

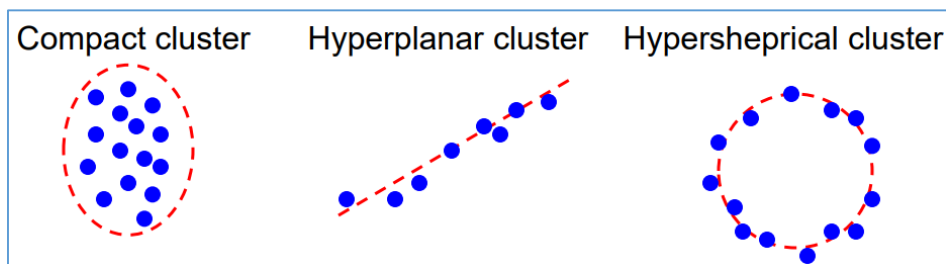
Clustering

- ✓ Classification 與 clustering 差別在於: **classification** 有用到 **class label**, 然 **clustering** 則無。
- ✓ clustering 中的 data 無 label, 因此為 **unsupervised** 的過程。也因為沒有 label, clustering 的結果無一定的好壞標準, 通常會依據問題或需求而定標準。
- ✓ dataset X , clustering 後分成 m 個 clusters。其中 clusters 的集合稱為 clustering。

$$\begin{array}{|l|} \hline C_i \neq \phi, \forall i \\ C_1 \cup C_2 \cup \dots \cup C_m = X \\ \hline \mathcal{R} = \{C_1, C_2, \dots, C_m\} \\ \hline \end{array}$$

- ✓ (cluster representative)為了要降低資料運算量, 通常將一個 cluster 用較少的 **parameters** 表示。常見的兩種表示方式分別為: **point representative (compact cluster)**與 **shell representative (shell cluster)**。Point representative 可以用 cluster 中所有 sample 點的平均(mean point)表示, 也可以用實際在 cluster 中的 sample 點表示(找到一個點使得相異度加總最小或是相似度加總最大)。

$$\begin{array}{l} \text{mean point } m_p: m_p = \frac{1}{n_C} \sum_{y \in C} y \\ \text{mean center } m_c \text{ (one of the samples in the cluster):} \\ m_c = \arg \min_{z \in C} \sum_{y \in C} d(y, z) \text{ or } m_c = \arg \max_{z \in C} \sum_{y \in C} s(y, z) \end{array}$$



- ✓ 在分群的過程中, 會遇到要將相似的 **patterns** 分在同一個 cluster; 將相異的 **patterns** 分在不同 clusters。多相似多相異, 這些皆需要判斷的標準, 我們稱為 **proximity measure**。當使用不同的 measure, 產生出 cluster 的形狀也有所不同。**Data** 使用上述的 **measure** 後, 就會成為 **relational data**。

- ✓ 相異度測量(dissimilarity measure)。須滿足以下兩個條件: 1. 自身相異度最小
2. 兩點之間的相異度不變。常見的方法有: Weighted lp distance、Generalized l₂ norm、Hamming distance、Edit distance。

$$\begin{aligned} d(\mathbf{x}, \mathbf{y}) &\geq d(\mathbf{x}, \mathbf{x}) = d_0, \quad \forall \mathbf{x}, \mathbf{y} \in X \\ d(\mathbf{x}, \mathbf{y}) &= d(\mathbf{y}, \mathbf{x}), \quad \forall \mathbf{x}, \mathbf{y} \in X \end{aligned}$$

It is a metric dissimilarity measure if it also satisfies

$$d(\mathbf{x}, \mathbf{y}) = d_0 \Leftrightarrow \mathbf{x} = \mathbf{y}$$

$$d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z}), \quad \forall \mathbf{x}, \mathbf{y}, \mathbf{z} \in X$$

- Weighted lp distance:

$$d_p(\mathbf{x}, \mathbf{y}) = \left[\sum_{i=1}^l w_i |x_i - y_i|^p \right]^{1/p}$$

l₂ with w_i=1: Euclidean distance; l₁: Manhattan norm。

- Generalized l₂ norm: 其中 B 為 symmetric、positive definite matrix。

$$d(\mathbf{x}, \mathbf{y}) = \left[(\mathbf{x} - \mathbf{y})^T B (\mathbf{x} - \mathbf{y}) \right]^{1/2}$$

- Hamming distance: 常用於離散的點。Hamming distance 就是兩 vectors 相異處的個數。例如: 101 與 111 的 hamming distance 為 1。
 - Edit distance: 用於資料為字串形式。Edit distance 就是可以使得兩字串相同的總改變量。例如: kitten 與 sitting 的 edit distance 為 3。
- ✓ 相似度測量(similarity measure)。須滿足以下兩個條件: 1. 自身相似度最大 2. 兩點之間的相似度不變。常見的方法有: cosine similarity、Tanimoto measure。

$$\begin{aligned} s(\mathbf{x}, \mathbf{y}) &\leq s(\mathbf{x}, \mathbf{x}) = s_0, \quad \forall \mathbf{x}, \mathbf{y} \in X \\ s(\mathbf{x}, \mathbf{y}) &= s(\mathbf{y}, \mathbf{x}), \quad \forall \mathbf{x}, \mathbf{y} \in X \end{aligned}$$

It is a metric similarity measure if it also satisfies

$$s(\mathbf{x}, \mathbf{y}) = s_0 \Leftrightarrow \mathbf{x} = \mathbf{y}$$

$$s(\mathbf{x}, \mathbf{y}) + s(\mathbf{y}, \mathbf{z}) \leq [s(\mathbf{x}, \mathbf{y}) + s(\mathbf{y}, \mathbf{z})]s(\mathbf{x}, \mathbf{z}), \quad \forall \mathbf{x}, \mathbf{y}, \mathbf{z} \in X$$

- Cosine similarity: 所有 vectors 皆已 normalized 為前提。其中 $\mathbf{x}^T \mathbf{y}$ 代表 \mathbf{x} 與 \mathbf{y} 內積的結果。

$$s_{\text{cosine}}(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|}$$

- Tanimoto measure: 所有 vectors 皆已 normalized 為前提。

$$s_T(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 - \mathbf{x}^T \mathbf{y}}$$

- ✓ 對於一個 clustering 演算法，我們需要一個 function 作為要 assign sample \mathbf{x} 給哪個 cluster 使用，也就是找出 point 與 cluster 之間的距離。它們之間的距離可以定義成: sample \mathbf{x} 與 cluster 最遠(或最近)的距離、sample 與 cluster 中所有點距離的平均等等。

$$\begin{aligned}\mathcal{P}_{\max}(\mathbf{x}, C) &= \max_{\mathbf{y} \in C} \mathcal{P}(\mathbf{x}, \mathbf{y}) \\ \mathcal{P}_{\min}(\mathbf{x}, C) &= \min_{\mathbf{y} \in C} \mathcal{P}(\mathbf{x}, \mathbf{y}) \\ \mathcal{P}_{\text{avg}}(\mathbf{x}, C) &= \frac{1}{n_C} \sum_{\mathbf{y} \in C} \mathcal{P}(\mathbf{x}, \mathbf{y})\end{aligned}$$

- ✓ 前面提到: 通常我們會用 cluster representative 來表示一個 cluster。因此，point 與 cluster 之間的距離就可以定義成 point 與 cluster representative 之間的距離。

Point representatives \mathbf{m} : $\mathcal{P}(\mathbf{x}, C) = \mathcal{P}(\mathbf{x}, \mathbf{m})$

Shell representatives (examples):

Hyperplanes: $d(\mathbf{x}, C) = |\mathbf{a}^T \mathbf{x} + a_0| / \|\mathbf{a}\|$

Hyperspheres: $d(\mathbf{x}, C) = |\|\mathbf{x} - \mathbf{c}\| - r|$

- ✓ 兩 clusters 之間的距離，可以定義成如上述提到的 point 與 cluster 距離、或是兩 cluster representatives 之間的距離。

$$\mathcal{P}_{max}(C_1, C_2) = \max_{x \in C_1, y \in C_2} \mathcal{P}(x, y)$$

$$\mathcal{P}_{min}(C_1, C_2) = \min_{x \in C_1, y \in C_2} \mathcal{P}(x, y)$$

$$\mathcal{P}_{avg}(C_1, C_2) = \frac{1}{n_{C_1} n_{C_2}} \sum_{x \in C_1} \sum_{y \in C_2} \mathcal{P}(x, y)$$

$$\mathcal{P}_{mean}(C_1, C_2) = \mathcal{P}(m_{C_1}, m_{C_2})$$

✓ Sequential clustering:

- BSAS: 每加入一個 sample 進 cluster，就會更新 cluster representative。根據設定的 θ 值，判斷要加入哪一個 cluster。

Basic Sequential Algorithmic Scheme (BSAS):

```

m ← 1
Cm ← {x1}
For i = 2 to N
    k ← argmin d(xi, Ck)
    If (d(xi, Ck) > θ) AND (m < q)
        m ← m + 1
        Cm ← {xi}
    Else
        Ck ← Ck ∪ {xi}
        Update cluster representatives if necessary
    End
End

```

- MBSAS: 一開始根據 θ 值建立初始的 clusters。有了這些初始的 clusters，將 sample 加入哪一個 cluster 的方法與 BSAS 完全相同。

Modified BSAS (MBSAS)

```

m ← 1
Cm ← {x1}
For i = 2 to N
    k ← argmin d(xi, Ck)
    If (d(xi, Ck) > θ) AND (m < q)
        m ← m + 1
        Cm ← {xi}
    End
End
For i = 1 to N
    If xi is not in any cluster
        k ← argmin d(xi, Ck)
        Ck ← Ck ∪ {xi}
        Update cluster representative if necessary
    End
End

```

- ✓ 在 BSAS 與 MBSAS 中不同的 θ 值會有不同的 clustering 結果。如何選出一個較好的結果？找到一個 clustering，它對於 θ 值影響最小。例如 $\theta = 1 \sim 5, q=6$ ，對於資料順序為 328716，並且使用 MBSAS， $\theta = 2 \sim 4$ 都是分成兩群，因此分成兩群是相對穩定的結果。
- ✓ (refinement by merging) 兩 clusters 可能非常相近，將它們合併成一個 cluster 會比較好，我們可以定義：相似程度超過門檻值時，就將它們合併。
(refinement by reassignment) 計算 clusters 的 representatives，並根據 representatives 從新 assign 到合適的 cluster 裡。
- ✓ **Hierarchical clustering:** 常用在 clustering relational data，會產生 nest effect (對於 agglomerative clustering 而言，當兩個 samples 被分配到同一個 cluster，則它們最終仍會在同一個 cluster 裡；對於 divisive clustering 而言，當兩個 samples 被分配到不同的 clusters，則它們最終仍會被分隔)。兩種常見的型態為 agglomerative clustering、divisive clustering。Agglomerative clustering 起始為一個 cluster，每次兩兩 clusters 合併直到所有 sample 都合併；division clustering 起始為一個 cluster，其中包含著所有的 samples，每次將 cluster 一分為二。
- ✓ 對於 agglomerative clustering，要如何選擇每次要合併的兩個 clusters？如使用相異度量測(dissimilarity measure)，則選擇的 clusters 為相異度最小的。常見的方法有 single link 與 complete link。Single link 定義的相異程度為兩 clusters 間的最小相異度；complete link 則為兩 clusters 間最大的相異度。當 clusters 被合併成另一個較大的 cluster 後，相異程度的判斷仍不變。以 single link 為例：假設 C_i 與 C_j 合併成 C_q ，則任一個 cluster C_s 與 C_q 的相異程度為 (C_s 與 C_i 最小相異度) 與 (C_s 與 C_j 最小相異度) 的最小值。

$$d_{\min}(C_1, C_2) = \min_{x \in C_1, y \in C_2} d(x, y) \quad (\text{single link})$$

$$d_{\max}(C_1, C_2) = \max_{x \in C_1, y \in C_2} d(x, y) \quad (\text{complete link})$$

$$d_{\min}(C_q, C_s) = \min\{d_{\min}(C_s, C_i), d_{\min}(C_s, C_j)\}$$

$$d_{\max}(C_q, C_s) = \max\{d_{\max}(C_s, C_i), d_{\max}(C_s, C_j)\}$$

- ✓ Single link 容易形成狹長型的 clusters。缺點為 cluster 的頭尾兩點距離遠。Single link (complete link) works better for elongated clusters (compact clusters). 除 single link 與 complete link 以外，也有其他的方法例如：WPGMA、UPGMA。

Weighted pair group method average (WPGMA)

$$d(C_q, C_s) = \frac{1}{2} [d(C_s, C_i) + d(C_s, C_j)]$$

Unweighted pair group method average (UPGMA; average-link)

$$d(C_q, C_s) = \frac{n_i}{n_i + n_j} d(C_s, C_i) + \frac{n_j}{n_i + n_j} d(C_s, C_j)$$

- ✓ 用暴力法(brute-force)產生 agglomerative clustering 需要 $O(N^3)$ ，而 single link 方法等同於建立 MST，因此只需要 $O(N^2)$ ；complete/average link 可以用 heap 建立，需花費 $O(N^2 \log N)$ 。
- ✓ 對於 divisive clustering，要如何選擇 split 成兩個 clusters？我們可以找相異度最高(或相似度最低)的兩個 clusters，然而在實作上此方法非常耗費時間。
- ✓ **C-Means algorithm:** clusters 的數量與 cluster representative (prototype)需給定。常見的方法有: HCM、FCM、PCM。

- **HCM (Hard/crisp C-Means):** 每一個 sample 被分配到唯一一個 cluster。其中 cost function 可以表示為: (d_{ij}^2 為 sample i 與 cluster j 的距離)

- $C = \sum_{i=1}^N \sum_{j=1}^C d_{ij}^2 \delta_{ij}$, subject to $\forall i, \sum_{j=1}^C \delta_{ij} = 1, \delta_{ij} = 0 \text{ or } 1$

它是一個 alternating optimization 過程。每次都會將 sample x 分配到最接近最合適的 cluster (固定 d_{ij}^2 降低 C)，接著再更新 prototype (固定 δ_{ij} 降低 C)，重複上述兩步驟直到 stopping criteria (例如 prototype 於兩次 iteration 間不再改變)。

- **FCM (Fuzzy C-Means):** 每一個 sample 可以被分配到多個 clusters，其隸屬於各群的程度為[0,1]之間，且所有程度加總為 1。它的 cost function 可以表示為:

- $J = \sum_{i=1}^N \sum_{j=1}^C d_{ij}^2 u_{ij}^q$ where $q > 1$, subject to $\forall i, \sum_{j=1}^C u_{ij} = 1$

越大的 q 值，能接受 fuzzy 的程度越高。

- Recalculate the memberships (u_{ij})

$$u_{ij} = \frac{1}{\sum_{k=1}^C \left(\frac{d_{ij}}{d_{ik}} \right)^{\frac{2}{q-1}}}$$

- Recalculate the prototypes (v_j)

This is for point prototypes only: $\rightarrow v_j = \frac{\sum_{i=1}^N (u_{ij})^q x_i}{\sum_{i=1}^N (u_{ij})^q}$

- **PCM (Possibilistic C-Means):** 將 FCM 中“隸屬程度加總為 1”的條件移除，也就是 sample 可以隸屬於多個 clusters，且程度加總允許大於 1，另一方面，也可以不隸屬任一個 cluster。對於 dataset 中有 outliers，FCM 會使得 prototype 外移，然而 PCM 可以讓 u_{ij} 變小，使 outliers 對 cluster 影響程度降至最小。它的 cost function 可以表示為：

$$J = \sum_{i=1}^N \sum_{j=1}^c u_{ij}^q d_{ij}^2 + \sum_{j=1}^c \eta_j \sum_{i=1}^N (1 - u_{ij})^q$$

- Recalculate the memberships (u_{ij})
$$u_{ij} = \frac{1}{1 + \left(\frac{d_{ij}^2}{\eta_j} \right)^{\frac{1}{q-1}}}$$
- Recalculate the prototypes (v_j)

Same as FCM \longrightarrow

$$v_j = \frac{\sum_{i=1}^N (u_{ij})^q x_i}{\sum_{i=1}^N (u_{ij})^q}$$

- 總結 HCM、FCM、PCM。
 - HCM 是(收斂)速度最快的。
 - HCM 是 FCM 中的一個特例($q=1$)。
 - HCM 相較於 FCM 更可能收斂到 local optima。
 - PCM 對於 outliers 的影響最小，對於 initialization 最敏感。
- **C-Medoid clustering:** 與 C-Means 不同的是: prototype 為 cluster 中真實的點，也就是找到 cluster 中的一個點與 cluster 中的其他點相異程度加總最小。

$$m_c = \arg \min_{z \in C} \sum_{\substack{y \in C \\ y \neq z}} d(y, z)$$

- Gustafson-Kessel algorithm (FCM 的一個變形): 與 FCM 唯一不同之處為: prototype 以 normal distribution 取代。

$$J = \sum_{i=1}^N \sum_{j=1}^c u_{ij}^q d_{ij}^2$$

$$d_{ij}^2 = d_{GK}^2(\mathbf{x}_i, \boldsymbol{\theta}_j) = \frac{1}{|\Sigma_j|^{1/l}} (\mathbf{x}_i - \mathbf{c}_j)^T \Sigma_j^{-1} (\mathbf{x}_i - \mathbf{c}_j)$$

The update equations:

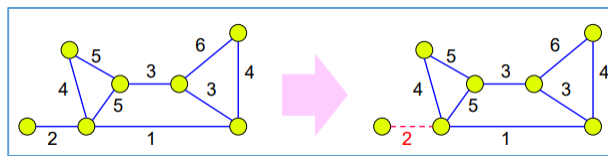
$$u_{ij} = \frac{1}{\sum_{k=1}^c \left(\frac{d_{ij}}{d_{ik}} \right)^{\frac{2}{q-1}}}$$

$$\mathbf{c}_j = \frac{\sum_{i=1}^N (u_{ij})^q \mathbf{x}_i}{\sum_{i=1}^N (u_{ij})^q}$$

$$\Sigma_j = \frac{\sum_{i=1}^N (u_{ij})^q (\mathbf{x}_i - \mathbf{c}_j)(\mathbf{x}_i - \mathbf{c}_j)^T}{\sum_{i=1}^N (u_{ij})^q}$$

- ✓ **Density-based Spatial Clustering of Application with Noise (DBSCAN):** 一種 **nonparametric (no prototype) clustering** 演算法。對於高密度點的會自成一類，低密度點的會被視為 **noise**。類似 **single link agglomerative clustering**，但此方法不會連接低密度點。首先，定義指定半徑 **r** 與密度閾值 **t**。每次拜訪尚未拜訪過的點，以此點為中心向外擴展半徑為 **r** 的圓形區域，如該區域包含到的點(自身點除外)之個數超過密度閾值 **t**，則將它們自成一類，同時包含到的其他點也向外擴展，找到滿足閾值的所有點，並將它們納入同一類。因此，只需要給定 **r** 與 **t** 就能分群。其餘不滿足的點則視為 **noise**。DBSCAN 的困難點在於: **t** 與 **r** 的選定以及找 neighbors 的效率。
- ✓ **Graph-cut based clustering:** 把 samples 描述成 **weighted graph**，其中 edges 上的 **weight** 為兩 samples 相似程度。透過移除 edges，可以將 graph 分成兩個(多個) **subgraph**。edges 上的 **weights** 為相似度，因此移除 edges 的 **wieght** 總和越小越好(如 **wieghts** 代表相異度，則移除 **wieght** 總和越大越好)。但這種方法容易產生不平均的 **subgrpahs** 或是產生 **isolated**

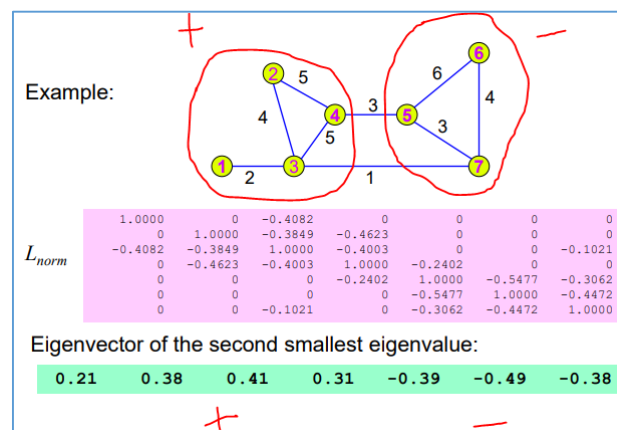
samples。我們可以透過 Normalized cut 找到最合適的 cut，使得 $ncut(A,B)$ 最小。



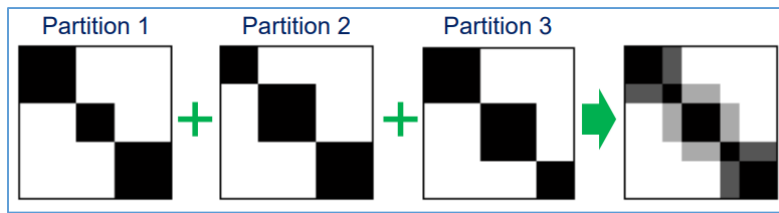
$$\text{Normalized cut: } ncut(A, B) = \frac{w(A, B)}{w(A, V)} + \frac{w(A, B)}{w(B, V)}$$

- V is the original set of vertices, and A and B are the two subsets of vertices.
- $w(A, B)$ is the total weights of edges that connect vertices between A and B .

- ✓ **Spectral clustering:** 因為“最小化 $ncut(A,B)$ ”的問題不容易。我們可以透過尋找 normalized Laplacian matrix 特徵向量的方式找到合適的切割。normalized Laplacian matrix (NLM) $L_{norm} = I - D^{-1/2} A D^{-1/2}$ 中 A 為 weighed 相鄰矩陣(有連接才有值，且對角線的值為 0)， D 為對角矩陣(D_{ii} 為連接到 vertex i 的 weight 總和)。因為 NLM 的最小特徵值為 0，故找第二小的特徵值。相對應的特徵向量，並依照其 component 的正負值來分群。例如第二小的特徵值對應到的特徵向量為(0.2, 0.3, 0.4, -0.3, -0.3, -0.4, -0.3)，前三個點、後四個點各自分在不同群。時間複雜度為 $O(N^3)$ 。



- ✓ 透過簡易的 clustering 演算法可以將 dataset 裡的 samples 分群。Cluster ensemble 就是將分群的結果再做分群(將某些 clusters 合併成更大的 clusters)。此時 similarity measure 就會是兩 samples 在同 cluster 但在不同 partition 的 likelihood。根據相似度量測，我們可以使用前述的 spectral clustering 或是 hierarchical agglomeration 做 cluster ensemble。Co-association (affinity) matrix 是相似度量測結果對於不同的 partition 常見的表示方法。



✓ Clustering 中無使用到 label，因此對於 clustering 的結果沒有直觀的判斷標準(例如正確率)。一些常見的 cluster validity 方法: Dunn Index、DB Index、XB Index。

- **Dunn Index:** 用於 hard (crisp) clustering。Dm 的值越大，代表 clustering 結果越好。其中，分子為兩 clusters 間相異度的最小值 $\min(d(C_i, C_j))$ ；分母為 cluster 中兩 samples 中最遠的距離 $\max(\text{diam}(C_k))$ 。例如:有兩個 clusterings 結果: $\{1,2\}, \{3,5\}, \{7,8\}$ 與 $\{1,2,3\}, \{5,7,8\}$ 。前者 $D_m = (3-2)/(5-3) = 1/2$ ；後者 $D_m = (5-3)/(8-5) = 2/3$ 。因為 $2/3 > 1/2$ ，所以後者 clustering 結果較好。

$$D_m = \min_{\substack{i,j \\ i \neq j}} \left[\frac{d(C_i, C_j)}{\max_k [\text{diam}(C_k)]} \right] \quad \text{The larger, the better.}$$

Dissimilarity between clusters: $d(C_i, C_j) = \min_{x \in C_i, y \in C_j} d(x, y)$

Cluster diameter: $\text{diam}(C_k) = \max_{x, y \in C_k} d(x, y)$

- **Davies-Bouldin (DB) Index:** 其中 S_i, S_j 以及分母的算法，如前述 weighted lp distance 雷同。Vi 為 cluster i 的 cluster representative (prototype)。q 為 fuzzy factor，q=1 為 hard clustering；q>1 為 fuzzy clustering。與 FLD 算 class separability 很像，分母如 within class variance；分子如 between class variance。因此 DBm 值越小，代表 clustering 結果越好。

DB index for hard clusterings:

$$DB_m = \frac{1}{m} \sum_{i=1}^m \max_{j \neq i} \left[\frac{s_i + s_j}{\|v_i - v_j\|_q} \right]$$

The smaller, the better.

where

$$s_j = \left[\frac{1}{n_j} \sum_{x \in C_j} \|x - v_j\|^r \right]^{1/r}$$

DB index for fuzzy clusterings:

Use

$$s_j = \left[\frac{\sum_{i=1}^N u_{ij}^q \|x_i - v_j\|^r}{\sum_{i=1}^N u_{ij}^q} \right]^{1/r}$$

- **Xi-Beni (XB) Index:** 用於 fuzzy clustering。XB 值越小，代表 clustering 結果越好。

$$XB = \frac{\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^c u_{ij}^2 \|x_i - v_j\|^2}{\min_{j \neq i} \|v_i - v_j\|^2}$$

The smaller, the better.

- ✓ 對於 fuzzy partition 的 crisp 程度。可以透過 partition coefficient 或是 partition entropy coefficient 判斷。它們的值越大，代表 partition 越 crisp (越不 fuzzy)。

Partition coefficient:

$$PC = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^m u_{ij}^2$$

The larger, the better.

Partition entropy coefficient:

$$PE = \frac{-1}{N} \sum_{i=1}^N \sum_{j=1}^m (u_{ij} \log u_{ij})$$

The larger, the better.

- ✓ 當我們有正確的 partition 資訊時，就可以對 clustering 結果好壞做評估。常見的方法有: Adjusted RAND Index (ARI)、Normalized Mutual Information (NMI)、clustering accuracy (Hungarian Algorithm)。

- **RAND Index:** 為 ARI 的簡易版。R 值越高，代表 clustering accuracy 越高。R = $\frac{a+b}{\binom{n}{2}}$ 。其中，n 為 samples 的總數、a 為 sample pairs 在同一個 cluster 不同 partition 的總數、b 為 sample pairs 在不同 clusters 不同 partition 的總數。例如:有兩個 clusterings 結果: {1},{2,3}與{1,2},{3}。R=(0+1)/3=1/3;

- **Adjusted RAND Index (ARI):** 為 RI 的調整版。根據下表的公式可以算出 ARI。ARI 值越高，代表 clustering accuracy 越高。其中 nij 代表 sample 在第一個 partition 中被分配到 cluster i 且同時在第二個 partition 中被分配到 cluster j 的個數。

	Y_1	Y_2	\dots	Y_s	$Sums$
X_1	n_{11}	n_{12}	\dots	n_{1s}	a_1
X_2	n_{21}	n_{22}	\dots	n_{2s}	a_2
\vdots	\vdots	\vdots	\ddots	\vdots	\vdots
X_r	n_{r1}	n_{r2}	\dots	n_{rs}	a_r
$Sums$	b_1	b_2	\dots	b_s	

the adjusted index is:

$$ARI = \frac{\sum_{ij} \binom{n_{ij}}{2} - [\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}] / \binom{n}{2}}{\frac{1}{2} [\sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2}] - [\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}] / \binom{n}{2}}$$

- Normalized Mutual Information (NMI): NMI 值越高，代表 clustering 結果越好。

[\[Here is an example\]](#)

$$NMI(Y, C) = \frac{2 \times I(Y; C)}{[H(Y) + H(C)]}$$

- 1) Y = class labels
- 2) C = cluster labels
- 3) H(.) = Entropy
- 4) I(Y;C) = Mutual Information b/w Y and C

- Hungarian Algorithm: 找到最低的分配成本。[slide reference]

Hungarian Method - Example 1

	Job1	Job2	Job3	Job4
Crane1	4	2	5	7
Crane2	8	3	10	8
Crane3	12	5	4	5
Crane4	6	3	7	14

min=2

min=3

min=4

min=3

	Job1	Job2	Job3	Job4
Crane1	2	0	3	5
Crane2	5	0	7	5
Crane3	8	1	0	1
Crane4	3	0	4	11

min: 2 0 0 1

- 1) Find the minimum of each row, and subtract from each row the min value.
- 2) Find the minimum of each column, and subtract from each column the min value.
- 3) Find the min number of vertical/horizontal lines required to cover the Zeros in the matrix.

	Job1	Job2	Job3	Job4
Crane1	0	0	3	4
Crane2	3	0	7	4
Crane3	6	1	0	0
Crane4	1	0	4	10

	Job1	Job2	Job3	Job4
Crane1	0	0	3	4
Crane2	3	0	7	4
Crane3	6	1	0	0
Crane4	1	0	4	10

If the number of lines is equal to m , the optimal solution is available among the Covered Zeros. Otherwise, proceed to Step four.

$3 \neq m=4$

- 4) Find the min of uncovered values. Then, subtract the min from all the uncovered values and add it to the corner points. Then go back to step 3.

go back to step 3: $3 \neq m=4$

	Job1	Job2	Job3	Job4
Crane1	0	0	0	1
Crane2	3	0	4	1
Crane3	9	4	0	0
Crane4	1	0	1	7

	Job1	Job2	Job3	Job4
Crane1	0	0	0	1
Crane2	3	0	4	1
Crane3	9	4	0	0
Crane4	1	0	1	7

min=1

minimum number of lines required = $4 = m$ ✓
So, we are at the optimal table

Start the assignment from the row or column that has minimum number of zeros.

	Job1	Job2	Job3	Job4
Crane1	0	1	0	1
Crane2	2	0	3	0
Crane3	9	5	0	0
Crane4	0	0	0	6

	Job1	Job2	Job3	Job4
Crane1	0	1	0	1
Crane2	2	0	3	0
Crane3	9	5	0	0
Crane4	0	0	0	6

2 2 2 1

Hidden Markov Model (HMM)

- ✓ 之前我們所學的 classifier (如 perceptron、Naïve Bayes 等等)都是假設 samples 之間互相獨立。然而對於許多實際應用上是不好的假設，因此有了 **context-dependent classification** 的議題產生，像是分析 speech、music、gesture 之問題等等。如 **samples 之間都有關聯性**，對於每個 sample 我們可以透過 bind 在此 sample 前所有的 observations (所謂的 combo-feature-vectors)，來得到更多的資訊加以分析？然而，這會出現一些問題，例如：原本的 samples 維度低，但加上 observation 後維度就提高許多(curse of dimensionality)、不同的 class 可能有不同的 dimension、太多的 redundancy 等等。Markov process 是一個有效率且很直觀的方法，可以用來記錄 samples 之間的關係。
- ✓ Markov process 的假設為：在 time t 的 observation，只跟前 n 個 observations 有關係，稱之為 **n-order Markov process**。
- ✓ Deterministic Markov process 為每個 state 由前一個 state 決定 (例如：紅綠燈 traffic light)；Non-deterministic Markov process 與 deterministic 的概念相同，但對於 **states 之間存在轉移的機率** (如 state i 有多少的機率會轉移到 state j)。 **a_{ij} 為 state i 到 state j 的機率**、A 為所有 state 的轉移矩陣、對於所有的 state 到另一個 state 的機率加總為 1。

$$a_{ij} = P(j|i) \quad A = [a_{ij}] \quad \sum_j a_{ij} = 1$$

- ✓ 對於 state 無法直接可見的，我們稱之為 **Hidden Markov process**。一個常見的例子為：behind-the-curtain coin-tossing。在這個例子中我們可以看到一連串硬幣正反面的結果，然而，無法得知是在哪個 state 下所產生的結果(每個 state 都有機會產生正反面結果)。要**建立一個 Hidden Markov Model (HMM)** 需要給定 state 的數量、states 之間的轉移機率 a_{ij} 、state 初始的機率 $P(i)$ 以及在給定 state 情況下所產生的 observation 機率 $P(x|j)$ 。於許多實際應用中，observation 與 state 並不一定相同，以 text analysis 為例：observation 為獨立的 word，而 state 為 noun、verb、adjective 等等。
- ✓ HMM 又分為 **discrete (離散)** 與 **continuous (連續)**。如果 observation 為 continuous，在建立 discrete HMM 時就會產生一些失真。對於建立 continuous HMM，它

的 observables 就會建立在 continuous space 上， $P(x|j)$ 就會是一個 pdf。雖然 continuous HMM 會比較接近實際的應用，但計算上需耗費較多的時間。

✓ HMM 使用的範疇包含 evaluation、recongntion、decoding、training problem。

➤ **Evaluation problem:** 給定 HMM 與 obervation sequence，評估此 observations 由 HMM 產生的機率。 $P(X|S) = \sum_I p(X|I,S)P(I|S)$ 其中， S 為 HMM、 X 為 observation sequence、 I 為 state sequence。假設一個 HMM 為 fully connected 且有 K 的 state，observation sequence 長度為 t ，如以暴力解(brute-force)，則需要遍巡 K^t 個長度為 t 的 path (時間複雜度為 $O(tK^t)$)。此方法耗費太多時間，我們可以用動態規劃來降低計算量。定義 $\alpha_t(j)$ 為給定 HMM 下，從起始時刻到 time t ，observation sequeece 為 $x_1x_2...x_t$ ，且在 time t 時停在 state j 的機率。首先，需要初始化 $\alpha_1(j)$ ，其關係式如下列圖示。[forward procdeure]

$$\alpha_t(j) = p(x_1x_2...x_t, i_t = j | S) \quad \alpha_1(j) = p(x_1 | j)P(j)$$

當到達 time $t+1$ 時，其狀態與 time t 時的有關係。下列等式右方，左邊中括號項為在 time t 時所有 state 轉移到 state j 的機率總和；右邊項為 state j 產生下一個 observable (也就是 x_{t+1}) 的機率。

$$\alpha_{t+1}(j) = \left[\sum_i \alpha_t(i)P(j|i) \right] p(x_{t+1} | j)$$

當到達最終時刻(time T)時，obervation sequence 由 HMM 產生的機率就是 time T 時所有 state 的機率加總。

$$P(X|S) = \sum_j \alpha_T(j)$$

- **Recognition problem:** 為 **evaluation problem** 延伸。給定多個 HMMs (每個 HMM 代表不同的 class) 與 observation sequence，評估此 observation 由哪個 HMM 產生的機率較高。以下圖為例，計算 $P(HHH | HMM1)$ 。

$$\alpha_1(\text{Left}) = 0.5 \times 0.8 = 0.4 ; \quad \alpha_1(\text{Right}) = 0.5 \times 0.2 = 0.1。$$

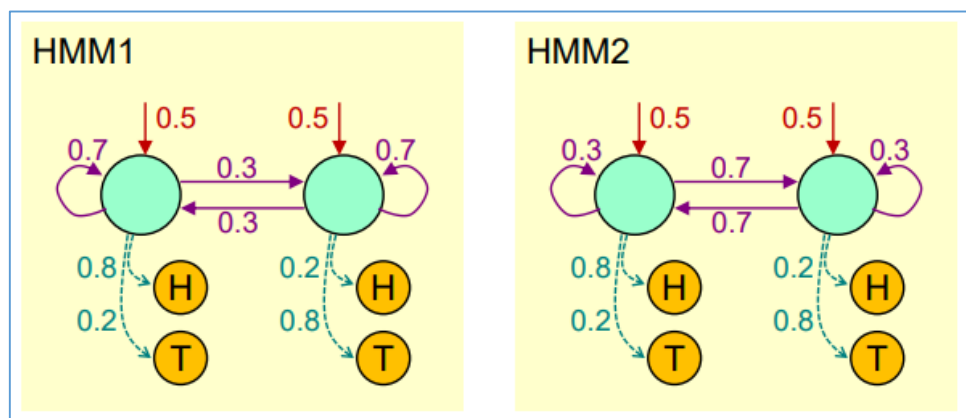
$$\alpha_2(\text{Left}) = (0.4 \times 0.7 + 0.1 \times 0.3) \times 0.8 = 0.248$$

$$\alpha_2(\text{Right}) = (0.4 \times 0.3 + 0.1 \times 0.7) \times 0.2 = 0.038$$

$$\alpha_3(\text{Left}) = (0.248 \times 0.7 + 0.038 \times 0.3) \times 0.8 = 0.148$$

$$\alpha_3(\text{Right}) = (0.248 \times 0.3 + 0.038 \times 0.7) \times 0.2 = 0.0202$$

$$P(X|S) = \alpha_3(\text{Left}) + \alpha_3(\text{Right}) = 0.148 + 0.0202 = 0.1682$$



Observations X	$P(X \text{HMM1})$	$P(X \text{HMM2})$
H H H	0.0168	0.0096
H H H H H H	0.0348	0.0084
H H H T T T	0.0218	0.0101
H T H T H T	0.0084	0.0348
H T H H H T	0.0129	0.0176

- **Decoding problem:** 給定 HMM 與 observation sequence，找到最有可能的 state sequence。 $I^* = \arg \max_I P(I | X, S)$ 如用暴力法(brute-force)，時間複雜度與 evaluation problem 提及的相同。同理，我們可以使用動態規劃方法解決，一個常見的演算法 **Viterbi algorithm**。其方法與 evaluation problem 很像，不同的是：每移動到下一步時，只選擇機率最大的那一步。因此，到達最後一步時也是選擇機率較大的結果。

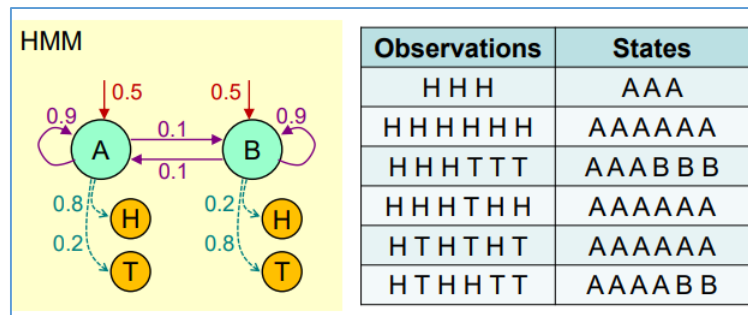
$$\delta_1(j) = p(x_1 | j)P(j) \quad \delta_{t+1}(j) = \left[\max_i \delta_t(i)P(j | i) \right] p(x_{t+1} | j)$$

我們可以發現算出來機率與問題本身機率不同，但因為兩機率之間只差一個倍數 $P(X | S)$ ，且此倍數不變。因此找 $P(I, X | S)$ 的最大值就等同於找 $P(I | X, S)$ 的最大值。

$$P(I^*, X | S) = \max_j \delta_T(j)$$

$$P(I | X, S) = \frac{P(I, X, S)}{P(X, S)} = \frac{P(I, X, S)}{P(X | S)P(S)} = \frac{P(I, X | S)}{P(X | S)}$$

以下圖為例，找出 **HTHHTT** 最可能的 state sequence。



H $\rightarrow \alpha_1(A) = 0.5 * 0.8 = 0.4$; $\alpha_1(B) = 0.5 * 0.2 = 0.1$ ($0.4 > 0.1$, 紀錄 **A**)

T $\rightarrow \alpha_2(A) = (\alpha_1(A) * 0.9) * 0.2 = 0.072$ ($\alpha_1(A) * 0.9 > \alpha_1(B) * 0.1$, 紀錄 **A**)

$\alpha_2(B) = (\alpha_1(B) * 0.9) * 0.8 = 0.072$ ($\alpha_1(B) * 0.9 > \alpha_1(A) * 0.1$, 紀錄 **B**)

H $\rightarrow \alpha_3(A) = (\alpha_2(A) * 0.9) * 0.8 = 0.05184$ ($\alpha_2(A) * 0.9 > \alpha_2(B) * 0.1$, 紀錄 **A**)

$\alpha_3(B) = (\alpha_2(B) * 0.9) * 0.2 = 0.01296$ ($\alpha_2(B) * 0.9 > \alpha_2(A) * 0.1$, 紀錄 **B**)

H $\rightarrow \alpha_4(A) = (\alpha_3(A) * 0.9) * 0.8 = 0.037348$ ($\alpha_3(A) * 0.9 > \alpha_3(B) * 0.1$, 紀錄 **A**)

$\alpha_4(B) = (\alpha_3(B) * 0.9) * 0.2 = 0.0023328$ ($\alpha_3(B) * 0.9 > \alpha_3(A) * 0.1$, 紀錄 **B**)

T $\rightarrow \alpha_5(A) = (\alpha_4(A) * 0.9) * 0.2 = 0.00672264$ ($\alpha_4(A) * 0.9 > \alpha_4(B) * 0.1$, 紀錄 **A**)

$\alpha_5(B) = (\alpha_4(B) * 0.9) * 0.8 = 0.00298484$ ($\alpha_4(B) * 0.9 > \alpha_4(A) * 0.1$, 紀錄 **A**)

T $\rightarrow \alpha_6(A) = (\alpha_5(A) * 0.9) * 0.2 = 0.0012100752$ ($\alpha_5(A) * 0.9 > \alpha_5(B) * 0.1$, 紀錄 **A**)

$\alpha_6(B) = (\alpha_5(B) * 0.9) * 0.8 = 0.00214900848$ ($\alpha_5(B) * 0.9 > \alpha_5(A) * 0.1$, 紀錄 **B**)

到最後一步時，逆推每一步最大機率所記錄的 state，因此會得 **AAAAAB**。

- **Training problem:** 給定多個 observation sequences，用來估計 HMM，其中包含的 components 有: states 之間的轉移機率 a_{ij} 、state 初始的機率 $P(i)$ 以及給定 state 情況下所產生的 observation 機率 $P(x|j)$ 。一個常見的方法 **Baum-Welch re-estimation method (EM algorithm)**。首先，定義 visit 的機率(給定 HMM 與 observation sequence，time t 時在 state i 的機率)，state 轉移機率(給定 HMM 與 observation sequence，time t 時在 state i 且 time $t+1$ 時在 state j 的機率)。

probability of "visit": $\gamma_t(i) = p(i_t = i | X, S)$

probability of "transition": $\xi_t(i, j) = p(i_t = i, i_{t+1} = j | X, S)$

$\beta_t(i)$ 與前述定義的 $\alpha_t(i)$ 不太一樣。 $\beta_t(j)$ 定義成給定 HMM 且 time t 時在 state i 的情況下，從 time t 到終止時刻 observation sequence 為 $x_{t+1}x_{t+2}\dots x_T$ 的機率。

$$\beta_t(i) = p(x_{t+1}x_{t+2}\dots x_T | i_t = i, S)$$

$$\alpha_t(j) = p(x_1x_2\dots x_t, i_t = j | S)$$

$\beta_t(i)$ 的關係是定義如下。

[backward procedure]

$$\beta_T(i) = 1, \quad \forall i$$

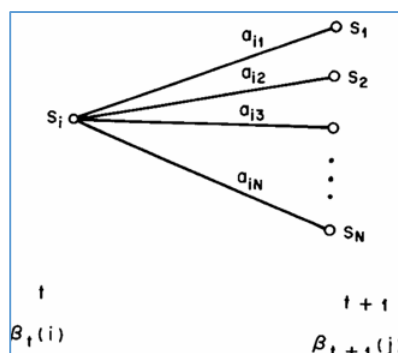
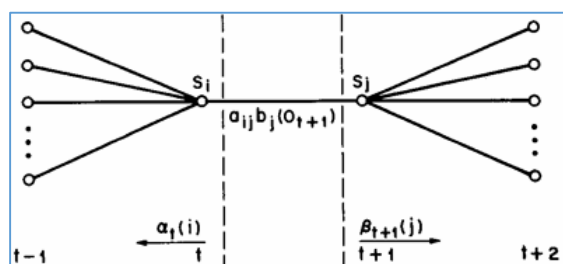
$$\beta_t(i) = \sum_j \beta_{t+1}(j) p(x_{t+1} | j) P(j | i), \quad T > t \geq 1$$

$$\gamma_t(i) = \frac{\alpha_t(i) \beta_t(i)}{\sum_j \alpha_t(j) \beta_t(j)}$$

$$\xi_t(i, j) = \frac{\alpha_t(i) P(j | i) p(x_{t+1} | j) \beta_{t+1}(j)}{\sum_{i'} \sum_{j'} \alpha_t(i') P(j' | i') p(x_{t+1} | j') \beta_{t+1}(j')}$$

$$r_t(i) = p(i_t = i | X, S) = \frac{p(i_t = i, X, S)}{p(X, S)} = \frac{p(i_t = i, X | S) p(S)}{p(X | S) p(S)} = \frac{p(i_t = i, X | S)}{p(X | S)} = \frac{\alpha_t(i) \beta_t(i)}{\sum \alpha_t(i) \beta_t(i)}$$

我們可以隨意初始化 HMM 的參數，並用上述的公式算出 $\alpha_t(i)$ 、 $\beta_t(i)$ 、 $\xi_t(i)$ 、 $\gamma_t(i)$ ，再將求出的結果帶入下列的公式更新 HMM 參數。重複相同的步驟直到收斂。



$$P(i) = \gamma_1(i)$$

$$P(j | i) = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)}$$

$$P(x | i) = \frac{\sum_{1 \leq t \leq T} \gamma_t(i)}{\sum_{1 \leq t \leq T} \gamma_t(i)} \quad x_t = x$$