feature | vector space = feature matrix = x
1. continuous numeric | label = ground truth = y
2. categorical feature

**Distance base** (NN-centroid, KNN, kD tree)
fast, least memory, scalability | sensitive to outlier

**Preprocess** การเก็บ Data ไว้ใน | NN-centroid (mean vector) | KNN | KD tree

① null/nan
1.1 ลบ row
1.2 แทนค่า mean ← continuous
   แทนค่า mode ← categorical
1.3 close fit ← nearest val

② outlier ข้อมูลผิดปกติ
2.1 plot & erase (visually)
2.2 ข้อมูลมากกว่า SD ออกนอกเส้น mean (statistically)

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

③ Data smoothing
3.1 smoothing
sort → split → replace (bin)

3.2 Feature scaling

Decimal scaling $v' = \frac{v}{10^j}$ ; $v' = a + \frac{v}{c_i}$ ; $j = $ min val

minmax normalization
$$v' = \frac{v - min_{old}}{max_{old} - min_{old}}(max_{new} - min_{new}) + min_{new}$$

z-score $v' = \frac{v - mean_{old}}{SD_{old}}$ (mean vector)

④ Feature selection
4.1 Feature selection ลบ column ที่ไม่ดี
4.2 feature aggregation รวม/แปลง ได้ mean vector

$$corr() = r = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum(x_i - \bar{x})^2 \sum(y_i - \bar{y})^2}} ← continuous$$

NN-centroid (mean vector)
1.1 หา mean แต่ละ label ของ feature
2. ทำนายกับกลุ่มของ label ด้วย

$$\sqrt{\sum(avg \; col_i - \bar{x}_i)^2} \; \leftarrow code \; ner$$
$$\sqrt{(x_{bi} - x_i)^2 + (y_{bi} - y_i)^2}$$

1. หา avg col;- sample col;²
2. เรียงลำดับ ได้
3. ระยะใกล้ = predict
hyperparameter: จำนวน $k \rightarrow$

**KNN**
1. หา mean แต่ละ label ของ feature
2. ระยะทางทุกจุด ได้
3. ถ้า len >= leaf size: loop
1-NN: มีโอกาส outlier
KNN: หลายจุด หาระยะทางที่ใกล้
เลือก 1 ได้เลย

hyper: k
leaf size

**KD tree**
1. sort data
2. หา median แบ่ง x,y,z data
2.1 ถ้า median หา if else ของ data
3. ถ้า len >= leaf size: loop
  ↙ else
4. KNN √ ได้


x = 6
y = 4 → (2,3)(4,1) | (1,1)(3,7)(5,4)
y = 8 → (6,8)(7,7)(8,8)

**D tree** (hierarchical structure)
explainable, ez implement | brute force ในการ decision
branch


Outlook
  sunny → humid → high → no / normal → yes
  rain → windy → true / false → yes / no
internal node
leaf node

root = ค่า เกิดมากสุด
entropy มากสุด

**Entropy** = วัดความไม่บริสุทธิ์ $0 \leq$ ดีสุด $1 = $ แย่สุด

$$H(S_1) = -\left[\frac{2}{5}\log_2(2/5) + 3/5 \log_2(3/5)\right]$$

**Gini Index** = ความเป็นไปได้ที่ทายถูก/ผิด
$$Gini(S_1) = 1 - \left[\frac{2}{5}^2 + (3/5)^2\right] = 0.36 \; ต่ำสุดดี$$

entropy สูง
entropy ต่ำลง ด้วย

**information gain**
$$IG(-,y) = 0.953 - \left[\frac{4}{8}H(A) + \frac{1}{8}H(B) + \frac{3}{8}H(C)\right]$$

K = ?
① entropy ถ้าน้อย = 0.
② information ข้อมูล

**Boosting (XGBoost)**
สร้าง weak (depth น้อยๆ/เดาผิดเดาถูก)

point | $y_i$ | $\hat{y}_i$

$z\sum(\hat{y}_i - y_i)\cdot x_i$
$\hat{y}_i - y_i$

epoch รอบการ train
ทำให้ผิดน้อยลง ปรับพารามิเตอร์

**Bagging (Random forest)**
+ได้แม่นยำ + explainable
→ data แตกต่าง ได้ดี

$\hat{y}_i = y_{0i}$ | $\hat{y}_1 = y_{0i} + hv$
lastest decision boundary
$\alpha = $ avg err

$L \; sum(err) / num \; err$

$x_{i+1} = w_0$
$z\sum(\hat{y}_i - y_i)\cdot x_i$

| LR | x | y | w1 | w0 | G(w1) | G(w0) | err | ŷ | z∑(ŷᵢ-yᵢ)·xᵢ | ε(err) | Loss | Epoch |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0.2 | 0.1 | 0.5 | 0.5 | 0.5 | | 0.700 | 0.05 | 0.55 | 0.175 | | | |
| | 0.5 | 0.7 | | | | 1.00 | 0.05 | 0.09 | | | 0.165 | 1 |
| | 0.8 | 0.5 | | | | | 0.40 | | | | | |
| | 0.1 | 0.5 | 0.36 | | -0.112 | -0.164 | 0.24 | | 0.382 | | | |
| | 0.5 | 0.7 | 0.3 | | -0.22 | -0.592 | 0.48 | | 0.418 | | 0.083 | 2 |
| | 0.8 | 0.5 | | | 0.088 | | 0.59 | | | | | |
| | 0.1 | 0.5 | | | | | 0.85 | | | | | |
| | 0.5 | 0.7 | | | -0.090 | -0.093 | 0.61 | | 0.26 | | 0.181 | 3 |
| | 0.8 | 0.5 | | | 0.224 | | 0.77 | | | | 0.060 | |

$$x_{i+1} = x_i - \alpha \cdot \nabla f(x_i) \quad G(w_i)$$

**Mutual information** ← เกือบ = entropy

| | Outlook | | temp |
| | sunny | overcast | rain |
|---|---|---|---|
| temp hot | 2 | 2 | 0 | 4 |
| cold | 1 | 1 | 2 | 4 |
| mild | 2 | 1 | 3 | 6 |
| | 5 | 4 | 5 | 14 |

$$P(outlook=sunny, temp=hot) = \frac{2}{14}\cdot\log_2\left(\frac{2/14}{\frac{5}{14}\cdot\frac{4}{14}}\right)$$

⑤ class imbalance

import pandas as pd

df = pd.read_csv('file name')
df = df.loc[(df['a']>0) & (df['b']<=10)]
df = df.loc[[df.col 1 ? col 2 ?]]
df = df.drop(columns=['col 1 ? 'col 2 '])
df.isna().sum()     df = df.dropna()

df = df.fillna(co)/df.fillna(df.mean()/df.fillna(df.median())
df = df.fillna(df.select_dtype(include='object').mode().iloc[0]) ← categorical | continuous
df['col...'] = df['col...'].replace({'ค่าเก่า':ค่าใหม่, 'ค่าเก่า2':ค่าใหม่2})
df = pd.DataFrame(att, columns=['col 1 ? col 2 ']) ← att to df
corr = df.corr()/df.corr()['col ~'] | df = df['cor.index[:5]']
df = df.sort_values(by='col ~', ascending=False) ← sort new
  ↘ น้อยไปมาก

# Logistic model

$$w_{t+1} = w_t - \alpha G(w_t) \cdot \frac{1}{1+e^z}$$

$\varepsilon(o_i - y_i)A$    $\varepsilon(o_i - y_i)B$    $\varepsilon(o_i - y_i)$    $\varepsilon(err)$

err: $-[Y\log_2(\hat{Y}) + (1-Y)\log_2(1-\hat{Y})]$

| A | B | Y | W1 | W2 | WO | Z | sigmoid | err | G(W1) | G(W2) | G(WO) | BCE loss | epoch |
|---|---|---|----|----|----|----|---------|-----|-------|-------|-------|----------|-------|
| 0 | 0 | 0 |    |    |    | 0.3  | 0.57 | 1.23 |      |      |      |      |   |
| 0 | 1 | 0 | 0.3 | 0.3 | 0.3 | 0.6 | 0.65 | 1.50 | 0.36 | 0.36 | 1.58 | 4.72 | 1 |
| 1 | 0 | 0 |    |    |    | 0.6  | 0.65 | 1.50 |      |      |      |      |   |
| 1 | 1 | 1 |    |    |    | 0.9  | 0.71 | 0.49 |      |      |      |      |   |
| 0 | 0 | 0 |    |    |    | -0.02 | 0.5 | 0.99 |      |      |      |      |   |
| 0 | 1 | 0 | 0.23 | 0.23 | -0.02 | 0.21 | 0.55 | 1.16 | 0.16 | 0.16 | 1.21 | 4.03 | 2 |
| 1 | 0 | 0 |    |    |    | 0.21 | 0.55 | 1.16 |      |      |      |      |   |
| 1 | 1 | 1 |    |    |    | 0.44 | 0.61 | 0.72 |      |      |      |      |   |

---

```python
from sklearn.preprocessing import MinMaxScaler, StandardScaler
scaler = MinMaxScaler(feature_range=(-1,1))   # ← ถ้า ...   LabelEncoder
scaler = StandardScaler()   # ← ถ้า ... math | df = pd.DataFrame(scaler.fit_transform(df), columns=df.columns)
labEn = LabelEncoder() | df['col'] = labEn.fit_transform(df['col'])
```

```python
import matplotlib.pyplot as plt        # import pickle
import seaborn as sns
plt.figure(figsize=(8,6))
sns.heatmap(df.corr(), annot=True)
plt.title('Title name')
plt.show()

pickle.dump(model, open('~.sav/pkl,wb'))   # ← save model
top.to_csv('~.csv', index=False)   # ← save.csv
```

```python
from sklearn.metrics import accuracy_score
acc = accuracy_score(y_test, model.predict(x_test))   # ← categorical
```

```python
Y-pred = model.predict(x_test)

from sklearn.metrics import r2_score
r2 = r2_score(y_test, y-pred)
print(r2)   # ← continuous
```

```python
from sklearn.model_selection import train_test_split
xTrain, xTest, yTrain, yTest = train_test_split(x, y, test_size=0.2, random_state=42)
```

```python
from sklearn.model_selection import GridSearchCV
parameters = {'n_neighbor': [1,20]}   # ← range
model = GridSearchCV(model, parameters)
model.best_params_    # ← วิธี fit
model.best_score_
```

```python
df['Date'] = pd.to_datetime(df['Date'])
df.set_index('Date', inplace=True)   # ← ดึง Date มาทำ index
df['Day'] = df['Date'].dt.day
df['Week'] = df['Date'].dt.isocalendar().week
df['Month'] = df['Date'].dt.month
df['Year'] = df['Date'].dt.year
# str to datetime
```

```python
df['col'].shift(1)
df['col'].shift(-1)
```

```python
from sklearn.neighbors import KNeighborsClassifier
model = KNeighborsClassifier(n_neighbors=5)   # ← dis, ind = knn, KNeighbors([xtest], n_neighbors=5) | print(dis, ind)
model = NearestCentroid()                      # NearestCentroid
```

```python
import xgboost as xgb
model = xgb.XGBClassifier(objective='multi:softmax'   # ← multi:softmax
n_estimators=10, max_depth=5, num_class=2)  ?
```

```python
from sklearn.linear_model import LinearRegression
model = LinearRegression()   | print('w1:', model.coef_)
model = fit(xtrain, ytrain)  | print('wo:', model.intercept_)
```

```python
from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier(max_depth=2, criterion='entropy')  ?
a = clf.best_estimator_       class_name = ['<=sok?','>sok?']
for i in range(6):
    plt.figure(figsize=(20,10))
    plot_tree(a.estimators_[i], feature_names=xTrain.columns)
    plt.title(f'tree{i+1}')
    plt.show()
```

```python
from sklearn.tree import DecisionTreeClassifier, plot_tree
model = DecisionTreeClassifier(max_depth=2, criterion='entropy')  ?
plot_tree(model, feature_names=xtrain.columns, class_names=
import matplotlib.pyplot as plt
plt.figure(figsize=(20,10))
plot_tree(model)
plt.show()
```

```python
from sklearn.metrics import mean_absolute_error, mean_squared_error
print('MAE:', mean_absolute_error(ytest, y-pred))   # ← 67...
print('MSE:', mean_squared_error(ytest, y-pred))    # ← outliner
```