

* Logistic Regression *

Binary classification.

We will learn binary class classification problem / multiclass we will learn later.

** calling logistic regression + 'regression' isn't correct
It actually a classification problem not regression one.

Linear regression

Polynomial regression

Logistic classification.

* Hypothesis Representation *

Denote

$P(y=0|x;\theta) \Rightarrow$ Probability that $y=0$ for given input values of x which is parametrised by θ .

Let's choose σ function for our σ {There are some} σ
A function which is such that
for x at all range of input values yields
result $\in (0,1)$ And it is calculus friendly.

The underlying assumption behind the sigmoid function is \Rightarrow

$$\ln\left(\frac{P}{1-P}\right) = \sum_{i=0}^K b_i x_i = \theta^T x$$

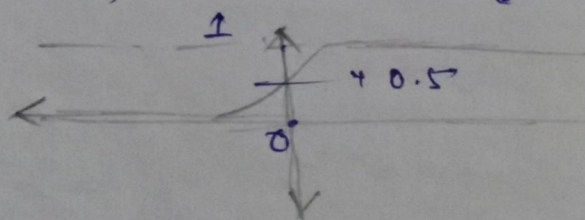
← odds ratio

$$h(x_i) = \frac{P}{1 + \exp(-\theta^T x)}$$

} The answer of sigmoid function is expressed in the form of probability.

* Decision Boundary *

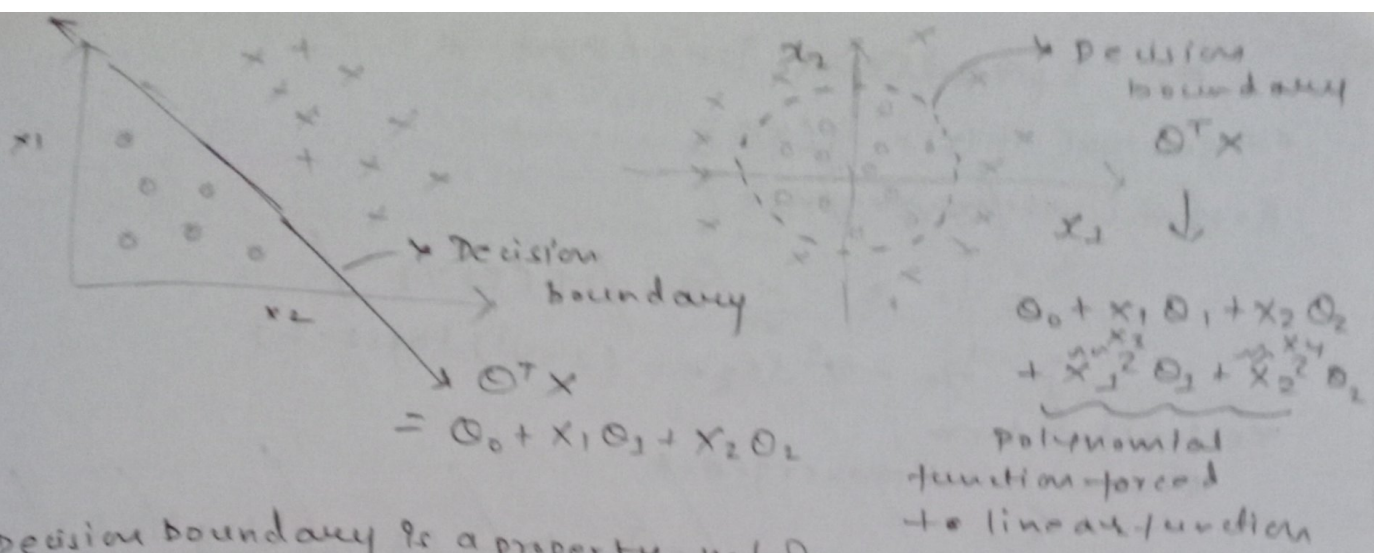
Every sigmoid function is such σ



\therefore when $h(x_i) = 0.5$
 \downarrow
 $x = 0.$

Predict $\Rightarrow y = \begin{cases} 0 & \theta^T x \leq 0 \\ 1 & \theta^T x \geq 0. \end{cases}$

This theory only applies when you find θ parameters.



"Decision boundary is a property, not of training set, but of the hypothesis under the parameters"

* Finding The Parameter Value (θ) for the most appropriate fit.

* Gradient Descent approach *

Minimise the cost of error.

$$J(\theta) = \sum_{i=1}^m (g(\dots)) \times \frac{1}{2m} = \frac{1}{m} \sum_{i=1}^m \text{cost}(h, y)$$

what to use ***

cost(h, y)

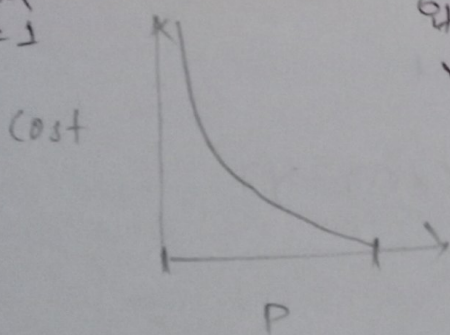
probability prediction

$$\begin{cases} -\log(h(x)) & y=1 \\ -\log(1-h(x)) & y=0 \end{cases}$$

It is the product of some superior statistical math

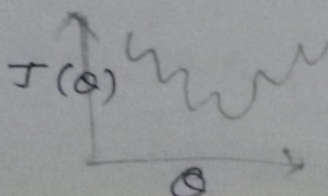
But let's concern ourselves only with intuition: \rightarrow

When $y=1$



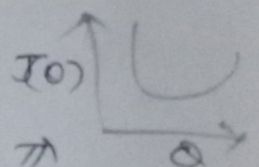
If you say no cancer when there actually a cancer then you have chosen death

Not $(h-y)^2$ because it will yield



But when cost is

$$\text{cost} = y * (-\log p) - (1-y) * \log(1-p)$$



* Gradient Descent *

** note that $h(x) = g(x\theta)$

Repeat until convergence

$$\theta_j = \theta_j - \frac{\alpha}{m} \sum (y_i - h_i) x_{ij}$$

vectorised form

Repeat until convergence

$$\theta = \theta - \frac{\alpha}{m} \sum (y - h) x^T$$

The vectorised cost function:

$$J(\theta) = \frac{1}{m} \sum \text{cost}(h, y)$$

$$J(\theta) = \frac{1}{m} \{ -y^T \log h - (1-y)^T \log(1-h) \}$$

The gradient descent:

Repeat until convergence

$$\theta_j = \theta_j - \frac{\alpha}{m} \frac{\partial}{\partial \theta_j} J(\theta)$$

Foundation

$$\theta_j = \theta_j - \frac{\alpha}{m} \sum (h - y_i) x_{ij}$$

% But remember that

$$h(x\theta) = y$$

also if log R
out to be this

}

Vectorised form of gradient descent

{

$$\theta = \theta - \frac{\alpha}{m} x^T (g(x\theta) - y)$$

}

feature Scaling

and

Normalisation

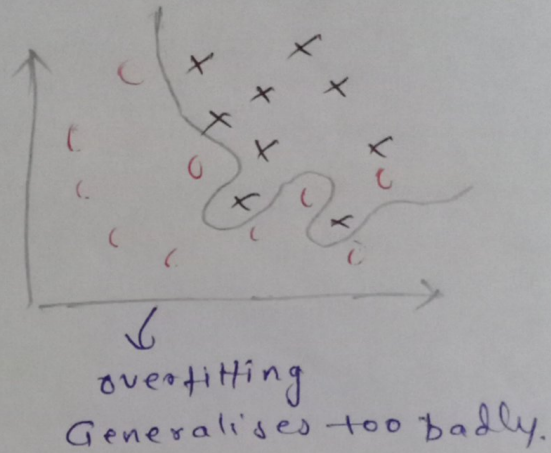
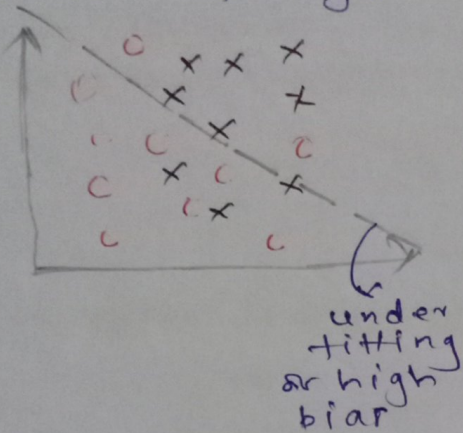


Don't forget me

* Overfitting *

"It makes accurate prediction for the example in training set, but it may not generalise the hypothesis which give good result while predicting new or unknown queries"

What is overfitting?



How to address overfitting?

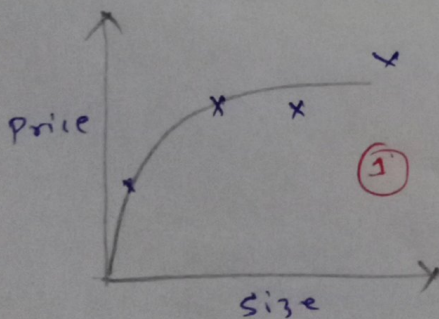
1. Reduce number of features { Also inc $x_3 = x_1 x_2^2$ }
 → Select which features to keep yourself
 → Model selection Algorithm
2. Regularisation.

* When to Use Regularisation? *

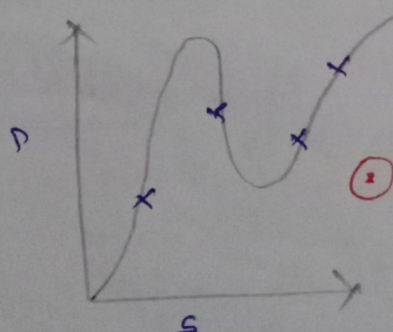
→ There are many features { including the made up ones }

Effect of regularisation

1 Perfect fit



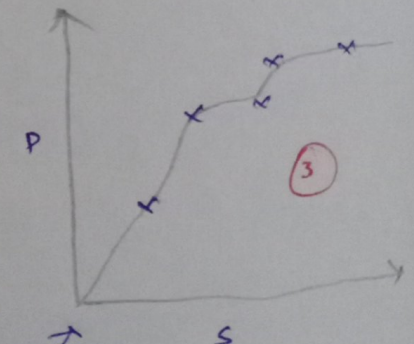
$$p = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2$$



Overfitting → Regularisation

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

Technically an extra term



How to reduce the effect of unnecessary term?
 It is to reduce the value of parameters.

So we make cost function pay greatly for each parameter

$$\therefore J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h - y)^2 + \sum_{i=1}^n (\lambda \theta_i^2) \quad \text{regularisation parameter}$$

* We want to reduce some parameter terms (θ) so that their amplification

Here θ_0 is not included, because it's not amplifying any parameter

undecisive ~~for~~ inputs/features decreases.

But / But / But their unwanted presence

* Improves the quality of ~~data~~ regression / Decision-boundary { see (3) >>> (1) }

* Gradient Descent equation *

{ % But before that note $J(\theta)$ is different for θ_0 and rest.

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

\Downarrow

$$\theta_0 := \theta_0 - \frac{\alpha}{m} \sum (h - y) x_0$$

$$\theta_j := \theta_j \left(1 - \alpha \frac{\lambda}{m} \right) - \frac{\alpha}{m} \sum_{i=1}^m (h - y) x_j^{(i)}$$

intuition.

$\lambda \downarrow$ so gradually the parameter is shrinking.

Regularisation:

Trade off between small error and simpler fit.

By adding more and more parameters I can actually reduce the squared error but at the same time is it actually worth it when compared to the penalty I am paying for it

Finally:

Improved Normal Equation \rightarrow

$$\theta = (X^T X + \lambda I)^{-1} X^T y$$

* Regularised Logistic Regression *

Again,

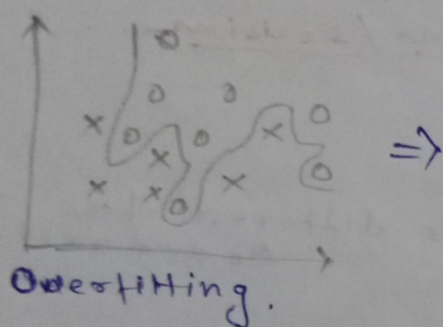
$$J(\Theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(h, y) + \frac{\lambda}{2} \sum_{i=1}^{n+1} \Theta_i^2 \quad \lambda \text{ Regularization parameter.}$$

Gradient Descent,

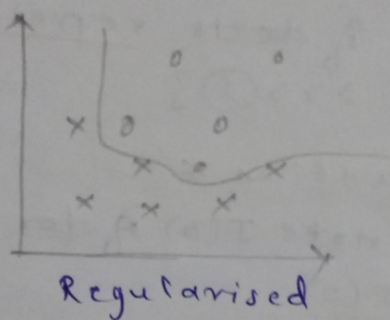
$$\{ \Theta_0 := \Theta_0 \left(1 - \alpha \frac{\lambda}{m} \right) - \alpha \frac{1}{m} \sum (h - y)$$

$$\{ \Theta_j := \Theta_j \left(1 - \alpha \frac{\lambda}{m} \right) - \alpha \frac{1}{m} \sum (h - y) x_j$$

Same as before for linear regression



\Rightarrow



Found A NEW NAME for parameter Θ , "WEIGHTS"

$$X^T X^{-1} (X^T X + \lambda I)$$