

Week 3 : Classification

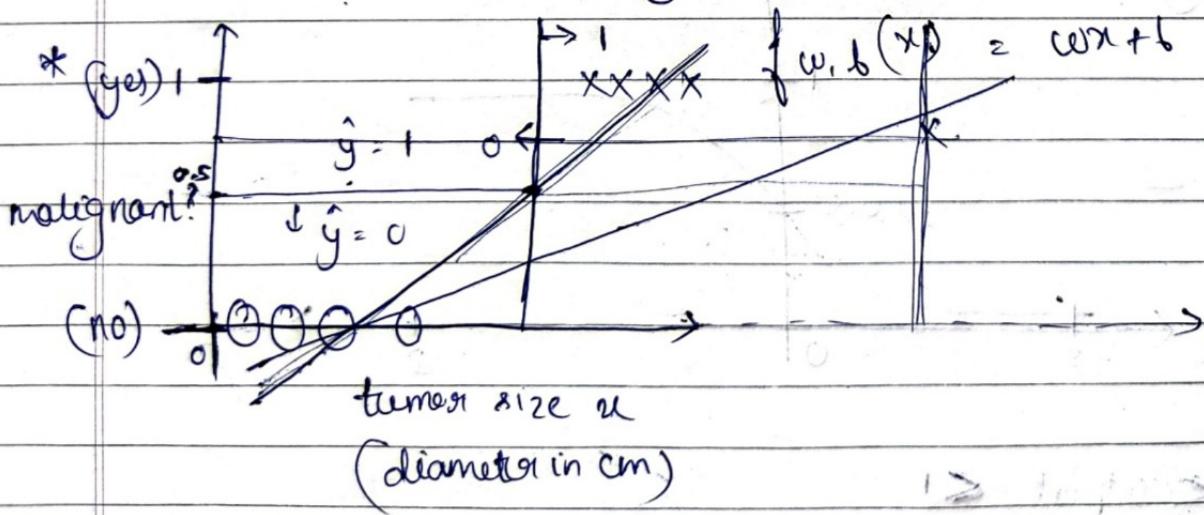
⇒ Classification with logistic regression

Video 1: Motivation

Question	Answer "y"	
Is this email spam ?	no	yes
Is the transaction fraudulent ?	no	yes
Is the tumour malignant ?	no (negative) (absent)	yes 1 (true) (presence)

* y can be only one of 2 values.

* "binary classification"
class = category



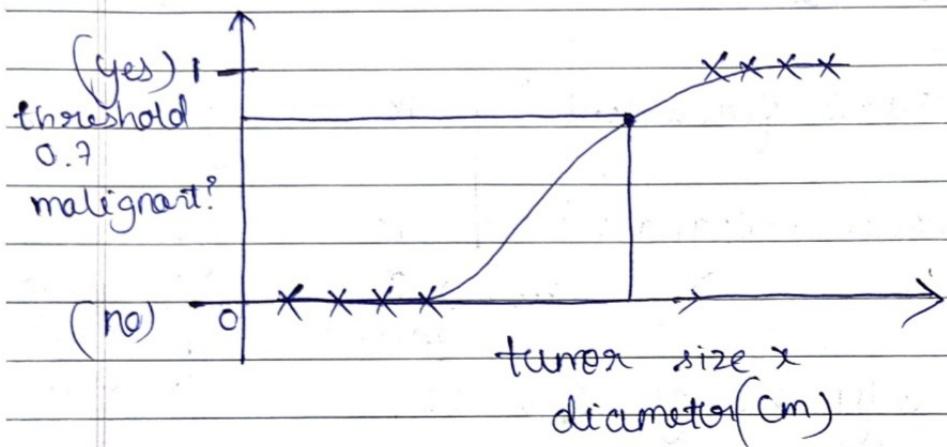
0 benign
x malignant

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

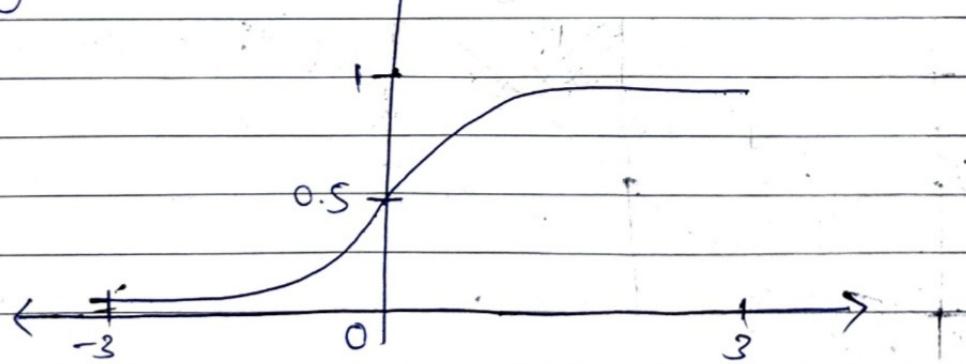
$$\text{if } f_{w,b}(x) \geq 0.5 \rightarrow \hat{y} = 1$$

Logistic regression is used to solve binary classification problems where the output is either 0 & 1

Ideo2: Logistic Regression



Sigmoid function / logistic function



Output ≤ 1

$$g(z) = \frac{1}{1+e^{-z}} \quad 0 < g(z) < 1$$

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

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

↓

z

↓

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

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

Interpretation of logistic regression output.

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

probability that the class is 1

Example:

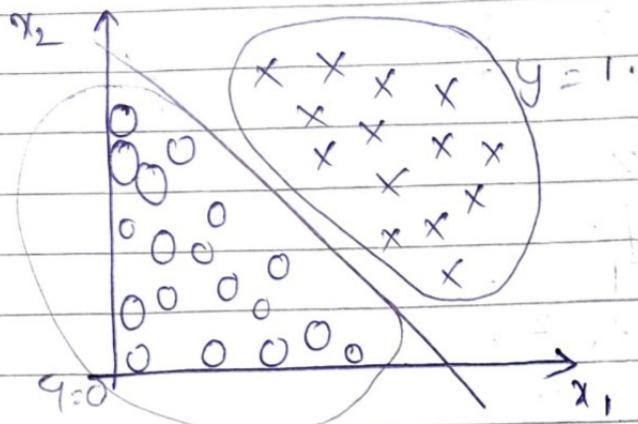
x is "tumor size"

y is 0
or 1

$$f_{\vec{w}, b}(\vec{x}) = 0.7$$

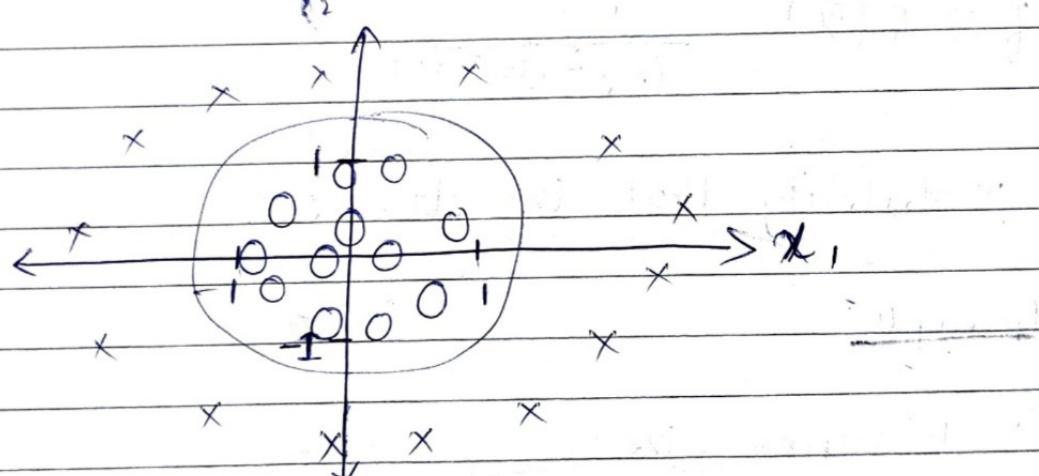
70% chance that y is 1

Videos: Decision Boundary



$$f(\vec{w}, b) (\vec{x}) = g(z) = g(w_1 x_1 + w_2 x_2 + b)$$

Non linear decision boundaries



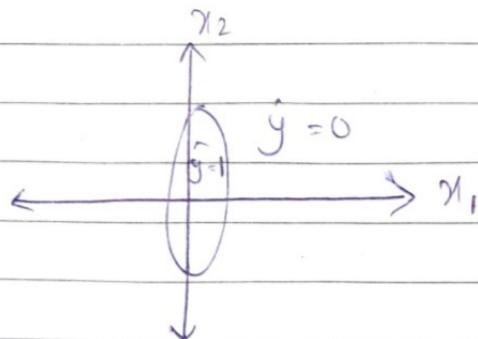
$$z = w_1 x_1^2 + w_2 x_2^2 + b$$

$$f_{\vec{w}, b} (\vec{x}) = g(z) = g(w_1 x_1^2 + w_2 x_2^2 + b)$$

decision boundary $= \frac{x_1^2 + x_2^2 - 1}{x_1^2 + x_2^2 - 1} = 0$

$$x_1^2 + x_2^2 < 1 \\ (1 = 0)$$

Non-Linear Decision Boundaries



$$f_{\vec{w}, b}(\vec{x}) = g(z) = g(w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_1 x_2 + w_5 x_2^2)$$

⇒ Cost function for Logistic Regression

Training Set

tumor size (cm)	patient's age x_n	malignant ↓
10	52	1
2	73	0
5	55	0
12	49	1

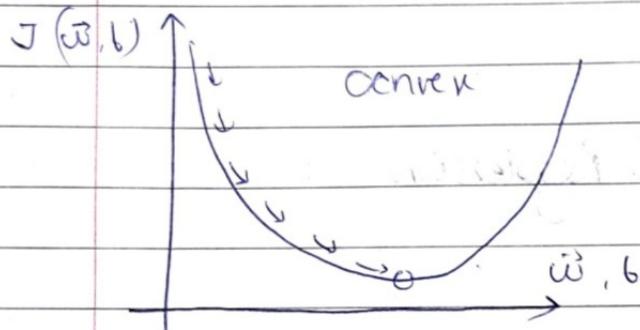
⇒ target y is 0 or 1.

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

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

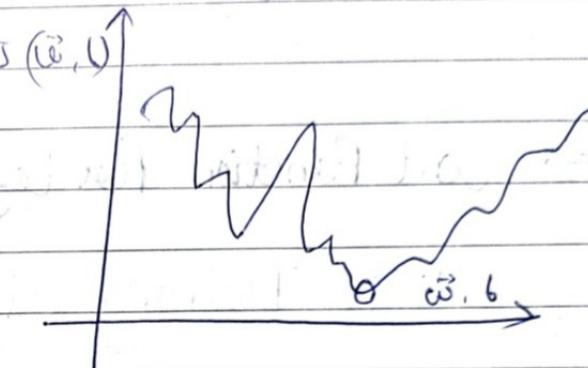
linear regression

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



logistic regression

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

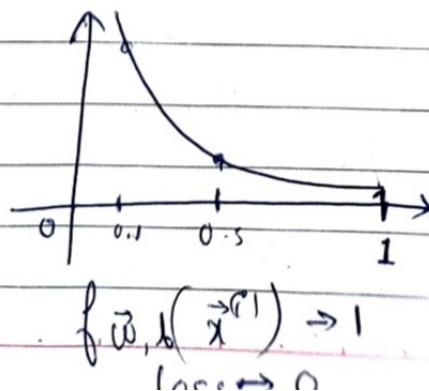
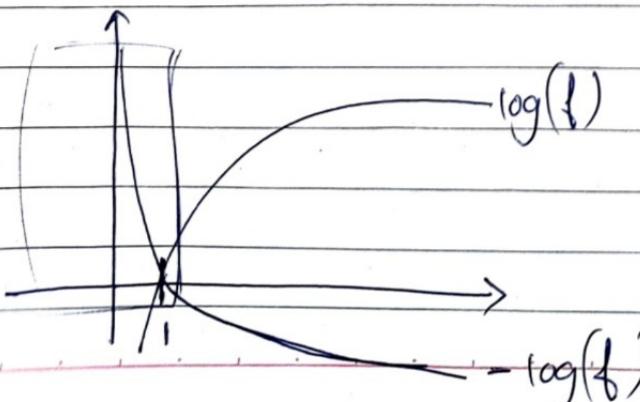


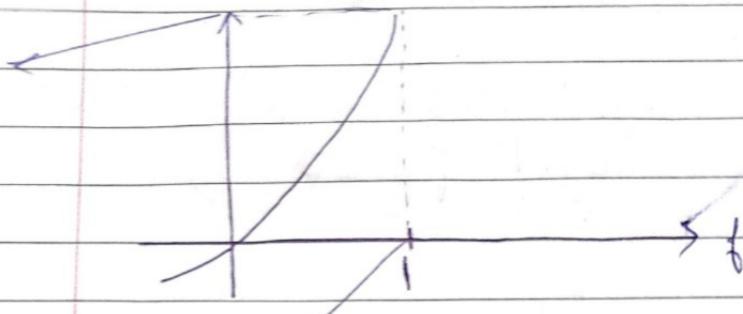
$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m \left[\frac{1}{2} f_{\vec{w}, b}(\vec{x}^{(i)} - y^{(i)})^2 \right]$$

loss on a single training example
function of x and true label y

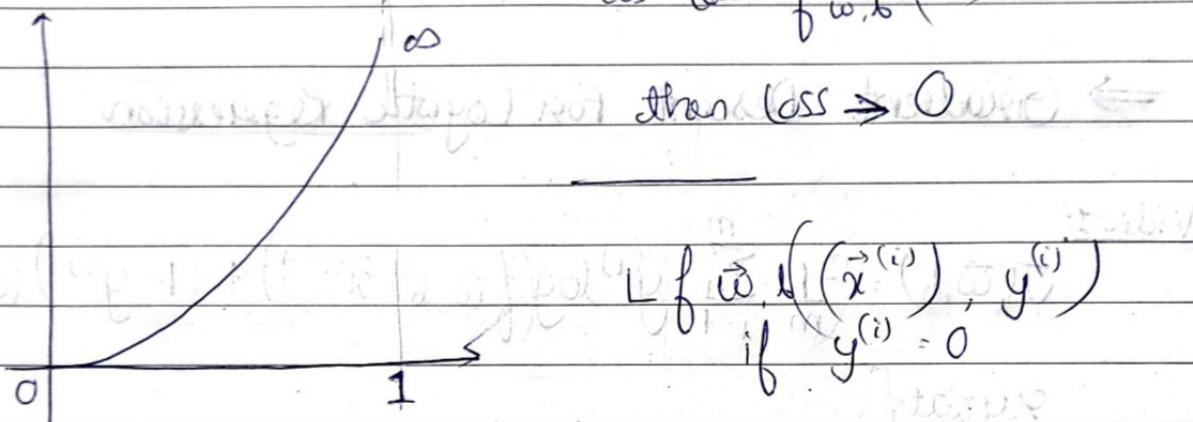
$$L = L(f_{\vec{w}, b}(\vec{x}^{(i)}, y^{(i)}))$$

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





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



Due to this cost function will be convex & will reach global minimum.

Video 2: Simplified cost function for Logistic Regression

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

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

loss

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

Cost function:

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

\Rightarrow Gradient Descent for Logistic Regression

Video 1:

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

repeat

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

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

3. Compute step by step loss function calculation

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

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

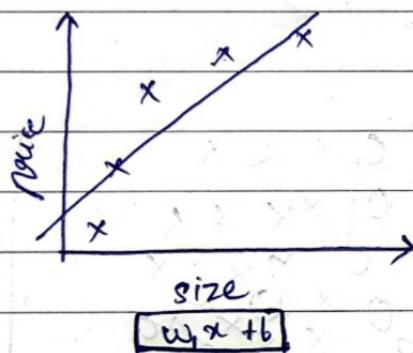
Linear regression: $f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$

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

⇒ The problem of overfitting

Video 1: The problem of overfitting

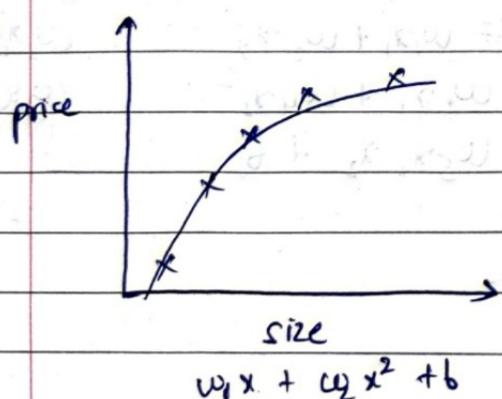
- * Regularisation - minimises overfitting
- * Regression example:



underfit: does not fit the training set well

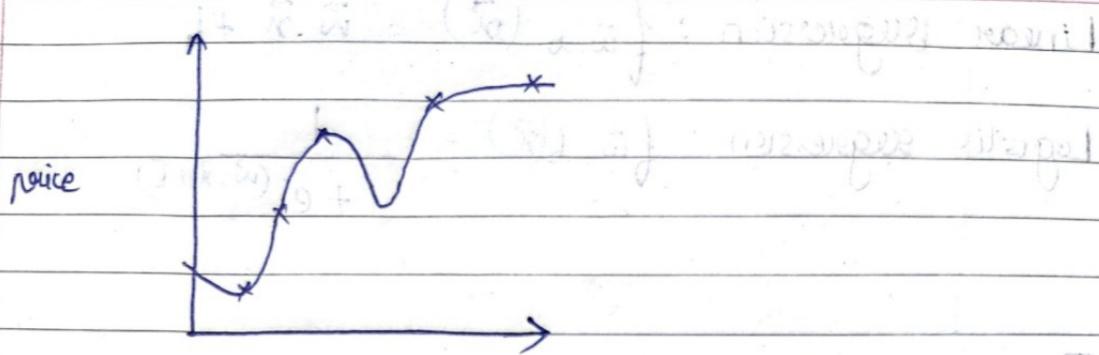
high bias

- * here there is preconception that the housing prices are going to a completely linear function of size.



generalisation

fits training set pretty well

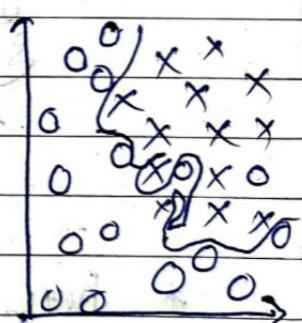
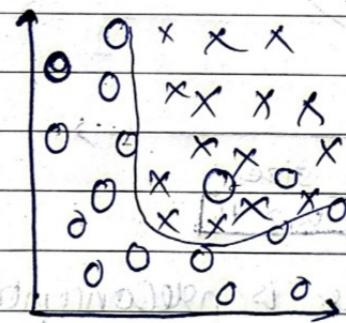
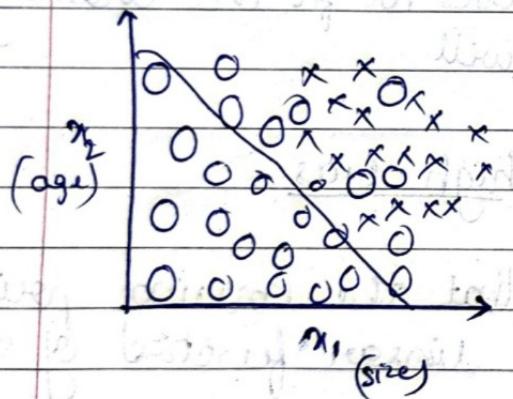


$$w_0 + w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + b$$

overfit / high variance

- * fits the training data too well.

Classification



$$z = w_0 x_0 + w_1 x_1 + b$$

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

$$z = w_0 x_0 + w_1 x_1 + b$$

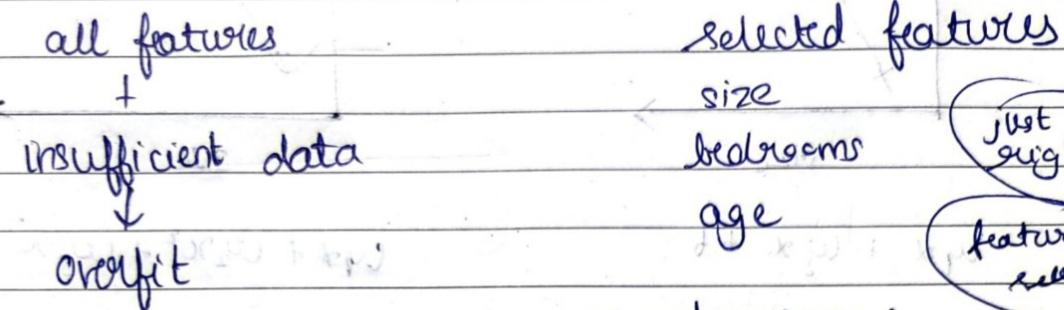
$$z = w_0 x_0 + w_1 x_1 + w_2 x_1^2 + w_3 x_1 x_2 + w_4 x_2^2 + w_5 x_2 + b$$

$$z = w_0 x_0 + w_1 x_1 + w_2 x_1^2 + w_3 x_1 x_2 + w_4 x_2^2 + w_5 x_2 + b$$

g is the sigmoid function

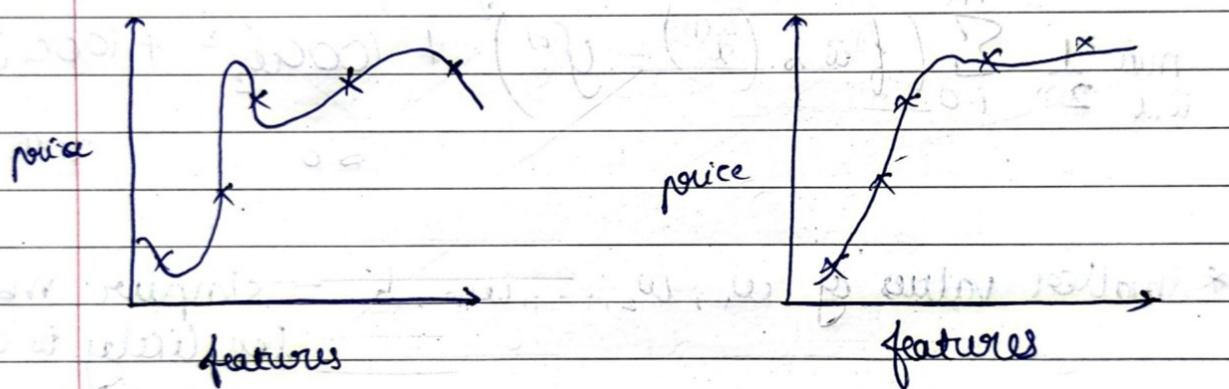
Video 2: Addressing Overfitting

- * collect more training data
- * use fewer features.



disadvantage:

(+) You will lose useful features could be lost



$$f(x) = 28x - 385x^2 + 39x^3 - 176x^4 + 100$$

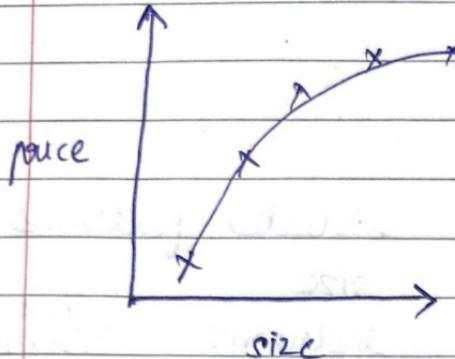
$$f(x) = 13x - 0.23x^2 + 0.000014x^3 - 0.0001x^4 + 10$$

smaller values

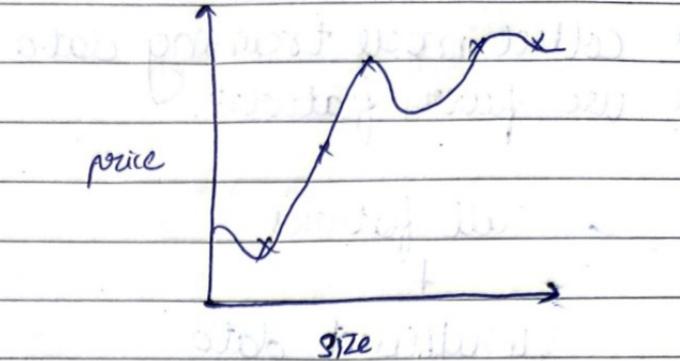
- * it makes very less difference if you regularise to or not

less regularization will result in overfitting

Video 3: Cost function with regularisation



$$w_0 + w_1 x + w_2 x^2 + b$$



$$w_0 + w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + b$$

make w_2, w_3, w_4 really small (≈ 0)

$$\min_{\vec{w}, b} \frac{1}{2m} \sum_{i=0}^m \left(f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)} \right)^2 + 1000 w_2^2 + 1000 w_3^2 + \dots$$

≈ 0 ≈ 0

* Smaller values of w_0, w_1, \dots, w_n, b - simpler model
less likely to overfit

* # regularise all the parameters

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

n = no. of features

λ = regularisation parameter

↑
makes
little
difference

like a, λ is a hyperparameter

÷ by $2m$ because the terms are scaled uniformly

* if $\lambda = 0$, we aren't regularising.
overfitting

* If $\lambda = 10^{10}$ (large large number)

$$f_{\vec{w}, b}(\vec{x}) \approx \begin{cases} w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + b \\ \approx 0 \quad \approx 0 \quad \approx 0 \quad \approx 0 \end{cases}$$

$$f(x) = b$$

does not fit



Video 4: Regularised Linear Regression

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

Gradient Descent

repeat {

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

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

don't have to regularise

Video 5: Regularised logistic regression

$$Z = w_1 x_1 + w_2 x_2 + w_3 x_1^2 x_2 + w_4 x_1^2 x_2^2 + w_5 x_1^2 x_2^3 + \dots$$

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

Cost function

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