# Lecture 10

#### Introduction to bioinformatics (MVE510)

Autumn, 2020

Additional reading: Clustering. Naomi Altman and Martin Krzywinski. Nature Methods, 14(6) 2017 and Principal component analysis. Jake Lever, Martin Krzywinski, and Naomi Altman. Nature Methods, 14(7) 2017.

1

## Todays agenda

- Unsupervised methods and explorative data analysis
- Clustering
  - Hierarchical clustering
  - K-means clustering
- Techniques for dimension reduction
  - Principal component analysis

## Supervised and unsupervised methods

- Data analysis methods can be divided into <u>supervised</u> and <u>unsupervised</u> methods
- Supervised methods rely on *metadata*, i.e. data that provides information about the data. This can be any form of medical, biological or technical information about the samples.
- Linear models are a supervised method. Here we specifically look for differences associated with a covariate (i.e. metadata). This includes, for example, identification of differentially expressed genes between different groups of patients.

#### Unsupervised methods

- Unsupervised methods does not utilize metadata, i.e. <u>no additional</u> information about the samples are used.
- Unsupervised methods focus on the identification of <u>patterns in the</u> <u>data</u>. In bioinformatics this typically means patterns between samples or genes.
- Unsupervised methods are explorative and does <u>not rely on any</u> <u>specific hypotheses</u>. This means that we can identify almost any form of patterns.
- Another common application of unsupervised methods is dimension reduction, where high-dimensional data is visualized in e.g. two dimensions.

#### Unsupervised methods

- In this course we will look at two main methods for unsupervised data analysis
- Clustering
  - Aims find groups of samples or genes ('clusters') that are 'similar'.
  - Two important methods: Hierarchical clustering and K-means clustering.
- Principal component analysis (PCA)
  - Used for <u>exploration and visualization</u> of high-dimensional data
  - Identifies the dimensions with highest variability. These are used to represent the data at a lower dimension (e.g. in two dimensions).

- Hierarchical clustering builds a hierarchy describing the similarity between clusters
- Agglomerative: All samples starts as individual clusters which are merged together (aggregated) according to
  - A <u>distance measure</u> describing the separation between data points
  - A <u>linkage criterion</u> measuring the distance between clusters
- Divisive: All samples starts in a single large cluster which is then successively divided
- The number of clusters are then defined based on a minimum similarity threshold

#### Hierarchical clustering - distances

All clustering techniques are based on a distance metric that measures the similarity between two data points.

#### A few notations

We will use  $x = (x_1, ..., x_p)$  and  $y = (y_1, ..., y_p)$  to denote two data points defined in  $\mathbb{R}^p$ . Let d(x, y) denote the distance between x and y.

#### Euclidean distance

$$d(x, y) = \sqrt{\sum_{i=1}^{n} (x_i - y_i)^2}$$



#### Manhattan distance

$$d(x, y) = \sum_{i=1}^{n} |x_i - y_i|$$



Maximum distance (
$$L^{\infty}$$
)  
 $d(x, y) = \max_{i} |x_{i} - y_{i}|$ 



#### Correlation distance

$$\operatorname{corr}(x, y) = \frac{\operatorname{cov}(x, y)}{\sqrt{\operatorname{var}(x)\operatorname{var}(y)}}$$
 and  $d(x, y) = \frac{1 - \operatorname{corr}(x, y)}{2}$ 

where  $cov(x, y) = \frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})$ ,  $\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$  and var(x) = cov(x, x).

9

## The distance matrix

- Contains all <u>pair-wise</u> distances between the data points.
- The distance matrix is symmetric since d(x, y) = d(y, x).
- This matrix can be really big since it contains  $n^2$  elements (only n(n 1)/2 elements are actually needed).
- Some algorithms works instead with similarity matrices.

	1	2	3	4	5	6	7	8	9	10
1	0.00	5.74	5.04	4.60	4.11	4.63	4.13	5.70	4.59	3.33
2	5.74	0.00	3.49	6.29	5.96	5.30	6.62	7.22	5.68	4.60
3	5.04	3.49	0.00	5.68	5.70	4.26	4.62	6.25	5.94	4.94
4	4.60	6.29	5.68	0.00	3.73	4.98	5.95	6.44	3.71	3.18
5	4.11	5.96	5.70	3.73	0.00	5.45	5.50	4.97	4.33	3.76
6	• • •									

#### Distances – which distance should I choose?

The choice of distance depends on the question. What similarities are you interested in?



Euclidean distance d(x, y) = 59.4 d(x, z) = 360.0Correlation distance d(x, y) = 0.666d(x, z) = 0.097

#### Hierarchical clustering - distances

Some more notation

We will use  $x = (x_1, ..., x_p)$  and  $y = (y_1, ..., y_p)$  to denote two data points defined in  $\mathbb{R}^p$ . Let d(x, y) denote the distance between x and y.

Let X and Y denote two clusters with k and l members each, i.e.  $X = \{x^{(1)}, x^{(2)}, ..., x^{(k)}\}$  and  $Y = \{y^{(1)}, y^{(2)}, ..., y^{(l)}\}$ . Let D(X, Y) denote the distance between X and Y.

#### Algorithm: Agglomerative hierarchical clustering

- 1. Calculate a distance matrix containing the pair-wise distance between all data point
- 2. Let each data point be a cluster
- 3. Identify the most similar clusters X and Y according to a *linkage criterion*.
- 4. Merge X and Y into a new cluster. Update the distance matrix.
- 5. Goto 3 if the number of clusters are more than 2.

#### Linkage criteria

The linkage criterion calculates the similarity between clusters.

Single linkage

$$D(X,Y) = \min_{x \in X, y \in Y} d(x,y)$$

Complete linkage

 $D(X,Y) = \max_{x \in X, y \in Y} d(x,y)$ 



#### Linkage criteria

Average linkage

$$D(X,Y) = \frac{1}{kl} \sum_{\substack{\text{all pairs } (x,y), \\ x \in X, y \in Y}} d(x,y)$$



Centroid linkage

 $D(X,Y) = d(x_c, y_c)$ 



#### Algorithm: Agglomerative hierarchical clustering

- 1. Calculate a distance matrix containing the pair-wise distance between all data point
- 2. Let each data point be a cluster
- 3. Identify the most similar clusters X and Y according to a *linkage criterion*.
- 4. Merge X and Y into a new cluster. Update the distance matrix.
- 5. Goto 3 if the number of clusters are more than 2.













Example: Agglomerative clustering, Euclidean distance, complete linkage.



	1	23	456
1	0.00	5.39	9.43
23	5.39	0.00	6.71
456	9.43	6.71	0.00

23







	123	456		
123	0.00	9.43		
456	9.43	0.00		



#### Hierarchical clustering - dendrogram

**Cluster Dendrogram** 



# Hierarchical clustering in R can be done using the **dist** and **hclust** commands.

28

## Clustering of gene expression data (RNA-seq)

**Data**: Gene expression of 45,015 transcripts in 33 patients and 6 cell lines. Among the patients, 23 were from less aggressive group and 10 from more aggressive group. For each patient, a biopsy was sampled, the RNA extracted, sequenced using standard protocols, preprocessed and transformed and normalized so that the resulting data is approximately normal distributed. For each patient, the age and gender was also recorded.

Aim: Use clustering so investigate patterns between the patients.

Clusters were calculated in R using dist and hclust.



Cluster Dendrogram

Euclidean distance, complete linkage



**Cluster Dendrogram** 

Euclidean distance, single linkage



#### Cluster Dendrogram

йй

Manhattan distance, complete linkage

32

- A method for clustering *n* observations in to *s* clusters. Note that *s* is pre-specified!
- We want assign each sample to one of the *s* clusters such that the <u>variance within each cluster is minimized</u>. This means that we want to minimize



 K-means clustering implicitly assumes a Euclidean distance between points.

The variance within a cluster can be calculated by

$$Var(c_j) = \sum_{x \in c_j} \sqrt{\sum_{i=1}^n (x_i - \mu_i)^2}$$



where  $\mu_i$  is the *mean point* of  $c_j$ . Note that minimizing  $Var(c_j)$  is equivalent to minimizing the distance between all pairs of data points in the cluster, i.e.

$$\frac{1}{2} \sum_{x,y \in c_j} \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$



#### The naïve algorithm

Start with an initial set of k randomly selected cluster means

$$\mu = (\mu_1, \mu_2, ..., \mu_k).$$

Step 1: Assignment

Assign each data point to the closest cluster based on the Euclidean distance to the cluster mean.

Step 2: Update

Recalculate the cluster mean values.

Repeat these steps until no more assignments are done.



36



Since there are no changes in assignment, the algorithm has converged.

A drawback with K-means clustering is that it is often hard to decide the number of clusters (s).

- Method used for <u>exploration and visualization</u> of high-dimensional data
- The main principle: Identify the dimensions of the data that explains as <u>much variability as possible</u>. These dimensions are defined by principal components.
- The data can then be projected on a set of the to principal components (often two) which results in a reduction of the dimensionality
- PCA is often used to identify patterns in the data. However, PCA does not assign each data point to a cluster.







41

#### Some notation

We will, as before, use  $x = (x_1, ..., x_p)$  to denote a data point. Note that we will here regard x as a random variable.

Let  $y_1, y_2, \dots y_n$  denote *n* linear combinations of the coordinates of *x*, i.e.

$$y_{1} = a_{11}x_{1} + a_{12}x_{2} + \dots + a_{1n}x_{n} = a_{1}^{T}x$$
$$y_{2} = a_{21}x_{1} + a_{22}x_{2} + \dots + a_{1n}x_{n} = a_{2}^{T}x$$
$$\dots$$
$$y_{n} = a_{n1}x_{1} + a_{n2}x_{2} + \dots + a_{nn}x_{n} = a_{n}^{T}x$$

**—** 

The principal components will be identify be selecting  $\{a_{ij}\}$  according to the following rules

The first principal component is the linear combination  $y_1 = a_1^T x$ that maximizes the variability  $Var(y_1)$  under the condition that  $a_1^T a_1 = 1$ .

The <u>second principal component</u> is the linear combination  $y_2 = a_2^T x$ that maximizes the variability  $Var(y_2)$  under the condition that  $a_2^T a_2 = 1$ and is independent of  $y_1$ .

The *nth* principal component is the linear combination  $y_n = a_n^T x$ that maximizes the variability  $Var(y_2)$  under the condition that  $a_n^T a_n = 1$ and is independent of  $y_1, y_2, ..., y_{n-1}$ .

The principal components corresponds to the eigenvectors of the covariance matrix. They are therefore often calculated by singular value decomposition (SVD).

A principal component analysis results

- 1. The principal components  $y_1, ..., y_n$  constitutes a new orthogonal basis for the data. The number of principal components is the same as the dimension of the data (or, if lower, the number of samples).
- 2. The rotation matrix, which is formed from vectors  $a_1, a_2, \ldots, a_n$ . This matrix describe how data can be rotated to the new basis. The individual elements are sometimes called <u>loadings</u>. Note that the rotation matrix is not unique it can always be multiplied by -1.
- 3. The estimated variance for each principal component. These are <u>always</u> in decreasing order.

# In R, principal component analysis can be done using the **prcomp** command.

> res.pca<-prcomp(data.matrix, center=FALSE)</pre>

PC1 PC2 PC3 [1,] -1.084225 -1.4550113 -0.3345473 [2,] 1.559554 -1.6491774 0.1780445 [3,] -1.815806 -0.5476472 0.3526779

> res.pca\$sdev
[1] 1.8580773 1.6026168 0.3660624

# 3x5 matrix

# the rotation matrix defined by the loadings

# the transformed data calculated from the
# rotation matrix

# the standard deviation of the three principal
# components

## PCA of gene expression data (RNA-seq)

**Data**: Gene expression of 45,015 transcripts in 33 patients and 6 cell lines. Among the patients, 23 were from less aggressive group and 10 from more aggressive group. For each patient, a biopsy was sampled, the RNA extracted, sequenced using standard protocols, preprocessed and transformed and normalized so that the resulting data is approximately normal distributed. For each patient, the age and gender was also recorded.

Aim: Use PCA so investigate patterns between the patients.

The PCA was performed in R using **prcomp.** The principal components were plotted using the **plot** function.







Scree plot

50

#### Clustering and PCA – some remarks

- Clustering and PCA are explorative methods. They can be used to visualize data and formulate hypothesis but they does not result in any statistical test (and thus no measure of the significance).
- Many methods automatically center and scale the data. You may need to manually turn such features off.
- There is a wide range of different tools for clustering and PCA. The command **heatmap** can for example cluster and display the results as a heatmap and a dendrogram.

#### Summary – Unsupervised methods

- Unsupervised methods can be used to explore and visualize highdimensional data.
- Clustering is a technique used to find groups of data points ('clusters'). Common clustering methods includes hierarchical clustering and k-means clustering. A distance metric and linkage criterion is needed to perform hierarchical clustering.
- Principal component analysis can be used to reduce the dimension of the data. The principal components are identified based on the dimension on where the variability of the data is as high as possible.