

-- @

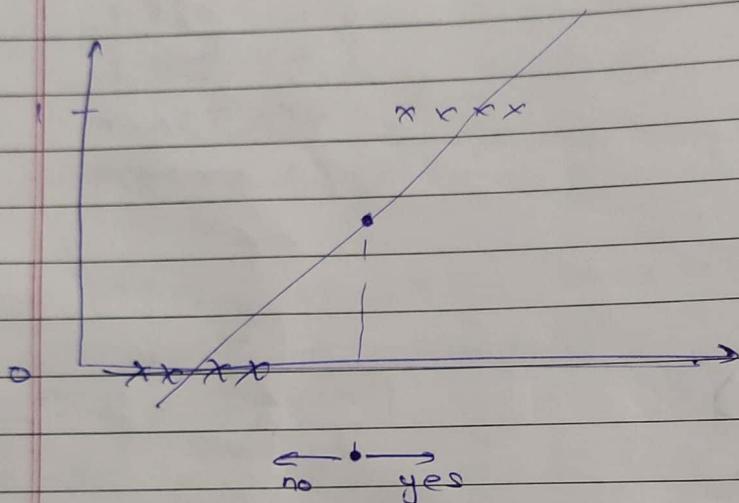
classification — — —

*

Classification.

$y \in \{0, 1\}$ binary classification

$y \in \{0, 1, 2, 3\}$ lattice



Threshold classifier output $h_0(x)$ at 0.5:

→ If $h_0(x) \geq 0.5$, predict $y=1$

→ else $y=0$.

In classification

$$y = 0 \quad y = 1$$

$h_0(x)$ can be >1 / <0

In logistic regression

$$0 \leq h_0(x) \leq 1$$

① Hypothesis Representation — — —

* Logistic Regression Model.

$$0 \leq h_{\theta}(x) \leq 1$$

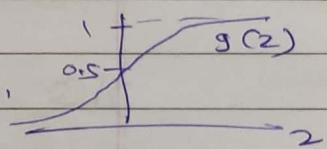
$$h_{\theta}(x) = \theta^T x \text{ previously}$$

$$h_{\theta}(x) = g(\theta^T x) \text{ now}$$

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

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

sigmoid function
logistic function.



• Interpretation of Hypothesis output.

$h_{\theta}(x)$ = estimated probability that $y=1$ on input x .

$$\text{e.g. if } x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorsize} \end{bmatrix}$$

$$\text{i)} g(z) \geq 0.5$$

$$h_{\theta}(x) = \theta^T x \quad y=1$$

' when $z \geq 0$

70% chance of being malignant from graph

$$h_{\theta}(x) = P(y=1 | x; \theta)$$

given parameters trained by

$$\text{i)} h_{\theta}(x) = g(\theta^T x) \geq 0.5$$

when $\theta^T x \geq 0$

$\nwarrow z$

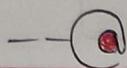
$$P(y=0 | x; \theta) + P(y=1 | x; \theta) = 1$$

$$\text{ii)} h_{\theta}(x) = g(\theta^T x) < 0.5$$

$g(z) < 0.5$

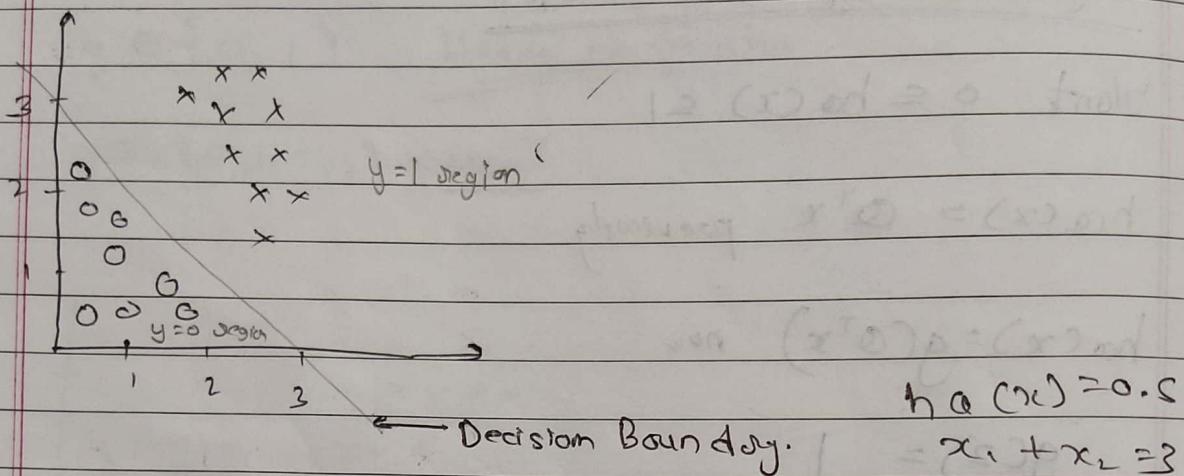
$\theta^T x \leq 0$

- i) predict $y=1$ if $h_{\theta}(x) \geq 0.5 \rightarrow \theta^T x \geq 0$
- ii) predict $y=0$ if $h_{\theta}(x) < 0.5 \rightarrow \theta^T x \leq 0$



Decision Boundary --

• Decision Boundary.



$$h_{\theta}(x) = g(\alpha_0 + \alpha_1 x_1 + \alpha_2 x_2)$$

If θ found out to be $\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$

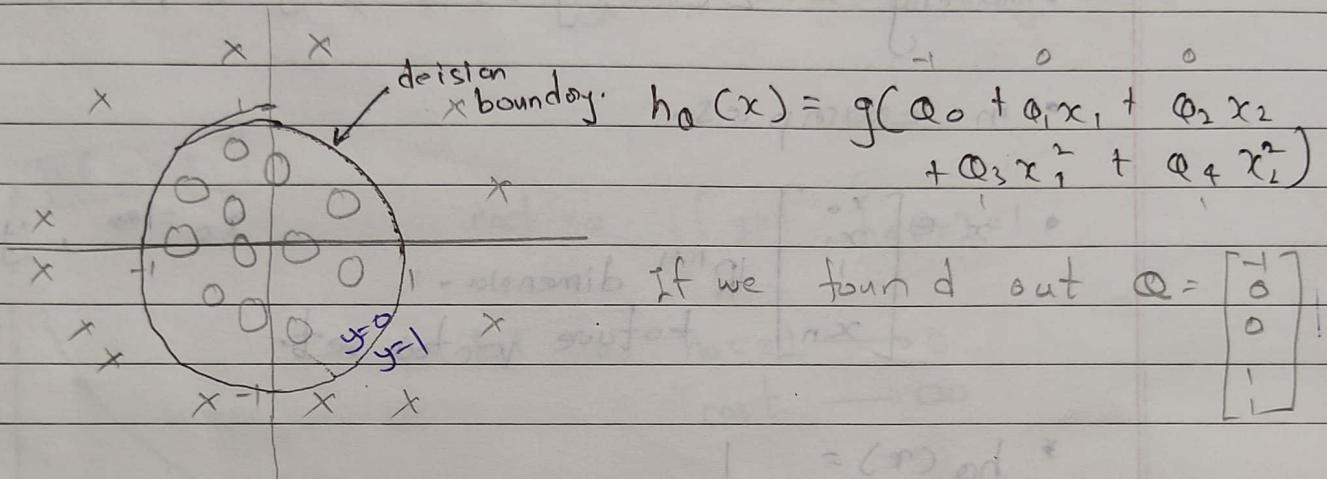
- predict. $y=1$ if $\theta^T x \geq 0$
 $-3 + x_1 + x_2 \geq 0$

- $y=0$ if $\theta^T x + x_2 < 0$

Decision Boundary is property of Hypothesis $h_{\theta}(x) = g(\theta^T x)$
and of parameters of hypothesis
and not property of Data set.

- Non-linear decision boundary

In linear region we added higher deg see term. Same here.



- Predict $y=1$ if $\alpha^T x \geq 0$

$$-1 + x_1^2 + x_2^2 \geq 0$$

$$x_1^2 + x_2^2 \geq 1$$

- We can add more terms like $x_1^2x_2, x_1^3x_2$ to get more complex decision boundary.

$$(C_0 - C_1) \alpha_0 + (C_0 + C_1) \alpha_1 = (C_0 - C_1) \alpha_0 + (C_0 + C_1) \alpha_1$$

comes out to be $C_0 - C_1$

so the final equation is $\alpha_0 - \alpha_1$

$\alpha_0 - \alpha_1$

---@ Cost function - - -

- Cost function: logistic regression.

Given:-

* Training set $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

- m examples

* $x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$ \mathbb{R}^{n+1} dimension
feature vector set.

$$x_0 = 1, y \in \{0, 1\}$$

$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

- How to choose parameter θ .

* cost fun before: $J(\theta) = \underbrace{\frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_\theta(x^{(i)}) - y^{(i)})^2}_{\text{cost } h_\theta(x^{(i)}, y)}$

$$\text{cost } h_\theta(x^{(i)}, y^{(i)}) = \frac{1}{2} \underbrace{(h_\theta(x^{(i)}) - y^{(i)})^2}_{\frac{1}{1 + e^{-\theta^T x}}}$$

making this is FINE,

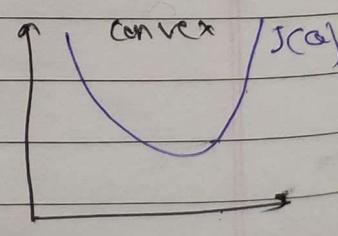
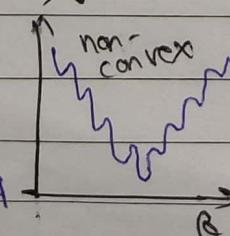
but this gives non-convex

θ , & that does not

guarantee global minima.

It could be case

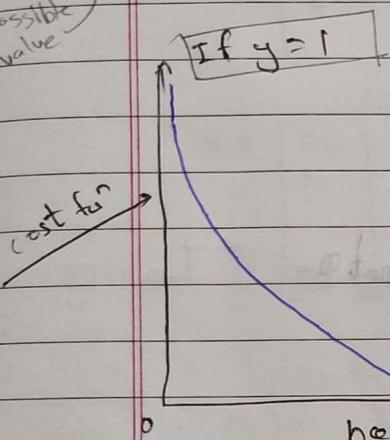
we got local minima



- Cost function: logistic regression.

$$\text{cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0. \end{cases}$$

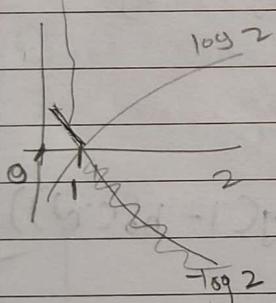
$h_{\theta}(x) \leq 1$
so it is 0
possible value



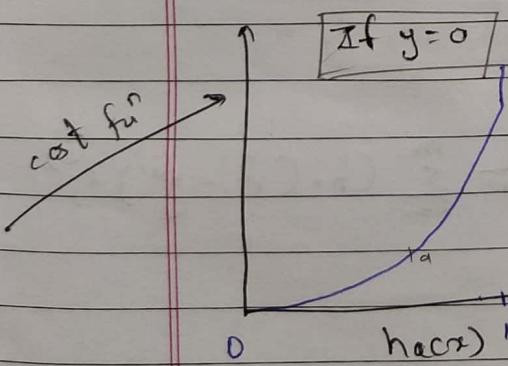
cost = 0, if $y=1$, $h_{\theta}(x)=1$

But as $h_{\theta}(x) \rightarrow 0$
cost $\rightarrow \infty$

Captures intuition that if $h_{\theta}(x)=0$ (predict $P(y=1|x; \theta)=0$), but $y=1$ we'll penalize learning algorithm by a very large cost.



We are wrong by 1 in probability because we told it is benign & found to be malignant, \therefore cost $\rightarrow \infty$.



Same here.

✓ If $h_{\theta}(x)=y$, then $\text{cost}(h_{\theta}(x), y)=0$ (for $y=0$ & $y=1$)

✓ If $y=0$, then $\text{cost}(h_{\theta}(x), y) \rightarrow \infty$ as $h_{\theta}(x) \rightarrow 1$

✓ If $y=0$, then $\text{cost}(h_{\theta}(x), y) \rightarrow \infty$ as $h_{\theta}(x) \rightarrow 0$

Regardless of whether $y=0$ or $y=1$, if $h_{\theta}(x)=0.5$, then $\text{cost}(h_{\theta}(x), y) \geq 0$ a case

$$\text{cost}(h_{\theta}(x), y) = -y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x))$$

$$\text{cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & \text{if } y=1 \\ -\log(1-h_\theta(x)) & \text{if } y=0. \end{cases}$$

Page No: 25
Date: / /

y can only be 0/1 simplified cost function and Gradient Descent

* logistic regression cost function:-

$$\text{cost}(h_\theta(x), y) = -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x))$$

→ Principal of maximum likelihood estimation. statistics.

That why this is cost fn, other fn can also work as cost fn by combining two fn.

• convex

- To fit parameters θ :

$$\min_{\theta} J(\theta) \rightarrow \text{Get } \theta. \text{ Get } \theta. \text{ Training}$$

- To make a prediction

$$f(x) = \text{output} \quad h_\theta(x) = \frac{1}{1+e^{-\theta^T x}} \quad P(y=1|x; \theta) \quad \text{Testing.}$$

- Gradient Descent.

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)})) \right]$$

Want to $\min_{\theta} J(\theta)$:

Repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j} \rightarrow \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

simultaneously update all θ_j

~~Same eqn of linear reg.~~

Repeat {

$$\alpha_j := \alpha_j - \alpha \sum_{i=1}^m (h_\alpha(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

simultaneously update all α_j)

}

$$\alpha = \begin{bmatrix} \alpha_0 \\ \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix}$$

logistic

$$h_\alpha(x) = \frac{1}{1 + e^{-\alpha^T x}}$$

linear

$$h_\alpha(x) = \alpha^T x$$

- * Check same way as in linear reg. for Is my gradient descent converging.

$$\uparrow J(\alpha) \rightarrow \frac{1}{m} \sum_{i=1}^m (h_\alpha(x^{(i)}) - y^{(i)})^2$$

↑
no of iterations

④ Advanced optimization --

* logistic regn : Advanced optimization

cost function $J(\alpha)$, want mince $J(\alpha)$.

- Given α , we have code that can compute.

$J(\alpha)$

$$\frac{\partial}{\partial \alpha_j} J(\alpha) \quad \text{-- for } j=0, 1, 2, \dots$$

Gradient descent

1 repeat

2

optimization algorithm

- Optimization algos:-

Gradient descent

conjugate gradient

BFGS

L-BFGS

Advantages:

- No need to manually pick L .
- often faster than gradient descent

Disadvantages:

• More complex.

• We don't completely know how it is working, hard enough

$$\text{eg. } \mathbf{Q} = \begin{bmatrix} Q_1 \\ Q_2 \end{bmatrix}$$

$$J(\mathbf{Q}) = (Q_1 - s)^2 + (Q_2 - s)^2$$

$$\frac{\partial}{\partial Q_1} J(\mathbf{Q}) = 2(Q_1 - s)$$

$$\frac{\partial}{\partial Q_2} J(\mathbf{Q}) = 2(Q_2 - s)$$

→ for $\min_{\mathbf{Q}} J(\mathbf{Q})$

$$Q_1 = s, Q_2 = s$$

Week 3, LRM → video → Hdu opt. for coding

Multiclass classification: one vs all --

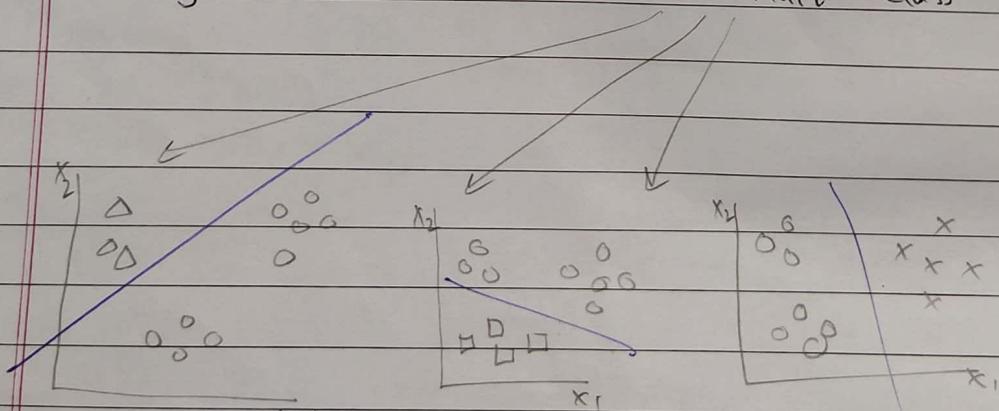
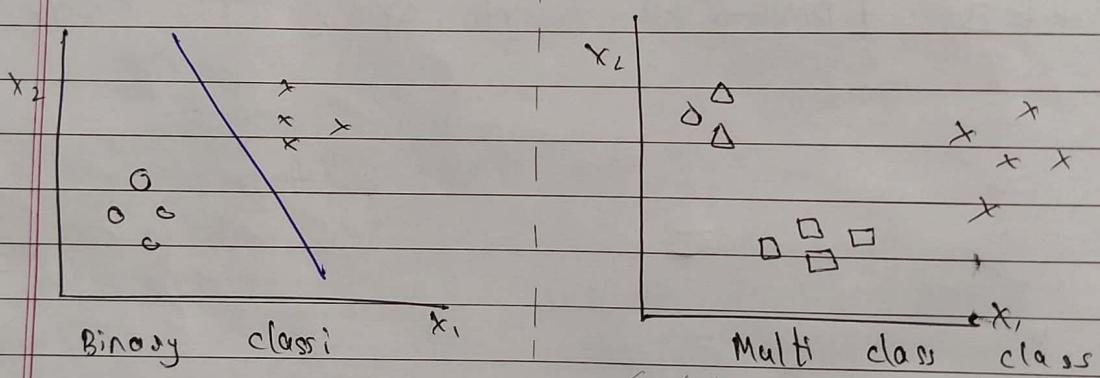
* logistic regression: Multi class classification:
one-vs-all

Multiclass classification:

Email folder tagging: Work, friend, family, Hobby
 y_{21} y_{22} y_{23} y_{24}

Disease cold, flu, Not ill
 y_{21} y_{22} y_{23}

Weather sunny, cloudy, rain, snow
 not matter 1, 2, 3, 4



$$\text{One-vs-rest } h_\alpha^{(1)}(x) \quad P(y=1|x; \alpha)$$

$$h_\alpha^{(2)}(x) \quad P(y=2|x; \alpha)$$

$$h_\alpha^{(3)}(x) \quad P(y=3|x; \alpha)$$

$$h_\alpha^{(i)}(x) = P(y=i|x; \alpha) \quad (i=1, 2, 3)$$

Ques K classes ($y \in \{1, 2, \dots, K\}$) using 1 vs all. How many different logistic regressions classifiers do you will you end up training? $\rightarrow K$.

$(i=1, 2, 3)$ we trained 3.

- One-vs-all.

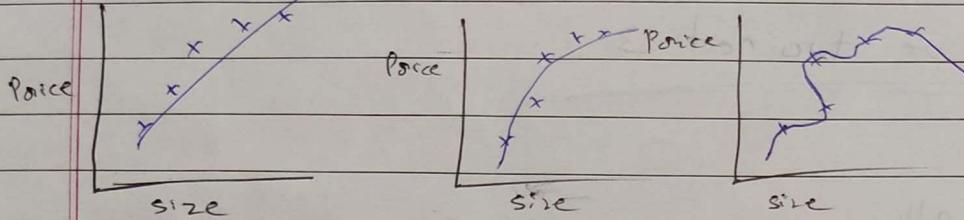
Train a logistic regression classifier $h_a^{(i)}(x)$ for each class i to predict the probability that $y=i$

On a new input x , to make a prediction, pick the class i that maximizes

$$\max_i h_a^{(i)}(x)$$

* Regularization : the problem of overfitting

e.g. linear regression (housing price)



Qotance

'underfit'

'High bias'

Qotance + Qout

'just right'

'overfit'

'High variance'

$Q_0 + Q_1 x + Q_2 x^2 + Q_3 x^3 + Q_4 x^4$

overfitting: If we have too many features, the learned hypothesis may fit the training set very well ($J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$) but fail to generalize to new examples (not able to predict price of new examples)

- Addressing overfitting:

options:

» Reduce number of features

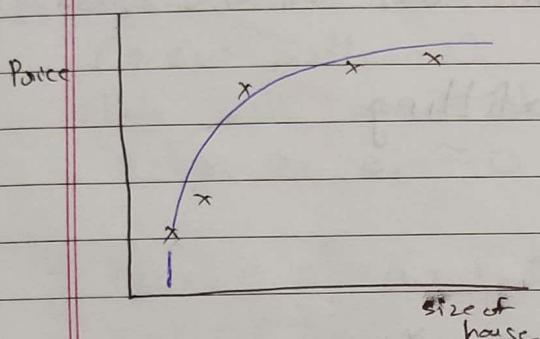
- manually select which feature to keep.
- Model selection algorithm (latter in course)

» Regularization

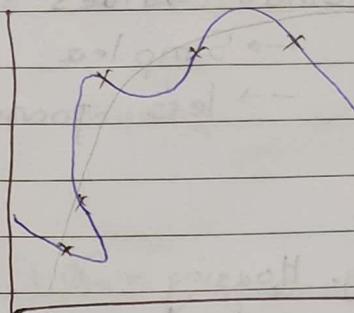
- keep all the features, but reduce magnitude values of parameters θ_j
- Works well when we have a lot of features, each of which contributes a bit to predicting y .

---@ cost function ---

- Intuition of regularization :-



$$\theta_0 + \theta_1 x + \theta_2 x^2$$



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Suppose we make θ_3, θ_4 really small

→ This can make $\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 \rightarrow \theta_0 + \theta_1 x + \theta_2 x^2$

Isn't it like removing x^3, x^4 terms or removing some features. Automatically / manually either maybe.

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \underbrace{1000\theta_3^2 + 1000\theta_4^2}_{\text{new terms added}}$$

$\theta_3 \approx 0 \quad \theta_4 \approx 0$

→ minimizing this cost function will surely make θ_3, θ_4

→ 1000, 1000 are just randomly big numbers, as we want to reflect cost caused by θ_3, θ_4 , so we choose big numbers.

• Regularization.

Small values for parameters $\theta_0, \theta_1, \dots, \theta_n$ all θ_i

- Simple hypothesis (smooth curve)
- less prone to overfitting

e.g. Housing.

- Features $x_1, x_2, x_3, \dots, x_{100}$

- Parameters $\theta_0, \theta_1, \dots, \theta_{100}$

• Here we don't know which parameter is to shrink
 $\theta_0 \rightarrow 0 / \theta_1 \rightarrow 0 / \theta_2 \rightarrow 0$ we don't know which $\theta_i \rightarrow 0$

• So we shrink all θ_i s. by modifying cost function of linear regn.

Regularization Parameter

Regularized linear regn.

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (\theta_0(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2$$

→ This will shrink all $\theta_i \rightarrow \theta_0, \theta_1, \theta_2, \theta_3, \dots, \theta_{100}$ term

• In practice Including θ_0 /not including does not make any difference in results.

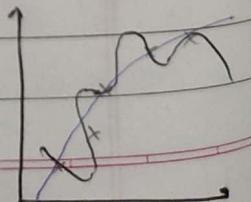
→ This is what we usually do

λ :- Regularization parameter. \rightarrow 2 goals

fit training set

shrink parameters

λ controls trade off b/w two goals.



Ans

$$\text{Regulation. } J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (\hat{y}_i - y_i)^2 + \lambda \sum_{j=1}^n \theta_j^2 \right] \quad \text{If } \lambda = 10^{\circ}$$

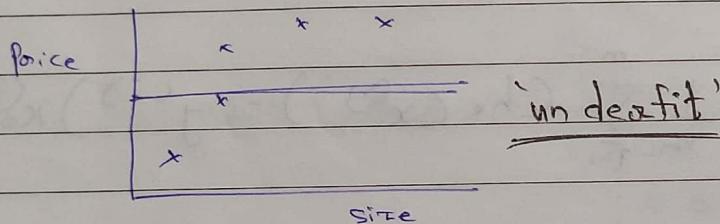
Ans Algo results in underfitting \rightarrow fails to fit even training set.

This will result shrinking too much,

$$\theta_1, \theta_2, \theta_3, \theta_4 \approx 0 \quad \underline{\text{all}}$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x + \theta_3 x + \theta_4 x$$

$$h_{\theta}(x) = \theta_0$$



Lattice \rightarrow choosing λ automatically in multiselection

--- (c) Regularized linear regression ---

⇒ Regularized linear regreⁿs.

- Gradient descent before (for linear regreⁿs)

Repeat {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (\text{h}_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad \text{--- } \theta_j = \theta_0, \theta_1, \dots, \theta_n$$

}

Just waiting do separately // no change.
updated for regularization.

Repeat {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (\text{h}_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (\text{h}_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \quad (j=1, 2, \dots, n)$$

}

We penalize $\theta_1, \theta_2, \theta_3, \dots, \theta_n$ only in regularization not θ_0 .

$$\boxed{\theta_j := \theta_j \left(1 - \alpha \frac{\lambda}{m}\right) - \alpha \frac{1}{m} \sum_{i=1}^n (\text{h}_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}}$$

$$1 - \alpha \frac{\lambda}{m} < 1$$

Its like $\theta_j \times 0.99$

- Normal eqn updated for regularization.

$$x = \begin{bmatrix} (x^{(0)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix} \quad m \times (n+1)$$

$$y = \begin{bmatrix} y^{(0)} \\ \vdots \\ y^{(m)} \end{bmatrix} \quad \mathbb{R}^m$$

$$\min_{\theta} J(\theta) \quad \frac{\partial J(\theta)}{\partial \theta_j} = 0$$

$$\theta = (X^T X + \lambda \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix})^{-1} X^T y$$

$(n+1) \times (n+1)$

e.g. $n=2$ $\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$

Non-invertibility (optional/advanced)

Suppose $m \leq n$

(#examples) \leq (#features)

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

\curvearrowright non-invertible singular.

If $\lambda > 0$

This is invertible matrix anyway.

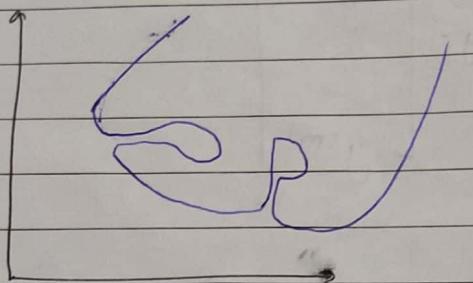
$$\theta = (X^T X + \lambda \begin{bmatrix} 1 & & \\ & 1 & \\ & & 1 \end{bmatrix})^{-1} X^T y$$

--- @ Regularized logistic regression ---

Page No: 37
Date: / /

★ Regularized logistic regⁿ.

• Regularized logistic regⁿ



$$h_{\phi}(x) = g(\phi_0 + \phi_1 x_1 + \phi_2 x_1^2 + \phi_3 x_1^3 + \phi_4 x_1^2 x_2^2 + \phi_5 x_1^2 x_2^3 + \dots)$$

• cost fun :-

$$J(\phi) = - \left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_{\phi}(x^{(i)}) + (1-y^{(i)}) \log (1-h_{\phi}(x^{(i)})) + \frac{\lambda}{2m} \sum_{j=1}^n \phi_j^2 \right]$$

$\phi_1, \phi_2, \dots, \phi_n$

• Gradient descent.

Repeat {

$$\phi_0 := \phi_0 - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_{\phi}(x^{(i)}) - y^{(i)}) x_0^{(i)} \right]$$

$$\phi_j := \phi_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_{\phi}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \phi_j \right] \quad (j=1, 2, \dots, n)$$

$$\left[\frac{\partial J(\phi)}{\partial \phi_j} \right] \text{ where } h_{\phi}(x) = \frac{1}{1 + e^{-\phi^T x}}$$

$$\bullet \text{ Plot } - \left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_{\phi}(x^{(i)}) + (1-y^{(i)}) \log (1-h_{\phi}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \phi_j^2$$

should decrease to work properly.