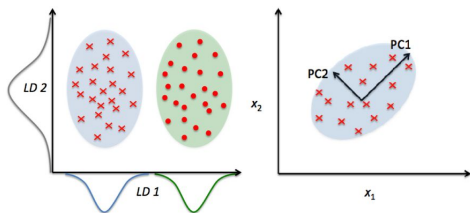# COMP90014

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
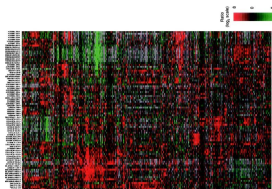- principal component analysis (PCA)
- singular value decomposition (SVD)

# 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
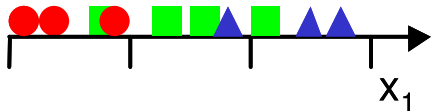- what is the dimensionality of this dataset?
- 27 000

thaliana      100

Arabidopsis
27,000
27,000

27,000

Wikimedia Commons; Zembutsu *et al.*, 2002. Cancer Res 62 (2): 518–527.
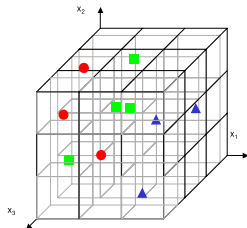
# High dimensionality has many costs


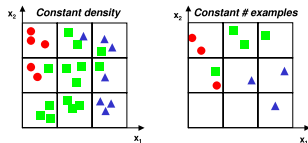
- data interpretation and visualisation
- redundant and irrelevant features degrade performance of ML algorithms
- computational cost may become intractable

# 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**

Ricardo Gutierrez-Osuna, Wright State University

# Goal of dimensionality reduction



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

Австралиец via [Wikimedia Commons](Wikimedia Commons)

# 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

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

- 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

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

Feature extraction: creating new features by combining existing ones

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

# Feature selection

Forward Greedy Feature Selection

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

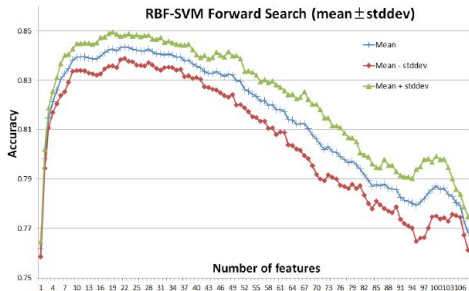Forward Greedy Feature Selection

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

这里是算法如何运作的:

1. **初始化**: 我们开始时没有任何特征被选中,所以特征集 $FS = \emptyset$,所有特征 $F = \{f_1, f_2, f_3, f_4\}$。

2. **第一轮迭代**:
   - 评估每个特征对模型性能的贡献。假设我们使用模型准确度作为评分标准。
   - 假设评分如下:$\text{score}(f_1) = 0.6, \text{score}(f_2) = 0.8, \text{score}(f_3) = 0.5, \text{score}(f_4) = 0.7$。
   - $f_2$ 有最高的评分,所以它被添加到特征集 $FS$,现在 $FS = \{f_2\}$,并且从 $F$ 中移除 $f_2$。

3. **第二轮迭代**:
   - 现在我们只评估剩下的特征 $f_1, f_3, f_4$ 与已选择的 $f_2$ 结合时的性能。
   - 假设新的评分如下:$\text{score}(f_1 + f_2) = 0.85, \text{score}(f_3 + f_2) = 0.82, \text{score}(f_4 + f_2) = 0.9$。
   - $f_4$ 与 $f_2$ 结合有最高的评分,所以 $f_4$ 被添加到 $FS$,现在 $FS = \{f_2, f_4\}$。

4. **第三轮迭代**:
   - 我们继续评估剩下的特征 $f_1, f_3$ 与已选择的特征 $f_2, f_4$ 结合时的性能。
   - 假设评分如下:$\text{score}(f_1 + f_2 + f_4) = 0.88, \text{score}(f_3 + f_2 + f_4) = 0.87$。
   - $f_1$ 有更高的评分,所以现在 $FS = \{f_2, f_4, f_1\}$。

5. **最后**:
   - 最后一个特征 $f_3$ 将被评估,但由于我们已经达到一个较高的性能得分,如果 $f_3$ 不能显著提高性能,它可能不会被添加到 $FS$ 中。
   - 假设 $\text{score}(f_1 + f_2 + f_4 + f_3) = 0.86$,比没有 $f_3$ 的得分还低,因此我们决定不添加 $f_3$。

最终,我们选择的特征子集是 $FS = \{f_2, f_4, f_1\}$。这个子集是基于贪婪算法迭代地评分和选择过程中得到的,它应该比其他可能的特征组合有更好的分类性能。

# Feature selection



RBF-SVM Forward Search (mean±stddev)

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

Nguyen & Allebach, 2017. 10.2352/ISSN.2470-1173.2017.12.IQSP-238

## 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

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

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

- 🍇 $y = f(x)$ is, in general, non-linear and problem-dependent
- 🍇 *e.g.* BMI $= \dfrac{\text{weight}}{\text{height}^2}$        BMI

We usually use linear projections $y = Wx$

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

# 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)
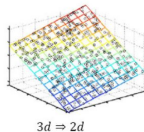Singular Value Decomposition (SVD)

# 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



$3d \Rightarrow 2d$

Applications:

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

# PCA

$x_2$
$x_1$

$x_2$
$z_1$
$z_2$
$x_1$

Mario Guarracino, ICAR-CNR

Principal Components PCs

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



PCA

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.
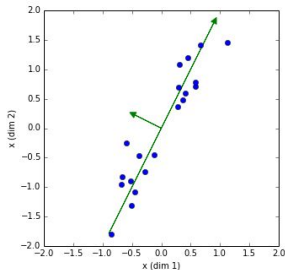- Equivalent to rotating the data onto new axes, then discarding some dimensions.

# 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
   PC1                                        PC2

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



k            k

k

PCA

# 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 + .....
- Used by PCA to define new axes which best explain the data in a smaller number of dims



Positive COV(x, y)   Zero COV(x, y)   Negative COV(x, y)

y

x

1. 2x+3. 4y
x
y
  PCA

# PCA Main Terminology

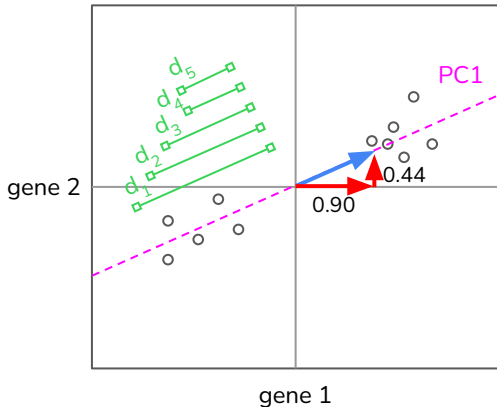(For a given PC)

**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×gene1 + 0.44×gene2
- A linear combination!

**Eigenvalue**

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

# 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* ← 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*

a.
b.
i .
   K
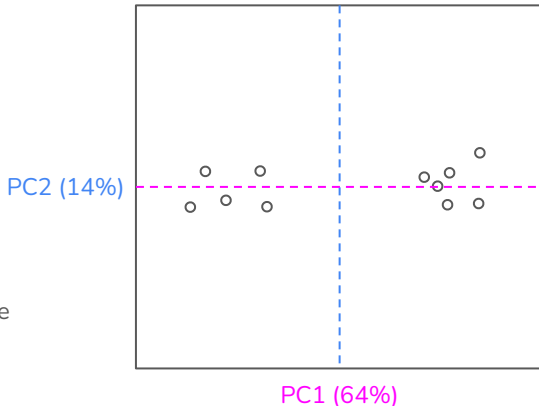Z              K

# 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?

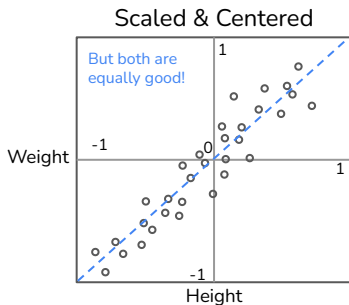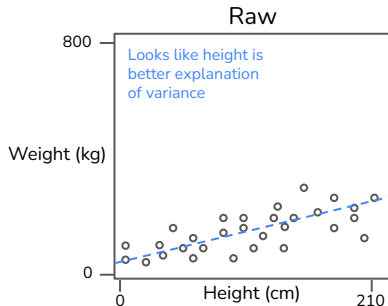- Min(num variables, num samples)
- Whichever is smaller



PC2 (14%)

PC1 (64%)

# 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"
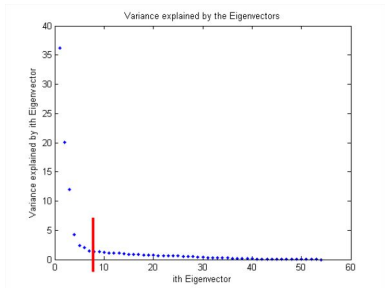- Why? Conceptually, PCA finds best fit lines which go through the origin.

### Raw



Looks like height is better explanation of variance

### Scaled & Centered



But both are equally good!

## Using PCA



Variance explained by the Eigenvectors

How many components to choose?

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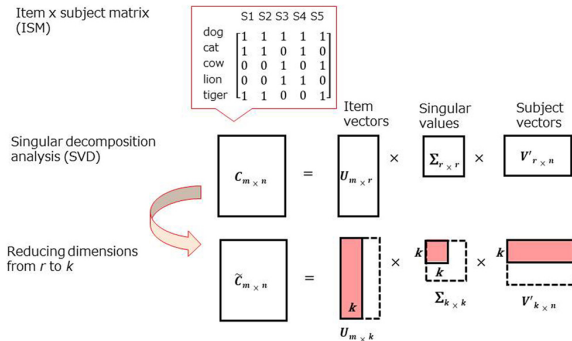- 🐾 *e.g.* use eigenvectors that cover at least *X* % of the variance

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

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

Alternatives?

# Singular Value Decomposition (SVD)



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

$$C = U * \underset{m \quad n}{\quad} * V\text{\textasciicircum}T$$

PCA                               SVD
PCA

SVD



(a)k=2   (b)k=12   (c)k=22   (d)k=52   (e)k=112
(f) k=202   (g)k=251   (h)k=260   (i)k=262   (j)k=264
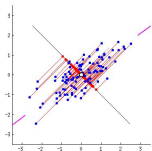
- can interpret PCA as a SVD on the centered covariance matrix
- SVD doesn't require you to calculate covariance matrix: more efficient

SVD

# Linear projections: PCA and SVD

PCA          SVD

PCA

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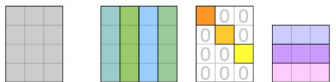See this animation.

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

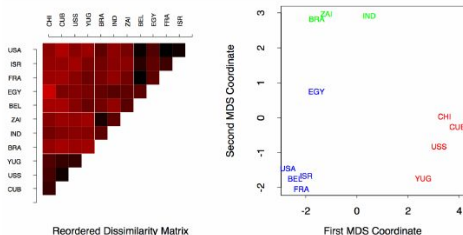🫐 StatQuest: Principal Component Analysis (20 minute video from Josh Starmer)

SVD

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

amoeba via stackexchange; Cmglee via Wikimedia Commons

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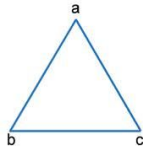- produces a lower-dimensional representation that *preserves pairwise distances or dissimilarities*  MDS
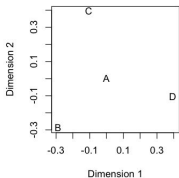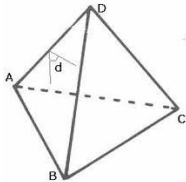- for visualising data
- sometimes called Principal Coordinates Analysis, but it's **not the same as PCA**!

# MDS

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

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

- 🫐 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

N

N−1

MDS          D

# Types of MDS

MDS " " stress "

lack-of-fit measure

xi xj

i, j

## Metric (distance-based) MDS

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

$$\text{Cost} = \sum_{i<j} \left( d_{ij} - \hat{d}_{ij} \right)^2$$
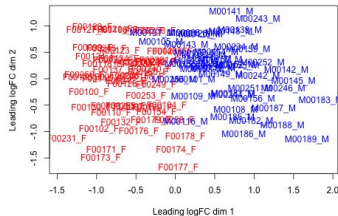
- actual cost/stress used may be more complex

## Non-metric MDS

- Try to maintain original ranking of pairwise distances

MDS

- 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

MDS            MDS

            MDS

        /        MDS

    MDS

# Thank you!

**Today:** Dimensionality Reduction I

**Next time:** Dimensionality Reduction II