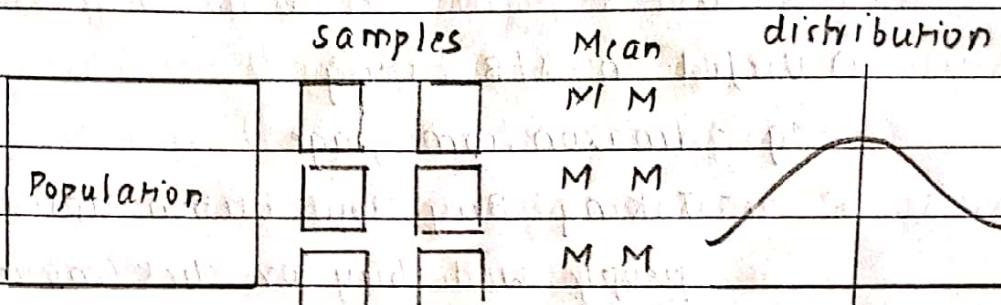


1. Central Limit Theorem



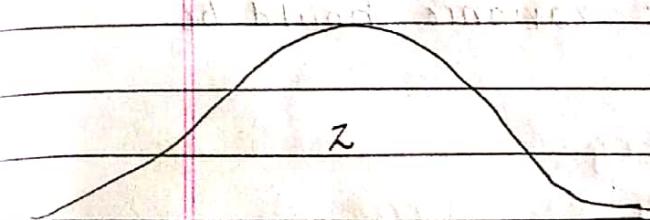
- Mean(sample) \sim mean(population)
- Sample SD $=$ SD of sampling distribution

$$= \frac{\sigma}{\sqrt{n}}$$

$$- SD = \frac{\sigma}{\sqrt{n}}$$

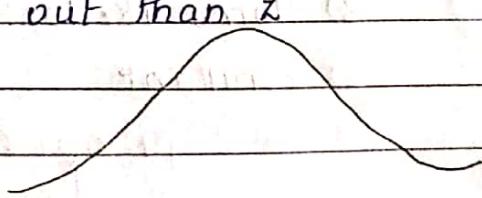
- If n is very large it causes uncertainty or huge spread of data from population

2. Z-test



T-test

- Slightly more spread out than Z



- Use when we don't know about SD of population

3. Hypothesis testing.

1) Useful in A/B testing

e.g.: - Amazon.com page

- Taking Buy Now button upside so people will buy or click more

0.2 0.20001

Apply hypothesis

2) At what point sample size is large enough to have correct test.

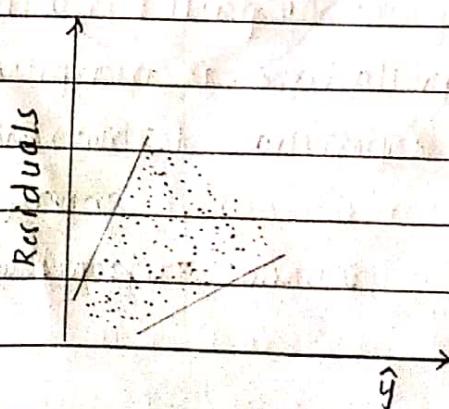
4. Linear Regression

ASSUMPTIONS:

1) No auto-correlation: values should not be dependent on past values

2) No SLR on TSA because TSA has auto-correlation

3) Homoskedastic: variance should be constant



4) Residuals: difference between actual and predicted values, should be normally distributed.

5) Hypothesis function:

$$h(x) = \theta^T x$$

Minimize error (actual - predicted value)

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m [y_i - \theta^T x_i]^2$$

$$6) \sqrt{(x)^2} = \pm 5, (\sqrt{x})^2 = x \quad \forall x \in \mathbb{R}$$

$$\sqrt{25} = \pm 5 \quad \sqrt{5} \times \sqrt{5}$$

$$= -1 \times 5$$

$$= -5$$

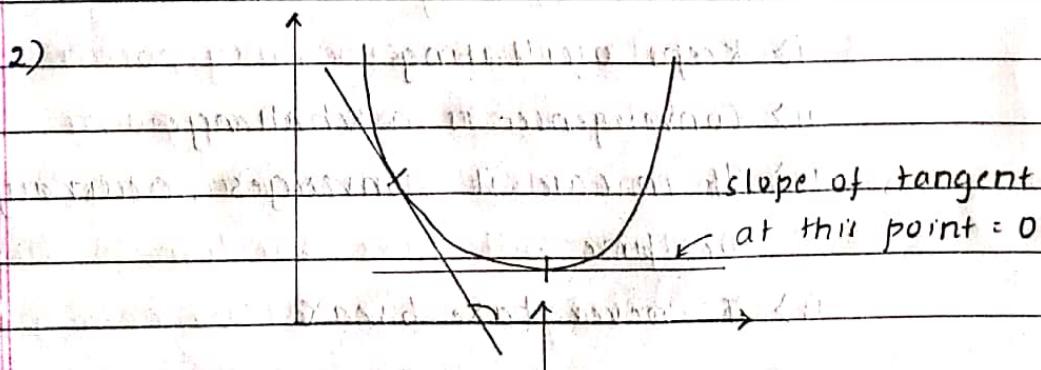
5. Gradient Descent: optimization algo used to find parameters (coefficients) of function

1) $\theta_i = \theta_i - \alpha \frac{\partial J(\theta)}{\partial \theta_i}$ that minimizes cost function.

Learning Rate

Gradient / slope

(Partial derivative)



Point of minimum (arrow)

operation of said algo implemented in

3) Determine slope of curve?

1) Take a point, draw a tangent

2) Slope of line = $\frac{y_2 - y_1}{x_2 - x_1}$

3) Since the line is constant, so taking 2 points is of no use

4) Slope of tangent = slope of curve at that point

5) Slope = $\frac{\partial J(\theta)}{\partial \theta}$ ← direction

6) Slope = $\frac{\partial J(\theta)}{\partial \theta}$ ← magnitude

7) α : control of jump of choosing points

6. a and tolerance

1) LdHT:

i) No high parameter

ii) Accuracy is low

iii) T stops too soon

2) HdLT:

i) Keeps oscillating

ii) Convergence is a challenge

iii) If in case it converges, accuracy will be fine

iv) T never take break

3) IdLT:

i) Extremely long time to converge

4) $H \approx HT$:

- i) Very likely that it will converge and accuracy is lower.

7. Scaling

8. Why to do scaling in LR?

→ i) Impact speed and learning phase.

- std ⁿ	\rightarrow 1) $H \alpha$
- Norm ⁿ	\rightarrow 2) $H \alpha$
- Min-Max	\rightarrow 3) $H \alpha$
- Unit vector	\rightarrow 4) $H \alpha$

- 3) We cannot have same α for both, but we can use only 1 of them.
- 4) Both will have different α s

8. Overfitting

- 1) Won't predict well on unseen data given model has performed well on training

- 2) What stops model from overfitting?

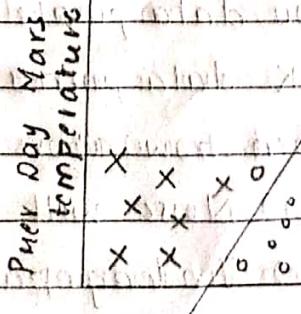
→ i) Predict rain in Mumbai

ii) We have less training data and a variable that

has nothing

to do with

our model



iii) Another thing that allows to make our model overfit is number of features.

iv) Degree of polynomial

overfitting: when a model learns detail & noise in training data to the extent that it's very imprecise

Date _____
Page 6

9. NOISE

1. Data that is unexplainable by given variable

2. What is the source of noise?

→ i) Eg: Data IP / Tracking error

- having cigarettes per day but lying about it

ii) Features of cigarette smokers having

lung cancer but 1 person has cancer

not by smoking because he lives in

an area with an asbestos factory nearby

3. Objective: Minimize error

- But causes above problem (i.e. noise in

data) because we have trained our

model to consider everything so it

will take each and every variable

10. Complexity of Model

1. The amount of data the model needs to store to save itself

2. Increase the size of training data to avoid overfitting

3. 1000 data points, 200 noisy points

10K data points, 2K noisy points

Then how above statement holds true?

→ i) Noise will take random to the model

ii) Finding pattern in 200 noise point

is easier than 2K noise points

4. Overfit: going away from signal

5. Reduce complexity

- Reduce features
- Reduce degree
- Reduce hidden nodes
- Reduce tree depth

(i) Whatever helps you to reduce model complexity will reduce noise

6. Increase training data

7. Update objective function (do that in regularization; Don't just minimize error or minimum error + something)

- * underfitting
 - Increase model complexity

II. Regularization

1) Lasso cost function:

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m [h(\theta_i) - y_i]^2 + \lambda (\sum |\theta_j|)$$

small difference increases value

2) Ridge cost function:

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m [h(\theta_i) - y_i]^2 + \lambda \sum \theta_j^2$$

Underfitting High bias

- increase model complexity

Chaitanya
Page No. 8

3) sparse: fairly spread out

4)	sparse data LASSO lot of features	RIDGE
----	--------------------------------------	-------

1. No. of observations < no. of features

1. Don't have any info. about features

2. Feature selection

2. Asymptotically to 0

3. Removes all of multi-collinear variables

3. If multi-collinearity is encountered use

(penalized variables)

Ridge, it will keep it.

4. Sparse data as it contains many zero features.

4. Ridge also does zero but rarely

5. Speed is much slower compared to Ridge

5. Speed is relatively higher

5) Lasso removes features so why it is not faster? OR Why Ridge is faster than Lasso?

→ i) Ridge is a continuous function so we can use differential function

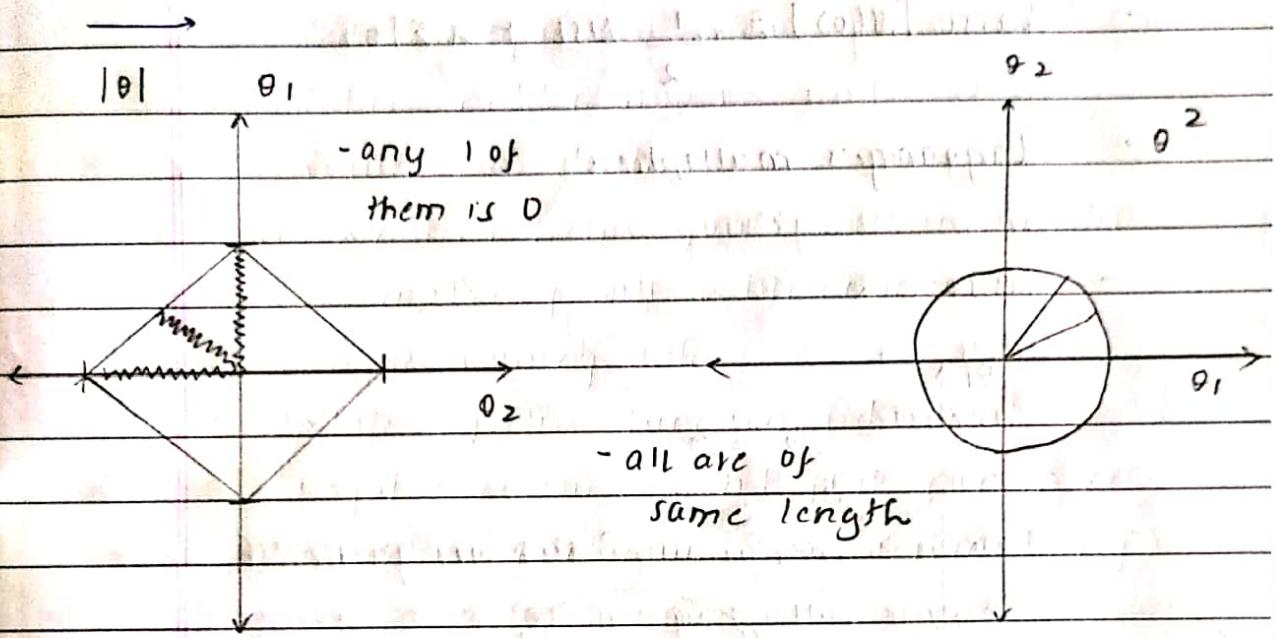
ii) Lasso cannot, we have to use conjugate differential equation

6) Kernel Function: Increases dimension;

spans the data it's easy to separate them

High Variance:
 - more data
 - Reduce model complexity

Q How does lasso removes features while ridge doesn't?



→ Elastic Net:

i) 1 extra parameter than ridge to tune (to bridge lasso and ridge)

ii) Sparse

$$\text{iii) } J(\theta) = \frac{1}{m} \sum_{i=1}^m [h(x_i) - y]^2 + \lambda_1 \sum |\theta| + \lambda_2 \sum \theta^2$$

→ Ridge:

i) less tuning

ii) Speed

→ overfit: difficult to separate minimum

Underfitting: model fits to weak

to fit our dataset

- use dataset or features or model complex

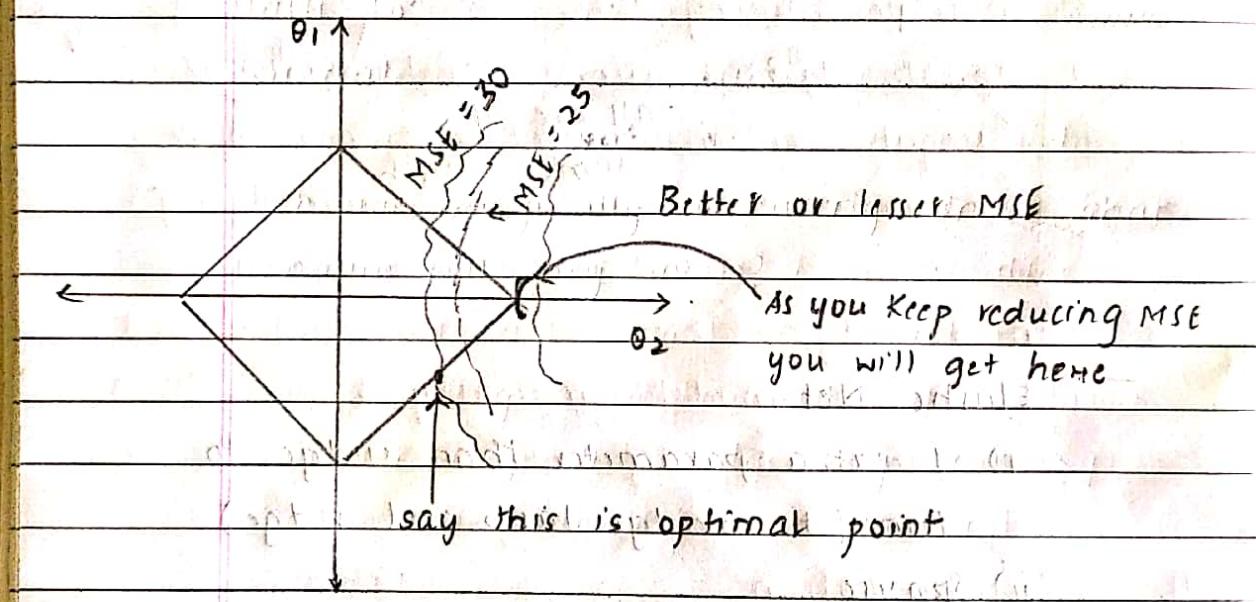
Date 08-12-2019
Page 10

$$\text{Ridge } J(\theta) = \frac{1}{2} \text{MSE} + \lambda \sum \theta^2$$

$$\text{Lasso } J(\theta) = \frac{1}{2} \text{MSE} + \lambda |\theta|$$

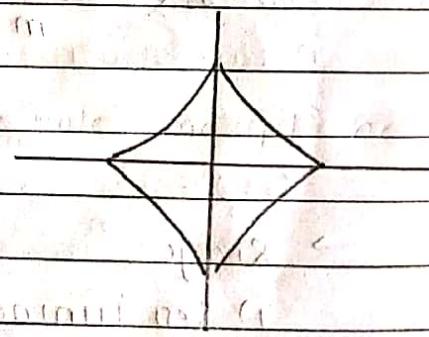
Lagrange's multiplier

- Lasso E



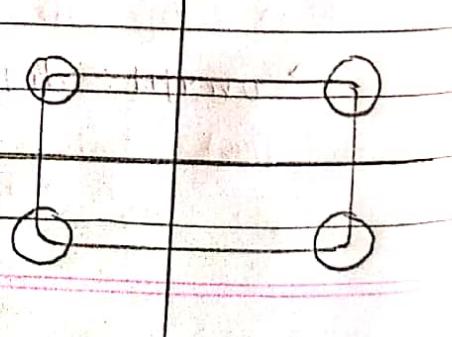
$$L_{0.5} = \lambda \sum |\theta|$$

$$= |\sqrt{\theta_1}| + |\sqrt{\theta_2}|$$



$$L_H = \lambda \sum \theta^2$$

Minimum MSE at
circled points



← Model is not complex enough to accurately capture relationships between features & target variable

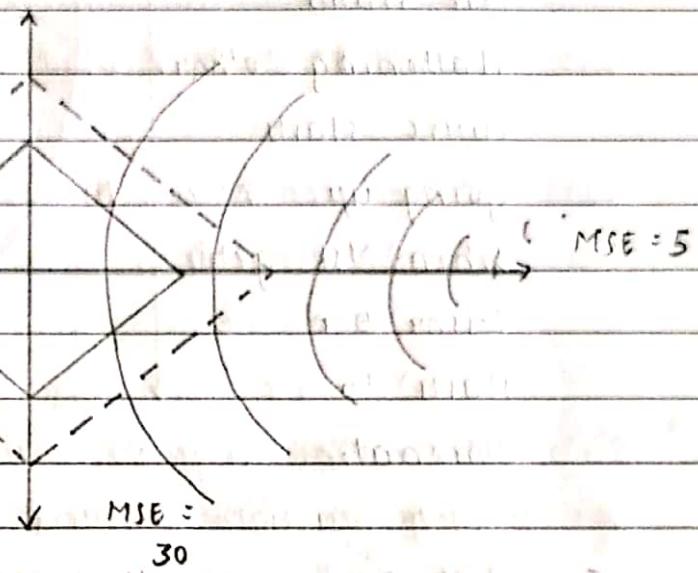
Class page N.....

- Lasso keeps them within check

Diamond

$$|\theta_1| + |\theta_2| = c$$

$$= C_2$$



- Regularization doesn't mean reducing θ at every iteration
- It means finding best way to have minimum error while keeping the model simple
- $\lambda(|\theta_1| + |\theta_2| = 20)$, $\lambda(|\theta_1| + |\theta_2| = 10)$

$$100 \quad 0.01$$

$$\text{Ridge} = \theta_1^2 + \theta_2^2$$

$$= 10000 \quad 0.01$$

$$\text{Lasso} = |\theta_1| + |\theta_2|$$

$$= 100 \quad 0.1$$

So ridge is more likely to reduce $\theta_1 = 10000$ than Lasso which is more likely to reduce both

$$J(\theta)$$

Regulatized cost function

UnRegulatized cost function

- As reduction in

MSE starts

flattening, other
curve starts

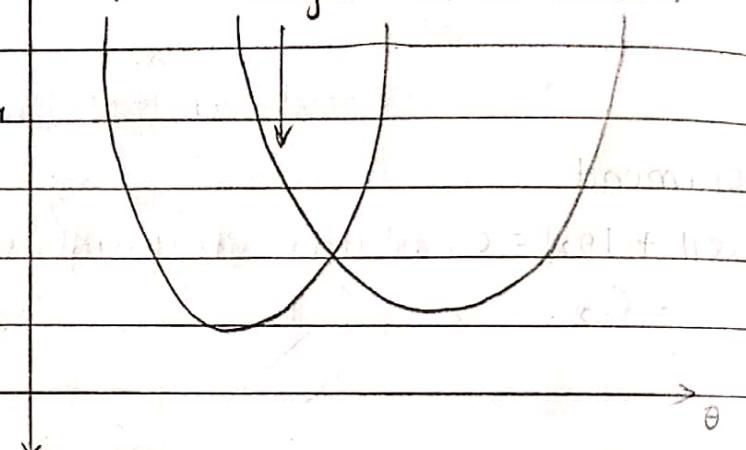
going up

When MSE gets

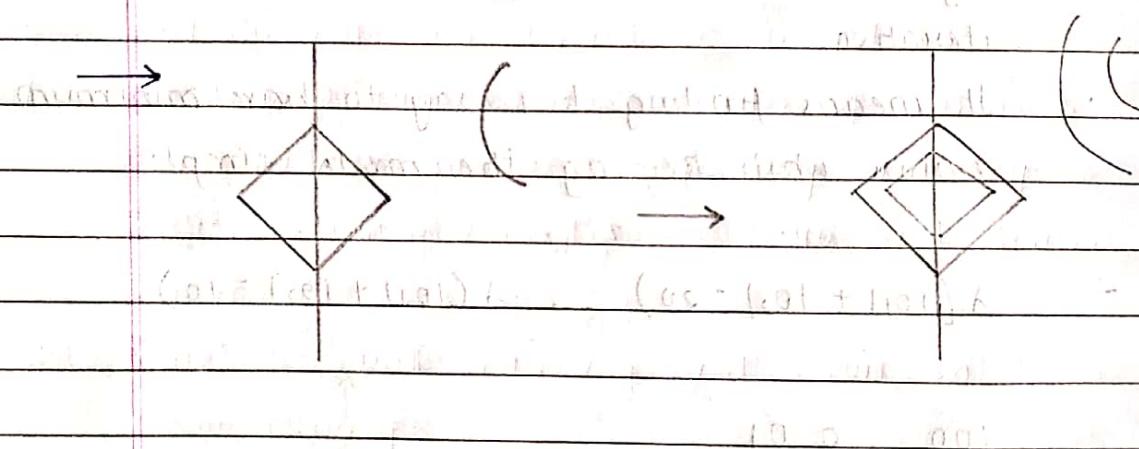
lesser θ

starts to

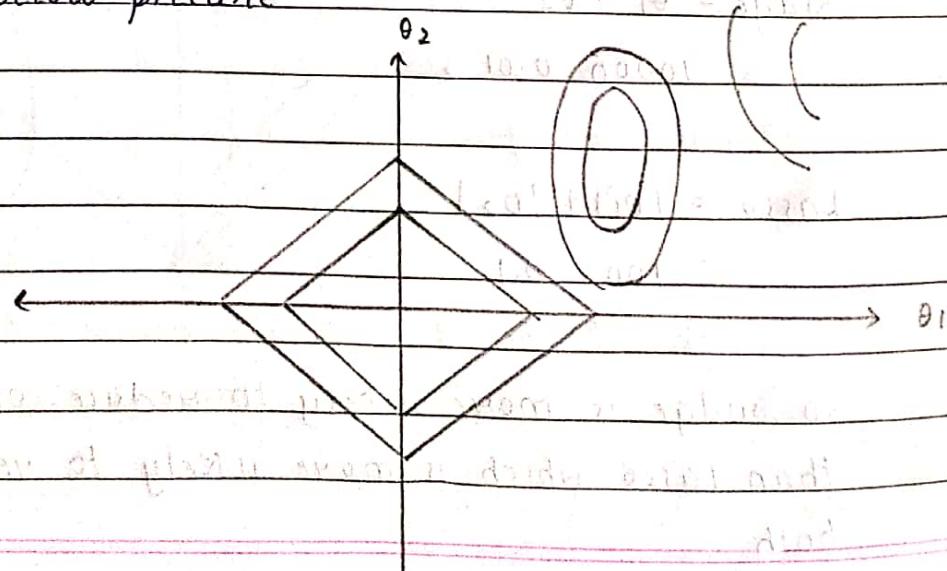
dominate



- Both are marginally increasing. Under what scenario will they meet?



- If θ_1 plays much bigger role, we get the below picture



- Lasso: In most cases we are likely to hit edge
- Ridge:
 - Less likely that θ_0 becomes 0
 - Less likely to hit edge
- Regularization sees if it can reduce 0 without impacting error

Role of λ :

1. How much diamond will expand in context to contour.
2. If λ is small, diamond will expand more, much more compared to contour.
3. If diamond is static, θ_0 are zero.
4. λ controls rate of expansion.
5. If $\lambda = 0$, focus on MSE, no effect from θ .
6. If λ is very high, optimal point will be very close to zero.
7. If λ is very low, minimize MSE.
8. If λ is very large, the error is high, so error contour will expand fast, so diamond won't expand, optimal point will be diamond close to (0,0).
9. Decide λ using cross validation.

Elastic Net

1.
$$\frac{1}{2} \text{MSE} + d_1 \sum |\theta| + d_2 \sum \theta^2$$
2. There are three things to consider:
 - i) diamond
 - ii) contour
 - iii) circle
3. Minimize large θ values
4. When multiple features are correlated to each other, lasso will keep 1 and remove the rest. Ridge will try and keep these variables.
5. Elastic net will depend on d_1 and d_2 .
6. E.N. will try to make cluster of variables; it will either go the ridge way and try to preserve the variables or lasso and remove them allowing the net to fit in given data.
7. Use manual training and cross validation for tuning.

AIC and BIC

1. Akaike Information Criteria and Bayesian Information Criteria both help to compare models.
2. BIC sort of assumes 1 out of say 10 models is absolute true model and assign values accordingly.
3. AIC is not as strict as much as BIC on parameters : BIC is more likely to pick a simple model. AIC might choose complex model.

4.

$$\text{AIC} \rightarrow 2K + \boxed{\ln \text{RSS}} \rightarrow \text{NMSE}$$

$$\text{BIC} \rightarrow K \cdot \ln(n) + \boxed{\ln \text{RSS}}$$

K : number of parameters

n : number of datapoints (observations)

If $n > e^2$, $\text{BIC} > \text{AIC}$ (BIC has or puts higher weightage on K than AIC)

5. Overfitting \rightarrow simpler model (observations increase, overfitting decreases)

Underfitting \rightarrow complex model

6. If sample size is small, it is shown AIC tends to pick better model.

7. $n < K^2$: always go for AIC (recommended)

Logistic Regression:

1) Sigmoid function: $\hat{y} = \frac{1}{1 + e^{-\theta^T x}}$

also known as logit link function or logitit transformation

2) Error: mis-classification

3) Output: 1 if either zero or 1 is true else 0

$\theta^T x = -\infty$, let $\hat{y} = 0$ if $x < 0$

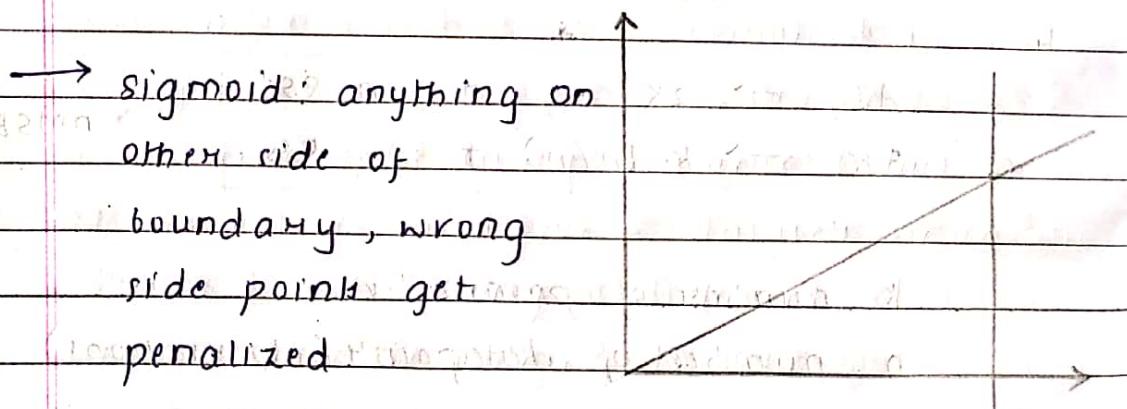
$\theta^T x = \infty$, $\hat{y} = 1$ if $x > 0$

4) Dependent variable is categorical.

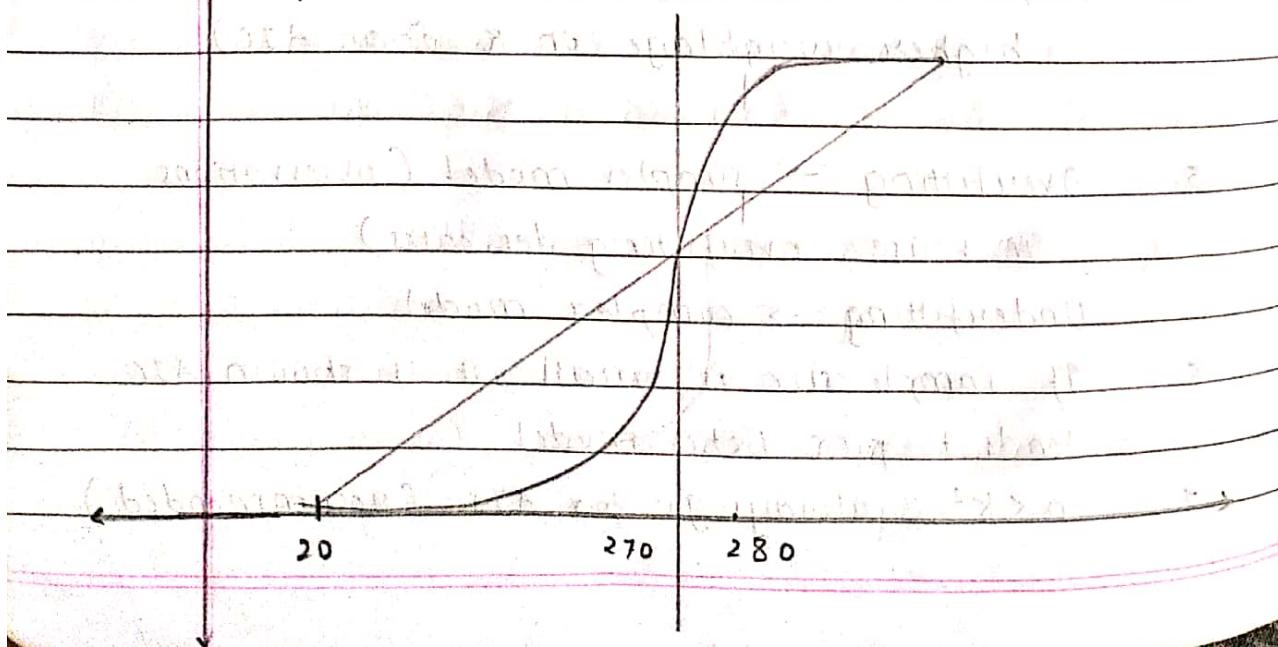
8. Why we use sigmoid in logistic regression?

to obtain symmetric output OR from 0 to 1

8. Why we cannot use below graph?



5) $\hat{y} = 0.002x - 0.04$

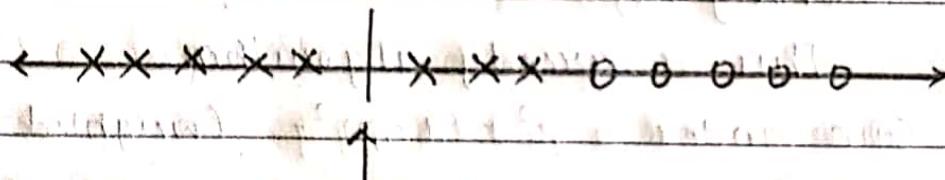


- let's say $x = 280$ (predicted value)

- As per linear R. : $0.52 / 0.66$

sigmoid : $0.66 / 0.85$

- 6) Sigmoid induces bigger errors for mis-classification and allows it to learn much faster.
- 7) Sigmoid gives much errors for points close to boundary.
- 8) Sigmoid pushes values either to 0 or 1, so very few in between values.



- and sigmoid will move this line much faster than linear because of high error
- 9) There is no arbitrary change in function for continuous function

- not differentiable

- cannot use GD and
other

Decision Trees

1. It keeps splitting till all datapoints belong to same class.

2. Goal is to make homogenous groups.

3. How to split?

→ We consider bunch of matrices

→ Try to optimize at every step with the assumption that it is optimized at every step.

This is greedy algorithm.

4. Gini index: $p^2 + (1-p)^2$ (weighted gini,

Entropy: $-p \cdot \log p$ ratio, etc.)

5. You decide split using variance reduction.

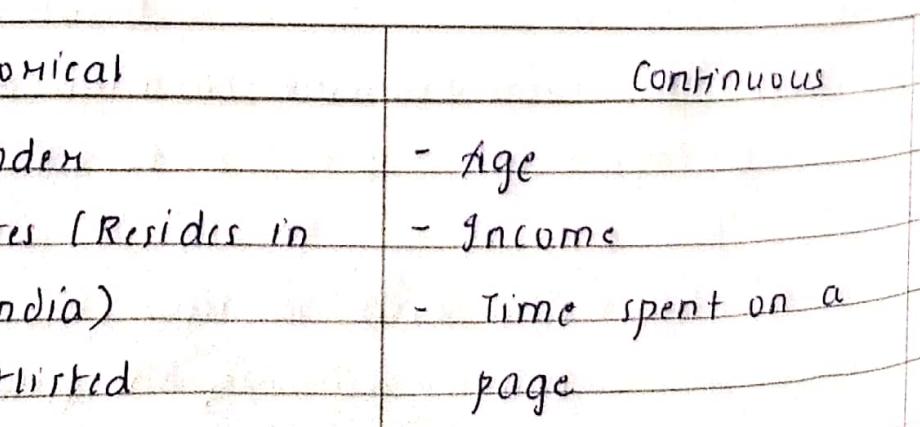
6. DT cannot be used for 2 datasets which are identical but belong to different classes because we have nothing to split.

So there is a hit with accuracy.

7. DT almost always overfits.

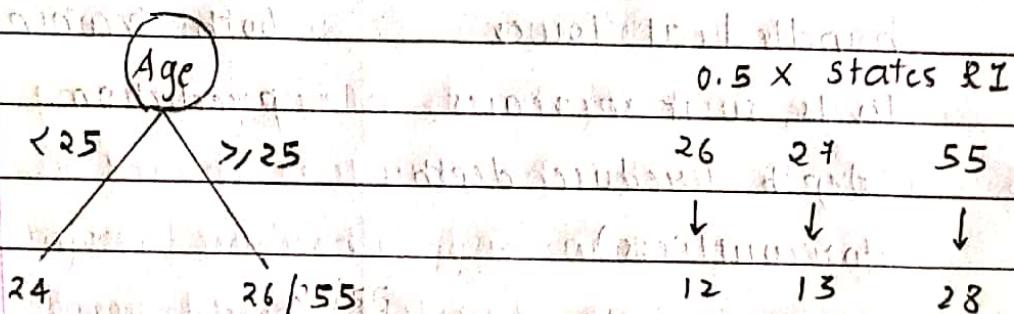
8. Gain is minimal so we cannot split.

9. Overall depth of the tree.



	Categorical	Continuous
10.	<ul style="list-style-type: none"> - Gender - states (Resides in India) - shortlisted 	<ul style="list-style-type: none"> - age - income - Time spent on a page

D.T. vs Logistic regression



Here, we are not able to differentiate between adult or middle aged person because of such a split.

	Decision Trees	Logistic
1)	Works better with categorical data	Works better with continuous data
2)	Can build almost any type of non-linear boundary / fit very complex function	less likely to overfit.
3)	Robust to outliers	Not robust
4)	OK with missing values (doesn't consider values for that feature but for others)	problems
5)	Better with class skew (Class imbalance)	Cannot handle much class skew
6)	Matters in ability to minimize error	and fine

7) Class skew will be handled at lower levels, so it increases depth (reduce depth for outliers.)

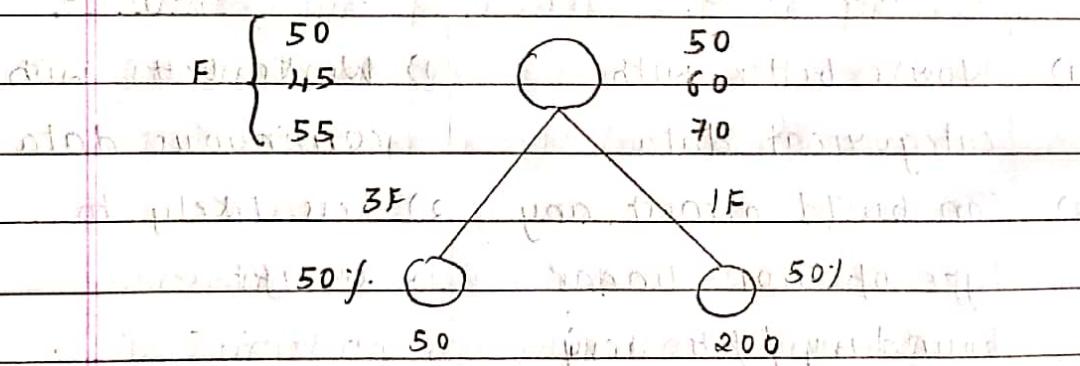
7) Faster than DT at both training and prediction.

12. How we use DT when target variable is continuous?

OR

How will you decide split?

→ 1) Gini index or entropy cannot be used.



2) So we use Reduction in variance for continuous target variable

- checks by how much variance has been reduced

- As much as variance has been reduced better it is

13. Suppose we want to check whether someone will buy the product or not.

→ 1) Check time spent on a page

2) Now the question arises whether where will we put DT boundary?

dtree	Gini index (default)	Information Entropy
-	Faster because we take square	log takes time

Ensembling

1. For binary classification, weak classifier is

Accuracy

0 0.5 1

2.

C1 is wrong

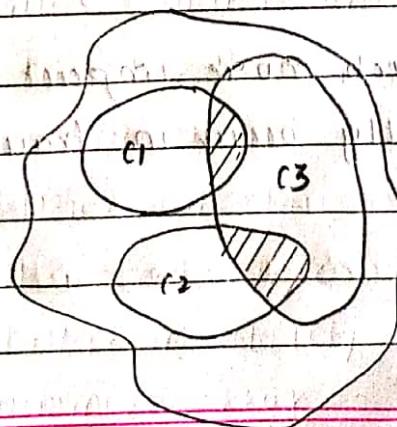
C2 is wrong

(3) is wrong

- This gets better and better

when we have more classifiers

- For overlapping if 2 classifiers vote in-
correct then the classifiers are wrong.



4. It would be correct only in shaded region else is incorrect.

So we are at disadvantage by taking ensemble.

5. For ensemble, the classifiers should be accurate by more than 50%. otherwise ensemble is useless i.e. each classifier should be better than random.

Bagging

1. Bootstrap aggregation

Sampling with replacement

2. How would bagging with DT work?

What type of model will we get if combine BA and DT?

→ Decision Tree

Random Forest

1. Sampling of features (decreases overfitting)

2. Sampling of records / datapoints (less

trained data for DT increases overfitting)

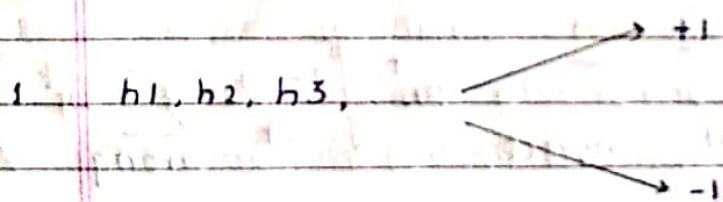
3. Majority vote / AM (reduces overfitting)

* Features have much more impact than data

* Noisy labels mostly occur in fraud detection

Random ForestDecision Trees

- | | |
|-----------------------------|--|
| 1. less likely to overfit | 1. Interpretable |
| 2. Handles noisy labels | 2. Speed (especially prediction speed) |
| 3. Tells feature importance | 3. Less likely to under-fit |

Boosting:

2. Boost misclassified datapoints in next classifier

$$\sum w_i = 1, \quad E = \sum w_i$$

wrong

Boost stronger classifier:

$$h(x) = \text{sign} [d_1 h_1(x) + d_2 h_2(x) + d_3 h_3(x) + \dots]$$

Start with, classifier error at that time

$$w_i = \frac{1}{N} \rightarrow \text{Pick } h_t \text{ that minimizes } E_t$$

ADA BOOST ✓

Pick d_t (weight you want to give to classifier)
 Calculate w_t ←

$$w_i^{t+1} = w_i^t \frac{e^{-\alpha h^t(x_i) y_i}}{z} \quad \text{if } h^t(x_i) y_i < 0 \quad \{ \text{Misclassified} \\ \text{then } -rc \}$$

$$\alpha^t = \frac{1}{2} \ln \frac{1 - Et}{Et}$$

$$\text{If correct, } w_i^{t+1} = w_i^t + e^{\ln \sqrt{\frac{Et}{1-Et}}}.$$

$$\text{If wrong, } w_i^t e^{\ln \sqrt{\frac{1-Et}{Et}}}$$

$$\frac{1}{Z} \left(\sqrt{\frac{Et}{1-Et}} \sum_{\text{correct}} w_i^t \sqrt{\frac{1-Et}{Et}} \sum_{\text{wrong}} w_i^t \right) = 1$$

$$Z = \sqrt{Et(1-Et)}$$

$$\sum_{\text{correct}} w_i^t = \frac{1}{Z} \sqrt{\frac{Et}{1-Et}}$$

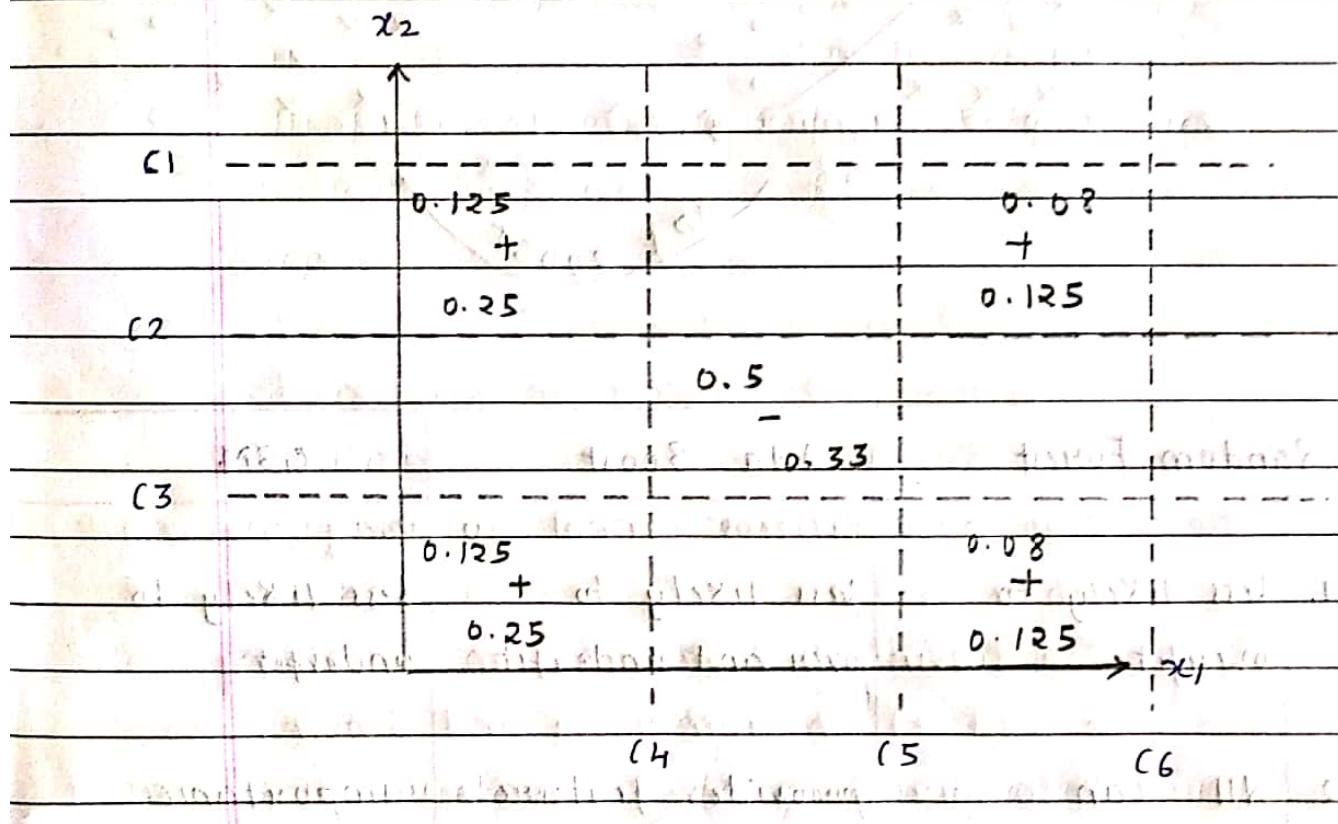
$$w_i^{t+1} = w_i^t \frac{Et}{1-Et} = w_i^t \cdot \frac{1}{2} \frac{1}{1-Et}$$

$$\sum_{\text{correct}} w_i^{t+1} = \frac{1}{2}$$

All correct and incorrect points get weightage

$$\text{of } \frac{1}{2}$$

If 80% of your points are correct, its collective weight will be $1/2$ and remaining 20% of incorrect points will be $1/2$



Gradient Boosting:

- Learn error in previous step

$$y_i = h(x_i) + \text{error}$$

$$\text{error}^t = h^{t+1}(x_i) + \text{error}^{t+1}$$

$$\text{error}^{t+1} = h^{t+2}(x_i) + \text{error}^{t+2}$$

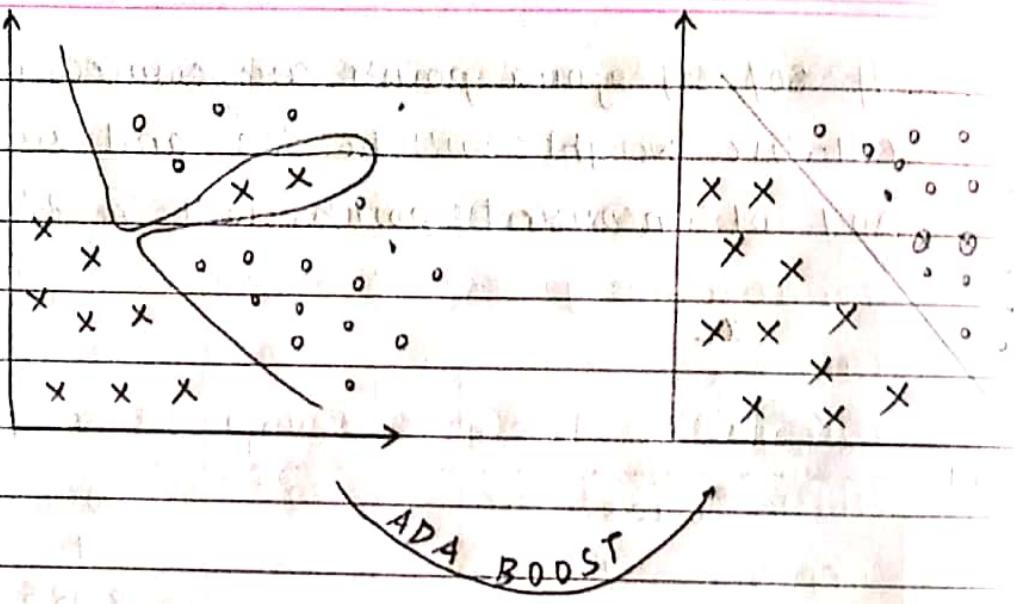
$$y_i = b^t(x_i)$$

$$y_i = b^t(x_i) +$$

$$h^{t+1}(x_i) +$$

$$h^{t+2}(x_i) +$$

Till error becomes very small



Random Forest	Ada Boost	GIBM
1. less likely to overfit	less likely to over and underfit	less likely to underfit
2. All can provide feature importance		
3. Works with noisy labels	Better with class skew	Better with class skew
4. Training speed when parallel	Prediction speed	Prediction speed
5. Longer to predict	Relatively fast	Relatively fast
6. Missing values, outliers	Lot of hyper parameters	Lot of hyper parameters.

XGBoost:

1. GBM + regularization + speed
2. Faster for parallel

21-12-2019

Naive Bayes (Generative Model)

1. $P(c, d)$ read as joint probability of c and d
2. In generative model what you learn is joint probability.
3. $P(c, d) = P(d|c) \cdot p(c) = P(c|d) \cdot p(d)$
4. $P(c|d) = \frac{P(d|c) \cdot p(c)}{P(d)}$

For 1 datapoint $p(d)$ is constant because $P(d)$ is fixed.

5. $\text{Max}[p(c|d)] = \text{Max}[p(d|c) \cdot p(c)]$
6. $P(d|c) = P(x_1, x_2, x_3, x_4, x_5 | c)$

(Joint probability of all features)

(and class)

$$= P(x_1, x_2, x_3, x_4, x_5 | c \cdot x_5)$$

$$= P(x_1, x_2, x_3 | c \cdot x_5 \cdot x_4)$$

7. In NB all these features are independent of each other.

$$= P(x_1|c) \cdot P(x_2|c) \cdot \dots \cdot P(x_5|c)$$

8. NB chooses class with highest probability.

Multinomial NB Classifier:

$$1. \text{ Cmap} = \underset{c \in C}{\operatorname{argmax}} P(x_1, x_2, x_3, \dots, x_n | c) \cdot P(c)$$

$$2. \text{ CNB} = \underset{c \in C}{\operatorname{argmax}} P(c_j) \prod_{x_i \in X} P(x_i | c)$$

Problems with Maximum Likelihood:

1. Zero probabilities cannot go away, no matter other evidence.

2. To solve this we perform smoothing.
- Add +1 in numerator and denominator
 - So our data never have zero. This is called Laplace smoothing.

$$P(w_i | c) = \frac{\text{count}(w_i)}{\text{count}(c) + 1}$$

Other types of smoothing:

(i) Good tuning

(ii) Mason A

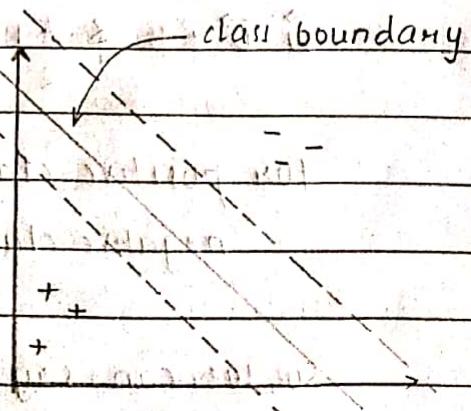
Advantages of NB

1. Very fast.
2. Training is also quick.
3. Low storage requirements.

4. Extremely robust to irrelevant features (IF)
IF cancel out each other without affecting features.
5. Least likely to overfit.
6. Very good in domains with many equally important features. (Decision trees suffers from fragmentation in such cases, when data is little)
7. NB takes probability and multiplies by 2 if 2 features are same.
8. Works really well on exceptionally high dimensional dataset.
9. Works good for text classification.
- Q. Why don't we use NB before logistic, linear...?
→ Discriminative classifiers (less likely to underfit)
works really well than generative model.

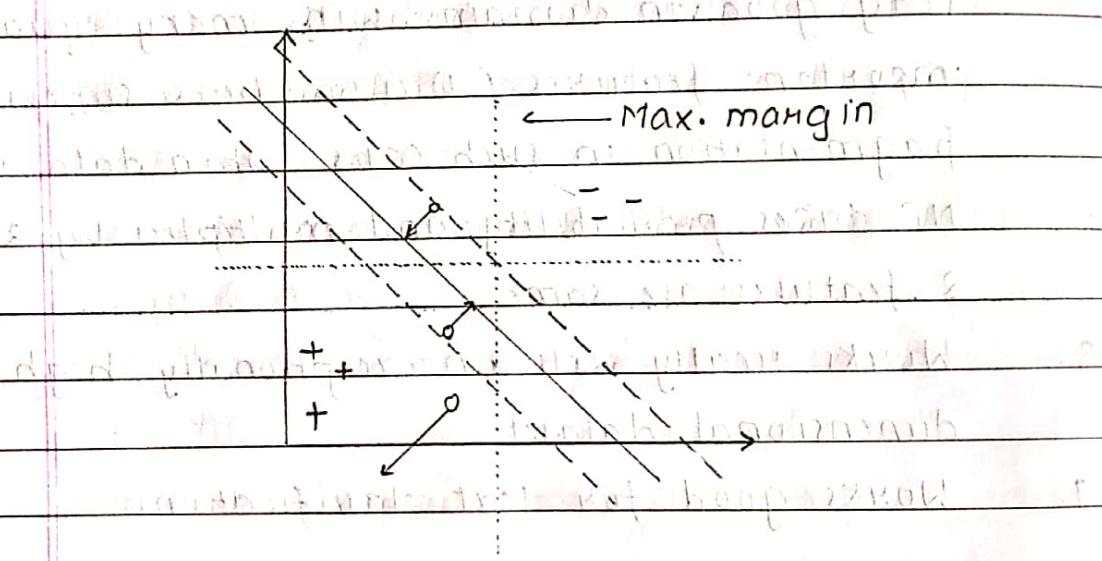
Support Vector Machines (SVMs)

1. Less likely to start overfitting
2. SVM can split these +ve and -ve classes on points in many ways
3. Logic puts boundary in such a way that all points are correctly classified.



4. SVM nuances further: The distance between boundary and closest points are high as possible.

5. Separation is as high as possible.



w : vector that allows to determine distance between margin and boundary of both sets

$$\vec{w} \cdot \vec{v} = c \text{ and } \vec{w} \cdot \vec{v} = c > 0$$

$$\vec{w} \cdot \vec{v} = c \quad \vec{w} \cdot \vec{v} + b > 0$$

$\vec{w} \cdot \vec{v} + b > 0$ for class 1, yet class 2

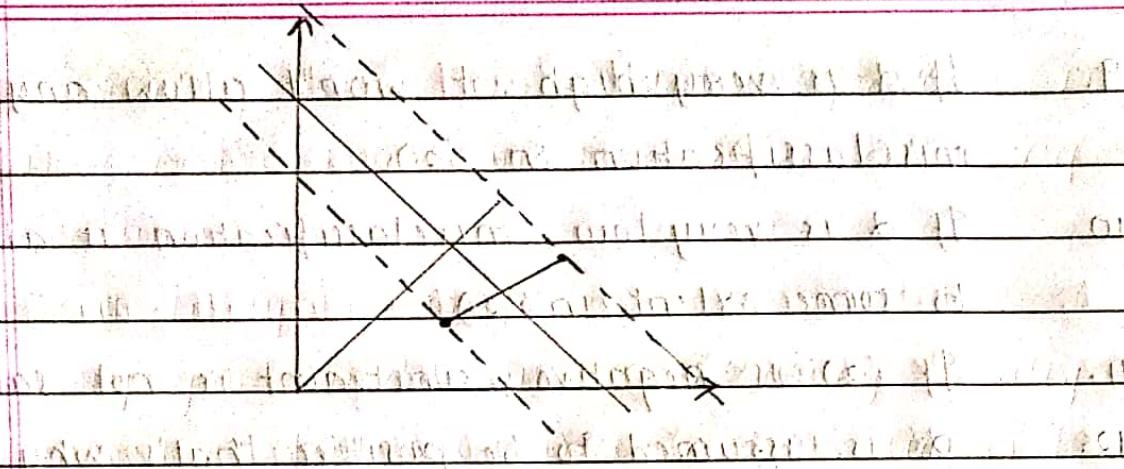
$$\vec{w} \cdot \vec{v} + b \leq -1$$

For positive class, $y = 1$

negative class, $y = -1$

so, we can say that, $y(\vec{w} \cdot \vec{v} + b) > +1$

$$y(\vec{w} \cdot \vec{v} + b) - 1 > 0$$



6. Objective:

i) Maximize width as much as possible

ii) Missclassified points

$$7. \vec{x}_p = (1-b) \frac{\vec{w}}{\|\vec{w}\|}, \vec{x}_n = (-1-b) \frac{\vec{w}}{\|\vec{w}\|}$$

$$\vec{x}_p - \vec{x}_n = \frac{\vec{w}}{\|\vec{w}\|} = \text{And this is what needs to be maximized.}$$

$$8. \text{Minimize } \frac{1}{2} \|\vec{w}\|^2 \text{ with } y_i (\vec{w}^\top \vec{v} + b) - 1 = 0$$

as constraint.

8. Lagrange's multiplier: Using come multiplier we will minimize entire term.

$$L = \underbrace{\frac{1}{2} \|\vec{w}\|^2}_{1} - \alpha \underbrace{[\sum y_i (\vec{w}^\top \vec{v} + b) - 1]}_{2}$$

α: how much importance to (2) w.r.t. (1)

9. If α is very high, it won't allow any misclassification in (2)
10. If α is very low, misclassification is allowed to some extent in (2)
11. If (2) is negative, constraint is not satisfied.
12. α is assumed to be positive that's why we take -ve sign.
13. We will have minimal value when derivative becomes zero.

$$\frac{\partial L}{\partial w} = \vec{w} - \sum \alpha_i y_i \vec{w} = 0$$

$$\frac{\partial L}{\partial b} = 1 - \sum \alpha_i y_i = 0$$

$$w = \sum \alpha_i y_i \vec{v}$$

$$\frac{\partial L}{\partial b} = \sum \alpha_i y_i = 0$$

$$0 = 1 - \sum \alpha_i y_i \vec{v} - \sum \alpha_j y_j \vec{v} - \sum \alpha_i y_i \cdot \sum \alpha_j y_j \vec{v}$$

$$\text{transferring } \sum \alpha_i y_i; b + \sum \alpha_i$$

$$\text{minimization} = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \vec{v}_i \vec{v}_j$$

In optimal scenario

it reduces to this

- Allows for special scenario:
- Class boundary is linear boundary
- To build Kernel function transforms to high dimension.

Dot product of every combination,

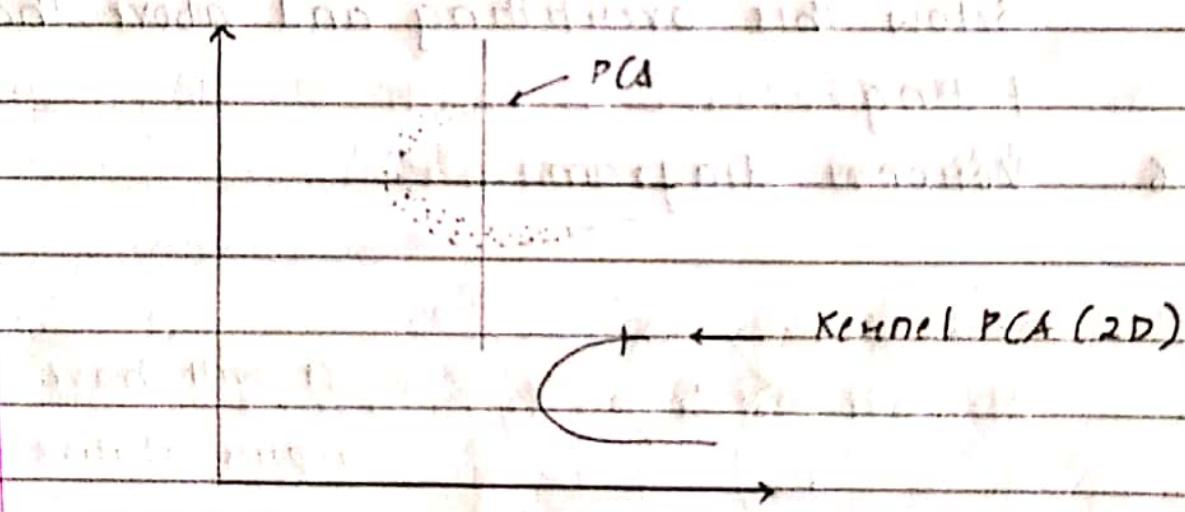
$$K \rightarrow (\vec{u} \cdot \vec{v} + 1)$$

14. SVMs are the only ones which can tolerate higher dimensions using Kernel.

15. Kernels are also used in PCA, but why?

→ PCA builds linear combinations of features

→ But PCA cannot build non-linear combinations, therefore we use Kernels.



Advantages:

- 1) less likely to overfit.
- 2) High dimensionality: works well with text

Disadvantages:

- 1) slow

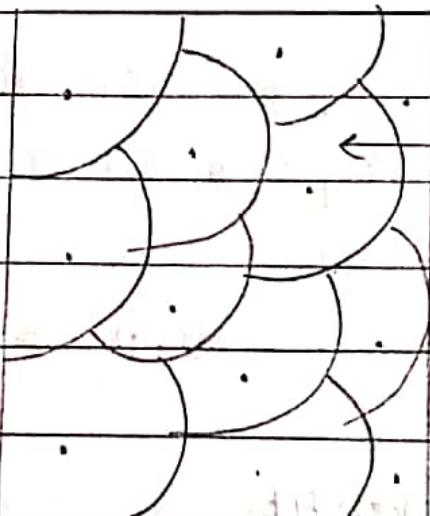
Neighbourhood Based Predictions:

KNN (K Nearest Neighbours)

1. Distances of each point from all points.
2. Low K → Overfitting
3. High K → Underfitting
4. Good with non-linear data
5. Value of $K = 2, 3, 4, 5$ are pretty good classifiers.

Below this overfitting and above that underfitting.

6. Voronoi diagram:



It will have
representative points
in each section.

representative

points of closest