# Multidimensional Scaling (MDS)

Note Title                                                                11/2/2010

MDS is a linear projection method. It is related to PCA. MDS and PCA can be used for non-linear projection — next lecture.

<u>Key Idea of MDS</u>: project to preserve the distances $|\underline{x}_i - \underline{x}_j|$ between the datapoints $\{\underline{x}_i : i=1 \text{ to } N\}$, i.e. $\underline{x}_i \to \underline{y}_i$ such that $|\underline{x}_i - \underline{x}_j| \approx |\underline{y}_i - \underline{y}_j|$, but the $y$'s have lower dimension.

This projection constraint is imposed on the dot products $\underline{x}_i \cdot \underline{x}_j \approx \underline{y}_i \cdot \underline{y}_j$ this will imply that $|\underline{x}_i - \underline{x}_j| \approx |\underline{y}_i - \underline{y}_j|$.

<u>Important Property</u>: We only need to know $|\underline{x}_i - \underline{x}_j| \stackrel{\Delta}{=} \Delta_{ij}$ in order to calculate $\underline{x}_i \cdot \underline{x}_j$. This will be useful for non-linear applications later. Also, sometimes only $|\underline{x}_i - \underline{x}_j|$ is specified.

<u>Result:</u> $\underline{x}_i \cdot \underline{x}_j = \dfrac{1}{2N} \sum_k \Delta_{ik}^2 + \dfrac{1}{2N} \sum_\ell \Delta_{\ell j}^2 - \dfrac{1}{2N^2} \sum_{\ell k} \Delta_{\ell k}^2 - \dfrac{1}{2}\Delta_{ij}^2$

provided $\sum_i \underline{x}_i = 0$ (Subtract $\dfrac{1}{N} \sum_{i=1}^{N} \underline{x}_i$ from data to ensure this)

<u>Proof</u>: $\Delta_{ij}^2 = |\underline{x}_i|^2 + |\underline{x}_j|^2 - 2 \underline{x}_i \cdot \underline{x}_j$

Let $T = \sum_i |\underline{x}_i|^2$, Note that $\sum_i \underline{x}_i \cdot \underline{x}_j = 0 = \sum_j \underline{x}_i \cdot \underline{x}_j$.

Thm $\sum_k \Delta_{ik}^2 = N|\underline{x}_i|^2 + T$, $\sum_\ell \Delta_{\ell j}^2 = N|\underline{x}_j|^2 + T$, $\sum_{\ell k} \Delta_{\ell k}^2 = 2NT$.

the result follows by substitution.

Define the <u>Gram matrix</u>
$$G_{ij} = \underline{x}_i \cdot \underline{x}_j = \sum_k \Delta_{ik}^2 + \dfrac{1}{2N} \sum_\ell \Delta_{\ell j}^2 - \dfrac{1}{2N^2} \sum_{\ell k} \Delta_{\ell k}^2 - \dfrac{1}{2}\Delta_{ij}^2.$$

Define an error function
$$\text{err}(\underline{y}) = \sum_{ij} \left( G_{ij} - \underline{y}_i \cdot \underline{y}_j \right)^2$$

$\underline{x}_i$ lies in $D$-dim space, $G_{ij}$ is $N \times N$ matrix
$\underline{y}_i$ is a vector in $d$-dim space     $d \ll D$
                                                    $d < N$

(2)

Minimize $\quad err(\underline{y}) = \sum_{ij} (G_{ij} - \underline{y}_i \cdot \underline{y}_j)^2$ .

Do spectral decomposition:

$$G = \sum_{\alpha=1}^{N} \lambda_\alpha \underline{v}_\alpha \underline{v}_\alpha^T$$

$\lambda_1 \geq \ldots \geq \lambda_n \geq 0$
eigenvalues of $\underline{\underline{G}}$

$\underline{v}_\alpha \cdot \underline{v}_\beta = \delta_{\alpha\beta}$
eigenvectors

Claim: optimal minimization
is $\quad \underline{y}_\alpha^i = \sqrt{\lambda_\alpha} v_i^\alpha$

ie. $\underline{y}_i = (\sqrt{\lambda_1} v_i^1, \sqrt{\lambda_2} v_i^2, \ldots, \sqrt{\lambda_n} v_i^N)$

$\underbrace{\qquad\qquad}_{N\text{-dimension}}$

Proof. $\underline{y}_i \cdot \underline{y}_j = \sum_\alpha \sqrt{\lambda_\alpha} v_i^\alpha \sqrt{\lambda_\alpha} v_j^\alpha = \sum_\alpha \lambda_\alpha \underline{v}_\alpha \underline{v}_\alpha^T = G_{ij}$

This gives $err = 0$.

## We can reduce the dimension by

truncating $\underline{y}_i$ to $\quad \underline{y}_i = (\sqrt{\lambda_1} v_i^1, \ldots, \sqrt{\lambda_d} v_i^d)$

$\qquad\qquad\qquad\qquad$ for $d < N$

In this case $\quad G_{ij} \neq \underline{y}_i \cdot \underline{y}_j$

$$\underline{y}_i \underline{y}_j = \sum_{\alpha=1}^{d} \lambda_\alpha v_i^\alpha v_j^\alpha \quad, \quad G_{ij} = \sum_{\alpha=1}^{N} \lambda_\alpha v_i^\alpha v_j^\alpha$$

Hence $\quad G_{ij} - \underline{y}_i \underline{y}_j = \sum_{\alpha=d+1}^{N} \lambda_\alpha v_i^\alpha v_j^\alpha$

Claim $\quad \sum_{ij} (G_{ij} - \underline{y}_i \underline{y}_j)^2 = \sum_{\alpha=d+1}^{N} \lambda_\alpha^2$

Proof. $\quad \sum_{ij} \left( \sum_{\alpha=d+1}^{N} \sum_{\beta=d+1}^{N} \lambda_\alpha \lambda_\beta v_i^\alpha v_j^\alpha v_i^\beta v_j^\beta \right)$

$\qquad\qquad\qquad \sum_i v_i^\alpha v_i^\beta = \delta^{\alpha\beta} , \quad \sum_j v_j^\alpha v_j^\beta = \delta^{\alpha\beta}$

$\sum_{\alpha=d+1}^{N} \sum_{\beta=d+1}^{N} \lambda_\alpha \lambda_\beta \delta^{\alpha\beta} \delta^{\alpha\beta} = \sum_{\alpha=d+1}^{N} \lambda_\alpha^2$.

Hence we project to $d$-dimension

provided $\quad \sum_{\alpha=d+1}^{N} \lambda_\alpha^2$ is small, or $\quad \dfrac{\sum_{\alpha=d+1}^{N} \lambda_\alpha^2}{\sum_{\alpha=1}^{N} \lambda_\alpha^2}$ is small.

$\underline{y}_i = (\sqrt{\lambda_1} v_i^1, \ldots \sqrt{\lambda_d} v_i^d)$

(3)

# Relation between MDS & PCA?

Both linear. Both depend on eigenvectors/eigenvalue.

Recall PCA    ( subtract mean to ensure $\sum_i \underline{x}_i = 0$ )

$$\underline{x}_p = (\underline{x}.\underline{e}_1, \ldots, \underline{x}.\underline{e}_d)$$

the $\underline{e}$'s are eigenvectors of $K_{ab} = \frac{1}{N} \sum_{i=1}^{N} \underline{x}_i \cdot \underline{x}_i^T$

MDS    $\underline{y}_i = (\sqrt{\lambda_1} v_i^1, \ldots, \sqrt{\lambda_d} v_i^d)$

the $\underline{v}$'s are eigenvector of $G_{ij} = \sum_{i=1}^{N} \underline{x}_i \cdot \underline{x}_j$

Claim: the eigenvalues of $\underline{G}$ and $\underline{K}$ are the same. The eigenvectors are closely related.

Proof  Let $\underline{X}$ be an $N \times D$ matrix     $i = 1 \text{ to } N$    no. of pts
with elements $X_{ia}$              $a = 1 \text{ to } D$    space dimension

$a^{th}$ component of $i^{th}$ datapoint

Consider   $\underline{X} \underline{X}^T$    $N \times N$ matrix,  $\left(\underline{X} \underline{X}^T\right)_{ij} = \sum_a X_{ia} X_{ja}$

$\underline{X}^T \underline{X}$    $D \times D$ matrix ,  $\left(\underline{X}^T \underline{X}\right)_{ab} = \sum_i X_{ia} X_{ib}$

Both are square matrices  and both are positive definite, so they have positive eigenvectors.

$\underline{X} \underline{X}^T$ is used for MDS,  $\underline{X}^T \underline{X}$ is used for PCA

Suppose $\underline{e}, \lambda$ are an eigenvector, eigenvalue of $\underline{X}^T \underline{X}$

$$\underline{X}^T \underline{X} \, \underline{e} = \lambda \underline{e}$$

So    $\underline{X} \, \underline{X}^T \underline{X} \underline{e} = \lambda \underline{X} \underline{e}$

$$(\underline{X} \underline{X}^T)(\underline{X} \underline{e}) = \lambda (\underline{X} \underline{e})$$

So   $(\underline{X} \underline{e})$ is an eigenvector (un-normalized)
of $\underline{X} \underline{X}^T$ with eigenvalue $\lambda$.

Similarly if   $\underline{X} \underline{X}^T \underline{v} = \lambda \underline{v}$
then $(\underline{X}^T \underline{v})$ is an eigenvector (un-normalized)
of $\underline{X}^T \underline{X}$ with eigenvalue $\lambda$.

Conclusion → the two matrices have the same eigenvalues
and related eigenvectors.

(4)

Result → the truncation conditions for MDS and PCA are similar

$$\frac{\sum_{i=1}^{\hat{d}} \lambda_i^2}{\sum_{i=1}^{} \lambda_i^2} > \text{Threshold.}$$

MDS projects $\underline{x}_i$

to $\underline{y}_i = \left( \sqrt{\lambda_1} \, v_i^1, \dots, \sqrt{\lambda_d} \, v_i^d \right)$

PCA projects $\underline{x}$ to $\underline{y} = \left( \underline{x} \cdot \underline{e}_1, \dots, \underline{x} \cdot \underline{e}_d \right)$

where the $\underline{e}$'s and the $\underline{v}$'s are related (see previous page).

Note: The equivalence between the eigenvalues of $\underline{\underline{X}} \, \underline{\underline{X}}^T$ and $\underline{\underline{X}}^T \underline{\underline{X}}$ has computational importance. If $N \ll D$, faster to compute eigenvalues/vectors for $(\underline{\underline{X}} \, \underline{\underline{X}}^T)$, then convert to eigenvalues/vectors of $\underline{\underline{X}}^T \underline{\underline{X}}$

Note: to do PCA we have to know the data $\{\underline{x}_i\}$ but to do MDS we only need to know $\Delta_{ij}$. In some applications it is possible to specify $\Delta_{ij}$ but not the $\{\underline{x}_i\}$.

Deeper Understanding: This relationship between $\underline{\underline{X}} \, \underline{\underline{X}}^T$ and $\underline{\underline{X}}^T \underline{\underline{X}}$ can be used to prove SVD:

$$\underline{\underline{X}} = \underline{\underline{F}} \, \underline{\underline{D}} \, \underline{\underline{E}} \qquad \text{where } \underline{\underline{F}} \, \underline{\underline{F}} = \underline{\underline{I}} \quad \times \text{Identity.}$$
$$\underline{\underline{D}} = \begin{pmatrix} d_1 & 0 \\ 0 & \ddots \, d_v \end{pmatrix} \qquad \underline{\underline{E}} \, \underline{\underline{E}}^T = \underline{\underline{I}}$$

This is a generalization of the spectral decomposition

$$\underline{\underline{G}} = \sum \lambda_a \underline{v}_a \underline{v}_a^T \qquad \underline{\underline{X}} \text{ is } N \times D \quad N \neq D$$

to any matrix $\underline{\underline{X}} \rightarrow$ so $\underline{\underline{X}}$ is not square.

It follows that $\underline{\underline{X}}^T \underline{\underline{X}} = \left( \underline{\underline{E}}^T \underline{\underline{D}} \, \underline{\underline{F}}^T \right) \left( \underline{\underline{F}} \, \underline{\underline{D}} \, \underline{\underline{E}} \right)$
$$= \underline{\underline{E}}^T \underline{\underline{D}}^2 \underline{\underline{E}}$$

spectral decomposition → with $\lambda_1 = d_1^2, \lambda_2 = d_2^2 \dots$

and $\underline{\underline{X}} \, \underline{\underline{X}}^T = \underline{\underline{F}} \, \underline{\underline{D}} \, \underline{\underline{E}} \, \underline{\underline{E}}^T \underline{\underline{D}} \, \underline{\underline{F}}^T = \underline{\underline{F}} \, \underline{\underline{D}}^2 \underline{\underline{F}}^T$

So SVD is like the square root of spectral decomposition.!