# K-Means (K clusters, unlabeled input)

Algorithm: Randomly choose K centroids
    Repeat:
        for each data:
            choose closest centroid
        for each cluster:
            update centroid with mean
    Until converage.

# Spectral (谱残) Clustering

Goal: Divide into 2 disjoint groups
* Good Cluster? Cut less edge & sides balanced
Minimize: $\phi(A) = \dfrac{|E|}{\min(vol(A), vol(B))}$
Notice: $vol(A) + vol(B) = 2m$, $vol$: degrees remained

Simplification
A: adjacency matrix
D: degree matrix
L: Laplacian matrix
* $L = D - A$
* $n \times n$ symmetric matrix
* $L \cdot \begin{pmatrix}1\\ \vdots\\ 1\end{pmatrix} = 0 = 0 \cdot \begin{pmatrix}1\\ \vdots\\ 1\end{pmatrix} \Rightarrow \lambda_1 = 0$
* Eigenvalues: $\lambda \geq 0$, $\lambda \in \mathbb{R}$
  Eigenvectors: $\vec{x} \in \mathbb{R}^n$, $(\vec{x_i}, \vec{x_j}) = 0$

| A | 1 | 2 | 3 |
|---|---|---|---|
| 1 | 0 | 1 | 1 |
| 2 | 1 | 0 | 0 |
| 3 | 1 | 0 | 0 |

| D | 1 | 2 | 3 |
|---|---|---|---|
| 1 | 2 | 0 | 0 |
| 2 | 0 | 1 | 0 |
| 3 | 0 | 0 | 1 |

| L | 1 | 2 | 3 |
|---|---|---|---|
| 1 | 2 | -1 | -1 |
| 2 | -1 | 1 | 0 |
| 3 | -1 | 0 | 1 |

Target: $\vec{x} = \underset{x_i \in \mathbb{R}}{\arg\min} \sum_{(i,j) \in E}(x_i - x_j)^2$ (for minimization)
Limit: $\sum x_i^2 = 0$, $\sum x_i = 0$ (for balance)
Finally: $\vec{x} = \underset{x_i \in \mathbb{R}}{\arg\min} \sum_{i=1}^{n}\sum_{j=1}^{n} L_{ij} x_i x_j$
Minimize: $x^T L x$ ($x^T x = 0$)
$\underset{x}{\min}\dfrac{x^T L x}{x^T x} = \underset{y}{\min}\dfrac{y^T \Lambda y}{y^T y}$, $\Lambda = P^T L P$ (正交变换)
$\quad = \underset{y}{\min}\sum_{i=1}^{n}\lambda_i y_i^2$
$\quad = \lambda_{min}$ (最小非零特征值)

比例 $y = \varepsilon_i$, $x = P^{-1}y$ 即为所求

# Support Vector Machine

Line L: $w \cdot x + b = 0$, Point A
$d(A, L) = \dfrac{|wA + b|}{||w||}$
Split: L: $w \cdot x + b = 0$
Perdiction: $\text{sign}(w \cdot x + b) \overset{def}{=\!=} y$
Confidence: $(w \cdot x + b) \cdot y$ ($> 0$)
Margin: $\gamma \overset{def}{=} \dfrac{(w \cdot x + b)}{||w||} \cdot y$
We define $|w \cdot x + b| = 1$ (Reasonable, don't care)
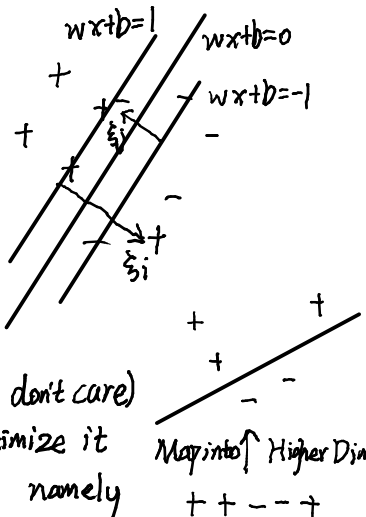Thus $\gamma = \dfrac{1}{||w||}$, we should maximize it
We may as well minimize $||w||$, namely
$\underset{w}{\min}\dfrac{1}{2}||w||^2$ s.t. $\forall i: y_i(w \cdot x_i + b) \geq 1$
Not seperable? Slack variables $\xi_i$
$\underset{w,b,\xi_i > 0}{\min}\dfrac{1}{2}||w||^2 + C\sum_{i=1}^{n}\xi_i$ s.t. $\forall i: y_i(w \cdot x_i + b) \geq 1 - \xi_i$
C is regularization parameter. Or: Map into higher dimension

*(figure: SVM separating hyperplanes $wx+b=1$, $wx+b=0$, $wx+b=-1$ with $\xi_i$ slack, + and − points)*

Map into↑ Higher Dimension
++ −− +

# C4.5 (Classification, Trained)

Which attribute differs best? 「Gain Ratio」determines
For $\{Y_i\}$ with m classes and distribution $\{P_i\}$:
  Entropy (熵) $H(P) = -\sum_{i=1}^{m} p_i \ln p_i = -\sum \dfrac{|Y_i|}{|Y|} \ln \dfrac{|Y_i|}{|Y|}$
  Information $Info(Y) = H(P)$ (越大越混乱)
After classification by attribute A with k classes:
$\quad Info_A(Y) = \sum_{j=1}^{k}\dfrac{|Y_j|}{|Y|} Info(Y_j)$
$\quad GainRatio(A) = \dfrac{Gain(A)}{SplitInfo(A)} = \dfrac{Info(Y) - Info_A(Y)}{Info(A)}$
For example: $Info(out) = -\dfrac{4}{14}\ln\dfrac{4}{14} - \dfrac{10}{14}\ln\dfrac{10}{14} = 0.940$
$\quad Info_{wind}(out) = -\dfrac{8}{14}(\dfrac{2}{8}\ln\dfrac{2}{8} + \dfrac{6}{8}\ln\dfrac{6}{8}) - \dfrac{6}{14}(\dfrac{3}{6}\ln\dfrac{3}{6} + \dfrac{3}{6}\ln\dfrac{3}{6}) = 0.892$
$\quad SplitInfo(wind) = -\dfrac{8}{14}\ln\dfrac{8}{14} - \dfrac{6}{14}\ln\dfrac{6}{14} = 0.985$
$\quad GainRatio(wind) = \dfrac{0.940 - 0.892}{0.985} = 0.049$

Then build the tree: 1. Select maximal GainRatio attribute from pool
    2. Remove it from pool  3. Generate new node
    4. For each partition, if not completed, back to 1

# PageRank

$PR(c) = \dfrac{PR(A)}{L(A)} + \dfrac{PR(B)}{L(B)}$
* Random Surfur: $PR(u) = \alpha\sum_{v \to u}\dfrac{PR(v)}{L(v)} + \dfrac{1-\alpha}{N}$
Let $P_0 = \dfrac{e}{N}$, $P = (PR(1), PR(2), \cdots, PR(m))$, $S$: transition matrix
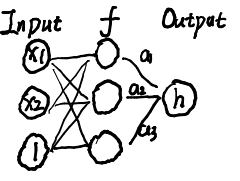We have $P_{i+1} = \alpha S P_i + (1-\alpha)P_0$
Example: $S = \begin{pmatrix}\frac{1}{2} & \frac{1}{2} & 0 & \frac{1}{2}\\ \frac{1}{3} & 0 & 0 & \frac{1}{2}\\ 0 & 0 & 0 & 0\\ \frac{1}{3} & \frac{1}{2} & 1 & 0\end{pmatrix}$, $P_0 = \begin{pmatrix}\frac{1}{4}\\ \frac{1}{4}\\ \frac{1}{4}\\ \frac{1}{4}\end{pmatrix}$

*(figure: directed graph A→B, A→C, B→D, C→D, D→B)*
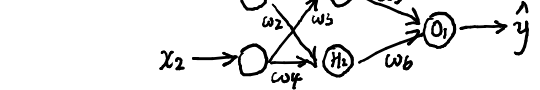
# Deep Learning (Neural Network)

前向传播 $a_1 = f(w_{11}x_1 + w_{12}x_2 + b_1)$
$\quad a_2 = f(w_{21}x_1 + w_{22}x_2 + b_2)$
$\quad h(x) = f(w_{11}a_1 + w_{21}a_2 + b')$
Or we say: $z_i^{(n+1)} = \sum w_{ij}x_j + b_i^{(n)}$
$\quad a_i^{(n)} = f(z_i^{(n)})$

激活函数BP: Sigmoid: $f(z) = \dfrac{1}{1+e^{-z}}$, $f'(z) = f(z)(1-f(z))$
$\quad$ ReLU: $f(x) = \max\{0, x\}$.

反向传播: $x_1 \to$ ... $E = \dfrac{1}{2}(y-\hat{y})^2 = \dfrac{1}{2}(y - O_1)^2$
$O_1 = f(O_1)$, $do_1 = \dfrac{\partial E}{\partial O_1}\cdot\dfrac{\partial O_1}{\partial O_1}$
$\therefore do_1 = -(y-O_1)\cdot f(O_1)(1-f(O_1))$

$O_1 = w_5 H_1 + w_6 H_2 + b$, $dw_5 = \dfrac{\partial E}{\partial O_1}\cdot\dfrac{\partial O_1}{\partial w_5} = do_1 \cdot H_1$
$db = \dfrac{\partial E}{\partial O_1}\cdot\dfrac{\partial O_1}{\partial b} = do_1$, $dw_6 = \dfrac{\partial E}{\partial O_1}\cdot\dfrac{\partial O_1}{\partial w_6} = do_1 \cdot H_2$
$dH_1 = \dfrac{\partial E}{\partial O_1}\cdot\dfrac{\partial O_1}{\partial H_1} = do_1 \cdot w_5$, $dw_1 = \dfrac{\partial E}{\partial H_1}\cdot\dfrac{\partial H_1}{\partial w_1} = dH_1 \cdot x_1$
同理有 $dH_2, dw_2, dw_3, dw_4$
Update: $w_i \mathrel{+}= \eta \cdot dw_i$, $\eta$: Learning Rate

*(figure: Input f Output network; Input–Hidden–Output with $\omega_1 ... \omega_6$, $H_1, H_2, O_1, \hat{y}$)*

# KNN (Supervised, Classification)

Algorithm: Find distance to each point
    Sort all the distanc
    Take closest K points
$\hat{Y}(x) = \dfrac{1}{K}\sum_{x_i \in N_k(x)} y_i$
With K-D Tree optimization: $O(\log N) \sim O(N)$

# Naive Bayes (Data Mining)

$P(B|A) = \dfrac{P(AB)}{P(B)} = \dfrac{P(A|B)P(B)}{P(A)}$
$\hat{v} = \underset{v_j \in V}{\arg\max}\, P(v_j | a_1 a_2 \cdots a_n) = \underset{v_j \in V}{\arg\max}\, P(a_1 a_2 \cdots a_n | v_j)P(v_j)$
$\quad = \underset{v_j \in V}{\arg\max}\, P(v_j)\prod P(a_i | v_i)$ (独立性)