# Data mining : Supervised Machine learning

machine learning

Supervised | unsupervised

**Supervised**
requires data
Sceintist to train
the algorithms
with Label input
with desired outputs


More accurate

**unsupervised**
do not required
any superviser
not Label or any
thing It will identif
the hidden patterns
it wesfful for Anglysi
g data.


Less accurate.
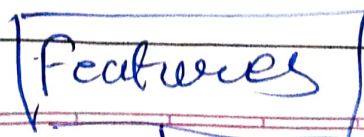
## workflow form Data to decision.

first we collect data may be it
will stored in sql table or No
sql

```
+-------------+
|    data     |
+-------------+
      |
      v
```

after importing data we are gonna find
insight of the data means like pattern
, what kind of data , what should we
predict the (y variable)

```
   ~~feature~~
+-------------+
|  insights   |
+-------------+
      |
      v
```

Next we are going to find some features
(what the important features) features
selection and extraction. we can
called as feature engeneering or
we can use various dimension re-
duction technique like PCA, TSNE etc.

```
+-------------+
|  Features   |
+-------------+
      |
      v
```

by domain knowledge we can know which are important feature and which are not -

```
┌──────────────┐
│   Domain     │
│  Knowledge   │
└──────────────┘
        |
        └──────→
```

Next we should select the right model (keras sklearn Pytorch) if it is supervised learning we sho uld go for sklearn and if it is Neural Network we should go for sklearn and Pytorch.
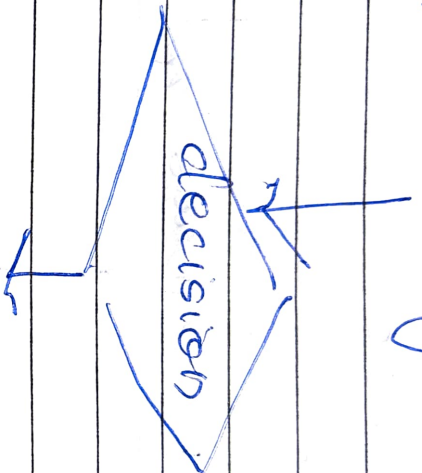
```
┌──────────┐
│  model   │
└──────────┘
      |
      └──────→
```

After that we should get the predic tion out of it. (by domain expert we can know if prediction is proper or not.

```
┌──────────────┐
│  Prediction  │
└──────────────┘
```

then take decision if prediction are
not proper do it again or not


decision

take feedback from domain expert

feedback

again to data

and these process may be
finely because delay prediction are
giving accurate for now or is not

Data Nuances

* Feature noise
  - value may not be accurate
  - it may be Extreme value, Random
    Anomaly

# how to remove Feature noise
=> by Binning.

A. but what is binning

binning is used to deal with Noisy data
we should smooth the data.

there are three technique of binning
1) equal Partitioned bin
2) Bin mean
  - we will before equal partioned bin
    to take out mean and we insert
    mean in Bin for example
    Bin₁ - 9,9,9
    Bin₂ - 22,22,22
    Bin₃ - 29,29,29.

3) Bin Boundaries.

we will take min max values and
thees (Boundaries values) and thos
who are in middles we will change
values thoes who are closest to th
boundaries. for example ( we will refe
to equal partitioned Bin).

Bin1 : 4, 4, 15
Bin2 : 21, 21, 24
Bin3 : 25, 25, 34

* Missing features
=> missing values

* Non - Normal feature Distribution

=> Highly skewed or kurtosis.

* Heterogeeous features.
=> ranges, scales, distribution
example. Age, temperature, Blood press

\* Multimodality

For to solve this we will use standardization.

\* multi modality issue features

Mix of Numeric symbolic series test.

log transformation

⇒ it reduces or remove the skewness of our original data there are two more method two remove skewness is

① log square transformation

② box cox transformation