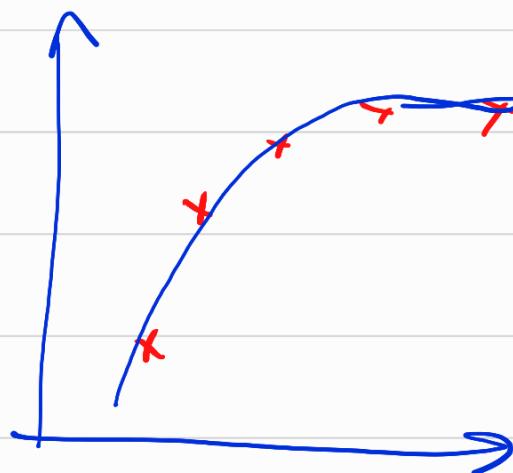


underfit

- high bias

✗

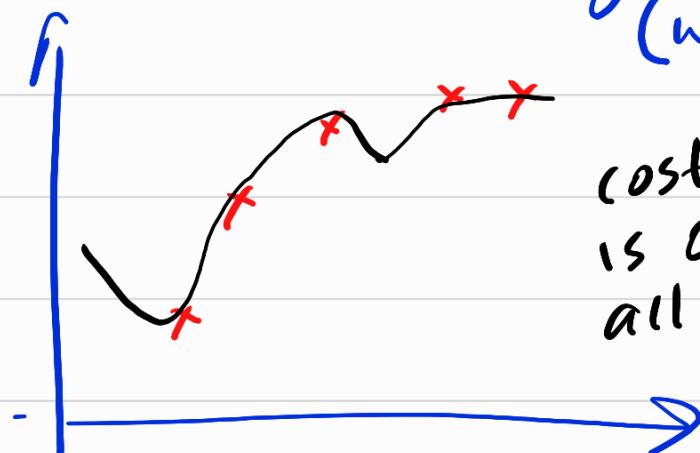


✓ $w_1 x + w_2 x^2 + b$

fits well

generalization

(well on brandex ex)



cost
is 0
all fits to
model

$$w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + b$$

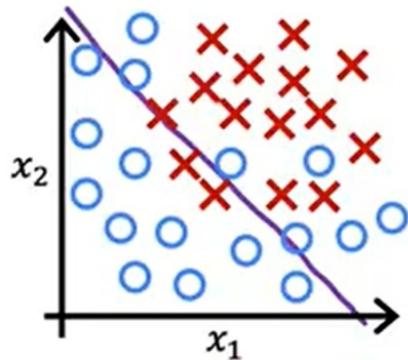
✗

* Overfit

not good for new examples
high variance

for classification — Overfitting

Classification

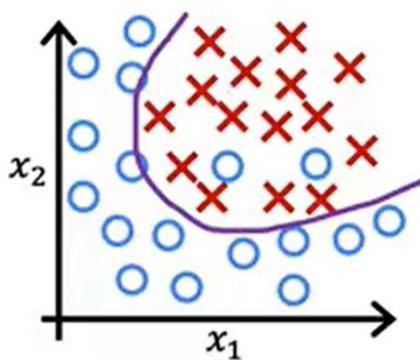


$$z = w_1 x_1 + w_2 x_2 + b$$

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

g is the sigmoid function

underfit high bias

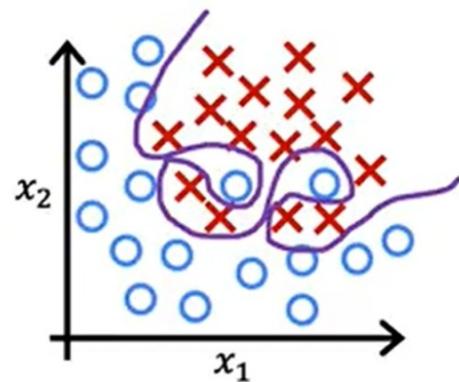


$$z = w_1 x_1 + w_2 x_2$$

$$+ w_3 x_1^2 + w_4 x_2^2$$

$$+ w_5 x_1 x_2 + b$$

just right



$$z = w_1 x_1 + w_2 x_2$$

$$+ w_3 x_1^2 x_2 + w_4 x_1^2 x_2^2$$

$$+ w_5 x_1^2 x_2^3 + w_6 x_1^2 x_2^6$$

$$+ \dots + b$$

10:50

Address ;

① Collect more training examples



② Select features to include/exclude

all features

+

insufficient data

} overfit



③ Regularization

- reduce size of parameters
(eliminate feature $x^q \rightarrow 0$)

↓
small values for w_j

reducing $b \rightarrow$ makes low diff.

Cost function with Regularization

$$\min \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(x^i) - y^i)^2 + 1000 w_3 + 1000 w_4$$

to minimize w_3 & $w_4 \rightarrow 0$

$$J(\vec{w}, b) = \dots + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

lambda > 0

(regularization parameter)

scale both similarly

Regularization

small values w_1, w_2, \dots, w_n, b

simpler model

less likely to overfit

$$w_3 \approx 0$$

$$w_4 \approx 0$$

size x_1	bedrooms x_2	floors x_3	age x_4	avg income x_5	...	distance to coffee shop x_{100}	price y
$w_1, w_2, \dots, w_{100}, b$							$n = 100$

n features

$$J(\vec{w}, b) = \frac{1}{2m} \left[\sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \underbrace{\frac{\lambda}{2m} \sum_{j=1}^n w_j^2}_{\text{"lambda" regularization parameter}} + \underbrace{\frac{\lambda}{2m} b^2}_{\lambda > 0} \right]$$

regularization term

can include or exclude b

not needed

if λ is very large

$$f_{\vec{w}, b}(\vec{x}) = \underbrace{w_1 x_1 + w_2 x_2 + \dots + b}_0$$

$$f(x) = b \rightarrow \text{underfit}$$

To implement Gradient Descent ;

repeat {

$$w_j = w_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m \left[(f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)} \right] + \frac{\lambda}{m} w_j \right]$$

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

} simultaneous update



same term
for logistic reg.

More Derivative Examples

$w = 3$

$J(w) = w^2 = 9$

$w \uparrow 0.001$

$J(w) = J(3.001) = 9.006001$

$\frac{\partial}{\partial w} J(w) = 6$

$J(w) \uparrow 6 \times 0.001$

$w = 2$

$J(w) = w^2 = 4$

$w \uparrow 0.001$

$J(w) = J(2.001) = 4.004001$

$\frac{\partial}{\partial w} J(w) = 4$

$J(w) \uparrow 4 \times 0.001$

$w = -3$

$J(w) = w^2 = 9$

$w \uparrow 0.001$

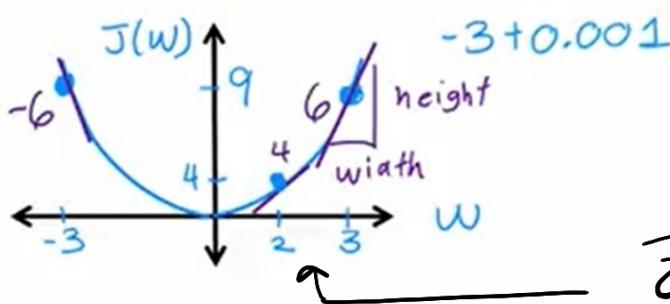
$J(w) = J(-2.999) = 8.994001$

$\frac{\partial}{\partial w} J(w) = -6$

$J(w) \downarrow 6 \times 0.001$

$J(w) \uparrow -6 \times 0.001$

↓ 0.006



$\frac{\partial}{\partial w} J(w) = 2w$

1

$y = x^2$
 $\frac{dy}{dx} = 2x$

$$w=2 \quad w^2=4 \quad 2w=4$$

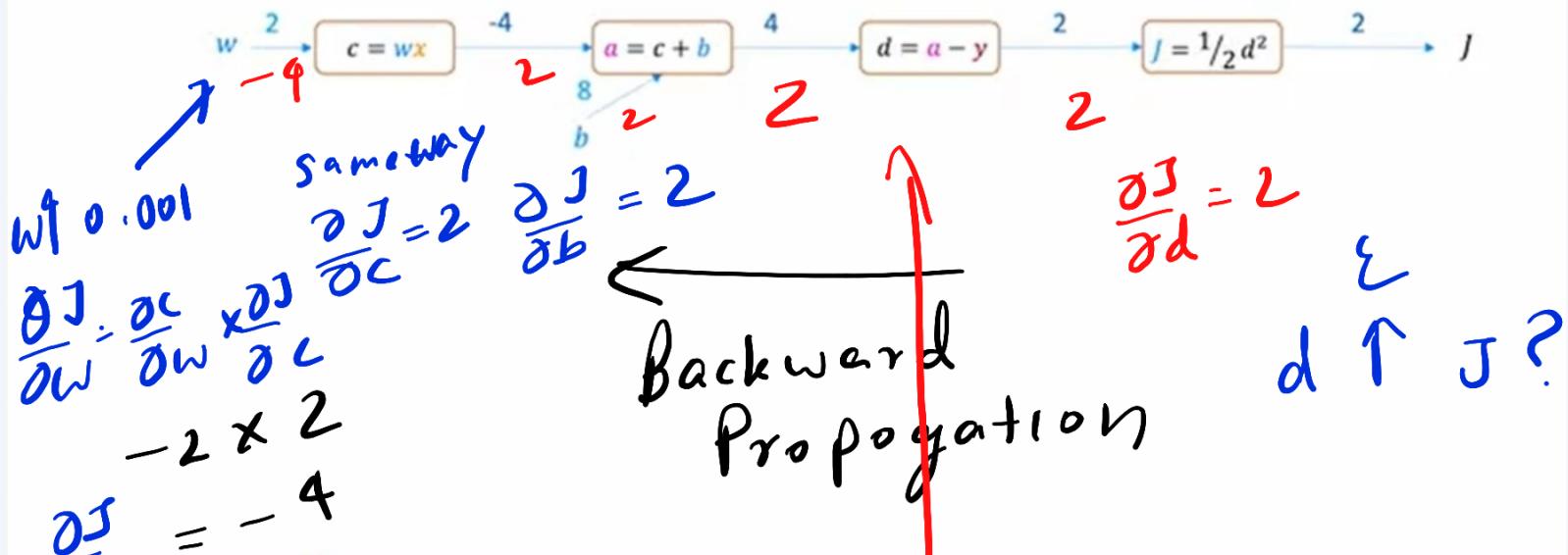
$$\rightarrow J(w) = 4.004001 \quad \uparrow 4 \times \epsilon_{\parallel}$$

$\frac{\partial}{\partial w} \rightarrow$ If w has many variables

sq error
Forward prop

Computing the Derivatives

$$w=2 \quad b=8 \quad x=-2 \quad y=2 \quad a=wx+b \quad J=\frac{1}{2}(a-y)^2$$



If w goes up by ϵ

B.P. compute derivatives in N+P steps rather than $N \times P$

if a goes up by ϵ

if $d \rightarrow 0.001$
but $J \rightarrow \frac{\partial J}{\partial d} = 2$
 $\uparrow 0.002$

$$\frac{\partial J}{\partial a} = 2 \parallel$$