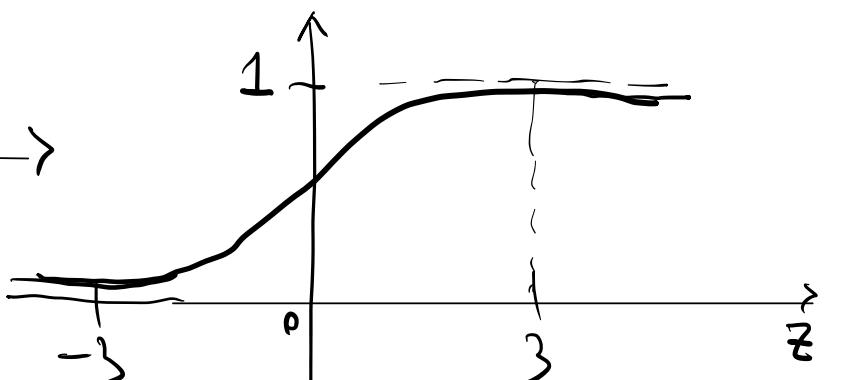
If  $f_{w,b}(x) < 0.5 \rightarrow \hat{y} = 0$

If  $f_{w,b}(x) \geq 0.5 \rightarrow \hat{y} = 1$

## Logistic Regression.



want outputs b/w 0-1



sigmoid

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

sigmoid (logistic) function  
outputs 0 - 1

$$0 < g(z) < 1$$

$$f_{\vec{w}, b}(\vec{x}) = \vec{w}\vec{x} + b$$

$\vec{z} = \vec{w}\vec{x} + b$

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

0 - 1

$$= \boxed{\frac{1}{1+e^{-(\vec{w}\cdot \vec{x} + b)}}}$$

$$f_{\vec{w}, b}(\vec{x}) = 0.7 \quad \text{%.70 chance that } y \text{ is 1}$$

%.30 chance that  $y$  is 0

$$P(y=0) + P(y=1) = 1$$

$$f_{\vec{w}, b}(\vec{x}) = \underbrace{P(y=1 | \vec{x}; \vec{w}, b)}$$

Decision Boundary.

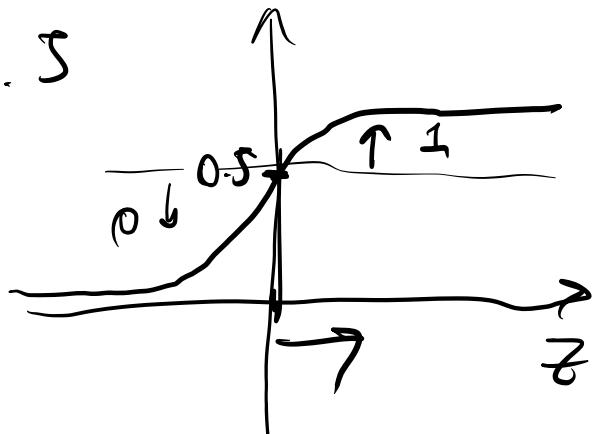
$$\frac{1}{1 + e^{-(\vec{w}\vec{x} + b)}} = P(y=1 | \vec{x}; \vec{w}, b) = 0.7$$

is  $f_{\vec{w}, b}(\vec{x}) \geq 0.5 \xrightarrow{\text{threshold}}$

yes  $\Rightarrow \hat{y} = 1, \hat{y} = 0$

when is  $f_{\vec{w}, b} \geq 0.5$

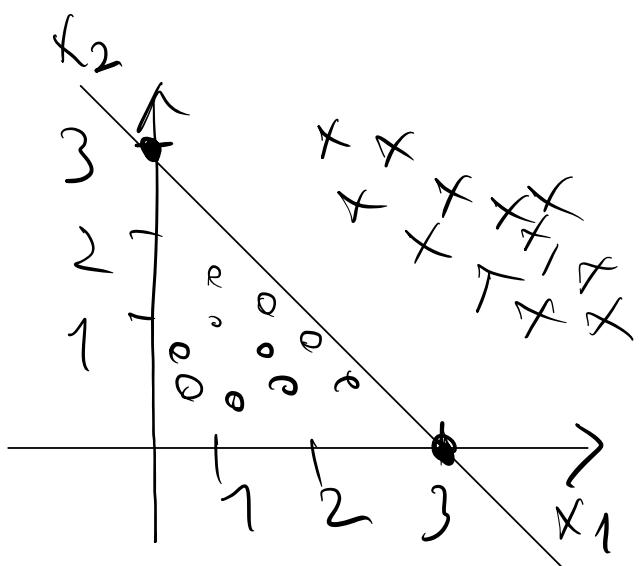
$$g(z) \geq 0.5$$



$$\begin{matrix} z \\ = \end{matrix} \geq 0$$

$$\vec{w}\vec{x} + b \geq 0 \Rightarrow \hat{y} = 1$$

$$\vec{w}\vec{x} + b < 0 \Rightarrow \hat{y} = 0$$



$$f_{\vec{w}, b} = g(z)$$

$$g(z) = g(w_1 x_1 + w_2 x_2 + b)$$

↓      ↓      ↓

1      1      3

boundary.

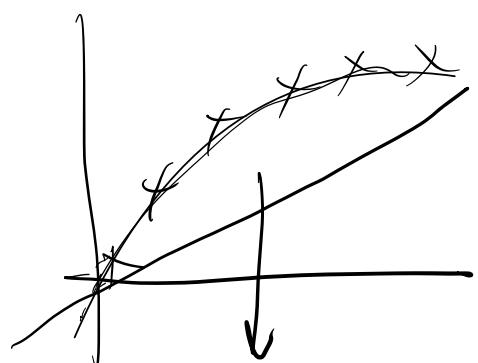
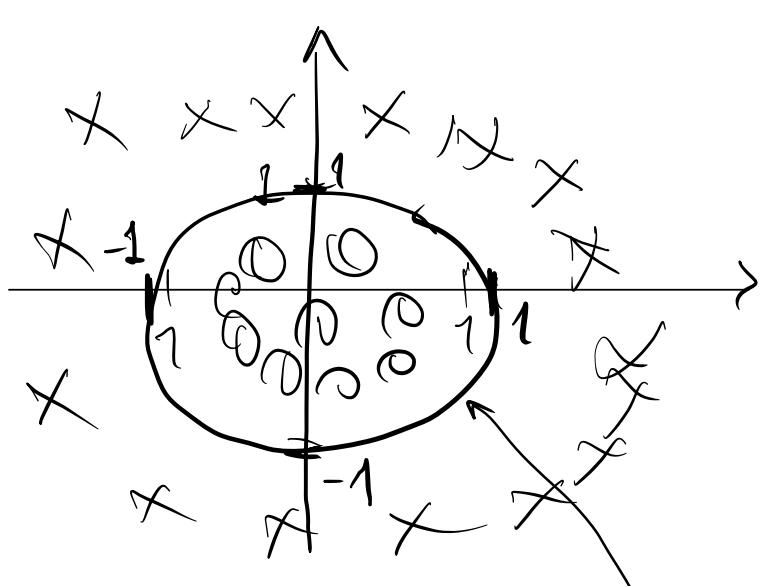
decision

$$z = \vec{w} \cdot \vec{x} + b \leq 0$$

$$z = x_1 + x_2 - 3 = 0$$

$$\begin{matrix} x_1 + x_2 = 3 \\ 0 \end{matrix}$$

Non-linear decision boundaries.



$$wx + b$$

$$wx + wx^2$$

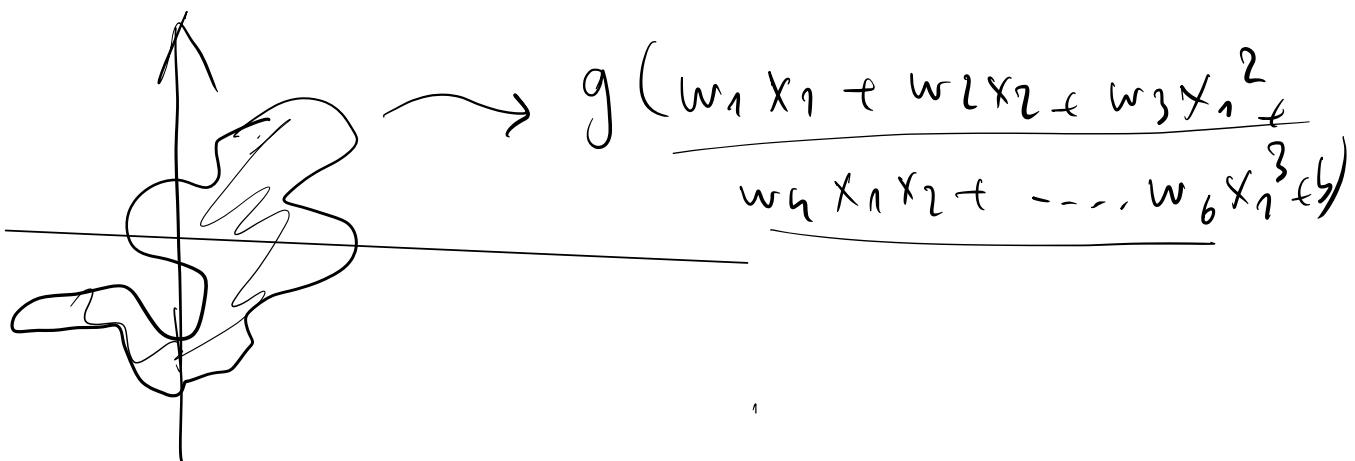
$$g(w_1x_1^2 + w_2x_2^2 + b) \quad \begin{matrix} 1 & -1 \end{matrix}$$

$$x_1^2 + x_2^2 \leq 1$$

$$y = 0$$

$$z = x_1^2 + x_2^2 - 1 = 0$$

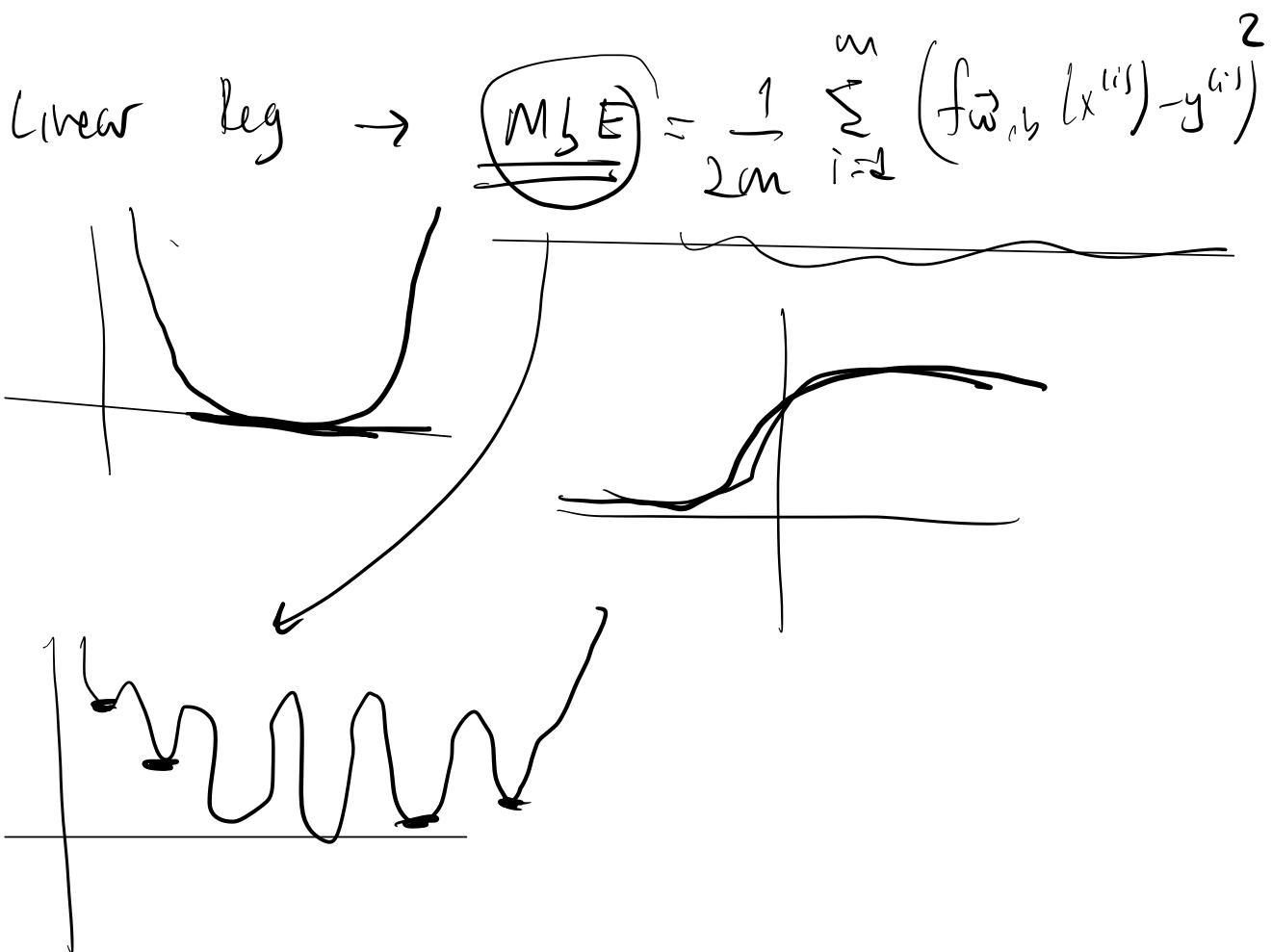
$$x_1^2 + x_2^2 \leq 1$$



$$f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w}\vec{x} + b)}}$$

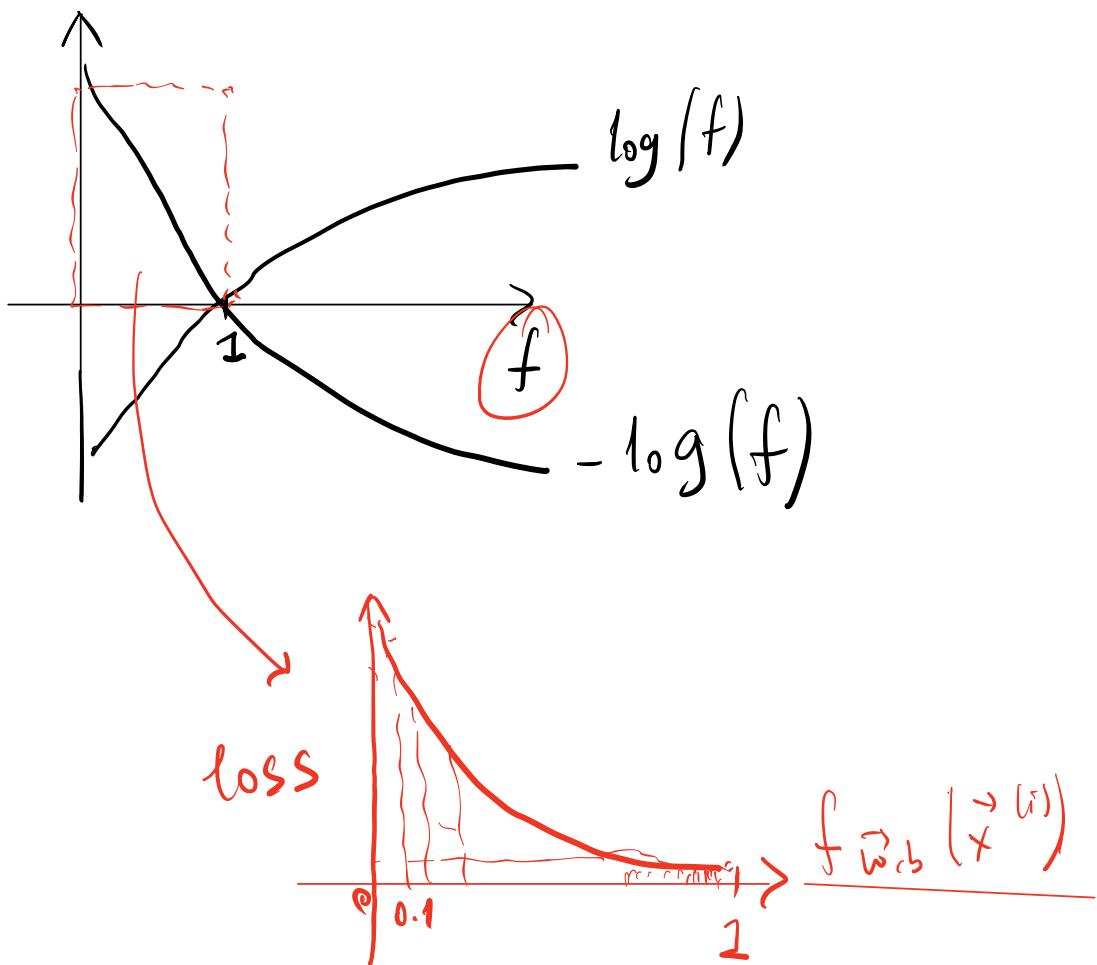
$$\vec{w} = [w_1, \dots, w_n] \text{ and } b \text{ u?}$$

predict malignant.



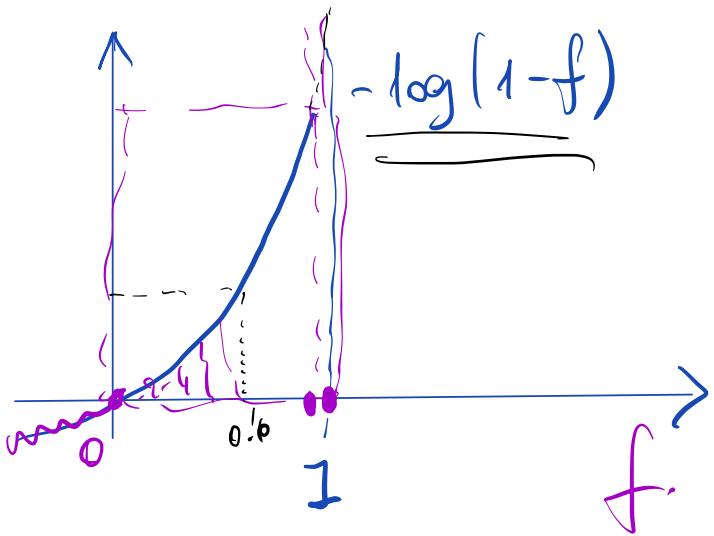
Logistic loss function. for only one fr.  
sample

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \begin{cases} 1 & -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) \text{ if } y^{(i)} = 1 \\ 2 & -\log(1-f_{\vec{w}, b}(\vec{x}^{(i)})) \text{ if } y^{(i)} = 0 \end{cases}$$



$\underline{x}^{(i)}, \underline{y}^{(i)}$   
email    spam (1)

- $f_{\text{w}, b}(\vec{x}) \rightarrow 1$  then loss  $\rightarrow 0$
- $f_{\text{w}, b}(\vec{x}) \rightarrow 0$  then loss  $\rightarrow \infty$



$$y = 0$$

$$\begin{array}{c} 0.51 \geq 0.5 \rightarrow 1 \\ 0.9999 \geq 0.5 \rightarrow 1 \end{array}$$

$f_{\vec{w}, b}(\vec{x}) = 1$  then loss  $\rightarrow \infty$

$f_{\vec{w}, b}(\vec{x}) = 0$  then loss  $\rightarrow 0, \dots$

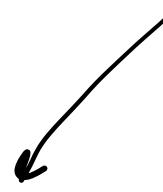
→ loss for one sample

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \underbrace{-y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)}))}_{\text{if } y^{(i)} = 1} + \underbrace{(1-y^{(i)}) \log(1-f_{\vec{w}, b}(\vec{x}))}_{\text{if } y^{(i)} = 0}$$

if  $y^{(i)} = 1 : -1 \log(f(\vec{x}))$

if  $y^{(i)} = 0 : -(1-f(\vec{x}))$

# of ex: M

$$J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m [$$


$$]$$

Repeat {

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$$

$$b = b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b)$$

} sum-up for.

$$J(\vec{w}) = -\frac{1}{m} \sum_{i=1}^m \left[ y_i \log(w^T x_i) + (1-y_i) \log(1-\sigma(w^T x_i)) \right]$$

$$\sigma = \frac{1}{1+e^{-z}}, \quad z = w^T x_i$$

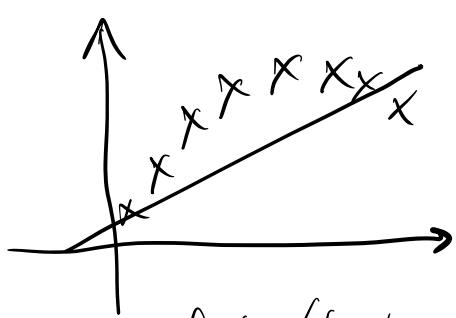
$$\textcircled{1} \quad \sigma'(z) = \frac{d}{dz} \left( \frac{1}{1+e^{-z}} \right) = \sigma(z) \cdot \sigma'(z) (1-\sigma(z))$$

$$\frac{\partial \mathcal{J}(\mathbf{w})}{\partial w_j} = \sum_{i=1}^m [y_i (\mathbf{w}^T \mathbf{x}_i) - (1-y_i) \sigma(\mathbf{w}^T \mathbf{x}_i)] \mathbf{x}_{ij}$$

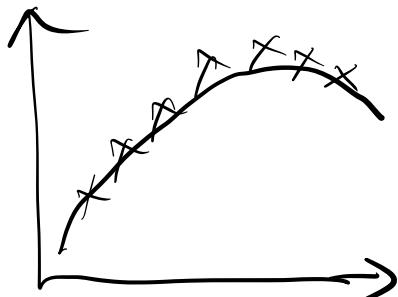
②

$$\frac{\partial \mathcal{J}(\mathbf{w})}{\partial w_j} = \frac{1}{M} \sum_{i=1}^m [y_i (\mathbf{w}^T \mathbf{x}_i) - (1-y_i) \sigma(\mathbf{w}^T \mathbf{x}_i)] \mathbf{x}_{ij}$$

### Overfitting



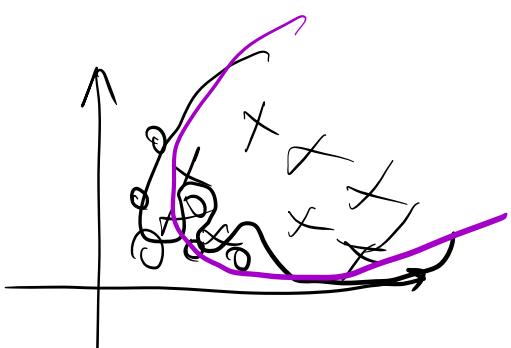
underfit / high bias.



generalize.



overfit.  
high variance

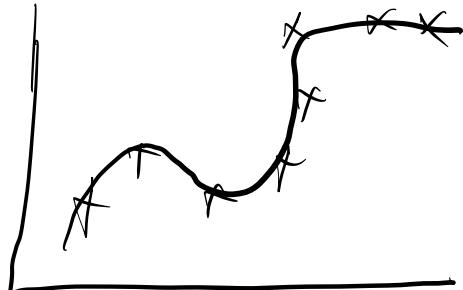
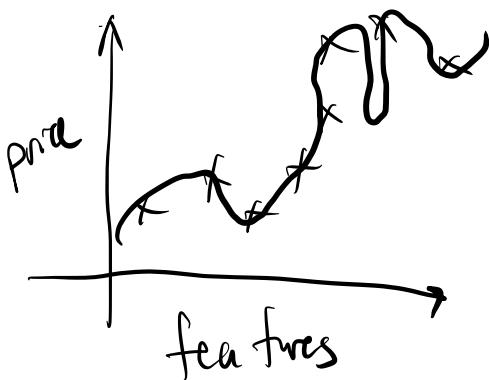


## Address overfitting.

- ① collect more training data
- ② select features to include/exclude  
~~~~~ feature selection.

1 - ----- 100

- ③ Regularization.



$$f(x) = 28x - 385x^2 + 39x^3 - 17x^4 + 100$$

$$f(x) = 13x - 0.23x^2 + 0.00000014x^3 - 0.0001x^4 + 10$$

# of factors

$w_1$  —————  $w_{100}$  ( $n = 100$ )

$$\text{min } J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

① fit data      ② keep  $w_j$  small

balance goal.

Lambda regularization term  $\geq 0$

$\lambda \downarrow \rightarrow$  overfit

$\lambda \uparrow \rightarrow$  underfit

$$\text{Ridge (L2)} = \text{cost func} = \text{OLS loss} + \lambda \sum_{j=1}^n w_j^2$$

$$\text{Lasso (L1)} = \text{cost func} = \text{OLS loss} + \lambda \sum_{j=1}^n |w_j|$$

feature selection, make zero.

$$w_1 x_1 + \frac{w_2 x_2}{0}$$

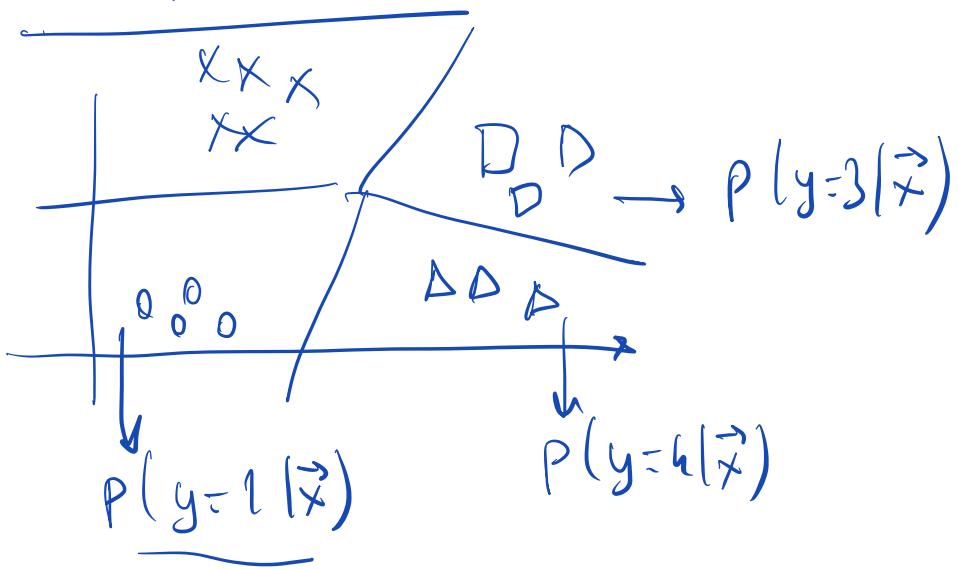
~~$w_1 x_1 + \frac{w_2 x_2}{0}$~~

# Elastic-Net Regularization.

$$\text{OLS Loss} + \frac{\rho}{p} \lambda \sum_{i=1}^n |w_i| + \frac{1-\rho}{2} \times \sum_{i=1}^n w_i^2$$

Lasso
Ridge

Multiclass



$$z = \vec{w} \cdot \vec{x} + b$$

$$u_1 = g(z) = \frac{1}{1+e^{-z}} = P(y=1|\vec{x}) \xrightarrow{z=0.70} P(y=0|\vec{x}) = 0.3$$

4 possible outputs.

Softmax

$$z_1 = \vec{w}_1 x + b_1 = q_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = P(y=1|x)$$

$$z_2 = \dots = P(y=2|x)$$

$$z_3 = \dots = -$$

$$z_4 = w_4 x + b_4 = q_4 = \frac{e^{z_4}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = -$$

# logistic Reg 2outbit

$$q(1) = P(y=1|x)$$

$$q(2) = P(y=0|x)$$

$$q_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2}}$$

$$\text{loss} = -y \log q_1 - (1-y) \log (1-q_1)$$

$$\text{if } y=1 \quad \text{if } y=0 \quad q_N = \frac{e^{z_N}}{e^{z_1} + \dots + e^{z_N}}$$

Crossentropy loss

$$= - \sum_{i=1}^N y_i \log(q_i)$$

$-\log q_i$  if  $y=1$

$$\text{loss}(a_1 - a_2) = \begin{cases} \dots & \dots \\ -\log_2 \text{if } y=2 \\ \dots & \dots \\ -\log_N \text{if } y=N \end{cases}$$