

Name: L.Narasimhanan, ID: 2019AIML618, Course: PCAM ZC211
 Q:01.

Regression EC-3R Q1

Given dataset with size: n

No. of features $\vec{x} = 10$ [vector of 10; dimension 10×1]

$$\Rightarrow x^{(1)} = \begin{bmatrix} x_1^{(1)} \\ \vdots \\ x_{10}^{(1)} \end{bmatrix}_{(10 \times 1)} \quad \text{where } {}^{(1)} \rightarrow \text{superscript 1 is the first sample out of } n$$

$$y = \begin{bmatrix} y_1^{(1)} \\ y_2^{(1)} \\ \vdots \\ y_n^{(1)} \end{bmatrix}_{(n \times 1)}, \text{ we add } x_0 = 1 \text{ as a standard feature.}$$

$$\Rightarrow x^{(1)}, \begin{bmatrix} x_0^{(1)} & x_1^{(1)} & \dots & x_{10}^{(1)} \end{bmatrix}^T$$

$$\text{thus, we have the data matrix } X = \begin{bmatrix} x_0^{(1)T} & \dots \\ x_0^{(2)T} & \dots \\ \vdots & \vdots \\ x_0^{(n)T} & \dots \end{bmatrix}_{(n \times 11)}$$

Simple regression model is given as $\hat{y} = \theta^T x$. where θ is the model parameters vector for each feature $\Rightarrow \theta = [\theta_0, \theta_1, \dots, \theta_{10}]^T_{(11 \times 1)}$

$$\Rightarrow \hat{y} = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_{10} x_{10} = \theta^T x \text{ (in matrix-vector notation)}$$

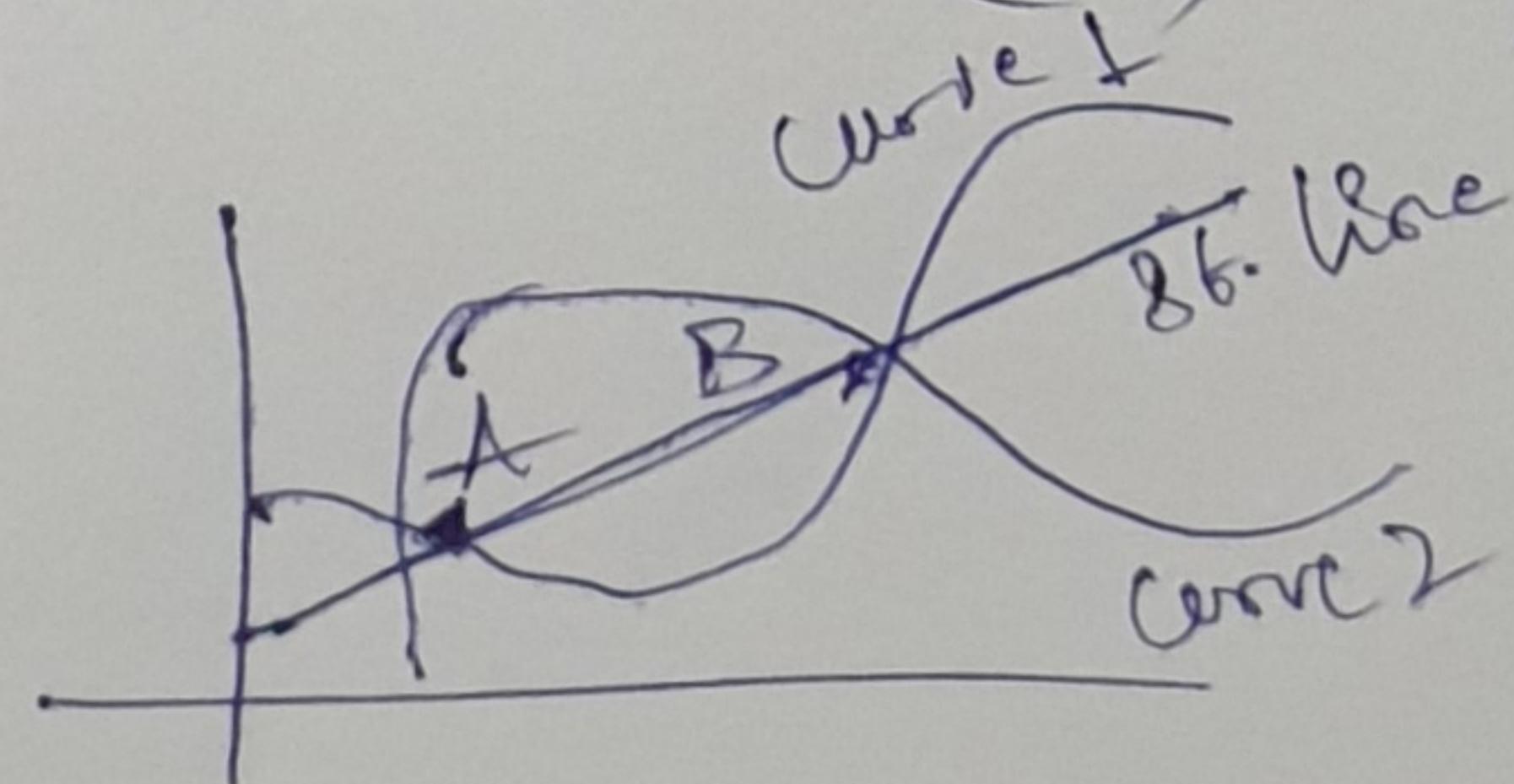
The empirical risk for the model is defined as:

$$J(\theta) = \frac{1}{2n} \sum_{i=1}^n (\hat{y}_i - y_i)^2, \quad \hat{y}_i = \text{fn. of } x_i \text{ parametrized by } \theta \\ \Rightarrow \hat{y}_i = h_\theta(x_i).$$

The goal of linear regression is to minimize the empirical risk

$$\Rightarrow \text{ERM (Emp. Risk minimization)} \text{ of } J(\theta) \Rightarrow \min_{\theta} J(\theta) := \frac{1}{2n} \sum_{i=1}^n [h_\theta(x_i) - y_i]^2$$

#1 If n is small ($n \leq 10$) \Rightarrow there is a chance of overfitting



there is a chance of fitting a curve of any shape can fit A & B

#2 n is large, overfitting can still occur with a large degree of polynomial chosen.

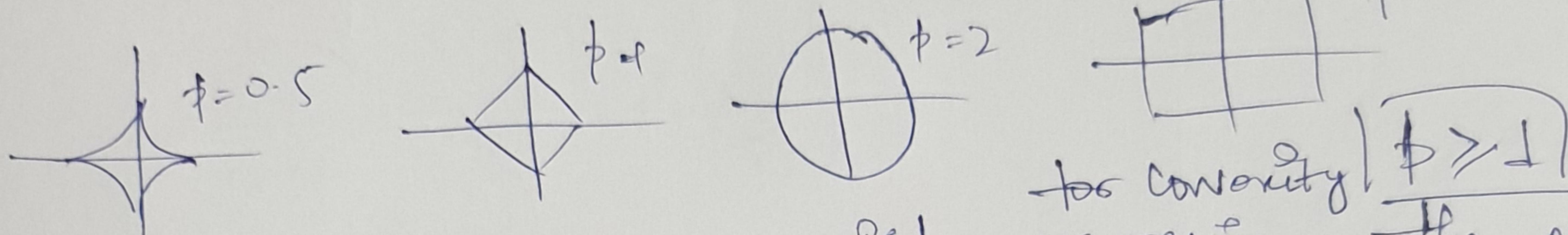
Hence we use regularization / shrinkage / stabilization to add bias to the model.

Name : Vigneshwaraman, ID: 2019A1ML612, Course: PCAM & C211
 Regression EC-3R C1

Q: Q Contd...

Hence, we minimize variance with added bias of the model by & called Regularized Empirical Risk minimization.

Regularizer: An ℓ_p norm is defined as $\left[\left(\sum_{i=1}^n \theta_i^p \right)^{\frac{1}{p}} \right]^p$ where p can range from $0 \rightarrow \infty$.



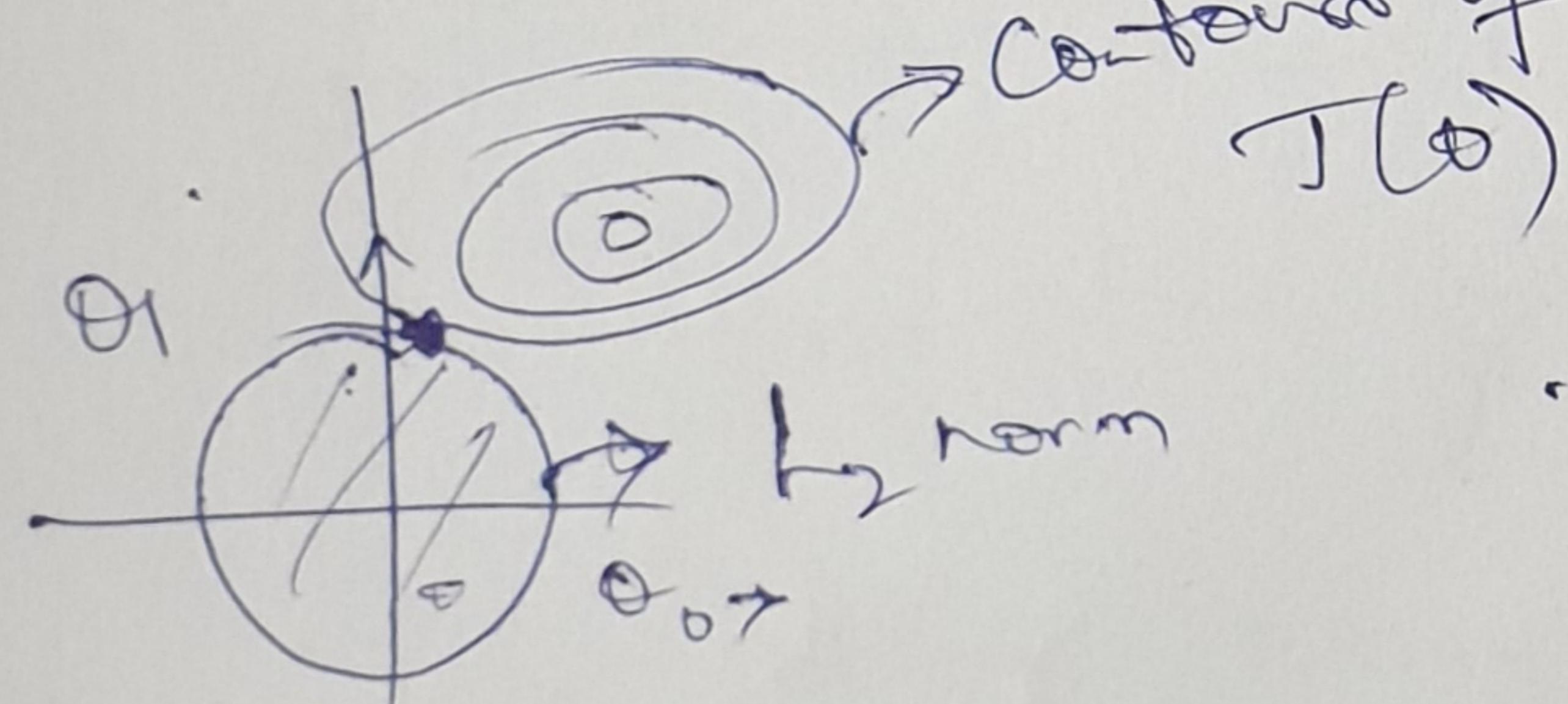
A) Ridge Regression: In Ridge regression, the norm ℓ_2 & the shape formed by the constraint is a circle.

& the shape formed by the constraint minimization fn.

$$\Rightarrow \sqrt{\sum_{i=1}^n \theta_i^2} = h_2 \quad \& \text{the constraint minimization fn.}$$

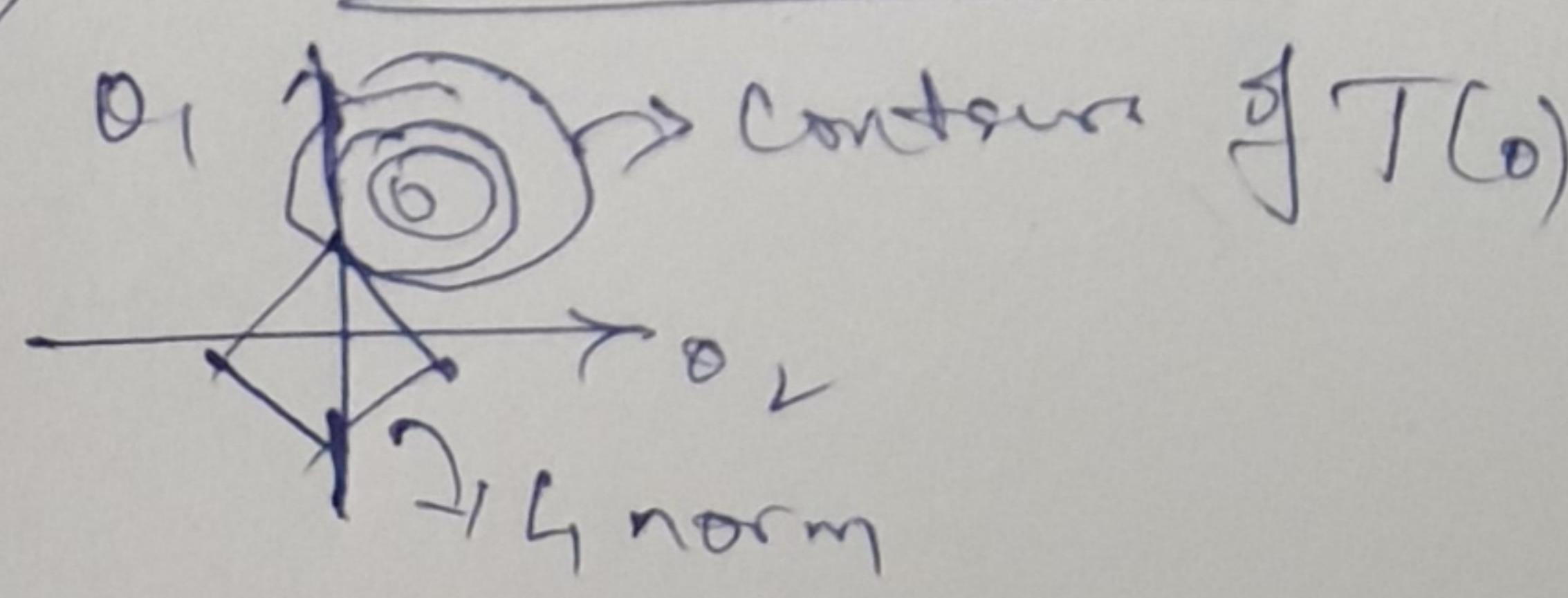
$$J(\theta) = \frac{1}{2n} \left[\sum_{i=1}^n (\hat{y}_i^{(i)} - y_i^{(i)})^2 + \lambda \sum_{i=1}^n \theta_i^2 \right]$$

\rightarrow we do not bias / regularize θ_0 & λ is the Lagrange multiplier for the constraint.



\therefore Note that the solution is when the contour of J touches the circle & the gradients are opposite.

B) Lasso Regression: Norm = 1 $\Rightarrow L = \left| \sum_{i=1}^n \theta_i \right| =$



$$J(\theta) = \frac{1}{2n} \left[\sum_{i=1}^n (\hat{y}_i^{(i)} - y_i^{(i)})^2 + \lambda \sum_{i=1}^n |\theta_i| \right]$$

again we do not penalize θ_0 / intercept term

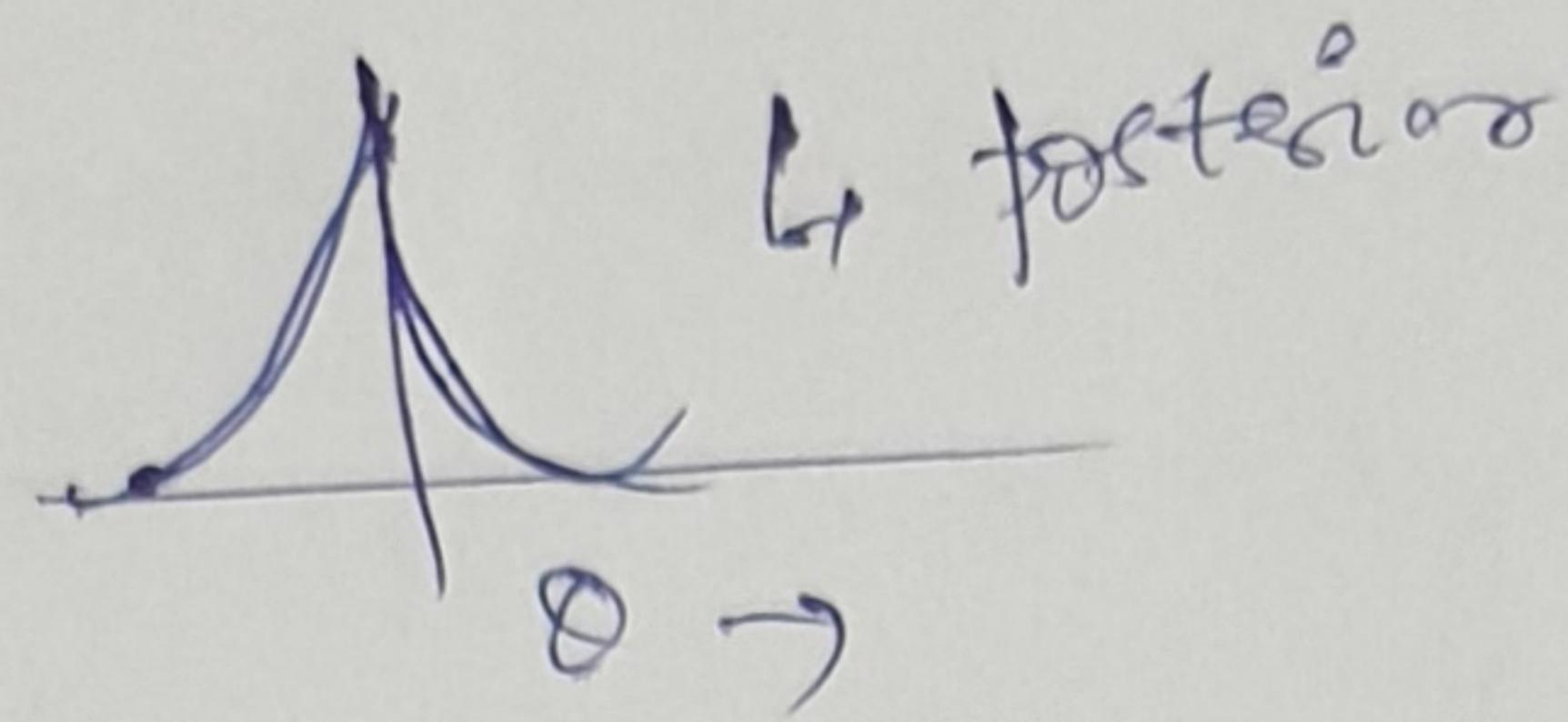
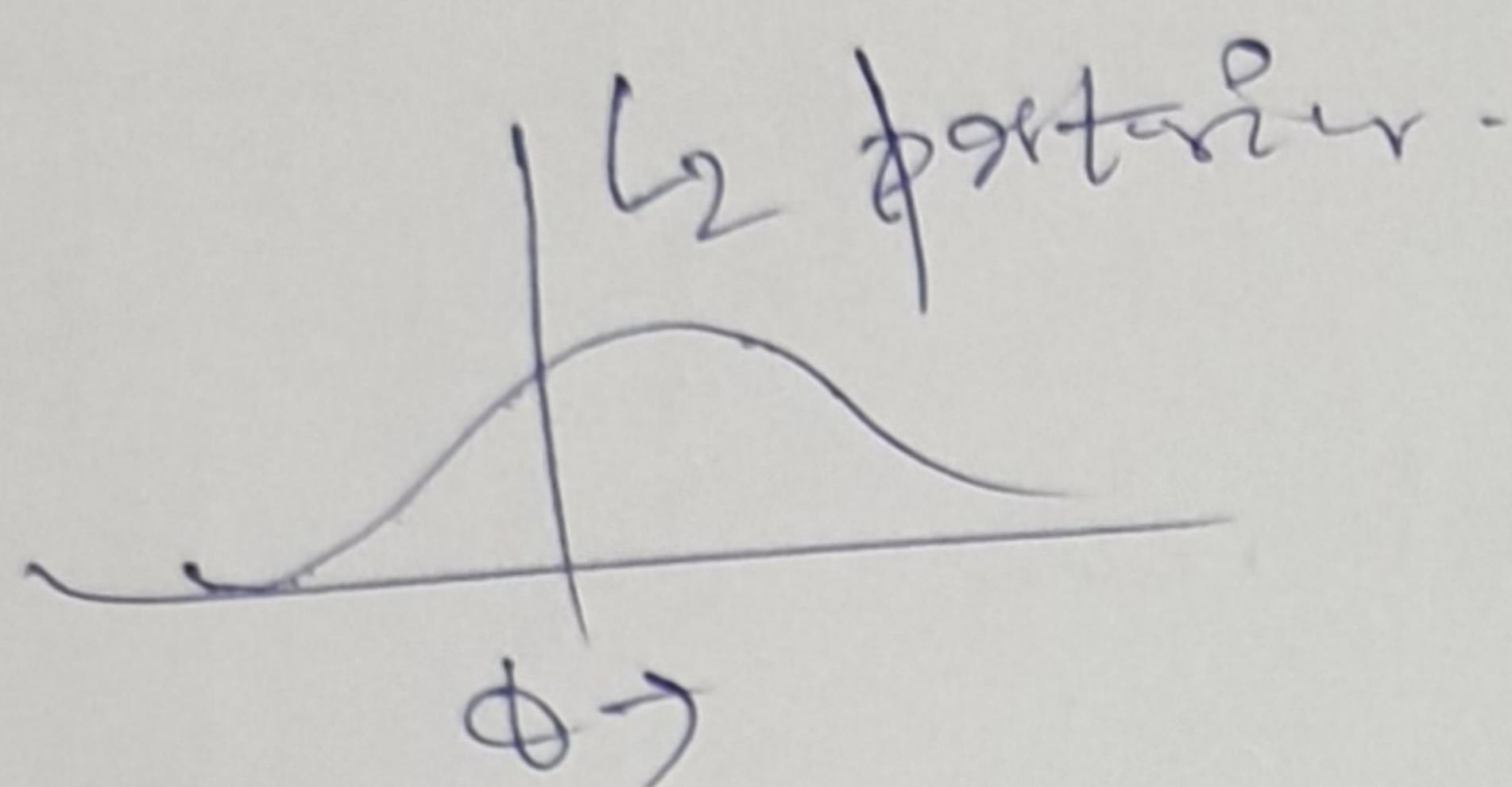
Name
Ayantha, ID: 2019AMLG18, Course: PCAMAC21

Registration No: BC3RCY

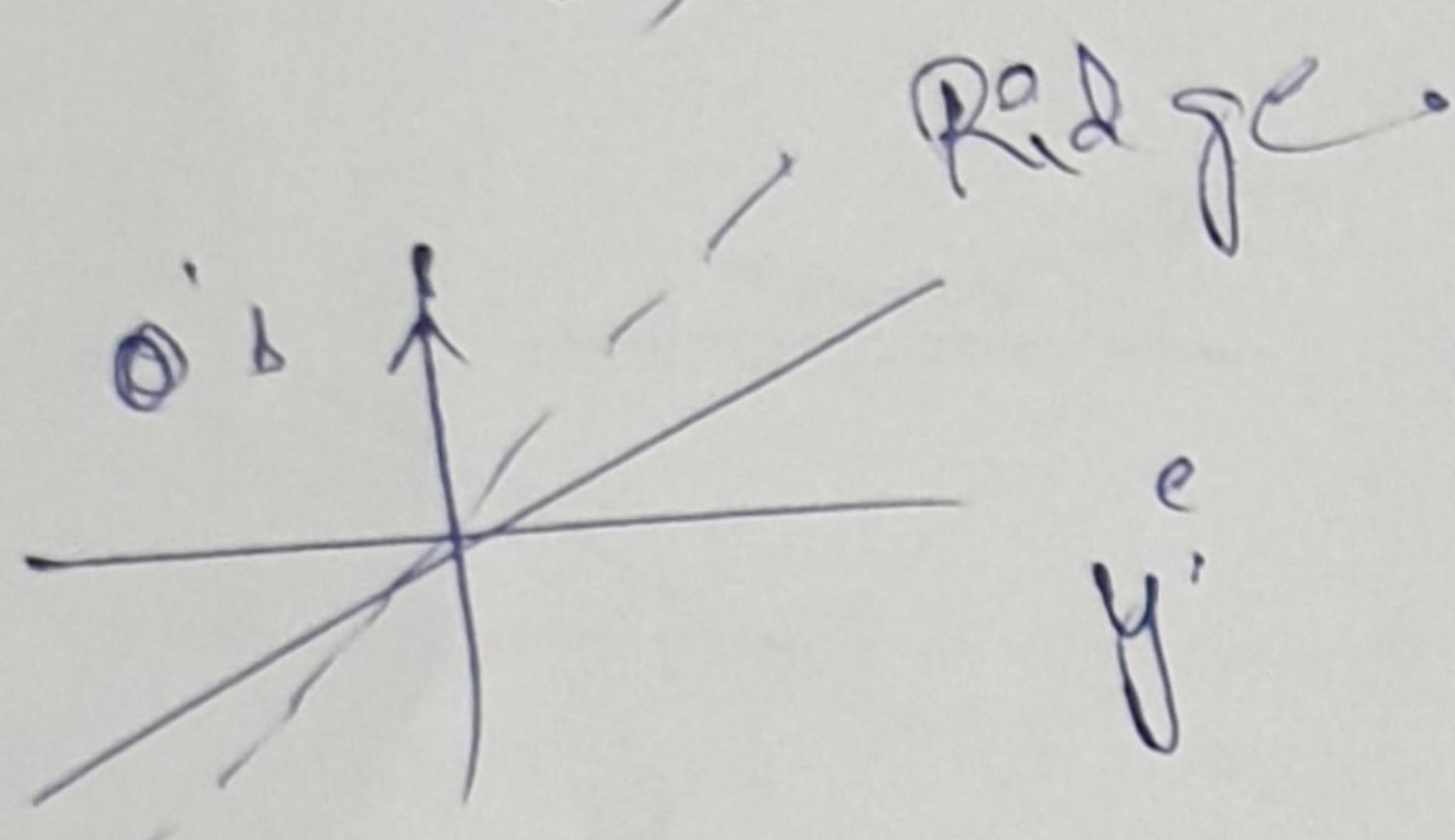
Q.9 Contd..

Condition: L_2 allows smooth solutions for convergence as it is differentiable & solutions are dense.

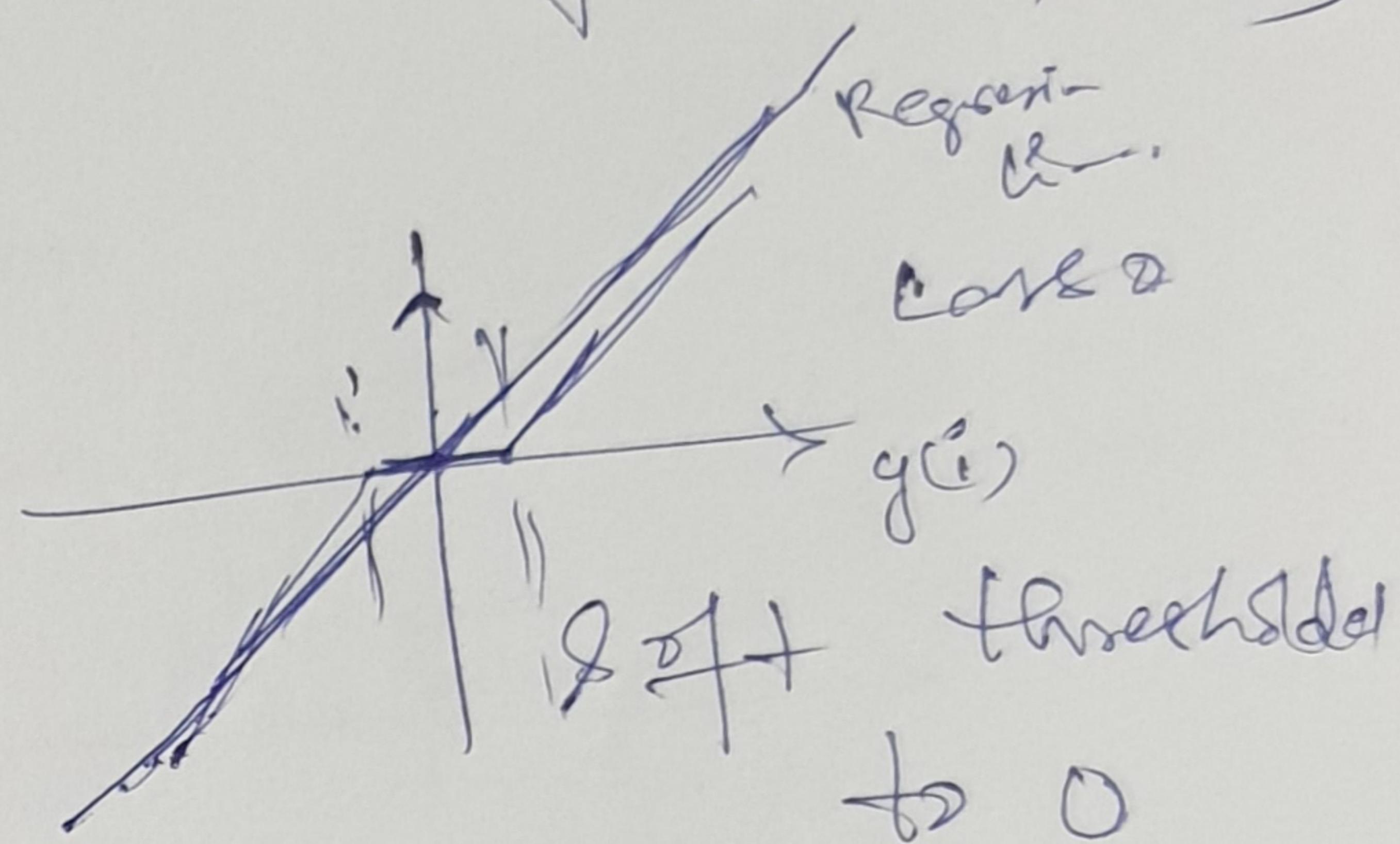
L_1 allows sparser solutions with some θ_i 's soft thresholded to 0.



{ posterior distributions
of L_2 & L_1 , note }



all θ_i 's
are
systematically
moved towards 0



This is because, from the regression diagrams in 7 & 8)
the anti alignment of Ridge lie with probability '0'
for $\theta_i \neq 0$.

while Lasso has
sharp turn & lie not differentiable.
i.e., allows many θ_i 's to be 0 while L_1 is still convex.

Name: Vigneshathamanu, ID: 2019-AIML 618, Grade ID: PGAM 2C20
 Regression EC-3RC1

Q.02.

Model 1

MAE: 900

R²: 68

RMSSE: 8000

Fscore: 23

Model 2.

1200

71

5690.

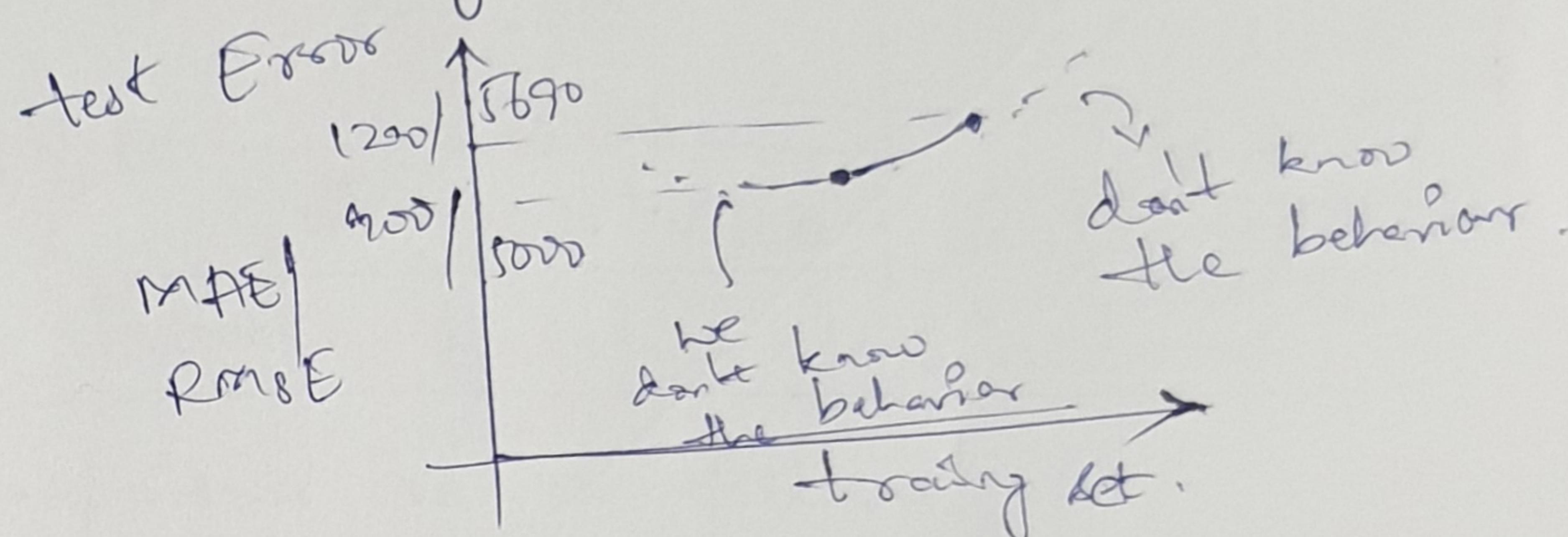
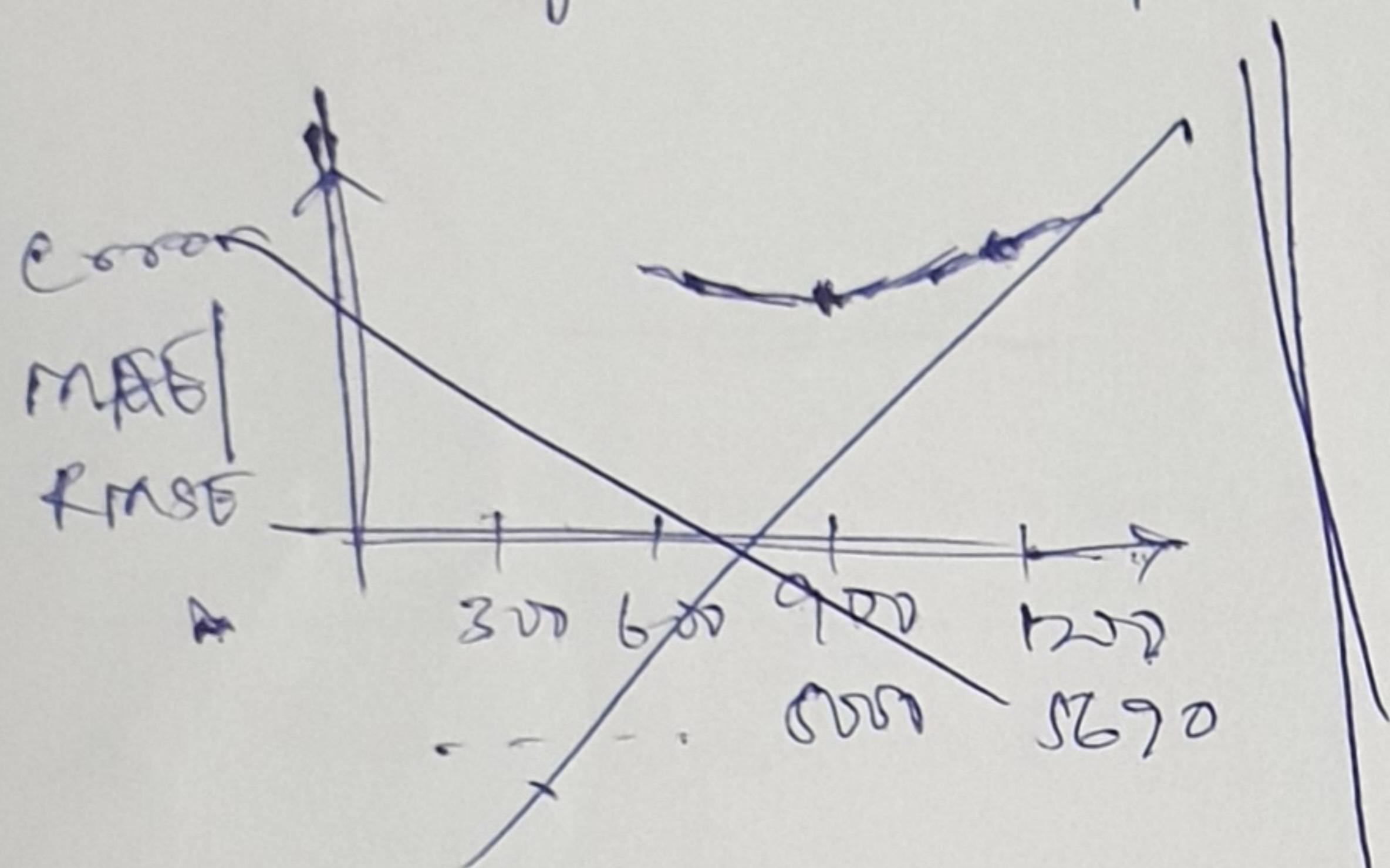
27.

Assumption

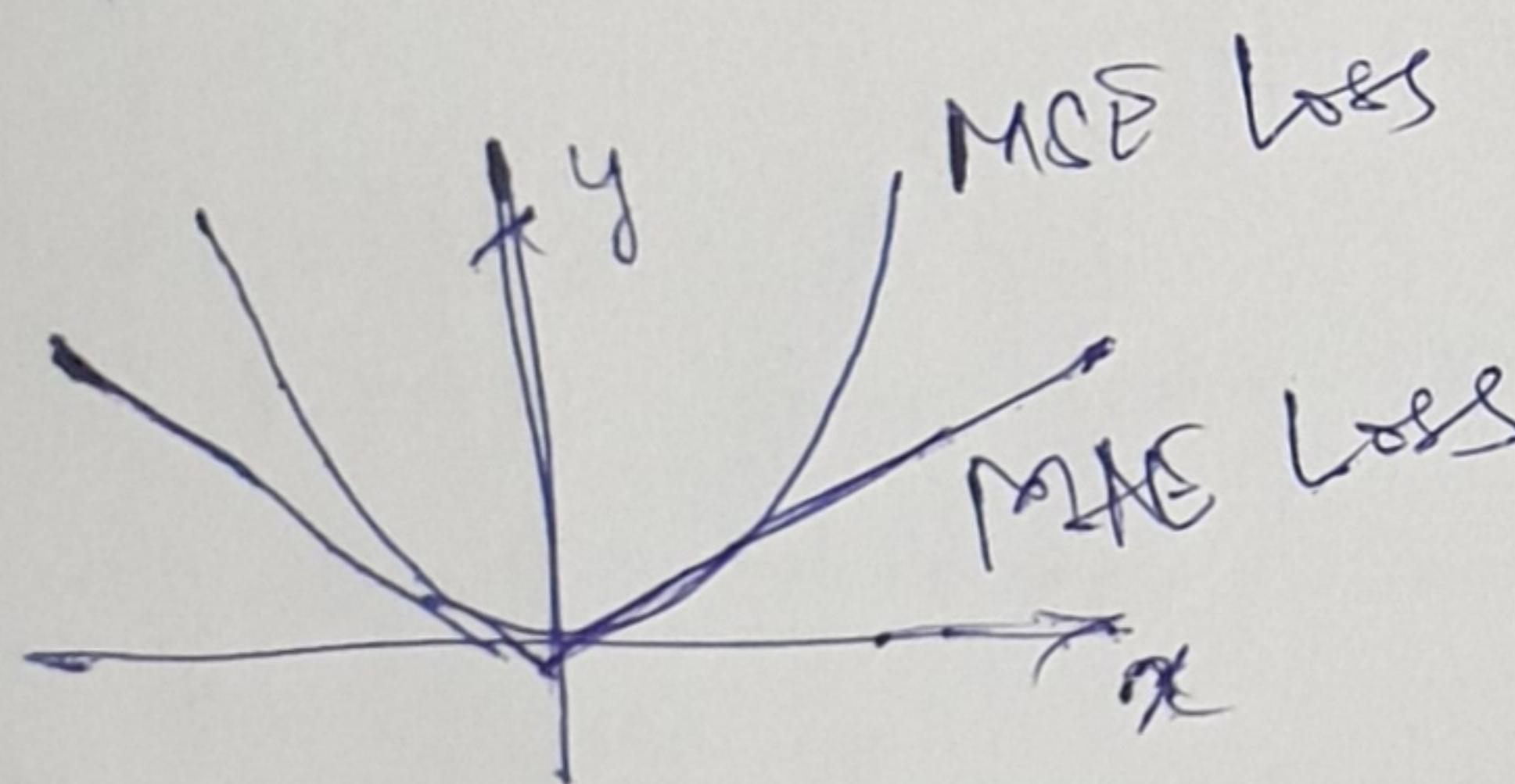
RMSSE to calculate
correct as there is

huge difference b/w MAE
& not in RMSSE

If we plot the model test error against MAE / RMSSE
we get the plot approximately.



→ With the info. I cannot claim one model over another
 → the deductions are below with regards to the data given.



→ These MAE & MSE plots are as shown on the left.

→ While MAE deviates less for outliers or wrong samples, MSE grows & penalizes wrong examples / outliers heavily.

CONCLUSIONS

MODEL - 1

1) Though R² (% variability explained by the model is lesser), so is the MAE & RMSSE
 → lower generalization error / risk

2) F-score > 1 \Rightarrow H₀ (no model parameters are related to response) is rejected

3) Choose this model to have lower risk & better generalization across other data

4) Slightly lower variance & higher bias model

MODEL - 2

1) R² is better than Model - 1 but the risk associated is more with MAE & RMSSE

2) F-score > 1 \Rightarrow H₀ - rejected.
 All model parameters are significant.

3) Choose this model when more variability needs to be explained & model with better performance & higher risk

4) Slightly greater variance & lower bias model

Q.03: 1 billion dataset, Algorithm comparison.

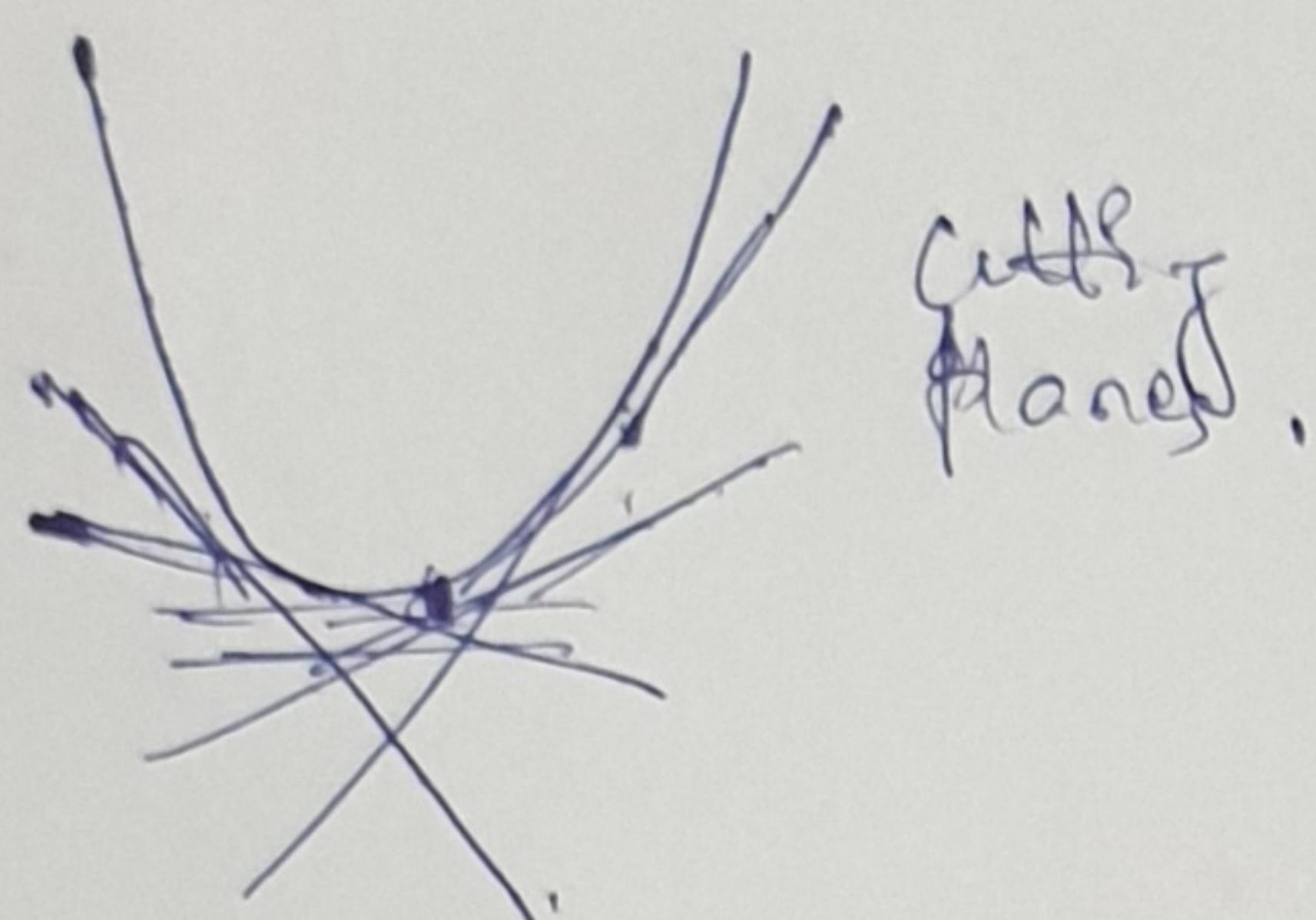
At the outset \rightarrow for $O(10^9)$, gradient methods will be too slow.

Prefer second order over first order gradient methods

A): Stabilized bundle method (Cutting plane) or LBFGS &

[BFGS will not work storing such huge data for memory.]

for calculating Hessian.

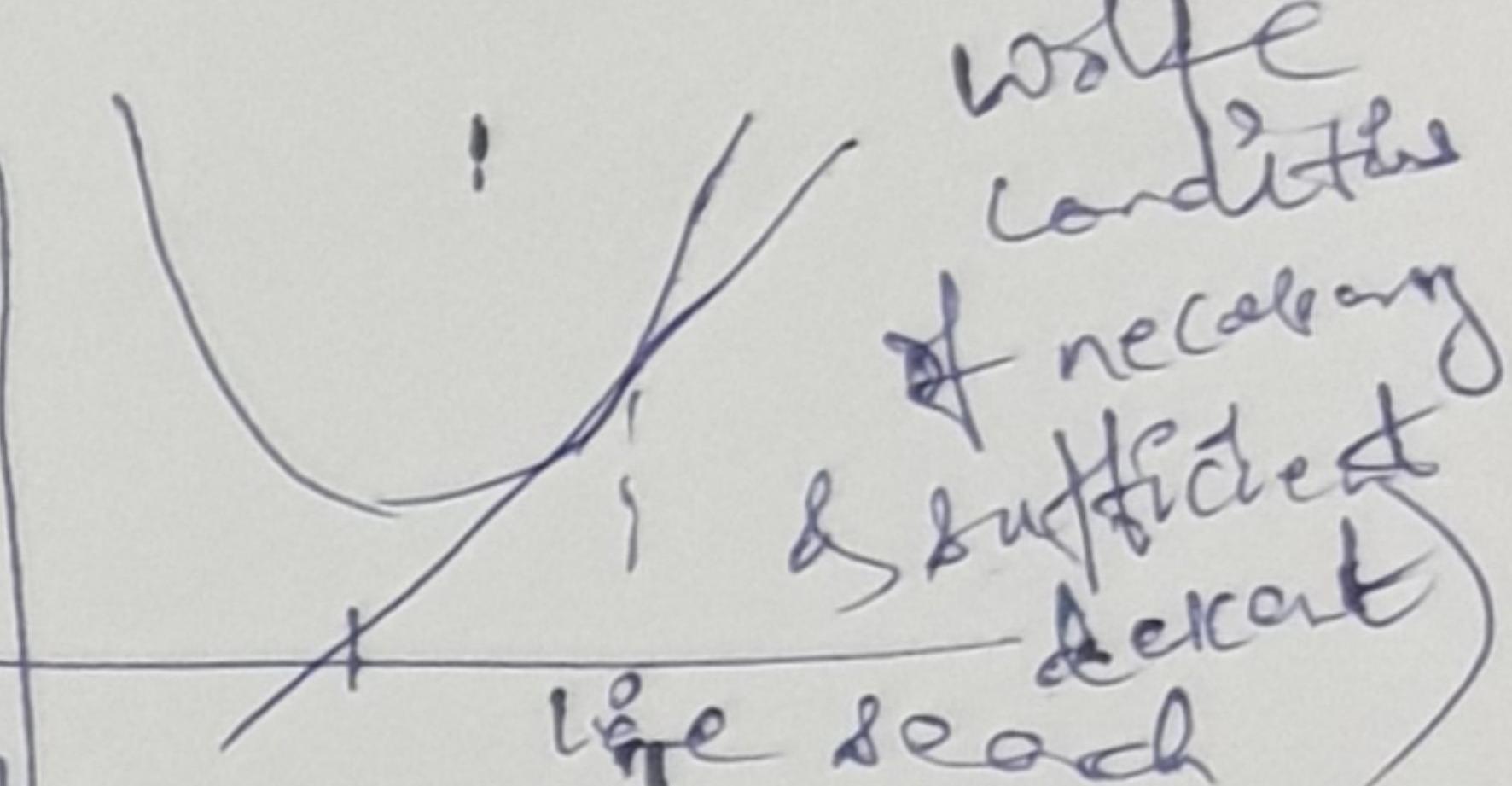


Newton-Raphson with line search (exact)

works pretty well too

$$g_{t+1} = \nabla f_t - \frac{f_t - f_0}{d_t}$$

$$\alpha_{t+1} = \alpha_t - \gamma H_t^{-1} g_t$$



line search obtained by

For the given comparison between batch & mini batch.

\rightarrow The algorithm & comparisons are as follows.

Vanilla (Batch G.D)

Mini-Batch G.D.

1) Algorithm, { for (N iterations or till convergence);

- \rightarrow calculate cost
- \rightarrow update parameter vector.
- \rightarrow check tolerance / convergence
- \rightarrow repeat }

1) Algorithm, { for (N iterations or till convergence).

- \rightarrow create shuffled batches (batch size is parameter)
- \rightarrow { for (each batch) ; update parameter vector }
- \rightarrow calculate cost, update parameter vector }
- \rightarrow check convergence
- \rightarrow repeat .

2) cost $m = \text{no. of samples.}$

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

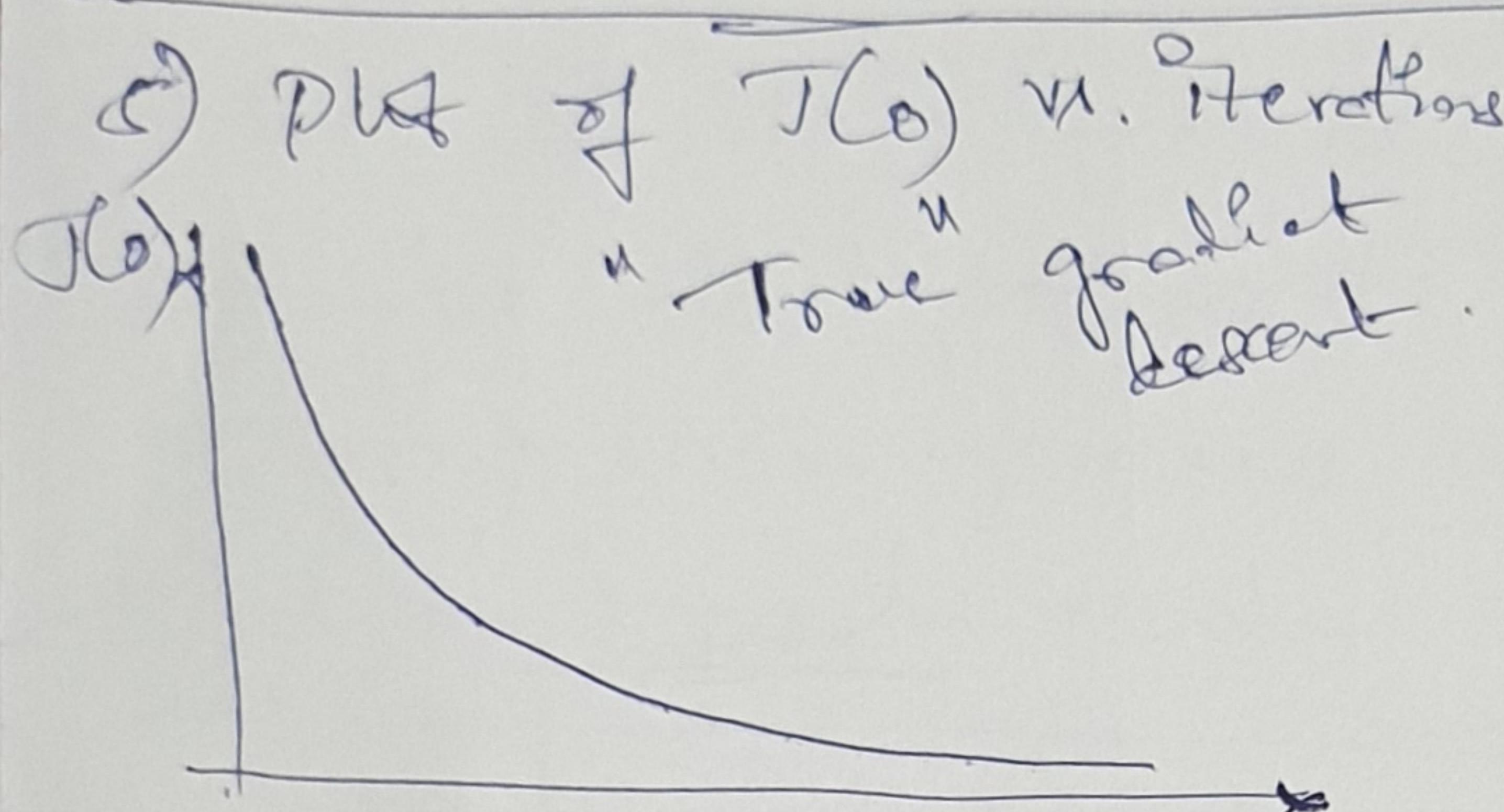
$$2) \text{Cost } J(\theta) = \frac{1}{2 \times \text{batch size}} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

3) See that batch descent iterates over all billion samples for a single step of descent or parameter calculation / update.

3) Mini-batch iterates only through the batch size, (typically chosen to be around 10-100 & can be tuned) for one step of param update / descent.

Q.03 Contd . . .

4) Calculating 1 billion is not feasible & need to use map-reduce & other techniques to parallel for it to work

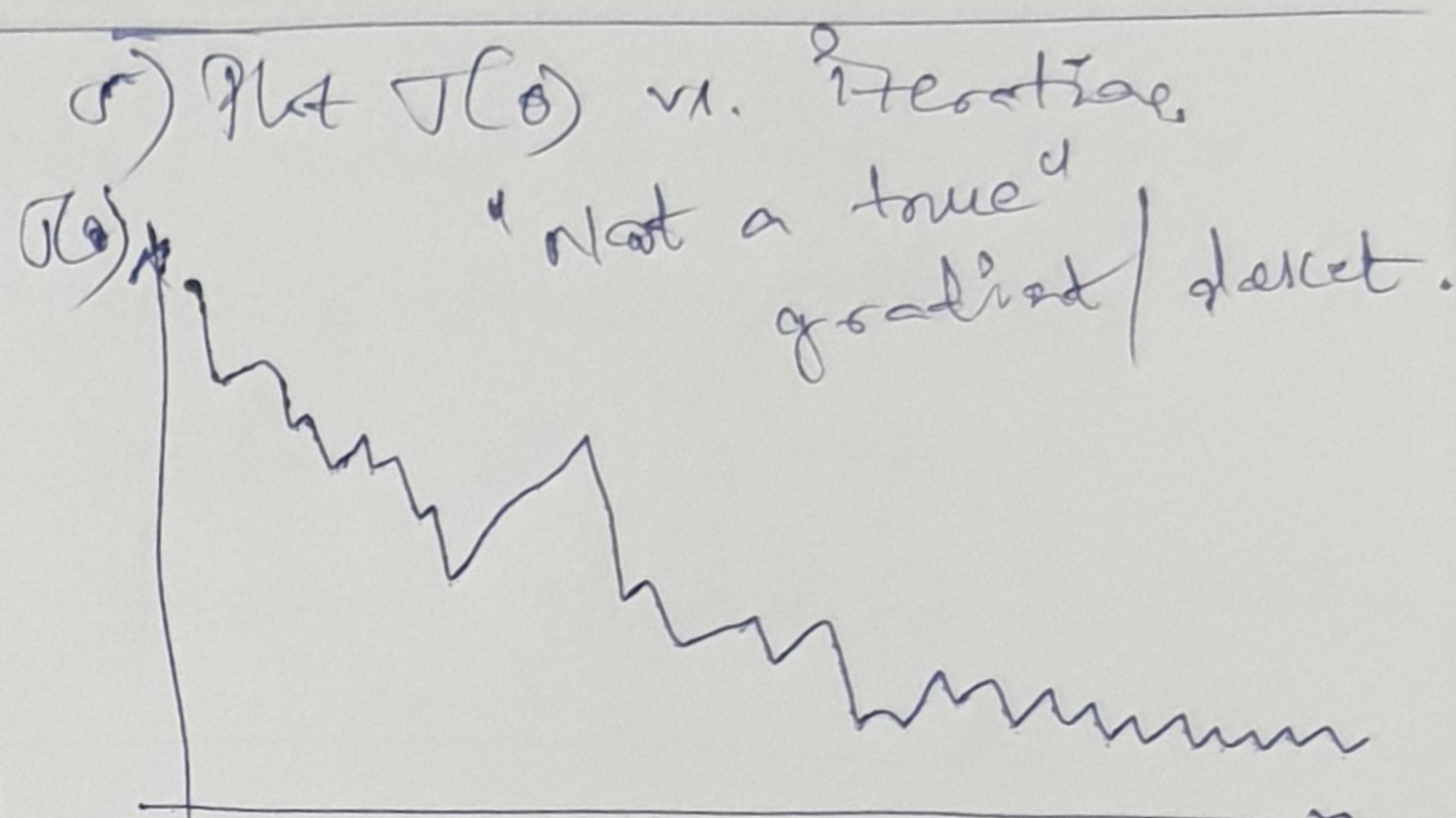


(d) Nearly worst choice of ' λ ' batch converges to the minimum but it can take forever to converge on 1 billion etc.
 May need $\lambda^{ten \text{ of thousands}}$ steps to converge $\approx O(10000)$

(e) Model created will have slightly different/better params than mini-batch (smooth convergence, see (5))

(f) This will be very slow to converge if at all & we will not use this for dataset of $O(10^{10})$

(a) Batch size can vary to get good convergence rates. Can also use map reduce for parallel batch-split-updates.



(c) Assuming correct ' λ ' choice, mini-batch converges faster compared to batch descent. Typically, chosen also good batch size, converges in 10-80 iterations (batch iteration)

(d) Model params will be slightly different as model, when it approaches minima, it oscillates as it is not a true gradient, but is typically near the minima, not much variation, see (5)

(e) Choosing dynamic ' λ ' to decay too as it approaches minima, can give very close result to batch descent while being very fast as well. Note this is typically not done.

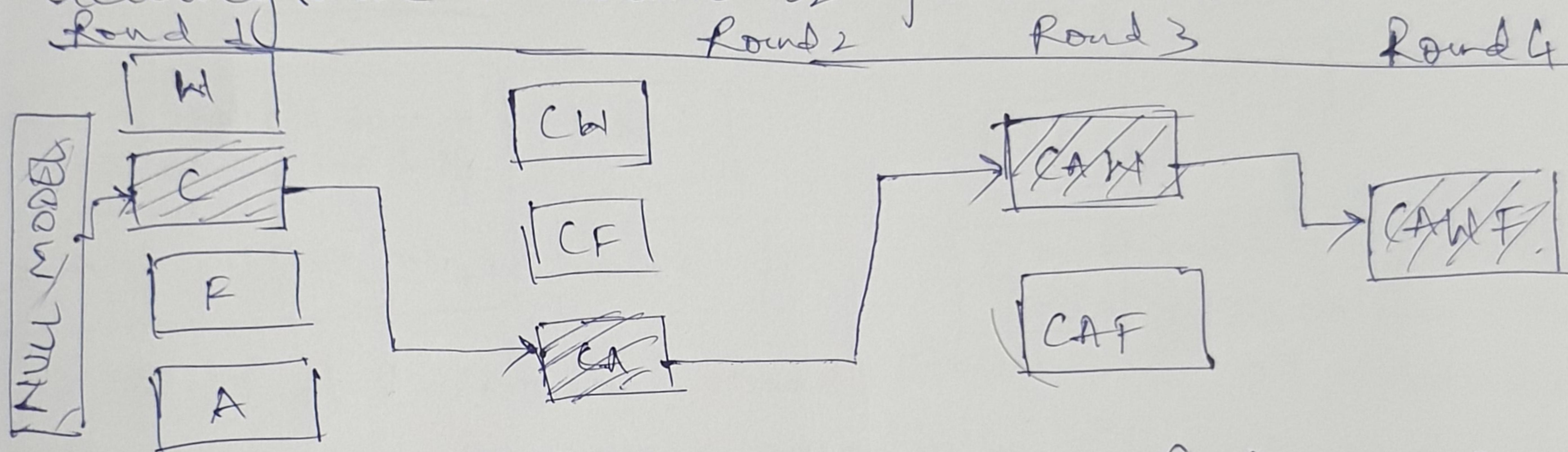
$$\lambda_t = \frac{\lambda_1}{\lambda_2 + t}$$

Q: 04

A) Let's say, Body weight = W, Caloric Intake = C, Fat Intake = F
 Age = A, Blood Cholesterol = Y

\Rightarrow features = W, C, F, A, Output = Y

Forward selection algorithm picks best feature starting with 1, 2 ... no. of features. Null model is first built assuming null hypothesis to be true & ~~then~~ ~~then~~ step by step new features are added till a good accuracy score is reached or feature threshold is reached.



\rightarrow In general for ~~1~~ 'n' predictors, $1 + 1 + (n-1) + \dots + n$ models are built = $1 + \sum_{i=1}^n i = (1+n)(n+1)/2$.

\rightarrow In this case, 4 features $\Rightarrow 1 + (4 \times 5)/2 = 11$ models (including null model)

\rightarrow The models are shown above with features at each step ^{round}.
 - Round 0: Null model
 - Round 1: C is the best feature, subject to some heuristic such as AIC, BIC or adjusted R^2 after adjusting for interactions of other predictors.
 - Round 2: C with A is best,
 - Round 3: C with AF is best,
 finally, CAWF are chosen

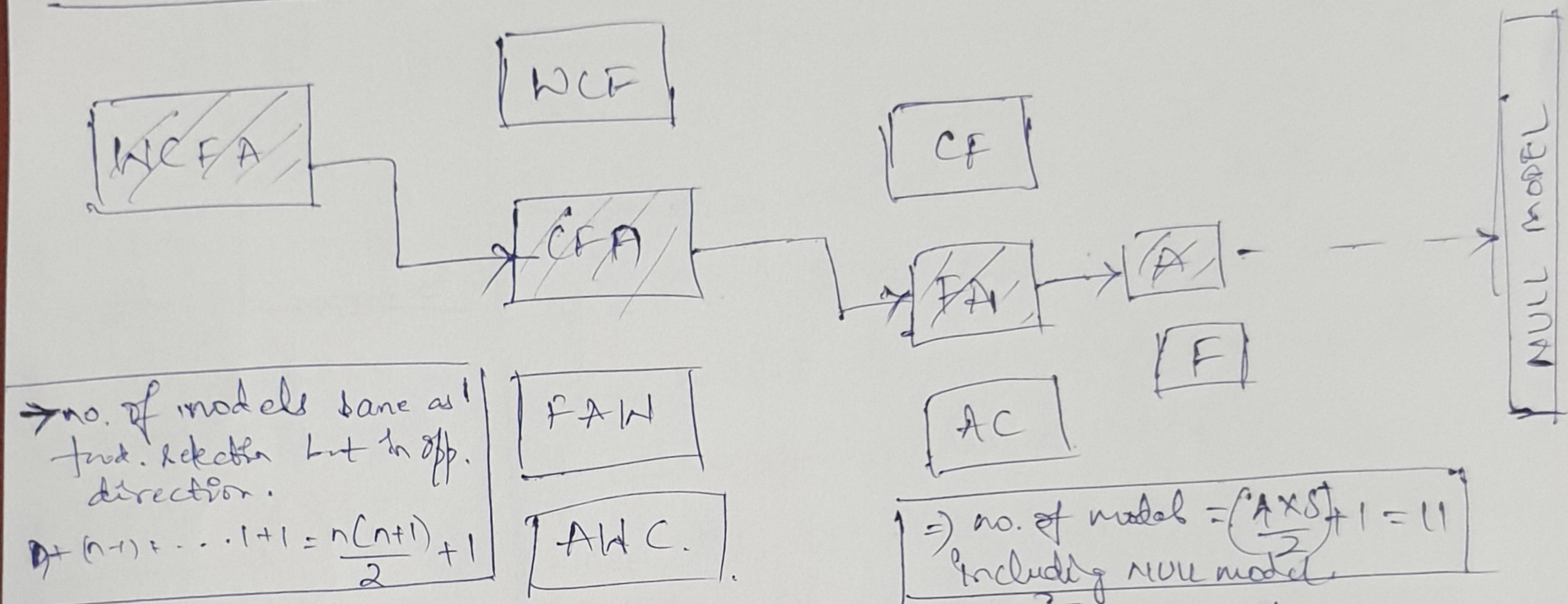
\rightarrow The above example was just one of the models that could have been chosen.

\rightarrow A best subset of 1 feature would be C, 2 \rightarrow CA, 3 \rightarrow CAF & so-on. No. of features needs to be decided.

Name: Vyarawathenam, ID: 2019AIPAL618 Core: PGAM 7E211
 Regression
 DD4 contd. -

B) Backward selection picks all features & then sequentially reduces each, based on how many features to keep or a accuracy score based thresholding such as AIC, BIC & adjusted R², etc.

Round 1 Round 2 Round 3 Round 4



- Opp. of fwd. selection. Model starts with all features dropping one by one sequentially till null model is reached.
- Ex. as shown, note, this is just one possible sequence.
- Round 1: WCFA (all), 2: CFA chosen as best feature (W is dropped)
- Round 3: FA (C is dropped), 4) F is dropped & AC chosen

Conclusion: Unless no. of features is chosen ~~not~~ to be a stopping criteria (like 3best, 2best, etc.) all features are chosen. ~~Process stops~~ Or a threshold accuracy score should be the stopping condition.

Fwd. selection ~~does not~~ be greedy & a feature once selected is fixed & cannot be dropped.

Bwd. selection can have max. k features where k < n when no. of features (n) > no. of samples (k).

Q.05: Examples of univariate & multiple regression.

Practical example of regression: Predict house price based on N features, $N=1 \dots n$

1) Simple regression / univariate

Let's assume that house price only depends on the house sq-ft.
 & area. Let's define 'Area' in (m^2) (metres²) as the feature
 & price of house as feature y which is the target.

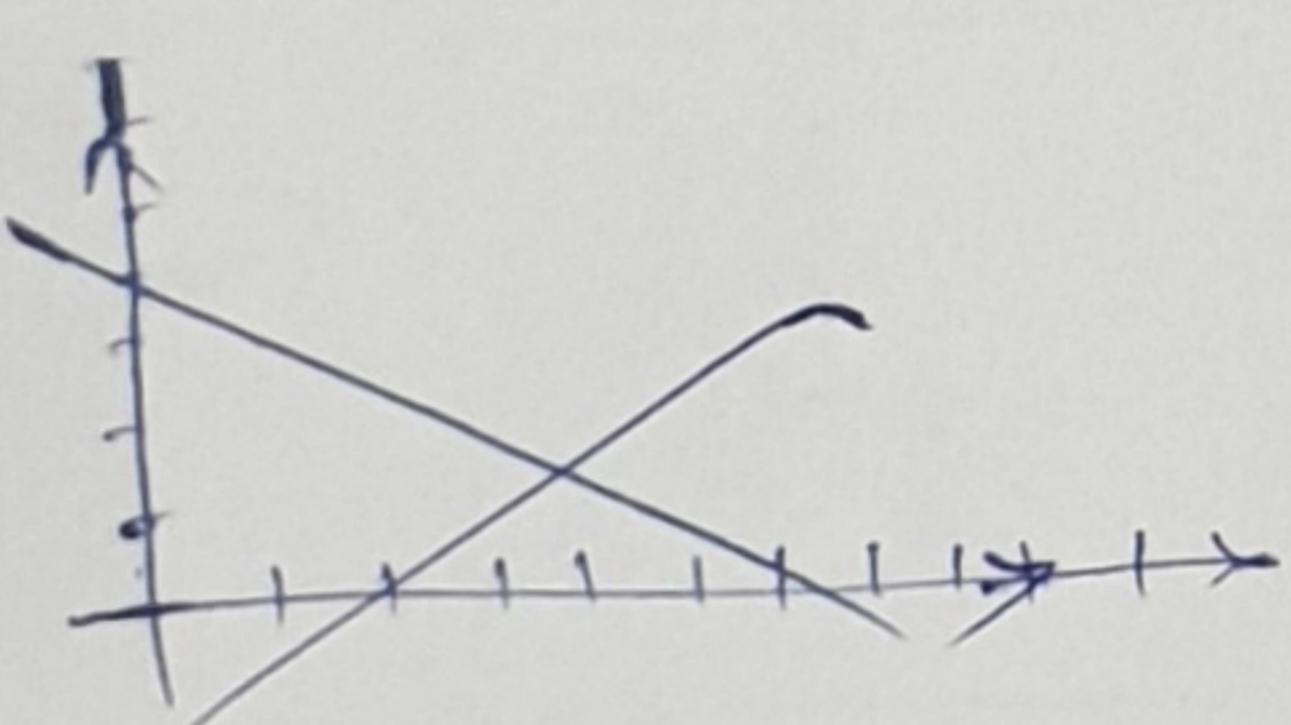
We have 5 samples, 1 feature.

Regression model will be defined as

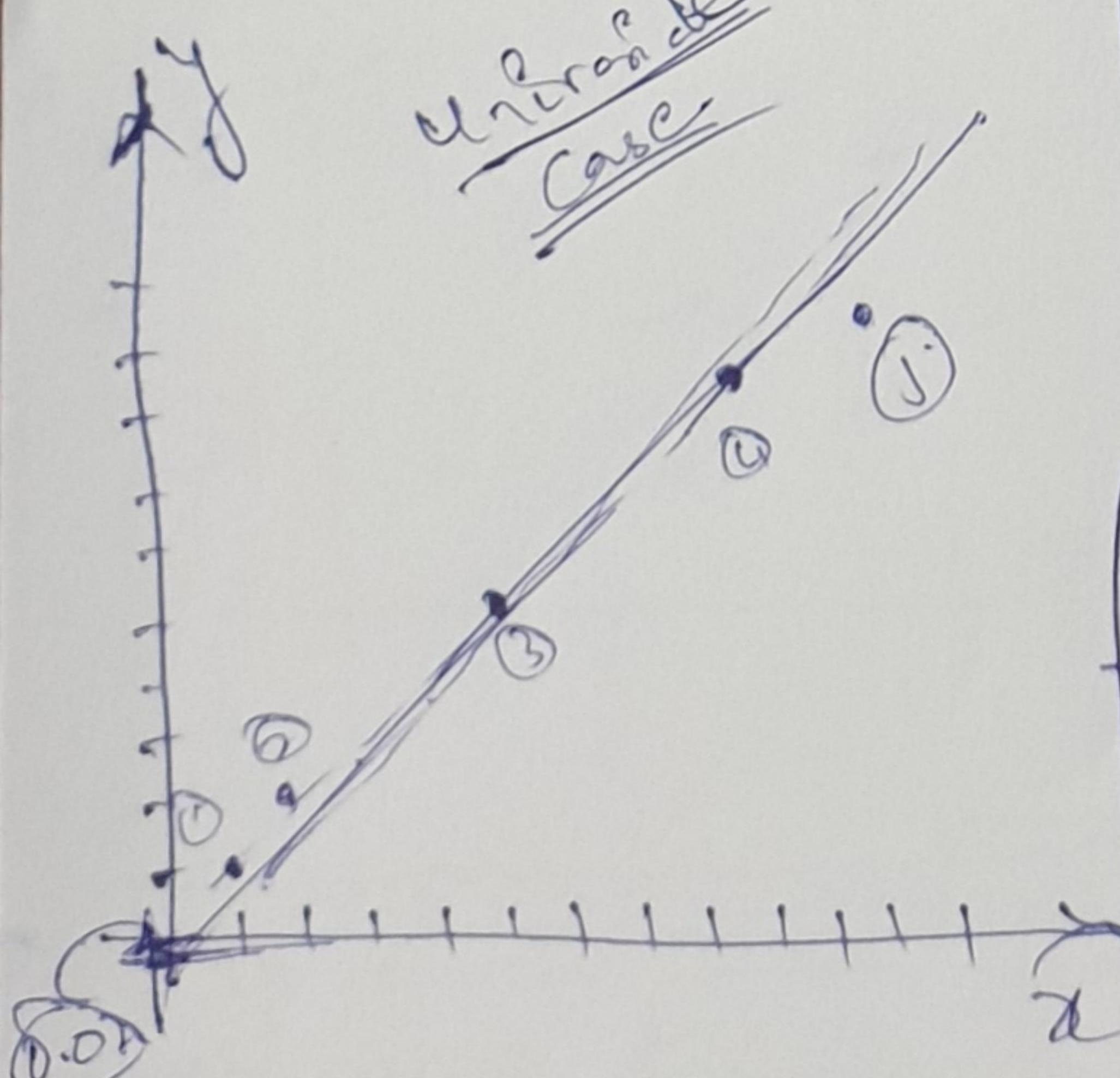
$y = \theta_0 + \theta_1 x$ where θ_0 & θ_1 are
 the model forms representing the
 y -intercept & slope respectively.

Since this is linear regression, model is the
 equation of a line.

x	y
1. 100	1000
2. 200	2000
3. 300	3000
4. 400	4000
5. 500	5000



~~univariate~~
 case



To find the params. normalize the
 slope & intercept to 1000 scale

$$\theta_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

x	y
1	1
2	2
5	5
8	7.8
9	9.2

$$\bar{x} = \frac{1+2+5+8+9}{5} = \frac{25}{5} = 5$$

$$\bar{y} = \frac{1+2+5+7.8+9.2}{5} = \frac{28}{5} = 5.6$$

$$\Rightarrow \theta_1 = \frac{(1-5)(1-5) + (2-5)(2-5) + (5-5)(5-5) + (8-5)(7.8-5) + (9-5)(9.2-5)}{(1-5)^2 + (2-5)^2 + (5-5)^2 + (8-5)^2 + (9-5)^2}$$

$$\Rightarrow \theta_1 = \frac{16 + 9 + 0 + (3 \times 2.8) + (4 \times 4.2)}{50} = \frac{50.2}{50} \approx 1.004$$

Name: Myron Wollenman, ID: 2019AHML618 Course: PGAM 2e2H

Regression EC-R3C1

Q.05. (Contd.)

$$\theta_0 = \bar{y} - \theta_1 \bar{x} = 8 - 8 \times 1.004 = 0.02 \Rightarrow \boxed{y = 0.02 + 1.004x}$$

$\Rightarrow y$ increases 1.004 times with every unit increase in x or area.

2) Multiple regression: Let's say house price depends also on no. of rooms, width, height & many other features. This becomes multiple regression. Here let's call

no. of rooms = x_2 Define $x_0 = 1$ for mathematical convenience

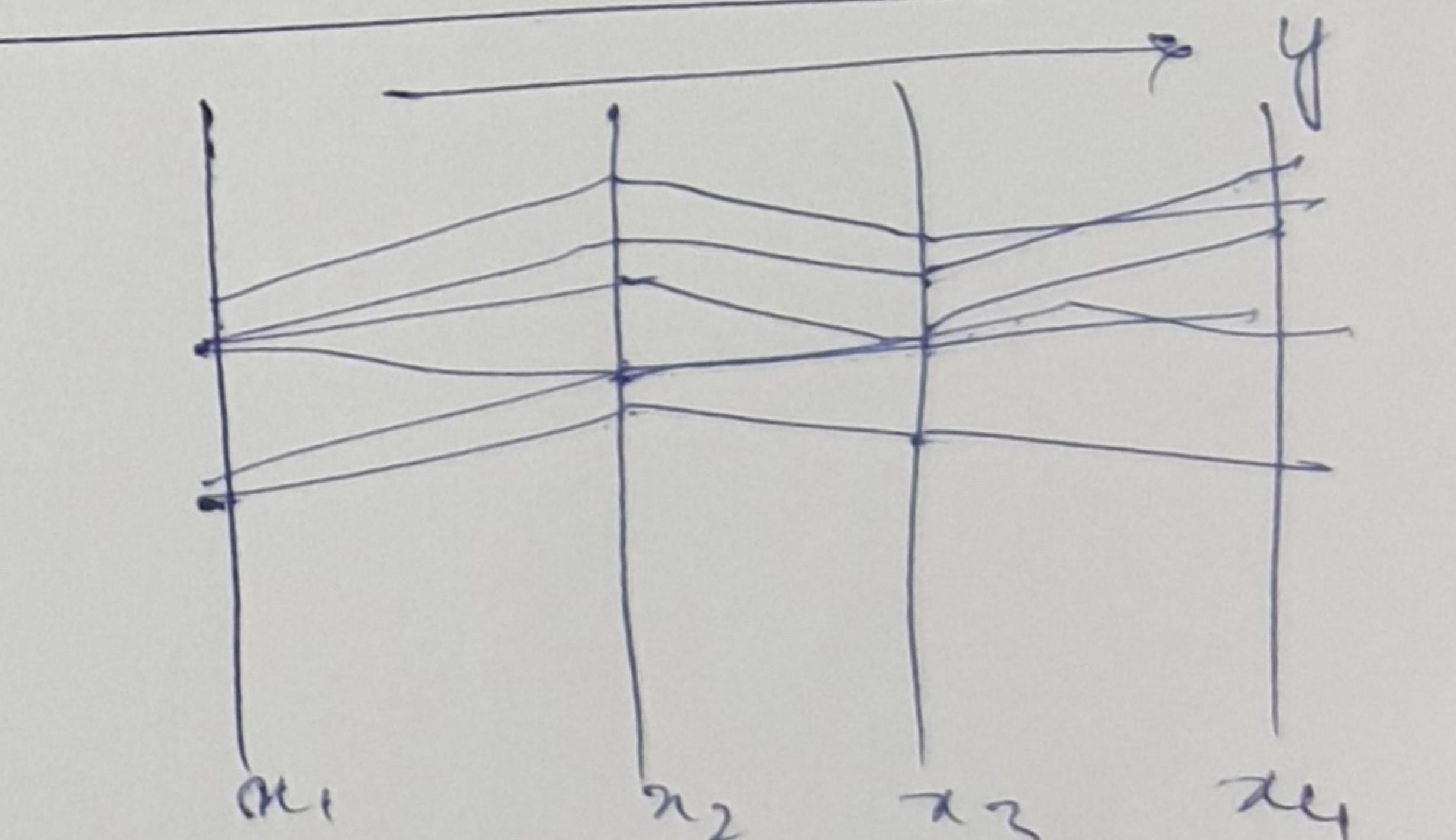
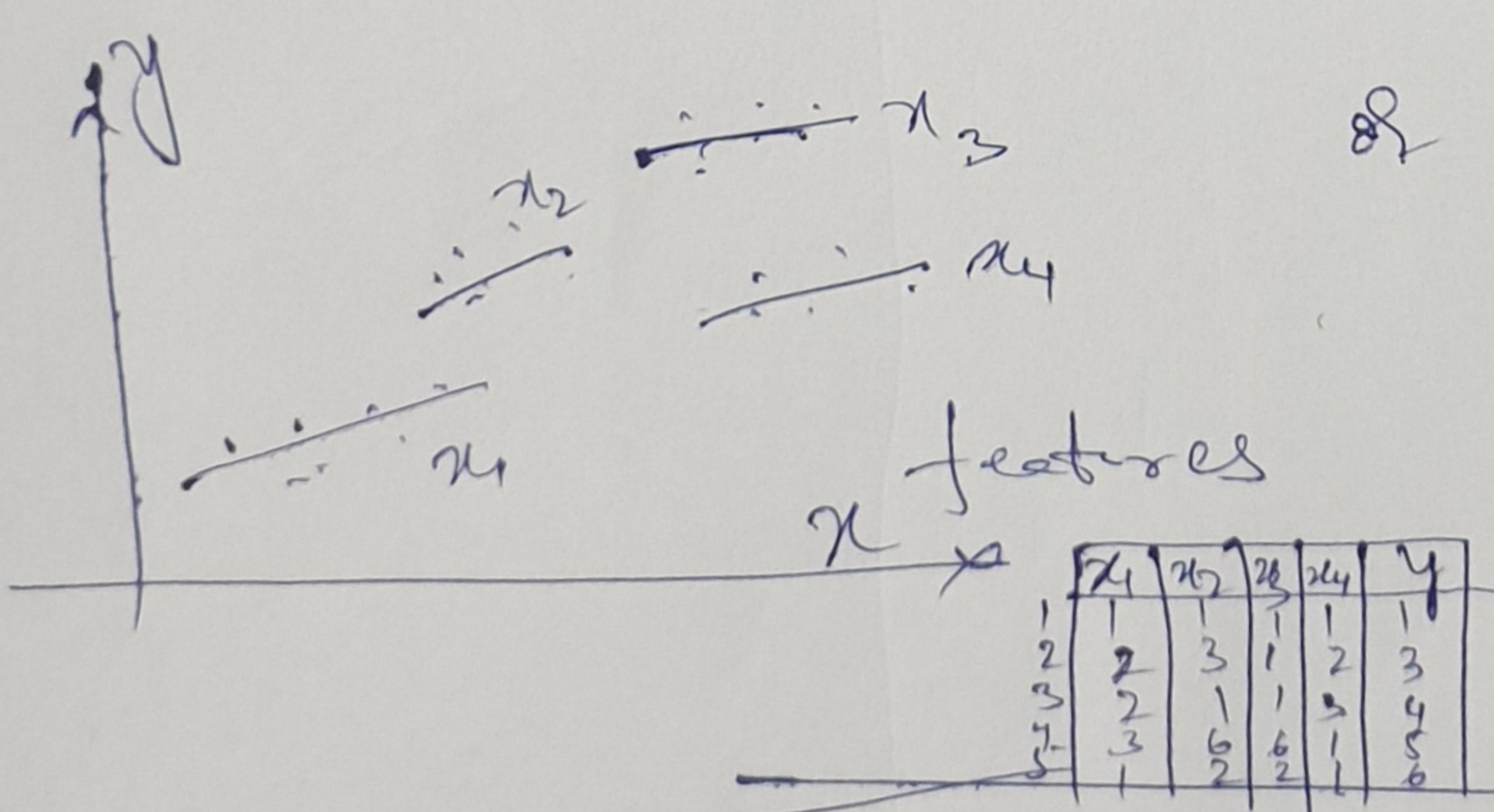
width = x_3 $\Rightarrow y = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4 \dots$

height = x_4 $\Rightarrow \vec{y} = \vec{\theta} \cdot \vec{x}$ or $\vec{y} = \vec{\theta}^T \vec{x}$ where.

area = x_1

$$y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix} \quad m = \text{no. of samples.} \quad x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix}_{(n+1) \times 1} \quad X = \begin{bmatrix} x_0^{(1)T} \\ \vdots \\ x_0^{(m)T} \end{bmatrix}_{(m+1) \times (n+1)}$$

$$\Rightarrow \vec{y}_{(m \times 1)} = \vec{\theta}^T \vec{x}, \text{ where } \vec{\theta} = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix}_{(n+1) \times 1} = \sum_{q=1}^n \theta_q x_q$$



why parallel coordinates.

It can be shown that the solution for $\vec{\theta}$ is $\boxed{\vec{\theta} = (\vec{x}^T \vec{x})^{-1} \vec{x}^T \vec{y}}$

$$\nabla_{\vec{\theta}} J(\vec{\theta}) = \frac{1}{2} (\vec{x}\vec{\theta} - \vec{y})^T \vec{W} (\vec{x}\vec{\theta} - \vec{y}) \text{ where } \vec{W} = \frac{1}{2} \sum_{i=1}^m (x^{(i)} - \bar{x})^T (x^{(i)} - \bar{x}) \text{ or } \boxed{\vec{\theta} = (\vec{x}^T \vec{x})^{-1} \vec{x}^T \vec{y}}$$

(Let deriving due to time constraint)

$$\Rightarrow \nabla_{\vec{\theta}} J(\vec{\theta}) = \frac{1}{2} (\vec{x}\vec{\theta} - \vec{y})^T (\vec{x}\vec{\theta} - \vec{y}) = \frac{1}{2} [2\vec{x}^T \vec{\theta} - 2\vec{x}^T \vec{y}] = 0 \Rightarrow \boxed{\vec{\theta} = (\vec{x}^T \vec{x})^{-1} \vec{x}^T \vec{y}}$$