

Traditional Machine Learning models

Dataset

Set of features or attributes $\rightarrow f_1 \ f_2 \ f_3 \dots f_n$
 Set of instances $\left\{ \begin{array}{cccc} i_{11} & i_{12} & i_{13} & i_{1n} \\ i_{21} & i_{22} & i_{23} & i_{2n} \\ \vdots & & & \end{array} \right.$

Suppose you have a property at a location L , having size S , having rooms R , having amenities A . . . Based on these, determine the rent price P .

1) Regression

determine whether the tumor is malignant or benign

2) Classification

For a tumor, given the size S of the tumor, location L of the tumor, comorbidities C of patient, age A , family history F ,

3) Clustering

The features are available but not the ~~labels~~ output or labels. Used to learn useful properties of the structure of data

for Supervised ML tasks we ~~have~~ require a set of labeled data

Given X predict Y .

Given X predict $P(Y)$
 $i.e P(Y|X)$

hypothesis $Y = f(X)$

Linear Regression

$$D = X =$$

$$\begin{bmatrix} (\text{cont}) & (\text{size}) & (\text{room}) & (\text{area}) \\ x_{11} & x_{12} & x_{13} & \dots \\ x_{21} & x_{22} & x_{23} & \dots \\ \vdots & \vdots & \vdots & \vdots \\ x_{n1} & x_{n2} & x_{n3} & \dots \\ & & & x_{nk} \end{bmatrix}$$

$$Y = \begin{bmatrix} \text{Set of Labels} \\ \text{(price)} \\ Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}$$

Input $X \in \mathbb{R}^{n \times k}$
 outputs $Y \in \mathbb{R}^n$

Objective

Given an input vector $\vec{x} \in \mathbb{R}^k$ predict a scalar value $Y \in \mathbb{R}$

Assume that the output is a linear function of the input features (Hypothesis)

$$\hat{Y}_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \dots + \beta_k x_{ik}$$

Let \hat{Y}_i be the predicted value for the i^{th} instance (i.e for giving \vec{x})
 $\hat{Y} = \theta^T \vec{x}$ where θ is the parameter vector $= [\beta_0, \beta_1, \dots, \beta_k]^T$

Given $\theta \in \mathbb{R}^{k+1}$ is the weight parameters
from the data set learn $\beta_0, \beta_1, \dots, \beta_k$ such that the Error

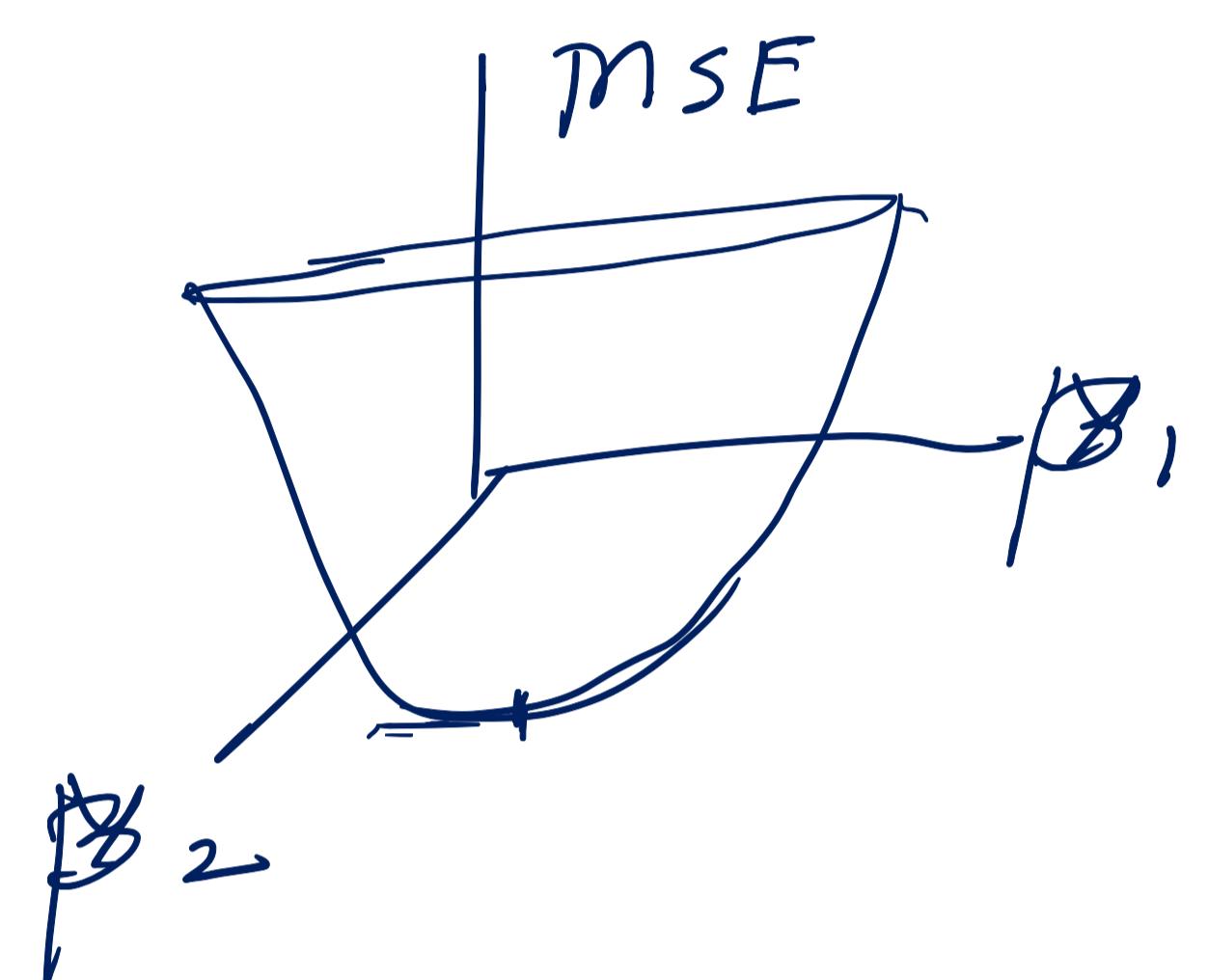
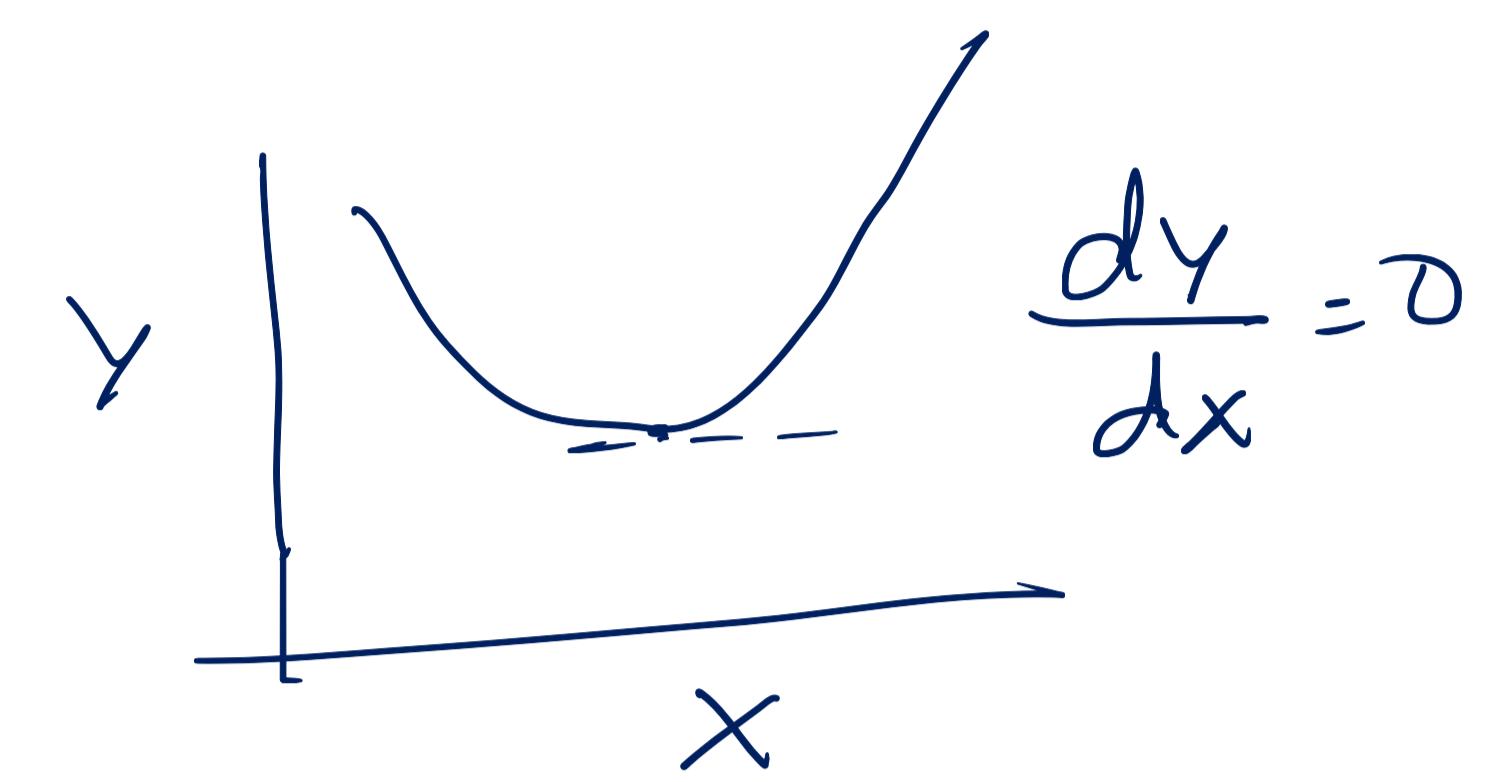
$$\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 \text{ is } \cancel{\text{minimized}} \rightarrow (\text{Mean Square Error (MSE)})$$

Divide the data into 2 parts, training set & test set
(usually 80-20 or 75-25)

$$\text{minimize} (\text{MSE} = \frac{1}{n} \|\hat{Y} - Y\|_2^2) \quad \dots \textcircled{1}$$

$$\nabla_{\theta} \text{MSE}_{\text{train}} = 0$$

$$\nabla_{\theta} \text{MSE} = \begin{bmatrix} \frac{\partial \text{MSE}}{\partial \beta_0} \\ \frac{\partial \text{MSE}}{\partial \beta_1} \\ \vdots \\ \frac{\partial \text{MSE}}{\partial \beta_k} \end{bmatrix} = 0 \quad (\text{for minima})$$



$$\hat{Y} = X\theta$$

Replace \hat{Y} by $X\theta$ in Eqn ①

$$\text{MSE} = \frac{1}{n} \|X\theta - Y\|_2^2$$

$$\nabla_{\theta} \text{MSE} = 0 \quad (\text{Derivation not done})$$

$$\theta = \frac{X^T Y}{X^T X} = \cancel{\theta} (X^T X)^{-1} X^T Y$$

for python Scikit learn provides all traditional ML package.

Ex. Approximating betweenness by Linear Regression.

Given a graph, ~~with degree of the nodes, extract~~
with the adjacency matrix, derive features of nodes
like degree of the nodes, average degree of neighbor
local clustering coeff and then use the same to predict
the betweenness.

Classification $\cdot X \in \mathbb{R}^k \quad Y \in \{0, 1\}$ or $Y \in \{C_1, C_2 \dots C_k\}$

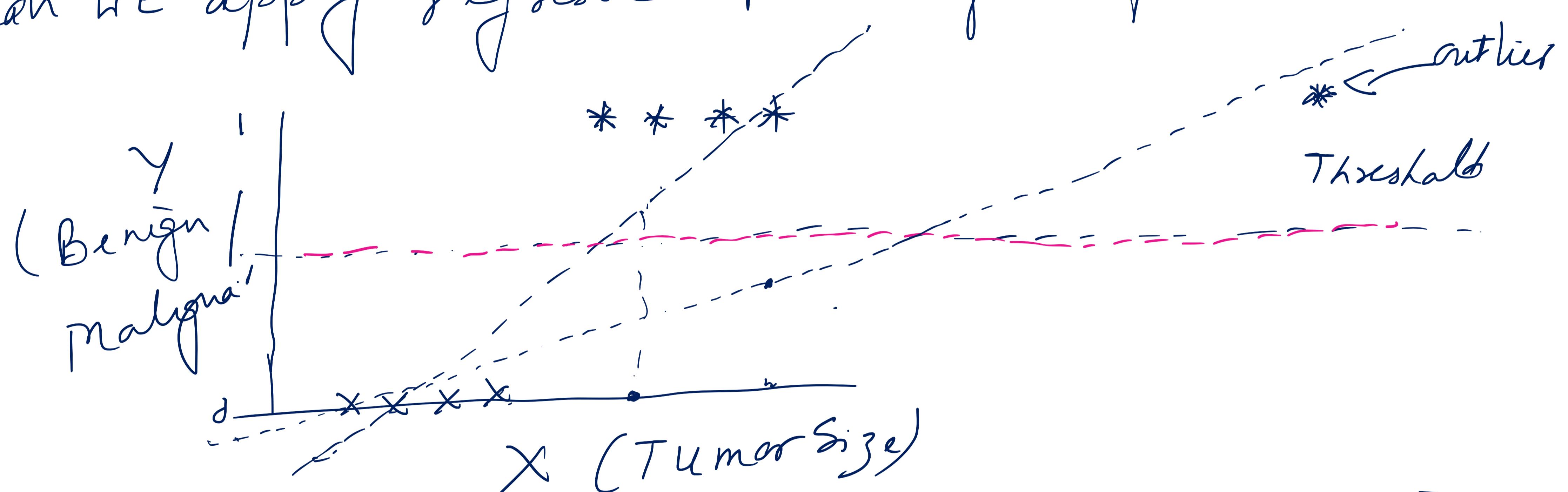
- Tumor :- Malignant / Benign
- Online Transaction :- Licit / Illicit
- Emails :- spam / non-spam
- Graphs - Link exists / not exists

When $Y \in \{0, 1\}$ it is called Binary classification

$Y \in \{C_1, C_2, C_3 \dots C_k\}$ it is called Multi-class classification.

for Regression we have $\hat{Y} = f(x)$

Can we apply regression techniques for classification



If we use the same Linear Regression principle $h_{\theta}(x) = \theta^T x$

We can use a threshold (say values of $h_{\theta}(x) > 0.5$)
then (+) we else (-) we class

Problems

- outliers can create errors and makes wrong learning
- Another problem is $h_{\theta}(x)$ may be much greater than 1 or lesser than 0.

- Fix - ~~θ~~ we want $0 \leq h_{\theta}(x) \leq 1$

\hookrightarrow Goal of Logistic Regression model

Logistic Regression Model

- for linear regression we have $h_{\theta}(x) = \theta^T x$
 $\theta^T x$ can be much larger than 1 or much lesser than 0.

For (Log Reg) we replace $\theta^T x$ by another function $g(\theta^T x)$

$g(\theta^T x)$ will have a property



So let $g(\theta^T x)$ be a sigmoid function given as $\frac{1}{1+e^{-\theta^T x}}$

Interpretation of the hypothesis's output

$h_\theta(x) = \text{Estimated prob that } Y=1 \text{ for input } x.$

$$\text{So If } x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumor size} \end{bmatrix}$$

Then $h_\theta(x) = 0.7$ indicates that there 70% chance of the tumor being malignant.

$$\text{So } h_\theta(x) = P(Y=1 | x; \theta)$$

\hookrightarrow Prob that $Y=1$ given input x , parameterized by θ .

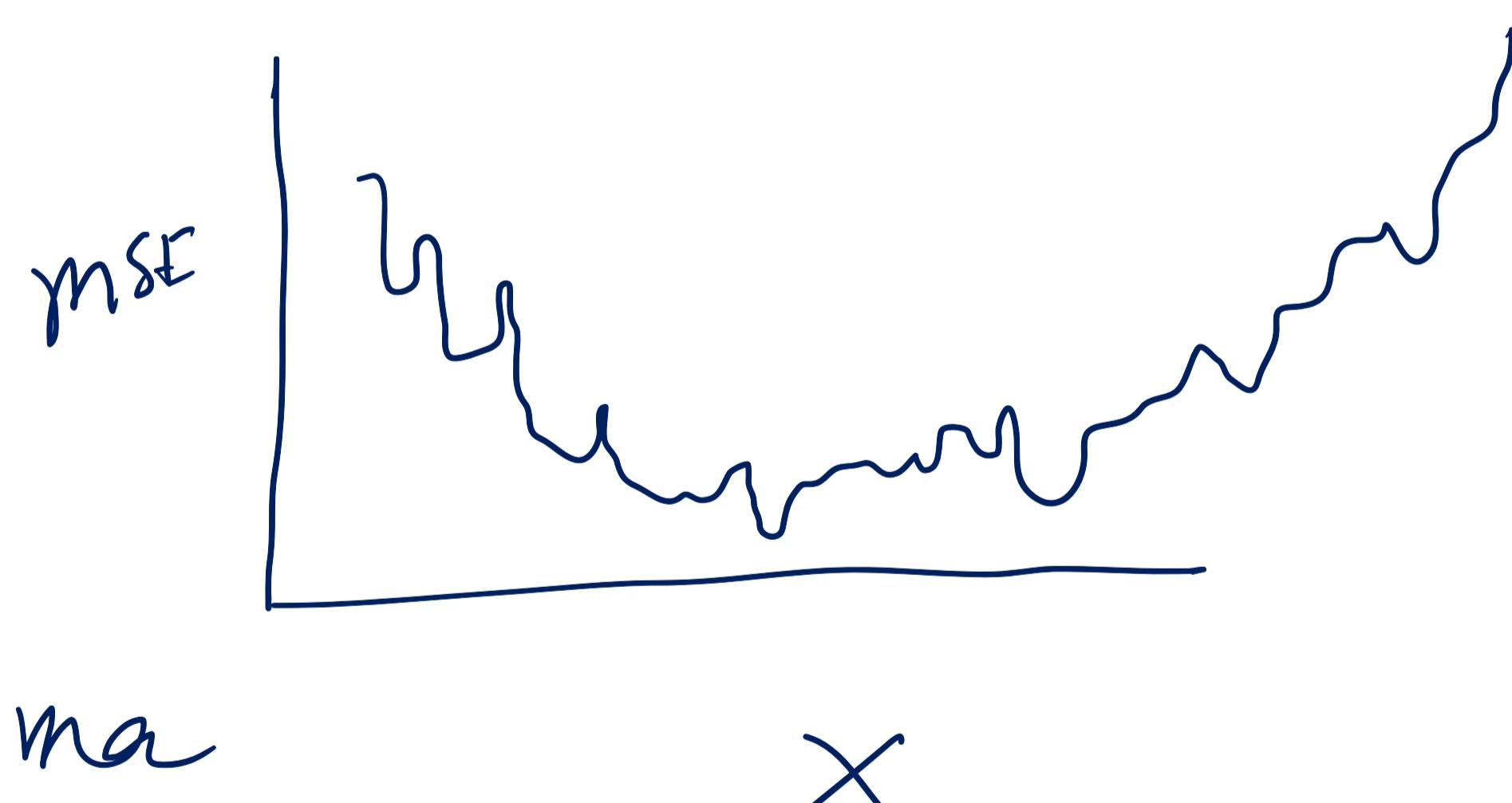
For linear regression to derive θ we minimized the

Mean Square Error (MSE)
ie $\sum_{i=1}^m (\hat{y}_i - y_i)^2$ \rightarrow when $\hat{y}_i = \theta^T x$
 MSE

for training set $\{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$

however when $\hat{y}_i = h_\theta(x_i) = \frac{1}{1+e^{-\theta^T x_i}}$ is exponential

the MSE



is not convex
ie will have lot of
local minima

So instead the following loss function is used

use Cross Entropy loss

$$\text{Cross Loss}(h_\theta(x), y) = -\log(h_\theta(x)) \text{ when } Y=1 \\ = -\log(1-h_\theta(x)) \text{ when } Y=0$$

Cross Entropy loss is given as

$$L = -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x))$$

$$\text{for } (x_1, y_1), (x_2, y_2), \dots, L_{\text{tot}} = \frac{1}{m} \sum_{i=1}^m L_i$$

Decision boundary of Logistic Regression

when $\hat{y} = 1$, when $h_\theta(x) \geq 0.5$ or $\theta^T x > 0$
 when $\hat{y} = 0$ when $h_\theta(x) \leq 0.5$ or $\theta^T x \leq 0$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_k \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_k \end{bmatrix}$$

$$\theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_k x_k > 0$$

This is a line, so the

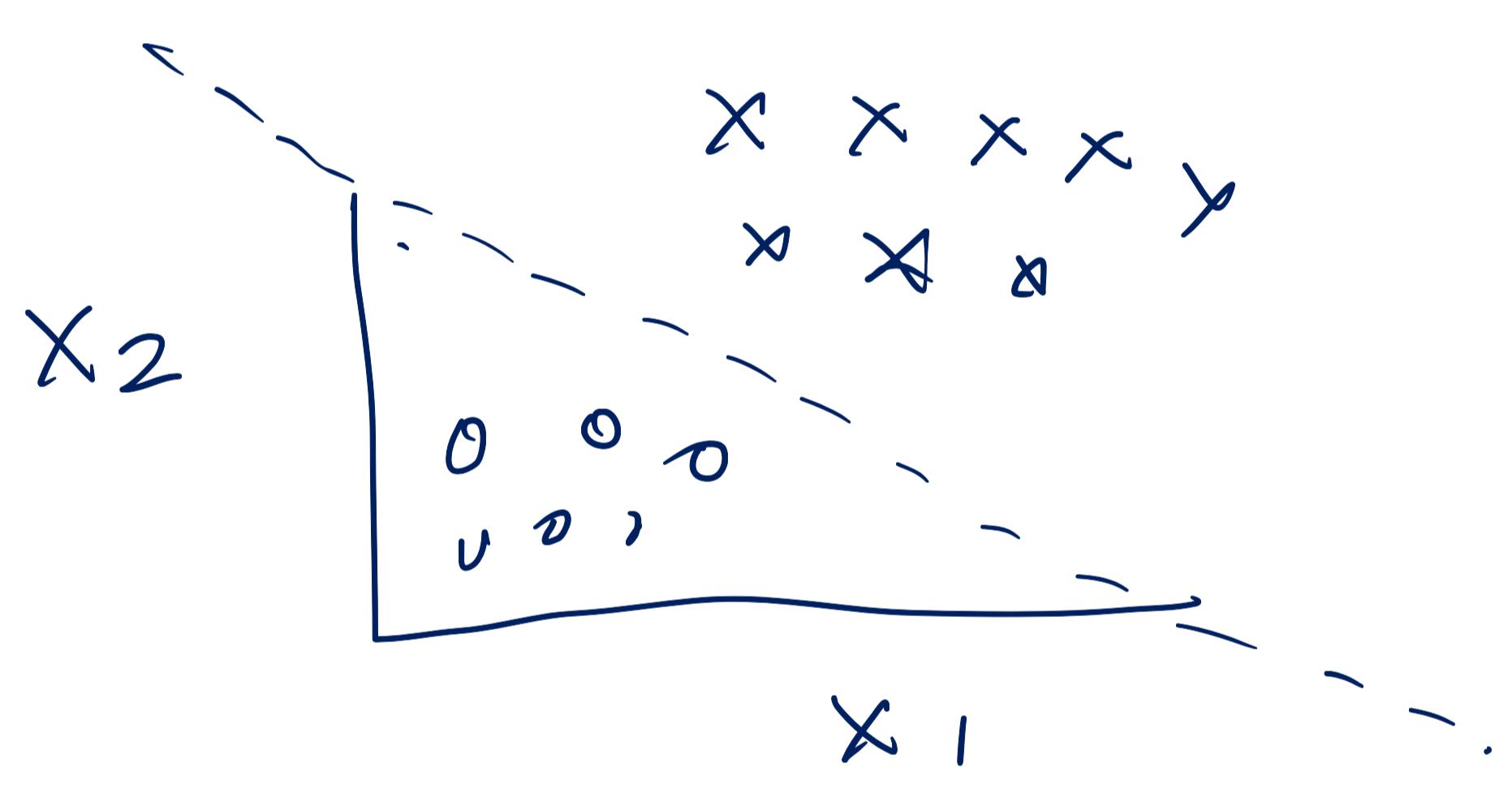
This is a hyperplane, so the decision boundary of logistic regression is a hyperplane

$$\text{Eg. } h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

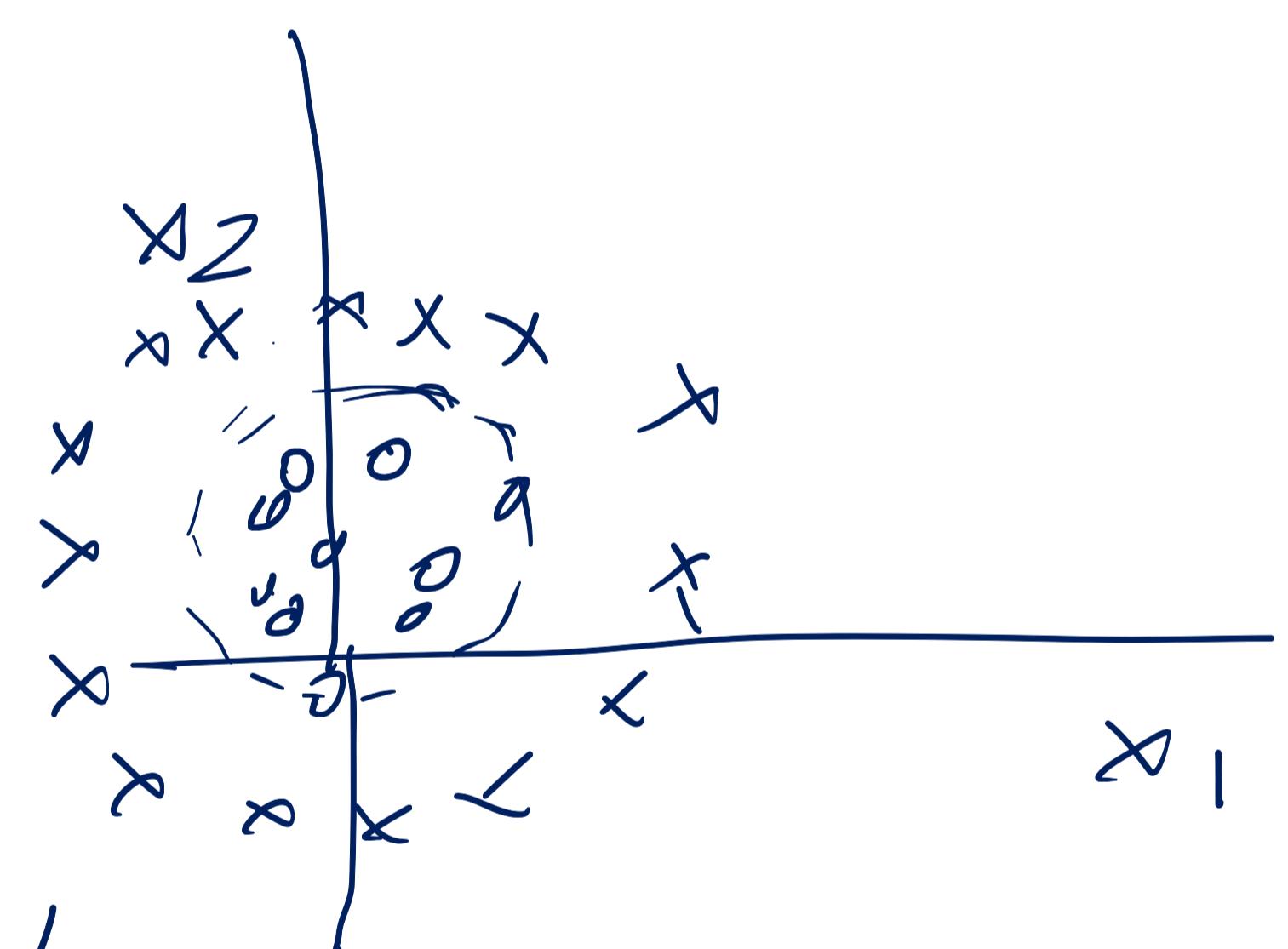
$$\text{Eg. } \theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

$$\text{Then } h_\theta(x) \geq 0.5 \text{ when } \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$$

$$\text{or } x_1 + x_2 \geq 1.3$$



Suppose data is



This is not
Linearly separable

for a non-linear separation we have to use different functions

$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \dots)}}$$

$$\text{if we use } \theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$$

$$h_\theta(x) \geq 0.5 \text{ when } \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 \geq 0$$

$$\text{or } x_1^2 + x_2^2 \geq 1$$

Support Vector machines

In Logistic Regression we had: $\hat{y} = \frac{1}{1 + e^{-\theta^T x}}$

objective when y true output $y = 1$ then $\hat{y} = 1$

$$\text{i.e. } h_{\theta}(x) = 1$$

$$\text{i.e. } \theta^T x \gg 0$$

when $y = 0$ then $\hat{y} = 0$

$$\text{i.e. } h_{\theta}(x) = 0$$

$$\theta^T x \ll 0$$

$$\text{loss func} = -(\gamma \log \frac{1}{1 + e^{-\theta^T x}}) - (1-\gamma) \log \left(1 - \frac{1}{1 + e^{-\theta^T x}}\right)$$

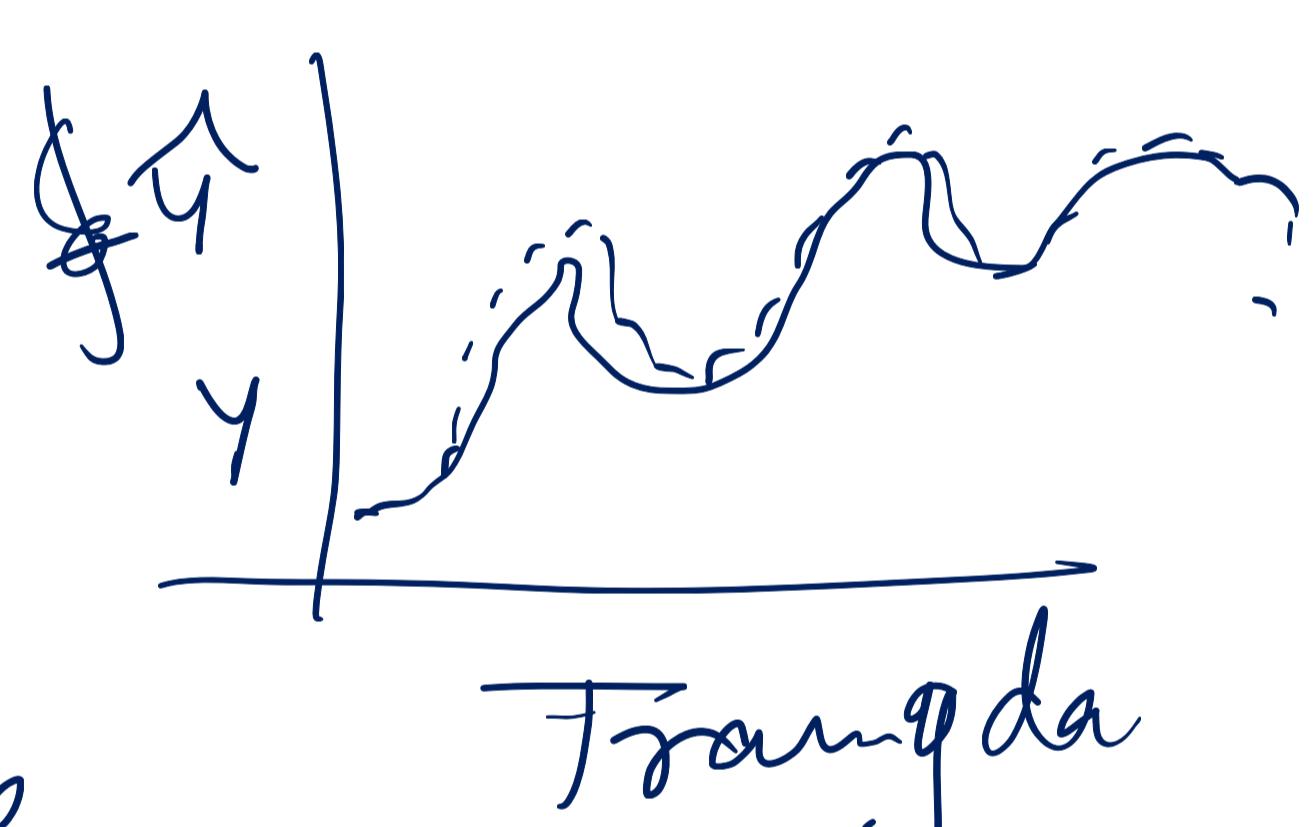
Generally for any loss func we also add a regularization term func $f(\theta)$

L2 regularization the $f(\theta) = \frac{1}{2} \sum_{i=1}^k \theta_i^2$ k is the no. of parameters

L1 regularization the $f(\theta) = \sum_{i=1}^k |\theta_i|$ overfitting

Regularized loss func = $L + \frac{1}{2} \sum_{i=1}^k \theta_i^2$

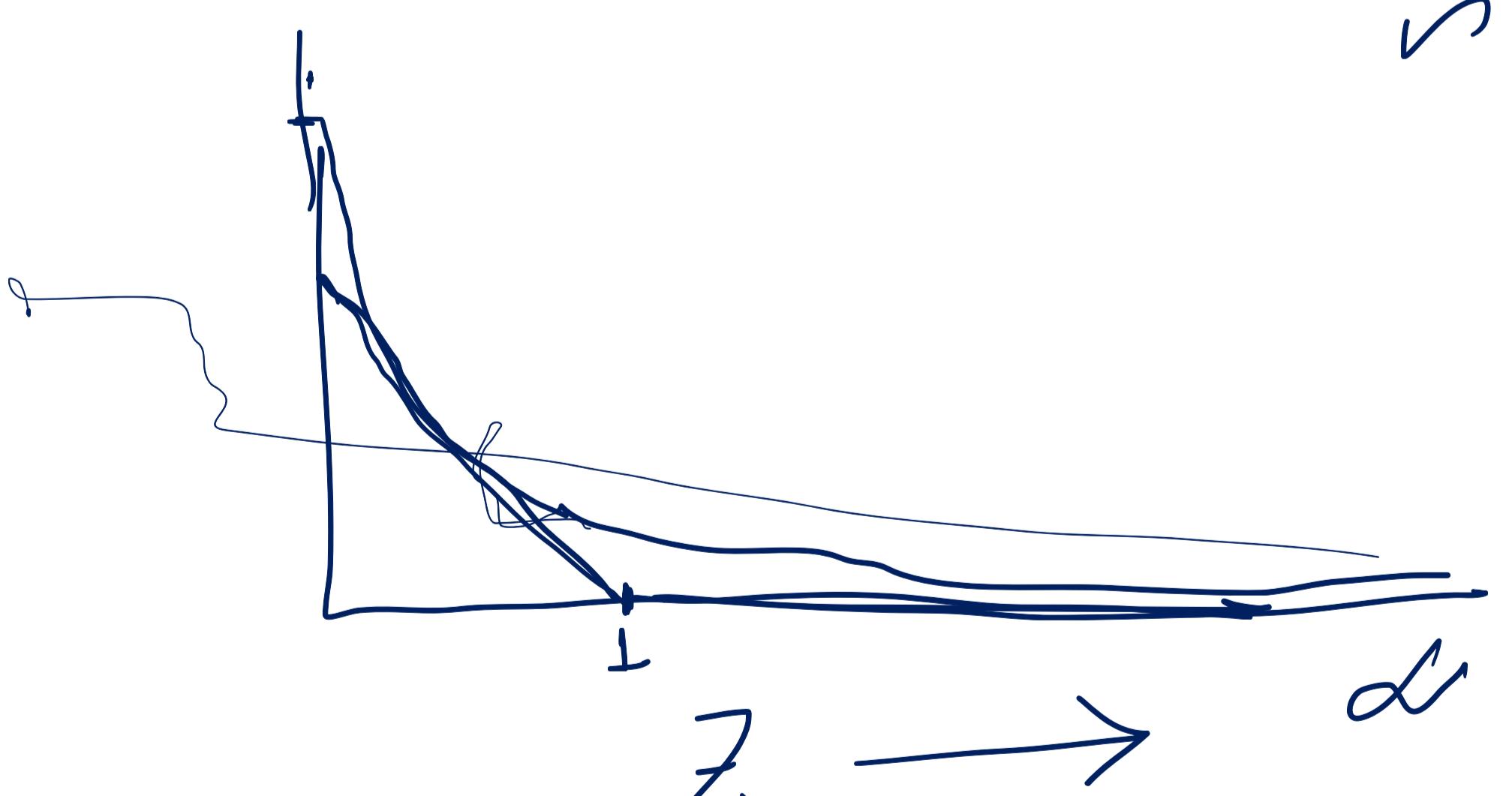
Regularization prevents the m/c learning model from overfitting.



For Logistic Regression

if $y = 1$ then in the loss func only $-(\gamma \log \frac{1}{1 + e^{-\theta^T x}})$ is active

$$\text{Let } \alpha \theta^T x = z$$



When $y = 0$ then the loss func only: $-(\log(1 - \frac{1}{1 + e^{-\theta^T x}}))$ is active

$$-\log \left(1 - \frac{1}{1 + e^{-\theta^T x}}\right)$$



For SVM we modify the cost func_s such that
as $y=1$ we have $\text{Cost}_1(z) = 0$ as $z \geq 1$

as $y=0$ we shall have $\text{Cost}_0(z) = 0$ as $z \leq -1$

Now loss func_s is $\frac{1}{m} \sum_{i=1}^m \left[Y_i \text{Cost}_1(z_i) + (1-Y_i) \text{Cost}_0(z_i) \right] + \frac{\lambda}{2} \sum_{j=1}^k \theta_j^2$

We may also use a weighted loss func_s that weights the regularization and the main loss func_n

$$\min_{\theta} \frac{\frac{1}{m} \sum_{i=1}^m [Y_i \text{Cost}_1(z_i) + (1-Y_i) \text{Cost}_0(z_i)] + \frac{\lambda}{2m} \sum_{j=1}^k \theta_j^2}{A + B} \quad - (4)$$

so the optimizat_e func_n can be written as

$$\min_{\theta} C A + B \quad \text{where } C \text{ is a constant that controls weight of } A \& B.$$

So what happens when C is kept very high
Then original minimizer of the loss func_n tends to make $A \approx 0$ and hence chooses θ accordingly

So when $Y_i = 1$ for $\text{Cost}_1(z_i)$ i.e. $\text{cost}_1(\theta^T x) = 0$

when $\cancel{\theta^T x \geq 1} \quad \theta^T x \gg 1$

the $Y_i = 0$ for $\text{Cost}_0(z)$ i.e. $\text{Cost}_0(\theta^T x) = 0$

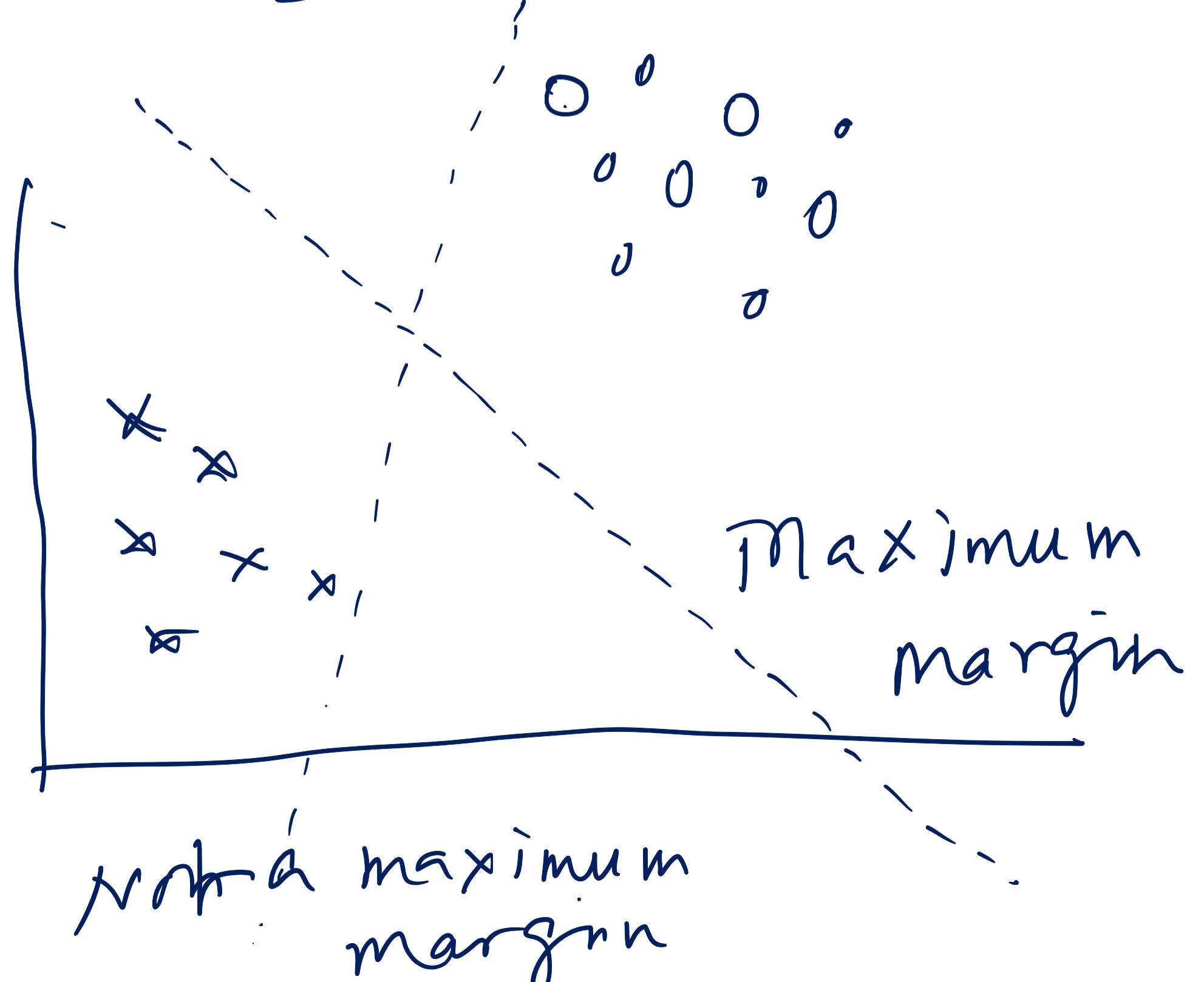
$\theta^T x \leq -1$

So the modified loss func_n in Eqn 4 can be written as

$$\text{Minimize} \frac{1}{2} \sum_{j=1}^m \theta_j^2 \quad \leftarrow \|\theta\|^2 \quad - (5)$$

s.t. $\theta^T x \geq 1$ whe $y=1$

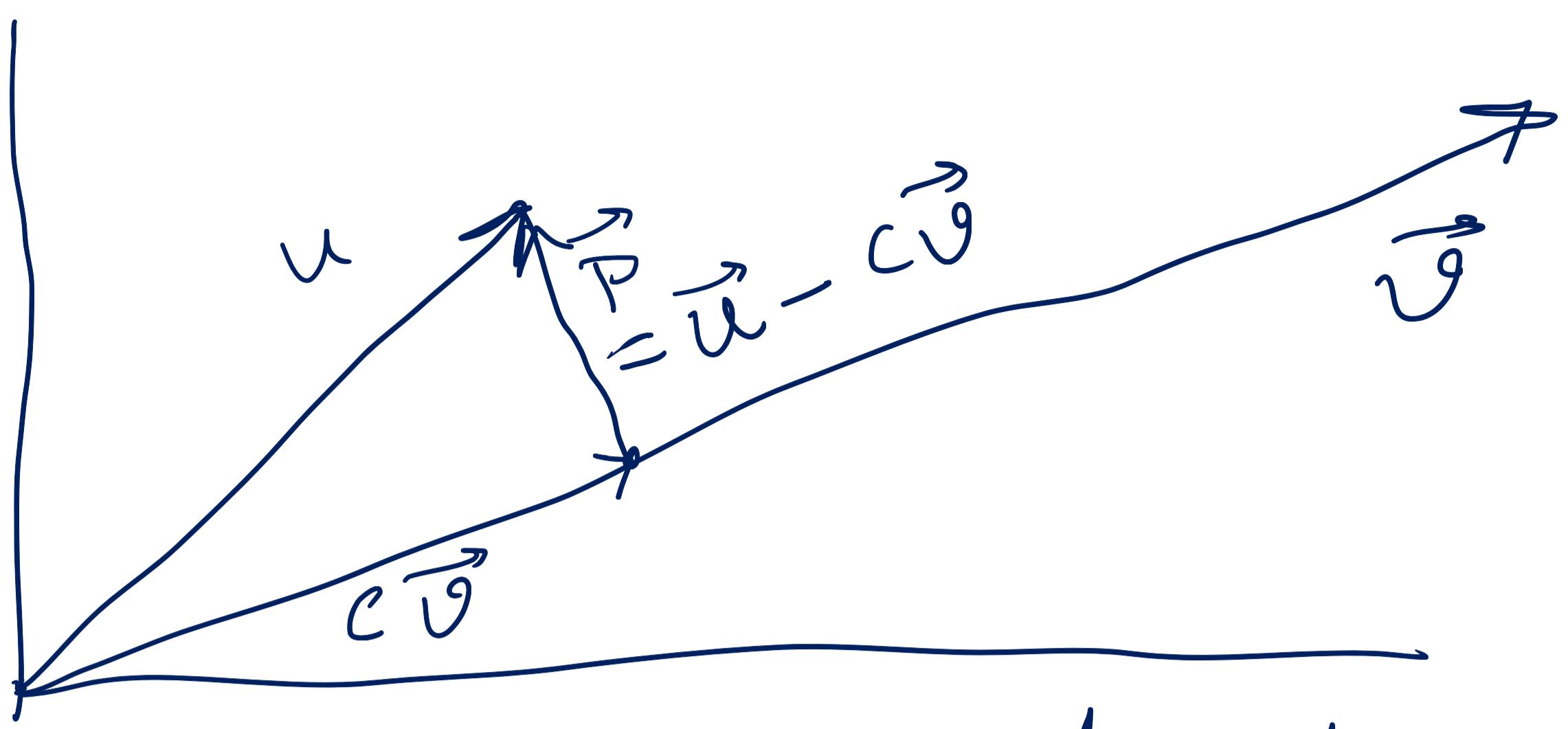
2 $\theta^T x \leq -1$ whe $y=0$



The maths behind the large margin

For 2 vectors \vec{u} & \vec{v} $u^T v$ can be represented as

$$\vec{\theta} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \\ \vdots \\ \theta_k \end{bmatrix}$$



$$c\vec{v} + \vec{p} = \vec{u}$$

$$\text{or } \vec{p} = (\vec{u} - c\vec{v})$$

\vec{p} is perpendicular to $c\vec{v}$

$$\text{So } \vec{p} \cdot \vec{c}\vec{v} = 0$$

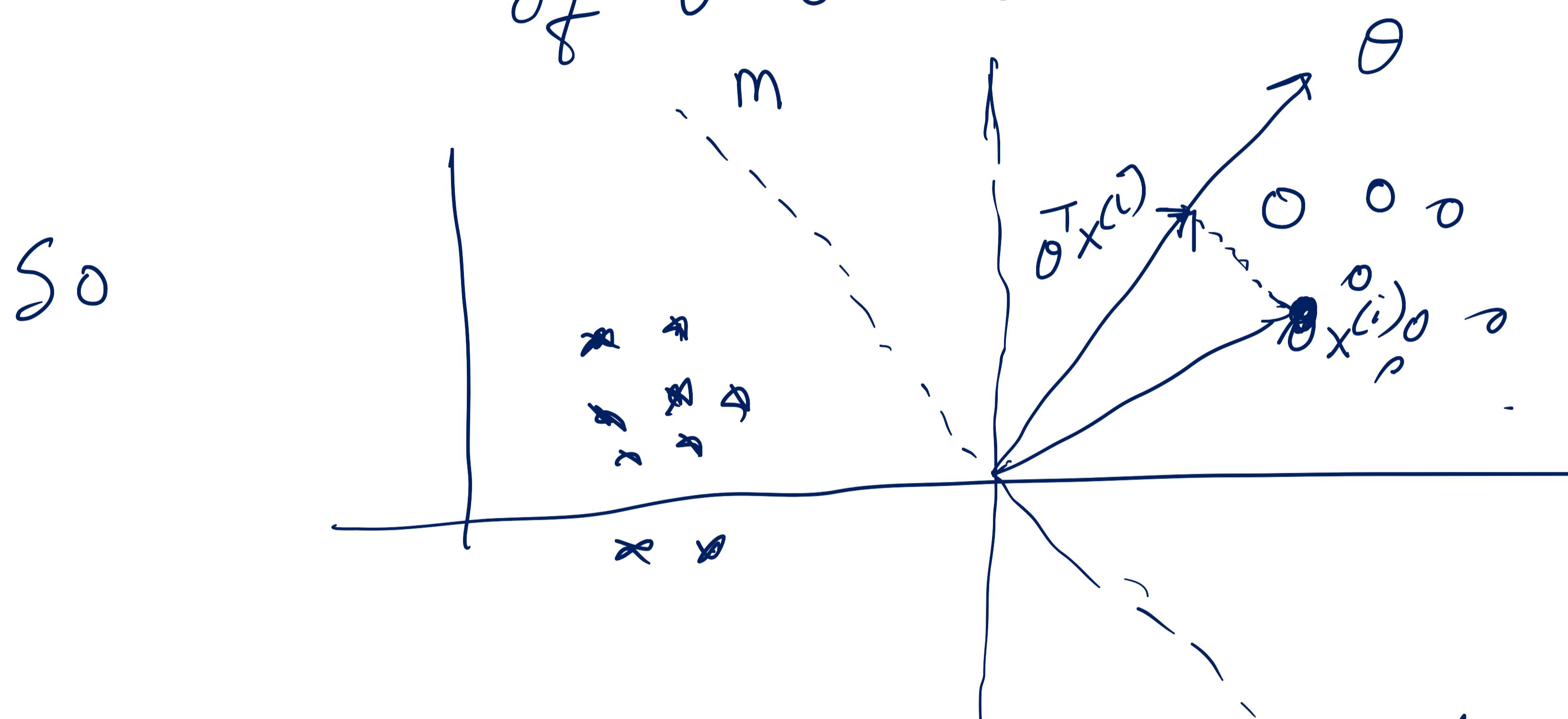
$$(\vec{u} - c\vec{v}) \cdot c\vec{v} = 0$$

$$\text{or } \ell(\vec{u} \cdot \vec{v}) = c^2 \|\vec{v}\|^2$$

$$\text{or } \vec{u}^T \vec{v} = c \|\vec{v}\|^2$$

$$v^T v$$

So $\vec{u}^T \vec{v}$ gives the magnitude or the length of the projection of \vec{u} on vector \vec{v}



So $\theta^T x^{(i)}$ gives the length of the projection of $\vec{x}^{(i)}$ on $\vec{\theta}$

If we consider $\vec{\theta}$ as normal to a margin m then if the margin is very close to $x^{(i)}$ then the projection length reduces and hence to have $c \|\theta\|^2 > 1$

$$\Rightarrow \theta^T x^{(i)} = c \|\theta\|^2 \quad \|\theta\|^2 \text{ must increase}$$

which is a conflict as

we minimize Egn 5.

For non-linear decision boundary.

Predict $y = 1$ if

$$\text{Say } \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \theta_5 x_2^2 > 0$$

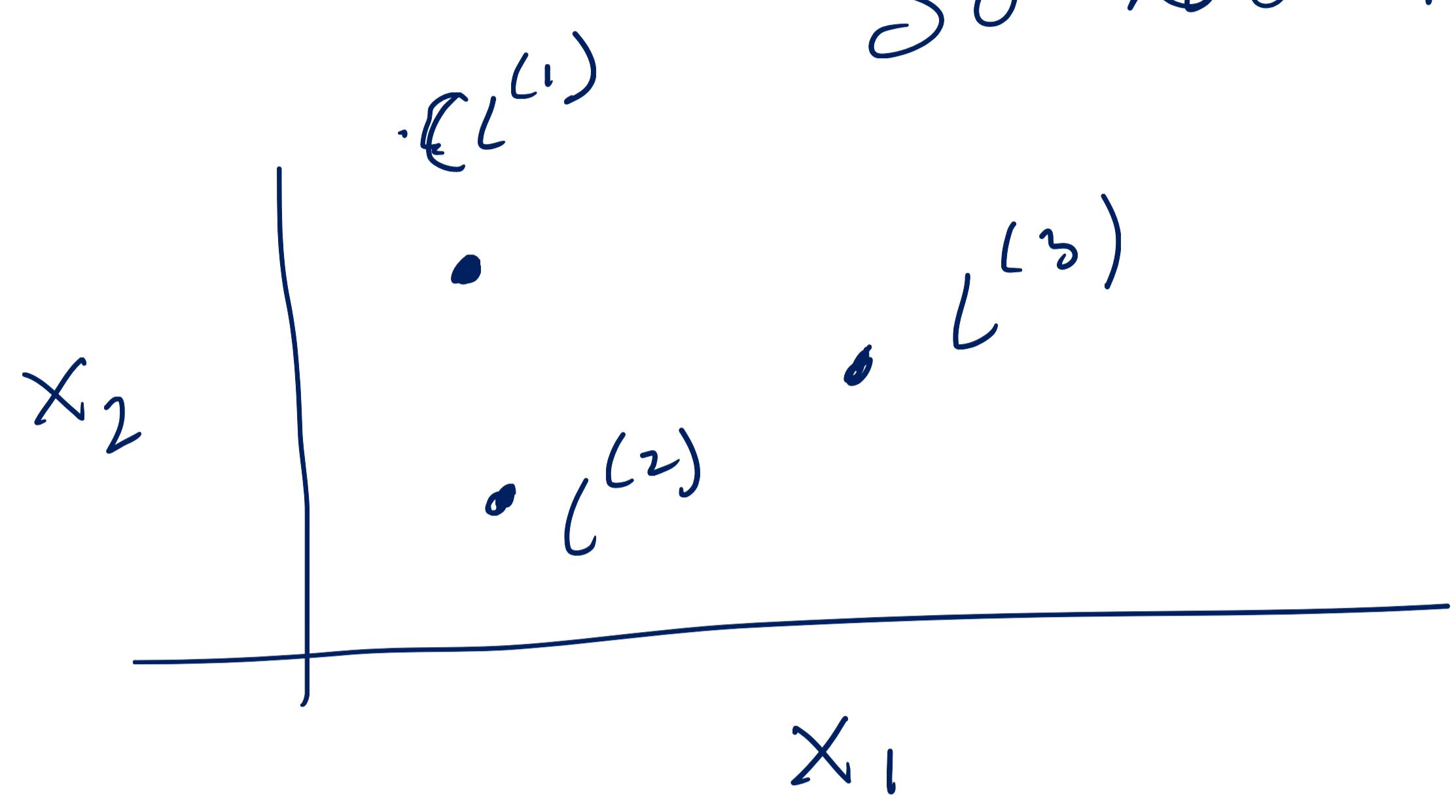
So hypothesis is now $h_{\theta}(x) = 1 \text{ if } \theta_0 + \theta_1 x_1 + \dots > 0$
= 0 \text{ otherwise.}

Other approach is create new features by combining feature i.e

let ~~x_1, x_2~~ $f_1 = x_1, f_2 = x_2, f_3 = x_1 x_2, f_4 = x_1^2, f_5 = x_2^2$
and apply the previous concept.

So how to create new features

choose few land mark points
and find the distance of the
points to the land marks.



Given x : $f_1 = \text{similarity}(x, l^{(1)})$

One similarity function can be

$$\exp\left(-\frac{\|x - l^{(1)}\|_2^2}{2\sigma^2}\right)$$

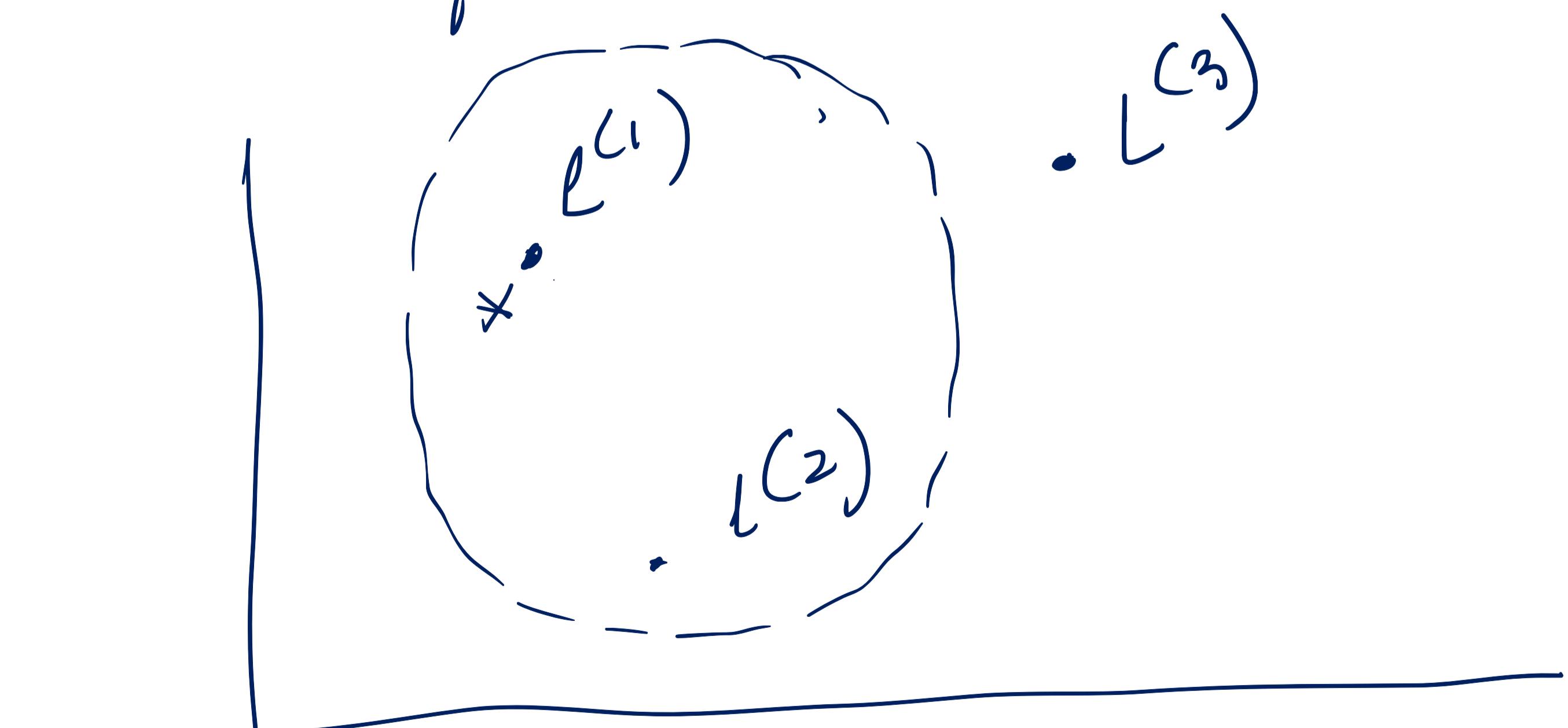
Similarly $f_2 = \text{similarity}(x, l^{(2)})$

$$\exp\left(-\frac{\|x - l^{(2)}\|_2^2}{2\sigma^2}\right)$$

So if $X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ $l^{(1)} = \begin{bmatrix} l_1 \\ l_2 \end{bmatrix}$

$$f_1 = \exp\left(-\frac{[(x_1 - l_1)^2 + (x_2 - l_2)^2]}{2\sigma^2}\right)$$

What is exactly does?



$$\text{If } \theta_0 = -0.5$$

$$\theta_1 = 1$$

$$\theta_2 = 1$$

$$\theta_3 = 0$$

$$\text{Then } h_\theta(f) = 1 \text{ if } \theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 > 0$$

or $f_1 + f_2 > 0$

What should be the landmarks?

In practice, every point is considered as a landmark