

## Auto ML :-

→ repetitive ml tasks  
can be automated.

→ Sagemaker Autopilot (AutoML)

## Model Building Challenges:-

→ steps involved in building ML model requires multiple iterations that can often result in increased time to market.  
(Ability to reduce the time to market).

→ Lack of ML skillsets (may be from the team)

→ Iterations (data transformations, differ algorithms, hyper-parameter tuning etc.) to find a good model.

→ Optimize scarce resources and skillsets.

AutonML:- automates the steps in ML workflow

## Benefits :-

- 1) Reduce the time to market by automating data transformation, feature engr, tuning
- 2) Non-data science folks can use the ML models
- 3) Iterate quickly using ML and automation to perform majority of the tasks in model building workflows.
- 4) AutonML + cloud train/tune model in parallel.

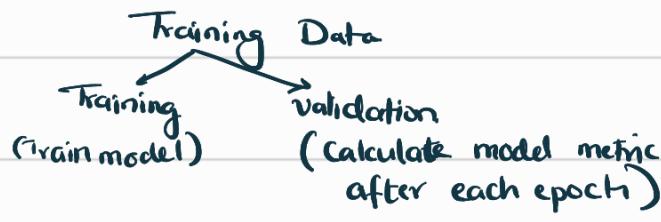
## Text Classification Data preparation :-

✗ No one-hot encoding.

Tokenization → stop words removal → convert into vectors ...

Also, we have to deal with imbalance.

Use synthetic data, resampling techniques,  
 choose different metric (say weighted F<sub>1</sub> score),  
 different algorithm (xgboost - grad descent  
 boosted decision trees)



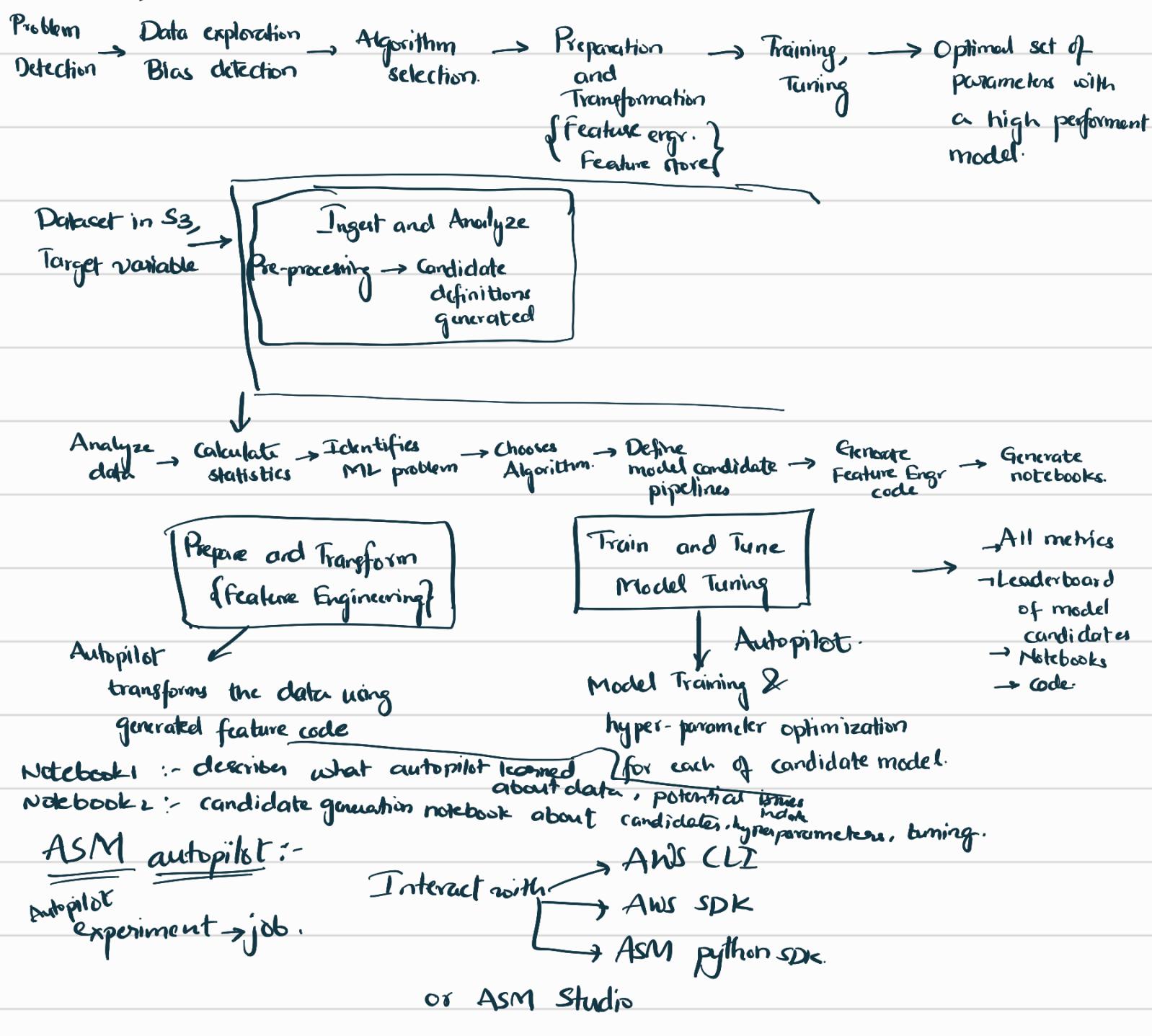
Training and tuning is time and resource consuming and hence distributed computing (scale out) plays an important benefit.

AutoML aims at automating the process of building a model.

- builds models without ML expertise
- models at scale without human intervention
- automate majority of work, then tweak  
(data cleaning, feature engr., feature selection)

Level of transparency and control is also be achieved in most cases. In some cases only the best performing model is returned, but it is hard to reproduce manually or explain.

### ASM autopilot:



ASM SDK on top of sklearn.

Text data → Features

## TFIDF

ASIM autopilot uses multi-column tfidf vectorizer and ASIM autopilot will automatically tune

MulticolumnTfidfVectorizer parameters.

How relevant a word is to a doc or collection of documents.

All text into a bag and sample from it without replacement. Form a vector of counts.

Then how important a word is to a document, relative to others in the corpus.

Word importance → statistical measure

- increases proportionally to number of times a word appears in the doc.
- decreases by the frequency of the word in the corpus.

t - term / word

d - document

D - corpus

$$tf(t, d) = \frac{f_{t,d}}{\sum_{t' \in D} f_{t',d}}$$

whole corpus →  $idf(t, D) = \log \left( \frac{|D|}{|\{d \in D : t \in d\}|} \right)$

say food appears in doc, 10 times of <sup>200 words</sup> and the word food is in 1000 docs of <sup>10<sup>6</sup> docs</sup>

then

$$idf(\text{food}, \text{corpus}) = \log \left[ \frac{10^6}{1000} \right] = 3$$

$$df(\text{food}, \text{doc}) = \frac{10}{200} = 0.05$$

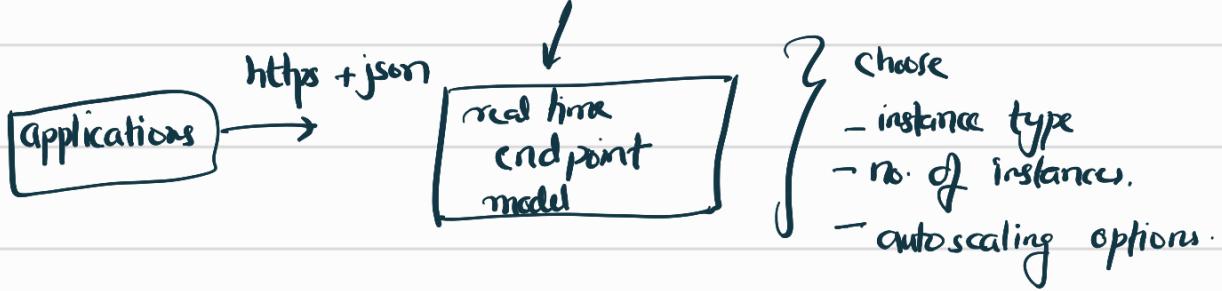
$$tfidf(\text{food}) = 0.05 \times 3 = 0.15$$

Once you find the model (best performant), if we need to deploy to production  
(prediction in real time)

- Optimized for low latency of model predictions.

Model serving stack → trained model + hosting stack

{ proxy, webserver with trained code }



## Inference Pipeline model

### Data Transformation Container

Container built from the model we selected and trained during the data transformer sections

### Algorithm Container

a container build from the trained model we selected above from the last HPO training job.

### Inverse Label Transform Container.

a container that converts intermediate numerical prediction to the label value.