

Week 03

01] Classification

* Classification is a type of supervised machine learning where the goal is to predict a label or category (class) given input data.

ex :-

- .) Email spam detection
- .) Fraud transaction detection

Types of classification

1] Binary Classification

- * Two possible outcomes
- * 0 (False) is the negative class (absence)
- * 1 (True) is the positive class (presence)

2] Multi - Class Classification

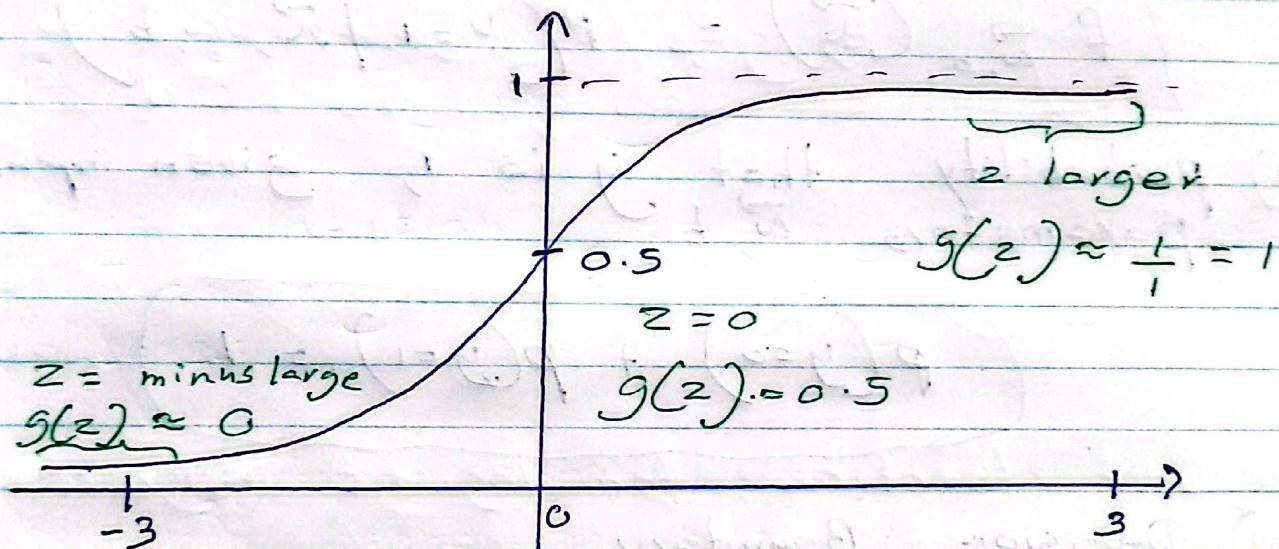
*) More than two categories

3] Multi - label classification

*) One item can belong to multiple categories
(A movie tagged as Comedy and romance)

02] Logistic Regression

- * Despite its name, Logistic Regression is actually used for classification, not regression.
- * It is a supervised learning algorithm used to predict discrete outcomes, like whether something belongs to class 0 or class 1.



want outputs between 0 and 1

Sigmoid function
(Logistic function)

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

$$f_{\vec{\omega}, b}(\vec{x}), \quad z = \vec{\omega} \cdot \vec{x} + b$$

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

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

Interpretation

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

"probability" that class is 1

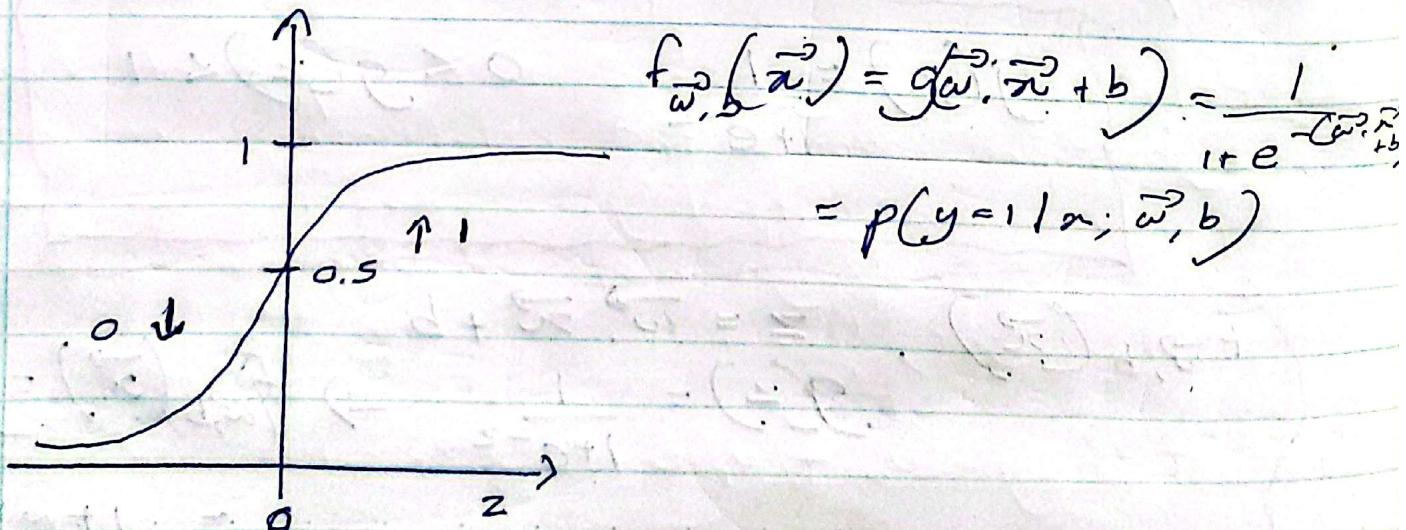
$$f_{\vec{w}, b}(\vec{x}) = P(y=1 | \vec{x}; \vec{w}, b)$$

- * probability that y is 1, given input \vec{x} , parameters \vec{w}, b

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

03) Decision Boundary

- * A decision boundary is an imaginary line (or surface) that separates different classes in your data.



IS $f_{\vec{\omega}, b}(\vec{x}) \geq 0.5$?

Threshold

yes: $\hat{y} = 1$

No: $\hat{y} = 0$

when is $f_{\vec{\omega}, b}(\vec{x}) > 0.5$?

$$g(z) \geq 0.5$$

$$z \geq 0$$

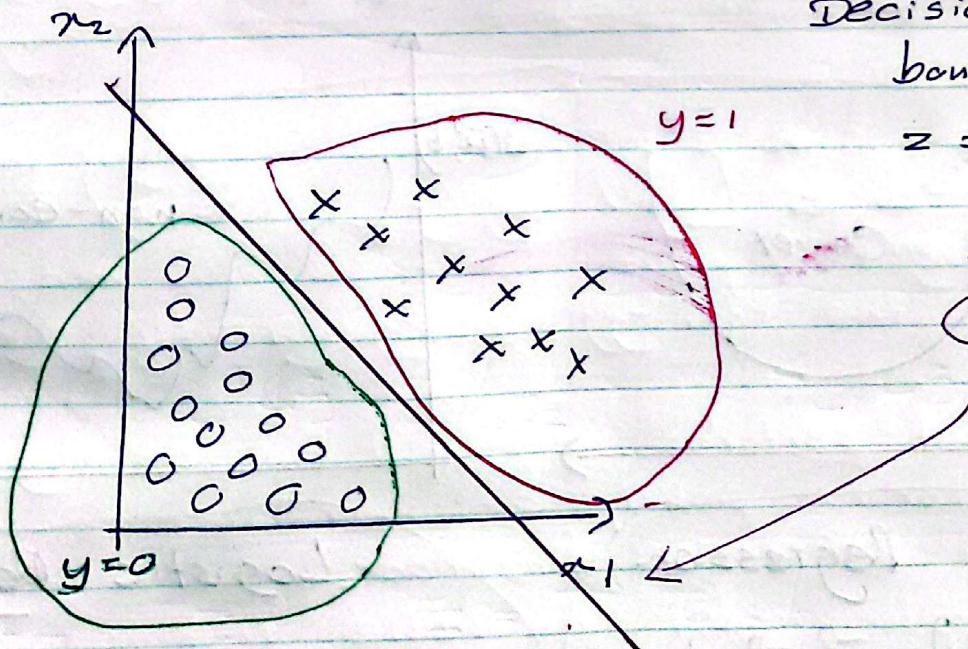
$$\vec{\omega} \cdot \vec{x} + b \geq 0$$

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

$$\hat{y} = 1$$

$$\hat{y} = 0$$

$$\Rightarrow g(z) = g(x_1 + x_2 - 3)$$



Decision

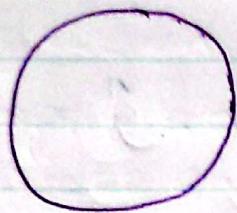
boundary \Rightarrow

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

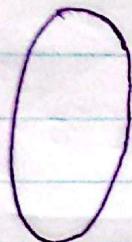
$$x_1 + x_2 - 3 = 0$$

$$x_1 + x_2 = 3$$

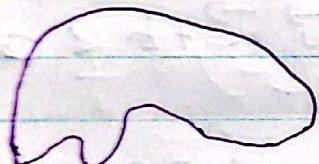
Non-linear Decision boundaries



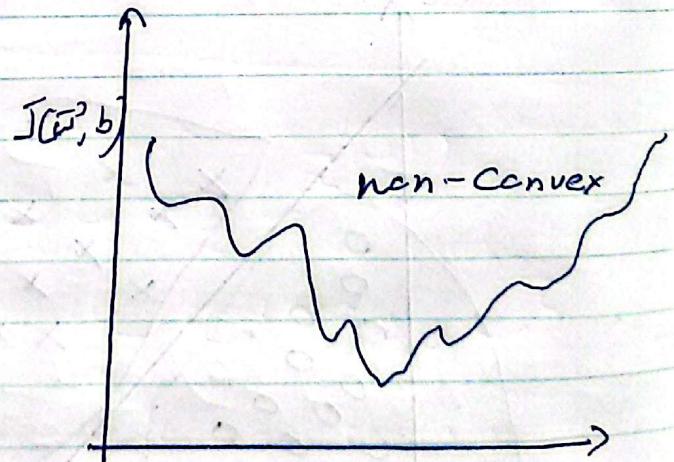
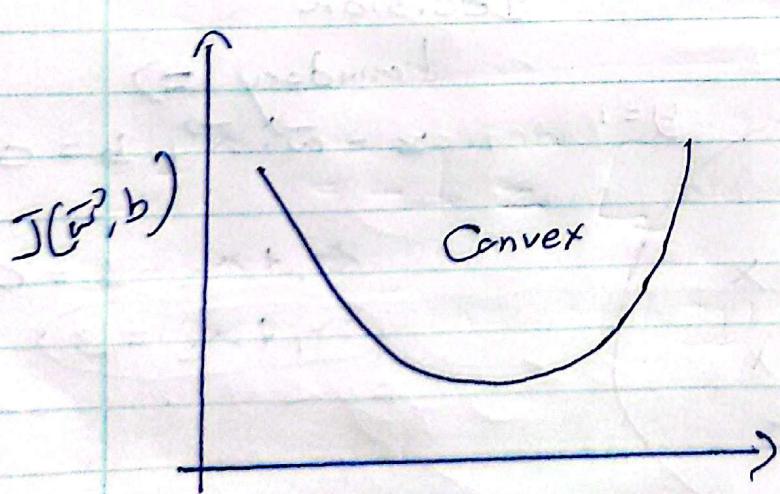
$$\Rightarrow g(\omega_1 x_1 + \omega_2 x_2 + b)$$



$$\Rightarrow g(z) = g(\omega_1 x_1 + \omega_2 x_2 + \omega_3 x_1^2 + \omega_4 x_1 x_2 + \omega_5 x_2^2 + \omega_6 x_1^3 + b)$$



04] Cost function for logistic Regression



Linear Regression

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

Logistic Regression

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

Cost

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

Loss

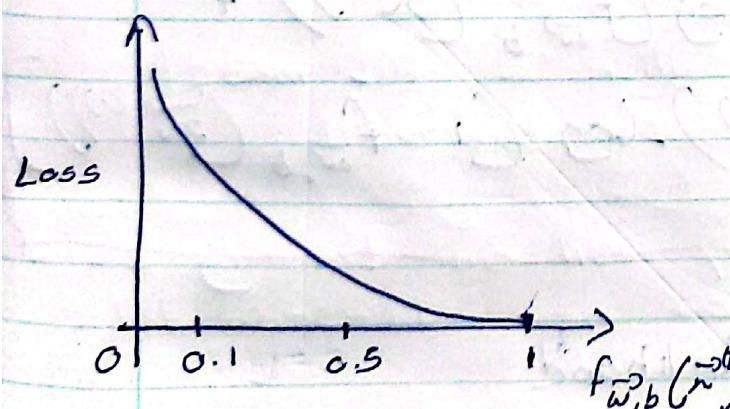
$$L(f_{\vec{\omega}, b}(\vec{x}^{(i)}), y^{(i)})$$

- *) Cost function is for model's all training set
- *) Loss function is only for one training example set

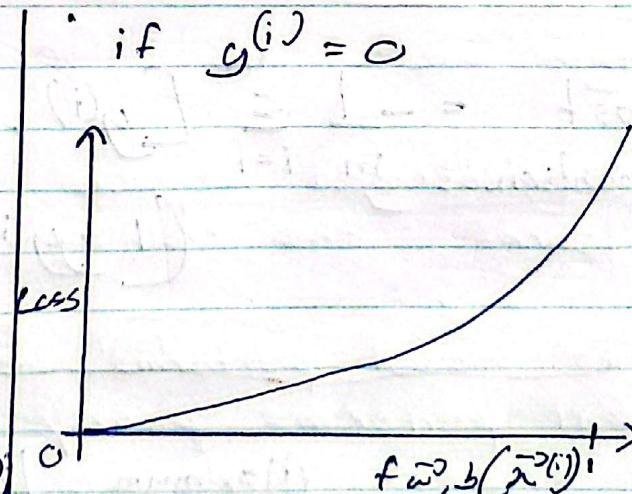
$$L(f_{\vec{\omega}, b}(\vec{x}^{(i)}), y^{(i)}) = \begin{cases} -\log(f_{\vec{\omega}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ -\log(1 - f_{\vec{\omega}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$

find $\vec{\omega}, b$ that minimize cost J

if $y^{(i)} = 1$



if $y^{(i)} = 0$



$f_{\vec{\omega}, b}(\vec{x}^{(i)}) \rightarrow 1$ then loss = 0

$f_{\vec{\omega}, b}(\vec{x}^{(i)}) \rightarrow 1$ then loss = ∞

$f_{\vec{\omega}, b}(\vec{x}^{(i)}) \rightarrow 0$ then loss = ∞

$f_{\vec{\omega}, b}(\vec{x}^{(i)}) \rightarrow c$ then loss = 0

* The further prediction $f_{\vec{w}, b}(\vec{x}^{(i)})$ is from target $y^{(i)}$, the higher the loss

Simplified Cost function for Logistic Regression:

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

Log loss (Convex)

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

Cost Function

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

Maximum likelihood

05) Gradient descent for logistic Regression

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

repeat

{

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

}

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

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

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

} simultaneous updates

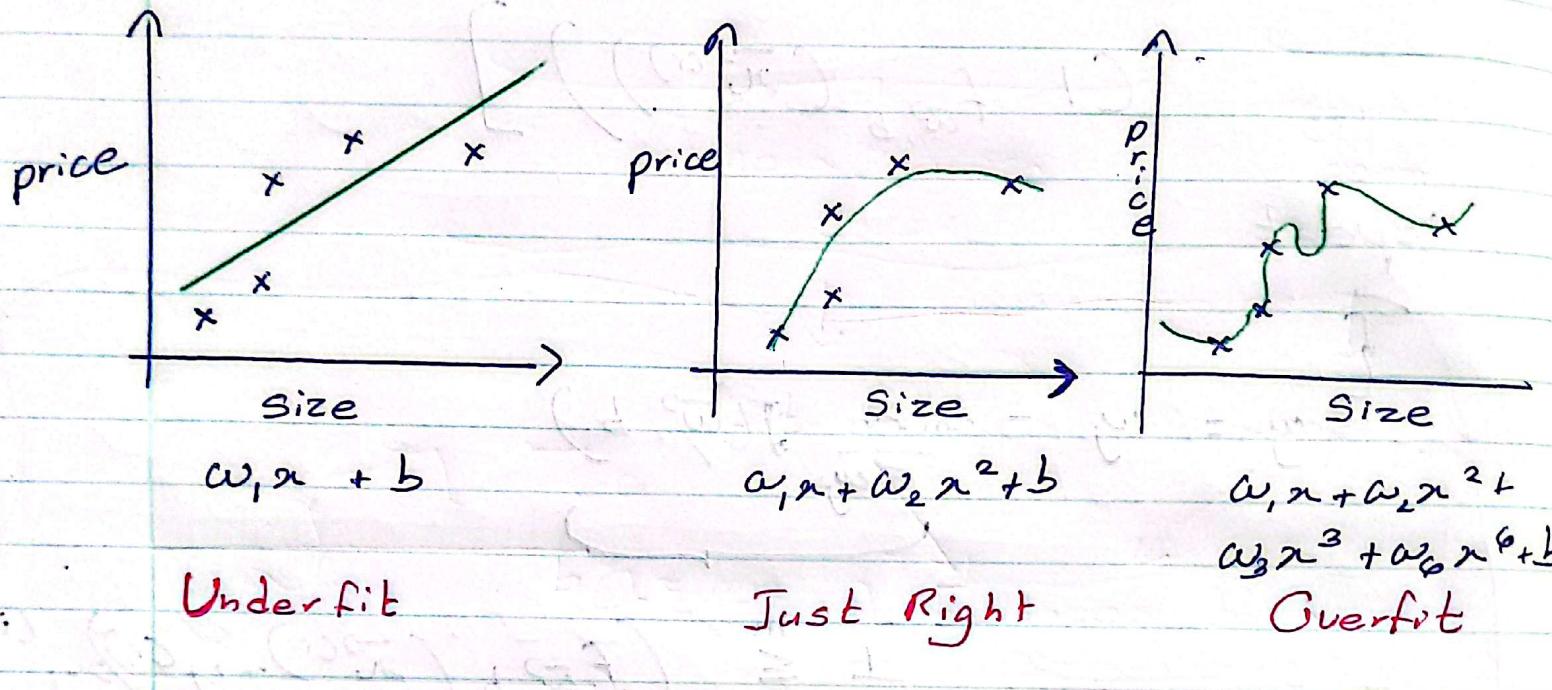
- * As gradient descent for both linear and logistic looks same it is the $f_{\vec{\omega}, b}(\vec{x})$ that changes.

linear regression = $\vec{\omega} \cdot \vec{x} + b$

logistic regression = $\frac{1}{1 + e^{-(\vec{\omega} \cdot \vec{x} + b)}}$

06] The problem overfitting

Regression Example



- *) Does not fit the training set well
- *) Fits training set pretty well
- *) Fits the training set extremely well

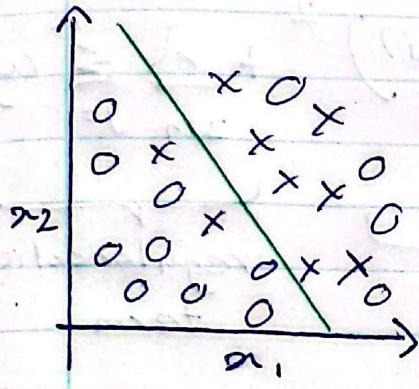
High bias

Generalization

High variance

- *) Overfitting happens when a ML model learns the training data too well.
- *) The model performs very well on training data But performs poorly on new un-seen data
- *) Too complex model, Too little training data, No regularization, Too many features can cause overfitting

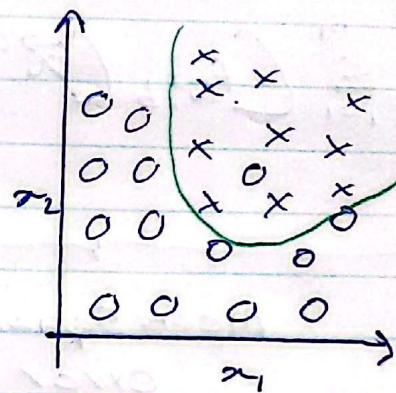
Classification Example



$$Z = \omega_0 x_0 + \omega_1 x_1 + \omega_2 x_2 + b$$

$$f_{\vec{\omega}, b}(\vec{x}) = g(Z)$$

g is the sigmoid function



$$Z = \omega_0 x_0 + \omega_1 x_1 + \omega_2 x_2$$

$$+ \omega_3 x_1^2 + \omega_4 x_2^2$$

$$+ \omega_5 x_1 x_2 + b$$

Just right

Underfit
High bias

$$Z = \omega_0 x_0 + \omega_1 x_1 + \omega_2 x_2$$

$$+ \omega_3 x_1^2 x_2 + \omega_4 x_1^2$$

$$x_2^2 + \omega_5 x_1^2 x_2^3 +$$

$$\omega_6 x_1^3 x_2 + \dots + b$$

Overfit

How to prevent Overfitting

1) Use more training data - More examples help the model generalize

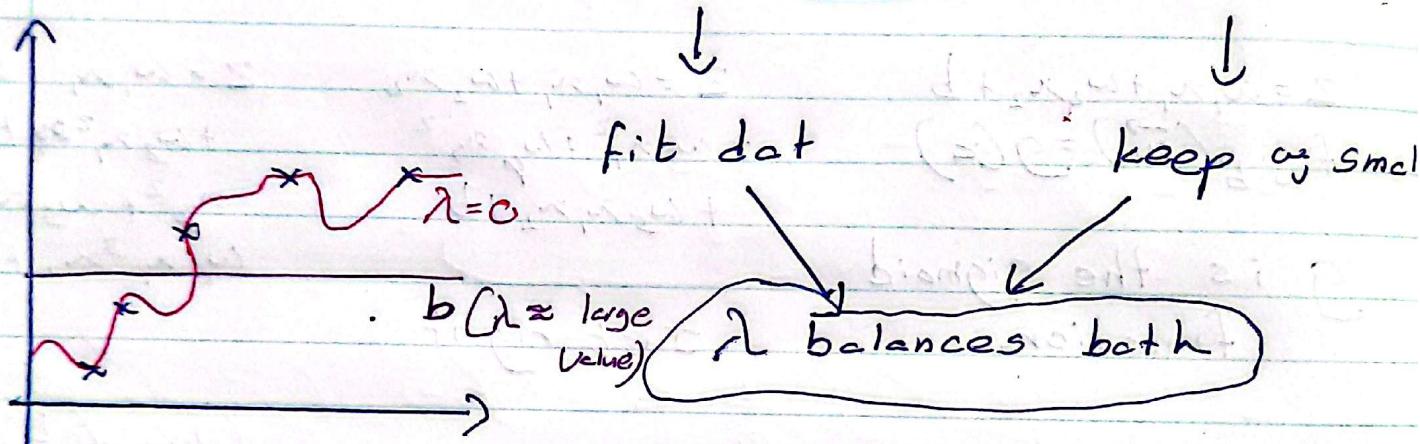
2) Regularization - penalizes large weights to prevent extreme fitting

3) Feature selection - Remove irrelevant features that cause noise

07) Cost function with regularization

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

Mean Squared error regularization term



- *) If $\lambda = 0$ overfit
- If λ enormous (like 10^{10}) underfit

Need to choose a λ in between

08) Regularized linear Regression

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

Gradient descent

repeat {

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

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

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

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

Don't have to
regularize b

3 Simultaneous update

Implementing Gradient descent

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

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

For logistic regression it looks same
but $f(\vec{x})$ changes ($\frac{1}{1 + e^{-x}}$)

$$\omega_j = \underbrace{\omega_j \left(1 - \alpha \frac{\lambda}{m} \right)}_{\text{new}} - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^m (f_{\vec{\omega}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j}_{\text{original part}}$$