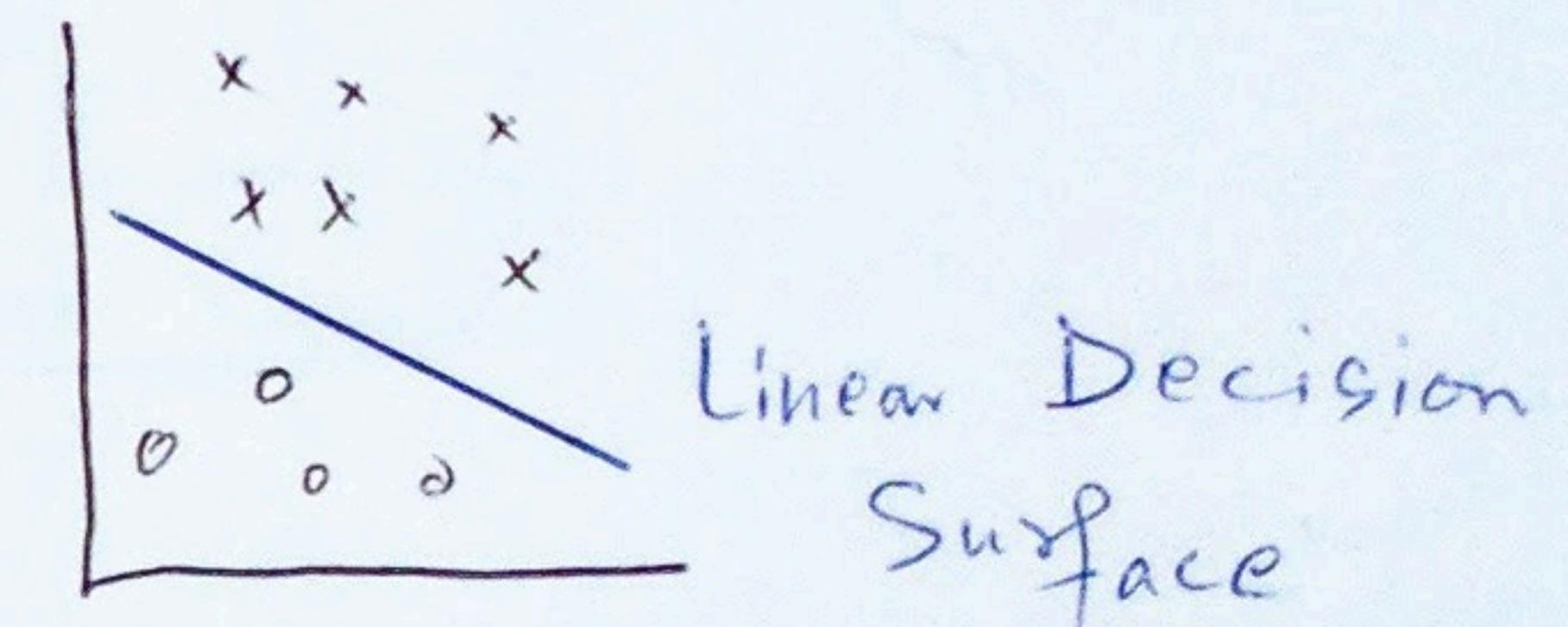
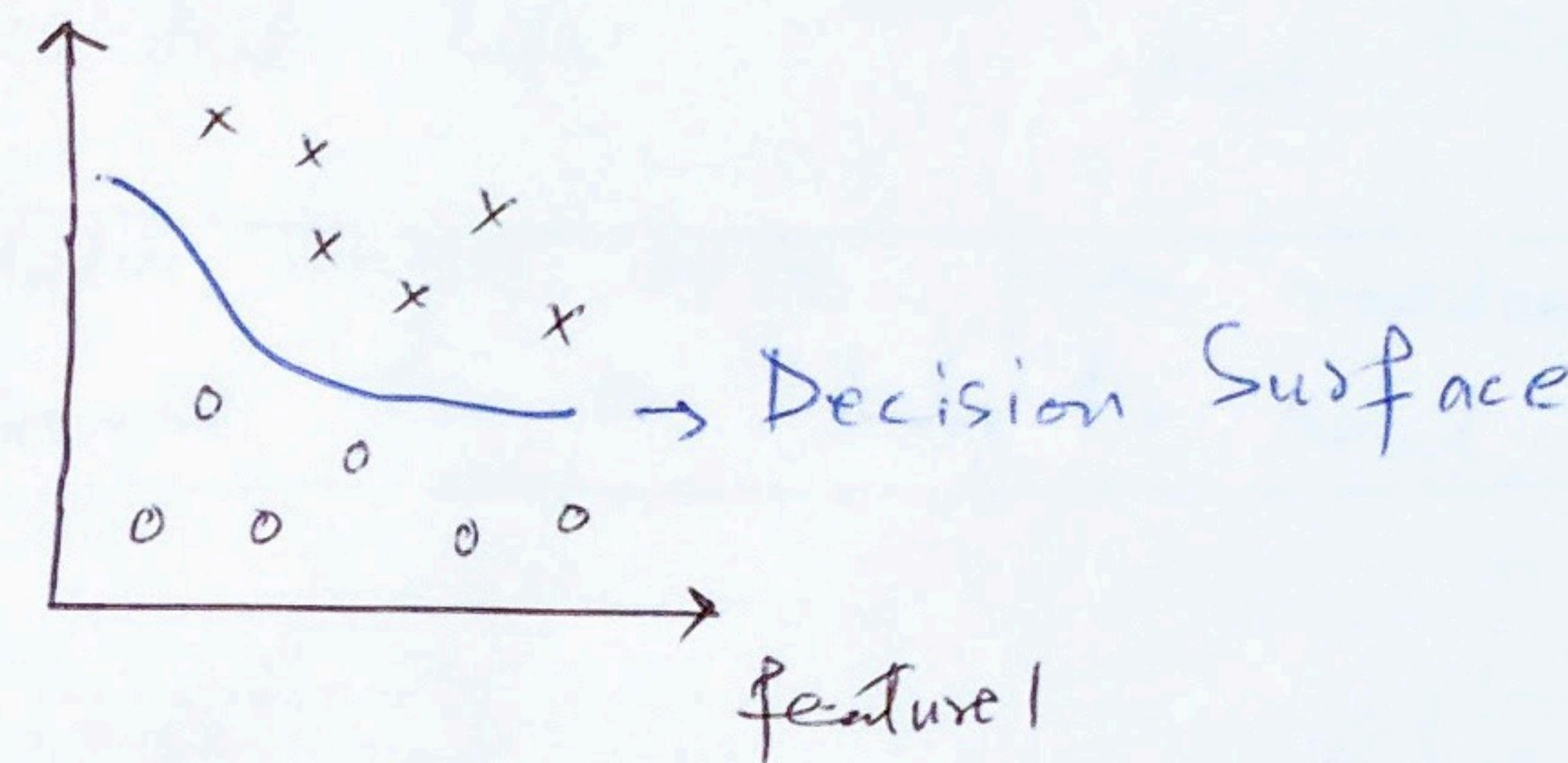


⇒ Acerous → animals w/ ~~no~~ horns.
 ⇒ Non-acerosus → animals w/ horns.
 feature
 & attribute

①

- * One of imp. things in ML is where to put (arrange) new data (or point) based on previous data.

feature 2



- * Naive Bayes is used to find decision surfaces.

* scikit-Learn:

- 1 - Create classifier.
- 2 - train it by fitting the features & Label data
- 3 - predict Label using test feature.

* Accuracy = $\frac{\# \text{ of points correctly classified}}{\text{all points in test set}}$

⇒ Bayes Rule:

Question:

$$P(c) = 1\%$$

$$\text{Sensitivity } P(\text{Pos}|c) = 90\%$$

$$\text{Specificity } P(\text{Neg}|>c) = 90\%$$

$$P(c|\text{Pos}) = ?$$

$$P(\text{Pos}|c) = 1 - P(\text{Neg}|c)$$

$$P(>c) = 0.1$$

$$P(c) = 1 - P(>c) = 0.99$$

Sol.

$$P(c|\text{Pos}) = \frac{P(\text{Pos}|c)P(c)}{P(\text{Pos})}$$

$$P(c|\text{Pos}) = \frac{P(\text{Pos}|c)P(c)}{P(\text{Pos}) + P(\text{Pos}|>c)P(>c)}$$

$$P(c|\text{Pos}) = \frac{0.9(0.01)}{0.9(0.01) + 0.1(0.99)}$$

$P(c|\text{Pos}) = 8\% \text{ or } 8.33\%$

8% chance of cancer given positive test

* return %

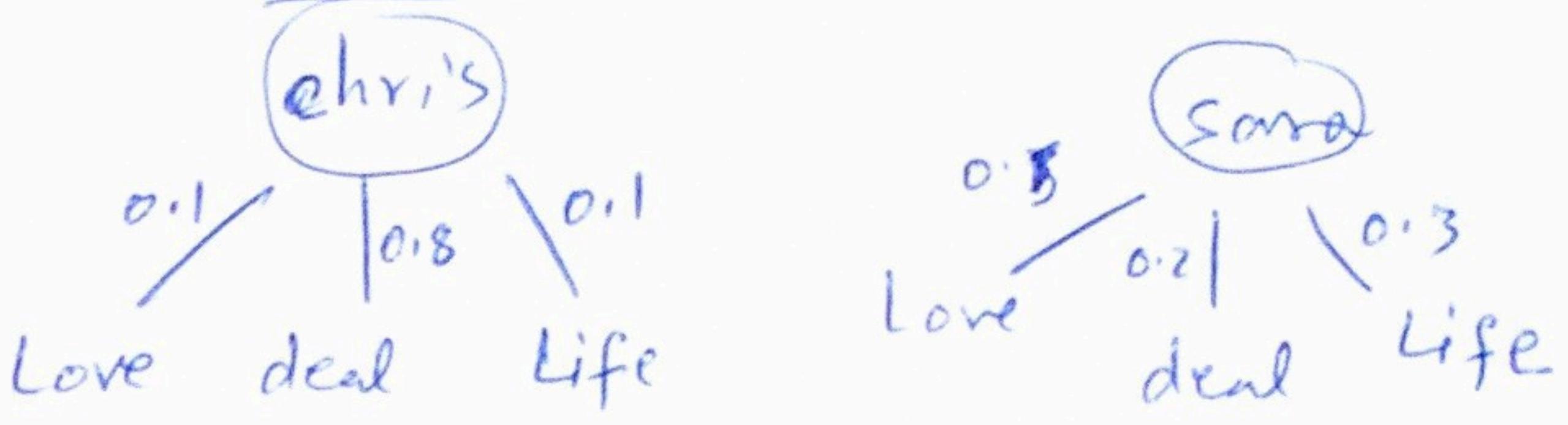
$x \neq y$

7 May (Night Shift)

(2)

→ Text learning:

Naive Bayes Rule:



+ Naive Bayes tells from random email who is likely to send this email

msg: Love Deal!

⇒ apriori:

$$P(\text{Chris}) = 0.5$$

$$P(\text{Sara}) = 0.5$$

• Sara is most likely to have sent this message

⇒ ~~apriori~~ New msg:

LIFE Deal

• More likely from Chris - As:

$$\text{Chris} = \frac{(0.1)(0.8)(0.5)}{\text{Life deal apriori}} = 0.04$$

$$\text{Sara} = \frac{(0.3)(0.2)(0.5)}{0.04 + 0.03} = 0.03$$

⇒ posterior:

$$P(\text{Chris} | \text{"Life Deal"}) = \frac{P(\text{"Life Deal"} | \text{Chris}) P(\text{Chris})}{P(\text{"Life Deal"} | \text{Chris}) P(\text{Chris}) + P(\text{"Life Deal"} | \text{Sara}) P(\text{Sara})} = \frac{0.04}{0.04 + 0.03} = 0.57$$

$$P(\text{Sara} | \text{"Life Deal"}) = \frac{P(\text{"Life Deal"} | \text{Sara}) P(\text{Sara})}{P(\text{"Life Deal"} | \text{Sara}) P(\text{Sara}) + P(\text{"Life Deal"} | \text{Chris}) P(\text{Chris})} = \frac{0.03}{0.03 + 0.04} = 0.428$$

* return *

⇒ Where ~~Naive~~ Naive Bayes Breaks:

(3)

→ New msg:

method-1: "Love Deal"

$$P(\text{chris} | \text{"Love Deal"}) = \frac{(0.1)(0.8)(0.5)}{(0.5)(0.2)(0.5)} = 0.04 \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{joint}$$

$$P(\text{Sara}, \text{"Love Deal"}) = (0.5)(0.2)(0.5) = 0.05$$

And:

$$\left. \begin{array}{l} P(\text{chris} | \text{"Love Deal"}) = \frac{0.04}{0.04 + 0.05} = 0.444 \\ P(\text{Sara} | \text{"Love Deal"}) = \frac{0.05}{0.04 + 0.05} = 0.555 \end{array} \right\} \text{posterior}$$

Method-2:

$$\begin{aligned} P(\text{chris} | \text{"Love Deal"}) &= \frac{P(\text{"Love Deal"} | \text{chris}) P(\text{chris})}{P(\text{"Love Deal"} | \text{chris}) P(\text{chris}) + P(\text{"Love Deal"} | \text{Sara}) P(\text{Sara})} \\ &= \frac{(0.1)(0.8)(0.5)}{(0.1)(0.8)(0.5) + (0.5)(0.2)(0.5)} = 0.444 \end{aligned}$$

* Actually ~~Method 1 &~~ Method 1 & Method 2 both are same!

⇒ Why Naive Bayes is Naive?



They are hidden,
always in supervised
Learning

* Every word gives us evidence either its person A or person B

* We multiply probabilities of all the words appearing in msg.

A $\boxed{\quad}$ Ratio
B $\boxed{\quad}$

↑
prior
probabilities

+ Called Naive, b/c it ignores the order of the words in message

In case of word order, so it won't differentiate b/w "Bull Chicago" & "Chicago Bull" although both have totally diff. meanings & "Chicago Bull" is actually 1 word not 2.

⇒ Project:

→ Based on text of email, classify which of two persons.

→ Naive Bayes is very good for text classification as it treats each word as a feature & where features are independent.

⇒ SVM [Support vector machine]

Invented by
Russian

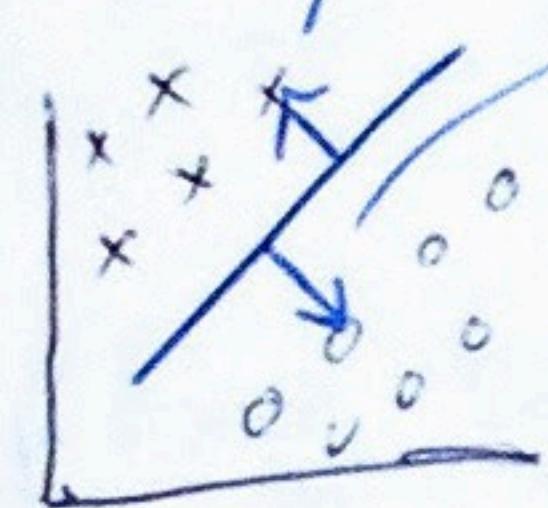
- * It is a supervised learning method like

Naive Bayes Rule.

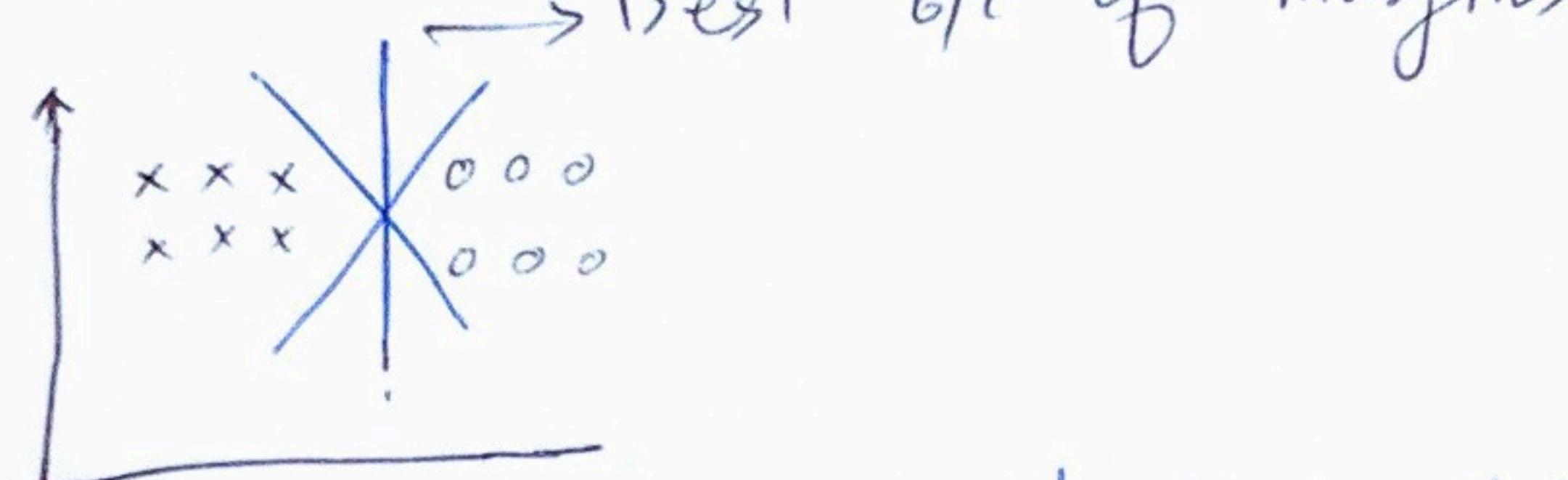
- * It takes data

2 sets/classes.

margin



Creates/forms boundary b/w
Best b/c of Margins



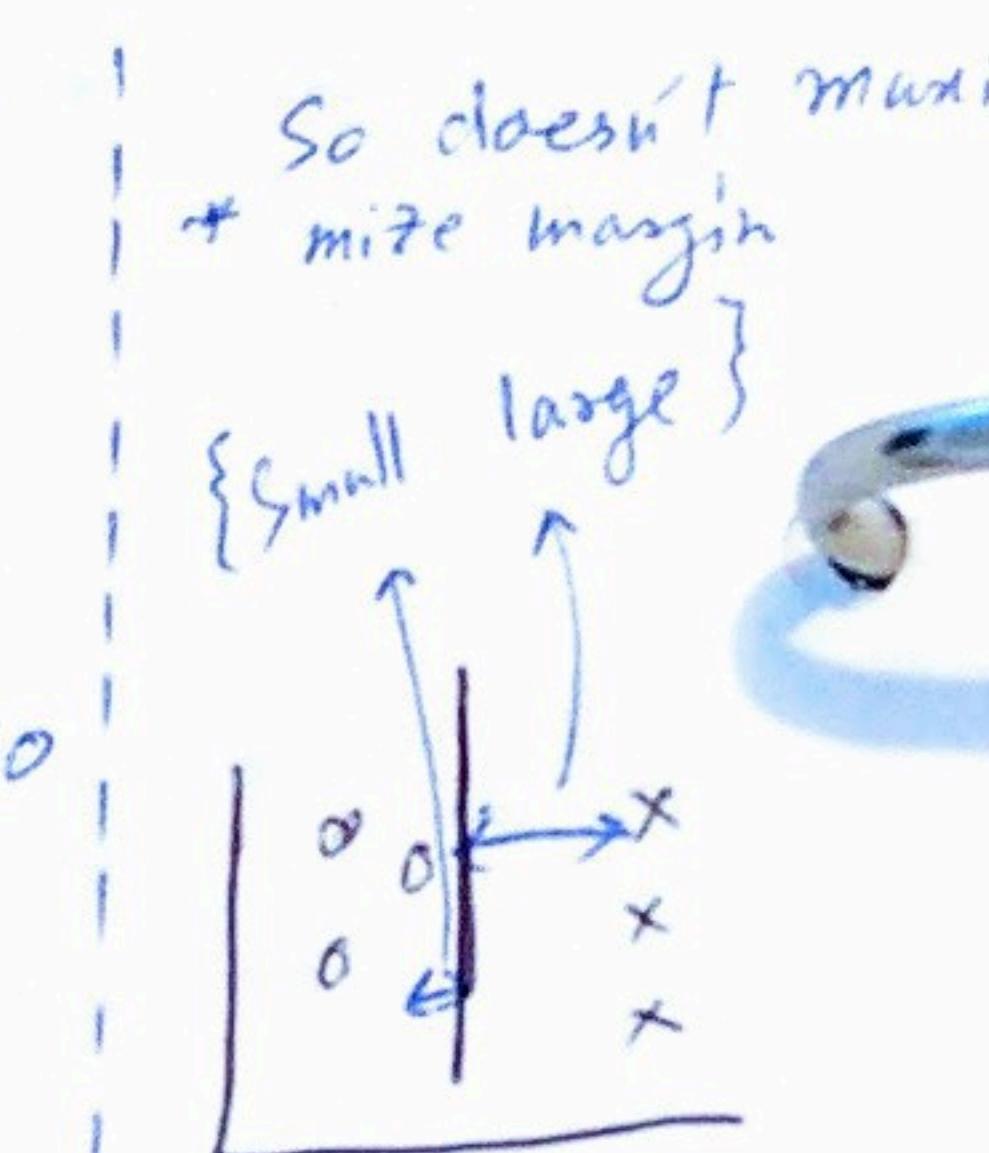
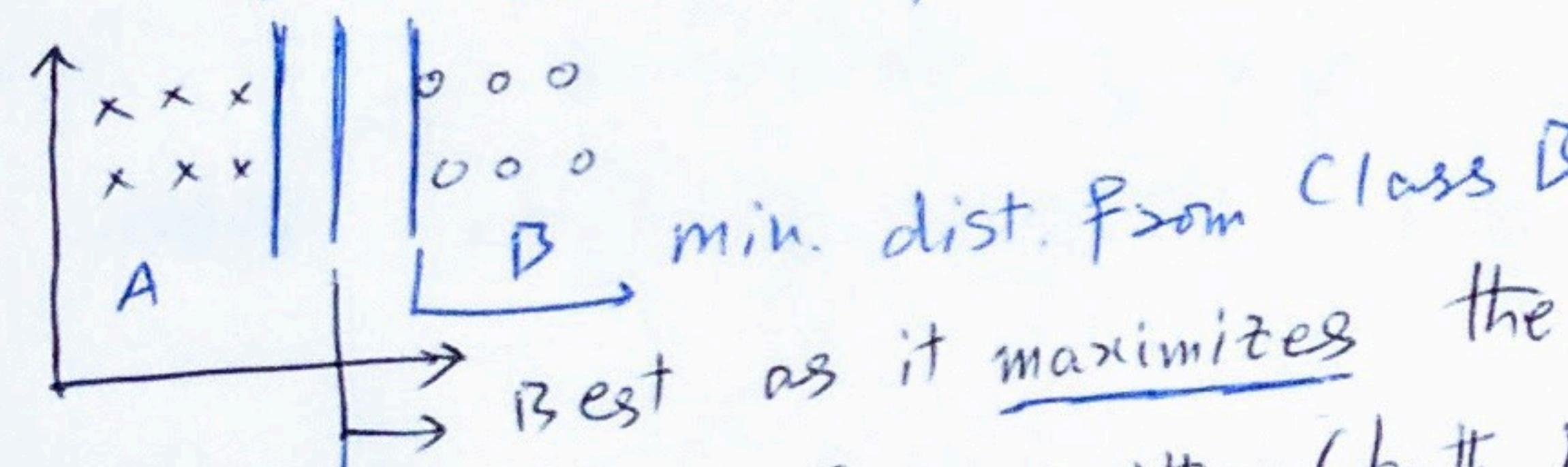
- * What makes a good separator?

Margin → max. distance to the nearest point from either of the two

min. dist.
from class A

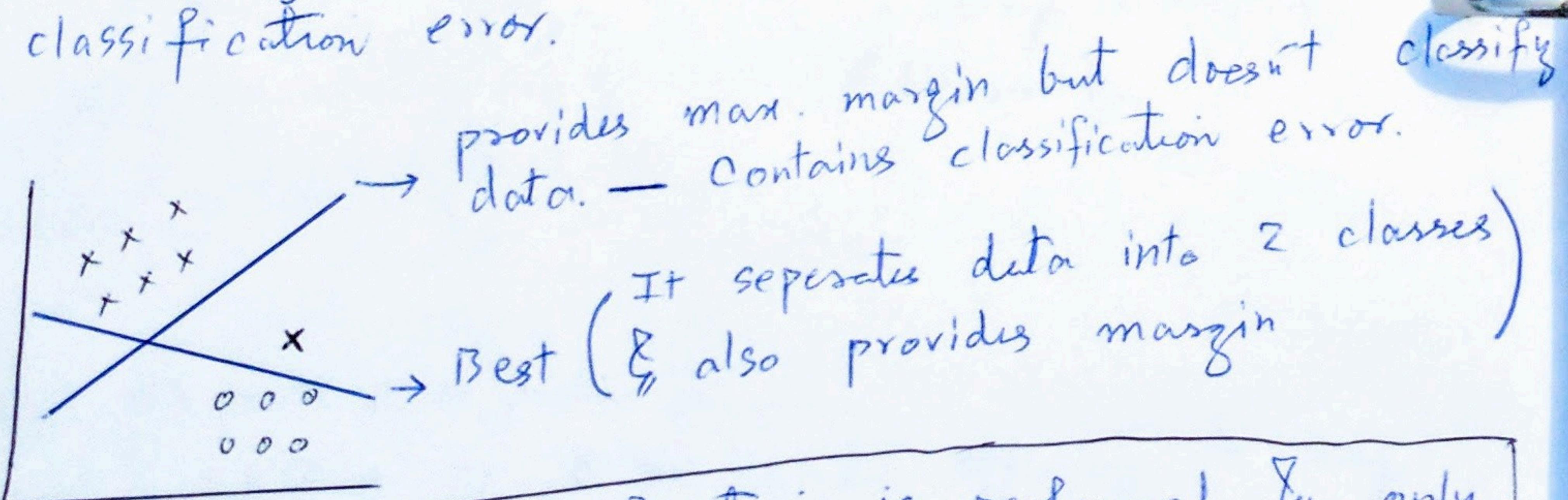
classes.

max. dist. from classes A & B

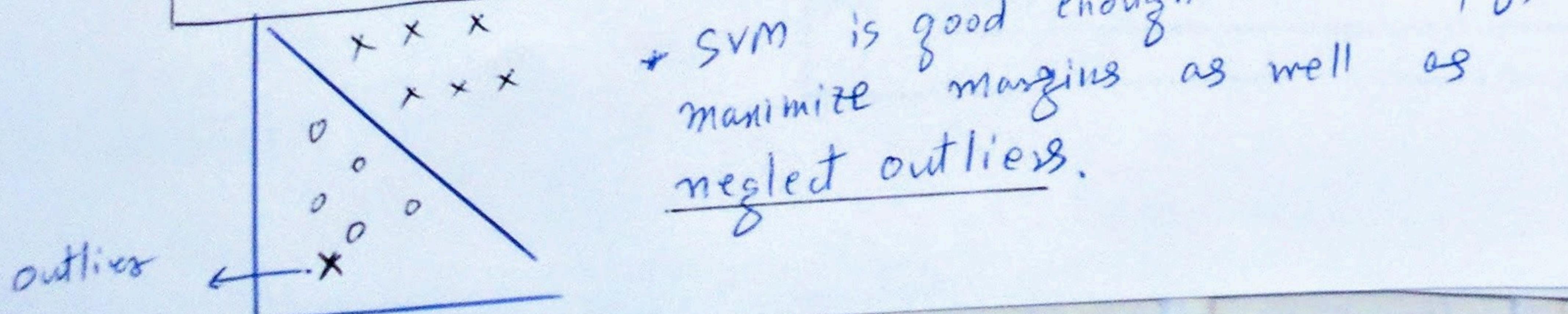


min. dist. from class B
Best as it maximizes the margin to nearest point from either (both) of ~~the~~ classes

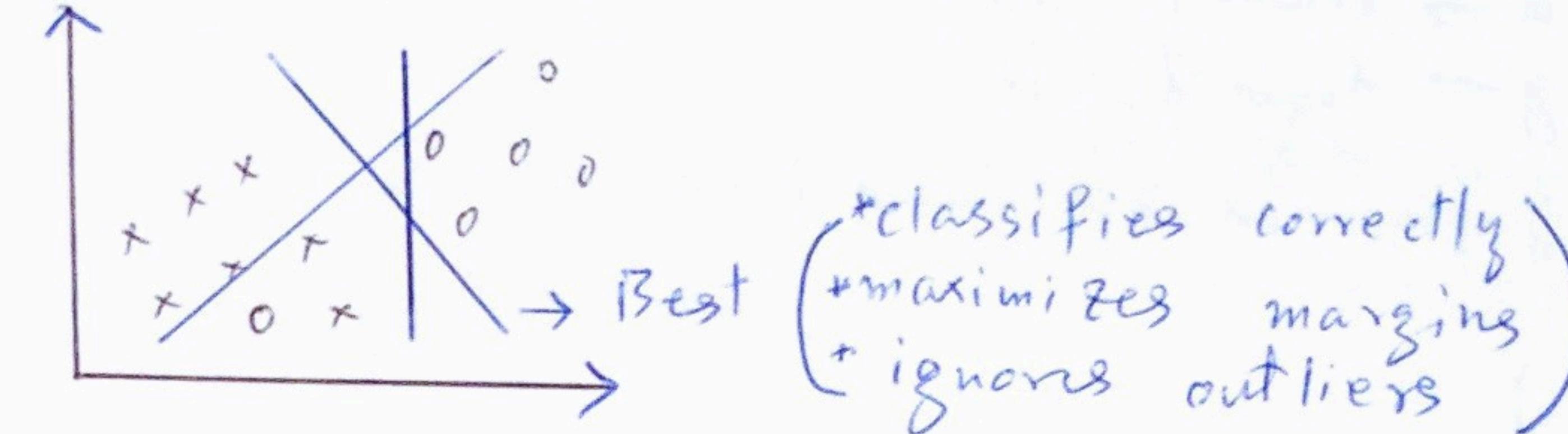
- * Max. margin on either side reduces classification error.



So first correct classification is performed & only then the margins are maximized.



SVM is good enough to classify, maximize margins as well as neglect outliers.



Adv. :

- * SVM is efficient in high dimensional spaces.

* Effective when dimensions are more than number of samples.

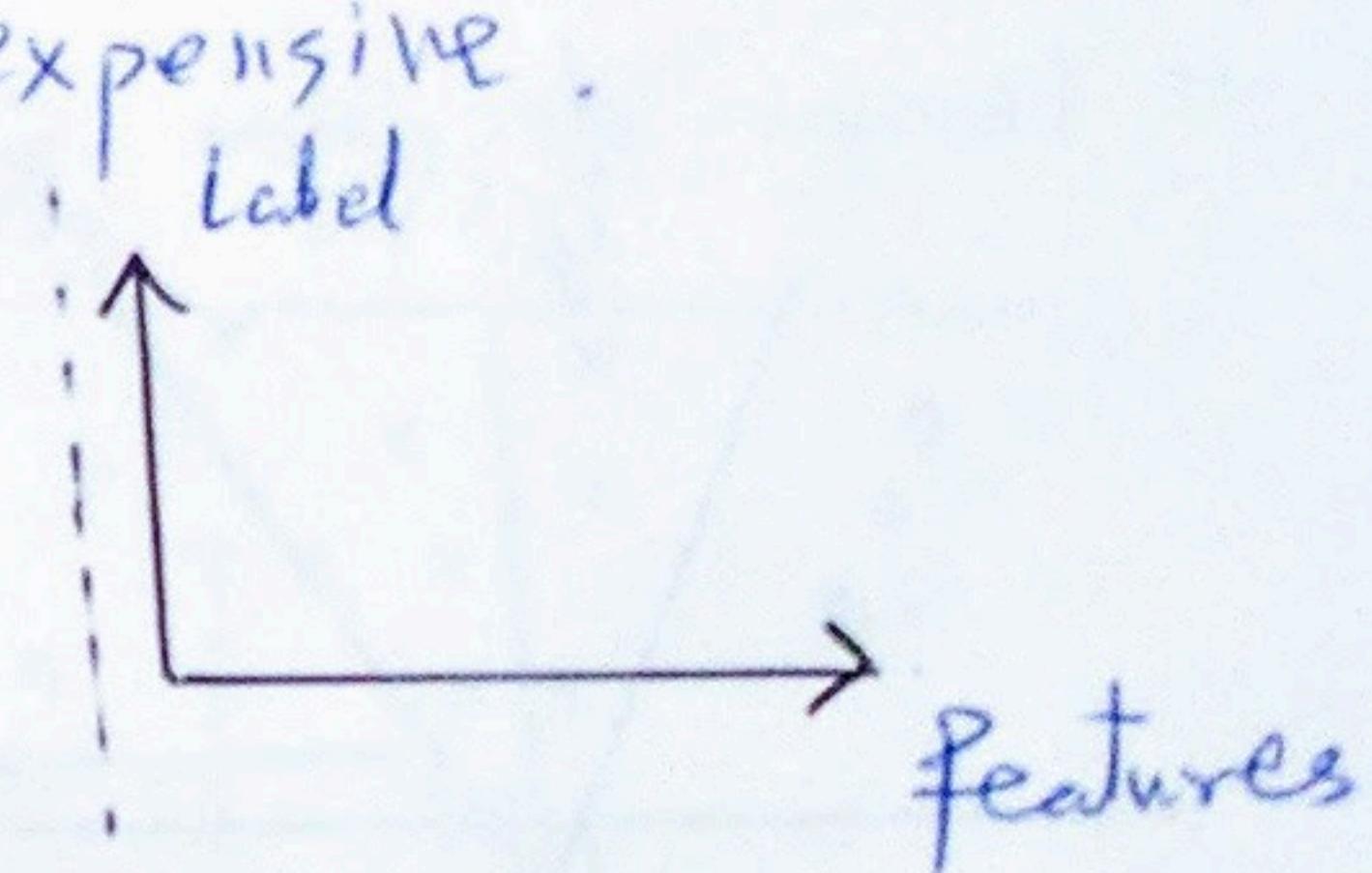
- * Diff. kernel functions can be used for decision function.

Disadv. :

- * If # of features is much greater than # of samples, it performs poorly — Hughes Phenomenon — Curve of Dimensions

* doesn't provide probability estimates. 5-fold cross validation is used for probability estimate which is expensive.

- * Margins make sure that there is as much separation among classes as possible.



Linear SVM will not work on this kind of data.

- * There is no linear Hyperplane (separator) b/w two classes.

* input for training → features & labels

* input for predict → test features & output → predicted labels

* accuracy score → compare test labels with predicted labels

* return q

④
* scikit-learn
Machine learning library.
Classification, Clustering, Regression
Naive Bayes, SVM,
Random Forest, KMeans ...

* Google Summer of Code

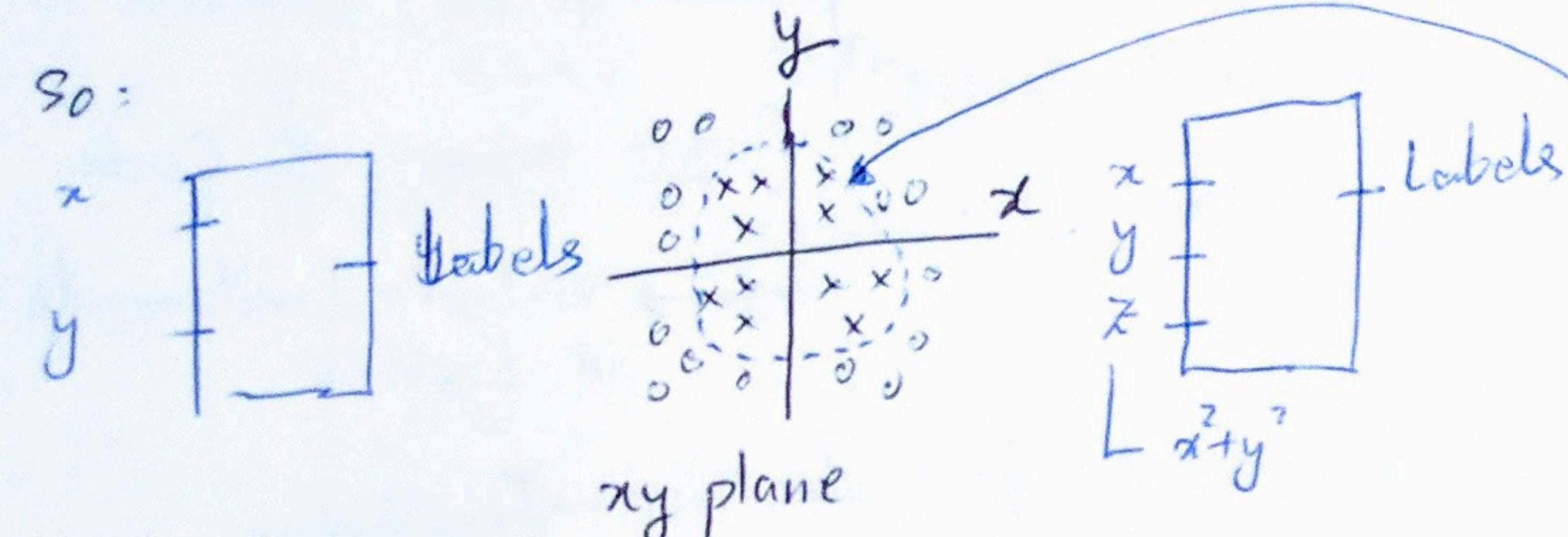
* VI → Virtual Instruments
in LabView

x → train-features.
 y → train-labels.
 y -predicted → predicted Labels

$$x^2 + y^2 = \boxed{?}$$

↳ non-negative (dist. from origin)

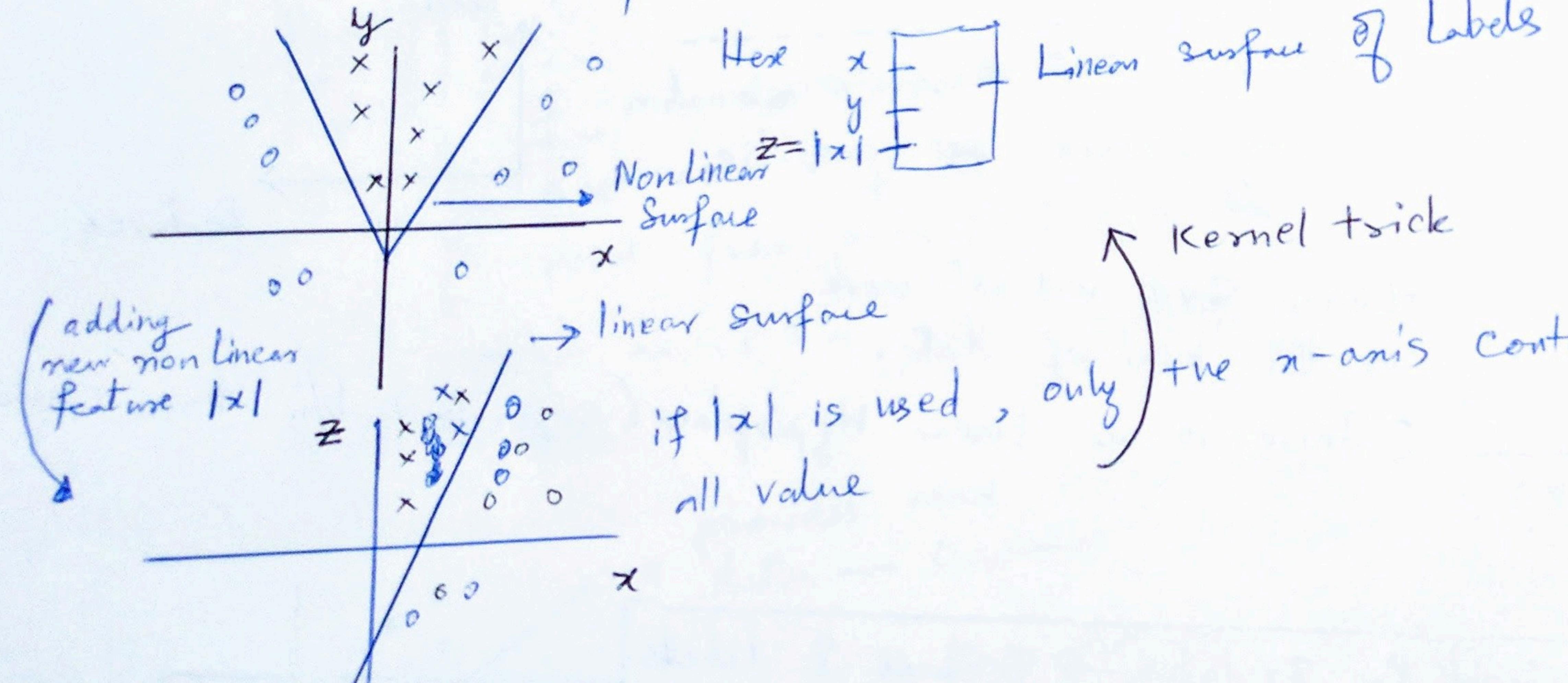
so:



Line in xz plane
 corresponds to circle in xy plane

$xy \rightarrow$ here its non Linear decision Surface

$xz \rightarrow$ linear decision Surface



⇒ Kernel Trick:

Low dimensional

x, y ↗ Kernels

non separable

Non Linear Separation

high dimensional

$x_1 x_2 x_3 x_4 x_5$

separable

Linear Separable Solution

③ We use kernel to go from low dimension to high dimension & find linear separation, using info of linear separation we can find nonlinear separation in original data

+ There are several SVMs in scikit-learn, we are using SVC (SVM classifier)

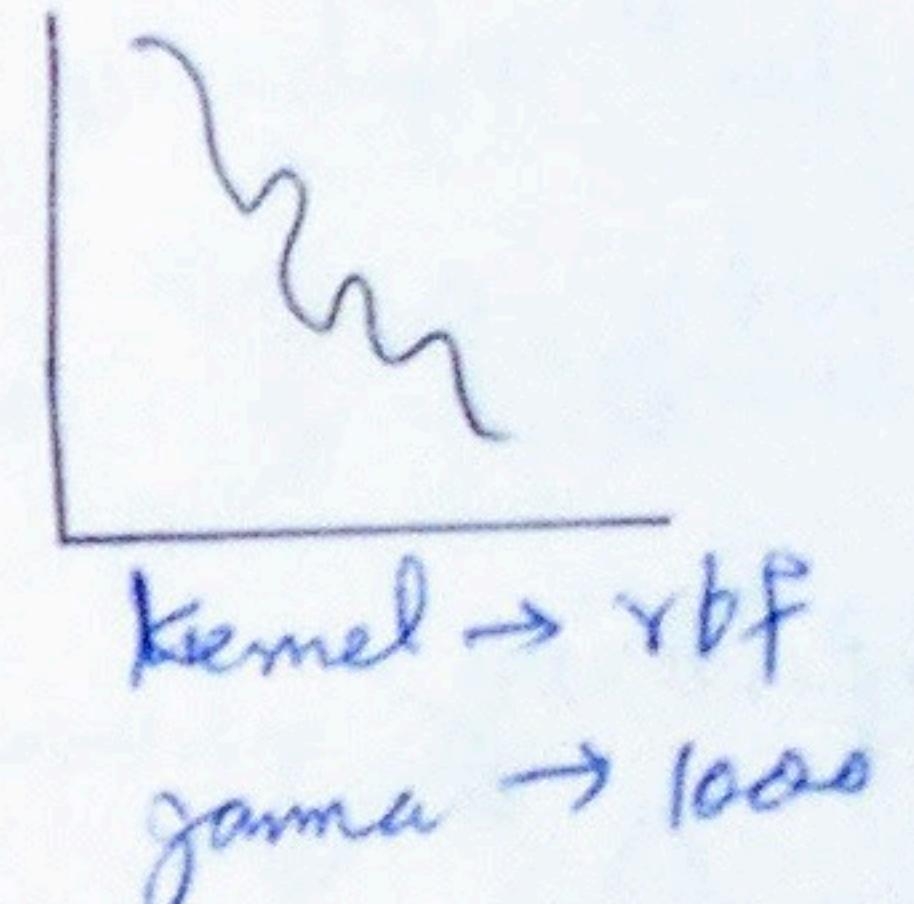
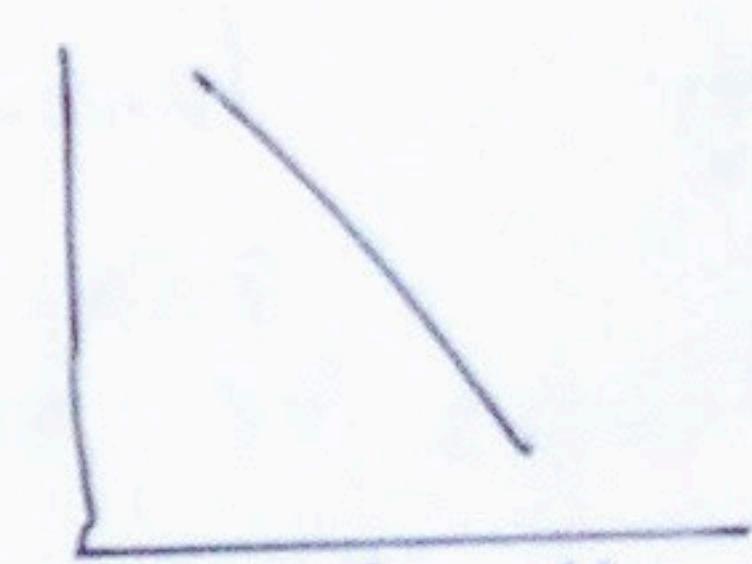
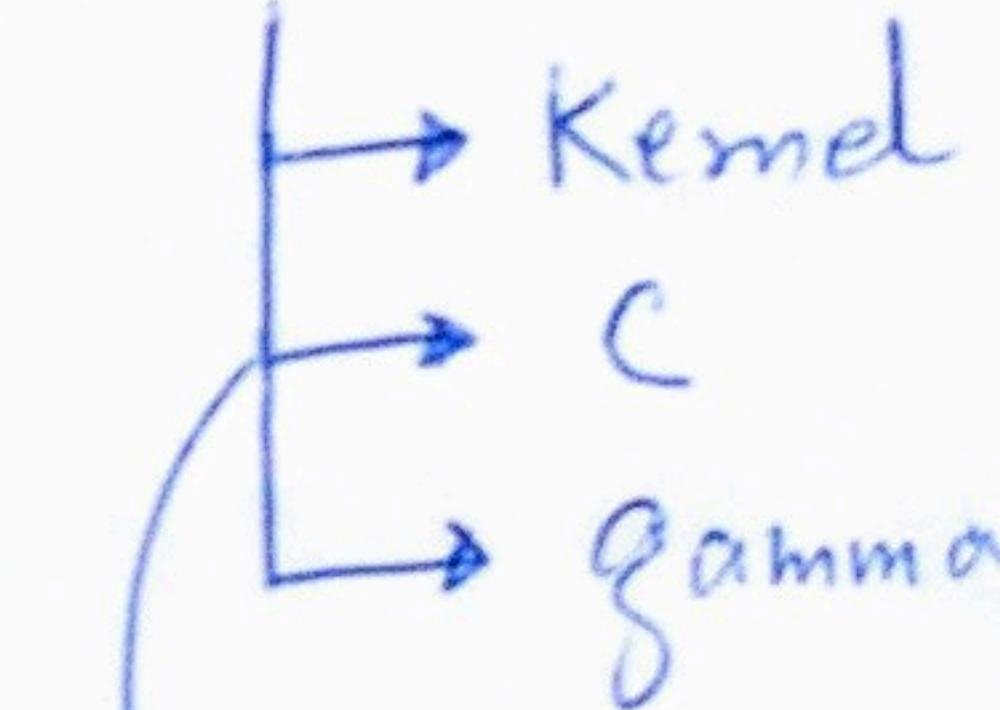
SVC supports these kernels:

- 1 - Linear
- 4 - Sigmoid
- 2 - poly
- 5 - precomputed
- 3 - rbf

* default is rbf

+ Parameters → passed when creating classifier before fit

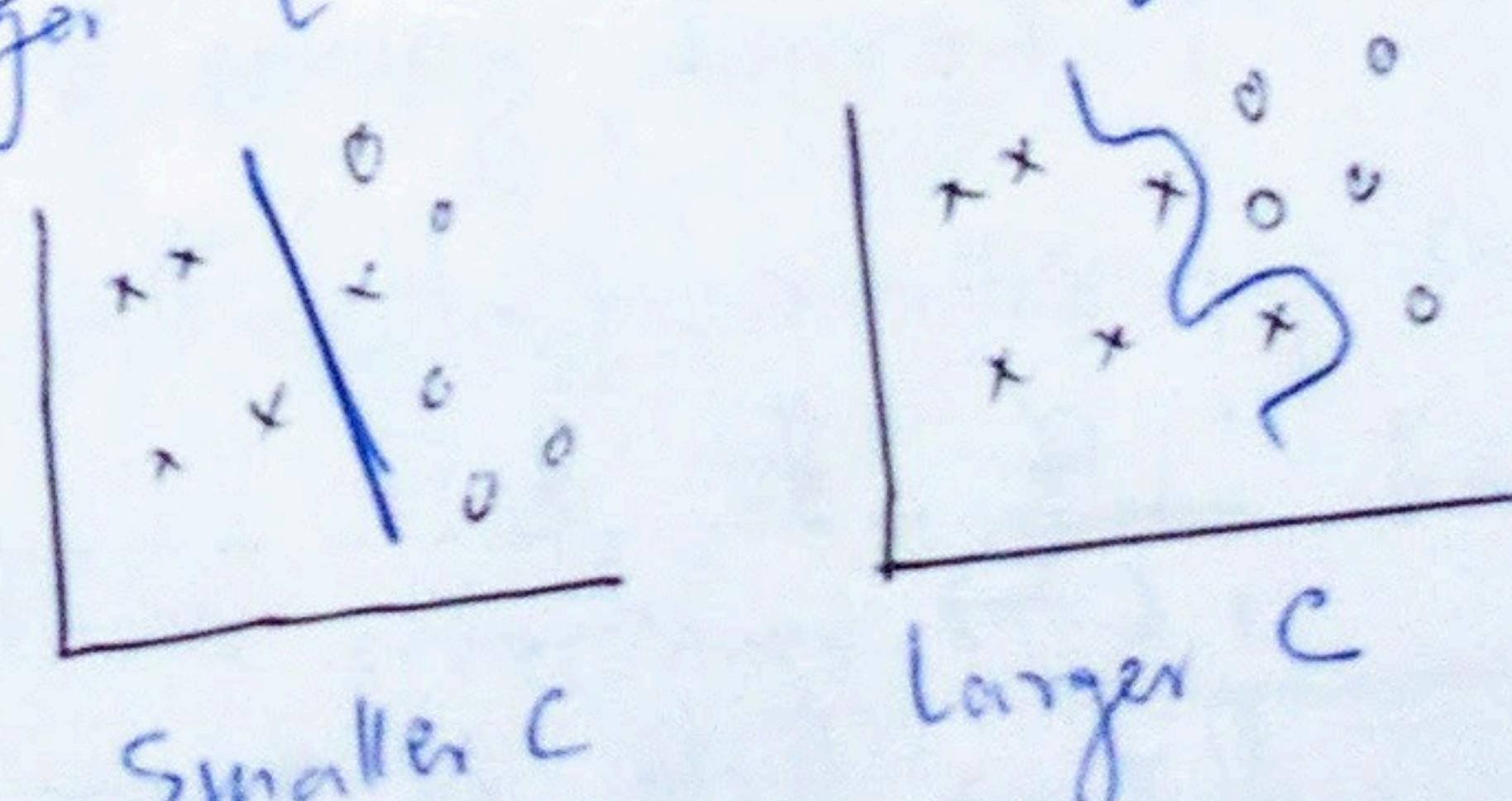
params for SVM



Kernel → rbf
gamma → 1000

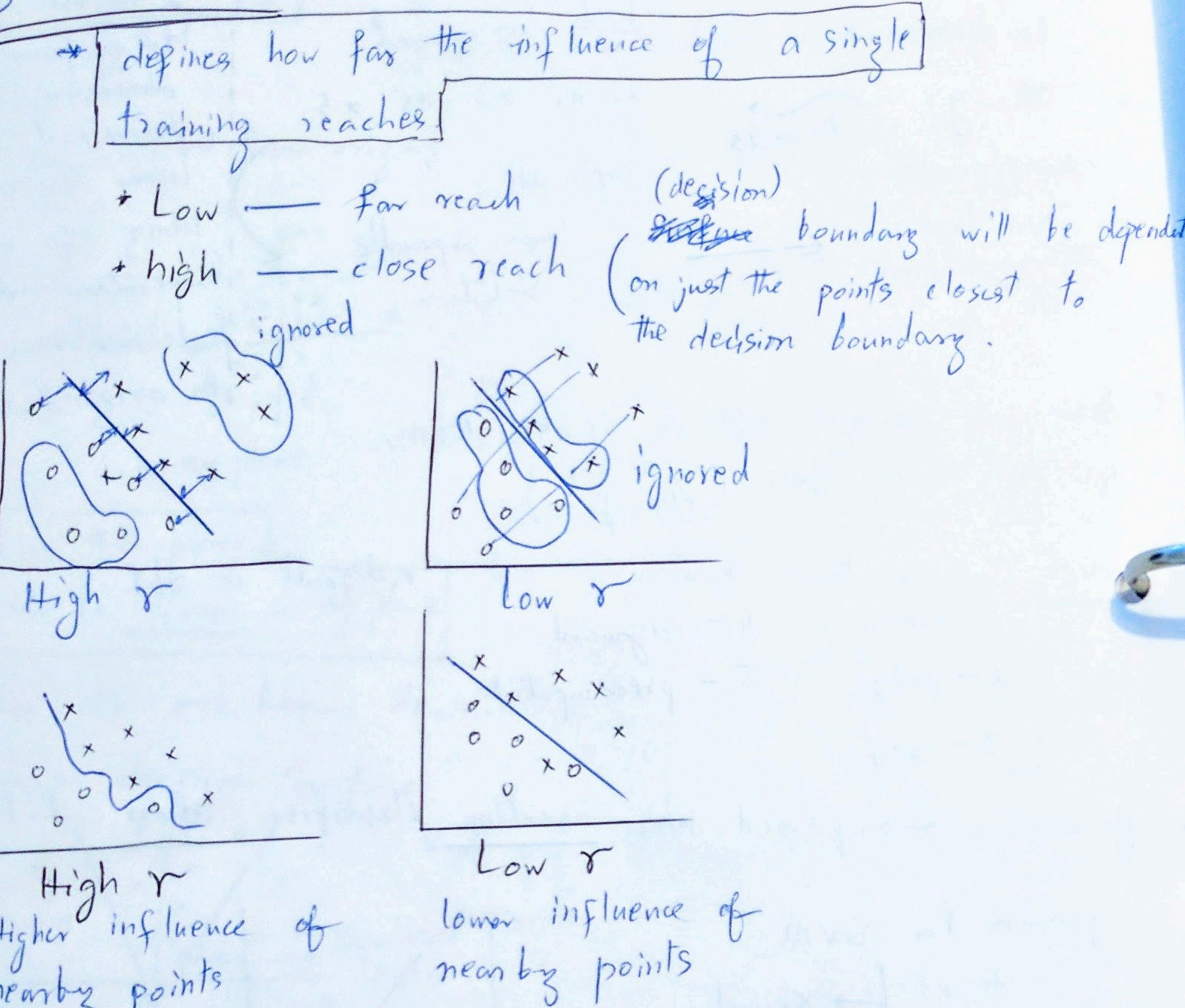
→ controls tradeoff b/w smooth decision boundary & classifying training points correctly.

- * smaller C → smooth -
- * larger C → correctly classifying but less smooth.

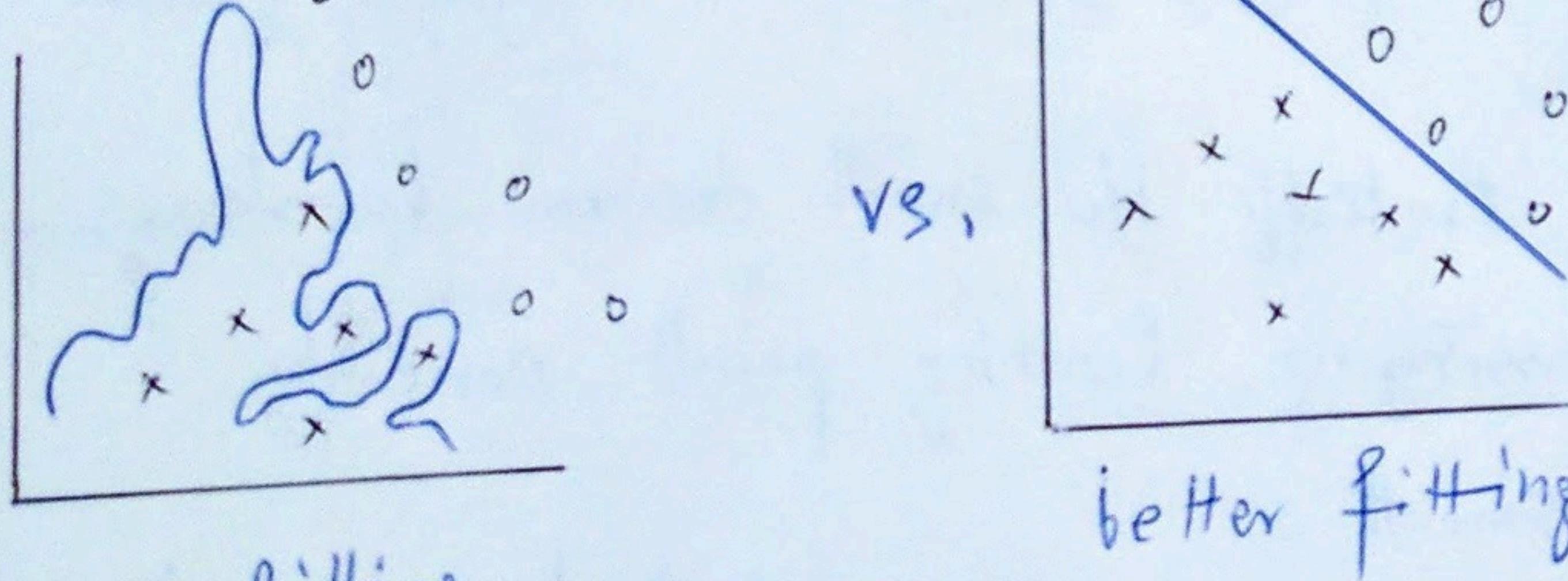


Smaller C
Larger C

→ Gamma in SVM:



→ Overfitting



- * when classifier takes the data too literally.
- * classifier is trained on same kind of training data again & again w/o variation. So its now overfit for that data & not generalized

Kernel, gamma & C
all controls the overfitting of the data:
 $C \rightarrow$ smoothness vs. best boundary
 $\gamma \rightarrow$ higher vs. lower influence of nearby points
Kernel \rightarrow Linear, Poly ...

Adv. of SVM:

Disadv. of SVM:

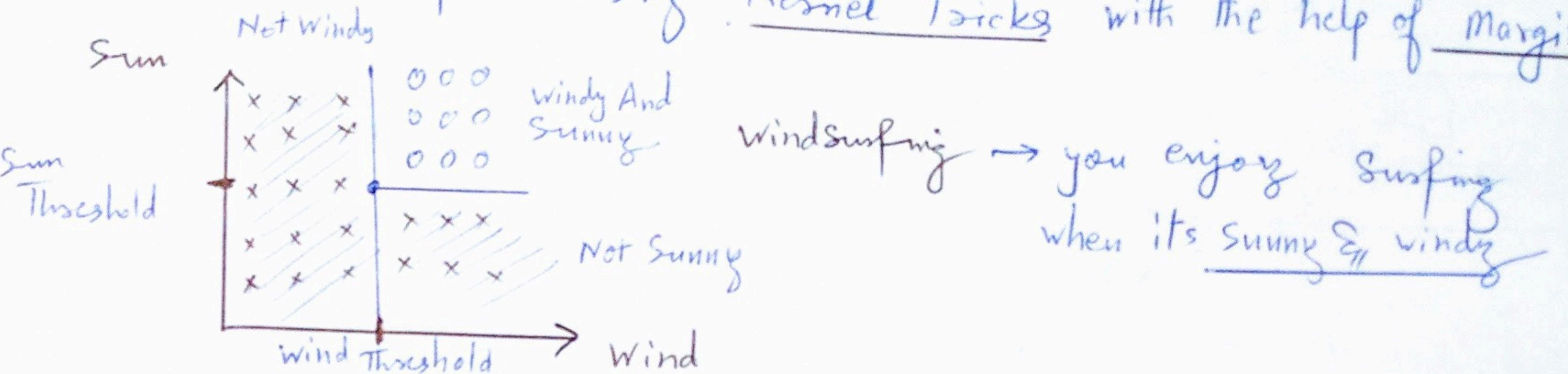
- * Noisy data, overlapping data \rightarrow doesn't work good.
- * With very large dataset & lots of feature, SVM can be slow. Naive Bayes can be a much better ~~classifier~~ for these situations as feature are independent & order doesn't matter in Naive Bayes.
- * Need to ~~perform~~ count "independent evidence" when classes are too overlapping i.e. noisy data
- * The more training data we have, the more is accuracy but slower is the data fitting — most of the time is consumed in training the classifier — Prediction takes very smaller time.
- * In real-time usage like voice recognition & transactions blocking etc., we prefer speed over accuracy.
- * naive bayes is usually way faster than SVM for large data sets.
- * Using combination of high C & soft kernel can give very good accuracy on smaller data set but at the expense of very complicated decision surface by the classifier will be much less generalized which is very bad.
- But if increase the data set size, classifier will be slow in training. — so right balance needs to be achieved.

* return q

- `svm.predict(test-labels
feature)`
- * It provides us predicted class for each data point of test-~~labels~~
features.
- * `gridCV` → sklearn tool that finds optimal parameters.

→ Decision Trees:

* In SVM, we were able to change from Linear to Non Linear Decision Surface using Kernel Tricks with the help of Margins



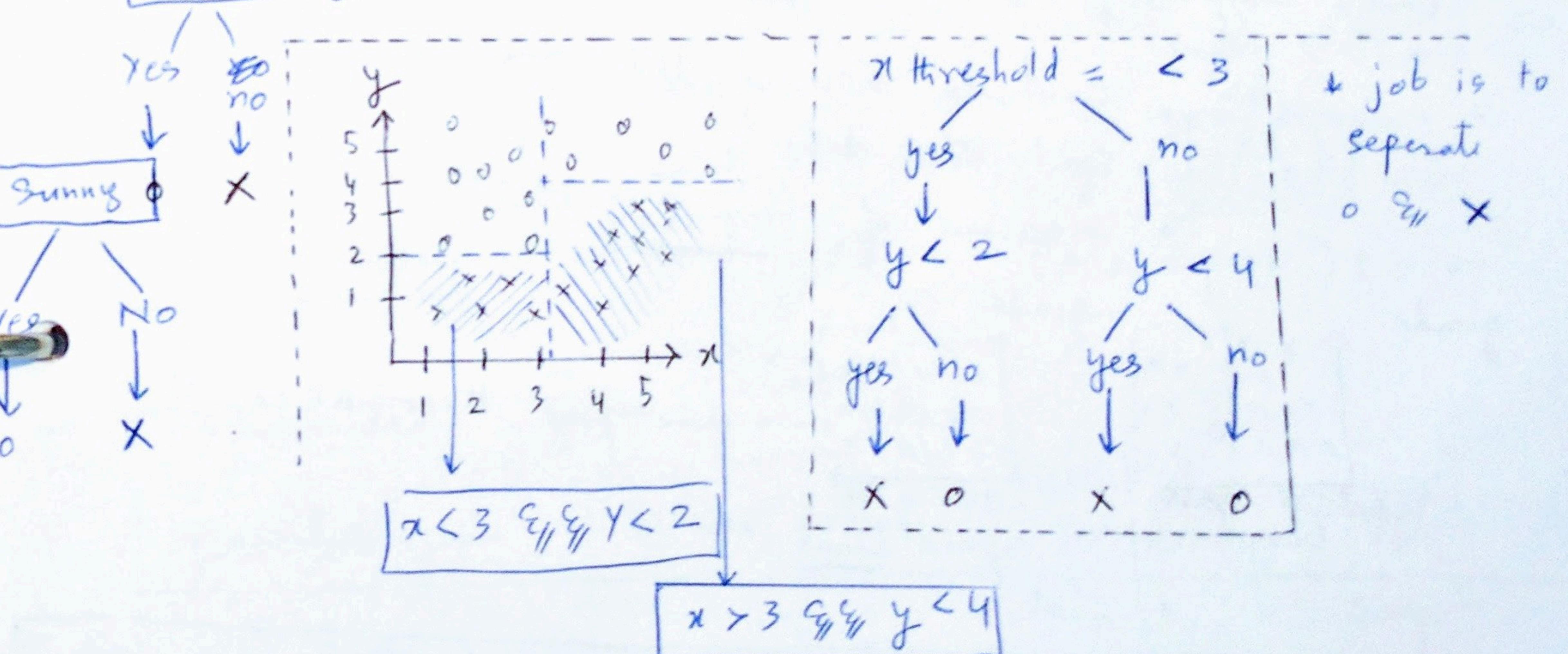
Windsurfing → you enjoy Surfing when it's Sunny & Windy

- * $ooo \rightarrow$ suitable condition for wind surfing
- * It's not linearly separable
- * Decision Tree allows to ask multiple linear Question one after another.

* Question 1:

Is it windy (with the threshold)

* So we classify data using data structure called Tree.

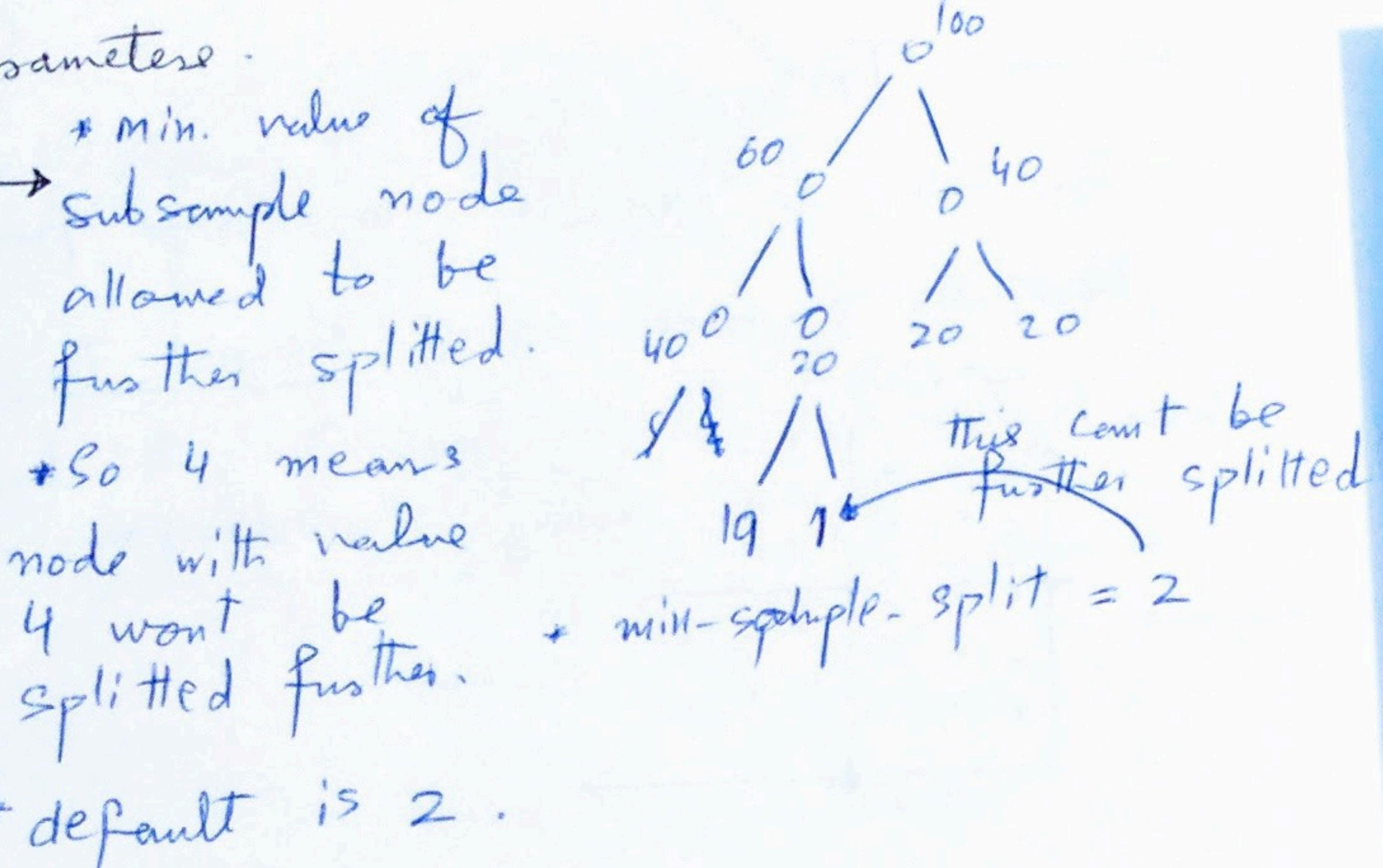


- * Cost of running tree is Logarithmic in the number of data point in training.

- * Uses white box model → any observable situation can be explained by boolean logic unlike artificial neural networks.
- * Can have overfit by creating over-complex decision trees.

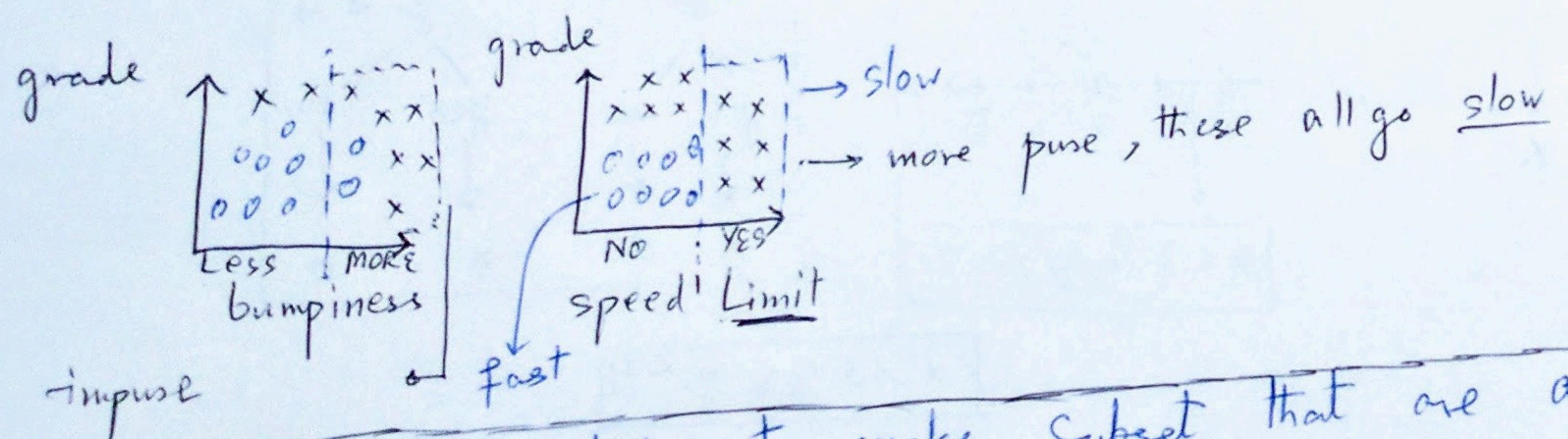
⇒ Decision Tree Parameters:

- * min-sample-split →
 - * min. value of subsample node allowed to be further splitted.
 - * So 4 means node with value 4 won't be splitted further.
 - * min-sample-split = 2
 - * default is 2.
- if min-sample-split = 2
Then it'll create this whole decision surface just to grab this point!



⇒ Data Impurity & Entropy:

- * entropy →
 - Controls how a DT decides where to split the data.
 - measure of impurity



- * so decision tree tries to make subset that are as pure as possible & that's how DT actually makes its decision

$$\text{Entropy} = \sum_i -P_i \log_2(P_i)$$

pure $\rightarrow 0 \leq \text{entropy} \leq 1$
impure
Fraction of examples in a given class.
2 classes can be 10, e...
sum over all classes available

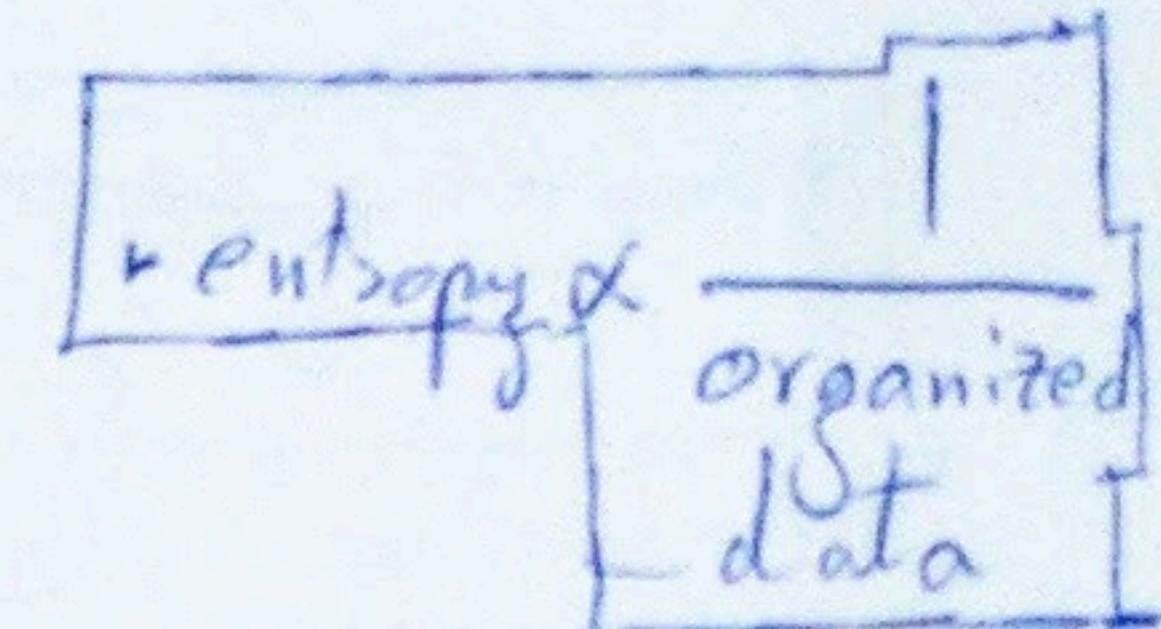
Cases for Entropy:

- * all examples are same class:

$$\text{entropy} = 0 \rightarrow \text{pure}$$

- * examples are evenly split b/w classes:

$$\text{entropy} = 1.0 \rightarrow \text{impure}$$



Lower entropy → More organized data
Higher entropy → Less organized data

grade	bumpiness	speed limit?	speed
steep	bumpy	yes	slow
steep	smooth	yes	slow
flat	bumpy	no	fast
steep	smooth	no	fast

$$+ p_{\text{slow}} = 0.5 \rightarrow \text{fraction of slow}$$

$$+ p_{\text{fast}} = 0.5 \rightarrow \text{fraction of fast}$$

so

$$\text{entropy} = \frac{0.5 + 0.5}{2} \cdot 1.0 = -0.5 \log_2(0.5) + -0.5 \log_2(0.5) = 1.0$$

impure state ↴

+ Most impure state is when samples are evenly split b/w 2 classes.

+ entropy of parent first(start) node (max impure) 1.0 ↴

⇒ Information Gain: (IG)

$$\text{information gain} = \text{entropy}(\text{parent}) - \frac{\text{Weighted average of entropy(children)}}{\text{children}}$$

- * decision tree maximizes information gain

$$\text{weight} = \frac{3}{4}$$

steep	ssff	flat	f	weight = 1/4
				(entropy = 0 i.e. pure i.e. this class has sample of same type)

pslow = $\frac{2}{3}$ ↴ 2slow, 1fast
pslow = $\frac{1}{3}$
 $\text{entropy} = \left(\frac{2}{3}\right) \log_2\left(\frac{2}{3}\right) - \frac{1}{3} \log_2\left(\frac{1}{3}\right) = 0.389 + 0.528 = 0.9183 \rightarrow \text{less organized (impure)}$

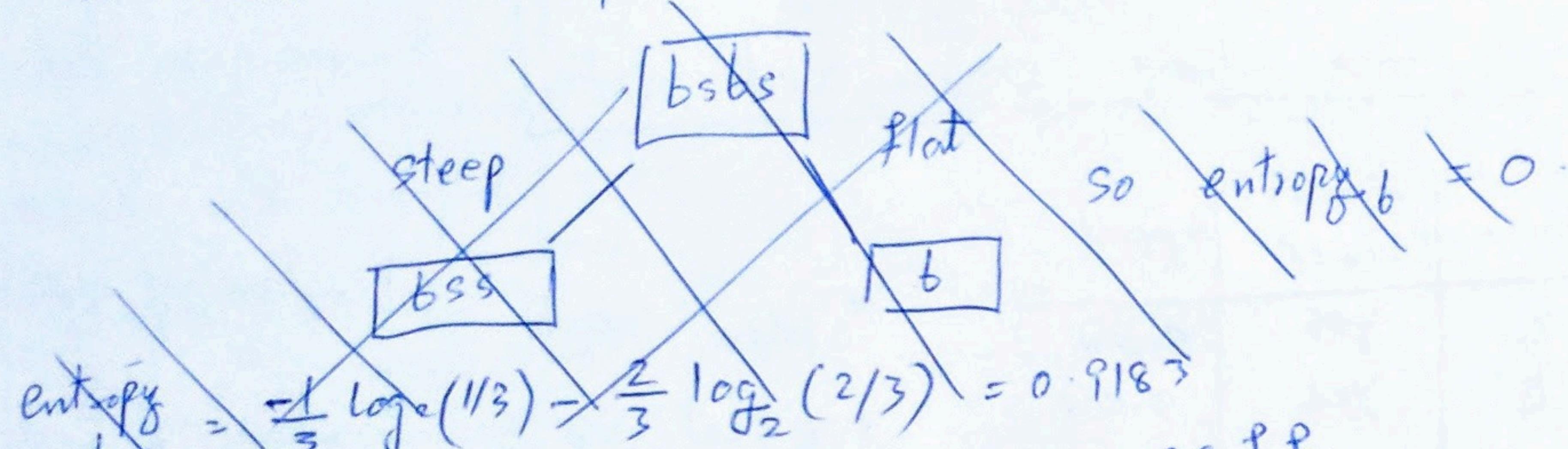
* return ↴

$$\text{So } I.G_I = (\text{parent entropy}) - [\text{weight}_1 \text{ (children entropy)} + \text{weight}_2 \text{ (children entropy)}]$$

$$I.G_I = 1.0 - \left[\frac{3}{4} (0.9183) + \frac{1}{4} (0.0) \right]$$

$$\boxed{I.G_I = 0.0817 \quad 0.08175}$$

Consider bumpiness



$$\text{entropy}_{bsbs} = \frac{2}{3} \log_2(1/3) + \frac{1}{3} \log_2(2/3) = 0.9183$$

$\text{entropy}_{ss} = 1.0$	$\text{entropy}_{ff} = 1.0$
$-0.5 \log_2(0.5) +$	$-0.5 \log_2(0.5) +$
$-0.5 \log_2(0.5)$	$-0.5 \log_2(0.5)$

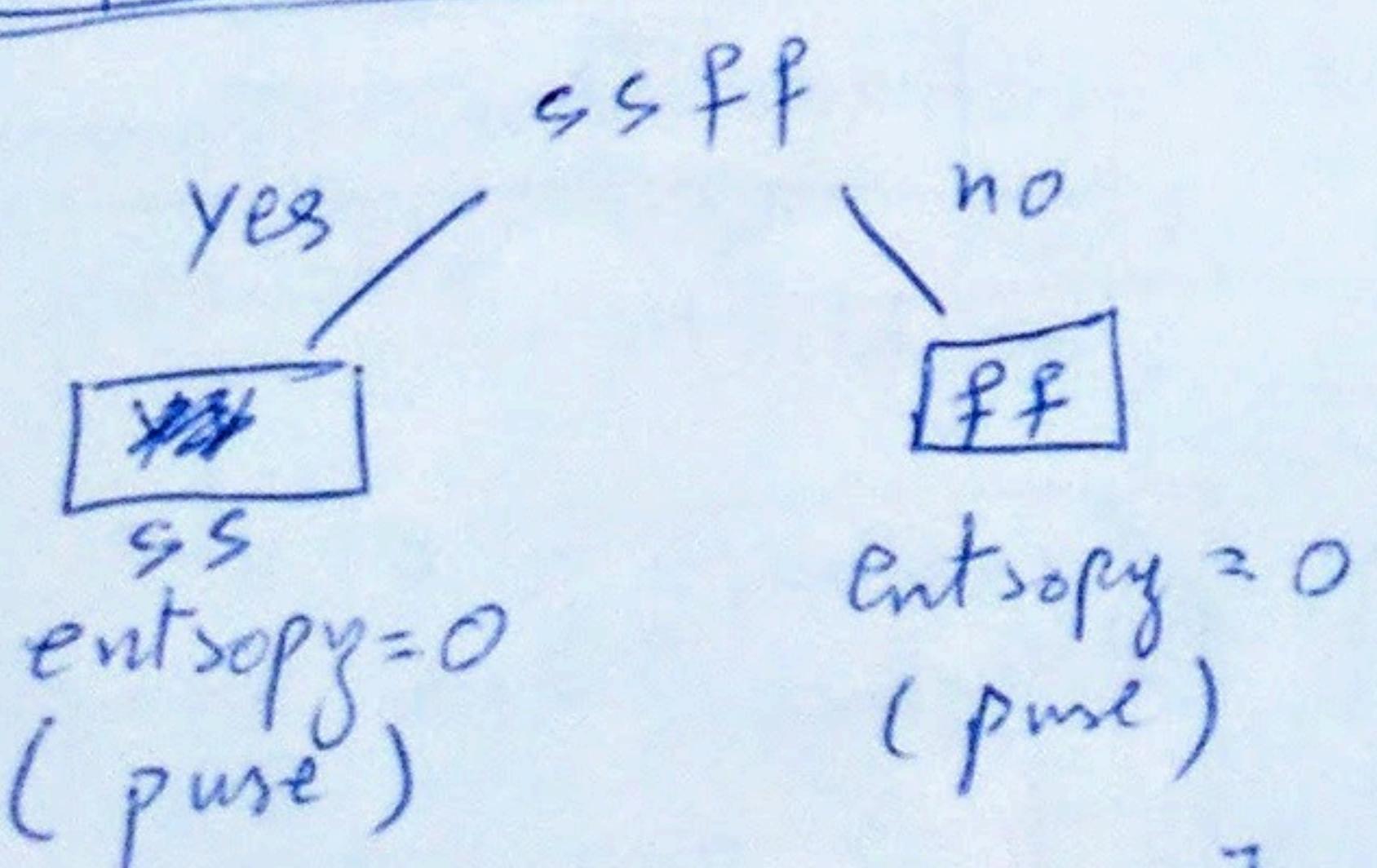
$$\text{And } I.G_I = 1.0 - \left[\frac{2}{3} (1.0) + \frac{1}{3} (1.0) \right]$$

weighted sum of entropy children

$$\boxed{I.G_I = 0}$$

we don't gain any information by splitting data in bumpiness for this example.

Consider Speed Limit:

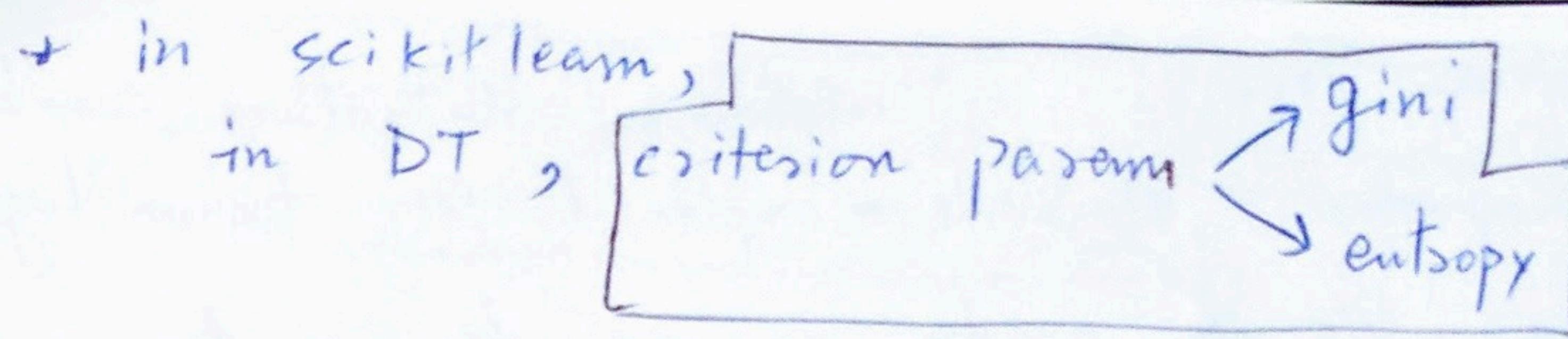


So:

$$I.G_I = 1.0 - \left[\frac{2}{3} (0) + \frac{1}{3} (0) \right]$$

We get max. info by splitting data based on Speed Limit

$$\boxed{I.G_I = 1.0}$$



(9)

(no Info) Worst $0 \leq \text{Inform. Gini} \leq 1$ Best (full Info)

(pure) Best $0 \leq \text{entropy} \leq 1$ Worst (Impure)

Information Gini $\propto \frac{1}{\text{entropy}}$

pure classification \Rightarrow Impure classif.

Class A	Class B
Apple Apple	Banana Banana

- * more Info
- * low entropy

Class A	Class B
Apple Banana	Apple Banana

- * no Info
- * high entropy

* for impure classification you get no information

* for pure classification, you get more information

Bias-Variance Dilemma:

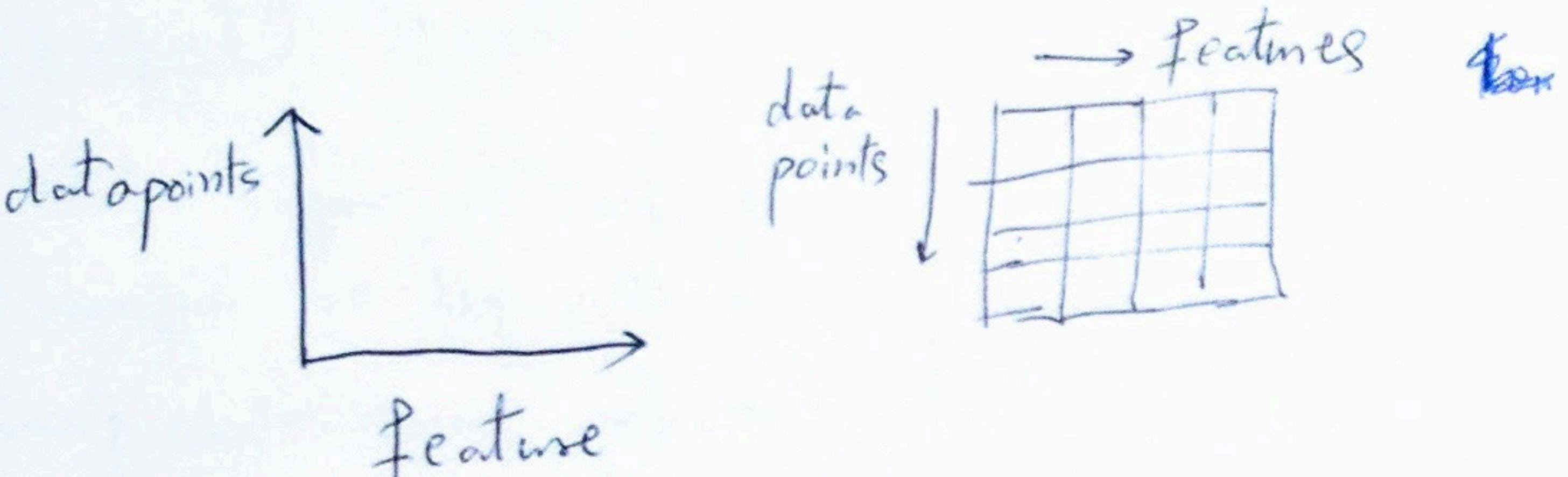
* Bias classifier is the one on which there is no effect of data. So whichever you train it, it'll never do things differently.

* Variance classifier is one which is very sensitive of the data, it can only replicate the stuff it has seen before. It'll react poorly for new data as it has very little bias to generalize itself to new stuff.

* So we need something in middle so it can learn from (good variance) but could also be generalized (good bias).

⇒ Tradeoffs with Decision Tree.

- * Can be easily interpreted, whitebox - explained using boolean logic
- * Prone to overfitting b/c of overly-complex decision tree with lots of features.
- * Decision Tree classifier can be used to new classifier (~~ensemble~~
ensemble classifier)



of feature & complex fit
assuming sample size is constant

⇒ AdaBoost.

- * Adaptive boost + meta algorithm.
- * ensemble → combines (weighted sum) several weak learners (decision trees) to form a robust classifier.

* adaptive → subsequent weak learners are favored over those instances misclassified by previous classifiers.

* Hughes phenomenon → with fixed number of training samples, prediction power & dimensions

Curve & Dimension
relation

→ means if you have lots of features but very small training sample, then prediction is weak as features are just increasing the dimension.

of Samples ↓ # of Features ↓

→ So, rows & cols should both grow to have good prediction.

* pruning → process of removing poorly performing weak classifiers

overfitted situations can be

+ Bias & ~~Data Dependence~~; Variance & Data Dependence

16

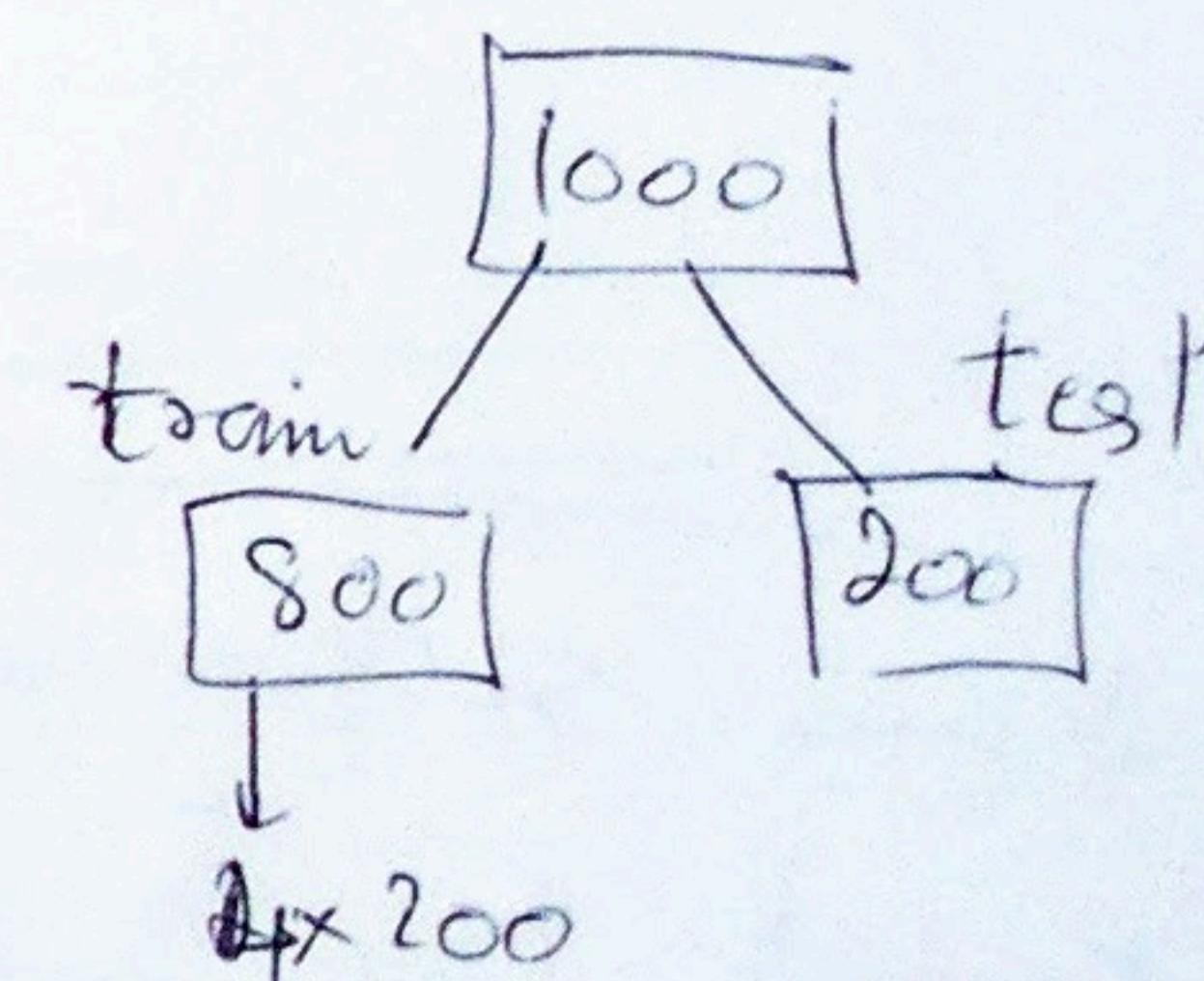
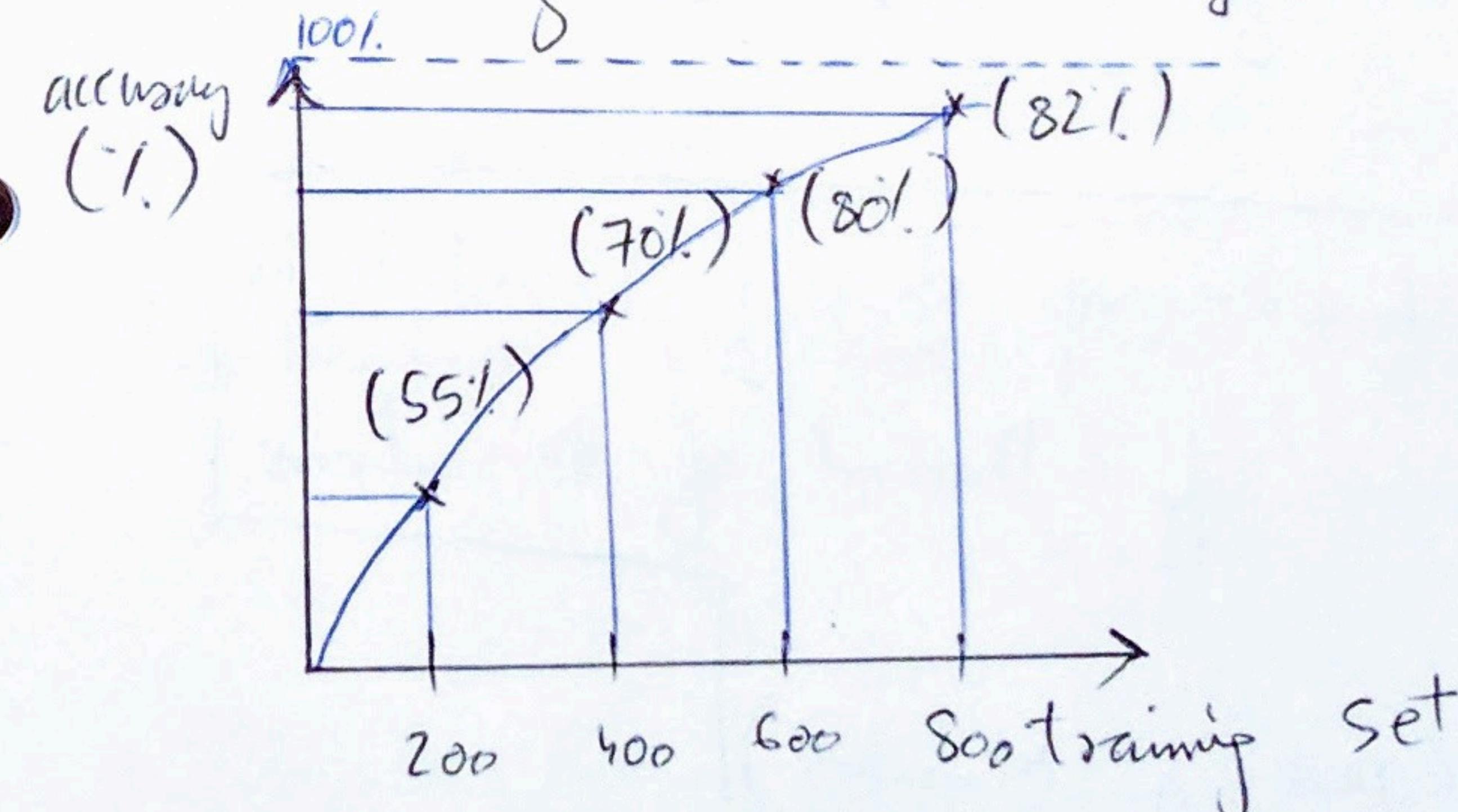
⇒ Enron fraud:

- * Enron was energy trading company with 101 bn \$ revenue.
- * went bankrupt b/c of clever economic & financial practices.
- * 10000 of emails (real) were released as corpus & we will use these emails to train our classifiers & get insight.

⇒ Datasets & Questions:

- * Person of Interest (POI):
 - indicted
 - settle w/o admitting guilty.
 - testified for immunity.

* Accuracy ⇔ Training set size



* Accuracy & Training set size

- * But we shouldn't have lots of features b/c it'll create Hughes phenomenon that's too much features w/o increasing the training set so just increasing dimensions & reducing prediction accuracy.

* More Data > Fine tuned Algorithm

* return q

Type of Data:

Numerical - numerical values (feature → number)

Categorical - limited # of discrete value (feature → category)

time series - temporal value (date timestamp)
↳ features

text - words → feature
↳ we convert words feature to number
using Bog of words formula etc.

- + salary → numerical feature
- + job title → categorical feature
- + timestamp on email → time series feature
- + email contents → text feature
- + # of emails sent → numerical feature
- + to/from field → can be categorical feature but fits better as text feature.
where text is email addresses.
text wrangling tools can be used.

x Jeffrey Skilling	→ CEO of Enron during fraud time
x Kenneth Lay	→ Chairman of Board of Directors.
x Andrew Fastow	→ CFO
+ Unfilled feature	→ NaN (Not a Number)

# of people	→ features			
	POI	Salary	expenses	...

Salary	to-messages	deferral-payments	total-payments	exercised-stock-options
bonus	restricted-stock	shared-receipt with POI	restricted-stock deferred	total-stock-value
expenses	loan-advances	from-messages	other	from-this-person-to-poi
poi	director-fees	defered-income	long-term-incentive	email
from-poi-to-this-person				

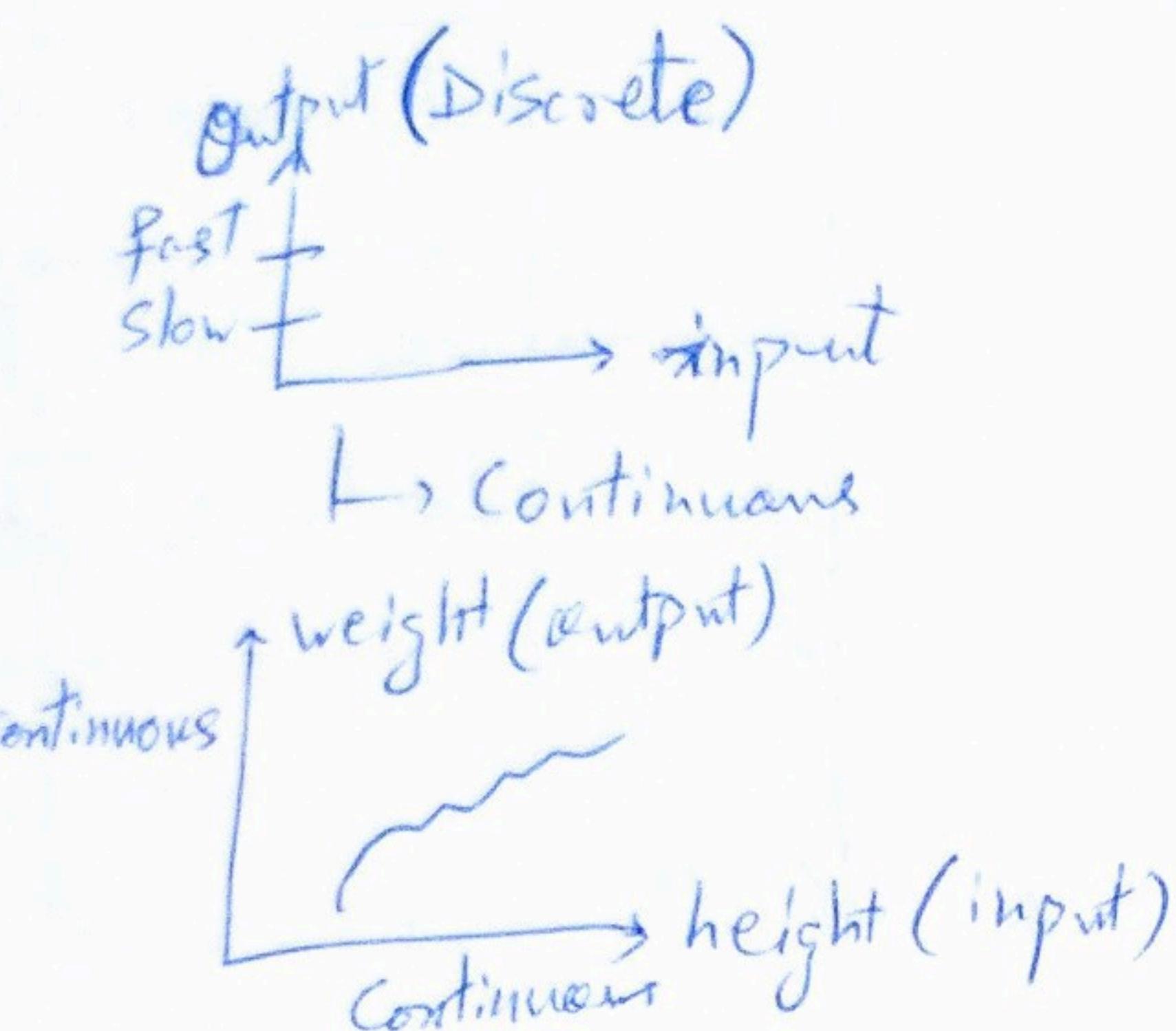
- + most non-POIs come from financial statements with complete data.
 - + POIs were added from Enron data with very incomplete data.
 - + classifier learner can be tricked into thinking that POIs are those which have several NaN fields — which is obviously not the case.
 - * So we should be very careful when introducing new feature while using data from several sources.
 - + 1 solution is to use features common to all sources.
- feature selection from Multiple sources.

⇒ Regression

+ Continuous Supervised learning:

↳ output.

* input is normally continuous but o/p can be discrete (binary) or continuous.



* Age → continuous

* weather → discrete (sunny or rainy)

* But most things we regard as discrete could be modelled as continuous. e.g. weather could mean the amount of light and this is continuous.

* person wrote email → discrete

* even though we can have thousands of authors but they are not continuous in time So #56 & #57 has no relation in time.

* phone number → discrete

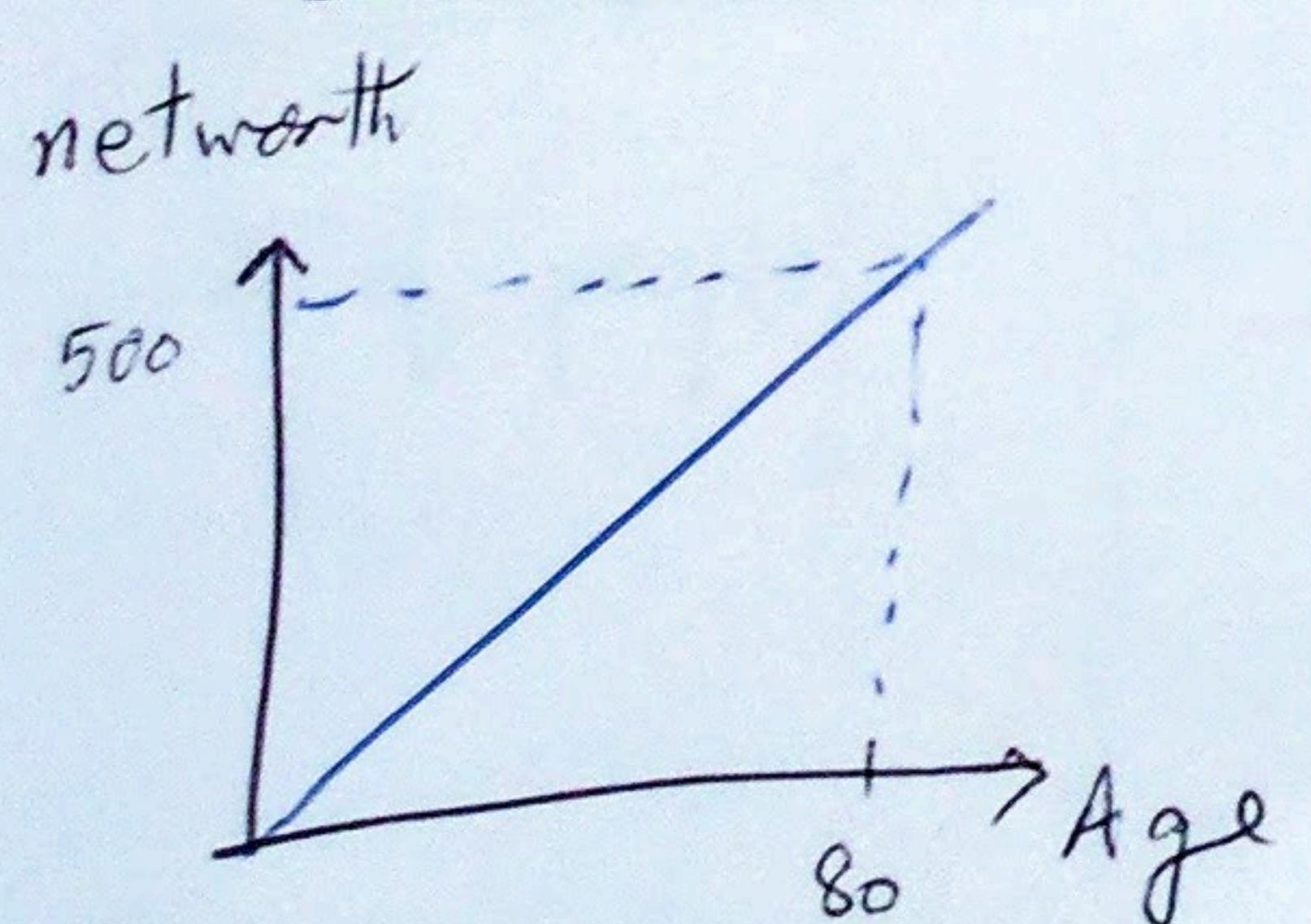
* income → continuous

Without order → discrete
With order → continuous

Regression → predicting
relation among data & finding best fit line.

+ Discrete
fast/slow

+ Continuous
speed in mph



$$\text{net worth} = m(\text{Age}) + b$$

$$m = \frac{500 - 0}{80 - 0} = 0.25.$$

$$b = 0$$

$$\text{So: net worth} = 0.25(\text{Age})$$

y-intercept → height when line intersects y-axis

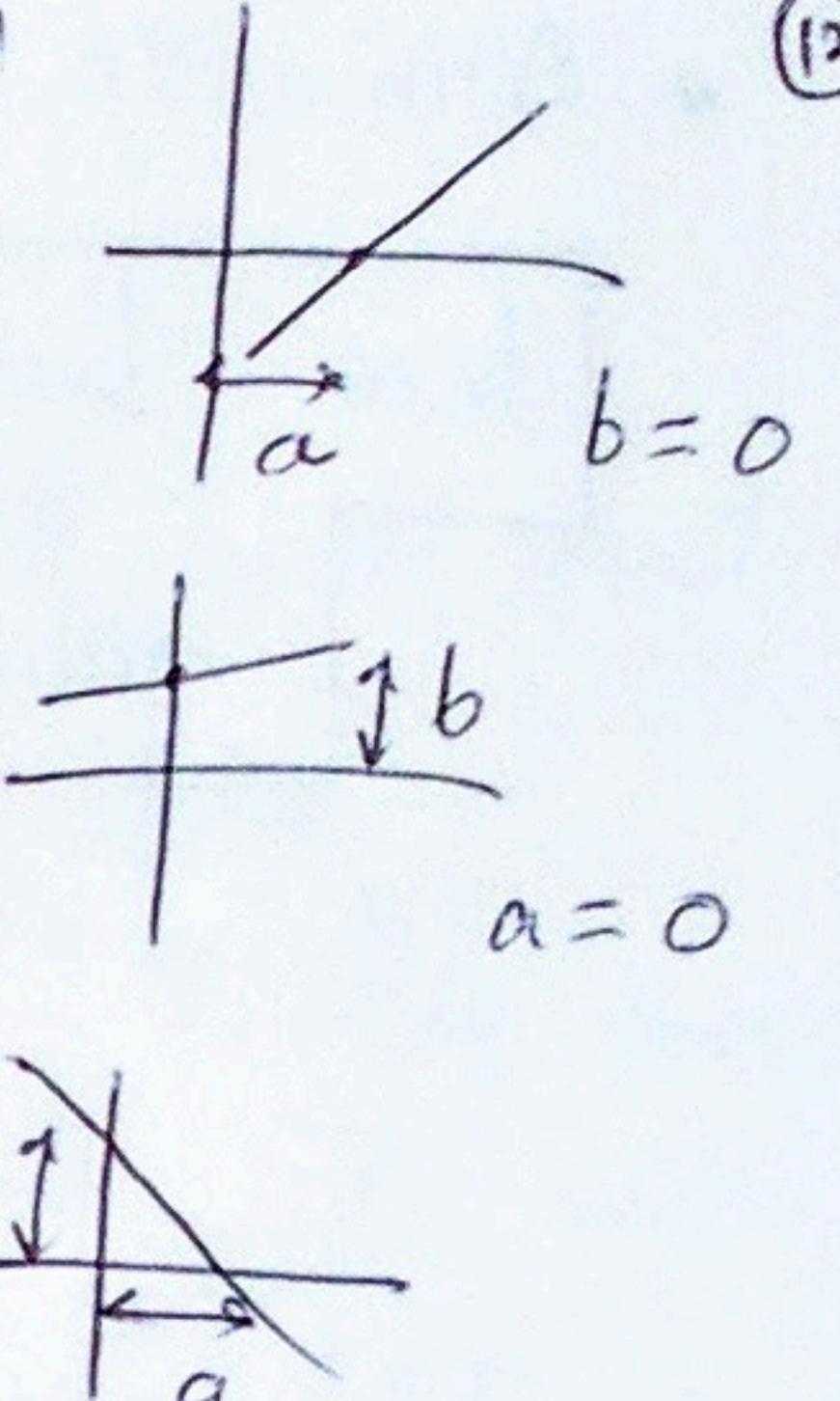
$$\text{Target (out)} = \text{Slope (Input)} + \text{y-intercept}$$

* Regression → predict output from input using:

$$y = mx + b$$

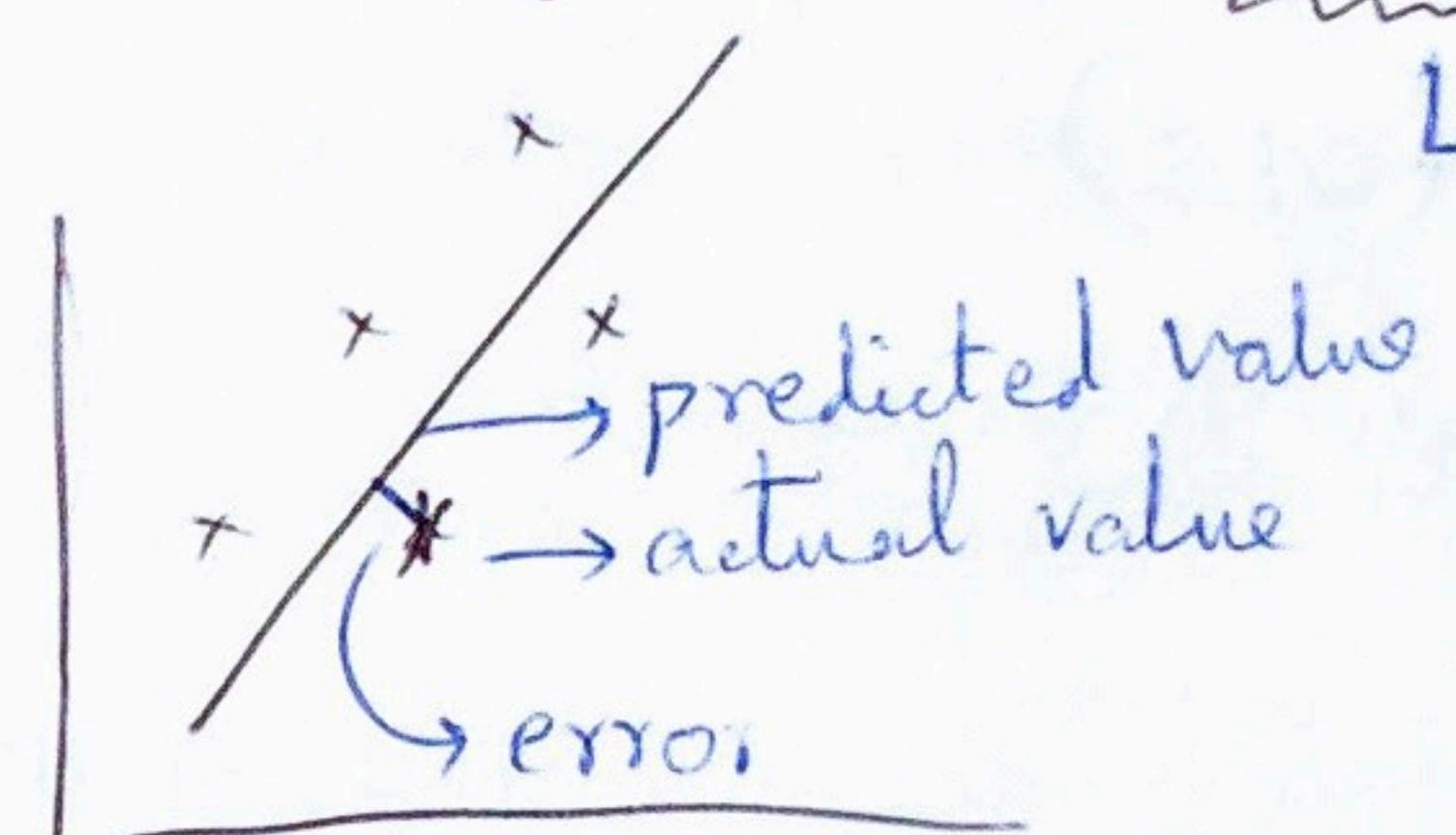
↳ Linear Regression

* We need to find m & b



Linear Regression Errors:

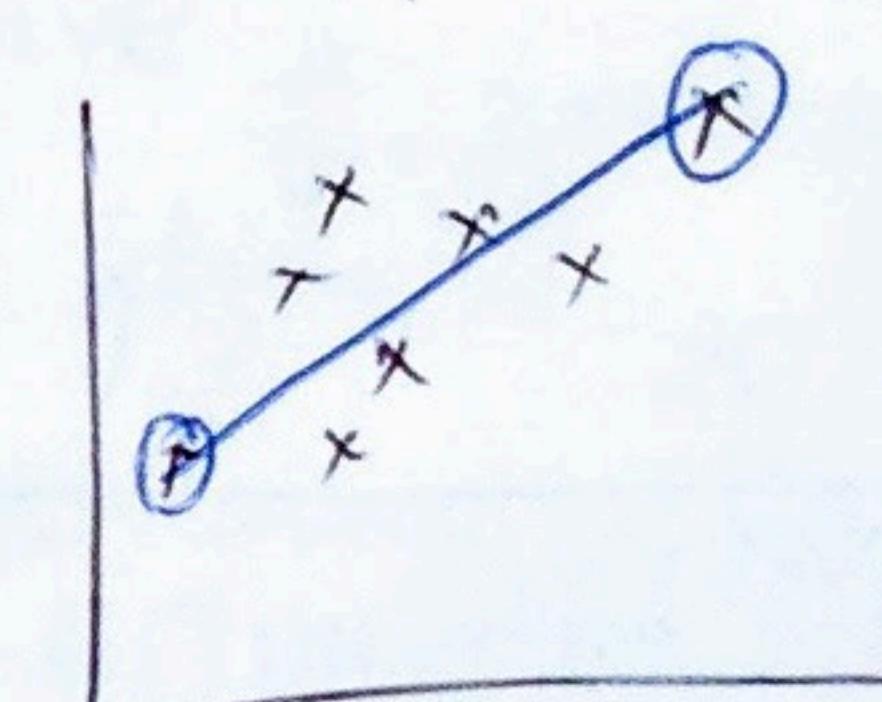
↳ actual - predicted



error → -ve means actual point is below predicted point.

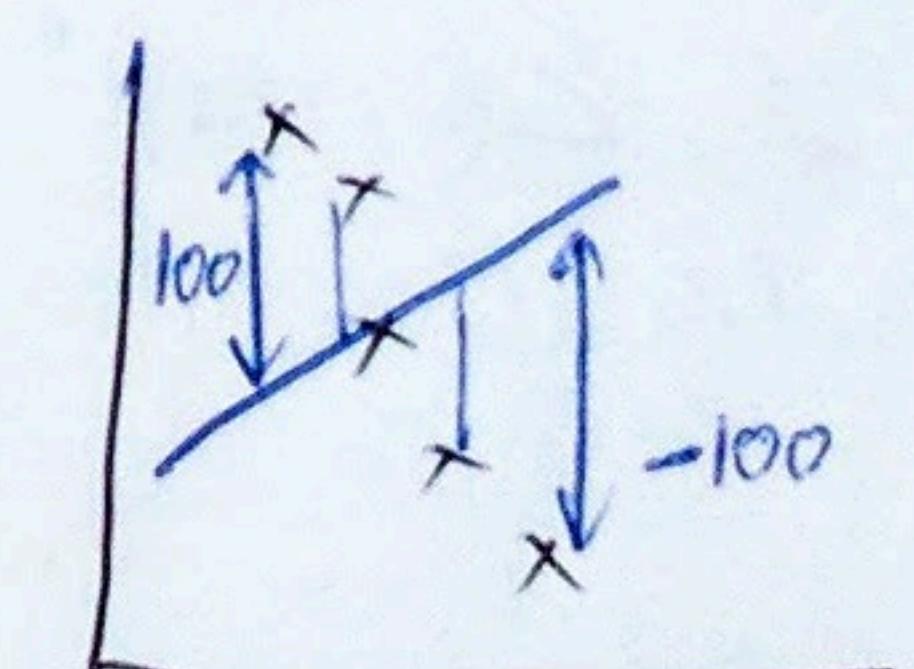
error → +ve means actual point is above predicted point.

* error on first & last data point:



Not a good fit

* Error on all data points



Not a good fit

* both errors cancel each other out & result is straight line.

Error = 0

even its not fitting at all.

* Minimize Sum of Square Error — $\sum(\text{error})^2$:

[the best regression is:]

$$\underset{\text{all training points}}{\text{minimizes } \sum (\text{actual} - \text{prediction})^2}$$

↑ prediction from regression
↑ training points
 $y = mx + b$

* Methods to minimize $\sum(\text{actual} - \text{prediction})^2$ OR SSE

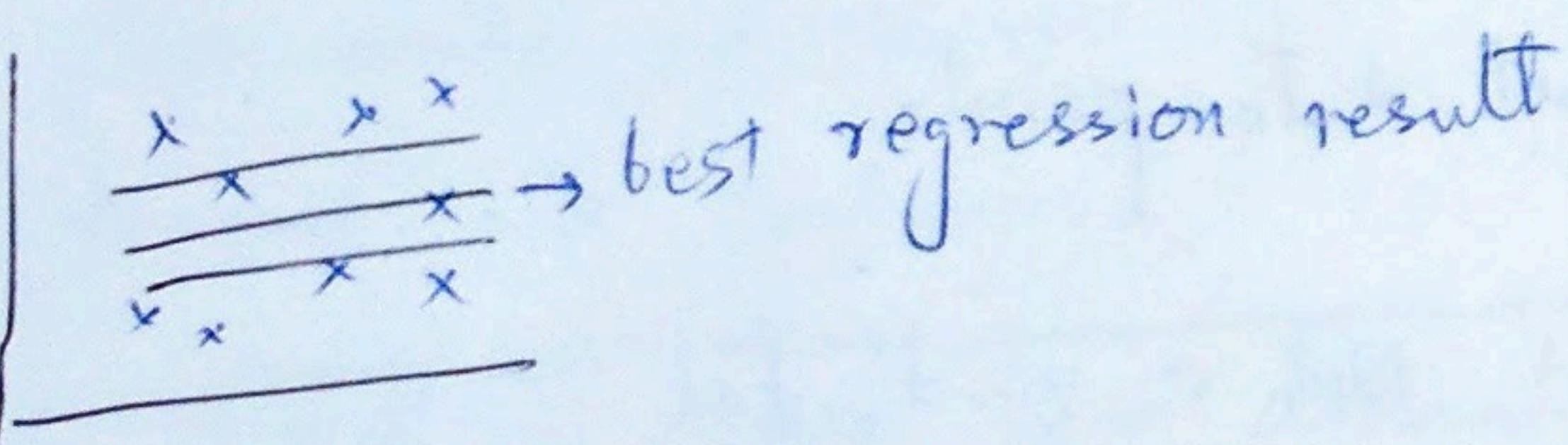
① Ordinary Least Square (OLS):
→ used in sklearn Regression.

② Gradient Descent:

$$y_i = y_i + \alpha (x_i - y_i)$$

$$y_i = y_i + \beta (f_{i-1} + f_i - 2(y_i))$$

$y_i \rightarrow \text{new}$
 $x_i \rightarrow \text{old}$
 $\alpha \rightarrow \text{data smooth}$
 $\beta \rightarrow \text{weight smooth}$



in terms of SVM, it has large margin on both sides.

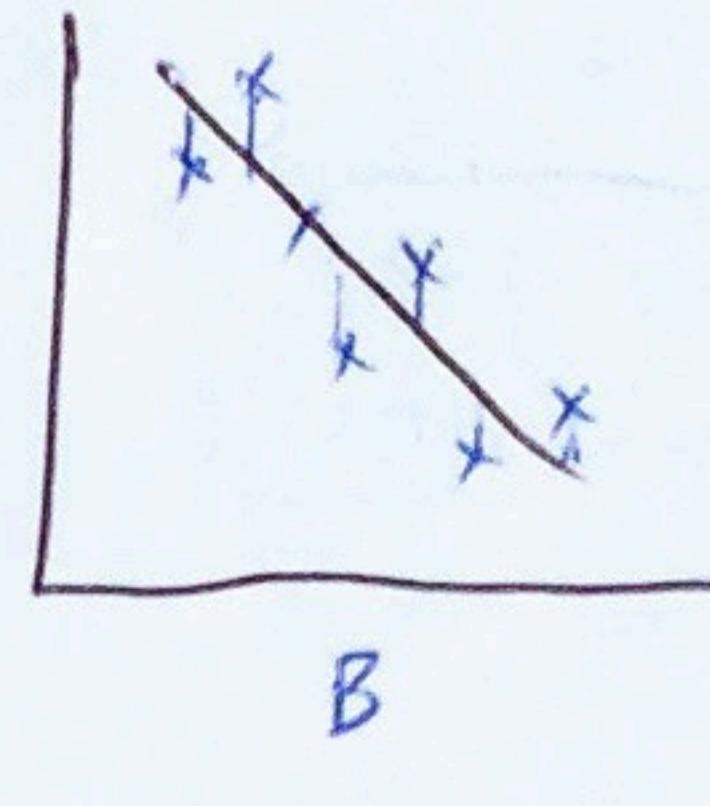
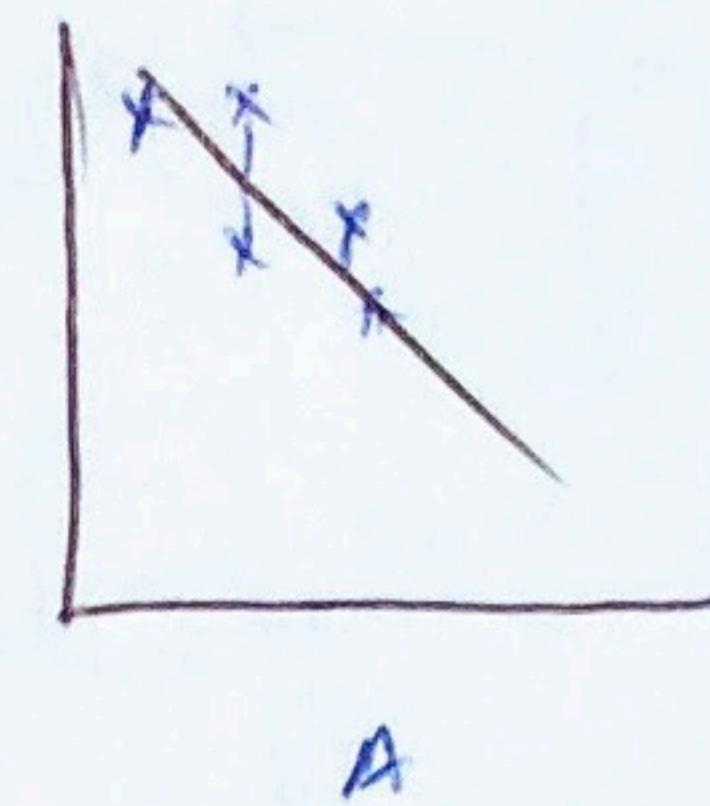
There can be multiple lines that minimize $\sum(\text{error})^2$
but only one line will minimize $\sum(\text{error})^2$

* Another reason for using SSE i.e. $\sum(\text{error})^2$ is that its much easier to code mathematically.

⇒ SSE is not perfect:

(13)

[larger SSE → Worst fit]



- * These both have about the same fit
- * But just B has more data points, it will have greater SSE.

In SSE, as you add more data, the SSE normally goes up even though the fit is not necessarily bad.

* SSE gives relation of test labels & predicted labels.

* r^2 gives relation b/w input & output of test data

⇒ r^2 of a regression:

+ how much Δoutput is explained by Δinput

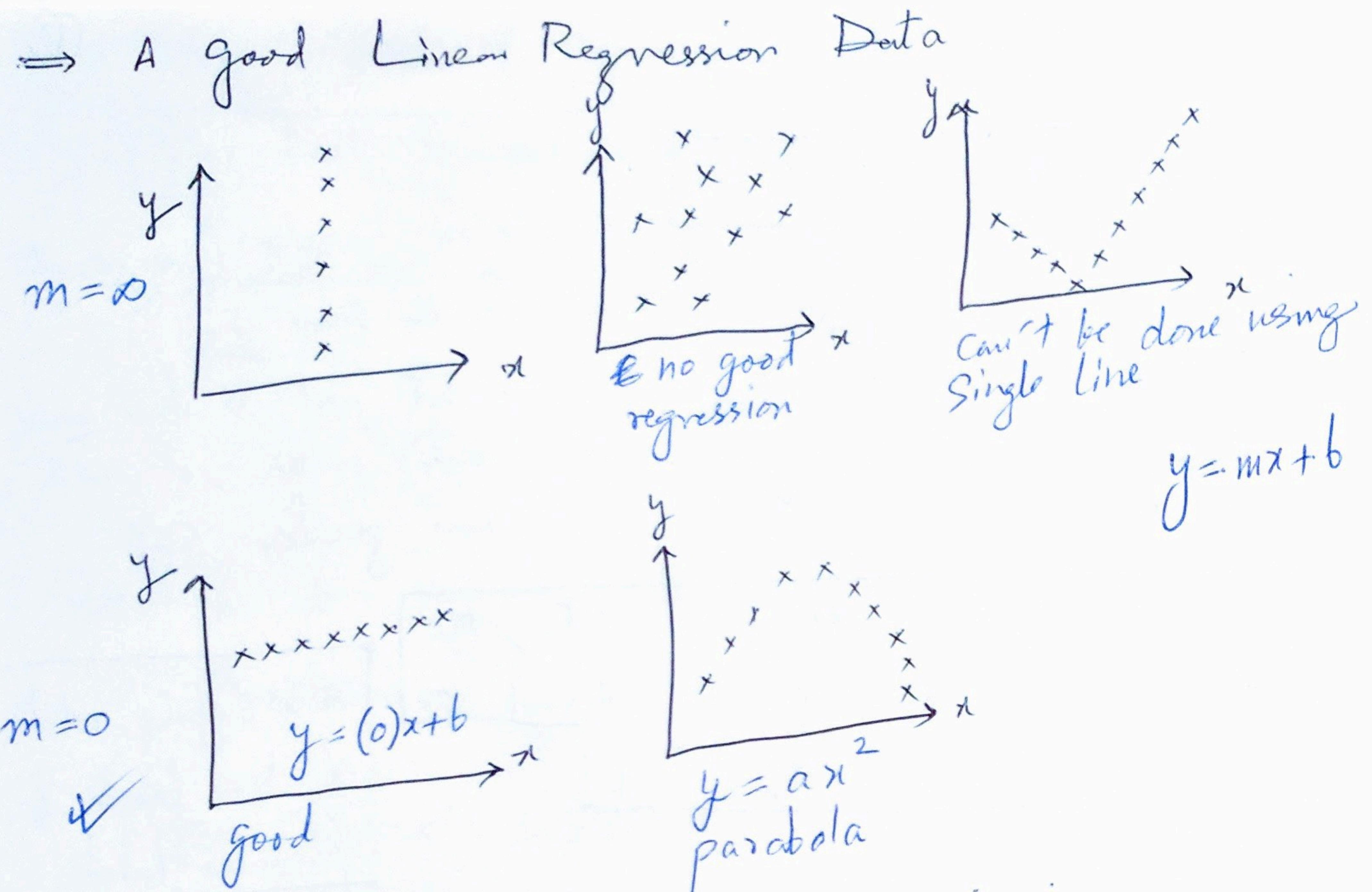
$$0 < r^2 < 1$$

not doing good job of capturing trend in data

good description / relation b/w input & output

+ r^2 can be in -ve when finding fit or regression for significantly different data set.

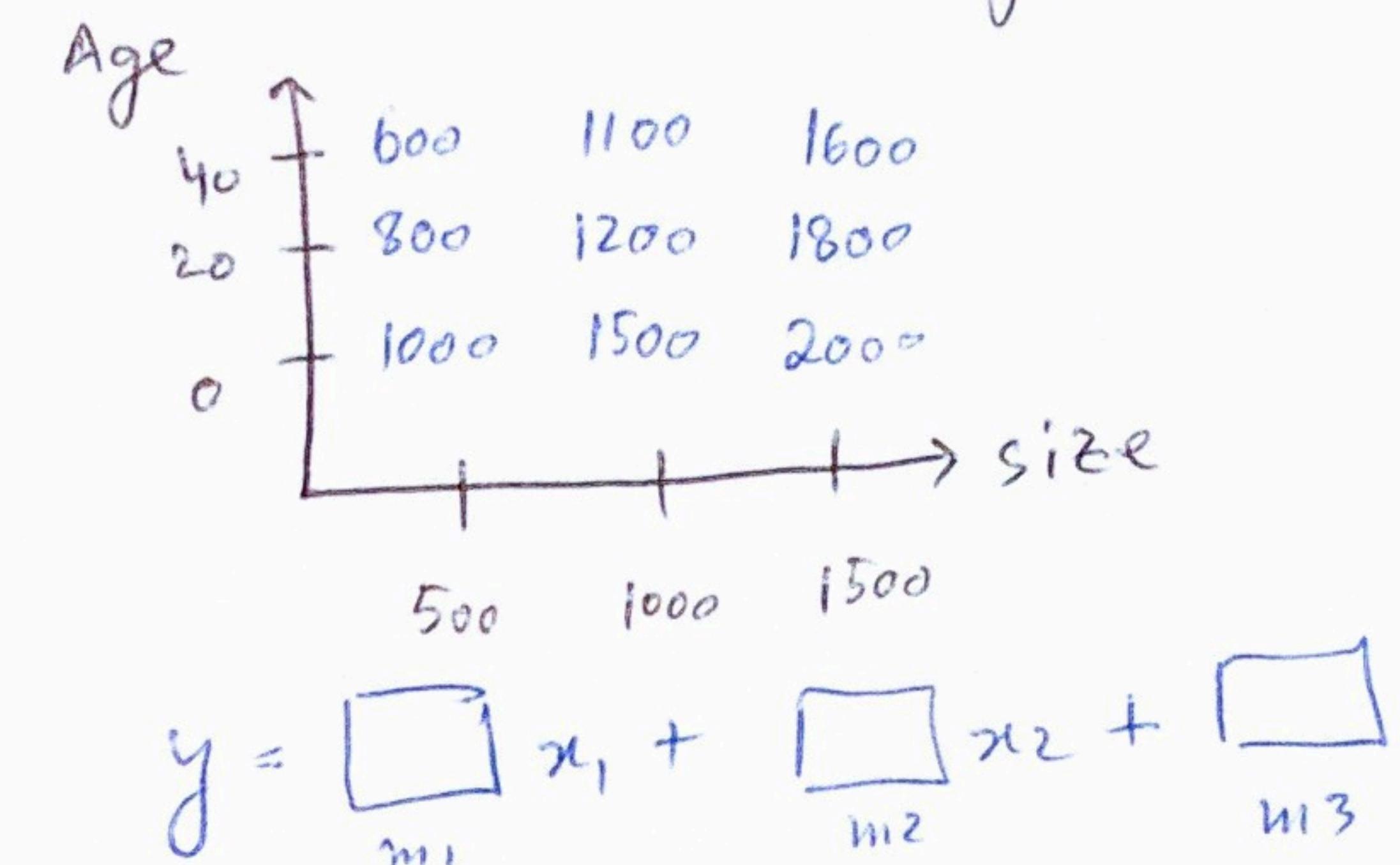
+ r^2 is independent of # of training points.



⇒ Comparing Classification & Regression:

property	Supervised Classification	Regression
output type	Discrete	Continuous (Numbers)
what are we find	Decision Boundary	Best line of Line
evaluation	accuracy score	* Sum of Squared Error * R Square * gradient descent

⇒ Multi Variate Regression:



Sol.

$$m_1 = \frac{1500 - 1000}{1000 - 500} = \boxed{1}$$

$$m_2 = \frac{1000 - 800}{0 - 20} = \boxed{-10}$$

And:

$$1000 = 1(500) + (-10)(0) + m_3 \\ \boxed{m_3 = 500}$$

* return q

height
weight
edu: } \rightarrow you

14

- * If outliers lie in training data, they change slope & intercepts.
- * If outliers lie in testing data, they affect negatively the regression score.
 - + train → for fit.
 - + test → for predict

\Rightarrow Outliers

target

target

feature

target

feature

\uparrow change b/c of outlier

very Large LS E (Least Sq. Error)

- * Outliers can be caused by  Sensor Malfunction
Data Entry Error

- * In robotics & SDC, we usually ignore outliers but in other fields like fraud detection, outliers are very highly looked for

The figure consists of six separate scatter plots arranged horizontally. Each plot has a vertical y-axis and a horizontal x-axis. The data points are represented by 'x' for outliers and 'o' for other data points.

- Plot 1:** Shows one outlier 'x' located at the top left of the plot area.
- Plot 2:** Shows one outlier 'x' located at the bottom left of the plot area.
- Plot 3:** Shows one outlier 'x' located at the top right of the plot area.
- Plot 4:** Shows one outlier 'x' located at the bottom right of the plot area.
- Plot 5:** Shows one outlier 'x' located at the top center of the plot area.
- Plot 6:** Shows one outlier 'x' located at the bottom center of the plot area.

Outliers Removal:

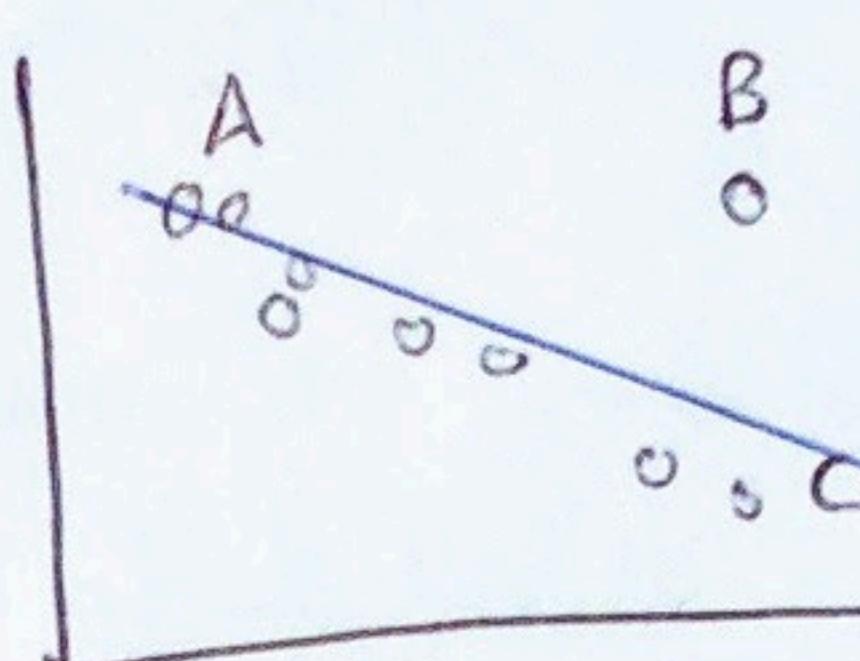
- ① - Train
- ② - Remove 10%
 remove points w/
 largest residual error.
- ③ - Retrain

Train

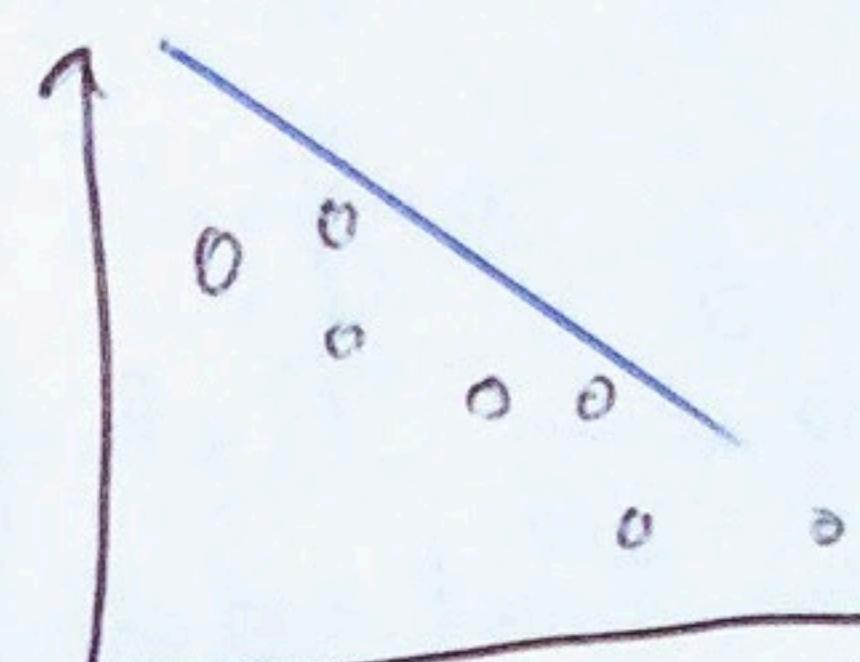
Remove 10%

Retrain

⇒ Residual Errors

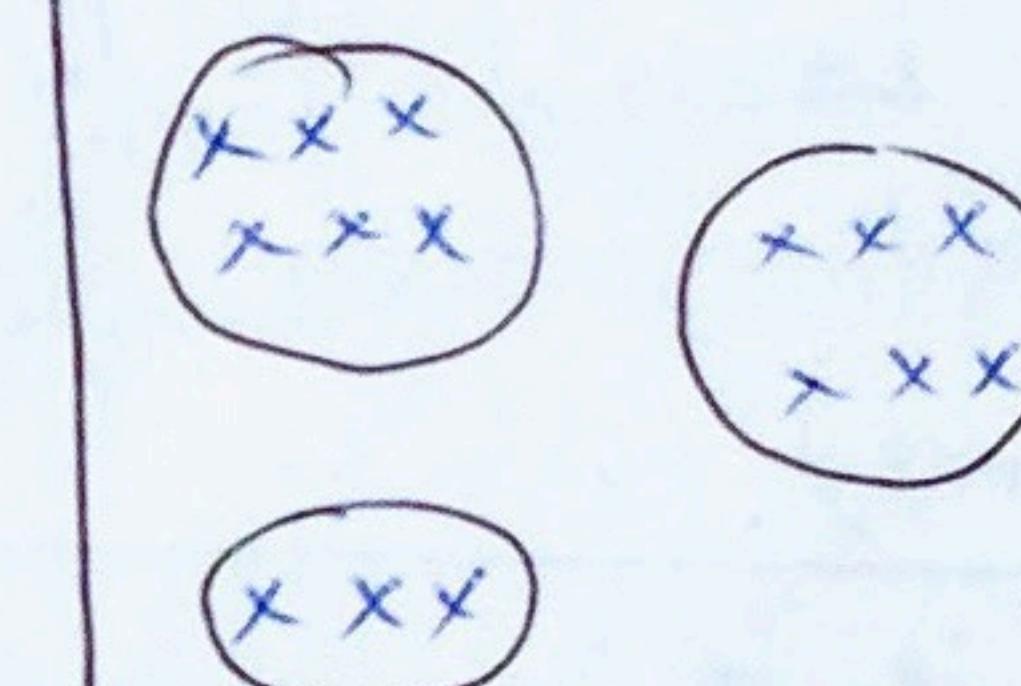


- * D) will have highest Residual Error after fit.

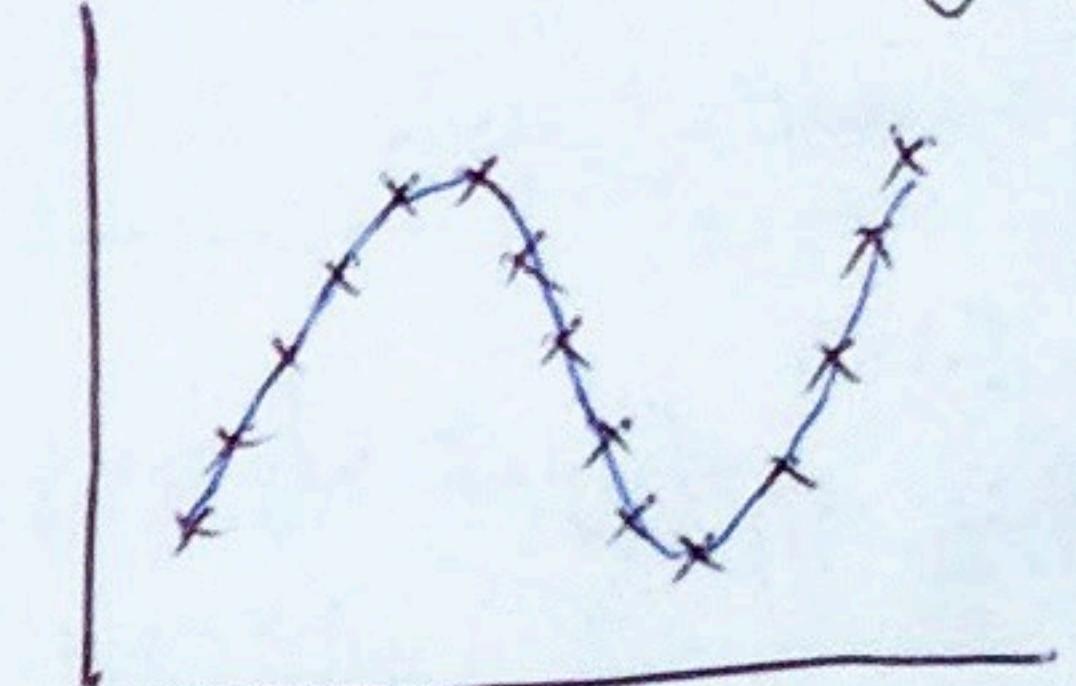


* Steeper now, as we have no outliers

⇒ Unsupervised Learning



Dimension Reduction



* Structuring data without Labels.

- with labeled data,
we don't generally need
Unsupervised

we didn't have classes but after looking at the data we can form the clusters & say them classes.

*K-Means:

R - means
+ Grouping a set of objects in clusters so that objects in same cluster are similar to each other than other clusters.

other than other clusters.

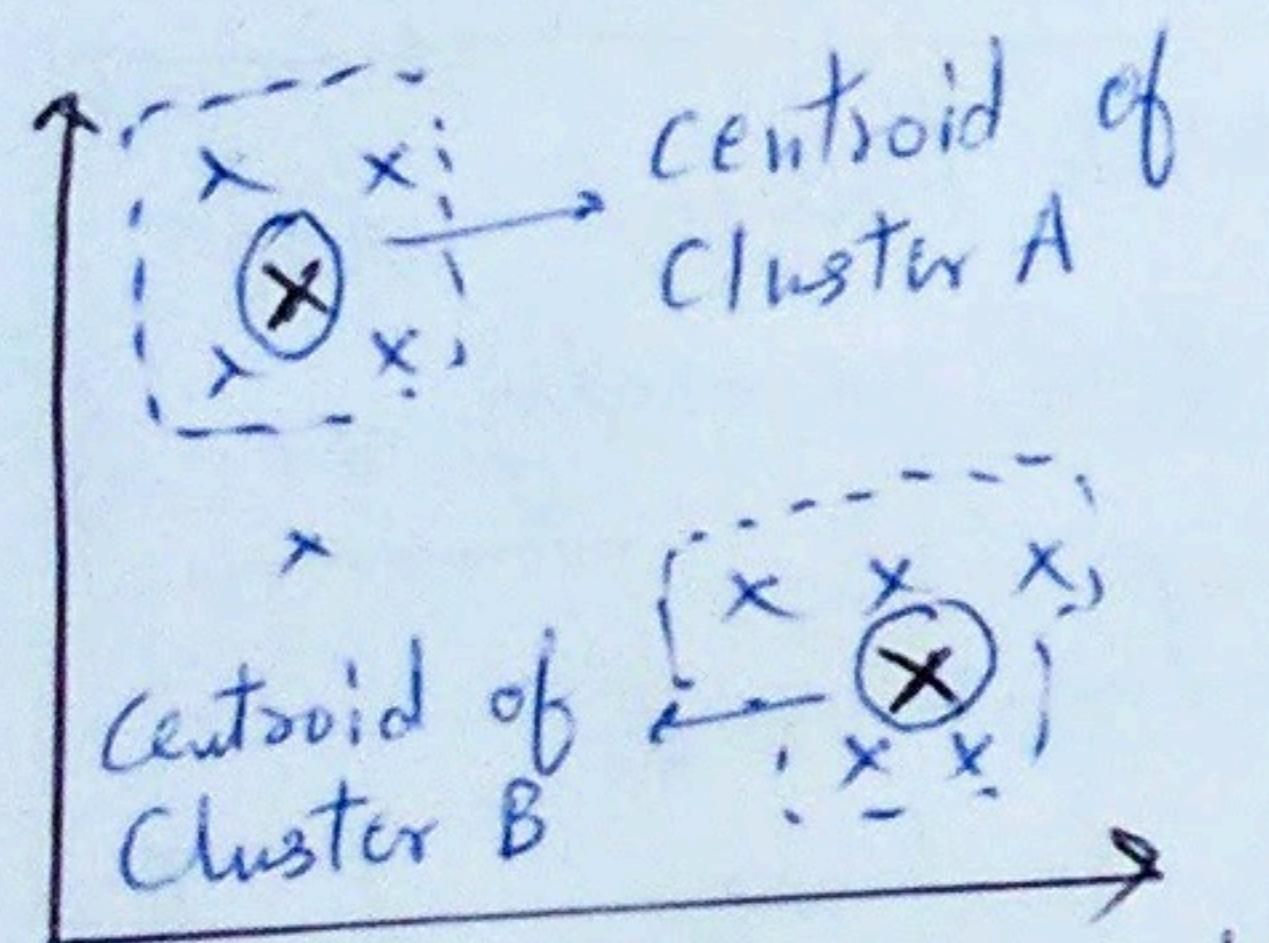
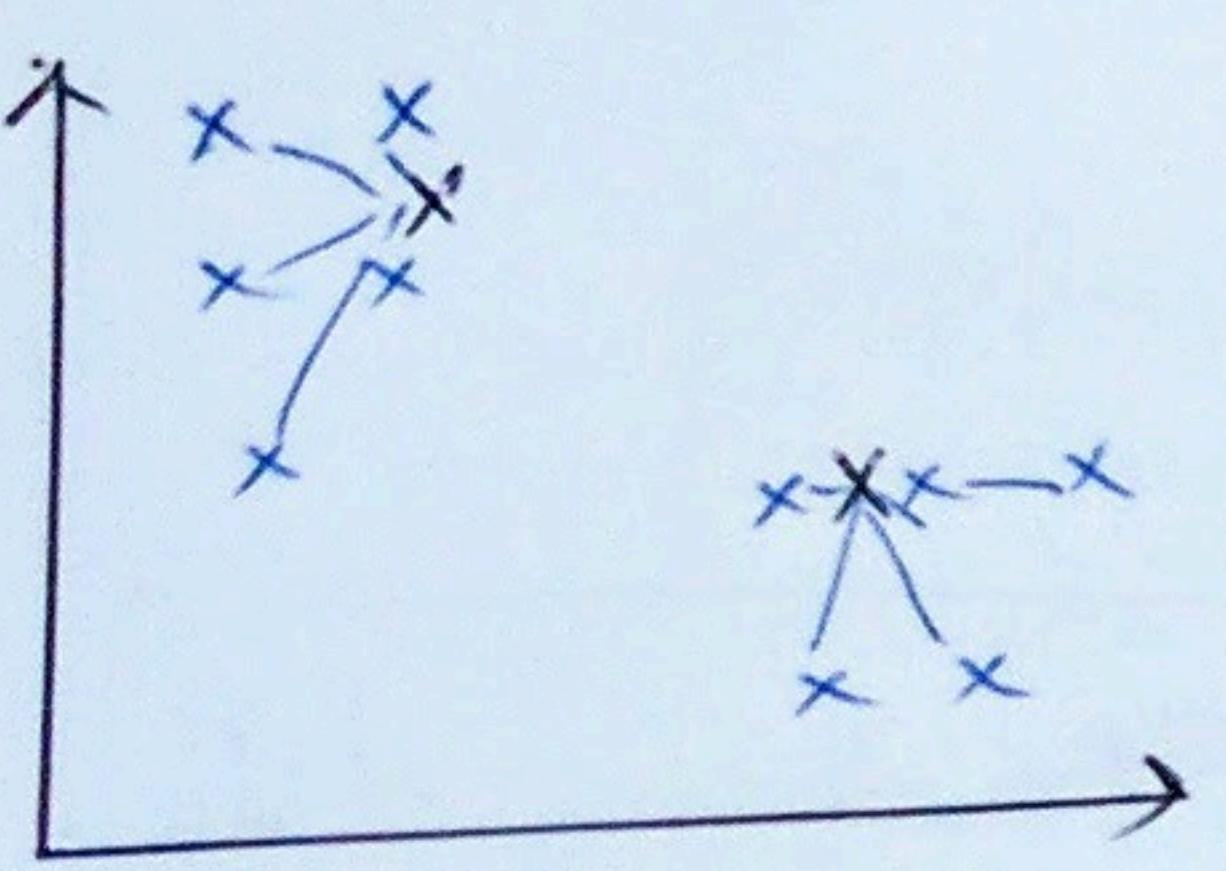
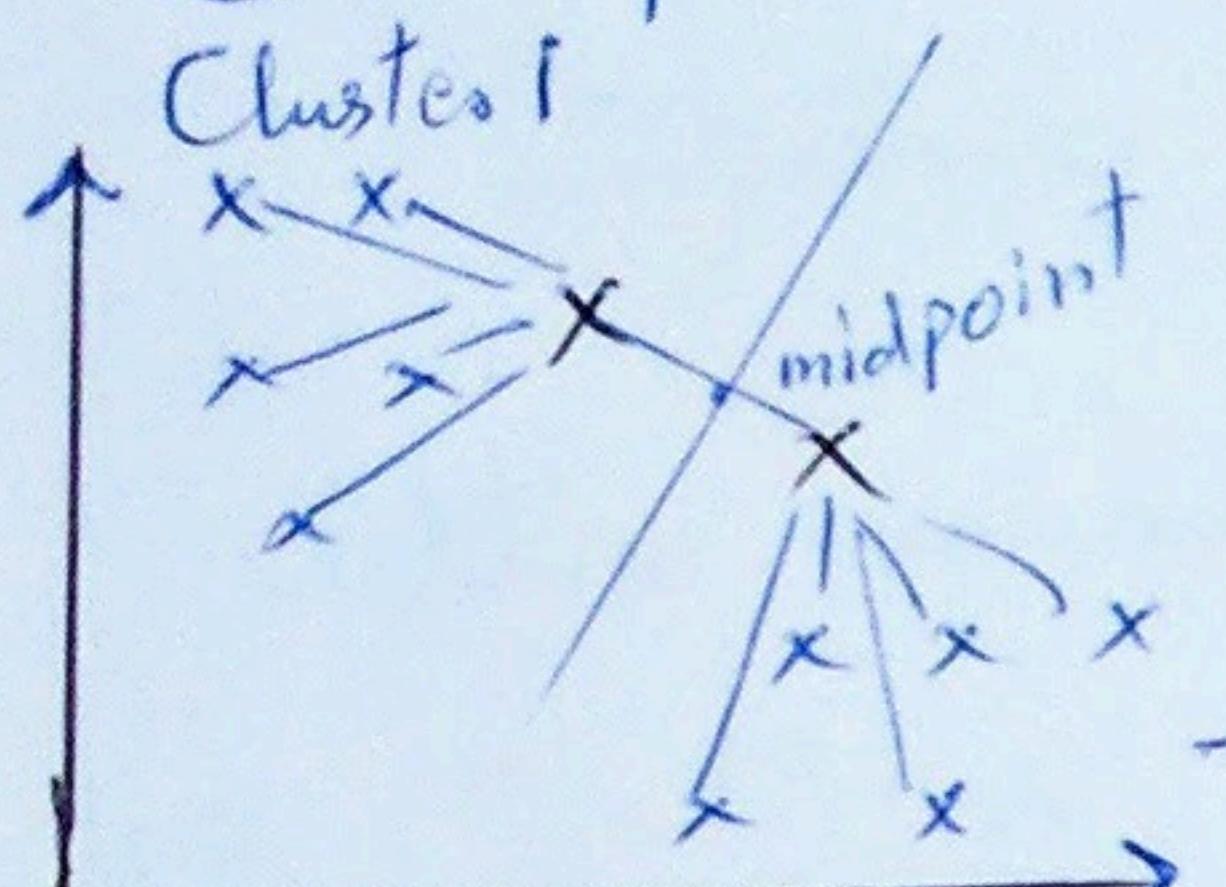
* K-means uses Centroid model for Clustering.

+ Kmeans (for Clustering) is different from
K Nearest Neighbours (for classification)

$$K \leq n \quad k \rightarrow \# \text{ of clusters} \quad n \rightarrow \# \text{ of data samples}$$

Algorithm for KMeans:

- ① - Assignment.
 - ② - Optimize.



- * Initial random positions of centroids are very important as they might end up with totally different clustering.

So: KMeans input \rightarrow n Samples \rightarrow how many clusters
 output \rightarrow K clusters

output \rightarrow K clusters
minimizes WCSS (Within Cluster Sum of Square)

In high dimension, Euclid distance becomes inflated

In high dimension, Euclid distance becomes inflated due to Curse of Dimensions & can be solved by dimension reduction algorithms e.g. PCA before clustering.

→ a.k.a LLOYD's algorithm.

Machine Learning

- * Algorithms which learn
from data.

- + Randomly place 2 points & form the line from their midpt.
 - + By making a model based on input data (Training features & Labels) & predicting output.
 - + with accuracy tested against test data.

- Change position of cluster point -
 - rubber bands tries to be as short as possible.
 - it Random forest &
 - Adaptive Boost corrects the overfitting problem in Decision Tree (makes it sometimes overcomplex)

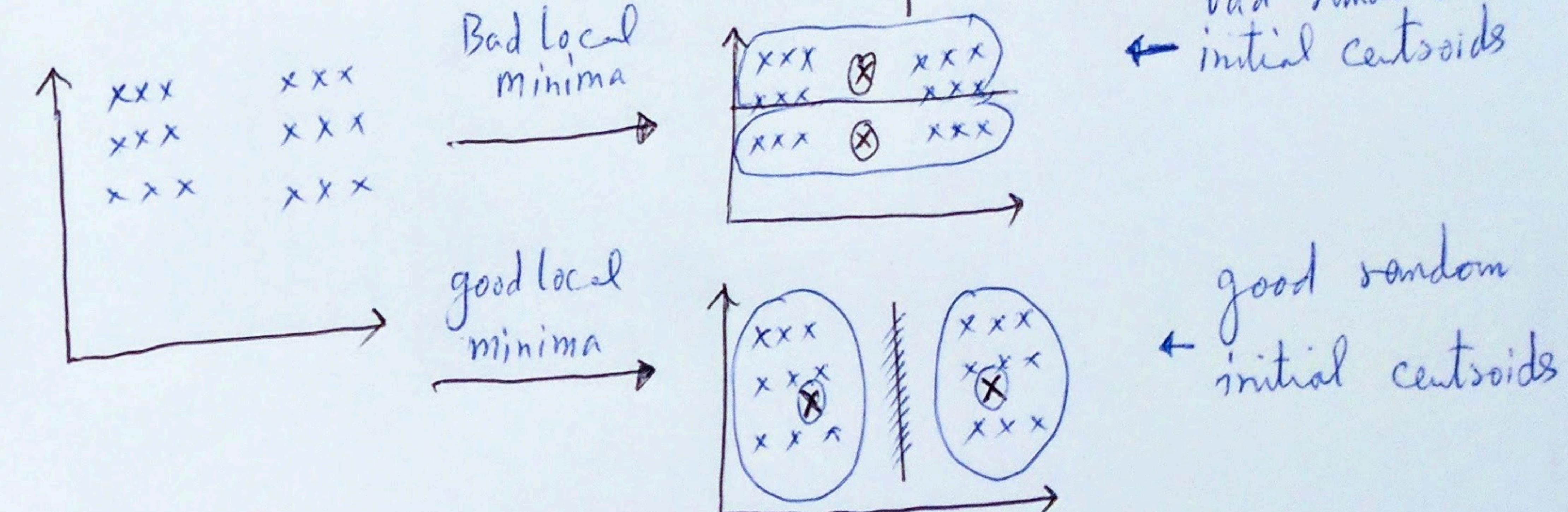
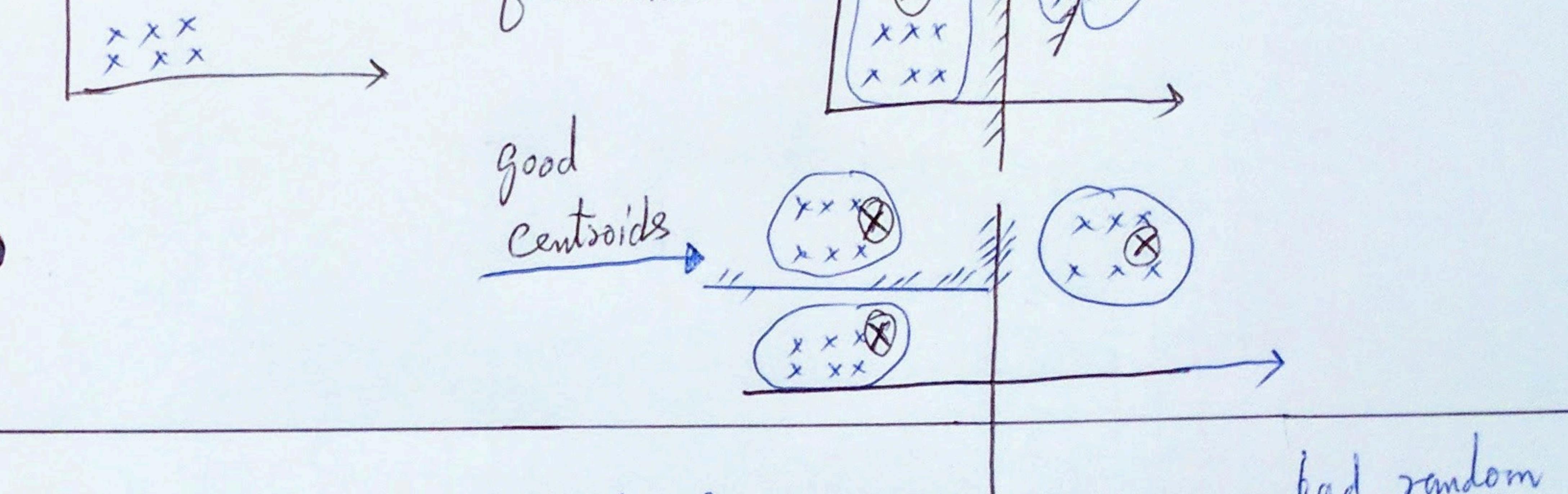
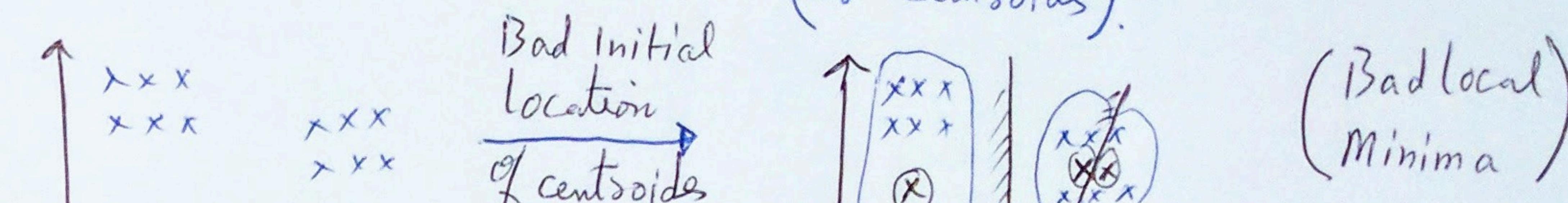
- + iterative optimization.
 - + This is centroid for both clusters.

- I keeps running until
centroid doesn't move very much

→ K-Means Challenges:

16

- * Uses Hill Climbing technique for local minima. So the clusterings depends on when you put the initial cluster centre (or centroids).



good local
minima

Bad location
minima

good random
← initial centroids

← bad random
initial centroids

* return a