# NEUR 608: Dimensionality Reduction

Bratislav Misic
McConnell Brain Imaging Centre
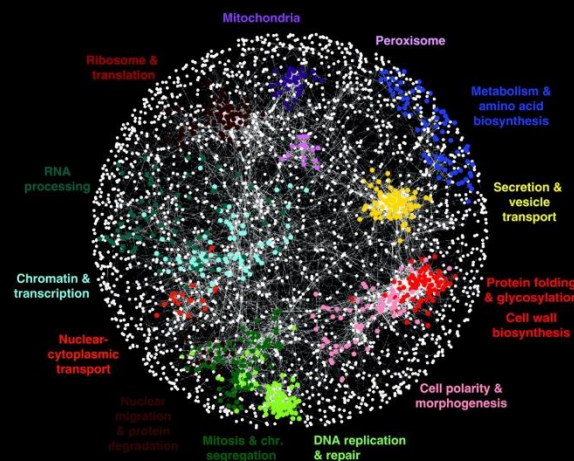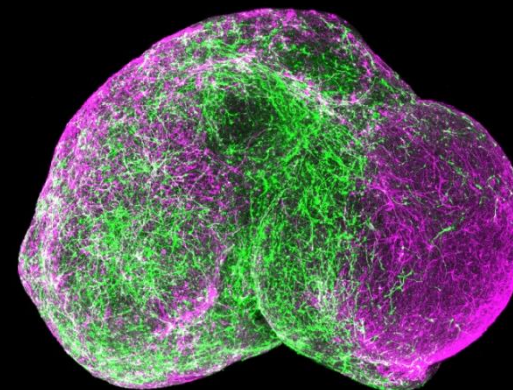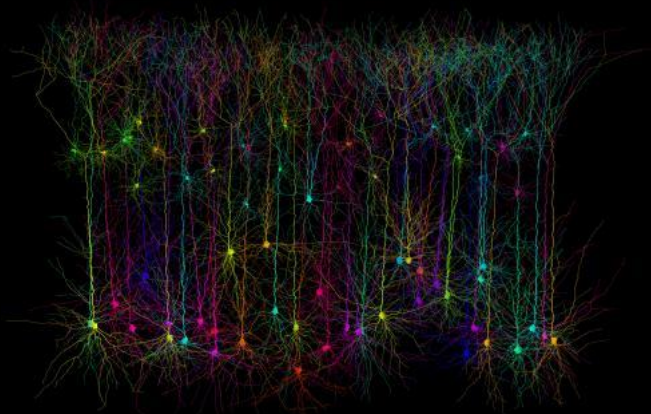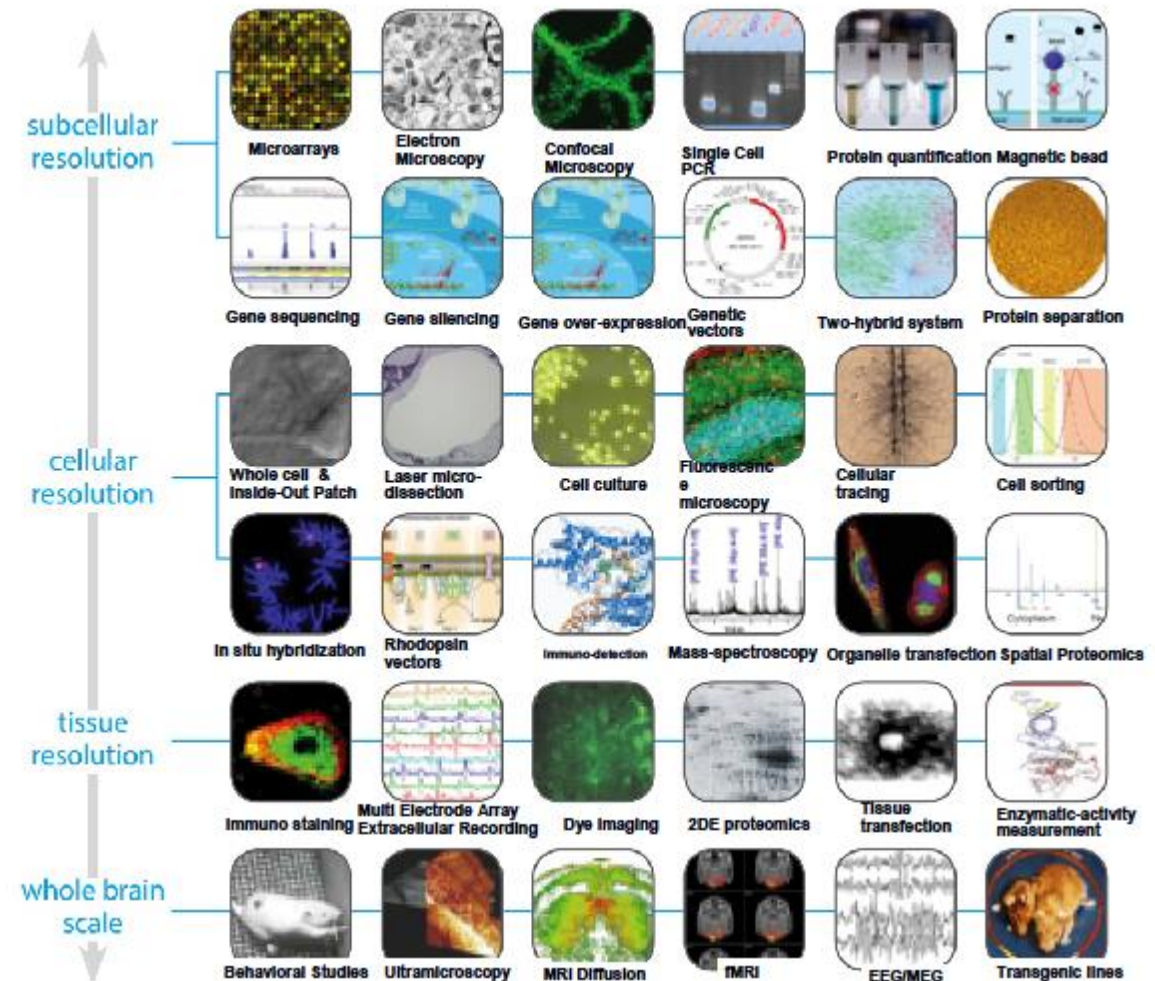Montréal Neurological Institute
McGill University

# Towards multivariate analysis

- **activation**: mapping individual elements to individual external variables

- **connectivity**: mapping individual structural or functional connections to individual external variables

- **networks**: mapping network attributes to individual external variables

- **multivariate systems**: mapping patterns of elements or connections to patterns of external variables
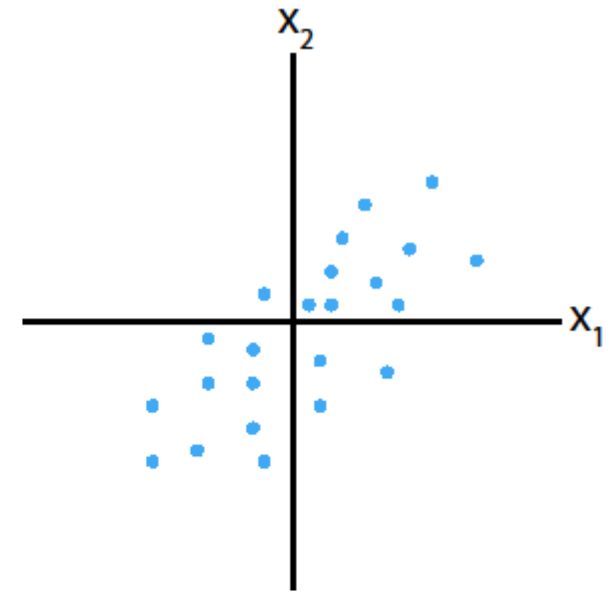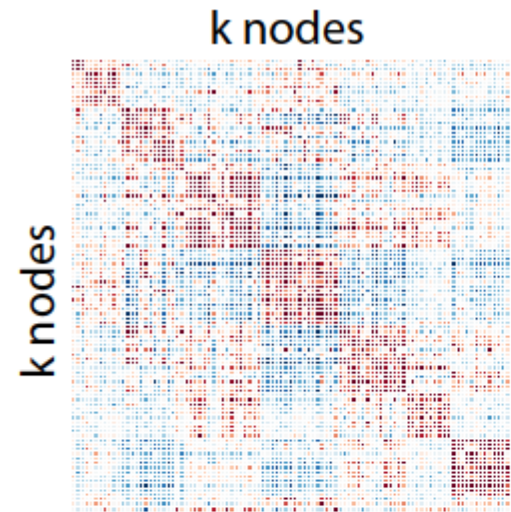


Misic & Sporns (2016) Curr Opin Neurobiol

- How to deal with more variables than observations?

- How to operationalize the network property?

- How to relate multiple data sets to each other?

- Goal #1: why dimensionality reduction?

    - linear: PCA, FA, ICA, MDS
    - nonlinear: kernel PCA, LLE, diffusion maps, t-SNE, autoencoders

- Goal #2: how is similarity expressed?
- Goal #3: is there an underlying generative model?
- Goal #4: statistical inference and reliability (and lack thereof)

- Goal #5: examples + demo

# Principal component analysis (PCA)

k nodes

k nodes

- original variables $x_1$ and $x_2$ are correlated but neither captures the dominant pattern of variance

- want a new variables $z_1$ and $z_2$ that
  (a) capture as much variance
  (b) are mutually uncorrelated

- need to find a rotation **u** to re-align the original axes (variables)

- **u** is chosen to maximize the variance of the new variables **z**

$X_2$

$X_1$

Hotelling (1933) J Educ Psychol

# Maximzing variance

- want to find a new variable $\mathbf{z} = \mathbf{Xu}$

- choose $\mathbf{u}$ to maximize $var(\mathbf{z})$

  under the constraint that $\mathbf{u'u} = 1$

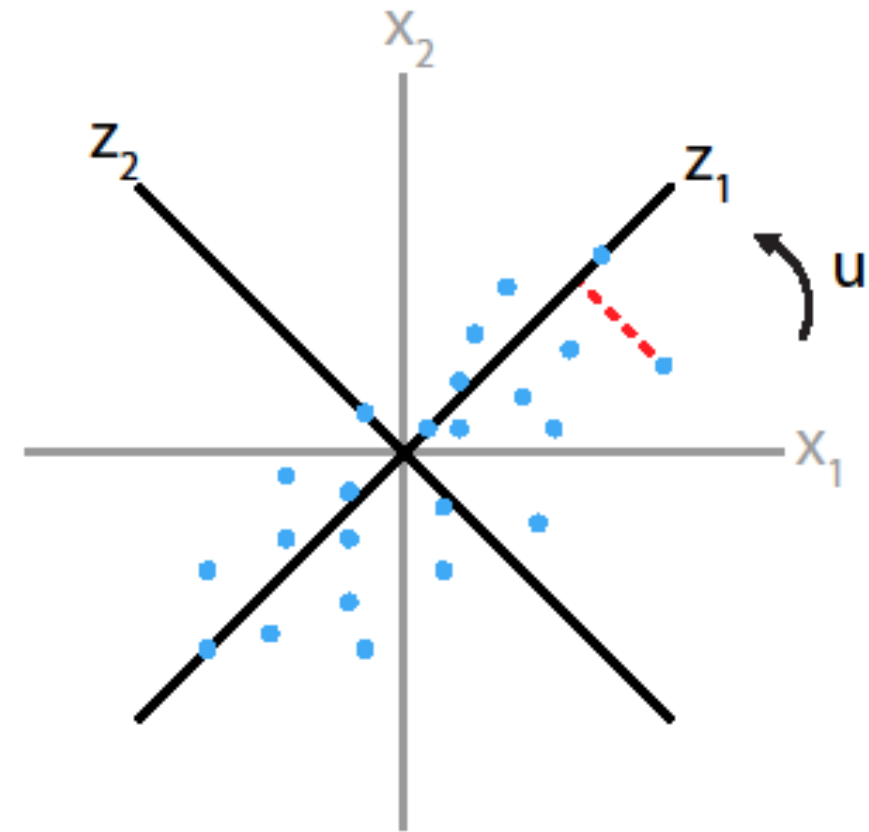- $var(\mathbf{z}) = \frac{1}{n-1}\mathbf{u'X'Xu} = \mathbf{u'Ru}$

  where $\mathbf{R} = \frac{1}{n-1}\mathbf{X'X}$

- define Lagrangian: $L = \mathbf{u'Ru} - \lambda(\mathbf{u'u} - 1)$

- find maximum: $\frac{\partial L}{\partial \mathbf{u}} = 2\mathbf{Ru} - 2\mathbf{u}\lambda = 0$

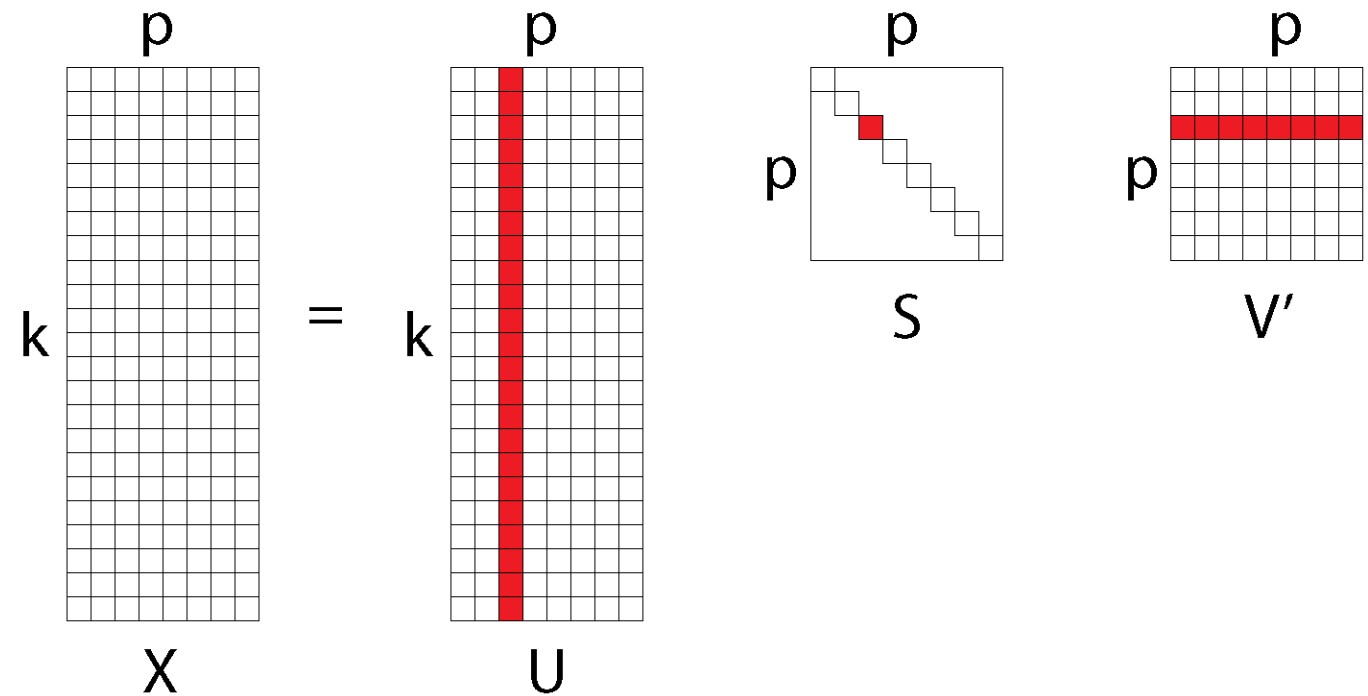  $$\mathbf{Ru} = \lambda\mathbf{u} \qquad (\mathbf{R} - \lambda\mathbf{I})\mathbf{u} = 0$$

- $\lambda$ = eigenvalue, $\mathbf{u}$ = eigenvector

- $var(\mathbf{z}) = \mathbf{u'Ru} = \mathbf{u'u}\lambda = \lambda$
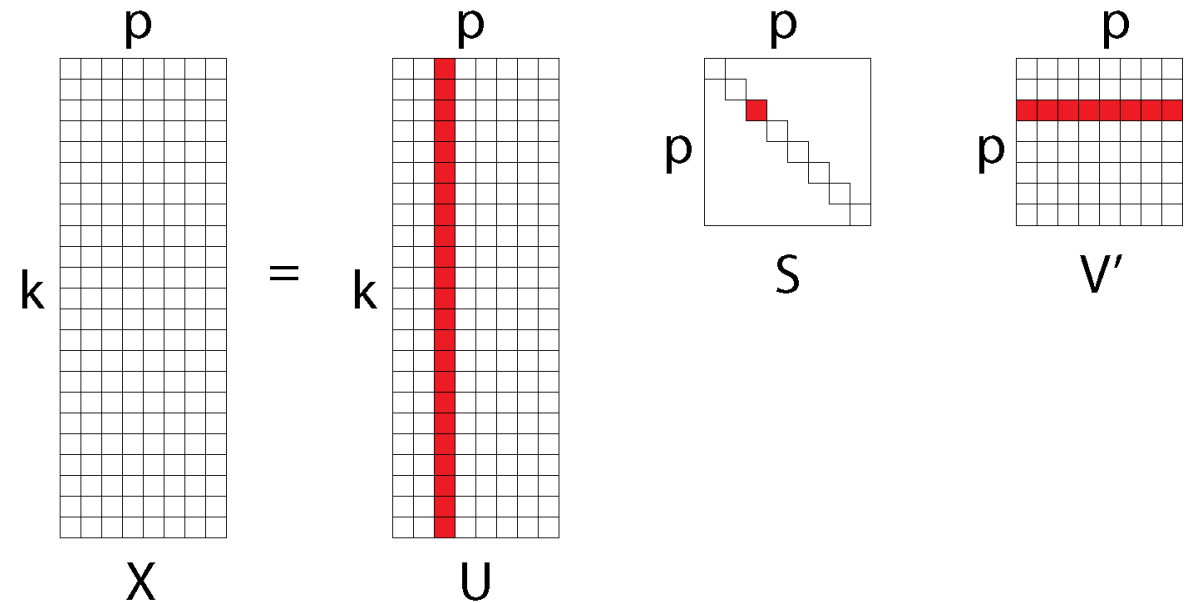
# Singular value decomposition (SVD)

- how can we factorize a data matrix?

  - $SVD(\mathbf{X}) = \mathbf{USV}'$

- this is a generalization of the spectral decomposition:

  - $EIG(\mathbf{X}'\mathbf{X}) = \mathbf{U\Lambda U}'$
  - $EIG(\mathbf{XX}') = \mathbf{V\Lambda V}'$

- **U** and **V** are orthonormal singular vectors; represent how you should weigh the original variables in **X**

- **S** is a diagonal matrix of singular values; represent how strongly paired the **U** and **V** are



Eckart & Young (1936) Psychometrika
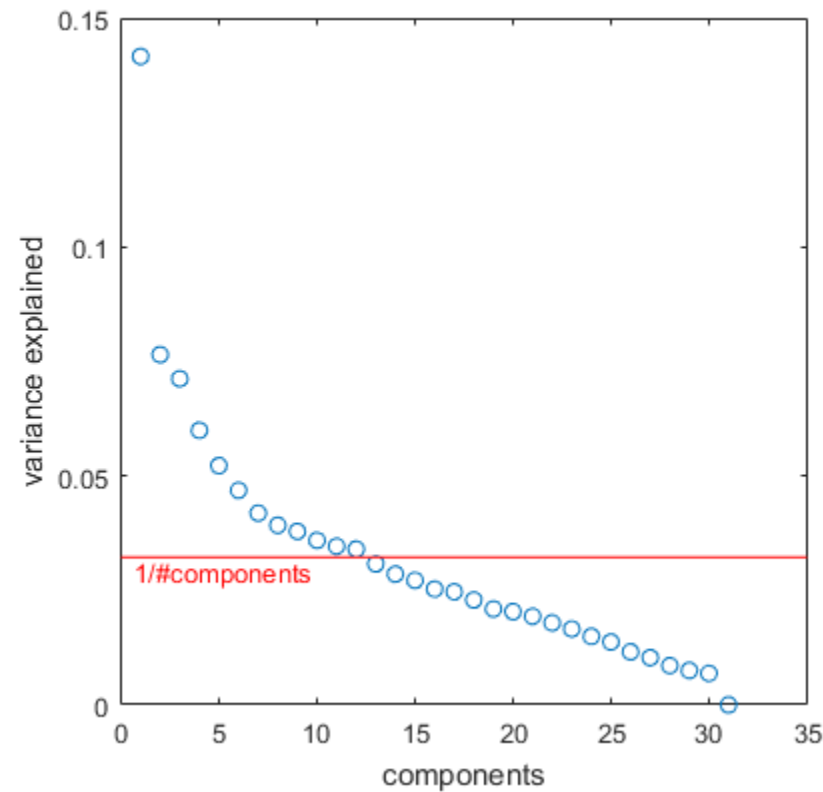
# Singular value decomposition (SVD)

- $SVD(\mathbf{X}) = \mathbf{USV}'$

- $EIG(\mathbf{X}'\mathbf{X}) = \mathbf{U\Lambda U}'$
- $EIG(\mathbf{XX}') = \mathbf{V\Lambda V}'$

- $\mathbf{X}'\mathbf{X} = (\mathbf{VS}'\mathbf{U}')(\mathbf{USV}') = \mathbf{VS}'(\mathbf{U}'\mathbf{U})\mathbf{SV}'$
  $= \mathbf{V}(\mathbf{S}'\mathbf{S})\mathbf{V}' = \mathbf{V\Lambda V}$

- $\mathbf{XX}' = (\mathbf{USV}')(\mathbf{VS}'\mathbf{U}') = \mathbf{US}(\mathbf{V}'\mathbf{V})\mathbf{S}'\mathbf{U}'$
  $= \mathbf{U}(\mathbf{SS}')\mathbf{U}' = \mathbf{U\Lambda U}'$

- eigenvector of $\mathbf{XX}'$ = left singular vector of $\mathbf{X}$

- eigenvector of $\mathbf{X}'\mathbf{X}$ = right singular vector of $\mathbf{X}$

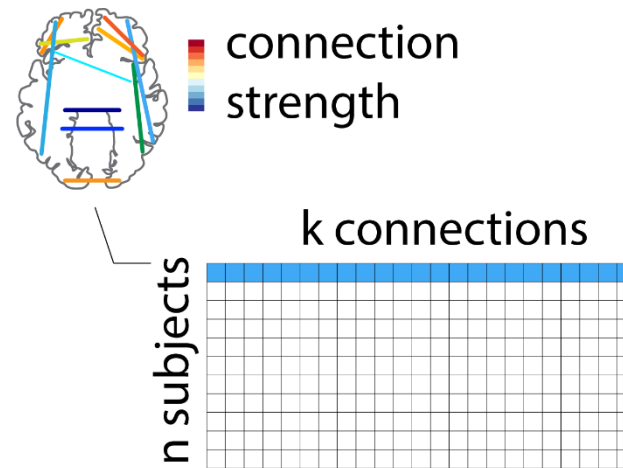- eigenvalue  = squared singular value



Eckart & Young (1936) Psychometrika
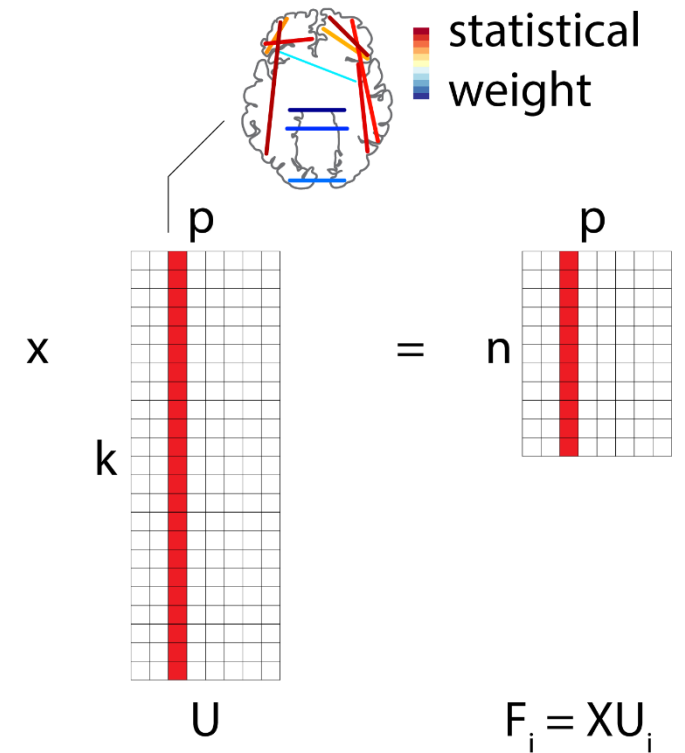
# How many components to retain?



- Kaiser criterion: drop all eigenvalues lower than 1
- "elbow" rule: look for the biggest drop-off in the scree plot

- weights that tell us how much each variable (connection, region, behaviour) contributes to the pattern

- but every participant is different – how do individual participants express the overall pattern?

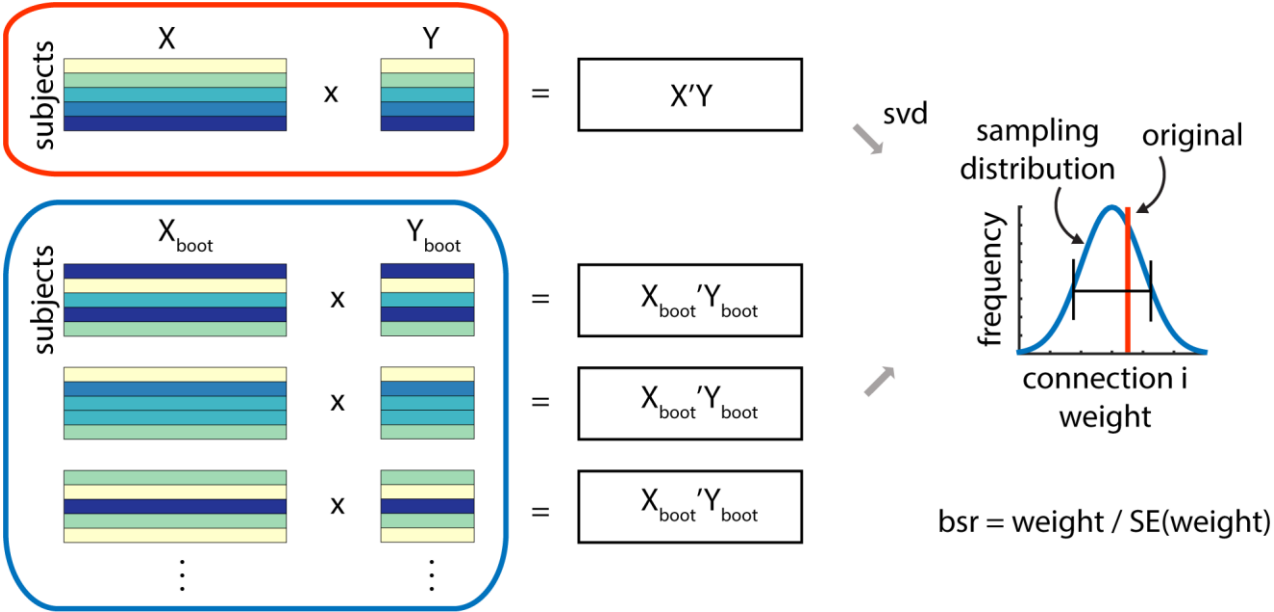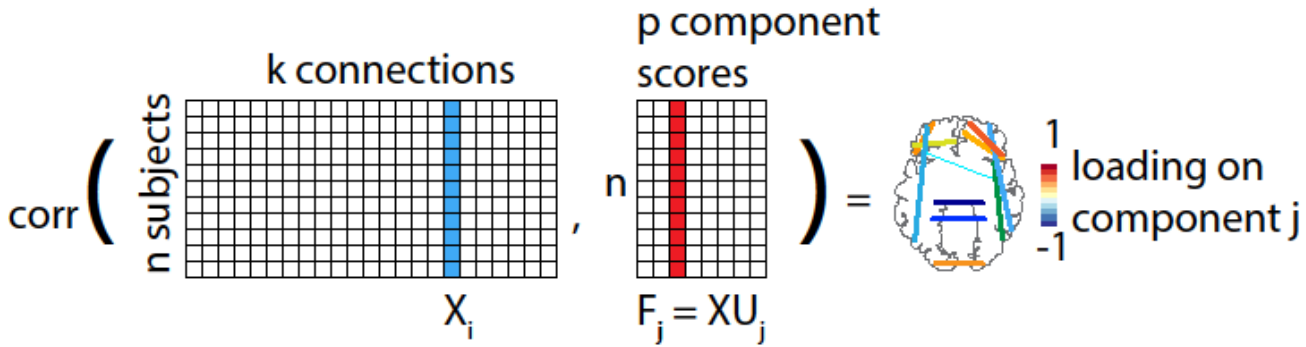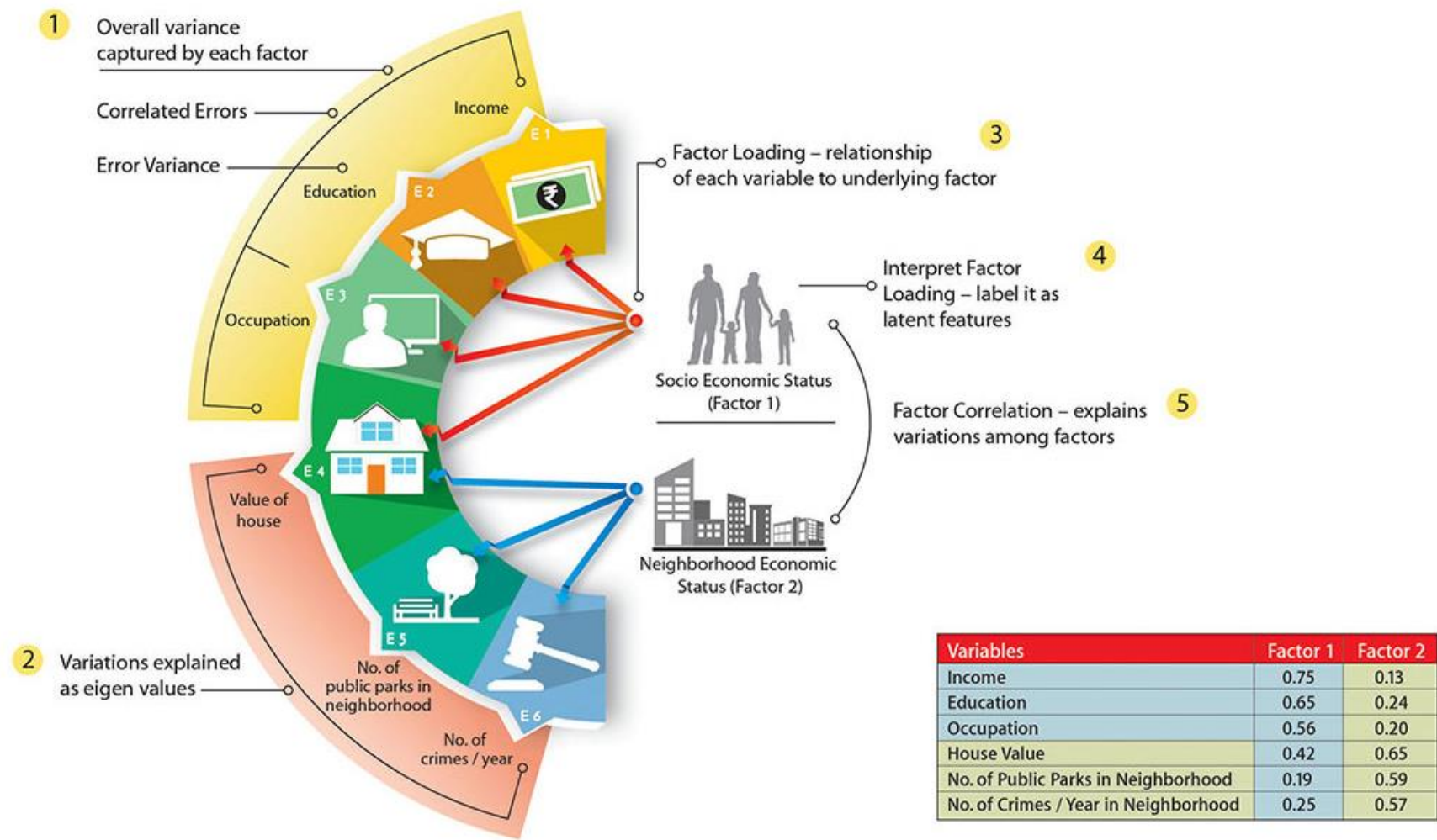- participant scores – project original data onto the latent variables

connection strength

statistical weight

k connections

n subjects

$X$

$p$

$k$

$U$

$\times$

$=$

$p$

$n$

$F_i = XU_i$

# Which variables are important?

- **singular vector weight**: how much should I weigh the variables?

- **loading**: does expression of the original variable correlate with expression of the latent variable?

- **"bootstrap ratio"**: which original variables are weighed highly but also stable across participants?

| Variables | Factor 1 | Factor 2 |
|---|---|---|
| Income | 0.75 | 0.13 |
| Education | 0.65 | 0.24 |
| Occupation | 0.56 | 0.20 |
| House Value | 0.42 | 0.65 |
| No. of Public Parks in Neighborhood | 0.19 | 0.59 |
| No. of Crimes / Year in Neighborhood | 0.25 | 0.57 |

- Assume that variance in variable $X_i$ comes from a set of latent factors $\eta_i$ and measurement error $\delta_i$:

  - $X_1 = \lambda_{1,1}\eta_1 + \lambda_{1,2}\eta_2 + \delta_1$
  - $X_2 = \lambda_{2,1}\eta_1 + \lambda_{2,2}\eta_2 + \delta_2$
  ...

- Unique factors $\delta_i$ are uncorrelated, so they only contribute to the diagonals of the covariance matrix

  - $var(X_i) = var(\lambda_{i,1}\eta_1 + \lambda_{i,2}\eta_2 + \delta_i)$

  - $var(X_i) = \lambda_{i,1}^2 + \lambda_{i,2}^2 + \theta_i^2$, where $\theta_i^2 = var(\delta_i)$

- The idea:

  - Estimate $\theta_i^2$, subtract them, and decompose the matrix
  - PCA: decompose correlation matrix **R**, with 1's on diagonal
  - FA: decompose matrix with diagonal elements $1 - \theta_i^2$

- How do we estimate $1 - \theta_i^2$ ("communalities")?

  - option #1: guess and re-adjust
  - option #2: estimate squared multiple
    correlation (SMC), i.e. regress $X_i$ on all $X_{j \neq i}$

- Procedure:

  1. calculate $\boldsymbol{R}$
  2. calculate $\boldsymbol{R}_{adj}$
  3. $EIG(\boldsymbol{R}_{adj})$
  4. are communalities stable?
     - no: recalculate $\boldsymbol{R}_{adj}$
     - yes: stop

- We want to improve interpretability: each factor has high loadings on only a few variables, and near zero for all others

  - orthogonal: factors remain uncorrelated (e.g. *varimax*)
  - oblique : factors may correlate (e.g. *promax, oblimin*)

- If we perform an oblique rotation, it is no longer obvious how to interpret factors:

  - structure loading
  - pattern loading



Orthogonal Rotation    Oblique Rotation

PSYC 4310/6310   Experimental Methods and Statistics   © 2014, Michael Kalsher

# Independent Components Analysis (ICA)



- Blind source separation

- Assumption: there exist a finite number of sources (independent components) that are mixed to produce observed neural activity

- Solution:

  - Any linear mixture of independent variables (e.g. voxels) will be more Gaussian than the original variables

  - Create new axes with maximally non-Gaussian projections

  - Many algorithms + objective functions

- Considerations:

    - Model selection
    - Starting point
    - Whitening
    - Algorithm
    - Subject- vs group-level
    - Spatial vs temporal independence

- Used for both preprocessing and for discovery

connection weight     transition probability

A Human     B Macaque monkey

C    D    E

1. Convert correlation matrix to a transition probability matrix

2. Compute diffusion operator (Laplacian):

3. Get eigenvectors and eigenvalues of Laplacian

- The resulting components reflect dominant modes of diffusion

- Points (e.g. voxels) are grouped by how close or accessible they are to each other by a random walker

Margulies et al. (2016) PNAS

# t-SNE



$$p_{j|i} = \frac{\exp(-||x_i - x_j||^2/2\sigma_i^2)}{\sum_{k \neq i} \exp(-||x_i - x_k||^2/2\sigma_i^2)} \quad (1)$$

$$p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n}$$

$$q_{ij} = \frac{(1 + ||y_i - y_j||^2)^{-1}}{\sum_{k \neq l} (1 + ||y_k - y_l||^2)^{-1}} \quad (2)$$

$$C = \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}} \quad (3)$$

- t-distributed stochastic neighbour embedding: visualizing data in 2 or 3 dimensions

- Input: data $X$ and perplexity σ

1. Estimate high-dimensional distance P
2. Estimate low-dimensional distance Q
3. Estimate distance between distances



van der Maaten & Hinton (2008) J Mach Learn Res

- Idea: input and output layers have the same number of nodes (i.e. equal to the number of variables), but hidden layers with fewer nodes

- The neural network has to efficiently encode and then reconstruct the input

- Typically feedforward, non-recurrent

- If activations are linear, autoencoders will approximate SVD (Bourlard & Kemp, 1988)

$X$

$\hat{X}$

Input Layer    Hidden Layer    Output Layer

- Overfitting

- Linear or nonlinear?

- Identifiable?

- Unique partitioning of variance

- Inference on individual variables

# Demo: PPMI

- Parkinson's Progression Markers Initiative (http://www.ppmi-info.org)

- Deep phenotyping in patients with Parkinson's disease

- n = 232 *de novo* patients

- p = 31 clinical, demographic and physiological variables

- Question: can we feasibly reduce this 31-dimensional dataset? Can we find dominant latent patterns?



Marek et al. (2011) Prog Neurobiol