



Table Of Contents

1. Steps in ML Project
2. Exploring Step 1 Data Collection
3. Exploring Step 2 Data Preparation
 - 3.1 Exploratory Data Analysis
 - 3.2 Data Preprocessing
 - Feature Imputation
 - Feature Encoding
 - Feature Normalization
 - Feature Engineering
 - Feature Selection
 - 3.3 Data Splitting

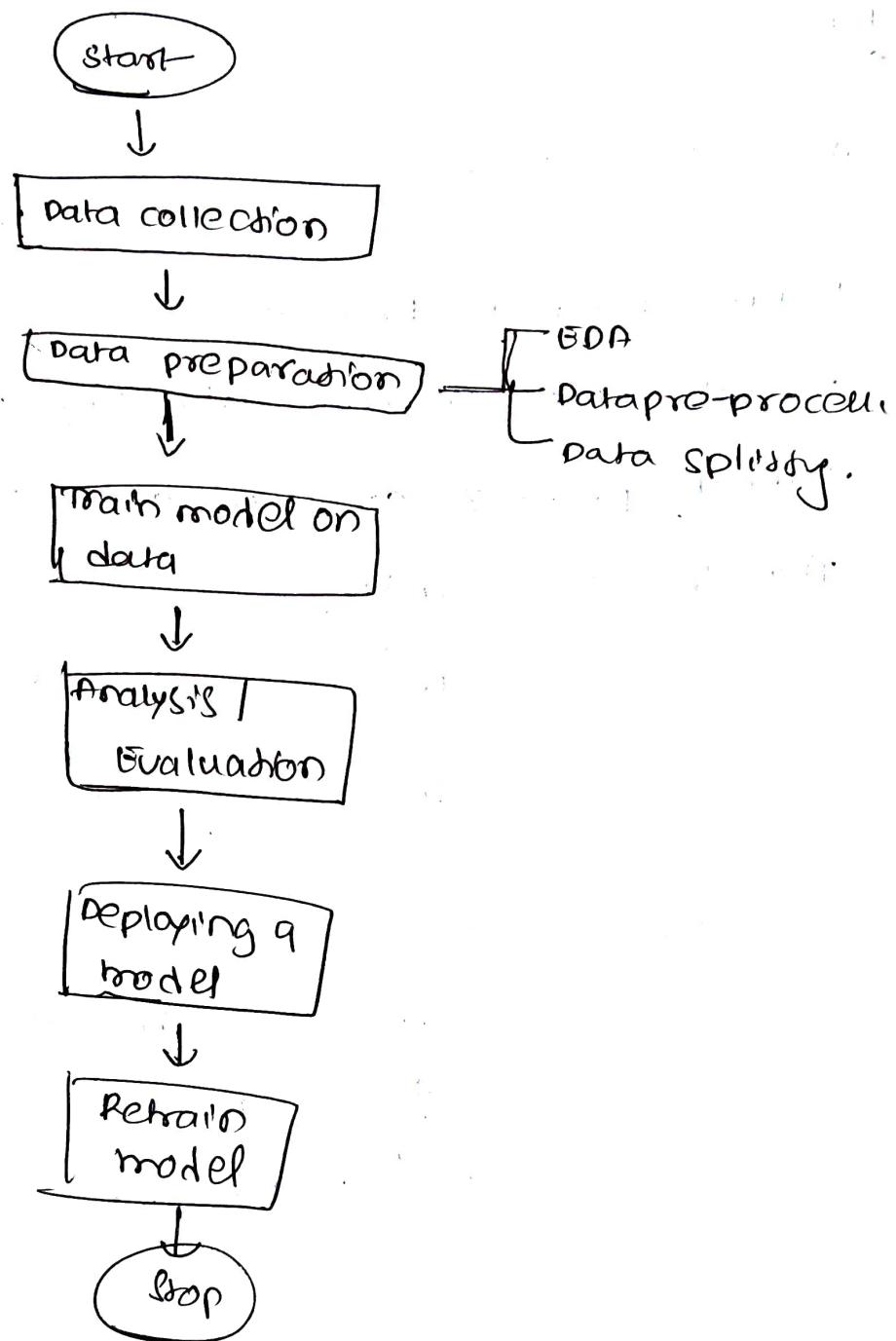
4. Exploring Step 3 - Train Model on Dataset

- 4.1 Types of Learning
- 4.2 Under Fitting and OverFitting
- 4.3 Regularization techniques
- 4.4 Hyperparameter Tuning

5. Exploring Step 4 - Evaluation of a Model

- 5.1 Evaluation Metrics
- 5.2 Confusion Matrix
- 5.3 Recall/Sensitivity
- 5.4 Precision
- 5.5 Specificity
- 5.6 F1 Score
- 5.7 AUC and ROC Curve
- 5.8 Analysis of a Model

ML Road map (Steps in machine learning project)



* Understand the process :-

Step ①

Data collection

a) Questions to ask?

- 1) what kind of problem are we trying to solve?
- 2) what data sources already exist?
- 3) what privacy concerns are there?
- 4) Is the data public?
- 5) where should we store the data?

b) understanding the types of data

i) structured data?

data which appear to be tabulated. Form
can contain diff types of data.
we
i) nominal / categorical.
ii) numerical
iii) ordinal
iv) time series

ii) unstructured data?

data with no rigid structure.

(Image, video, natural language text, speech,

Step(2)

Data preparation

It has 3 main process.

- a) Exploratory data analysis (EDA), understanding data.
- b) Data preprocessing, preparing your data, to be modelled.
- c) Data splitting.

a) EDA, understanding the data

- 1) what are feature variables (Input) and the target variables (output)?

[e.g. For predicting heart diseases, feature variable may be a person's age, weight, heart rate etc. and target variable is whether or not they have heart disease]

- 2) what kind of data do you have?
structured, unstructured, categorical / numerical.
- 3) Are there missing values? Should you remove them / fill them with feature imputation?
- 4) what are outliers?
How many of them are there?
Are they out by much ($3 + 2\sigma$)?
Why are they there?
- 5) Are those questions you could ask a domain expert about the data?

b) Data Preprocessing, preparing your data to be modelled

① Feature Imputation

(Filling missing values)

→ A ML model can't learn on data that isn't there.

i) Single Imputation: Fill with mean / median or column

ii) multiple imputation:

and fill with what your model finds.

iii) KNN (K-nearest neighbor):

Fill data with a value from another example which is similar.

iv) many more such as random imputation,
last observation carried forward (for time series),
moving window, most frequent.

∴ there are ways to handle missing values.

② Feature encoding

(Turning values into numbers)

→ A ML model requires all values to be numeric

i) OneHotEncoding:-

Turn all unique values into 1's or 0's & 1's where the target value is 1 & the rest are 0's.

e.g. car colors green, red, blue,

a green car's color feature would be [1, 0, 0]

& a red one would be [0, 1, 0]

& a blue one [0, 0, 1]

ii) Label Encoder :-

Turn labels into distinct numerical values

e.g.

If your target variables are diff.

animal as dog, cat, bird then 0, 1, 2 respectively

iii) Embedding encoding :-

Learn a representation amongst all the diff. data points.

e.g.

A language model is a representation of how diff. words relate to each other. Embeddings are also becoming more widely available for structured data.

③ Feature normalization (scaling) or standardization :-

when your numerical variables are on diff scales. (like size of land b/w 1000 - 20000 sq ft). In such cases some ML models don't perform well. Scaling / standardization help to fix this.

i) Feature scaling (normalization) :-

shifts your values so they always appear b/w 0 & 1.

This is done by subtracting min value & divide (value b/w 0 & 1 but numerical relationship is still there) with range.

ii) Feature standardization :-

Standardized all values so they have a mean of 0 & variance of 1.

This is done by subtracting mean & divide with SD.

values doesn't end up b/w 0-1 (st. var.).

Standardization is more robust to outliers than scaling.

④ Feature Engineering

→ Transform data into potentially more meaningful representations by adding in domain knowledge.

i) decompose

Like, turn a date (such as 2020-06-18 to 2016:
into hour-dr-day, day-dr-week, day-dr-month, ts-holt
etc.

ii) discretization

(Turning larger groups to smaller groups)

For numerical variables,

like age you may want to move it into bucket such as over-50 (under-50, 21-30, 31-40 etc.) this process is known as BINNING.

For categorical variables,

such as car color, this may mean combining colors such as (light-green); (dark-green), (lime green) into a single green

iii) crossing & interaction features

(combining two/more features)

→ depending on situation sometimes

adding / subtracting two features could help.

iv) Indicator features

like using
is-Paid-traffic
for Paid

[using other parts of data to
indicate something potentially significant]

+ Create a X-missing feature for whenever
a column X contains a missing value.

⑤ Feature Selection :-

- Selecting the most valuable feature or your dataset to model.
- Potentially reducing overfitting & training time (less overall data & less redundant data to train on) & improving accuracy.

i) Dimensionality reduction:-

A common dimensionality reduction method, PCA takes a larger no. of dimension (features) & uses linear algebra to reduce them to less dimensions. For example, say you have 10 numerical features. You could run PCA to reduce it to 3.

ii) Feature Importance :-

Fit a model to a set of data, then inspect which features were most important to the results, remove the least important ones. (can use "TPOT" for this)

⑥ Dealing with Imbalance :-

Does your data have 10,000 examples or one day, but only 100 examples the other?

- collect more data (if possible)
- use scikit-learn-contrib imbalanced-learn package
- use SMOTE (synthetic minority over-sampling technique) which creates synthetic samples of your minor class to try & level the playing field.

→ So this is all about ^{data} Preprocessing.

c) Data Splitting

1) training set (usually 90-80% of data):

model learns on this.

2) validation set (typically 10-15% of data):

model hyperparameters are tuned on this.

3) test set (usually 10-15%): (Dev Set)

model's final performance is evaluated on this.

④ Do not use this dataset to tune the model

a. why shouldn't we tune on test set?

A.

To do the regularization we have to take 100 diff. models & try on the test data. And so last we finds (suppose) the best hyperparameter that produce a model with low generalization error, say just 5%.

And you launch it to production & found it produces 15% errors.

④ Now, here, is what happened as you trying every hyperparameter on the same test data eventually your model learned the test data & fitted to it thus reducing generalization errors but not improving models.

This is called "overfitting".

→ So, to avoid we take validation set / cross-validation (more on)

④

Step③

Train model on dataset

It has 3 main process.

a) choose an algorithm

b) overfit the model

c) Reduce overfitting with regularization.

a) Algorithm can be choosed by understanding supervised, unsupervised (Explained previous)

b) Type of learning

i) Batch Learning

All of your data exists in a big static warehouse & you train a model on it.

You may train a new model once per month once you get new data. Learning may take a while & isn't done often. Runs in production without learning (can be retrained).

ii) Online Learning

Your data is constantly being updated. You constantly train new models on it. Each learning step is usually fast & cheap.

Runs in production & learns continuously.

iii) transfer learning

Take the knowledge one model has learned & use it with your own. Give you the ability to leverage SOTA (state of the art) models for your own problems.
Helpful if you don't have much data / want compute resources.

iv) Active learning

Also referred to as "human in the loop" learning. A human expert interacts with a model & provides updated labels for samples where the model is most uncertain about.

v) Ensembling

Not really a form of learning, more combining algorithms which have already learned in some way to get better results.

Underfitting

Happens when your model doesn't perform as well as you'd like on your data. Try training for a longer & more advance model.

overfitting

Happens when your validation loss (how your model is performing on the validation dataset, lower is better) starts to increase, or if you don't have a validation set,

* Happens when the model performs far better on the "training set" than on the "test set".
(e.g. 99% accuracy on training set, 67% accuracy on test set) → overfitting the training set.

Fix through various regularization techniques.

Regularization [Techniques to prevent/reduce overfitting.]

i) L1 (Lasso) & L2 (ridge) regularization.

→ L1 regularization sets unneeded feature coefficients to 0.

→ L2 constrains a model's features (won't set them to 0)

ii) Dropout:

Randomly remove parts of your model.

So the replacement has to become better.

iii) Early stopping:

Stop your model from training before the validation loss starts to increase too much / more generally, any other metric has stopped improving. Early stopping usually implemented in the form of a model callback.

iv) Data Augmentation

- manipulate your dataset in
artificial ways to make it harder to learn.
- For example, if you're dealing with images, randomly rotate, skew, flip & adjust the height of your images.
 - this makes your model have to learn similar patterns across different styles of the same image (harder).
 - use functions like `ImageDataGenerator` in keras or transforms in torchvision.

v) Batch Normalization

standardize inputs as well as adding two parameters (beta, how much to offset the parameter for each layer & epsilon to avoid division by zero) before they go into next layer. This often results in faster training speed since the optimizer has less parameter to update.

may be a replacement for dropout in some networks.

Hyperparameter Tuning

→ Run a bunch of experiments with different model settings & see which works best.

a) Setting a learning rate:

(often the most important hyperparameter)
Generally

High Learning rate = algorithm rapidly adapts to new data

Low learning rate = algorithm adapts slower to new data
(esp. for transfer learning)

i) Finding the optimal learning rate

Train the model for a few hundred iterations starting with a very low learning rate (10^{-6}) & slowly increase it to every large value (e.g. 1). Then plot the loss vs learning rate (using a log-scale for learning rate), you should see a U-shaped curve, the optimal learning rate is about 1/2 notch to the left of the bottom of u-curve. (P326 - Handwritten)

ii) Learning rate scheduling involves decreasing the learning rate slowly as model learns more (get closer to convergence). [see Adam optimizer]

iii) cyclic learning rates

Dynamically change the up & down b/w some thresholds & potentially speed-up training.

other hyperparameters you can tune

1) No. of layers (deep learning networks)

2) Batch size (How many eg. of data your model sees at once.)

3) No. of trees (Decision tree algos)

4) No. of iterations (How many times model goes through the data - depends on algorithm. (can use early-stopping))

→ Stochastic gradient descent

SGD

SGD

Step 4

Analysis / Evaluation

to estimate how well a model will generalize to out-of-sample data.

a) Evaluation metrics

i) Classification →

ROC & AUC curves

Accuracy

Precision

Recall

F1 score

+ confusion matrix.

ii) Regression

MSE (mean squared error)

MAE (mean absolute error)

R² (r-squared)

Task-based metric (create one based on your specific problem).

i) Classification accuracy

Percentage of correct predictions. ($\frac{Y_{\text{pred}}}{Y_{\text{true}}}$)

Null accuracy

Accuracy that could be achieved by always predicting the most frequent class.

→ It's a good way to know the min. we

should achieve with our model.

→ It gives a baseline to compare our

model's accuracy.

→ Classification accuracy is easy, but it doesn't tell you the underlying distribution or response class.

→ And it doesn't tell you what type of error your classifier makes.

a) confusion matrix

FP - False Positive

TN - True Negative

Predicted

		0	1
		TN	FP
Actual 0	TP	FN	
Actual 1	FP	TP	

They actually don't have Corona. But we predicted they have.

They actually have, but we say they don't.

TP% which are actually correct

TN% which are actually wrong

FP% they are actually wrong, but predicted as true.

FN% they are actually true, but predicted as False

Classification

accuracy

$$\rightarrow \frac{TP + TN}{TP + TN + FP + FN}$$

How often is the classifier correct?

talks about correct.

misclassification rate

$$\rightarrow \frac{FP + FN}{TP + TN + FP + FN}$$

Recall / sensitivity

[True positive rate]

→ How "sensitive" is the classifier to detecting positive instances?

→ when the actual positive is tested, how often is it correct?

Predicted Pos.

TP

$$\text{Recall} = \frac{TP}{TP + FN}$$

→ It depends on FN,

which were actually true, but predicted false.

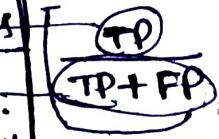
out of total positive ones, how many positive I am able to **Recall**

Factual Pos - TP + FN

Precision & Recall linked with TP

Precision

Really Pos
Total Pos:
actually



$$\frac{TP}{TP+FP} = \frac{\text{Really Pos}}{\text{Predicted Pos}}$$

→ when a positive value is predicted,
How often is the prediction correct?

→ How "precise" is the classifier when
Predicting positive individual.

out of
Predicted Pos,
How many
really Pos.

$$\text{Precision} = \frac{TP}{TP+FP}$$

→ It depends on FP,
which were actually
negative, but
predicted true.

Q) which metrics to focus on?

Eg ①

I = COVID 19 (+)

0 = - (Healthy)

		Pred	
		TN	FP
Act	0	FN	TP
	1	TP	FP

→ Healthy predicted as sick.

→ Sick predicted as healthy.

→ Now analyze the cost.

→ we know a sick predicted as healthy
could spread a lot more easily.

So,

$\text{cost of FN} > \text{cost of FP}$

Hence, It depends on FN, we go for "Recall".

Eg ②

I = Spam

0 = Not Spam

		Pred	
		TN	FP
Act	0	FN	TP
	1	TP	FP

→ Not spam predicted as spam.

→ Spam got into mail box.

→ we know if a not spam (imp. mail) goes to spam,
that could cause lot of damage.

So,

$\text{cost of FP} > \text{cost of FN}$

Hence, It depends on FP, we go for "Precision".

(29)

Sensitivity P

$$\frac{TN}{TN + FP}$$

False positive rate

$$\frac{FP}{TN + FP}$$

FB-Score

when you want to consider both PP & PN like we need both Sensitivity (recall) & precision,
for such cases we choose F-Beta.

General Formula

$$F_B = \frac{(1+\beta^2) \cdot \text{Precision} \times \text{recall}}{(\beta^2 \cdot \text{Precision}) + \text{recall}}$$

$$F_B = \frac{(1+\beta^2) \times TP}{(1+\beta^2) \times TP + \beta^2 \cdot FN + FP}$$

In beta, there is possibility for 3 ranges.

- 1) when FP & FN are equally important, like you can't compromise even a bit both are crucial, then we set beta=1 which is F1-score.

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

- a) when both are worse, but FP is drastic than FN. then we go for beta = 0.5 (usually 0-1)
- b) when both are worse, but FN is drastic than FP. then we go for beta = 2 (usually 1-10).
-

Note

we can also adjust classification threshold to modify performance & classified.

Roc curve

can see how specificity & sensitivity are affected by various thresholds, without actually changing the threshold.

AUC curve or (Area under curve)

If you randomly choose one + one, AUC represents the likelihood that your classifier will assign a higher predicted probability to the positive observation.

Analysis of model

1) Feature Importance

- Which feature contributed most to model
- Should some be removed?

2) Training / Inference time

- How long does the model take to train?
• Is this feasible?

- How long does inference take?
• Is it suitable for production?

3) using what-if tool

- what-if exchanged going to data?

• how does this effect the outcome?

4) Bias / variance trade-off

→ High bias results in underfitting &
a lack of generalization to new sample

→ High variance results in overfitting due
to the model finding patterns in the
data which are actually random

Step 5

Deploying a model

Put the model into production & see how it goes. Evaluation metrics in notebook are great but until it's production, you won't know how it performs for real.

→ Step 6

Retrain model

→ See how the model performs after serving based on various evaluation metrics & revisit the above steps as reqd.

→ You will find your model's prediction starts to lag or drift, as in when data source change or upgrade. This is when you will want to retrain it.

Stop / NO END

DO IT

So, basically the flowchart is -

