



**Melbourne Bioinformatics**

BIOINFORMATICS + DATA SERVICES + INFRASTRUCTURE, FOR LIFE SCIENCES TODAY



# COMP90014

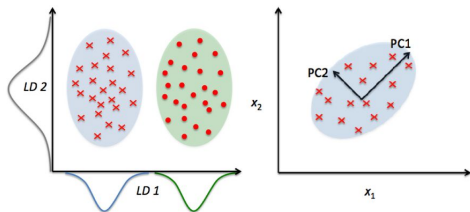
Algorithms for Bioinformatics

Week 9A: Dimensionality Reduction I

# Dimensionality Reduction

1. Why do we need dimensionality reduction?
2. Approaches for dimensionality reduction
3. Principal Component Analysis (PCA)
4. Multi-dimensional scaling (MDS)

# Why do we need dimensionality reduction?



## Feature selection

- 🍇 greedy feature selection

## Feature extraction

- 🍇 unsupervised methods
- 🍇 supervised methods
- 🍇 principal component analysis (PCA)
- 🍇 singular value decomposition (SVD)

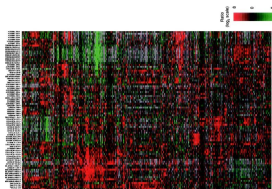
# Datasets have huge numbers of features (dimensions)



**Documents:** thousands of words

**Images:** thousands to millions of pixels

**Genomics:** thousands of genes, millions of DNA variants



e.g. a gene expression microarray

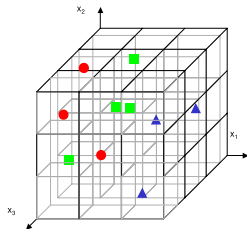
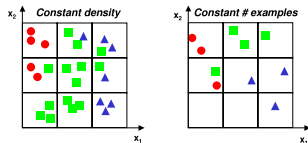
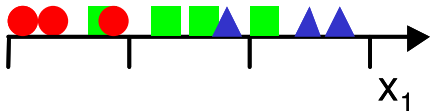
- 🍇 100 samples from *Arabidopsis thaliana*
- 🍇 data from 27 000 genes
- 🍇 what is the dimensionality of this dataset?
- 🍇 27 000

# High dimensionality has many costs



- data interpretation and visualisation
- redundant and irrelevant features  
degrade performance of ML  
algorithms
- computational cost may become  
intractable

# The curse of dimensionality



- as the number of dimensions increases the data become **sparse**
- an huge amount of data is needed to “cover” all the dimensions
  - number of data points needed **grows exponentially with the number of dimensions**

# Goal of dimensionality reduction



1. transform data from high-dimensional space into low-dimensional space
2. retain meaningful properties of the original data



# Dimensionality Reduction

1. Why do we need dimensionality reduction?
2. Approaches for dimensionality reduction
3. Principal Component Analysis (PCA)
4. Multi-dimensional scaling (MDS)



## Two approaches to perform dimensionality reduction

$$\mathbb{R}^N \rightarrow \mathbb{R}^M \quad (M < N)$$

- the goal is to find a low-dimensional representation of the data
- preserving (most of) the information or structure in the data

Feature selection: choosing a subset of all existing features

$$[x_1, x_2, \dots, x_N] \xrightarrow{\text{feature selection}} [x_{i_1}, x_{i_2}, \dots, x_{i_M}]$$

Feature extraction: creating new features by combining existing ones

$$[x_1, x_2, \dots, x_N] \xrightarrow{\text{feature extraction}} [y_1, y_2, \dots, y_M] = f([x_{i_1}, x_{i_2}, \dots, x_{i_m}])$$

# Feature selection

---

## Forward Greedy Feature Selection

---

```
1  $FS^{(0)} = \emptyset; F^{(0)} = \{f_1, f_2, \dots, f_n\};$   
   $i = 0; \text{opt} = 0; \text{iter} = 0;$   
2 while  $i < n$  do  
3    $k = \text{size}(F^{(i)})$   
4    $\text{max} = 0; \text{feature} = 0$   
5   for  $j = 1$  to  $k$  do  
6      $\text{score} = \text{eval}(F_j^{(i)})$   
7     if  $\text{score} > \text{max}$  then  
8        $\text{max} = \text{score}; \text{feature} = F_j^{(i)}$   
9   if  $\text{max} > \text{opt}$  then  
10     $\text{opt} = \text{max}; \text{iter} = i$   
11     $FS^{i+1} \leftarrow FS^{(i)} + \text{feature}$   
12     $F^{i+1} = F^{(i)} - \text{feature}$   
13     $i++$ 
```

---

## Supervised feature selection

- 🌸 optimise an objective function  
e.g. F1 metric

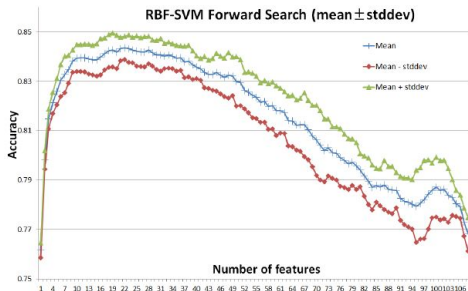
## Algorithms for supervised feature selection

- 🌸 Backward Greedy Feature Elimination
- 🌸 Forward Greedy Feature Selection

## Forward Greedy Feature Selection

1. choose the feature with the highest score
2. add feature to the pool
3. re-score other features together with feature in the pool
4. repeat the process

# Feature selection



$$[X_1, X_2, \dots, X_N] \xrightarrow{\text{feature selection}} [X_{i_1}, X_{i_2}, \dots, X_{i_M}]$$

## Time complexity

- forward / backward:  
 $O(n^2)$
- consider evaluation time:  
 $O(n^2 \times \text{eval})$
- might be unfeasible, if
  - you have a lot of features
  - it takes a long time to evaluate features
- stop after selecting / removing  $K$  features
- stop when objective function stops improving

# Feature extraction

Define a mapping function  $y = f(x)$  over all features:

$$[x_1, x_2, \dots, x_N] \xrightarrow{\text{feature extraction}} [y_1, y_2, \dots, y_M] = f([x_{i_1}, x_{i_2}, \dots, x_{i_m}])$$

🍇  $y = f(x)$  is, in general, non-linear and problem-dependent

🍇 e.g.  $\text{BMI} = \frac{\text{weight}}{\text{height}^2}$

We usually use linear projections  $y = Wx$

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} \xrightarrow{\text{linear feature extraction}} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1N} \\ w_{21} & w_{22} & \dots & w_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ w_{M1} & w_{M2} & & w_{MN} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix}$$

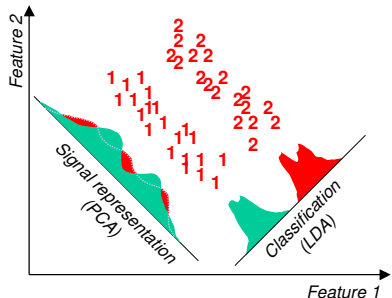
# Feature extraction

Two criteria to find the best mapping function  $y = f(x)$

Supervised

**Classification:** maximize separation among classes  
e.g. Linear Discriminant Analysis (LDA)

**Regression:** maximize correlation between projected data and target variable  
e.g. Partial Least Squares (PLS)



Signal representation

**Unsupervised:** retain as much data variance as possible  
Principal Component Analysis (PCA)  
Singular Value Decomposition (SVD)

# Dimensionality Reduction

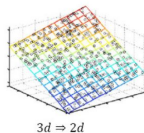
1. Why do we need dimensionality reduction?
2. Approaches for dimensionality reduction
3. Principal Component Analysis (PCA)
4. Multi-dimensional scaling (MDS)

# Principal Component Analysis

- ✿ widely used method for unsupervised, linear dimensionality reduction
- ✿ find a  $k$ -dimensional set of vectors (**components**) such that, if we project our data into this subspace, as much as possible of the variance is retained

**Input:** points scattered in multidimensional space

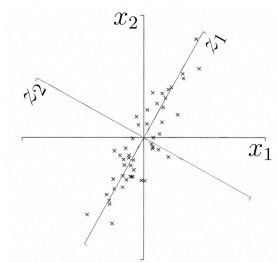
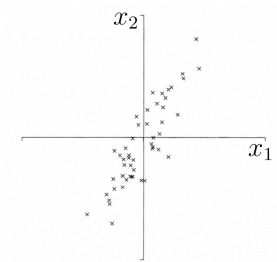
**Output:** a more compact representation of the data



## Applications:

- ✿ Data visualisation
- ✿ Data compression
- ✿ Signal processing
- ✿ Assist other learning algorithms

# PCA



Mathematical procedure:

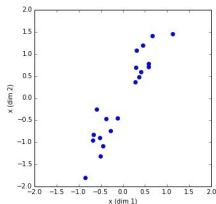
- transform a number of (possibly correlated) variables into a (smaller) number of uncorrelated variables
- the uncorrelated variables are called ***principal components***

Principal components (PC):

- first PC is the projection direction that maximizes the variance of the projected data
- second PC is the projection direction that is orthogonal to the first PC and maximizes variance of the projected data
- Spatial rearrangements may reveal relationships that were hidden in higher dimension space



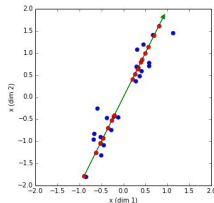
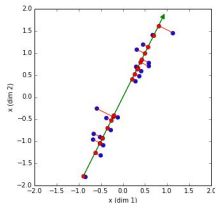
# PCA concept



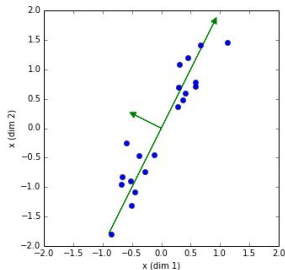
1. Find a line, such that when the data is projected onto that line, it has the **maximum variance**.

This is the same as **minimising** the orthogonal projection error.

- Projecting the data onto the PC(s) gives us an approximate representation in fewer dimensions.
- Equivalent to rotating the data onto new axes, then discarding some dimensions.



# PCA concept



2. Find another line, orthogonal to the first, that has maximum projected remaining variance.
3. Repeat until we have  $k$  orthogonal lines (PCs)

✿ The projected position of a point on the PCs is its coordinates in the  $k$ -dimensional space

PCs are ordered

- ✿ PC1 has the most variance
- ✿ PC2 has the second most, etc.

Principal components are uncorrelated

- ✿ Original data may contain correlated features
- ✿ Principal components are orthogonal

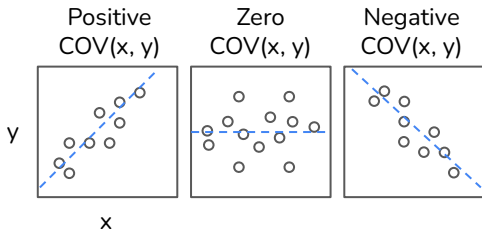
# PCA Main Terminology

## Covariance

- Measure of relationship between variables
- Positive = trend the same
- Zero = no relationship
- Negative = trend inversely

## Linear combination

- Multiplying each variable by a scalar
- Eg.  $1.2x + 3.4y + \dots$
- Used by PCA to define new axes which best explain the data in a smaller number of dims



# PCA Main Terminology

(For a given PC)

## Eigenvector

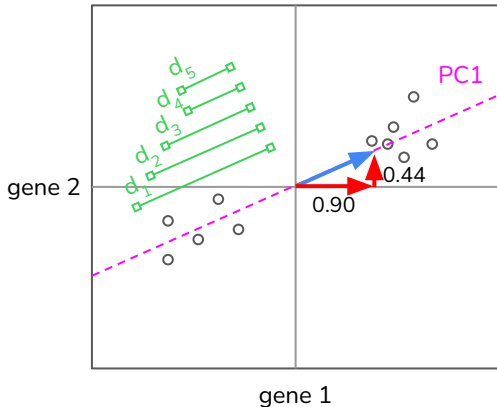
- Defines the line of best fit
- Unit vector (length = 1)
- Also known as singular vector

## Loading Scores

- Combination of each variable to create the PC
- $0.9 \times \text{gene1} + 0.44 \times \text{gene2}$
- A linear combination!

## Eigenvalue

- Measure of variance explained by this PC
- Average of SS(distances for PC)
- $\text{Eigenvalue} = (d_1^2 + d_2^2 + \dots) / n - 1$



# PCA algorithm

**Input:** Data matrix  $X$ , integer  $K$

**Output:** Projected matrix  $Z$

1. Calculate the covariance between each dimension of the data.
2. Find **eigenvectors** of the covariance matrix. This:
  - a. Finds basis vectors which are uncorrelated
  - b. Finds **eigenvalues**
    - i. Give a measure of variance along each eigenvector
3. Choose  $K$ : how many principal components to keep
4.  $Z \leftarrow$  Project data points onto  $K$  basis vectors

🧠 Principal Components chosen are the eigenvectors with largest eigenvalues

$$\text{Cov}_{xy} = \frac{\sum (x - \bar{x})(y - \bar{y})}{n - 1}$$

$$Ax = \lambda x$$

*For a matrix  $A$ , there is a vector  $x$  (eigenvector) such that the product of  $Ax$  is equal to a scalar multiplied by  $x$*

# PCA Output

## Output

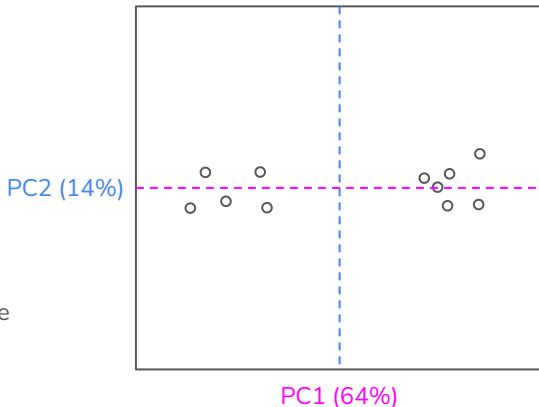
- Data is centered
- Axes are now the PCs

## Variation explained by each PC

- Can determine using Eigenvalues
- Eigenvalues: measure of variance
- Eg  $PC1=9$ ,  $PC2=2$ , ...,  $Sum=14$
- Variation explained by  $PC1 = 9/14 = 0.64$
- First 2 PCs explain 78% of the total variance
- The higher the better!

## How many components will be produced?

- $\text{Min}(\text{num variables}, \text{num samples})$
- Whichever is smaller



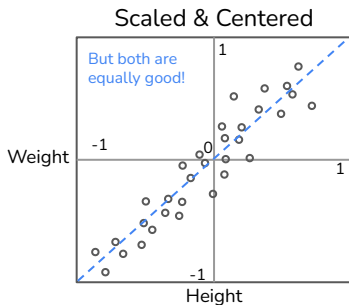
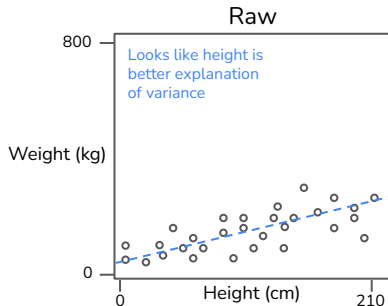
# PCA Preprocessing

## Scaling the data

- Variables should be on same scale
- RNAseq: Counts per million (CPM), Log CPM
- Log scale often useful for biological count data!

## Centering the data

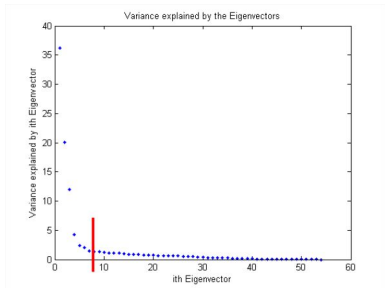
- PCA needs the data to be centered first
- Some libraries will do this for you, others won't
- Imagine centering about the "origin"
- Why? Conceptually, PCA finds best fit lines which go through the origin.



# Using PCA

How many components to choose?

- plot the proportion of variance explained by each eigenvector
  - eigenvalue divided by the sum of eigenvalues
- e.g. use eigenvectors that cover at least X % of the variance



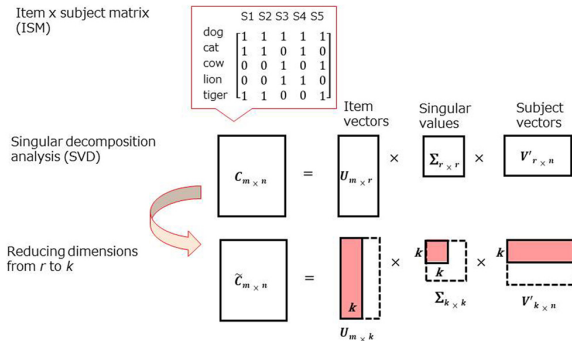
Time complexity ( $n$  dimensions and  $m$  observations)

- covariance matrix computation is  $O(n^2m)$
- eigenvalue decomposition is  $O(n^3)$
- total for PCA is  $O(n^2m + n^3)$

Alternatives?

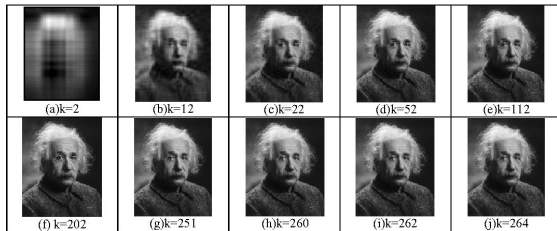


# Singular Value Decomposition (SVD)



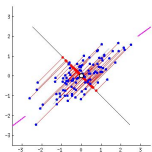
- data matrix  $C$
- factors  $U$ ,  $\Sigma$  and  $V^T$ 
  - $\Sigma$  (diagonal matrix) with singular values ( $\sigma$ )
  - singular values ( $\sigma$ ) are the square roots of the eigenvalues
- $U$  and  $V^T$  are orthogonal eigenvectors
- can use SVD to calculate PCA

# PCA and SVD



- 🍇 can interpret PCA as a SVD on the centered covariance matrix
- 🍇 SVD doesn't require you to calculate covariance matrix: more efficient

# Linear projections: PCA and SVD



$$\mathbf{X}_{n \times m} = \mathbf{U}_{n \times n} \mathbf{\Sigma}_{n \times m} \mathbf{V}^*_{m \times m}$$

**PCA:** Finds orthogonal vectors (components) to represent as much variance is as possible. Decomposing the covariance matrix,  $C = W\Sigma W^T$ . See [this animation](#).

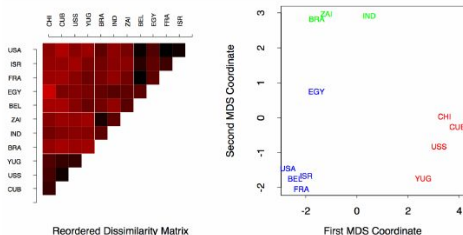
**SVD:** More efficient  
Can interpret PCA as a SVD on the centered covariance matrix

 [StatQuest: Principal Component Analysis](#) (20 minute video from Josh Starmer)

# Dimensionality Reduction

1. Why do we need dimensionality reduction?
2. Approaches for dimensionality reduction
3. Principal Component Analysis (PCA)
4. Multi-dimensional scaling (MDS)

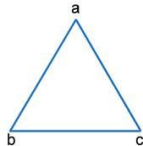
# Multi-Dimensional Scaling (MDS)



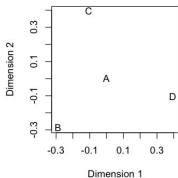
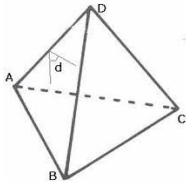
- 🍇 feature extraction
- 🍇 produces a lower-dimensional representation that ***preserves pairwise distances or dissimilarities***
- 🍇 for visualising data
- 🍇 sometimes called Principal Coordinates Analysis, but it's **not the same as PCA!**

# MDS

$$D(x_i, x_j) = \begin{matrix} & \begin{matrix} a & b & c \end{matrix} \\ \begin{matrix} a \\ b \\ c \end{matrix} & \begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix} \end{matrix}$$



$$D(x_i, x_j) = \begin{matrix} & \begin{matrix} a & b & c & d \end{matrix} \\ \begin{matrix} a \\ b \\ c \\ d \end{matrix} & \begin{pmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix} \end{matrix}$$



- start with a pairwise distance matrix or dissimilarity matrix
- we can represent three points that are equally-spaced in 3D **exactly in 2D**
- we can represent four points that are equally-spaced in 3D **exactly in 3D ...**
- ... but not in 2D**
- in general, we need  $N - 1$  dimensions to exactly represent pairwise distances between  $N$  samples

# Types of MDS

## Metric (distance-based) MDS

- 🍇 Minimise **stress**
- 🍇 Stress is the error in the distances (lack-of-fit measure)
- 🍇 Try to preserve distance between each vector  $x_i, x_j$

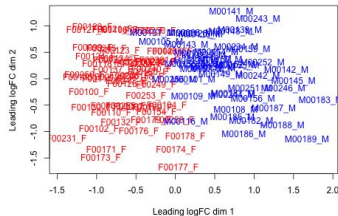
$$\text{Cost} = \sum_{i < j} (d_{ij} - \hat{d}_{ij})^2$$

- 🍇 actual cost/stress used may be more complex

## Non-metric MDS

- 🍇 Try to maintain original ranking of pairwise distances

# Interpreting low-dimensionality visualisation



- 🍇 clusters may represent real clusters
- 🍇 outliers may represent real outliers
- 🍇 the position of points is **not useful** in MDS
- 🍇 resulting dimensions are not ordered by importance/variance
- 🍇 can't reconstruct original data



# Thank you!

**Today:** Dimensionality Reduction I

**Next time:** Dimensionality Reduction II