

Olyse

✓

Machine Learning! "Field of study that gives computers the ability to learn without being explicitly programmed."

Algorithms:

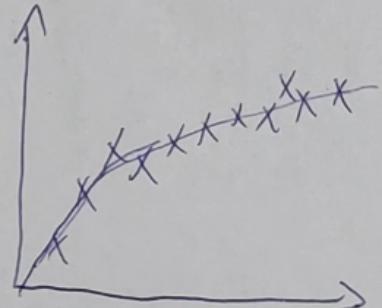
- Supervised Learning
- Unsupervised Learning
- Recommender systems
- Reinforcement Learning

Supervised Learning:

$$X \rightarrow Y$$

input              label

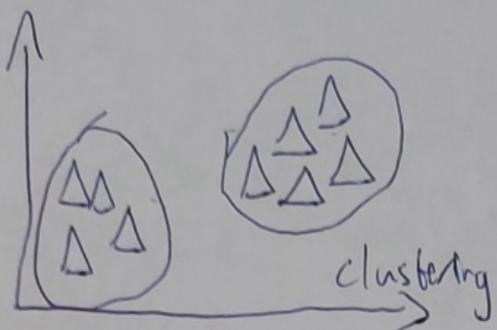
Regression: learning to predict a number out of infinitely many possible numbers



Classification) predicts categories, don't have to be numbers.

Categories are finite/limited

Unsupervised Learning: Find something interesting in unlabeled data



Think Google News

Anomaly Detection / find unusual data points  
/ outliers

Dimensionality Reduction

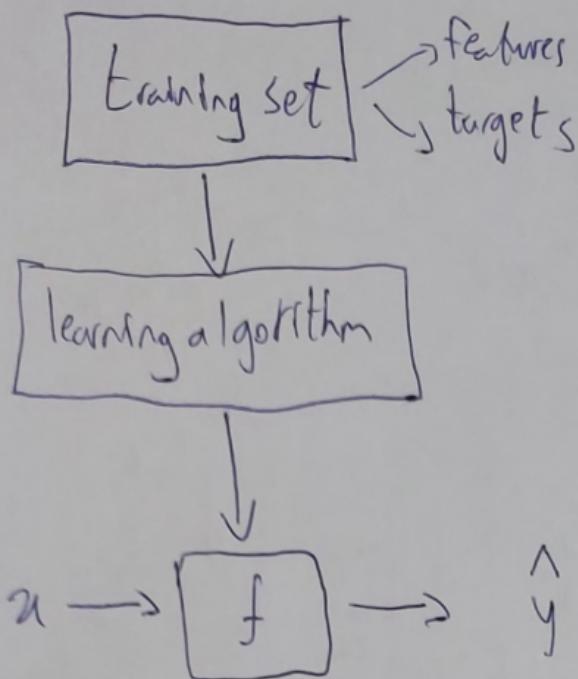
Compress data using fewer numbers

Notation:

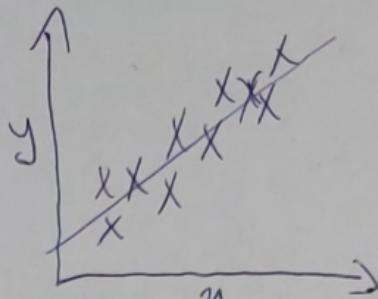
$x$  = "input" variable / feature

$y$  = "output" variable / target

$(x^{(i)}, y^{(i)})$  =  $i^{\text{th}}$  training example



feature      hypothesis/function  
 model  
 prediction  
 (estimated  $\hat{y}$ )  
 target



$$\text{If } f_{w,b}(x) = w_n x + b$$

and 1 variable

$\Rightarrow$  univariate linear regression

(Cost functions) used to find the best parameters for a model

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2 \quad // \text{Squared error cost function}$$

$m$  = number of training examples

Goal is to minimize  $J$  and find  $w$  and  $b$  corresponding to that

Gradient descent: algorithm to reach a minimum by taking the direction of steepest descent multiple times.

Repeat until convergence (simultaneously):

$$\begin{cases} w = w - \alpha \frac{\partial}{\partial w} J(w, b) \\ b = b - \alpha \frac{\partial}{\partial b} J(w, b) \end{cases}$$

$\alpha$  is the learning rate (determines how big our step will be)  
 $\alpha$  is a small positive number in  $(0, 1)$

Linear regression model

Cost function

$$f_{w,b}(n) = w_n + b$$

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

Gradient descent algorithm

$$\begin{aligned} w &= w - \alpha \boxed{\frac{\partial \mathcal{J}(w, b)}{\partial w}} \\ b &= b - \alpha \boxed{\frac{\partial \mathcal{J}(w, b)}{\partial b}} \end{aligned} \Rightarrow \begin{aligned} \frac{1}{m} \sum_{i=1}^m (f_{w,b}(n^{(i)}) - y^{(i)}) n^{(i)} \\ \frac{1}{m} \sum_{i=1}^m (f_{w,b}(n^{(i)}) - y^{(i)}) \end{aligned}$$

"Batch" gradient descent:

"Batch": Each step of gradient descent uses all the training examples

Multiple features:

$n_1$ Size	$n_2$ # Bedrooms	$n_3$ # Floors	$n_4$ Age	Price
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
...	...	...	...	...

$n_j$  = jth feature

$n$  = Number of features

$\vec{x}^{(i)}$  = features of the i<sup>th</sup> training example

$x_j^{(i)}$  = value of feature j in the i<sup>th</sup> training example

$$f_{\vec{w}, b}(\vec{u}) = w_1 u_1 + w_2 u_2 + \dots + w_n u_n + b$$

$\vec{w} = [w_1 \ w_2 \ w_3 \ \dots \ w_n]$  parameters of the model

$$\vec{u} = [u_1 \ u_2 \ u_3 \ \dots \ u_n]$$

$$f_{\vec{w}, b}(\vec{u}) = \vec{w} \cdot \vec{u} + b \quad \text{multiple linear regression}$$

To implement this, we use a neat trick called Vectorization.

Instead of manually typing the expression or using a for loop,  
just use  $f = \text{np.dot}(w, u) + b$

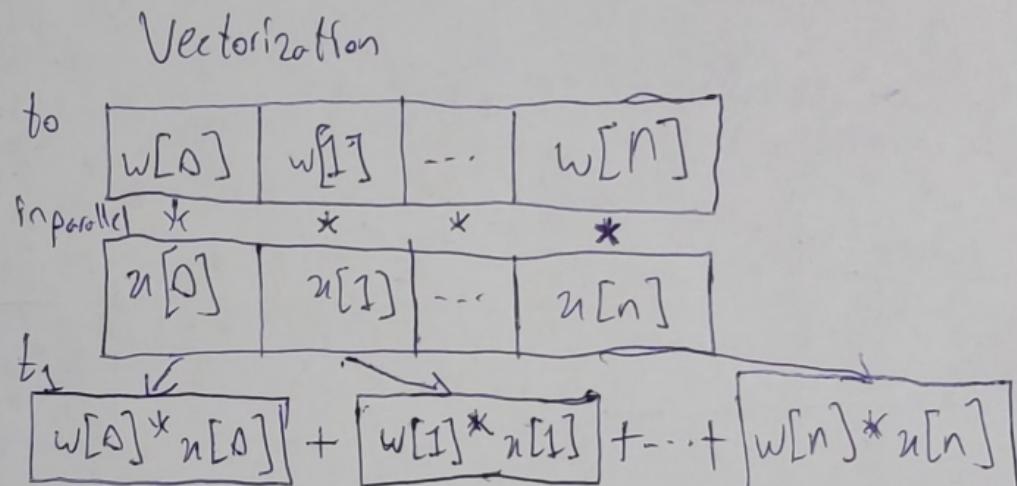
Without vectorization:

$$t_0: f + w[0] * u[0]$$

$$t_1: f + w[1] * u[1]$$

⋮

$$t_n: f + w[n] * u[n]$$



Utilizes computer hardware

Efficient + Scalable

Normal equation: Alternative to gradient descent for linear regression only. It is not generalized and can become less efficient for larger inputs.

Feature Scaling: Dividing by the max so that scaling  $\leq 1$

$$300 \leq u_1 \leq 2000 \quad 0 \leq u_2 \leq 5$$

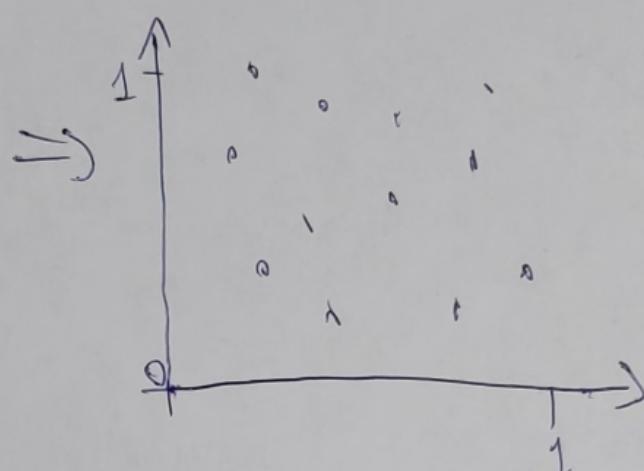


$$u_{1,\text{scaled}} = \frac{u_1}{2000}$$

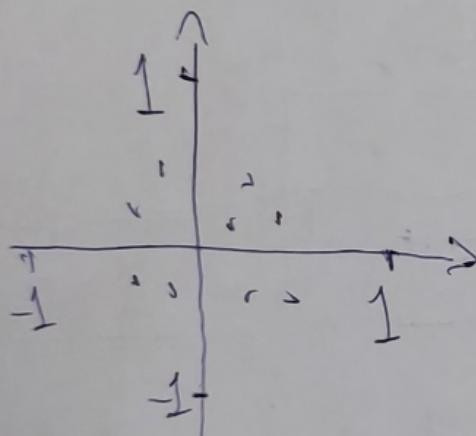
$$u_{2,\text{scaled}} = \frac{u_2}{5}$$

$$0.15 \leq u_{1,\text{scaled}} \leq 1$$

$$0 \leq u_{2,\text{scaled}} \leq 1$$



Mean normalization: rescaling and centering around 0.



$$\text{then form: } u_{\text{rescaled}} = \frac{u - \mu}{\text{max} - \text{min}}$$

$$Z\text{-score normalization: } u_{\text{scaled}} = \frac{u - \mu}{\sigma}$$

- Aim for about  $-1 \leq u_j \leq 1$  for each feature  $u_j$

Why should we use it?

Let  $\widehat{\text{price}} = w_1 u_1 + w_2 u_2 + b$

$\downarrow \quad \downarrow$   
size      #bedrooms

$u_1:$  ranges 300-2000       $u_2:$  ranges 0-5  
large      small

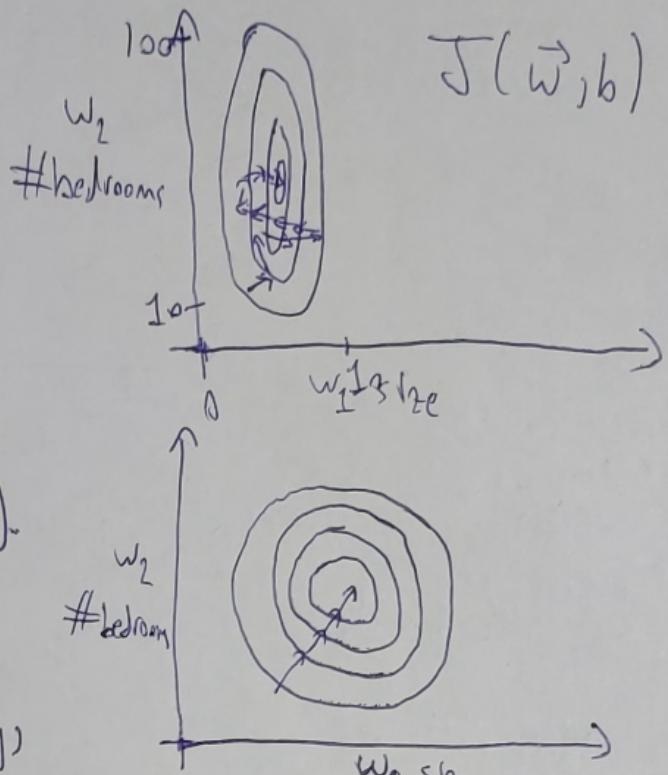
In a contour plot:

The contours form ellipses and are shorter on one side and longer on the other.

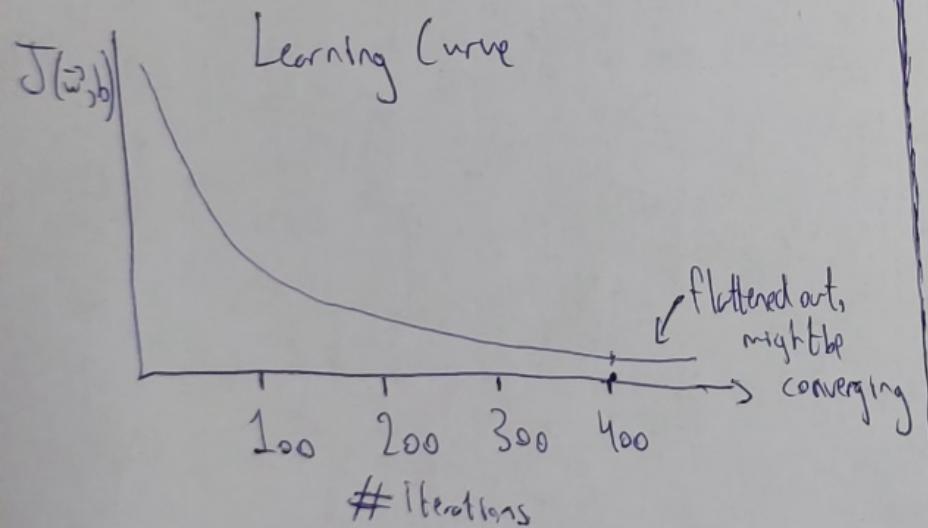
This is because a very small change to  $w_1$  can have a very large impact on  $J(w, b)$ . ( $w_1$  is multiplied by a very large number (size of house)).

It is the opposite for  $w_2$ .

Because the contours are tall and skinny, gradient descent may bounce back and forth for a long time before it can find the way to the global min.



Checking Gradient Descent for convergence:

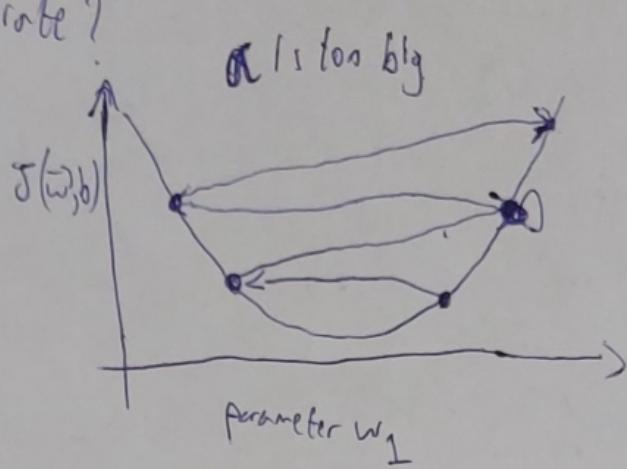
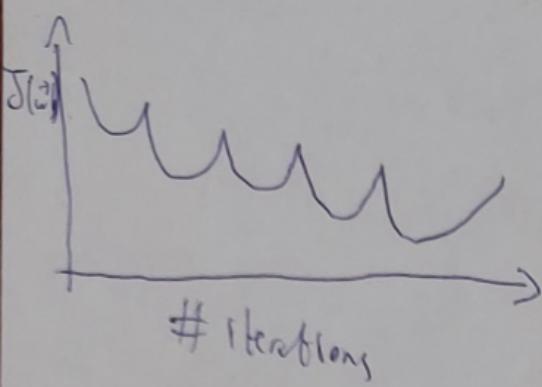


Automatic convergence test

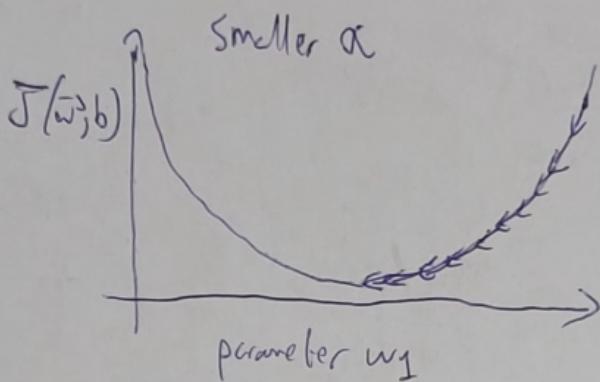
Let  $\epsilon = 10^{-3}$

If  $J(\vec{w}, b)$  decreases by  $\leq \epsilon$  in one iteration, declare convergence

How to choose the learning rate?



overshooting the minimum

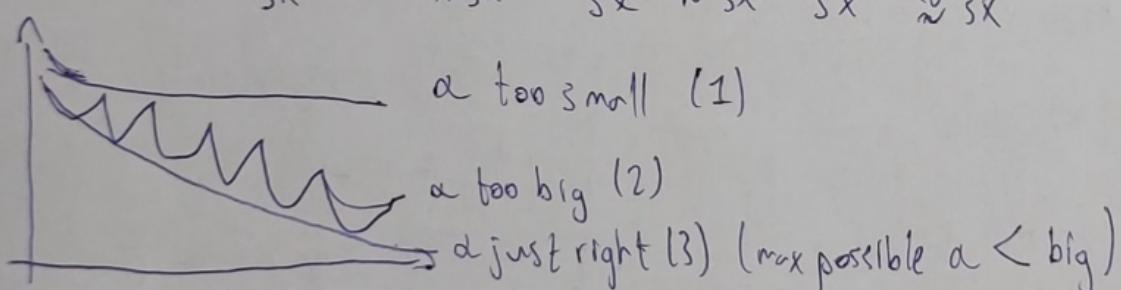


If the learning curve increases, then either  $\alpha$  is too big or there is a bug in the code.

Example:  $w_1 = w_1 + \alpha d_1$

Values of  $\alpha$  to try:

$$0.001 \xrightarrow{\approx 3x} 0.003 \xrightarrow{\approx 3x} 0.01 \xrightarrow{3x} 0.03 \xrightarrow{\approx 3x} 0.1 \xrightarrow{3x} 0.3 \xrightarrow{\approx 3x} 1$$



Polynomial Regression:

Just use Linear Regression with  $x^1, x^2, \dots, x^n$ .

# Classification

class = category

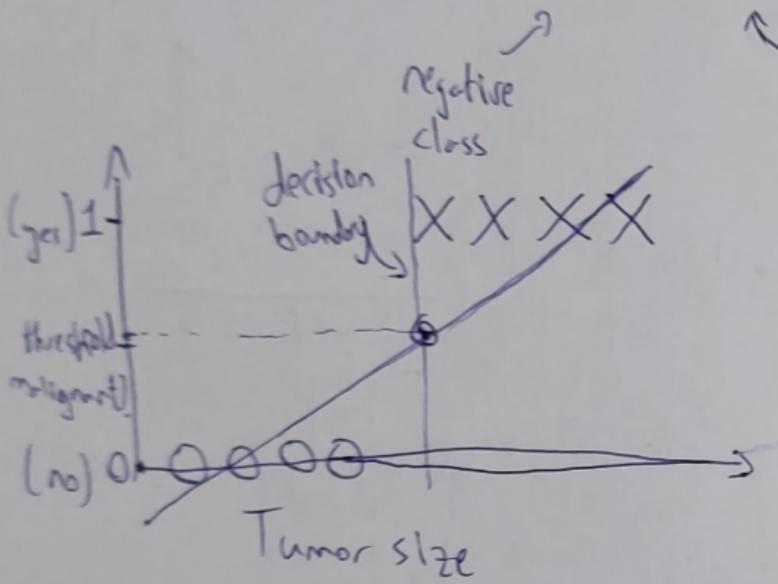
binary classification when you have 2 values for  $y$ .

no      yes

false      true

0      1

↑  
positive  
class

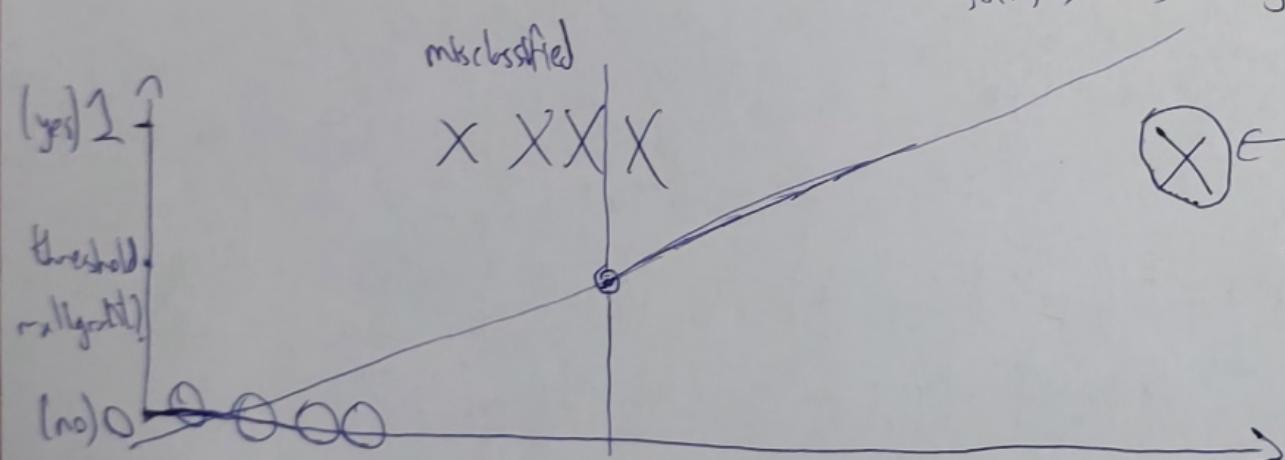


$$f_{w,b}(x) = w*x + b \text{ by linear regression.}$$

linear regression outputs numbers greater than 0, less than 1, and even greater than 1. Let's define a threshold (0.5).

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

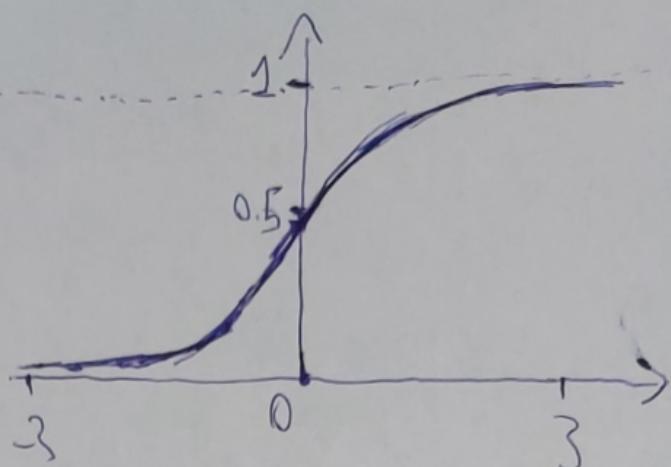
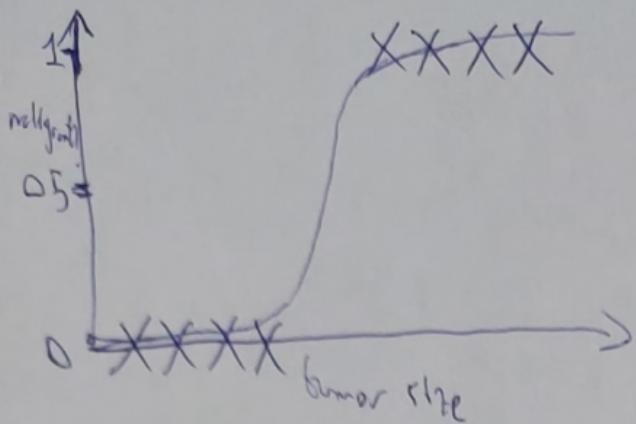
$$\text{if } f_{w,b}(x) \geq 0.5 \rightarrow \hat{y} = 1$$



Another example  
Would change  
the line of best fit

Therefore linear regression is not reliable

# Logistic Regression:



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

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

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

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

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

$$\text{So } g(\underbrace{\vec{w} \cdot \vec{x} + b}_z) = \frac{1}{1+e^{-(\vec{w} \cdot \vec{x} + b)}}$$

"logistic regression"

Interpretation of the output:  $f_{\vec{w}, b}(\vec{x}) = 0.7$

$\Rightarrow 70\%$  chance that  $y$  is 1

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

Probability that  $y$  is 1, given input  $\vec{x}$ , parameters  $\vec{w}, b$

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

Is  $f_{\vec{w}, b}(\vec{x}) \geq 0.5$ ?

Yes:  $\hat{y} = 1$       No:  $\hat{y} = 0$

When is  $f_{\vec{w}, b}(\vec{x}) \geq 0.5$ ?

$$g(z) \geq 0.5$$

$$z \geq 0$$

$$\vec{w} \cdot \vec{x} + b \geq 0$$

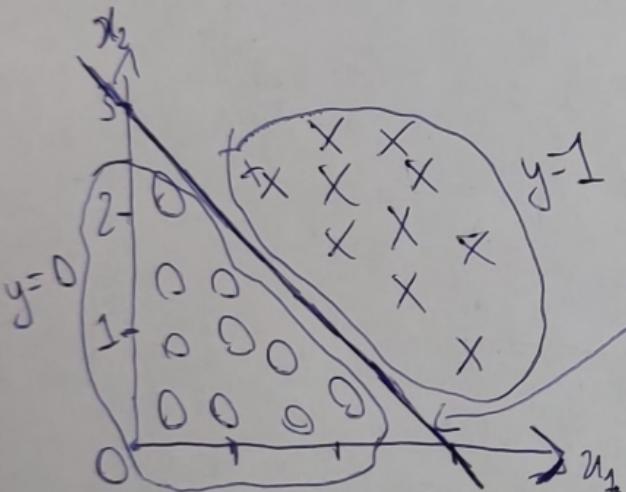
$$\hat{y} = 1$$

$$\vec{w} \cdot \vec{x} + b < 0$$

$$\hat{y} = 0$$

$$f_{\vec{w}, b}(\vec{x}) = g(z) = g(w_1 u_1 + w_2 u_2 + b)$$

1      1      -3



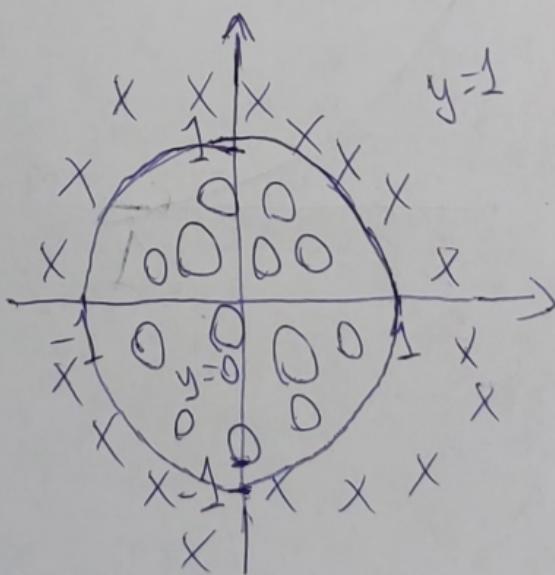
$$\text{Decision boundary: } z = u_1 + u_2 - 3 = 0$$

$$u_1 + u_2 = 3$$

$$f_{\vec{w}, b}(\vec{x}) = g(z) = g(w_1 u_1^2 + w_2 u_2^2 + b)$$

$$z = u_1^2 + u_2^2 - 1$$

$$\text{decision boundary: } z = 0 \Rightarrow u_1^2 + u_2^2 = 1$$

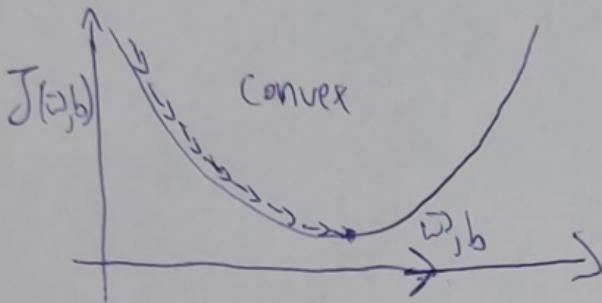


Squared Error cost for linear regression:

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

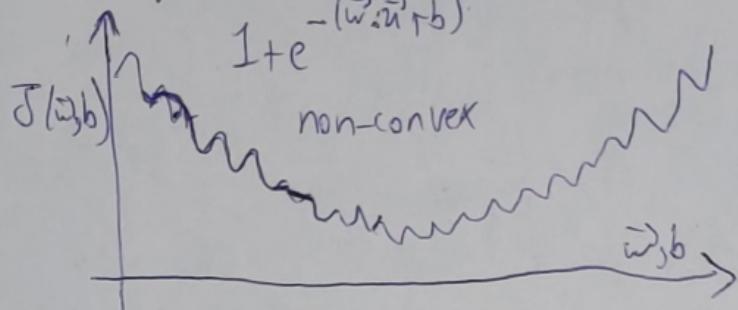
Linear Regression

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



Loss  $L(f_{\vec{w}}, b(\vec{x}^{(i)}), y^{(i)})$

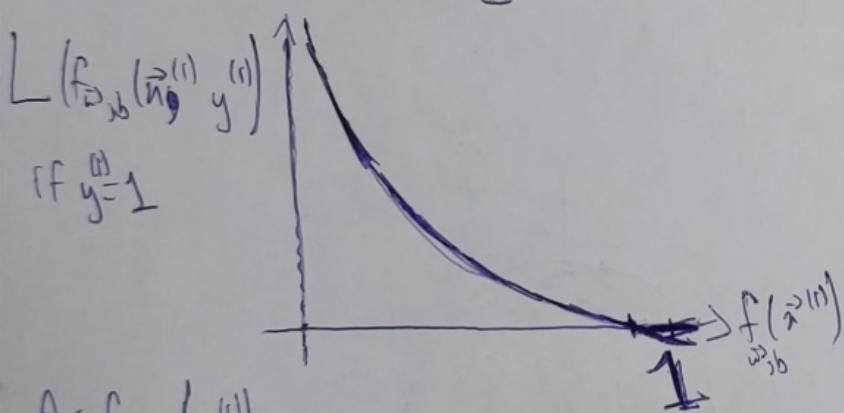
$$f_{\vec{w}}, b(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$



"Wiggly" with many local minima so gradient descent won't work. So squared error cost is not a good choice for logistic regression.

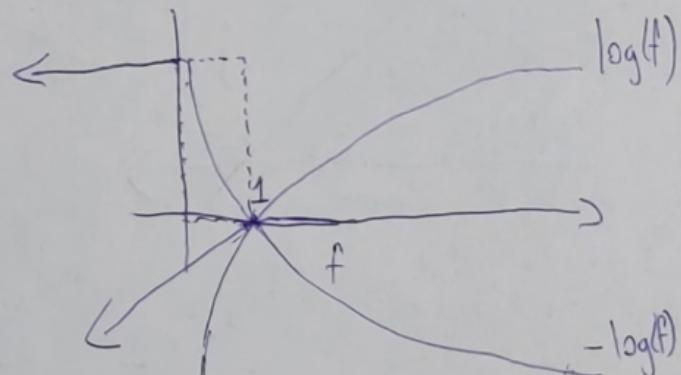
We need a new loss function.

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

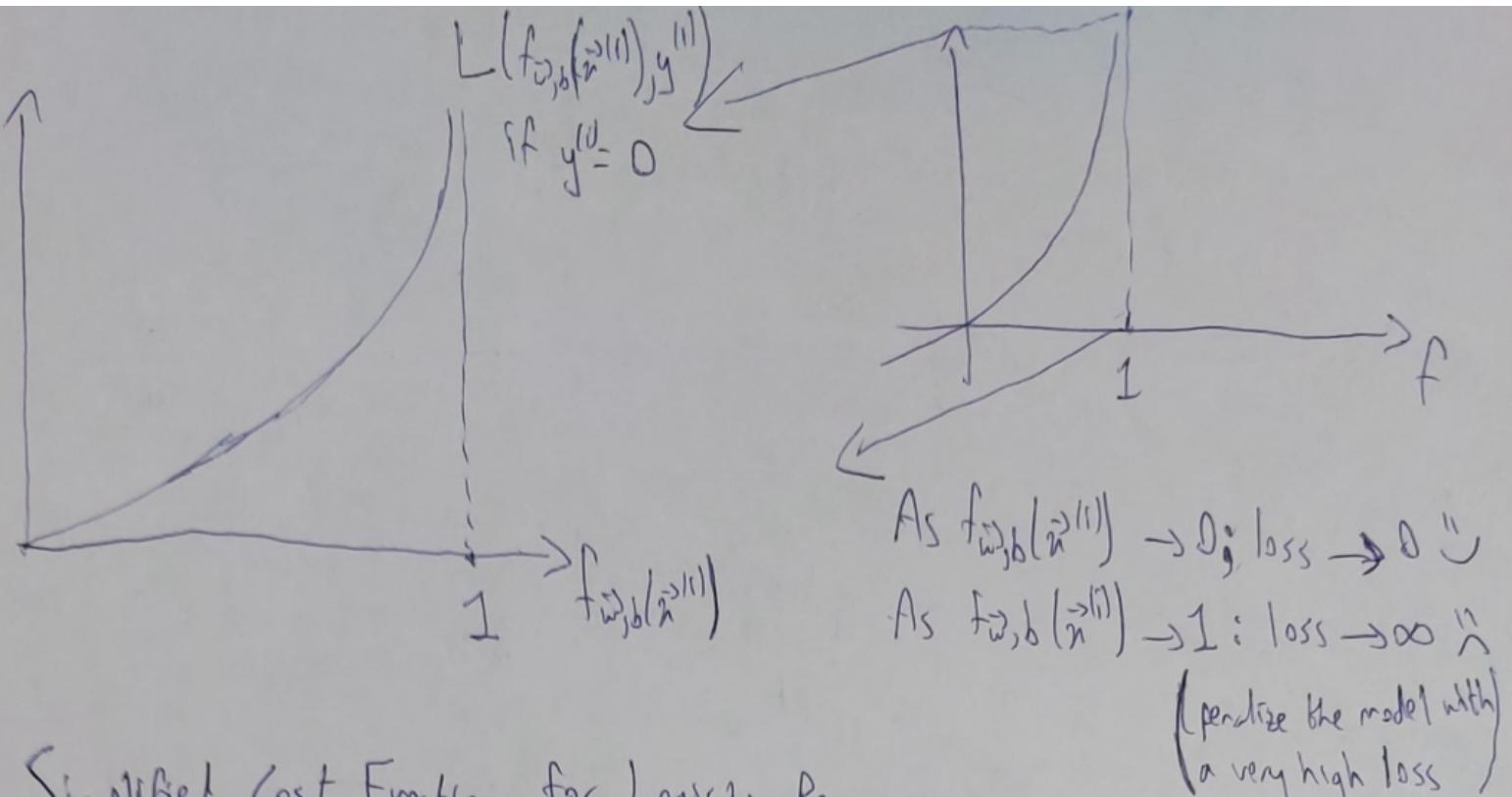


As  $f_{\vec{w}}, b(\vec{x}^{(i)}) \rightarrow 1$ : loss  $\rightarrow 0$

As  $f_{\vec{w}}, b(\vec{x}^{(i)}) \rightarrow 0$ : loss  $\rightarrow \infty$



$f$  = output of logistic regression  $E(0,1)$



Simplified Cost Function for Logistic Regression:

Because  $y$  can be either 0 or 1, and cannot take on values in between, we can rewrite the cost function as:

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = -y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) - (1-y^{(i)}) \log(1-f_{\vec{w}, b}(\vec{x}^{(i)}))$$

Therefore:

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

$$\frac{\partial}{\partial w_j} J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)} \quad \leftarrow j^{\text{th}} \text{ feature of training example } i$$

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

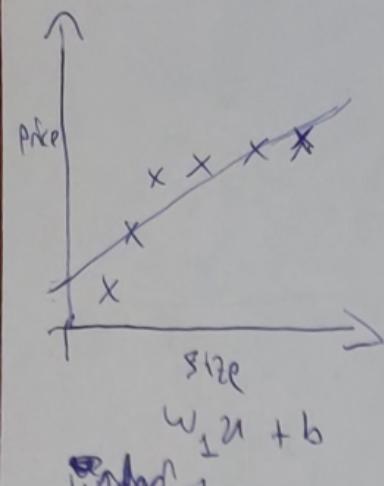
repeat {

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

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

} simultaneous updates

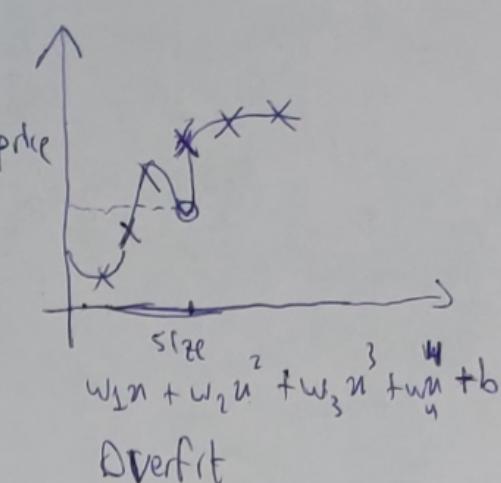
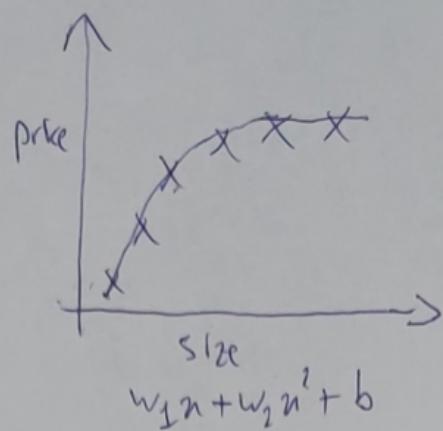
Looks the same as linear regression!



Does not fit the training set well

high bias; it has a preconception that the pattern is linear despite the data showing otherwise

"High bias"



Fits the training set extremely well (cost=0)

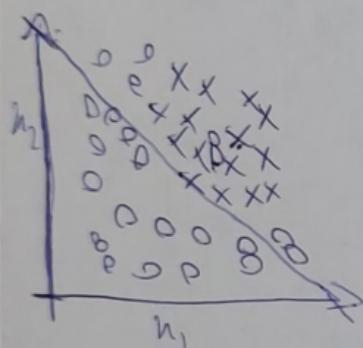
But it's wiggly and might predict data incorrectly such as a house of a larger size being cheaper than one with a smaller size.

It looks like the model has fit the training data too well and can't generalize to new examples

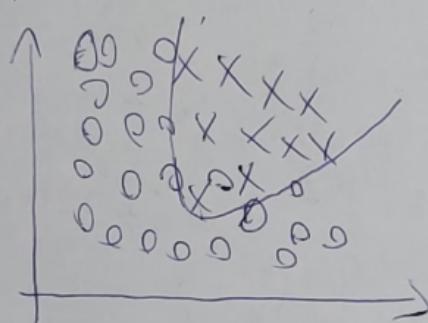
"High variance"

Two models with slightly different training sets might end up with totally different/highly variable predictions

Applies to classification as well

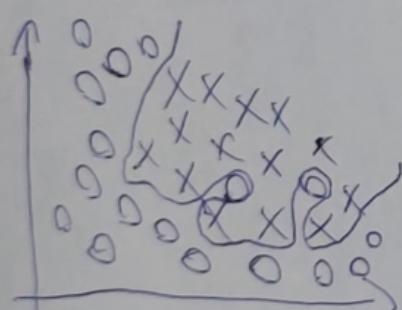


$z = w_1 n_1 + w_2 n_2 + b$   
underfit high bias



$z = w_1 n_1 + w_2 n_2 + w_3 n_1^2 + w_4 n_2^2 + w_5 n_1 n_2 + b$   
just right

Overfitting



$z = w_1 n_1 + w_2 n_2 + w_3 n_1^2 n_2 + w_4 n_1^2 n_2^2 + w_5 n_1^2 n_2^3 + w_6 n_1^3 n_2 + \dots + b$

## Addressing Overfitting:

Options:

1. Collect more data

2. Select features

- Feature selection

3. Reduce size of parameters

- Regularization

Regularization: penalize for high/large parameters in the cost function

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

can exclude

$\downarrow$

fit data

Regularization term  
keep  $w_j$  small

$\lambda$  balances both goals

$\lambda$ : regularization parameter

If  $\lambda = 0$ : overfits

If  $\lambda = 10^{10}$  (large number): underfits because model focuses too much on making  $w_j \approx 0$  so fit to b  
So we'll have a horizontal line.

Gradient descent:

repeat {

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

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

} simultaneous update

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

$$\frac{\partial}{\partial b} J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) \quad (\text{because we don't regularize } b)$$

$$\begin{aligned} \frac{\partial}{\partial w_j} J(\vec{w}, b) &= 1_{w_j} - \alpha \frac{\lambda}{m} w_j - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) u_j^{(i)} \\ &= w_j \left(1 - \alpha \frac{\lambda}{m}\right) - \underbrace{\alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) u_j^{(i)}}_{\text{usual update}} \end{aligned}$$

$$\alpha \frac{\lambda}{m}$$

$$0.01 \cdot \frac{1}{50} = 0.0002 \quad \text{then} \quad w_j \left(1 - \alpha \frac{\lambda}{m}\right) = 0.9998 w_j$$

This is what regularization does. This is how it works: by multiplying  $w_j$  by a number slightly  $< 1$ , effectively shrinking it every iteration.

End of Course One