

Decision Trees

Introduction and Geometric Intuition

Trainer: Dr Darshan Ingle

Trainer: Dr. Darshan Ingle



Example 1

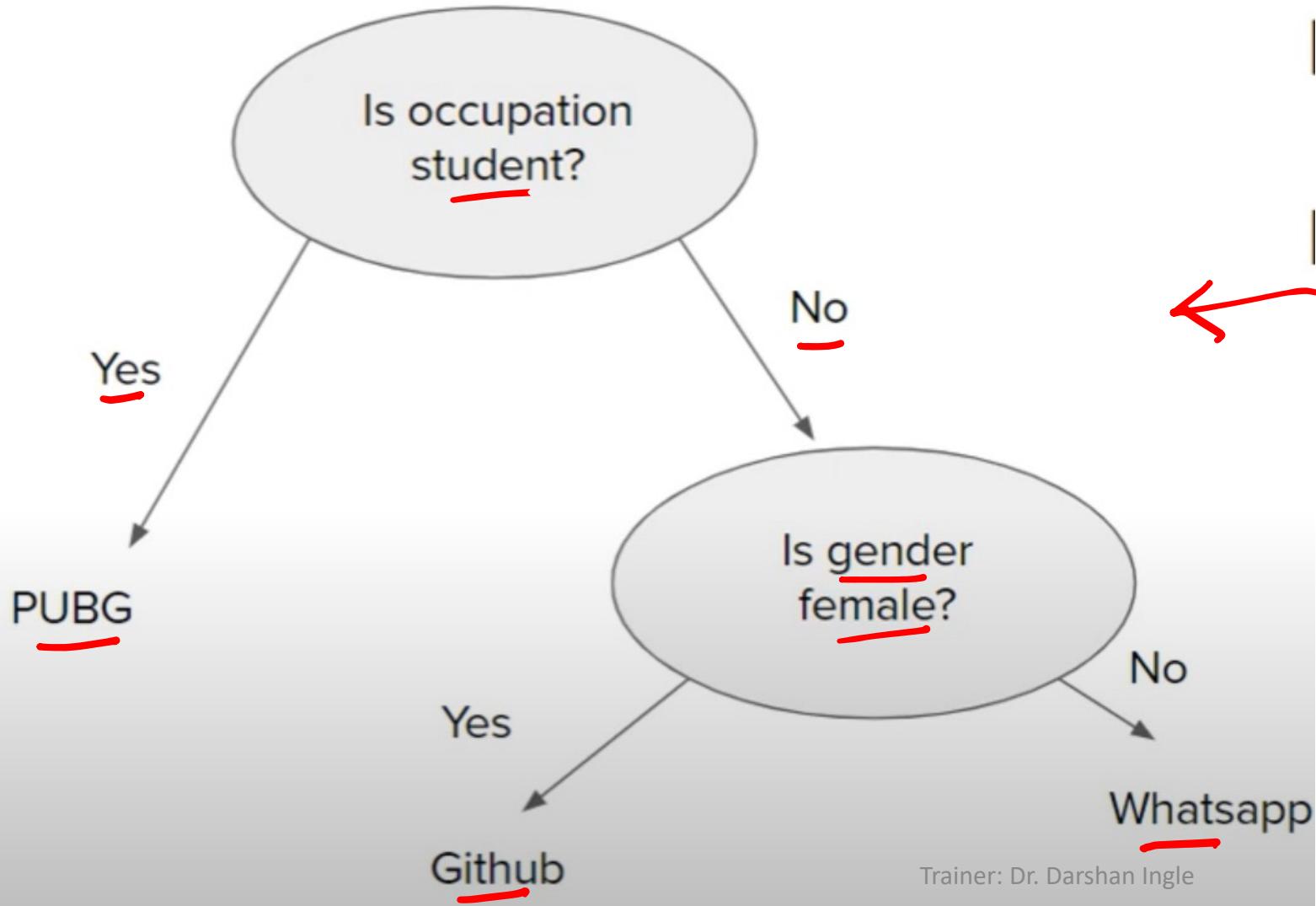
x_1 x_2

y

Gender	Occupation	Suggestion
F	Student	PUBG
F	Programmer	Github
M	Programmer	Whatsapp
F	Programmer	Github
M	Student	PUBG
M	Student	PUBG

```
If occupation==student  
    print(PUBG)  
Else  
    If gender==female  
        print(Github)  
    Else  
        print(Whatsapp)
```

Where is the Tree?

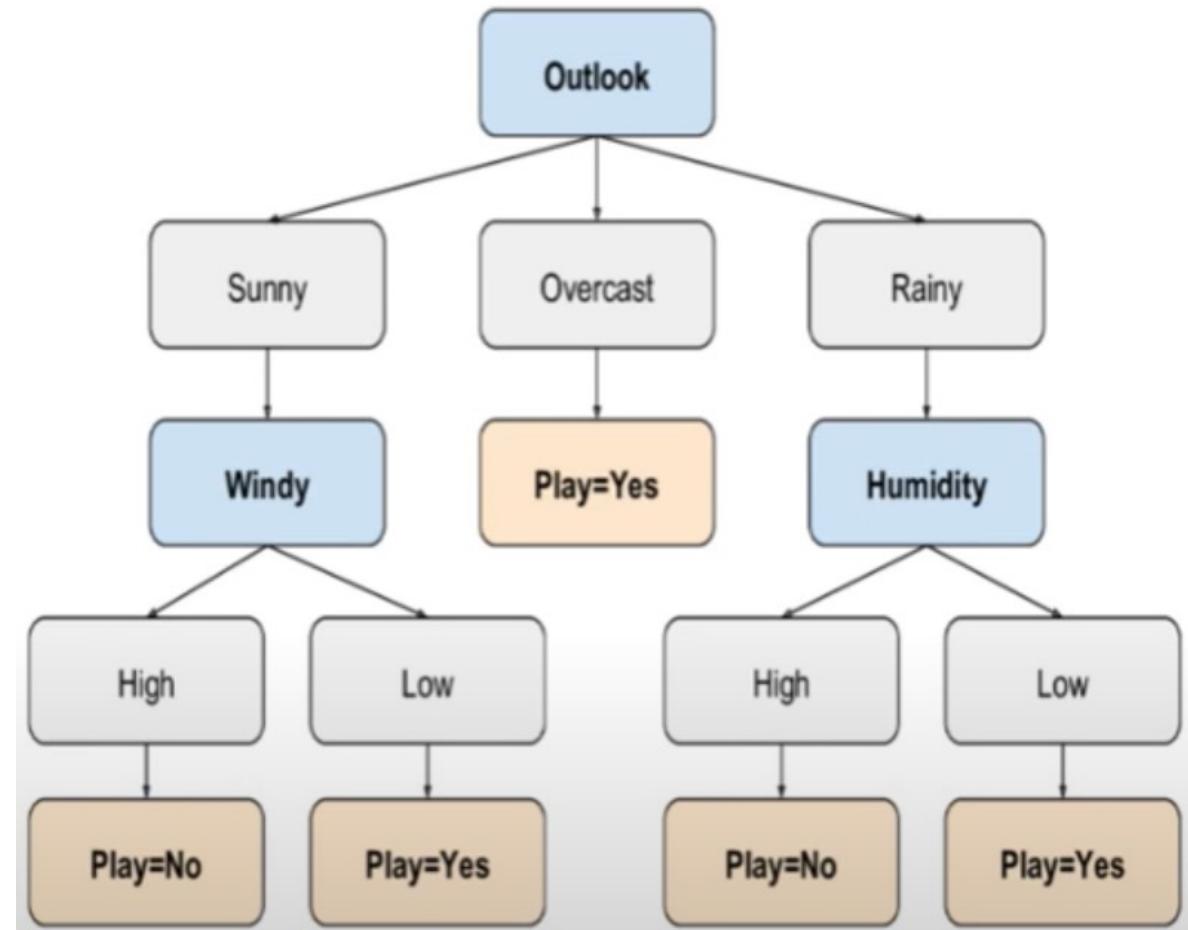


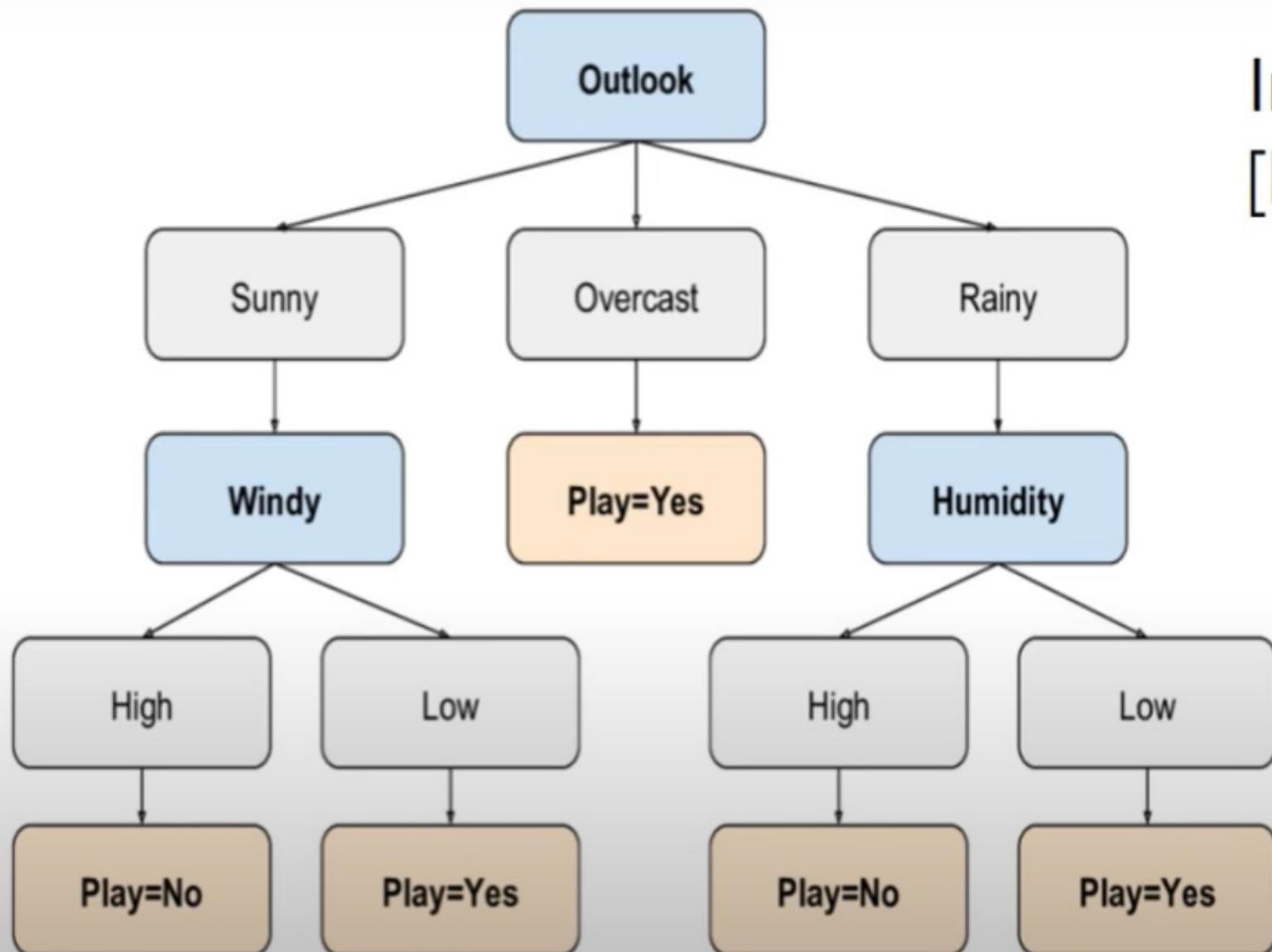
```
If occupation==student  
    print(PUBG)  
Else  
    If gender==female  
        print(Github)  
    Else  
        print(Whatsapp)
```

Example 2

Day	Outlook	Temp	Humid	Wind	Play?
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

y





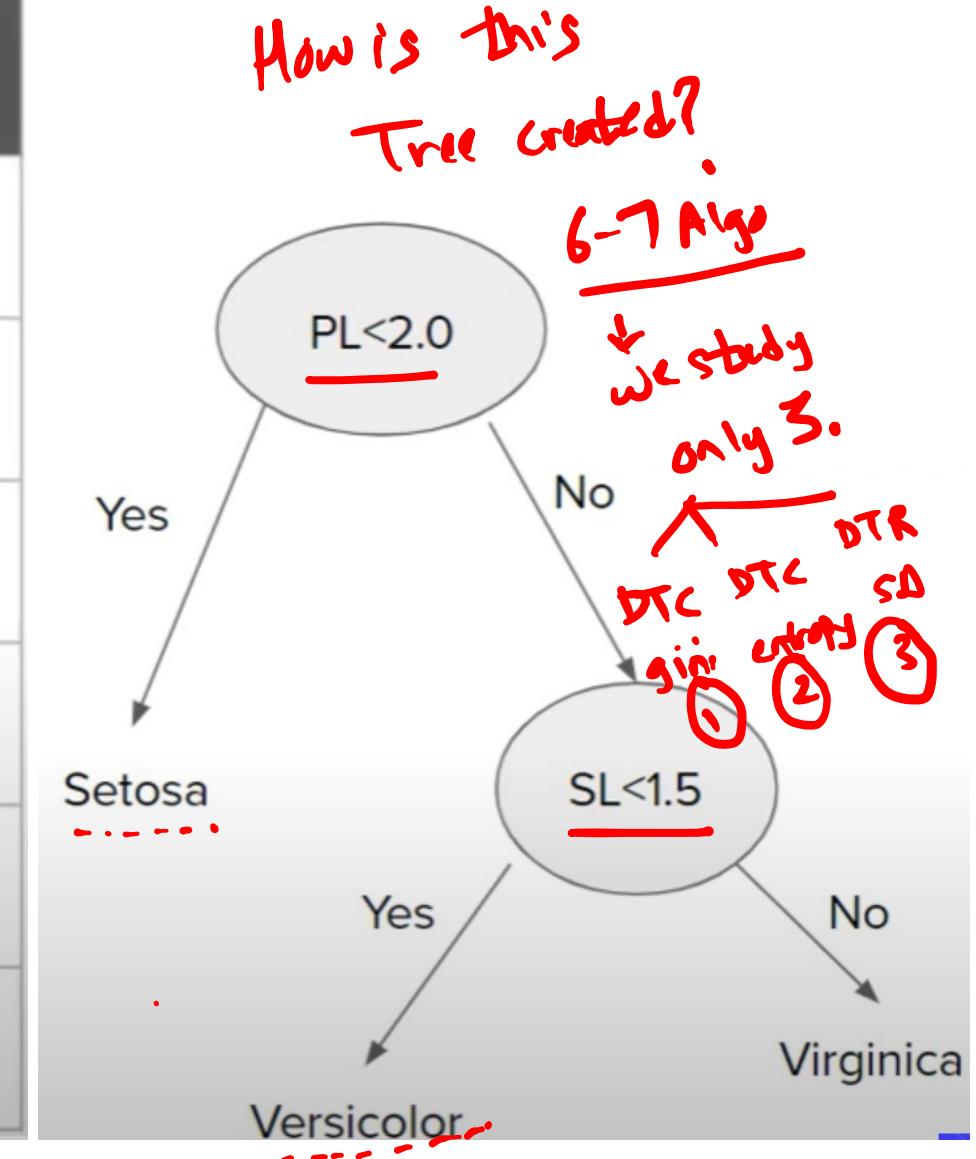
Input query point:
[Rainy, Mild, High, Strong]

What if we have numerical data?

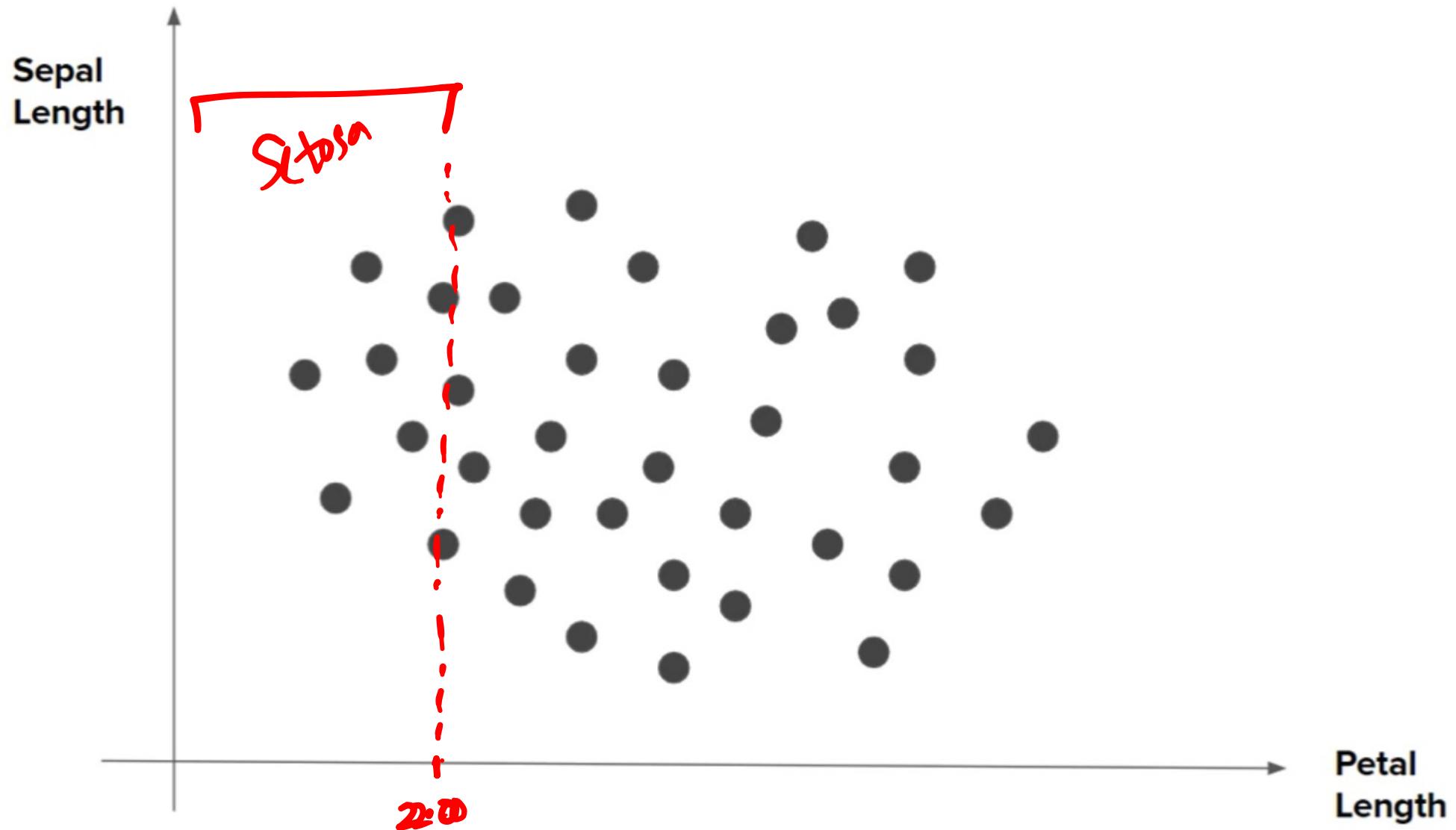
x-axis X_1 *y-axis X_2* *y*

Petal Length	Sepal Length	Type
1.34	0.34	Setosa
3.45	1.45	Versicolor
1.69	0.98	Setosa
2.56	1.79	Virginica
3.00	1.13	Versicolor
1.3	0.88	Setosa

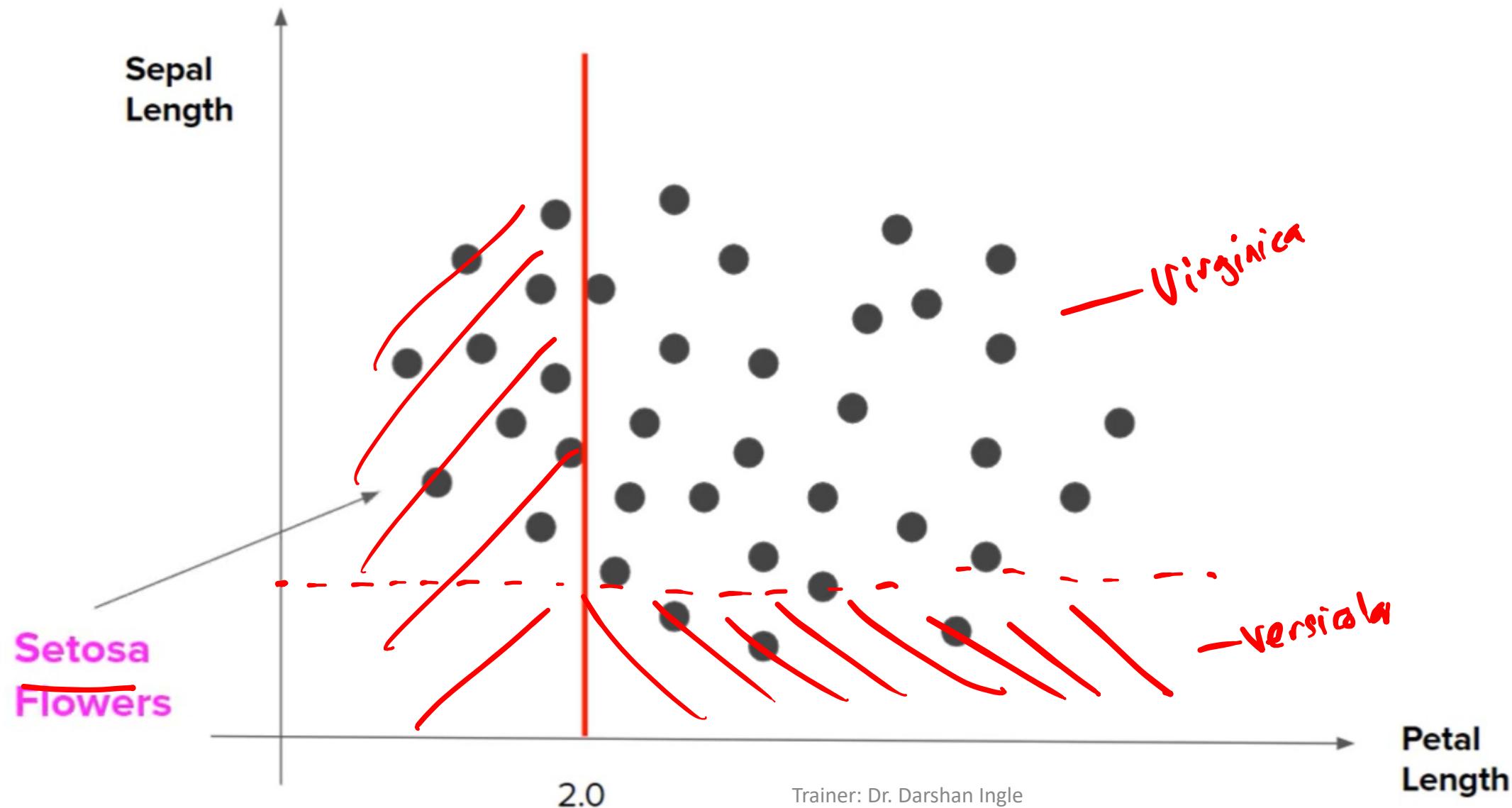
Trainer: Dr. Darshan Ingle



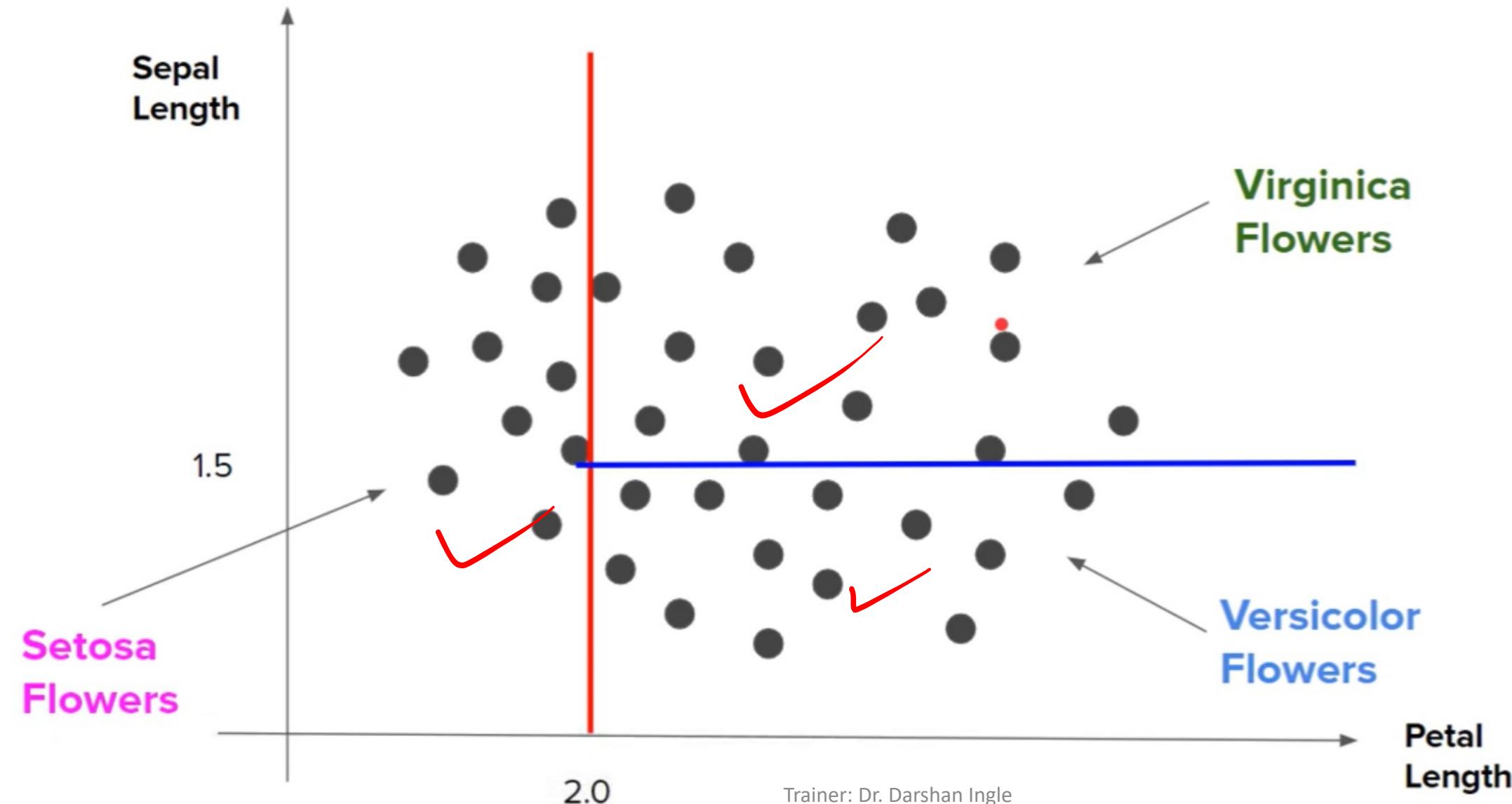
Geometric Intuition



Geometric Intuition



Geometric Intuition



Pseudo code

- Begin with your training dataset, which should have some feature variables and classification or regression output.
- Determine the “best feature” in the dataset to split the data on; more on how we define “best feature” later
- Split the data into subsets that contain the correct values for this best feature. This splitting basically defines a node on the tree i.e each node is a splitting point based on a certain feature from our data.
- Recursively generate new tree nodes by using the subset of data created from step 3.

Conclusion

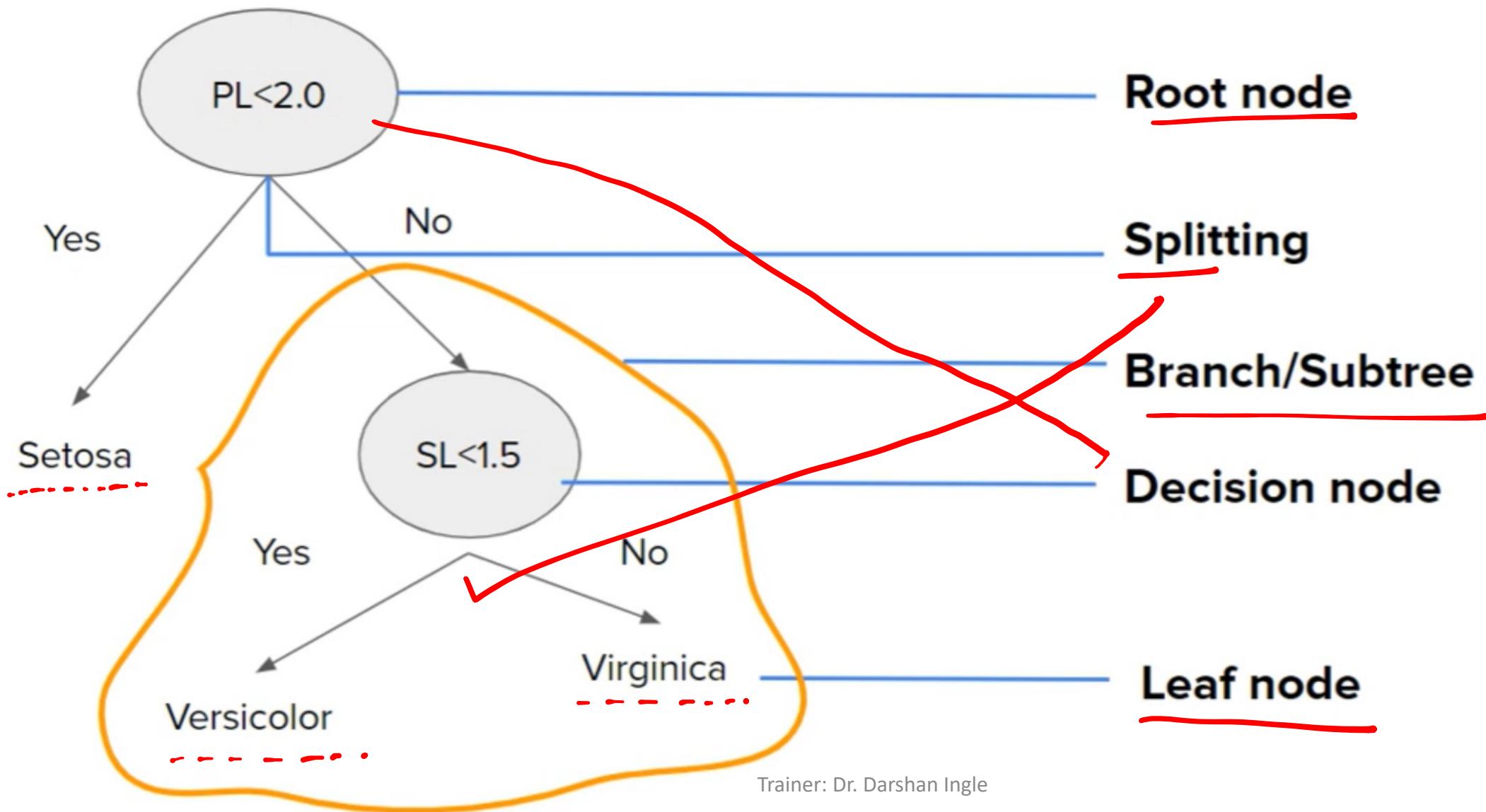
Programmatically speaking, Decision trees are nothing but a giant structure of nested if-else condition



Mathematically speaking, Decision trees use hyperplanes which run parallel to any one of the axes to cut your coordinate system into hyper cuboids

Hyperplane → 3D AD, ...

Terminology



Some unanswered questions

How to decide which column should be considered as root node?

How to select subsequent decision nodes?

How to decide splitting criteria in case of numerical columns?

P2 < 2.0

S2 < 1.5

Advantages



Intuitive and easy to understand

Minimal data preparation is required No Scaling is reqd. in D.T.

The cost of using the tree for inference is logarithmic in the number of data DS points used to train the tree

Disadvantages

Overfitting

Prone to errors for imbalanced datasets



Overshooting:

Rule! ① Difference b/w Train & Test accuracy < 5%, else it is overfitting.
② Industry Accepted Model, Test acc $\geq 85\%$.
 & Train acc \geq Test acc.
③ If Train Acc < 85%, model is underfitting.

Train Acc. Test Acc. Model

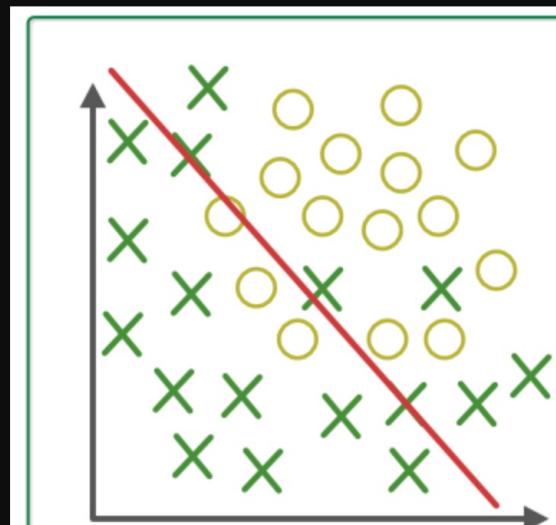
① 98%	96%	Good
② 98%	90%	OF
③ 92%	93%	Recheck prev. steps.
④ 86%	84%	UF
⑤ 60%	51%	UF

When he was given unseen data, he predicted wrongly on it & that too CONFIDENTLY.

OVERFITTING

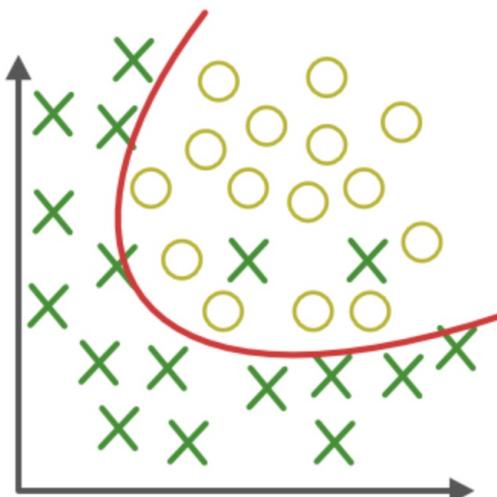
You are screen sharing Stop Share

Classification

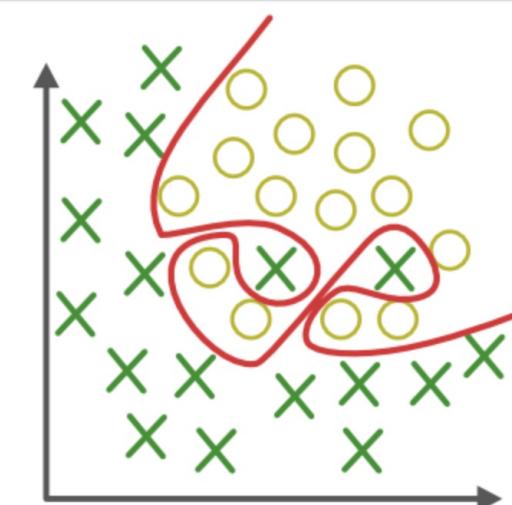


Under-fitting

(too simple to explain the variance)



Appropriate-fitting

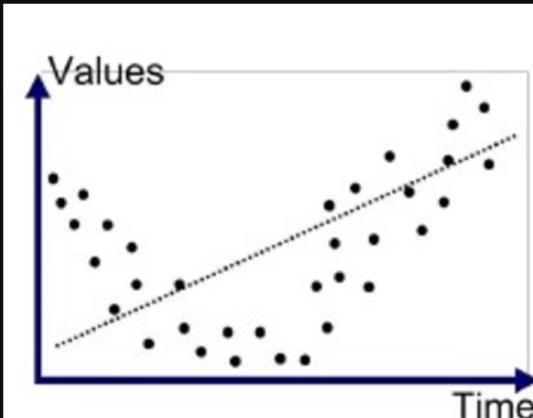


Over-fitting

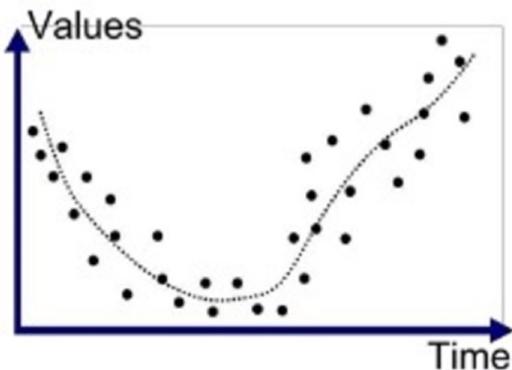
(force fitting--too good to be true)



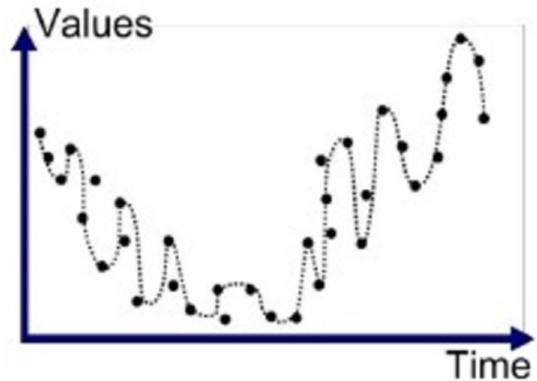
Regression:



Underfitted



Good Fit/Robust



Overfitted

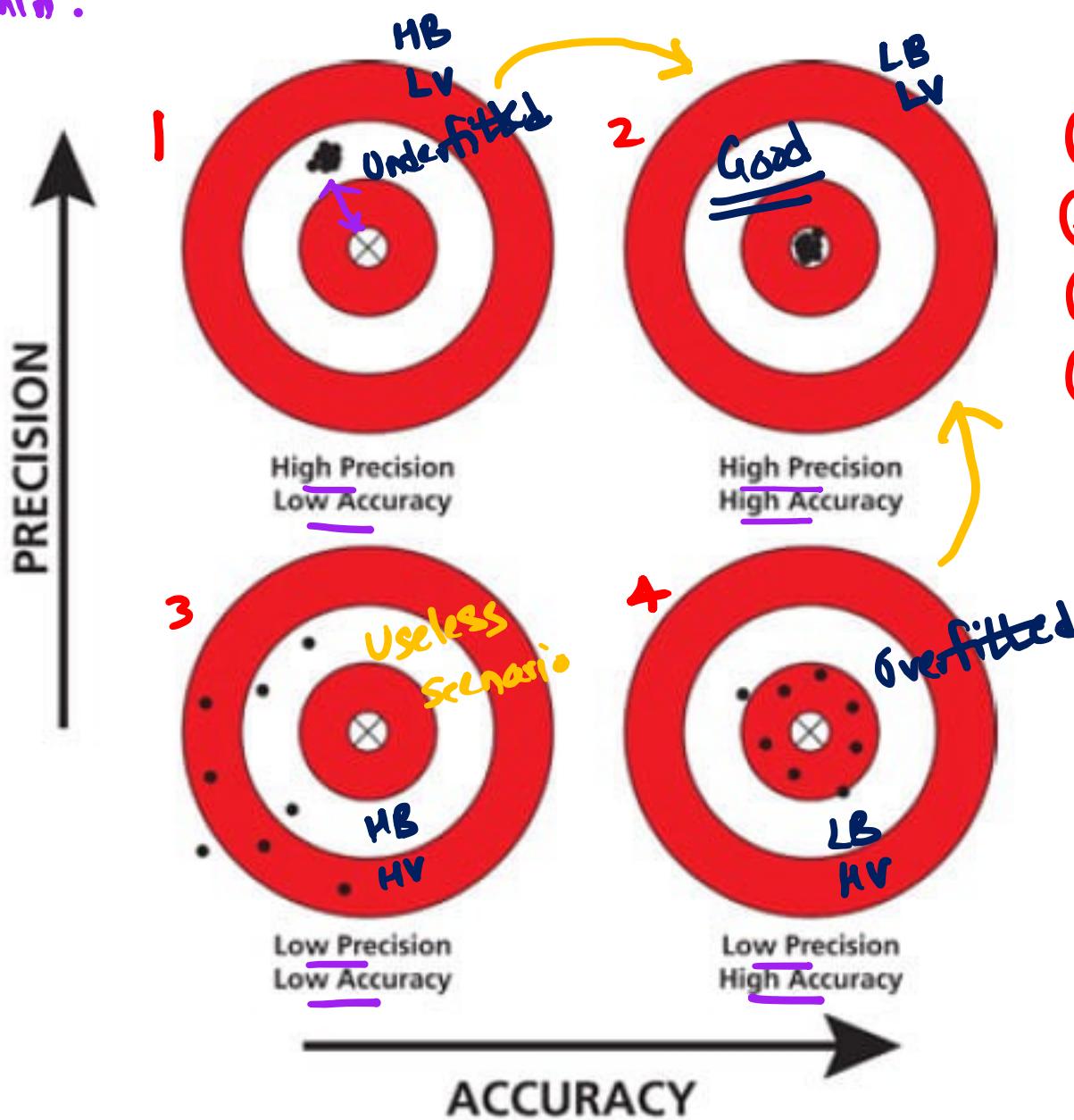


Precision: Spread is less/min.

Accuracy: Close to target / less error.

Bias: error

Variance: Spread of data



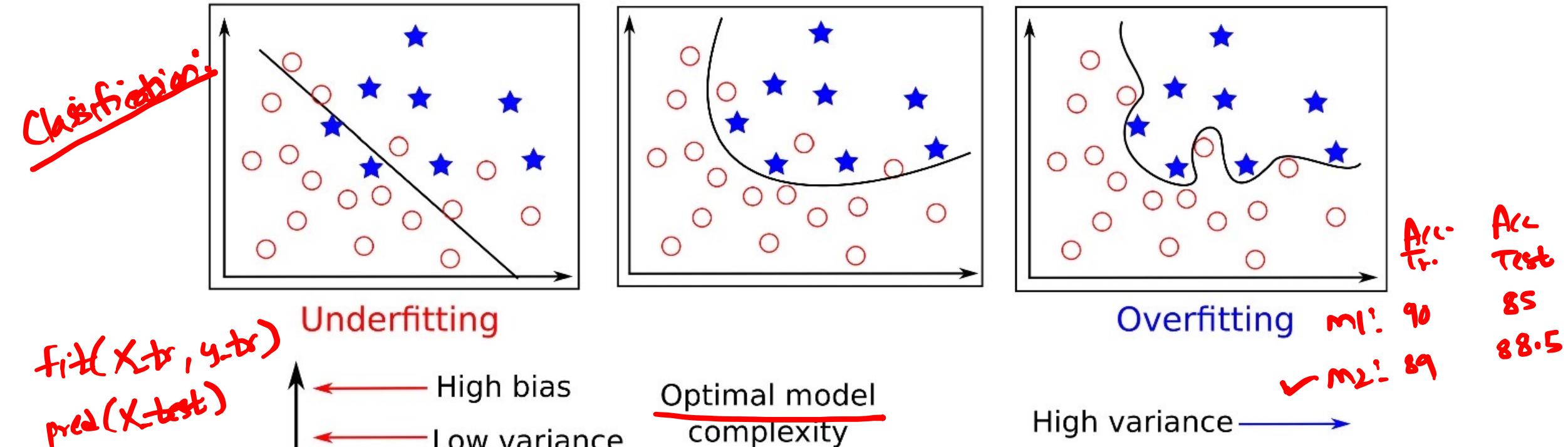
Arches: 20 arrows

① Arjun (Mahabharat)²

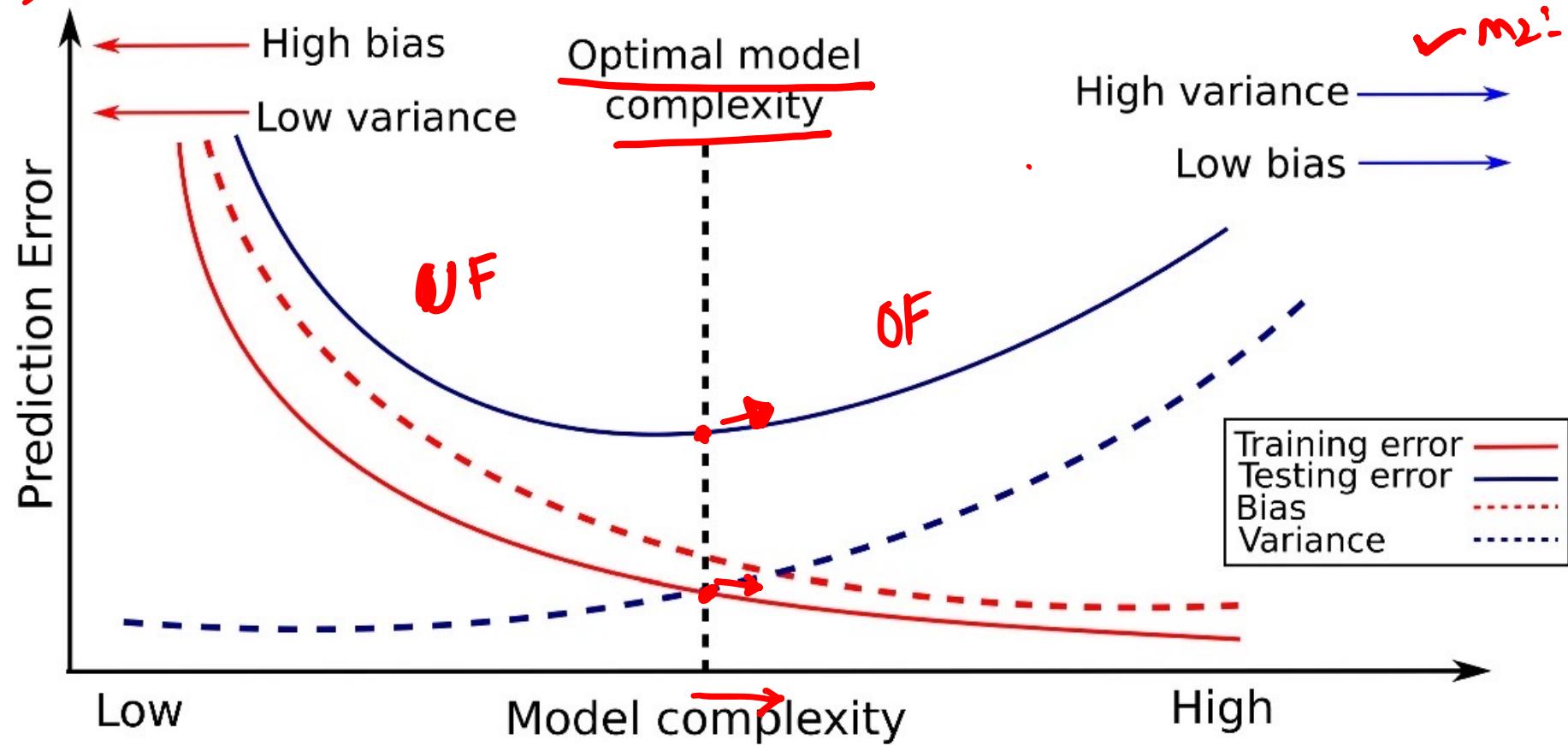
② Arjun with squint eyes¹

③ Darshan +

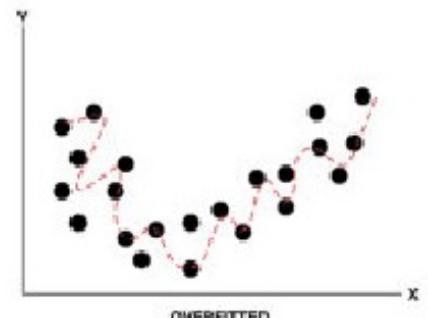
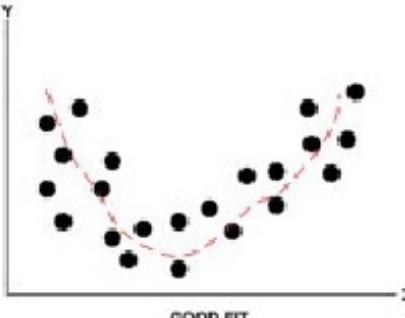
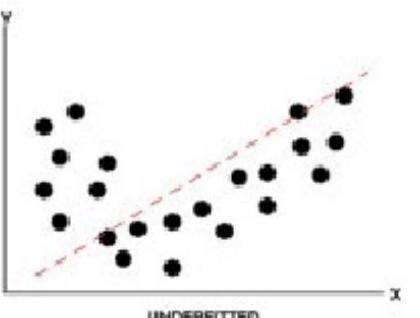
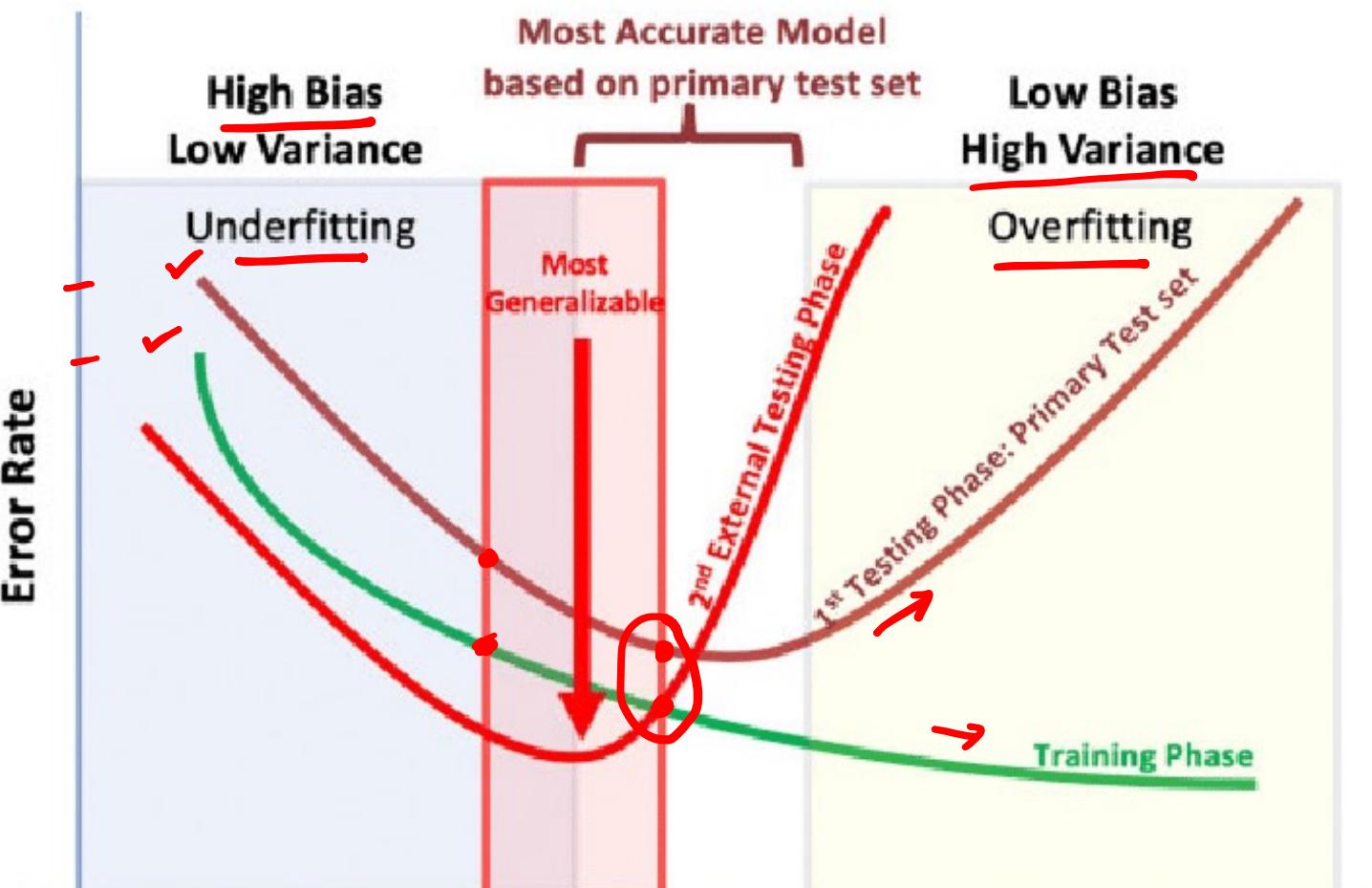
④ Darshan with squint eyes³



$\text{fit}(X_{\text{tr}}, y_{\text{tr}})$
 $\text{pred}(X_{\text{test}})$



Regression!



CART - Classification and Regression Trees

The logic of decision trees can also be applied to regression problems, hence the name CART

A fun example

<https://en.akinator.com/>



Trainer: Dr. Darshan Ingle

Thank You



Trainer: Dr. Darshan Ingle