# Cross - Industry

- Crisp - DM
- { regression
   { classification

## 1. Crisp - DM

cross- industry standard process
for data mining

Step I.  understand purpose of the data
        mining study

Step 2.  understand data.
  {1. perform statistical analysis
  {2. perform various types of visualizations

Step 3.  data consolidation { Collect
                              Select
                              Integrate
            ⇓
       data cleaning { Impute missing values
                      Reduce noise in data
                      Eliminate inconsistencies
            ⇓
     data Transformation { Normalize data
                           Discretize/aggregate data
                           Construct new attributes
            ⇓
       data reduction { reduce number of variables
                        reduce number of cases
                        balance ..

Step 4 — model building

Step 5 — testing & evaluation
        { regression → 偏差值
        { classification → 离群率
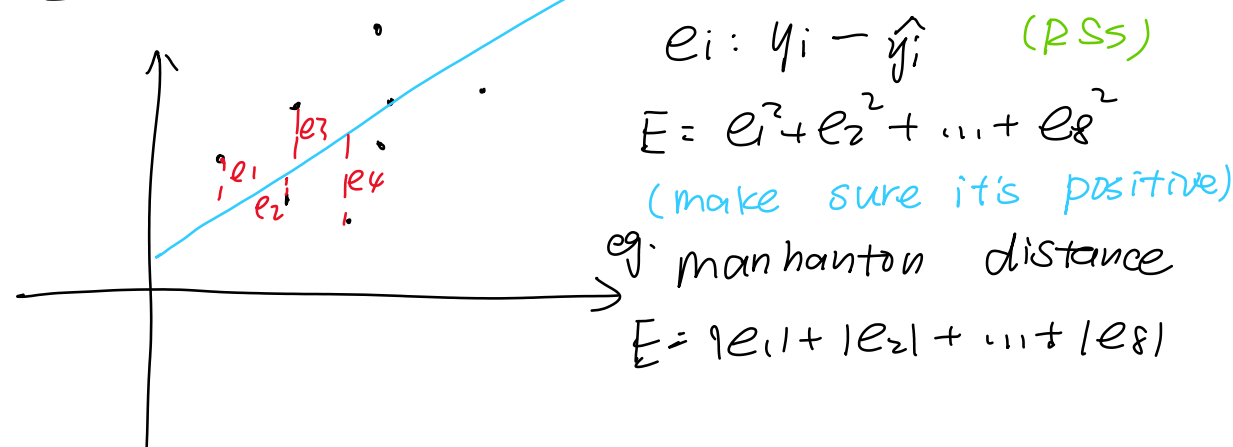        { other evaluation

Step 6 — deployment

# Simple Linear Regression

$Y$ 与 $\beta_0 + \beta_1 x$  ⇒  Sales 与 $\beta_0 + \beta_1 TV$
         └─ gradient

Training Data

$E = e_1^2 + e_2^2 + \cdots + e_8^2$  (square errors)



$e_i : y_i - \hat{y_i}$  (RSS)

$E = e_1^2 + e_2^2 + \cdots + e_8^2$

(make sure it's positive)

eg. manhanton distance

$E = |e_1| + |e_2| + \cdots + |e_8|$

Residual Sum of Square

useful predictors

$R^2 : 0 \sim 1$  好

| X1 | X2 | X3 | X4... | Y |
|----|----|----|----|----|
|    |    |    |    |    |

# F-STATISTICS

越接近1, 越无关, 越大引, 越相关

Prob: (F-statistic): 1.47e-17

MSE: Mean Square

$= \dfrac{1}{degrees-of-freedom} \displaystyle\sum_{i=1}^{n} (y_i - \hat{f}(x_i))^2$

Coding SCHEME (gender..)

gender  numeric { 0
                { 1

Ethnicity  numeric → { Caucasian 0
                      { Asian        1    ✗
                      { African American 2

→ { 1  if in Caucasian
   { 0  if not Caucasian