

# Pros and cons of ML algorithms

1

## SVM

Pros: ① Performs well in higher dimensions

② Best algorithm when classes are separable

③ Outliers have less impact

④ Suited for extreme case binary classification

Cons: ① Slow: for larger dataset, it requires a large amount of time to process

② Poor performance with overlapped classes.

③ Selecting appropriate hyperparameters is important.

④ Selecting the appropriate kernel function can be tricky.

## Naive Bayes

Pros: ① Real time prediction.

② Scalable with large datasets

③ Insensitive to irrelevant features

④ Multi class prediction.

⑤ Good performance

Cons: ① Independence of features doesn't hold.

② Bad estimator: Probability outputs from predict-proba are not to be taken too seriously.

③ Training data should represent population well.

## Logistic Regression

- Pros :
- ① Simple to implement
  - ② Effective
  - ③ Feature scaling not needed
  - ④ Tuning of hyperparameters not needed.

- Cons :
- ① Poor performance on non-linear data (e.g., image data)
  - ② Poor performance with irrelevant and high correlated features
  - ③ Not very powerful algorithm.

## Random Forest

- Pros :
- ① Reduced risk of overfitting
  - ② Good performance on imbalanced datasets.
  - ③ Can handle huge amounts of data
  - ④ Good handling of missing data
  - ⑤ Little impact of outliers.
  - ⑥ Useful to extract feature importance.

- Cons :
- ① Time consuming process
  - ② Requires more resources
  - ③ Appears as a black box

(3)

## Decision Trees

- Pros:
- ① Normalization or scaling of data not needed.
  - ② Handles missing values.
  - ③ Easy to explain
  - ④ Easy visualization
  - ⑤ Automatic feature selection.

- Cons:
- ① Prone to overfitting
  - ② Sensitive to changes in data.
  - ③ Higher time required to train.

## NN

- Pros:
- ① Simple to understand and implement
  - ② No assumption about data (e.g. in case of linear regression we assume dependent variable and independent variables are linearly related, in Naive Bayes we assume features are independent of each other etc.)
  - ③ Constantly evolving model.
  - ④ Multi-class problems can be solved
  - ⑤ One hyper-parameter to tune.
  - ⑥ No training involved ("lazy"). New training examples can be added easily.

(4)

- Cons :
- ① Slow for large datasets
  - ② Doesn't work very well on datasets with large number of features
  - ③ Scaling of data is required
  - ④ Doesn't work well on imbalanced data
  - ⑤ Sensitive to outliers
  - ⑥ Cannot deal well with missing values.
- XGBoost
- Expensive & slow:  $O(md)$ ,  $m = \# \text{ examples}$ ,  $d = \# \text{ dimension}$
  - To determine the nearest neighbor of a new point  $x$ , must compute the distance to all  $m$  training examples. Runtime performance is slow, but can be improved.
    - Pre-sort training examples into fast data structures
    - Compute only an approximate distance
    - Remove redundant data (condensing).

- Pros :
- ① Less feature engineering required (no need for scaling, normalizing data, can also handle missing values well)
  - ② Feature importance can be found
  - ③ Fast to interpret
  - ④ Outliers have minimal impact
  - ⑤ Handles large datasets.
  - ⑥ Good execution speed
  - ⑦ Good model performance
  - ⑧ Less prone to overfitting.

Cons :

- ① Difficult interpretation, visualization tough
- ② Harder to tune as there are many hyperparameters.

## K-Means

- Pros:
- ① Relatively simple to implement
  - ② Scales to large data sets
  - ③ Guarantees convergence
  - ④ Can warm start the position of centroids
  - ⑤ Easily adapts to new examples
  - ⑥ Generalizes to clusters of different shapes and sizes, such as elliptical clusters

- Cons:
- ① Choosing k manually
  - ② Being dependent on initial values
  - ③ Clustering data of varying sizes and density.
  - ④ Clustering outliers.
  - ⑤ Scaling with number of dimensions.