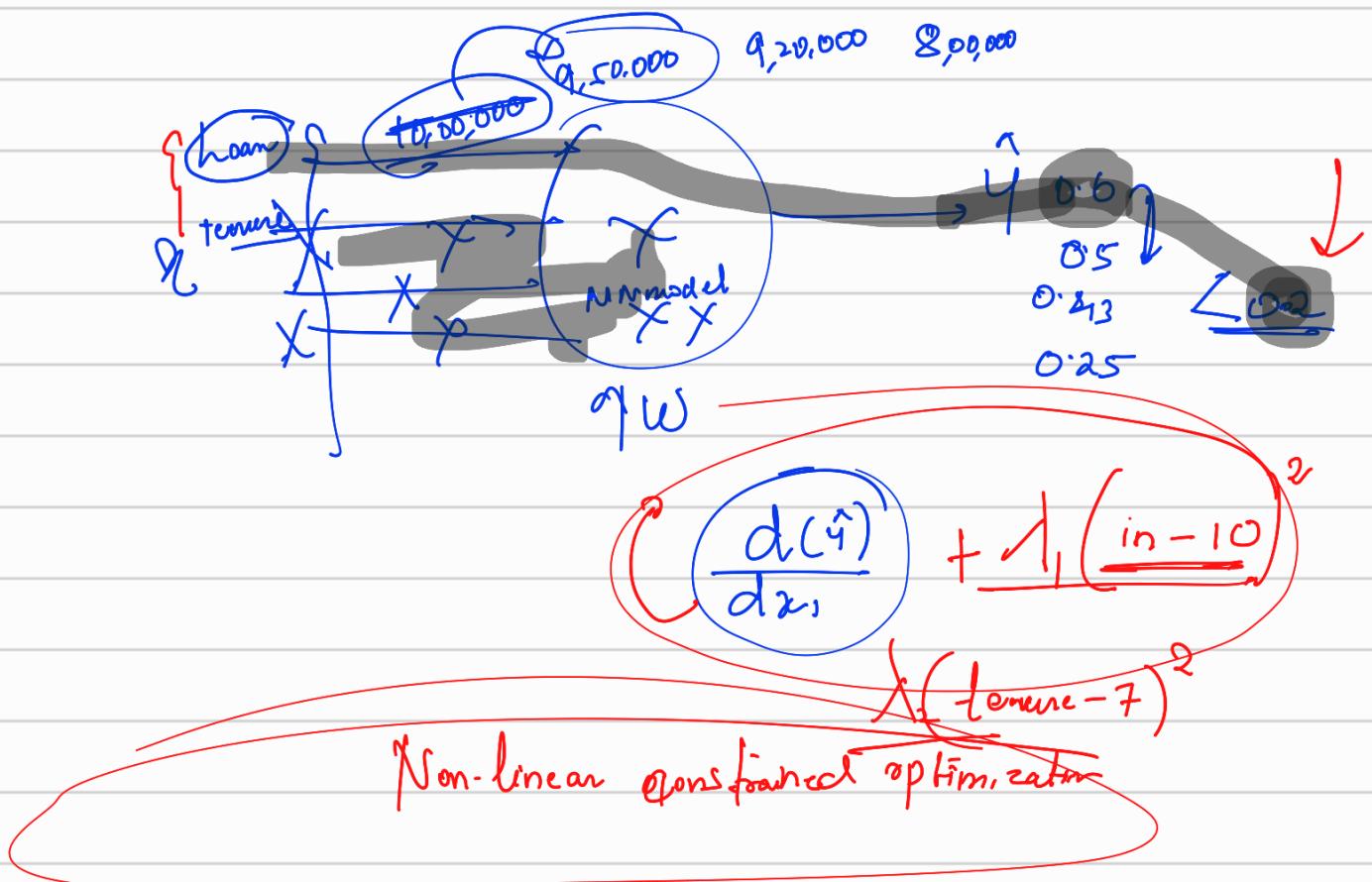


# Prediction Solved



## Model Explainability:

is a concept where in we are able to explain the prediction that we do, in the sense that we can atleast let the user know that what was the reason that the current sample got this prediction.

## Feature importance- Model Explainability?

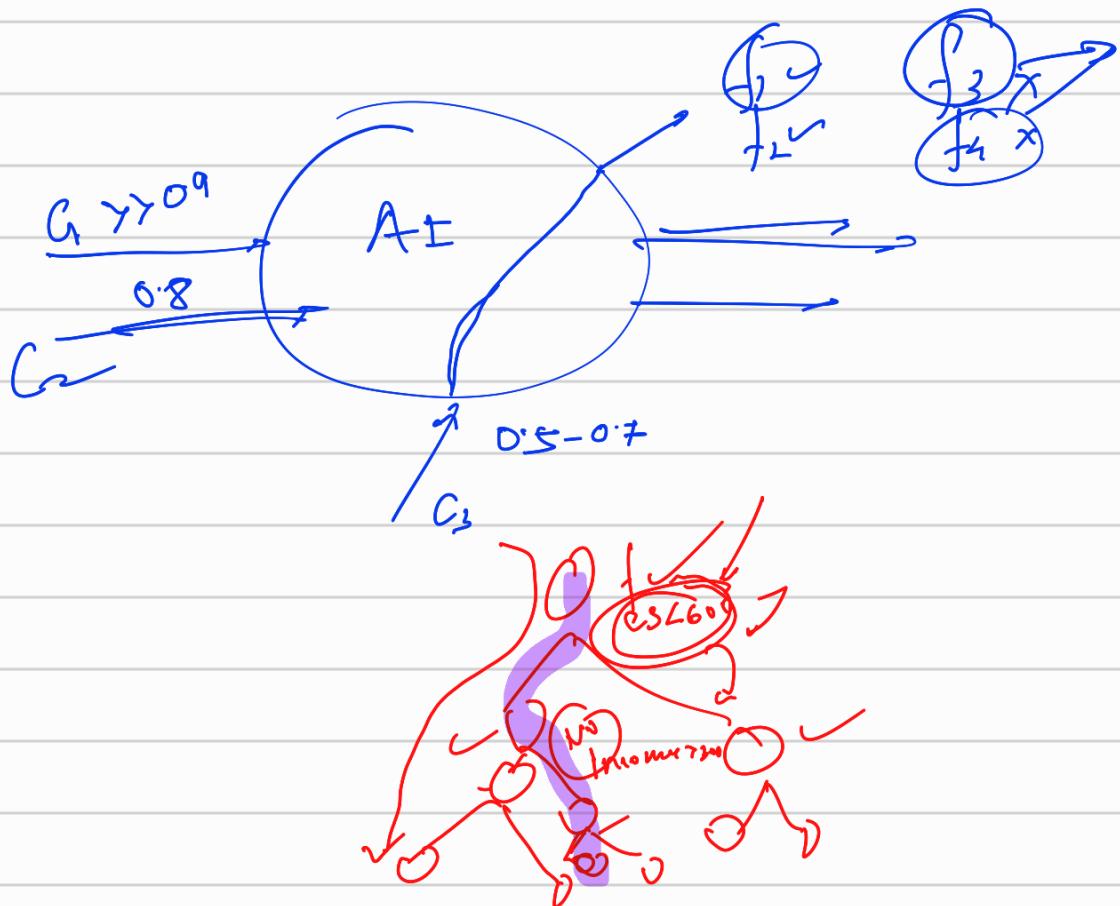
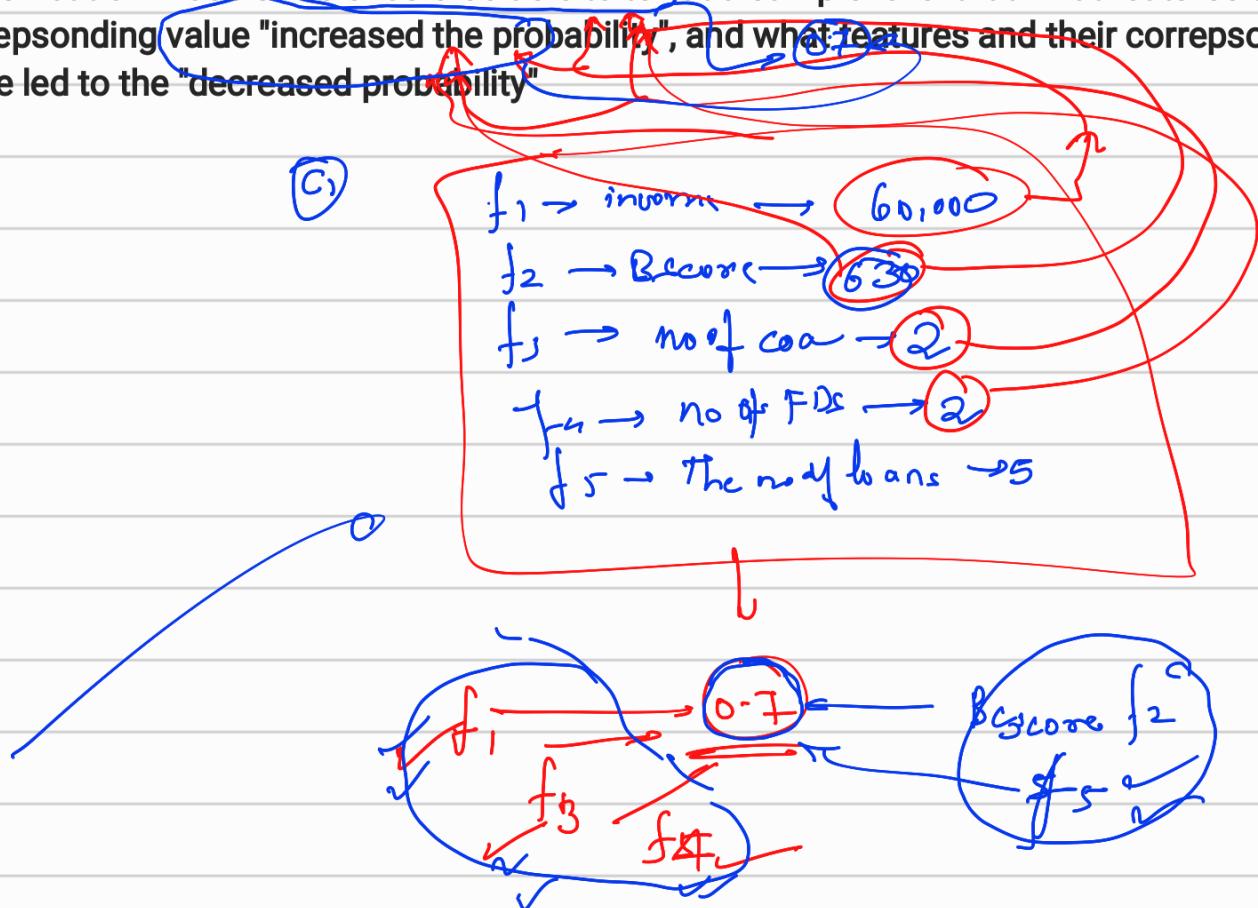
Feature importance from DT/XGB models can be considered as explainability however they

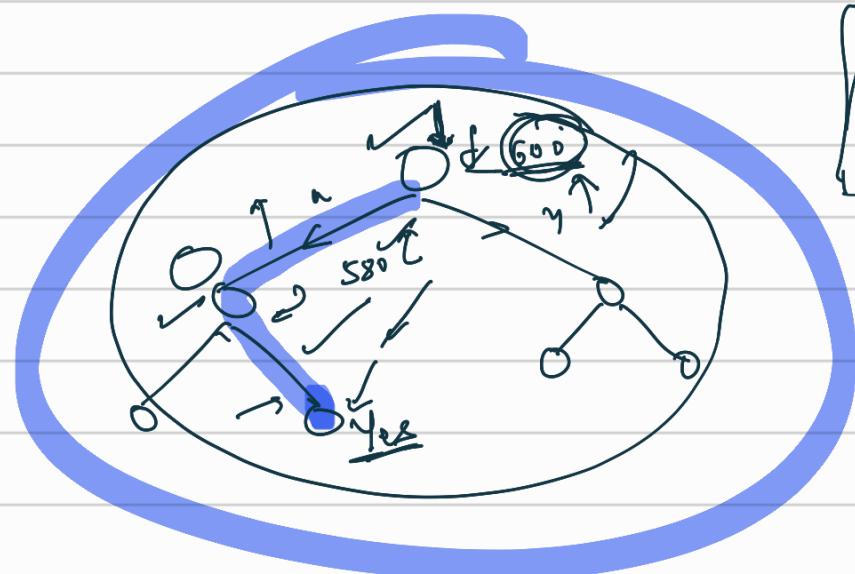
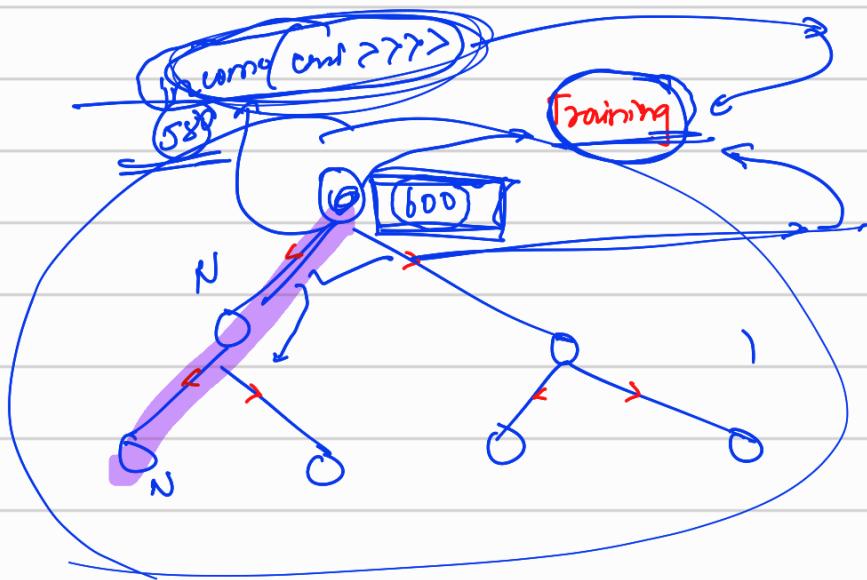
work on global level, they can help you to do feature selection

How do we define explainability?

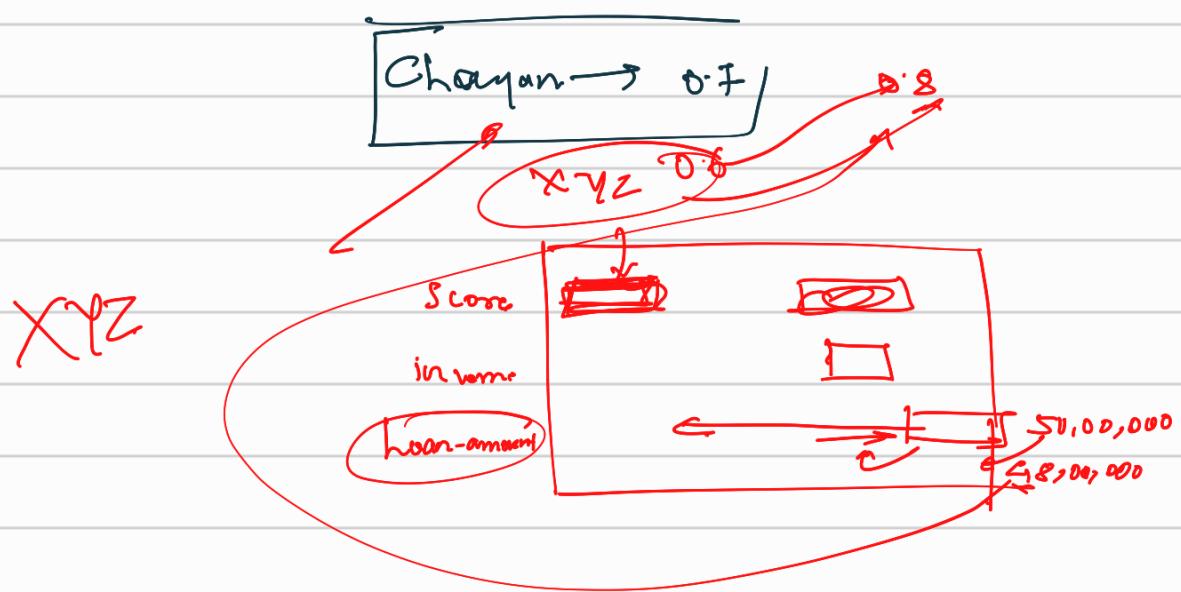
Explainability is during the inference time/prediction has to work at a sample level, what features led to this value.

Classification Models - we should be able to tell at a sample level that what features and its corresponding value "increased the probability", and what features and their corresponding value led to the "decreased probability"



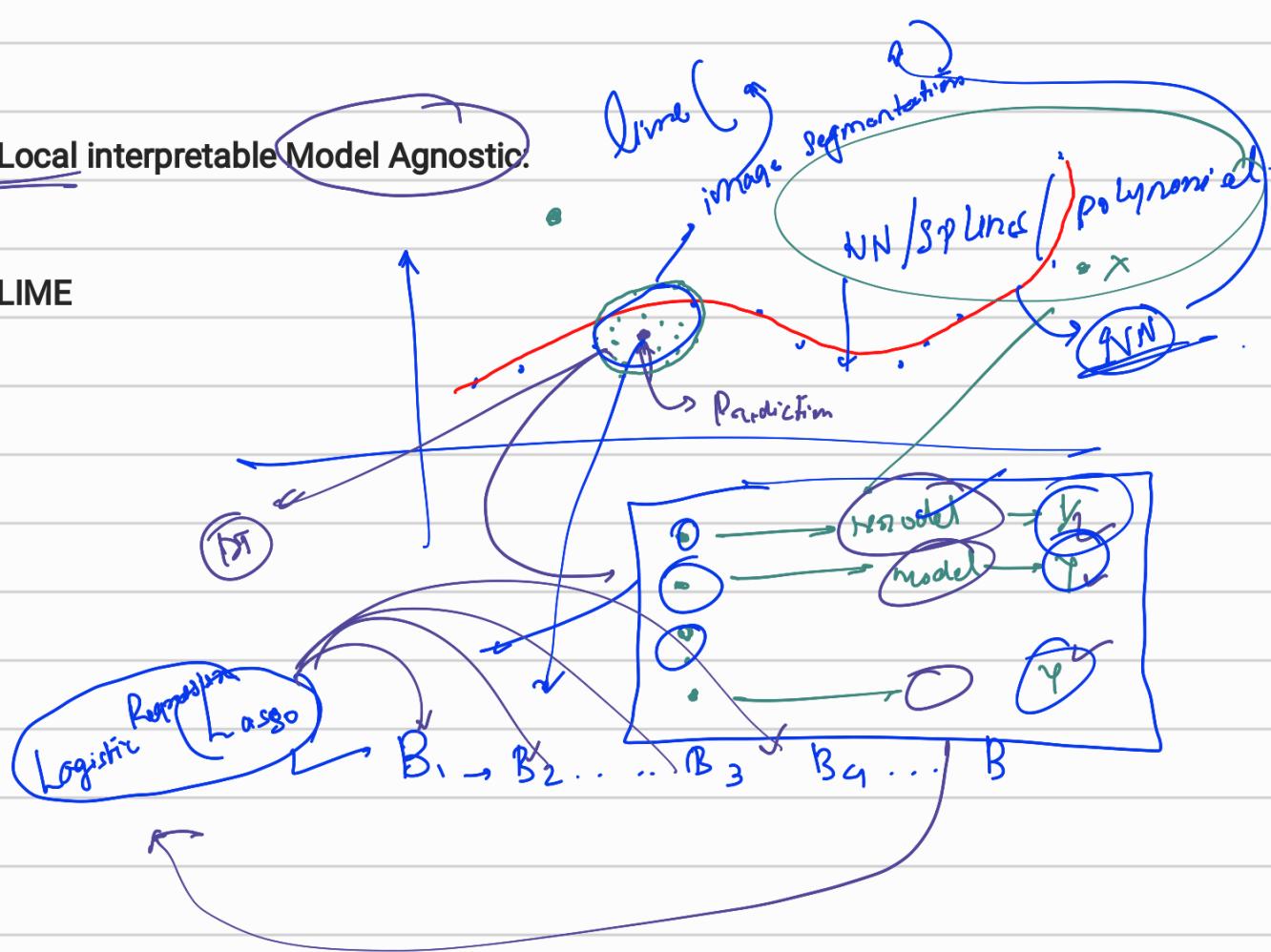


Local level / Sample level

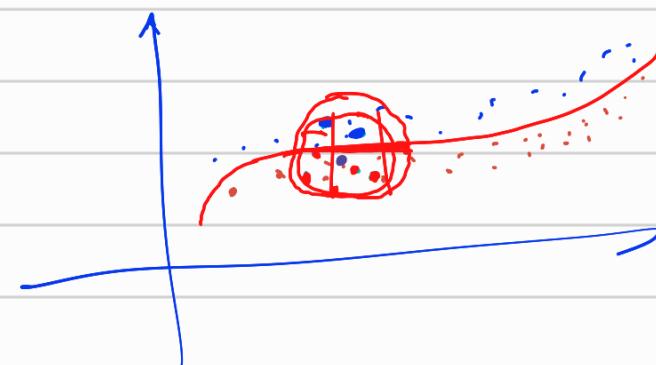


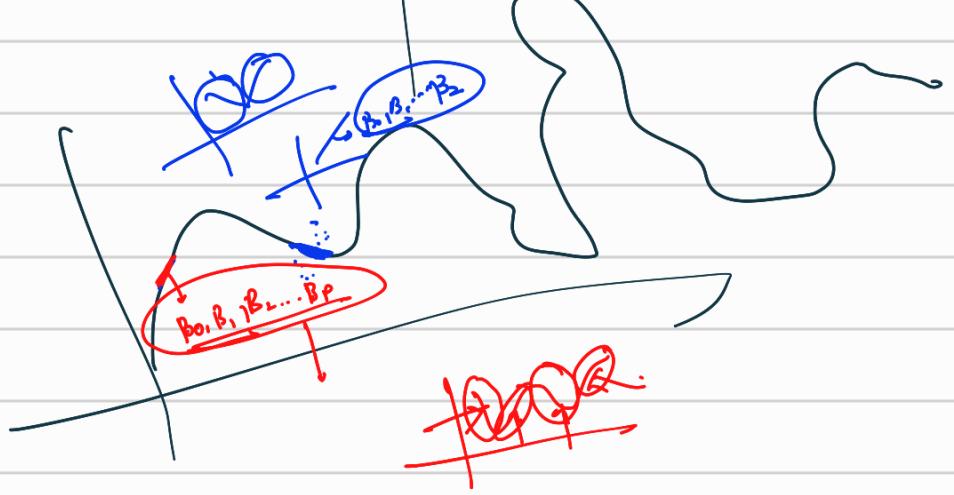
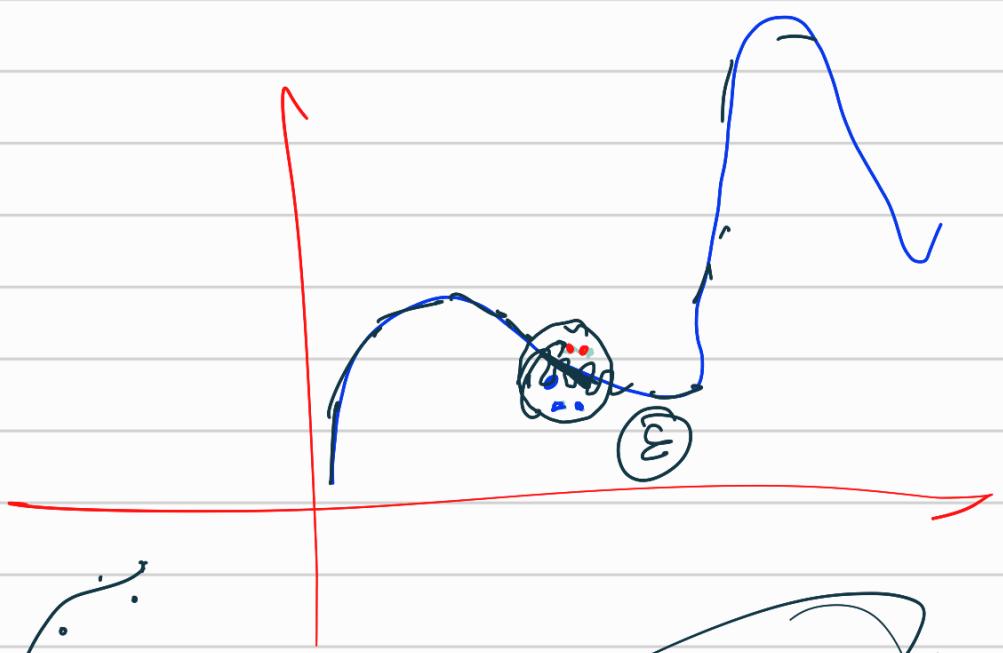
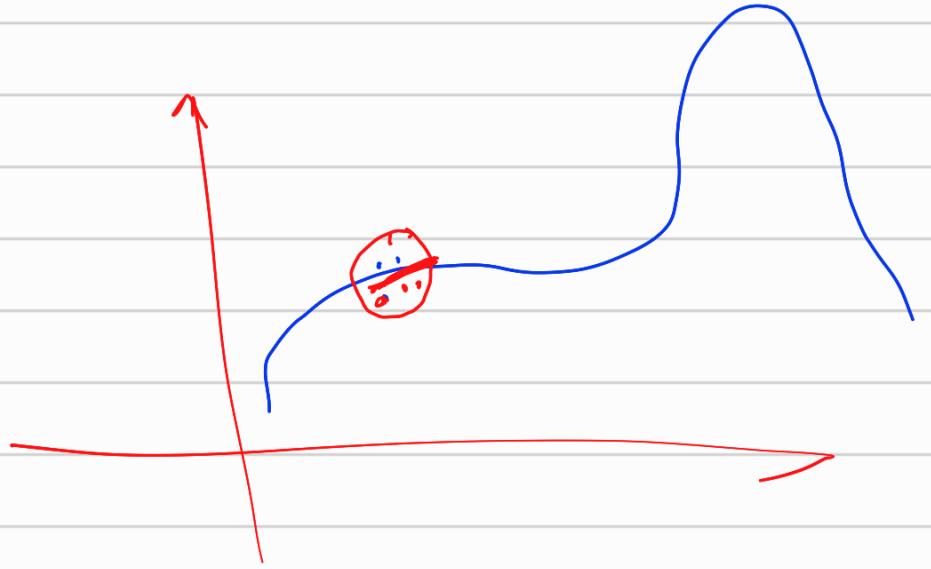
## Local interpretable Model Agnostic:

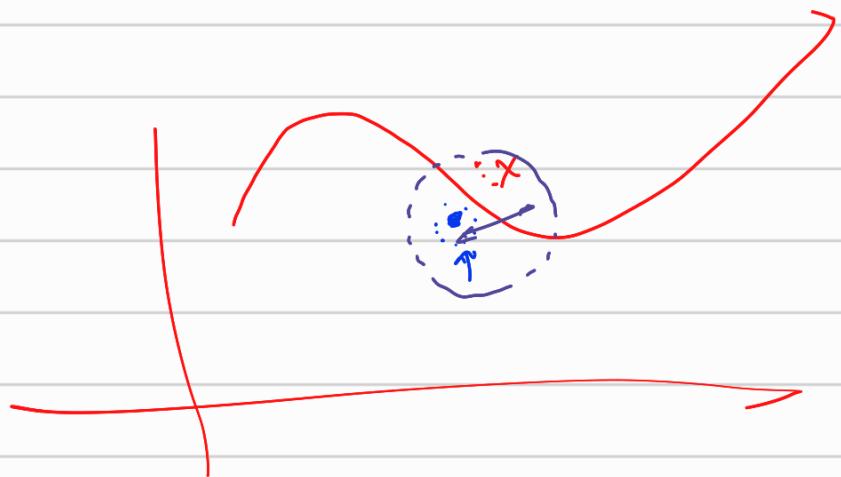
LIME



- We will be given a single point (purple Point) for prediction
- We will sample random points around the given point within a epsilon radius
- All these random points will be passed through the main model (red curve/NN mode)
- thereby we would have created a sample data of green points and their corresponding predictions
- This above data becomes a training data for me.
- I select a simple (explainable model like DT/Logistic Regression) to fit the above data
- The feature importance/coefficients become the explainability factors of the given purple point







Shapley '0 /

