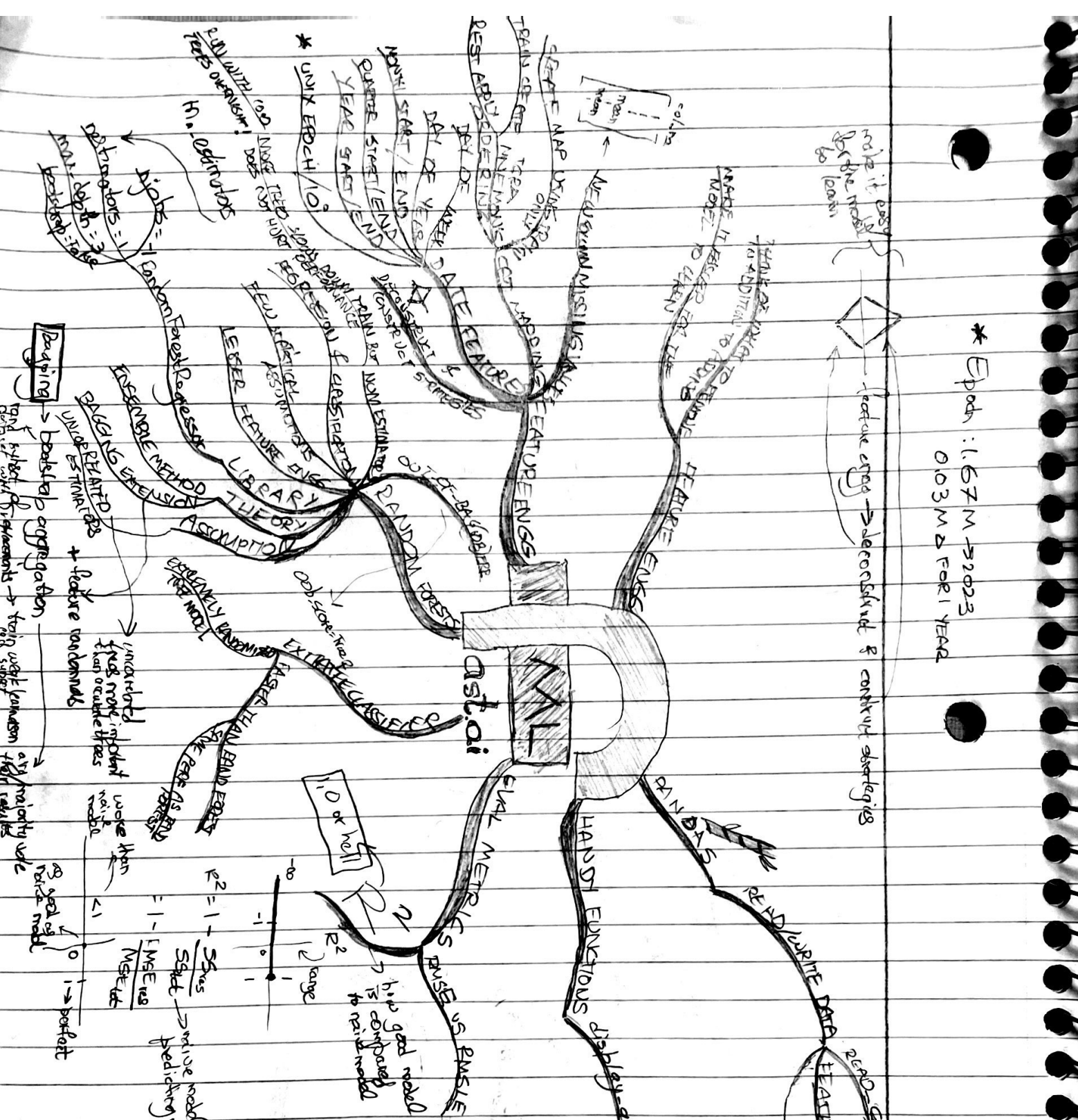


* Epoch : 1.67M → 2023
0.03M → FOR 1 YEAR

feature eng → decrease bias & control overfitting



RMSE vs RMSLE

RMSE
scale sensitive
outliers sensitive

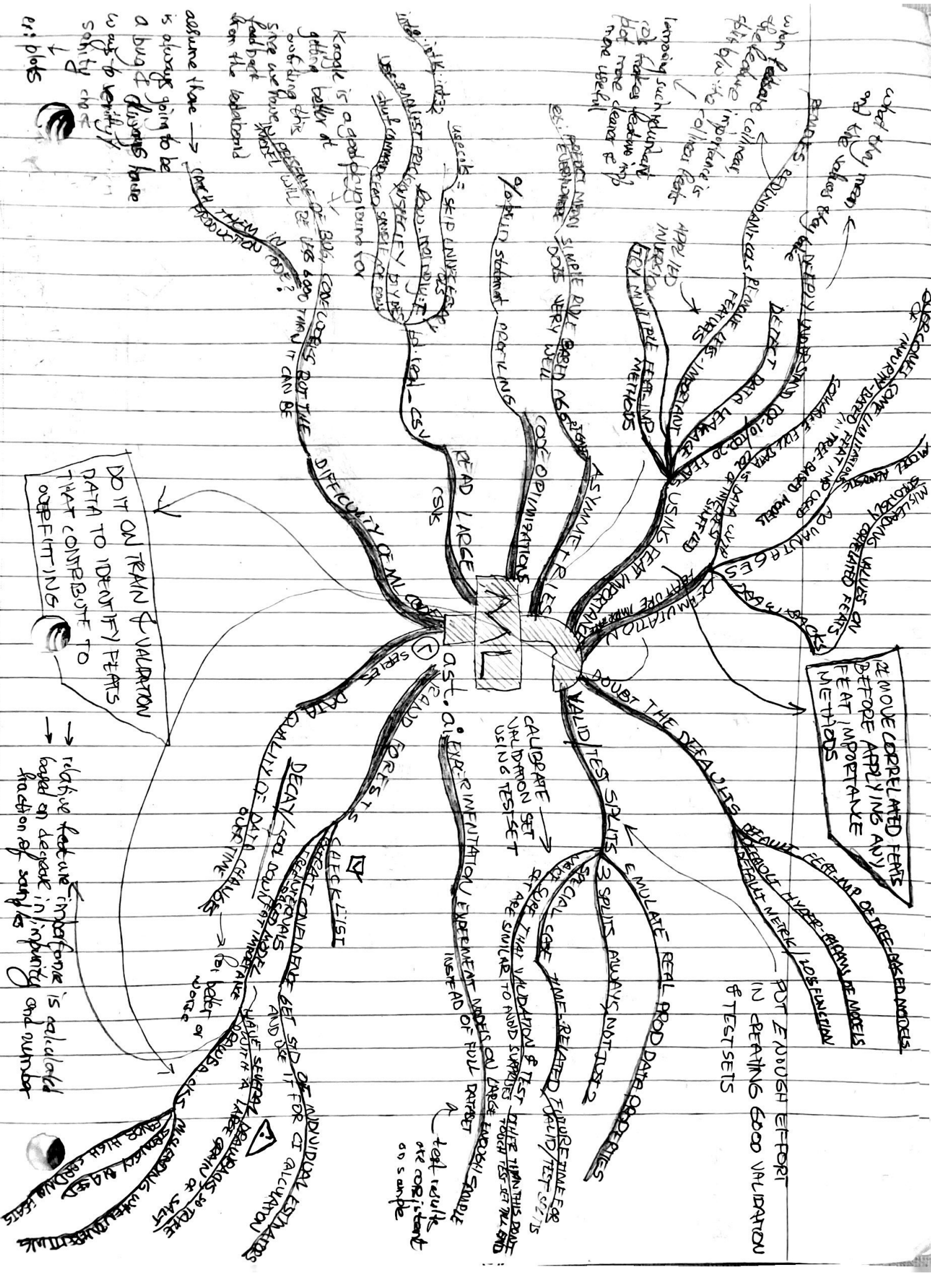
RMSLE
scale insensitive
penalize underfit more
more tolerant of outliers

$RMSE = \sqrt{\frac{1}{N} \sum (y_i - \hat{y}_i)^2}$

$RMSLE = \sqrt{\frac{1}{N} \sum (\log(\hat{y}_i) - \log(y_i))^2}$

For small x , $\log(x) \approx x$

$y \rightarrow y$ and small: RMSE & RMSLE is same
one is large: RMSE > RMSLE
both large: RMSE > RMSLE
(RMSE > RMSLE because RMSE is more sensitive)



when feature collinear, the feature importance is not clear
 removing such redundant features makes feature importance more useful

KAGGLE IS A GREAT PLACE TO FIND DATA
 GETTING BETTER AT FINDING THE BEST WILL BE USEFUL THAN IT CAN BE
 DON'T USE ALL DATA
 DISTRIBUTION CHANGE
 SYNTHETIC DATA
 CROSS-VALIDATION

DENSITY UNDERSTAND TOP-10 TO TOP-20 FEATURES USING SHAP OR SHAPLETS
 DETECT DATA LEAKAGE
 APPLY L2 REGULARIZATION TO PREVENT OVERFITTING
 INTERPRETABILITY OF SPLIT METHODS
 ASYNCHRONOUS FEAT IMPORTANCE
 DETERMINATION OF FEATURE IMPORTANCE
 PREDICT WITH SIMPLE RULES FOR HARD ASSESSMENTS
 EASY TO IMPLEMENT BUT VERY DULL
 % BUILD STATEMENT PORTALUS
 CODE OPTIMIZATIONS
 READ LARGE DATA

VALID TEST SPLITS
 3 SPLITS ALWAYS NOTICE
 EMULATE REAL WORLD DATA DISTRIBUTION
 CALIBRATE SET VALIDATION USING TEST SET
 SPECIAL CASE: TIME-RELATED FEATURE TIME FOR CALIBRATE SET VALIDATION & TEST SETS
 THESE TIME RELATED TEST SETS ARE SUITABLE TO HARD SPLITTING
 SIMULATE REAL WORLD DATA DISTRIBUTION
 3 SPLITS ALWAYS NOTICE
 EMULATE REAL WORLD DATA DISTRIBUTION

REMOVE CORRELATED FEATURES BEFORE APPLYING ANY FEAT IMPORTANCE METHODS
 BEWARE OF DEPENDENT VARIABLES
 HYPER-PLANE OF BEST FIT
 BEWARE OF MULTICOLLINEARITY
 DON'T DO OVERFITTING
 IN GETTING GOOD VALIDATION & TEST SETS

DECAY OF DATA QUALITY
 CHECKLIST FOR DATA QUALITY
 DECAY OF DATA QUALITY OVER TIME
 CHECKLIST FOR DATA QUALITY
 DECAY OF DATA QUALITY OVER TIME

EXPERIMENTATION/ EXPERIMENTATION
 ASSESS MODEL PERFORMANCE
 TEST RESULTS ARE ROBUST ON SAMPLE

TEST RESULTS ARE ROBUST ON SAMPLE
 CHECKLIST FOR DATA QUALITY

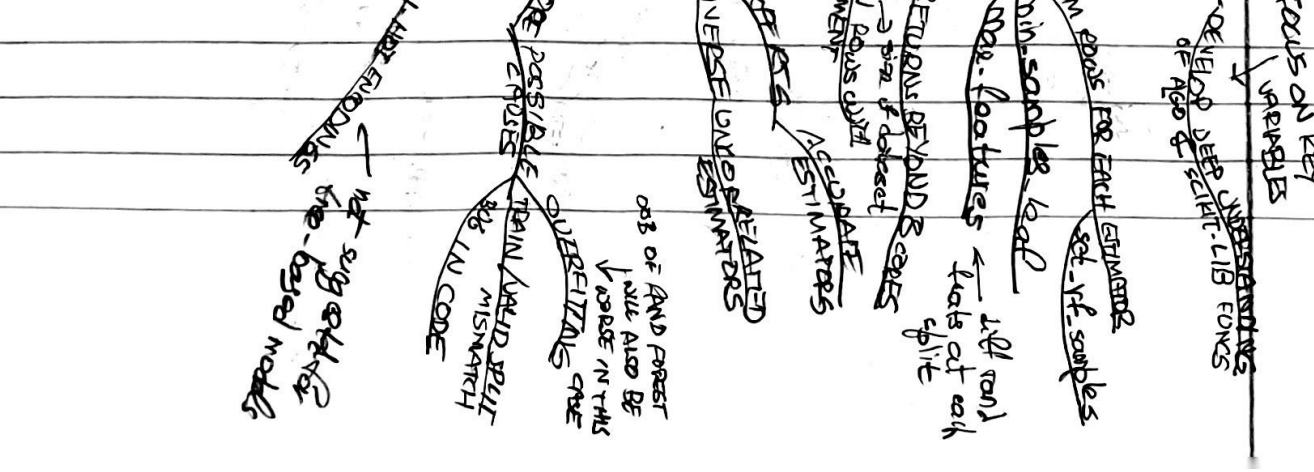
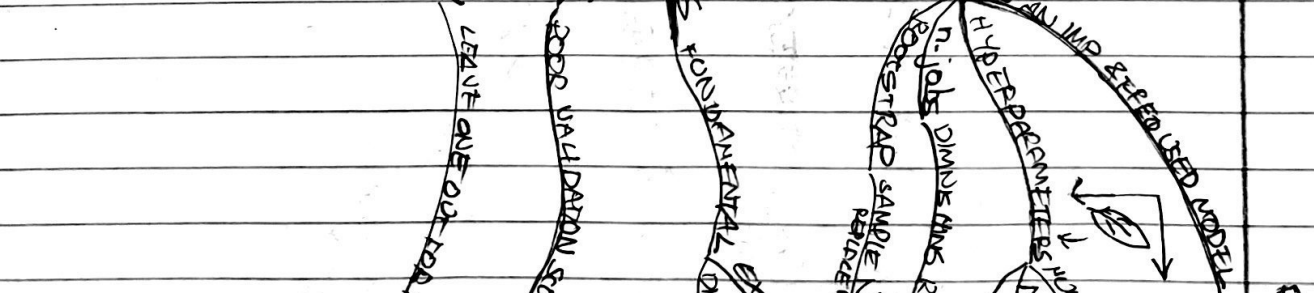
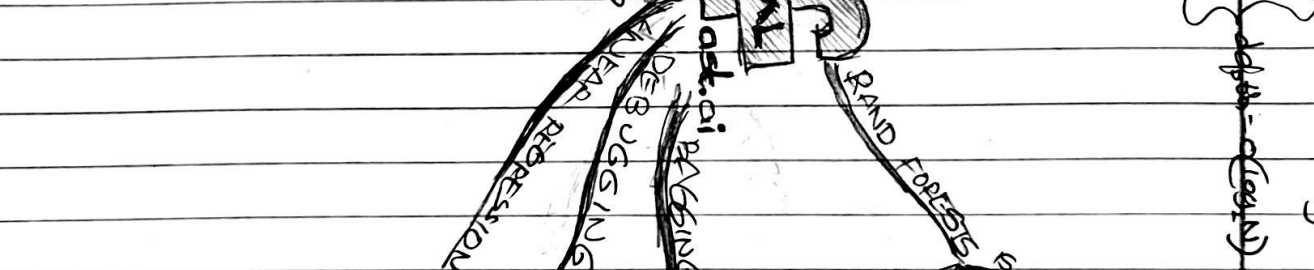
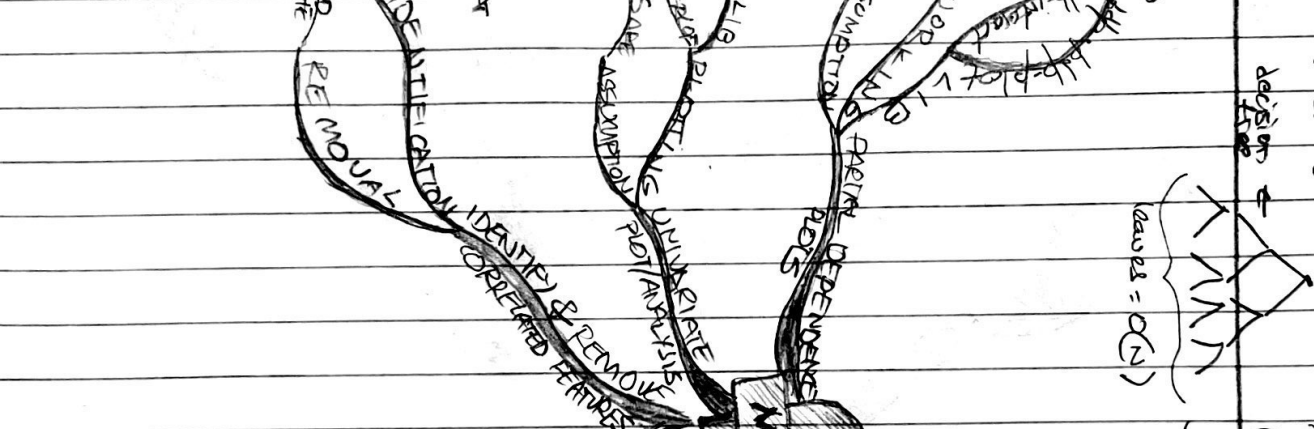
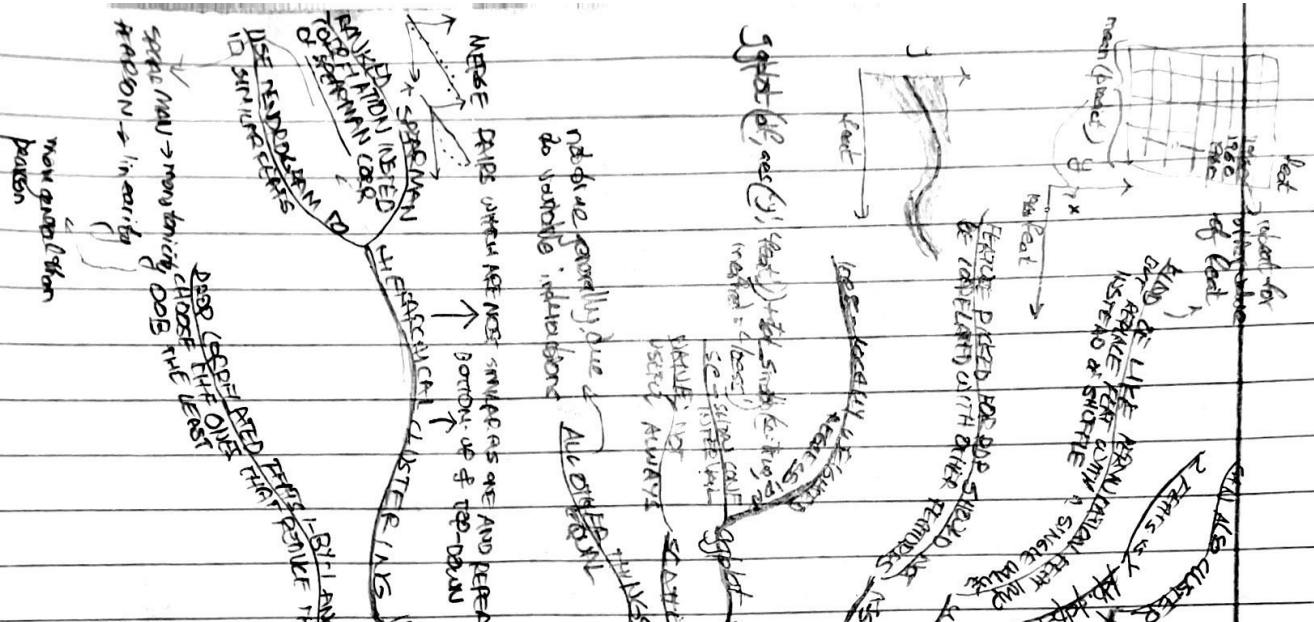
CHECKLIST FOR DATA QUALITY
 DECAY OF DATA QUALITY OVER TIME

CHECKLIST FOR DATA QUALITY
 DECAY OF DATA QUALITY OVER TIME

DON'T ON TRAIN & VALIDATION DATA TO IDENTIFY FEATS THAT CONTRIBUTE TO PROFITING

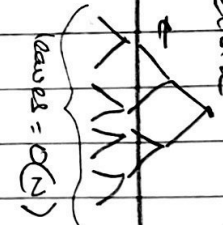
relative feature importance is calculated based on degree of importance and number fraction of samples

CHECKLIST FOR DATA QUALITY
 DECAY OF DATA QUALITY OVER TIME



$N = \text{dataset size}$

Decision Tree

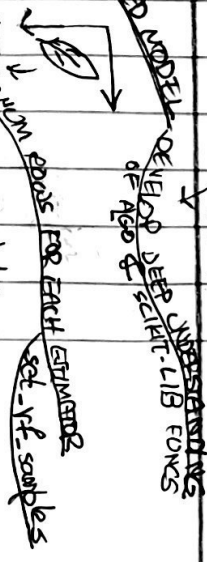


depth = $O(\log N)$

low feature eng

high classical

FOCUS ON KEY VARIABLES



MIN-SAMPLES LEAF

N-JOBS DIVIDE THIS RETURNS BEYOND SCORES

POSTSTRAP SAMPLE N ROWS WITH REPLACEMENT

ESTIMATE ESTIMATORS

LEAF ONLY

SPLIT

MSE OF AND FOREST

MAY BE ALSO BE

UPGRADE IN THIS

OVERFITTING

TRAIN AND TEST

MISMATCH

BIG IN CODE

MAY BE ALSO BE

UPGRADE IN THIS

OVERFITTING

TRAIN AND TEST

MISMATCH

BIG IN CODE