

## Creating a ML Program

① Load the data

② Pre-processing tasks

ML algo (sklearn) {

- ① Your data must be complete
- ② Your data must be strictly numeric
- ③ Your data must be in the form of numpy array.

① Statistical summary. (describe)

② Check for missing data (info)

③ Handle missing data (Imputation) →

- mean
- median
- mode
- default

④ Handling categorical data → Dummy variables

⑤ Handling ordinal data →

- Encode the data and convert the ordinal data to discrete numeric.
- \* Dummy variables (domain & statistics)

③ Separate your data as features and label. Ensure the data is in the form of numpy array.

④ ⇒ Handle Numerical data. (Feature Scaling)

⇒ This step is applicable for features.

This step ensure your feature columns follow a common statistical scale. (sklearn)

① Rescale your feature manually by defining a range (0 to 1) (-1 to 1) (MinMaxScaler)

⇒ ② Standardize your data (mean = 0, std = 1)  
(StandardScaler)

- ③ Normalize your data. (manual method.... write code using numpy & scipy)
- ① L1 normalization
  - ② L2 normalization
- euclidean ← L2 normalization
- manhattan ← L1 normalization
- ④ Feature Transformation
- ① Log Transformation
  - ② Trigonometric Transformation

The above step ensures your scale is solved, however your magnitude is unaffected.

Rule:- for ML, if your algorithm uses distance formula

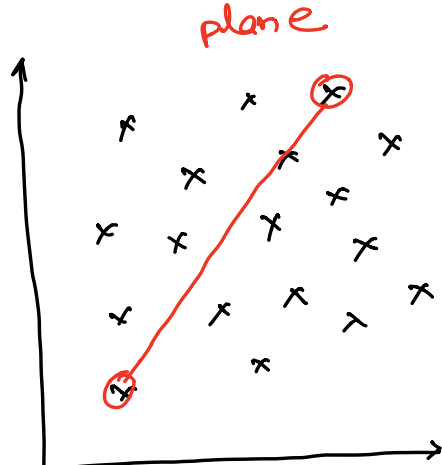
- len ← ① Euclidean \*
- msb ← ② Manhattan
- ③ Minkowski

- k-means
- k-NN
- \* svm (optional)
- [Gradient Descent]

for DL, feature scaling is mandatory for ANN.



① manhattan dist.



② Euclidean dist

magnitudes remain same

100	$\Rightarrow$	1
200		2
300		3
400		4

scaling is reduced

$$37 \Rightarrow 3.7 \times 10^1 \Rightarrow \underbrace{3.7e^{+1}}_{\text{exponential}}$$

$3.700 e^{+1}$