

AI - ENGINEERING BLUE - PRINT

- Dream in years
- Plan in months
- Evaluate in weeks
- Ship Daily

• Prototype for 1x

• Build for 20x

• Engineer for 200x

- What's required to cut the time-line in half
- What needs to be done to double the impact

Video-series - 1 :-

‘ Machine-learning Basics (Supervised Learning) ’

$$y = f(x) \rightarrow \text{It's a form of equation}$$

$y \rightarrow$ output
 $f()$ → function form
 $x \rightarrow$ Input

2) y - output y stock → millions
 accuracy - out of sample

time-series - continuous

$$\text{Mean Squared Error. (MSE)} = \frac{1}{n} \sum (y_i - \hat{y}_i)^2$$

$$\text{MAE (Mean absolute Error)} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Classification

$$\text{Accuracy} = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}}$$

$$\text{RMSE} = \text{Root Mean Squared Error}$$

$$= \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

$$\text{Precision} = \frac{\text{True Positive}}{\text{true positive} + \text{false positive}}$$

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

$$\text{Recall} = \frac{\text{True Negative}}{\text{true positives} + \text{false negatives}}$$

$$\text{MAPE} = \text{Mean Absolute percentage}$$

$$F_1\text{-score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$\text{Error (MAPE)} = \frac{1}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right| \times 100$$

ROC - AUC (Reciever Operating Characteristic)

Measures the ability of classifier to distinguish b/w classes.

A higher AUC indicates better model performance.

• Logarithmic Loss (Log loss) : $\rightarrow 0$ and 1

performance of classification model whose output $\rightarrow 0$ or 1

$$\text{Log Loss} = \frac{1}{n} \sum_{i=1}^n [y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)]$$

2. Test the model

6. Train the model

5. Build and Compile the model

- CNN, P. f.
- Pre-trained Models like VGG16, ResNet, or MobileNet

Independent
 (x)

education \rightarrow

Salary

1. Collect Images {Source}
2. Label Images \rightarrow cat \rightarrow non-cat
3. Pre-process Image
 1. Resize Images (128×128)

2. Normalization (scaling pixel values between 0 and 1 by dividing by 255)

\hookrightarrow can be fine tuned to specific task using transfer learning

2) $f(x)$

\rightarrow salary.

$$\text{Salary} = \beta_1 \text{Education} + \beta_0$$

\hookrightarrow but estimate.

1) Defining variables

(X) : Education level

• years of education

• highest degree

(Y) : Salary (Continuous variable)

2) Formulate the regression model

$$Y = \beta_0 + \beta_1 X + \epsilon$$



If family - background comes into picture

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$$

Y = dependent variable (salary)

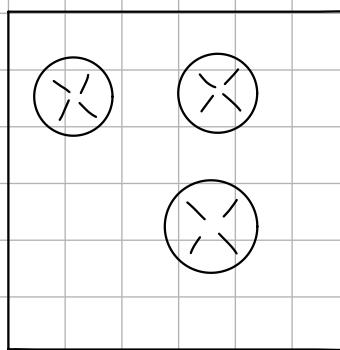
X = independent var (education)

β_0 = intercept (salary when edu is zero)

β_1 = slope-coefficient (the change in salary for each additional unit of edu)

ϵ → error-term

3) X-Inputs Optimization :- are central to ml, as they involve finding the best parameters for a model to minimize (or maximize) a certain objective function. → error or loss of model, and optimizing it means improving the model's performance on a given task



Case (i) :- Handing out coupons leads to ↑ customer spending or wrong?)

(benchmark)

RCT :- Randomized controlled trial bias spend
coupon (FOMO) cost

wrong
• customer who get coupon are likely to spend anyway

$$F_b \quad Y \quad T \quad P$$

$$y = \left(\frac{111111}{p-v-e} \right) \quad (1-q)$$

sign
↓
 $\frac{1}{20}$

1 - 8

wine-quality

(a)

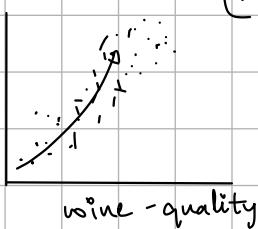
alcohol

texture

sulphur-content

nitrogen

alcohol



(Wrong!)

Bias → refers to a systematic error that occurs when a model consistently makes errors.

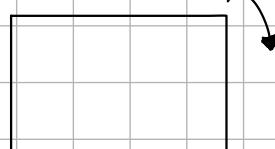
Unbiased → in a statistical sense, makes predictions that, on avg, are correct. It does not overpredict or underpredict the true values across many data samples.

Here lies
the dragon
in statistics

(Case iii)

risk of relapse

x wrong

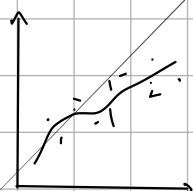


(Model) - NN - 99%

Ideal - scenario

$$y = f(x)$$

Discriminative Models

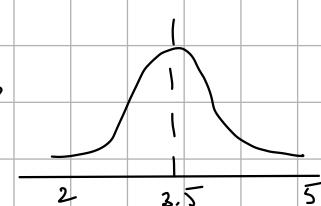
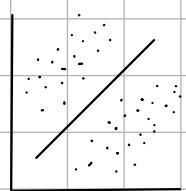


frequentist - setting

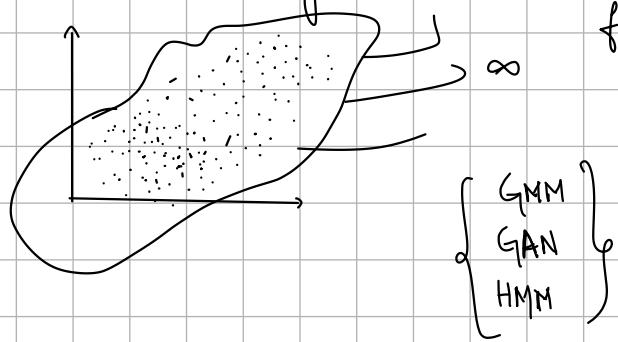
$$y = \beta_0 + \beta_1 x + \epsilon$$

Probabilistic setting
 $P(y|x)$

$$\beta \rightarrow$$



Generative Modeling



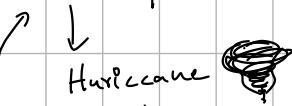
figuring out data

ML-pipeline

blue-print

(USA)

resource - equipment



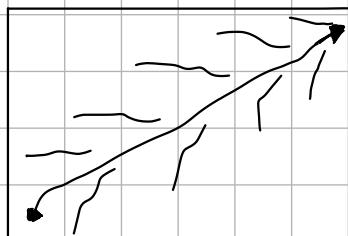
path of hurricane

we can pull the resource equipment

8 → 3 samples

1000 → 3997 → 1-7

① What is the business usecase ?

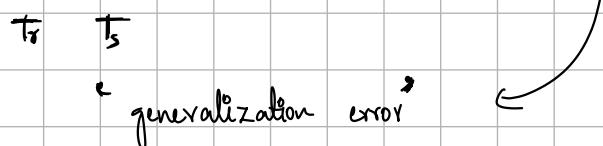


- i) Infinite data would I be able to solve the problem
 → ii) Given that I will have enough variation of the data

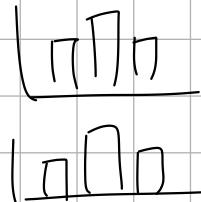
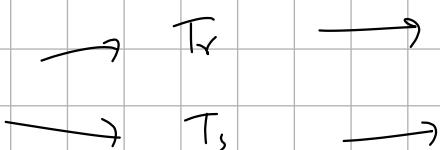
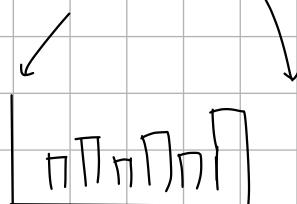
3) Define Error-metric $T_S \rightarrow T_S \rightarrow MAPE, MSE, RMSE$

classification → Accuracy, precision, recall, --- f2-score

Aim of modelling process to min OOS error.

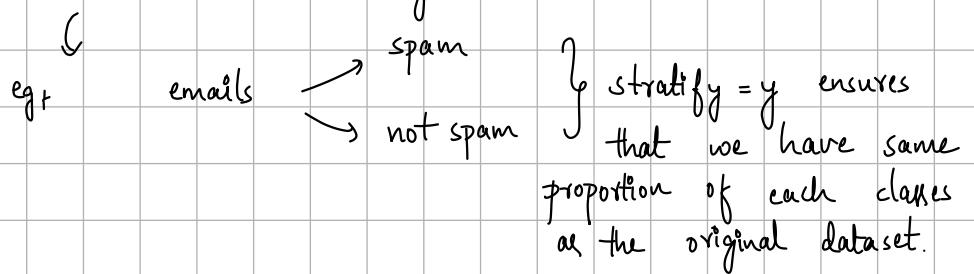


Balanced Split



2) Proper Data splitting :- 1) Stratified Sampling (Classification)

ensure distribution of classes in the training and test sets is similar to that in the original dataset.



Regression example : Random sampling → output variable is continuous
: Cross-validation → assessing model's performance on different subsets of the data which helps reduce generalization error.

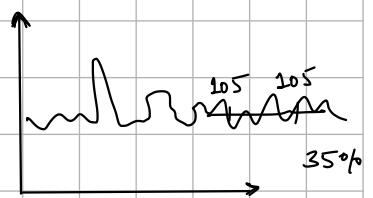
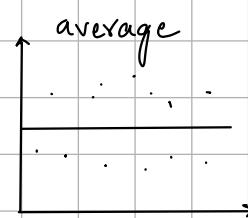
3) Determine the task → whether it is regression or classification

6) Setting the base-line

(no-covid) (covid) (baseline)
95% 5% (95% accuracy)

{95% accuracy}

min-accuracy
(17% MAPE)



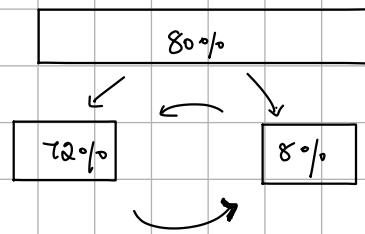
35% not beating your base line

if the client doesn't provide baseline, we set up our own baseline

(cold-storage)

Train	Test
80%	20%

7) Apply models on train and cv-set.



- * Logistic regression
- * Random Forest
- * SVM
- * Gradient Boosting Method

Cross-validation (CV)

8) Prediction interval }
Confidence interval }

1) Books to read → * Rob J Hyndman and George Athanasopoulos

Forecasting Principles and Practice (2nd ed)