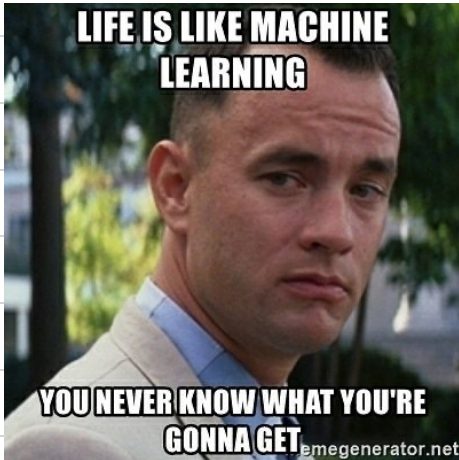


Session 12

MODEL INTERPRETABILITY

March 04,
2024



AGENDA

① TIME

② SHAP

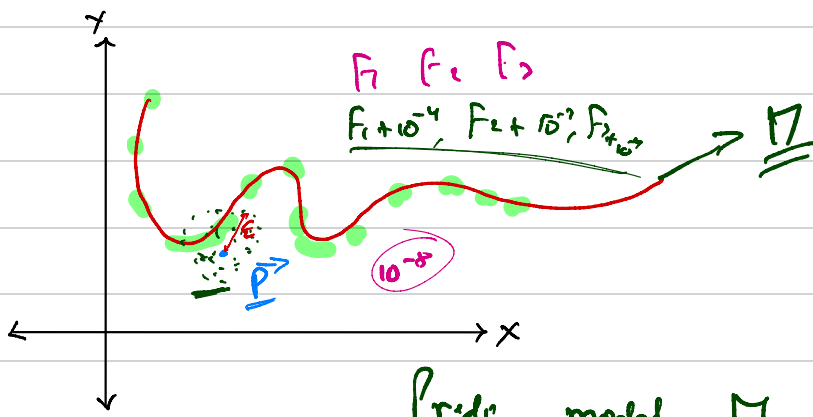
Why do we need it?

① Hedgefox: Mathematical model \rightarrow explain

<u>Feature Importance</u>	<u>Model Explainability</u>
look at global level <u>Entire dataset</u>	look at <u>local level</u> sample level

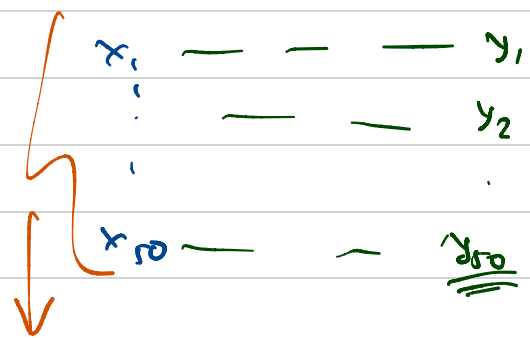
Home \rightarrow looks at local level
Shop \rightarrow " " global "

LIME



Predict model M on new generated Points

1, 10 \rightarrow some nois \rightarrow 50, 10



Fit a Linear model + lasso

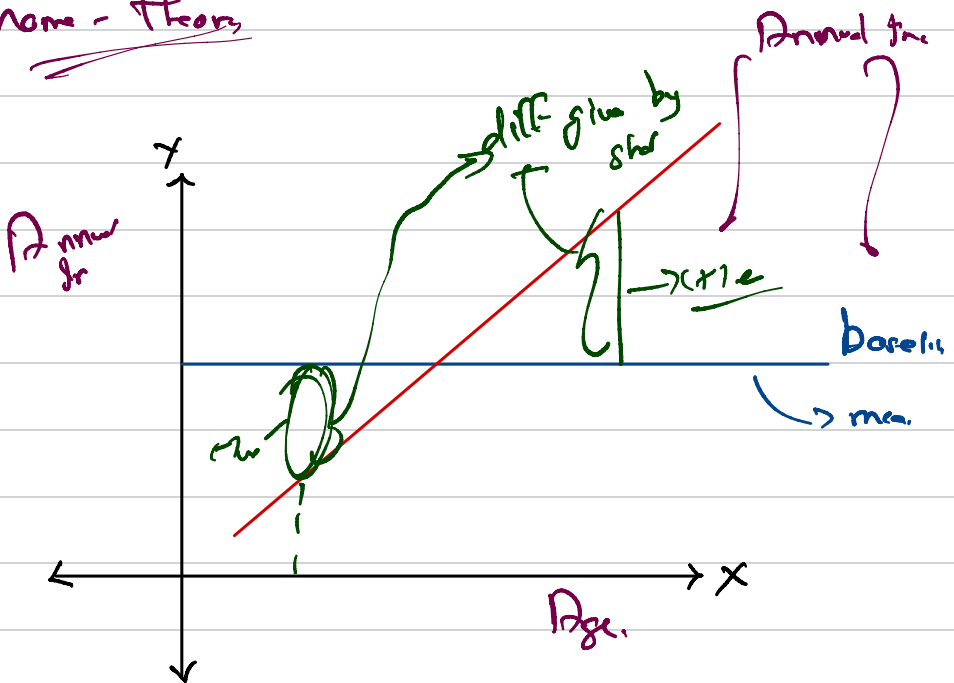
coeff of my linear model \rightarrow help me understand why my model

is giving that output

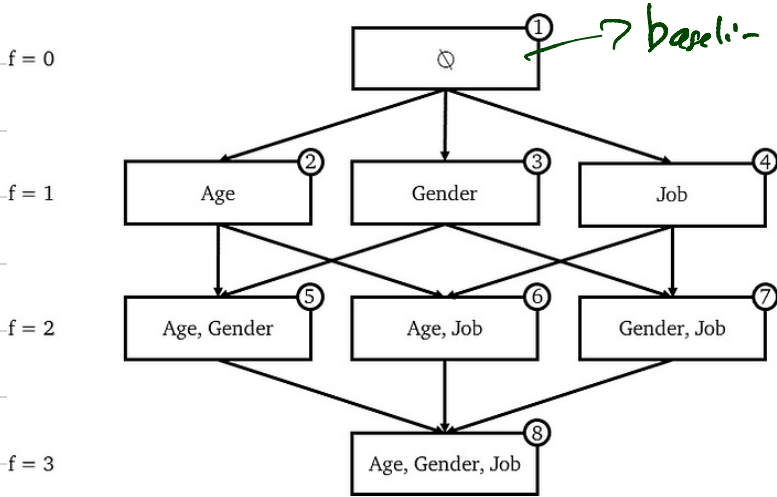
→ Points closer to blue point will be given higher weightage

SHAP

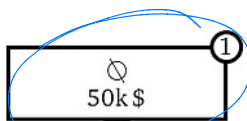
Game theory



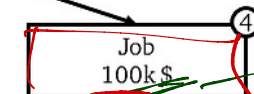
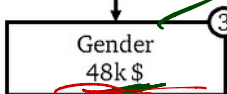
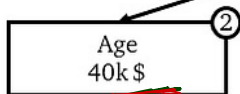
NN
mod



f = 0

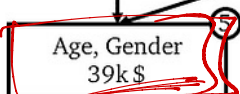


f = 1



50k

f = 2



f = 3



40k - 50k

-212

-1012

-9k -1512 } Age

-212

Age Gender Job } NN → Output
Predict

Split ← Train data mean - Age -
Gender -
Job

10

1 →

Age, Gen, Job
≡ Av., Av. → Output

45012, 40012

50% ave
↑
acron
res

Age Gen Del

avg av 11

Baseline + 201^ \rightarrow 201L

Shapley additive Exploratio