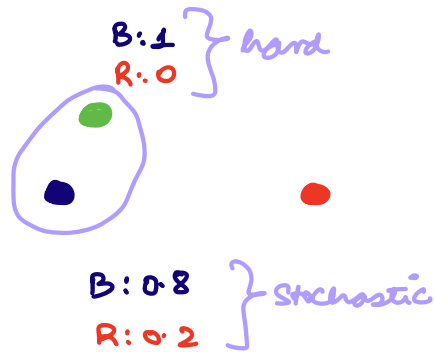
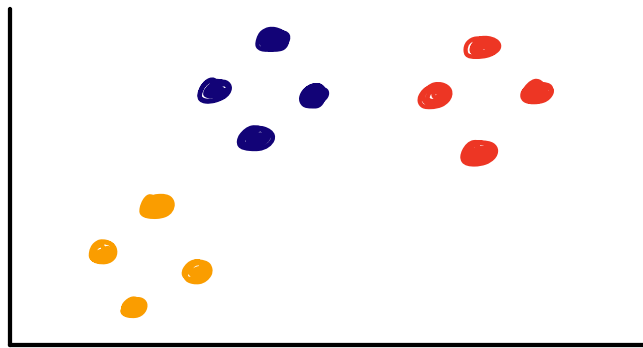


t-SNE  
 ↓  
 ↘ Stochastic Neighbour Embedding  
 ↙ Probabilities  
 † distribution



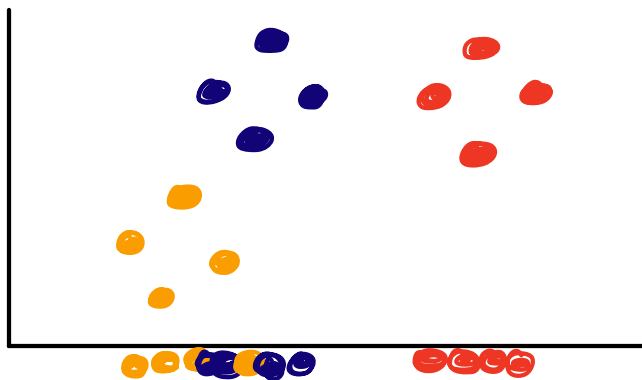
2D Scatter Plot



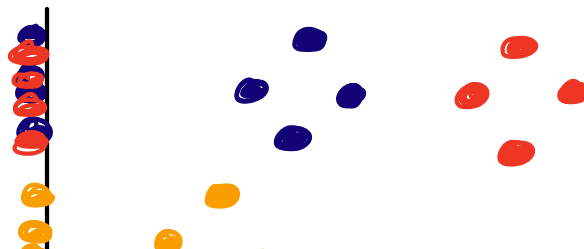
1D number line



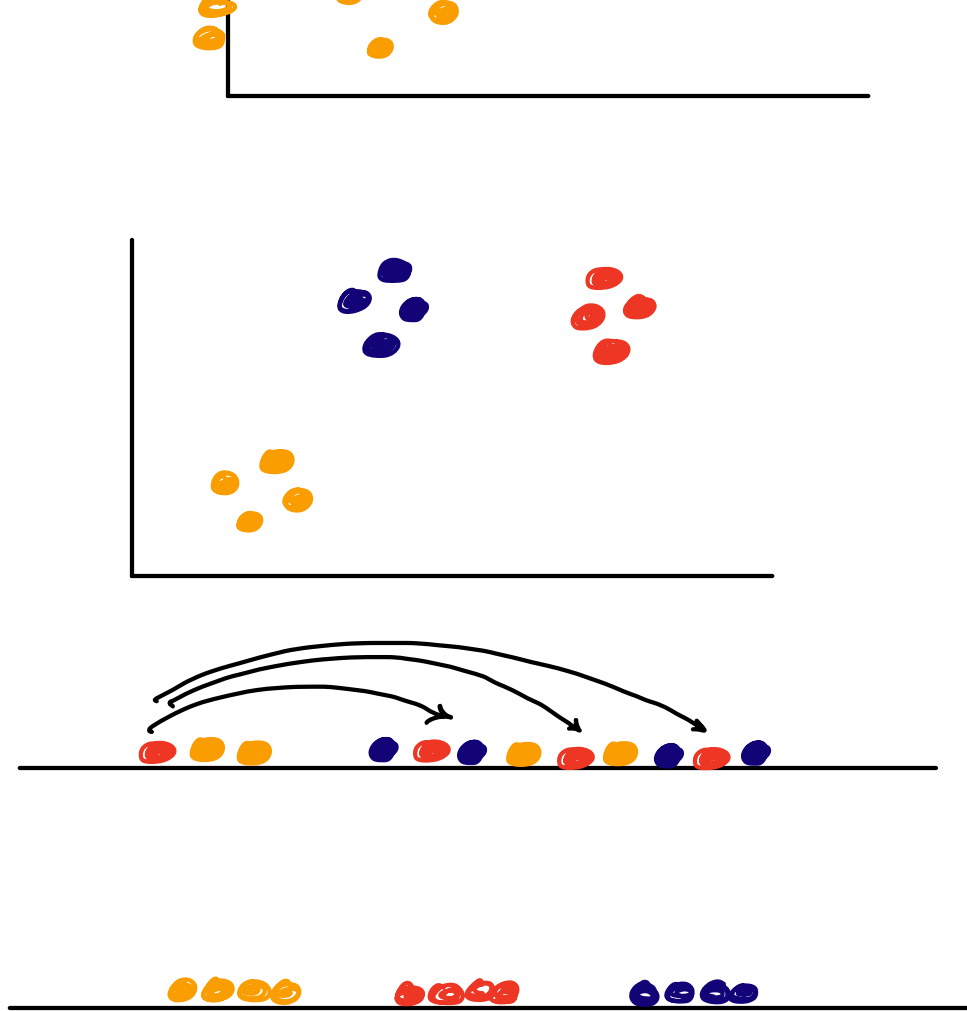
Spread (Projection on x)



Projection on y

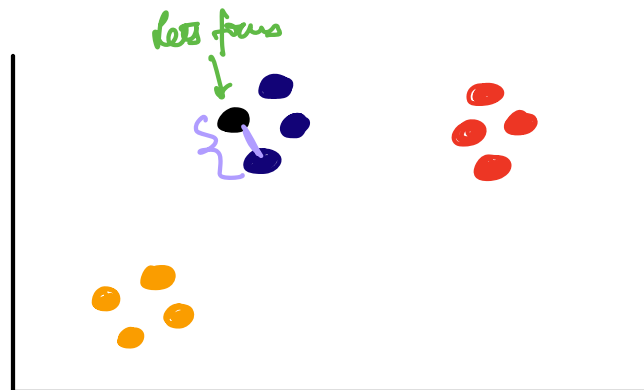


## Intuition

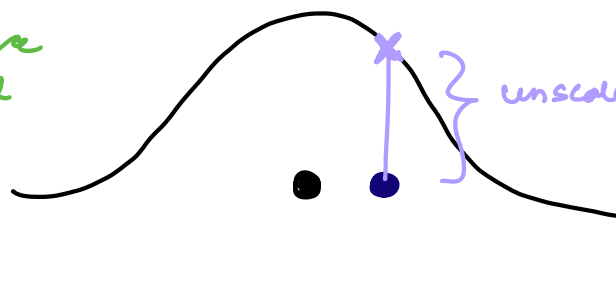


## Step 1:

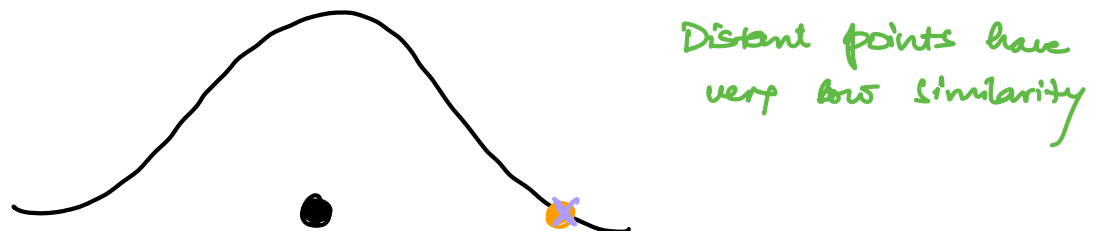
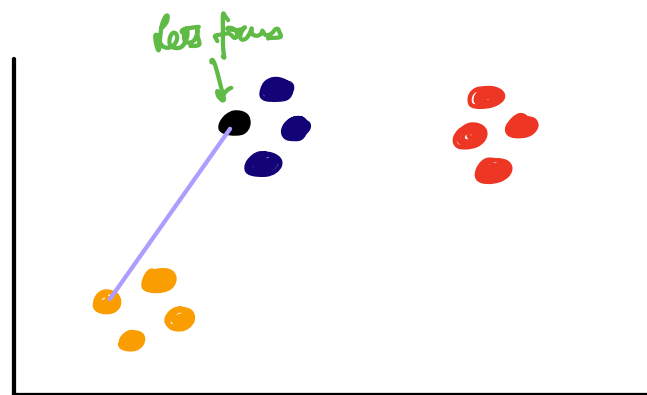
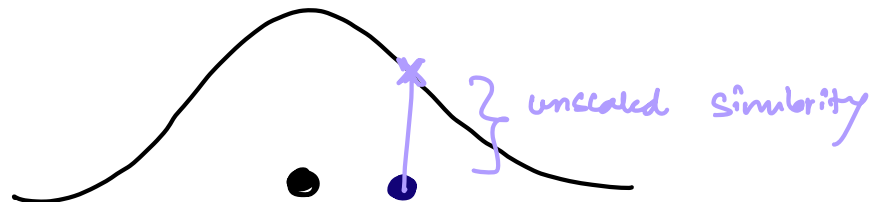
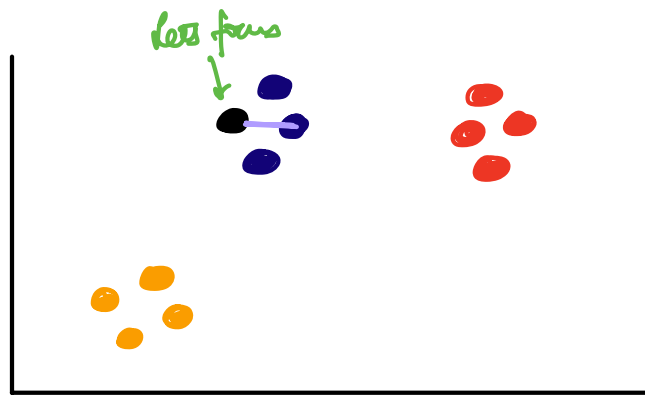
Determine the similarity of all the points in the scatter plot



Plot the distance on a normal curve that is centered on point of interest

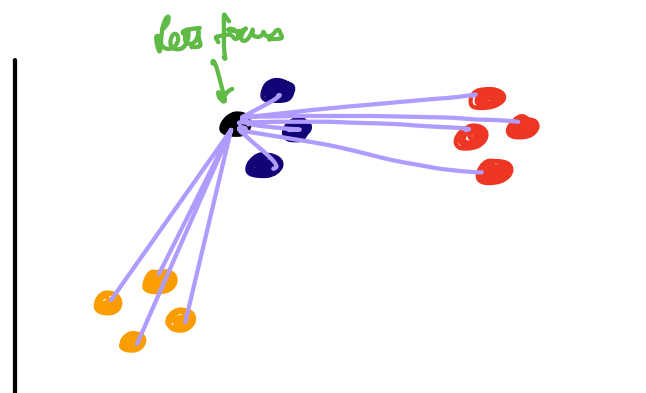


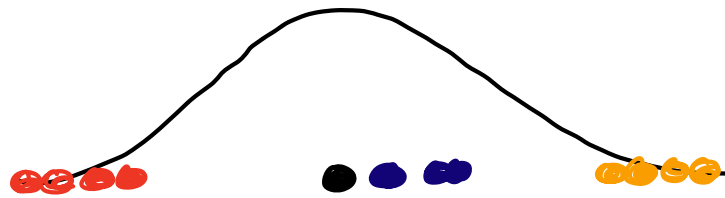
Close points have high similarity



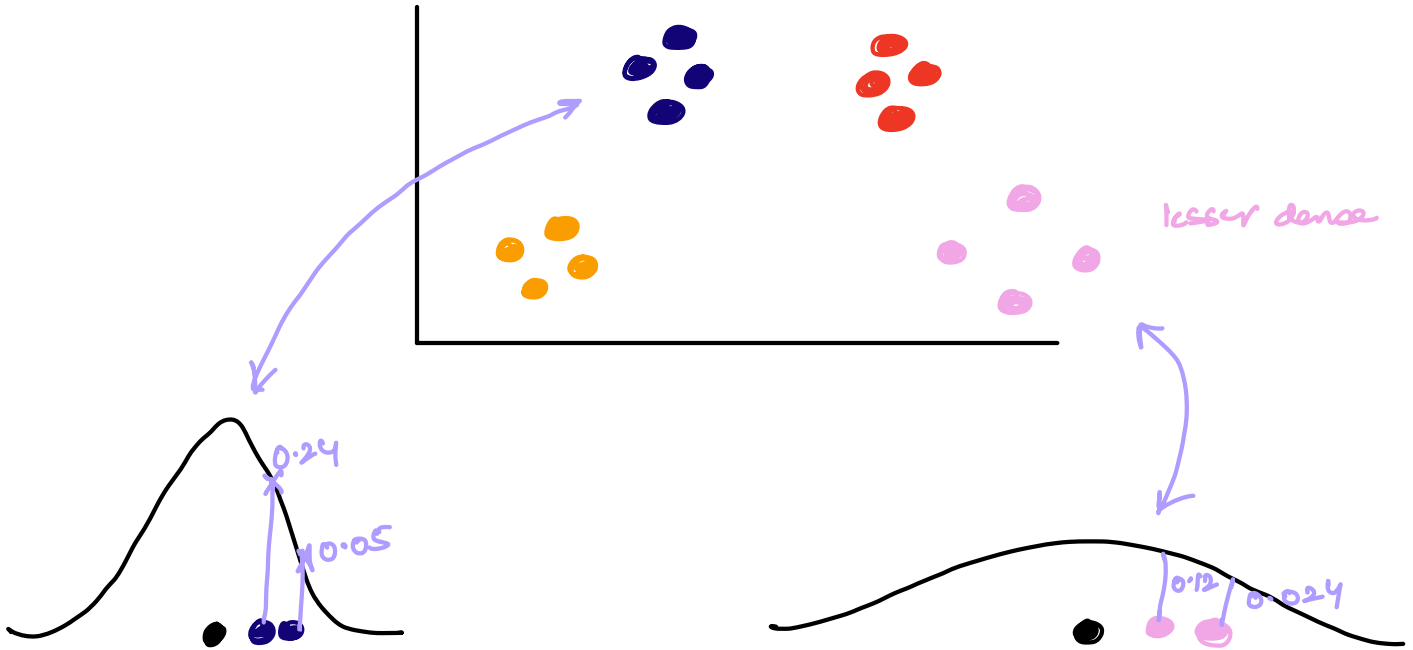
Distant points have very low similarity

We measure the distances b/w all the points & point of interest.





Unscaled  $\rightarrow$  Scaled ?

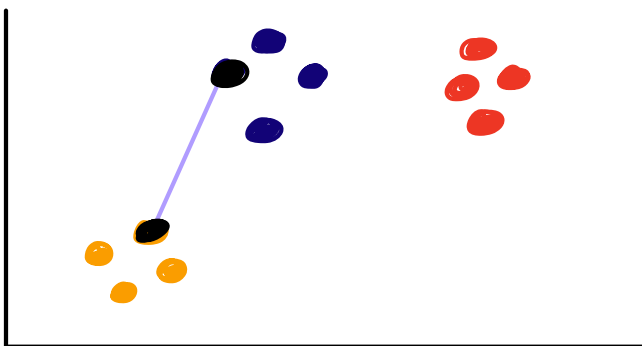


$$\frac{0.24}{0.24 + 0.05}, \frac{0.05}{0.24 + 0.05}$$

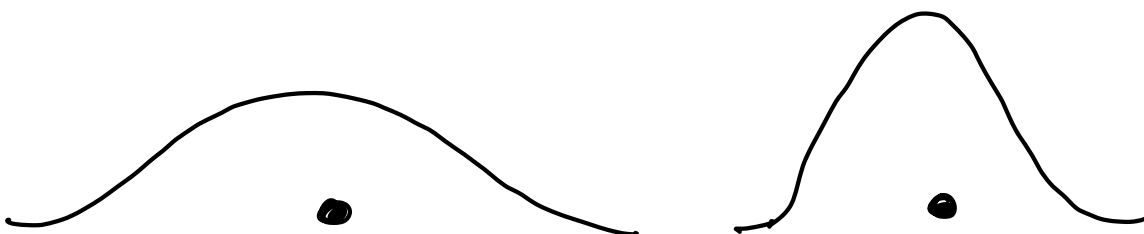
$$0.82, 0.18$$

$$\frac{0.12}{0.12 + 0.024}, \frac{0.024}{0.12 + 0.024}$$

$$0.82, 0.18$$

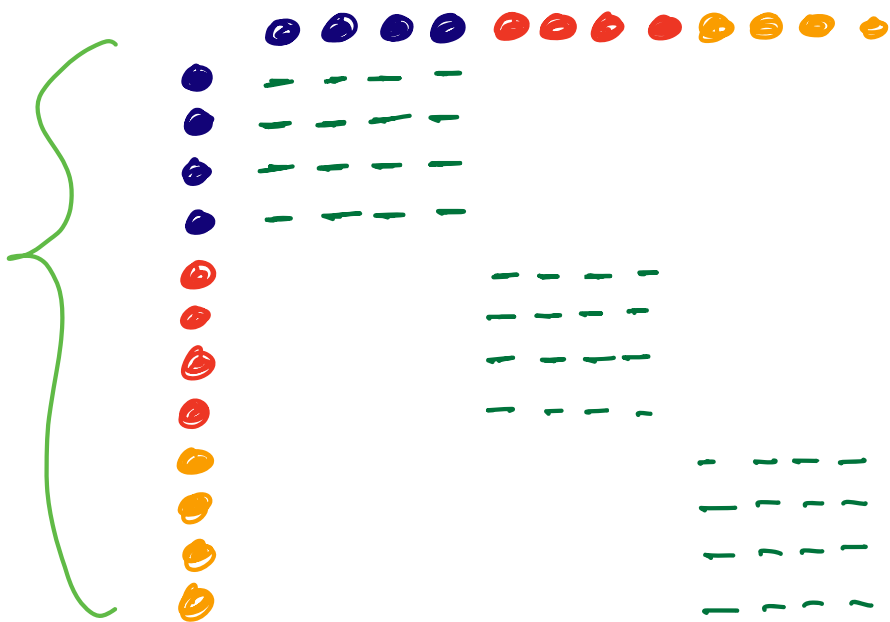


Keep each point in focus & calculate the distance.

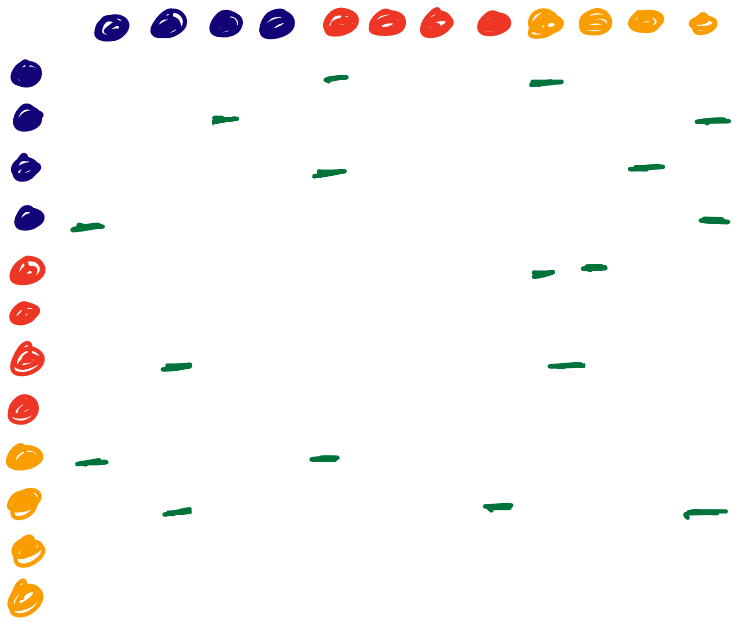
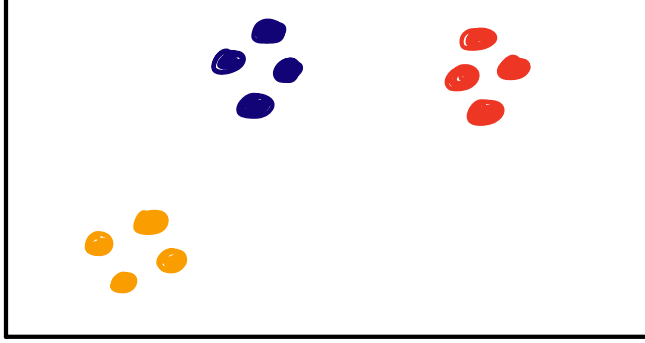


} Avg  
similarity  
score from  
both  
dr's

Similarity  
Matrix



— high similarity



# Fairness

Sensitive: Race, Age, Gender...

Group fairness

Individual fairness

males & females

2 individuals

$P(A|B)$ : Probability of A happening given that B has already happened.

$y$  = actual  $y$  / true  $y$

$\hat{y}$  = predicted  $y$

$G$  = groups

Group fairness:

- Statistical Parity / Demographic Parity:

$$P(\hat{y}=1 | G=f) = P(\hat{y}=1 | G=m)$$

	group	$\hat{y}$
1.	f	✓
2.	f	
3.	f	✓
4.	f	
5.	m	✓
6.	m	
7.	m	✓
8.	m	
9.	m	✓
10.	m	
11.	m	✓
12.	m	

$$\frac{2}{4} = \frac{4}{8} \quad \checkmark$$

males:  
 $P(\hat{y}=1 | G=m)$   
 0.4 → 0.6 → 0.8  
 20% 20% 20%  
 that's fine

	group	$\hat{y}$
1.	f	✓
2.	f	
3.	f	
4.	f	
5.	m	✓
6.	m	
7.	m	
8.	m	
9.	m	
10.	m	
11.	m	
12.	m	

$$\frac{1}{4}$$

$$\frac{6}{8}$$

not fair

5.	m	✓
6.	m	
7.	m	✓
8.	m	✓
9.	m	✓
10.	m	✓
11.	m	✓
12.	m	

### Equal Opportunity

(TPR should be same across both groups)

$$P(\hat{Y}=1 | Y=1, G=f) = P(\hat{Y}=1 | Y=1, G=m)$$

$$\frac{1}{2} \qquad \frac{2}{6}$$

	group	Y	$\hat{Y}$
1.	f	✓	✓
2.	f		
3.	f	✓	
4.	f		✓
5.	m	✓	
6.	m		
7.	m	✓	✓
8.	m	✓	
9.	m	✓	✓
10.	m		✓
11.	m	✓	
12.	m	✓	

### Equalized Odds:

TPR & FPR should be same across groups

$$P(\hat{Y}=1 | Y=1, G=f) = P(\hat{Y}=1 | Y=1, G=m)$$

and

$$P(\hat{Y}=1 | Y=0, G=f) = P(\hat{Y}=1 | Y=0, G=m)$$

### Overall Accuracy Equality:

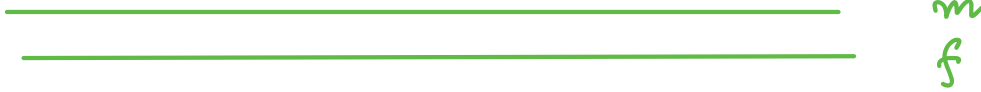
Acc. should be same across both the groups

$$\frac{TP_f + TN_f}{TP_f + TN_f + FP_f + FN_f} = \frac{TP_m + TN_m}{TP_m + TN_m + FP_m + FN_m}$$

Individual fairness:

### - Causal Independence (Causal Discrimination)

CGPA, Branch, Exp., Projects, Gender



} Results should be same for these 2 individuals.

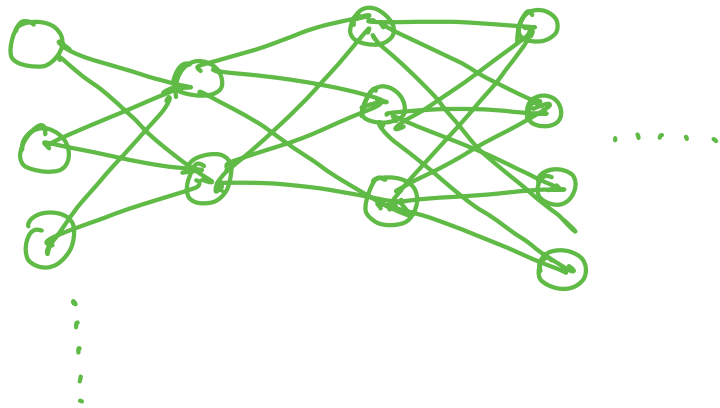
### - fairness through unawareness:

Sensitive attribute are excluded in the decision making process.

### Explainable AI

#### - Transparency : Functioning & Decision making process.

C, D, H  
(weights)



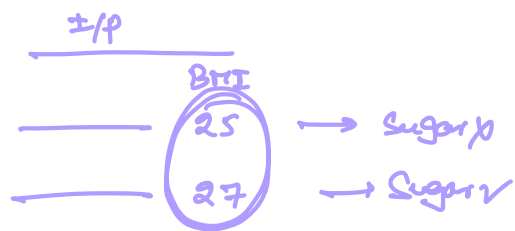
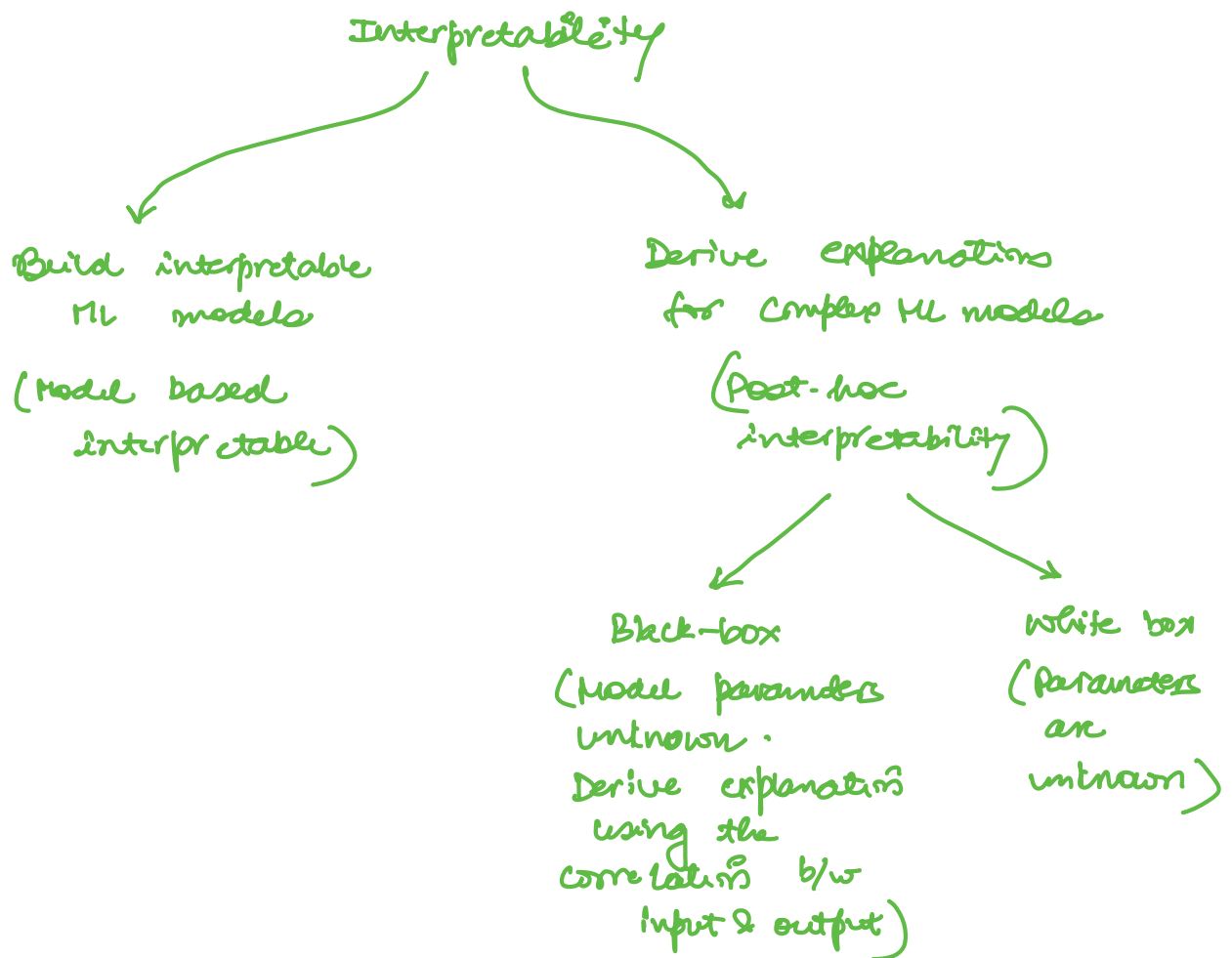
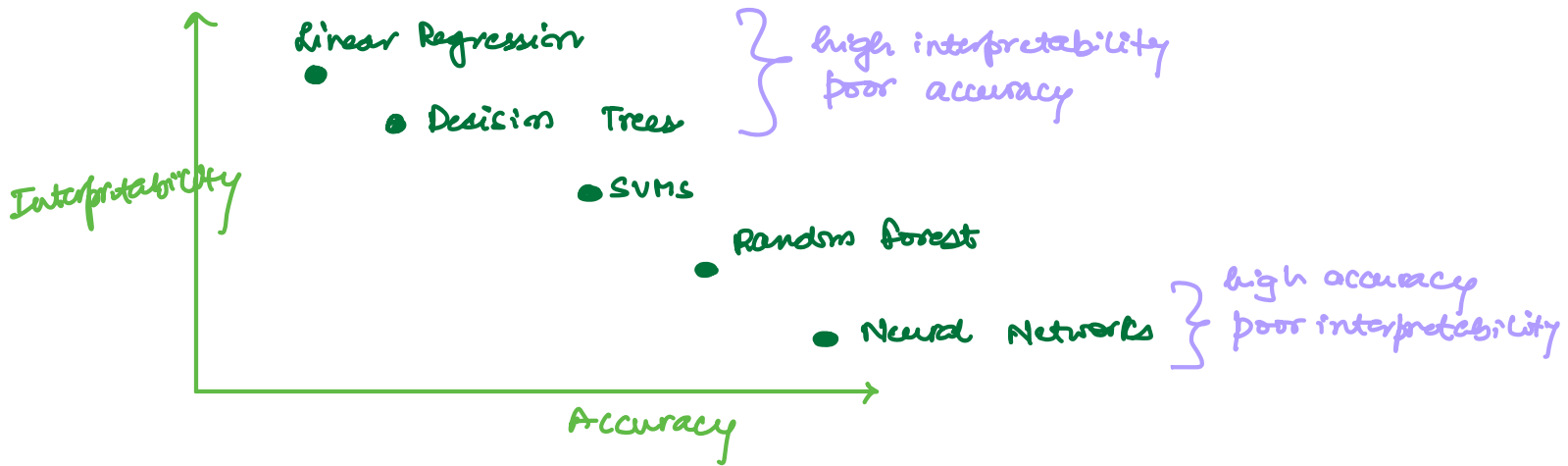
#### - Interpretability

Human understandable

#### - Accountability

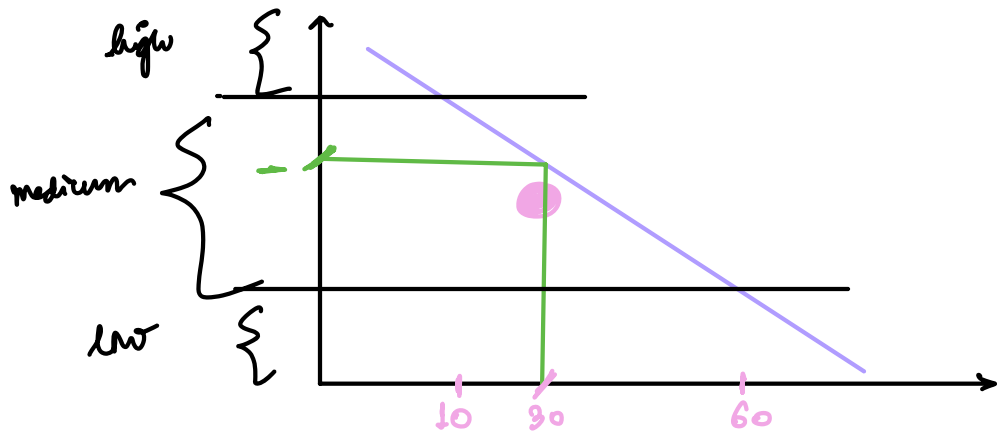
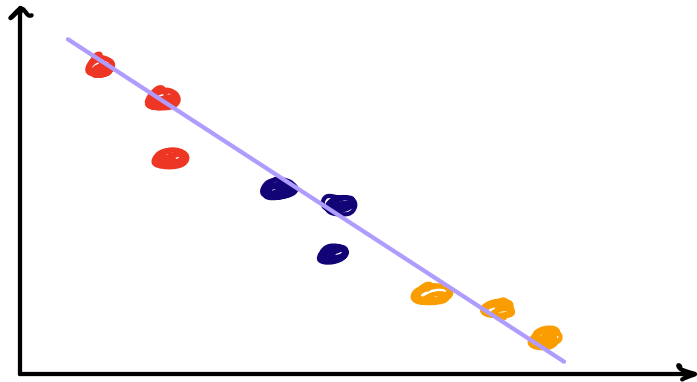
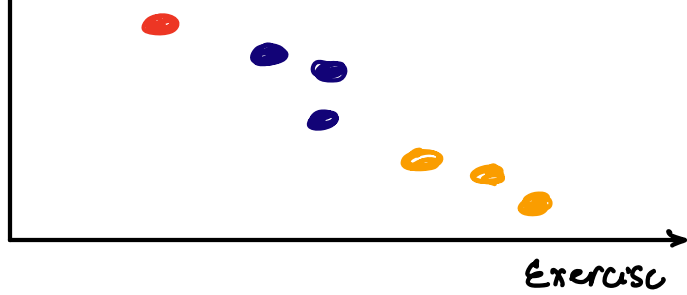


## Interpretability v/s Accuracy



## Linear Regression



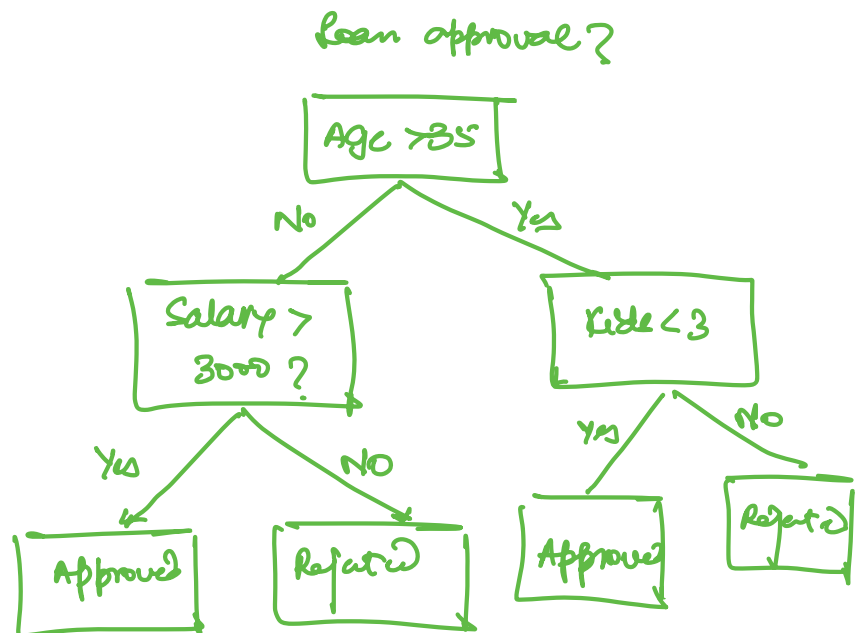


<10: high  
10-60: medium  
>60: low

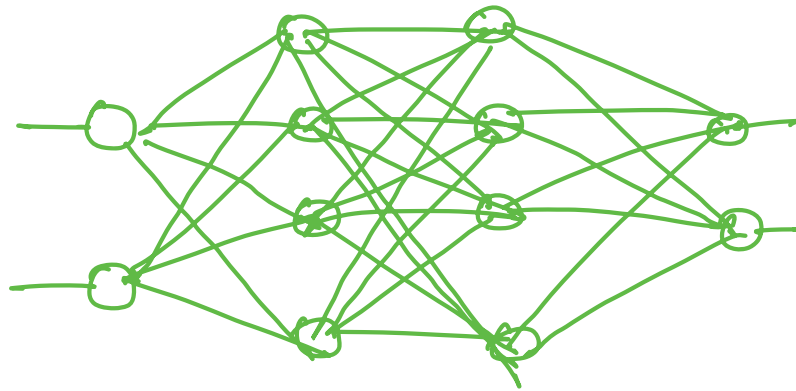
## Decision Trees

Age: 32  
Salary: 2800  
Title: 2

↓  
Reject

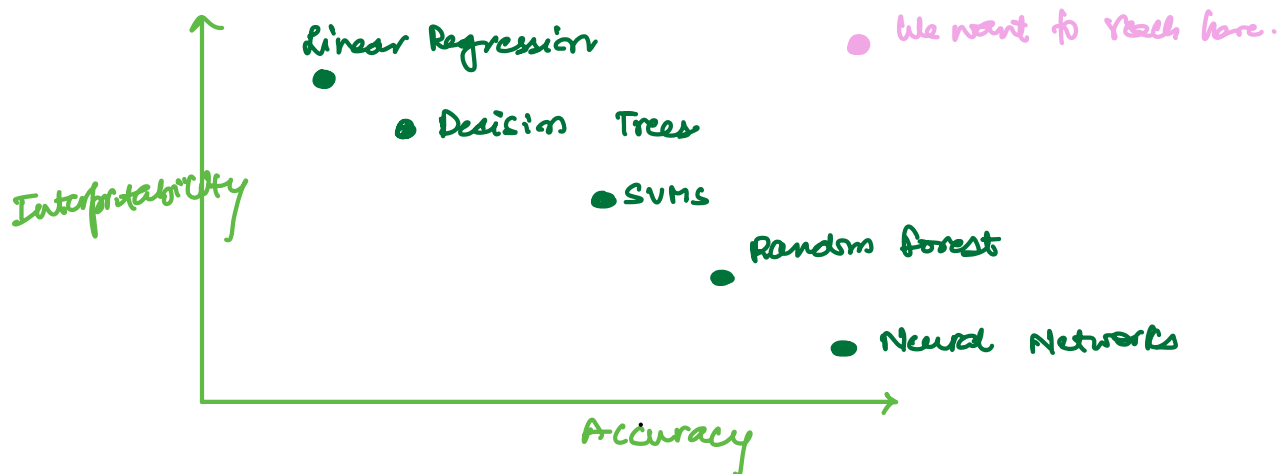


# Neural Networks



Tumor  
Detect

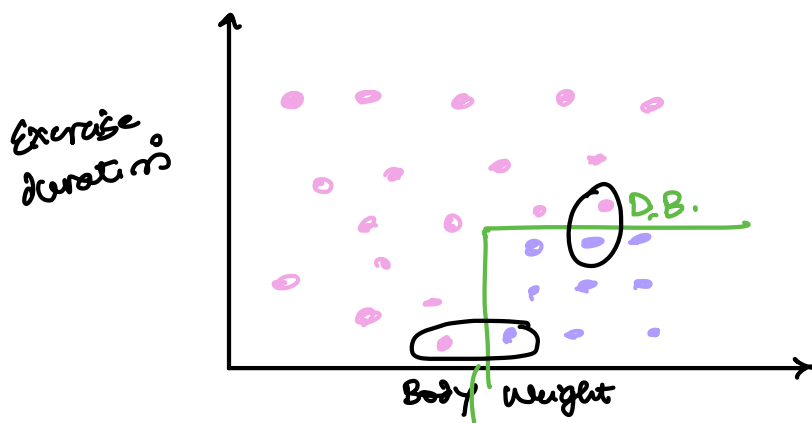
NN are not inherently interpretable. we don't know which feature got maximum importance.



## LIME

Local  
Interpretable  
Model-Agnostic  
Explanations

- local neighbourhood of the instance
- A human should be able to interpret
- Applicable to all models
- Explanations that help interpretation

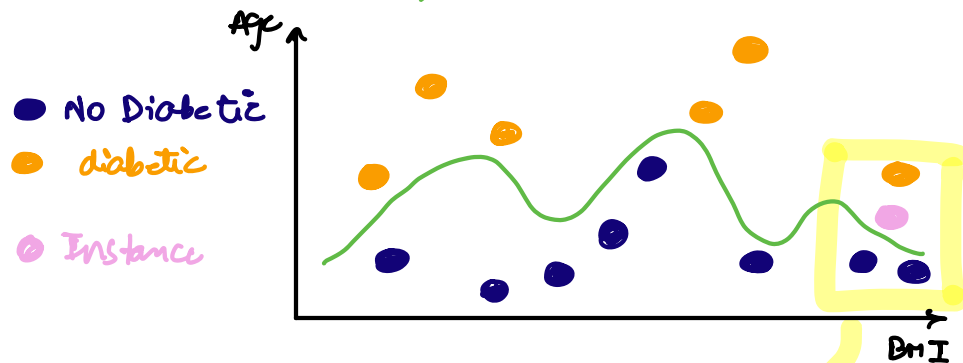


• healthy  
• unhealthy

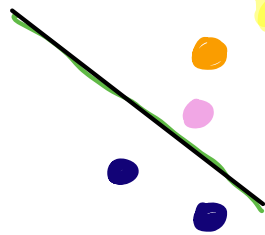
feature differentiating  
these 2 people?

LIME Implementation:

① Global model, local data point



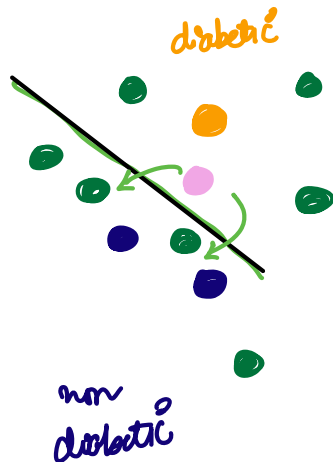
② Local neighbourhood



local linear model near the point

③

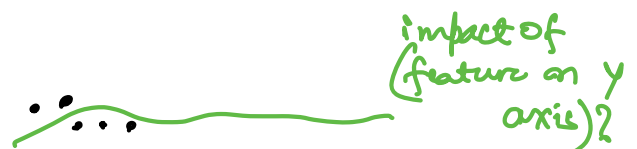
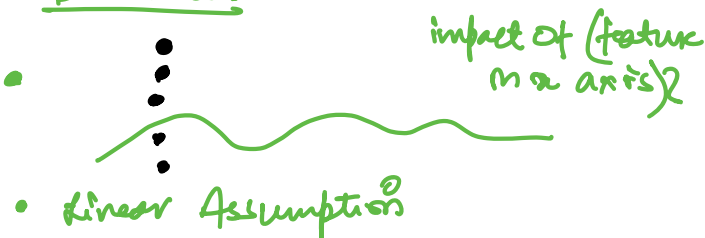
• New  
Random  
data points  
generate



④ features relevant to  
the data point

Age —  
BMI —  
Heart Disease —  
Gender —  
Cholesterol —

Drawback:

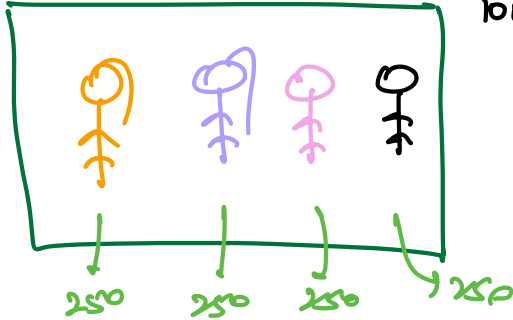


# SHAP

Shapley Additive explanations

Lloyd Shapley  $\rightarrow$  Game theory

team



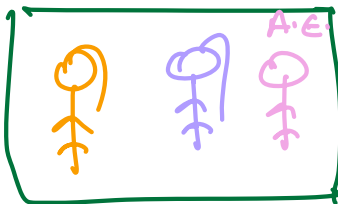
1000 R.C.

Distribution ?

Option 1:

Equal  
Distribution:

Option 2:



H.E. 2nd pos  
600 R.C.


$\rightarrow$  400 R.C.?

We don't account for  
interactions ?



Contributions ?



- Consider different subsets.
- Calculate individual contributions of  in each subset
- weighted avg of these individual contributions will give you the contribution of each player.

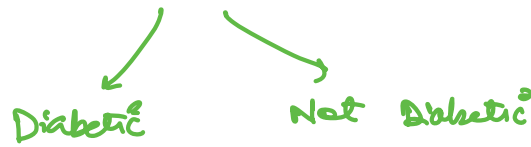
① Age

② Exercise Duration

③ BMI

④ Diabetic parents

⑤ Calorie Intake



Subset: Age Random Values BMI Random Values Calorie Intake

Subset: Random Values Random Values BMI Random Values Calorie Intake

ML model

Data point

weight

$M = \text{total \# features}$

$z' = \text{\# features in subset}$

$$\phi_i(f, x) = \sum_{z' \subseteq x'} \frac{|z'|! (M - |z'| - 1)!}{M!}$$

Shapley value for feature  $i$

Subsets

if features are:  
Age | BMI | Calorie Intake  
then the result from ML model

$$\left[ \underbrace{f_x(z')}_{60\%} - \underbrace{f_x(z' \setminus i)}_{15\%} \right]$$

BMI | Calorie Intake  
result from ML model

Drawback:

Computational Complexity is high

#Subsets:  $2^n$