

STATISTICS & ML WITH R

Intro to Machine Learning

2024

M. Chiara Mimmi & Luisa M. Mimmi

WORKSHOP SCHEDULE

- Modules

- 1. Intro to R and data analysis
- 2. Statistical inference & hypothesis testing
- 3. Modeling correlation and regression
- 4 Mapping causal & predictive approaches
- 5. Machine Learning
- 6. Extra topics:
 - MetaboAnalyst;
 - Power Analysis

- Each day will include:

- Frontal class (MORNING)
- Practical training with R about the topics discussed in the morning. (AFTERNOON)

MODULE 5 – LECTURE OUTLINE

- Intro to Machine Learning (ML)
- Classification of ML algorithms
- *Supervised* ML examples
 - Logistic Regression
 - Decision trees
- *Unsupervised* ML examples
 - K-means clustering
 - PCA
 - PLS-DA

Different goals of statistical modeling

1. **ASSOCIATION/CORRELATION** → observational studies
 - aimed at **summarizing or representing the data structure**, without an underlying causal theory
 - may help **form hypotheses** for explanatory and predictive modeling

2. **CAUSAL EXPLANATION** → experimental studies
 - aimed at **testing “explanatory connection”** between treatment and outcome variables
 - prevalent in “**causal theory-heavy**” fields (economics, **psychology**, environmental science, etc.)
 - **Note:**
 - ✓ The **same modeling approach** (e.g., fitting a regression model) can be used for **different goals**
 - ✓ While they shouldn't be confused, **explanatory power** and **predictive accuracy** are complementary goals: e.g., in bioinformatics (which has little theory and abundance of data), predictive models are pivotal in generating avenues for causal theory.

3. **EMPIRICAL PREDICTION** → algorithmic machine learning and data-mining modeling

Different goals of statistical modeling (*today!*)

1. **ASSOCIATION/CORRELATION** → observational studies
2. **CAUSAL EXPLANATION** → experimental studies
3. **EMPIRICAL PREDICTION** → algorithmic machine learning and data-mining modeling
 - aimed at **predicting new or future observations** (without necessarily explaining how)
 - relies on **big data**
 - prevalent in fields like natural language processing, **bioinformatics**, etc.. In **epidemiology**, there is more of a mix causal explanation & empirical prediction
 - **Notes:**
 - ✓ “Prediction” does not necessarily refer to future events, but rather to *future* datasets that were previously unseen to the algorithm

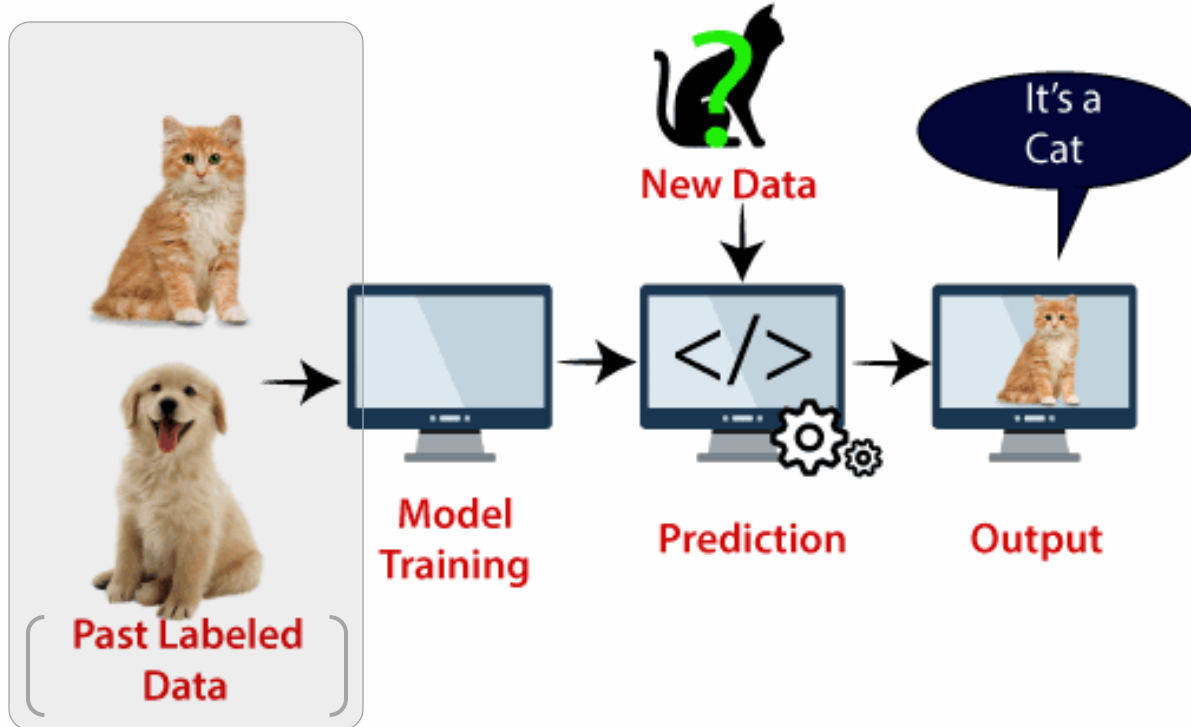
MACHINE LEARNING

Defining Machine Learning (ML)



“At its core, Machine Learning is just a “thing-labeler”, taking something and telling you what label it should get.”

(Cassie Kozyrkov)



Source: Image from <https://entri.app/blog/what-is-svm-algorithm-in-machine-learning/>

Defining Machine Learning (ML)

- **Machine Learning** is a broad and highly active research field. (In the life sciences, “*precision medicine*” is an application of machine learning to biomedical data)
- The **general idea** is to **predict** or **discover outcomes from measured predictors**, in problems like:
 - *Can we discover new types of cancer from gene expression profiles?*
 - *Can we predict drug response from a series of genotypes?*
 - *How do we classify a set of images/spectrometry outputs, etc.*
 - *Given various clinical parameters, how can we use them to predict heart attacks?*
- The **ML is a data-driven (inductive) approach**, where a machine **learns** the rules/patterns from a set of **training data** and (then) **validates** findings on a set of **testing data**
- In contrast with inferential statistics, **ML doesn't worry about assumptions on parameters** (probability distribution, error, correlation, etc.), **nor the causal nexus** between specific predictor(s) and response, **nor the data collection strategy**
- In contrast with standard statistics, **in ML the rules are not necessarily specified...** hence ML = a subfield of AI

Stylized comparison between statistics and machine-learning

	Standard (causal inference) Statistics	Machine Learning
Typical Goal	Explanation, uncovering causal relationships	Predicting an outcome as accurately as possible
Typical Task	Research based on a theory to identify the <u>causal effect</u> (better: pre-register your hypothesized model).	Try out and tune many different algorithms in order to <u>maximize predictive accuracy</u> in new and unseen test datasets.
Data generating process	Designed ex-ante based on study goal (e.g. randomized control trial, or observational study with statistical control variables)	Useful but not strictly necessary, and often not available
Parameters of interest:	Causal effect size and statistical significance, p-value of <u>treatment X</u> for outcome Y	Model's accuracy (%), precision/recall, sensitivity/specificity, in <u>predicting Y</u>
Dataset	Use ALL AVAILABLE DATA to calculate effect of interest (it was designed to be representative of a population).	It is critical to SPLIT THE DATA (usually 75% for training and 25% for testing the algorithms) leaving aside a sub-sample to test the model with unseen new data

Source: Adapted from <https://forloopsandpiepkicks.wordpress.com/2022/02/10/beginners-guide-to-machine-learning-in-r-with-step-by-step-tutorial/>

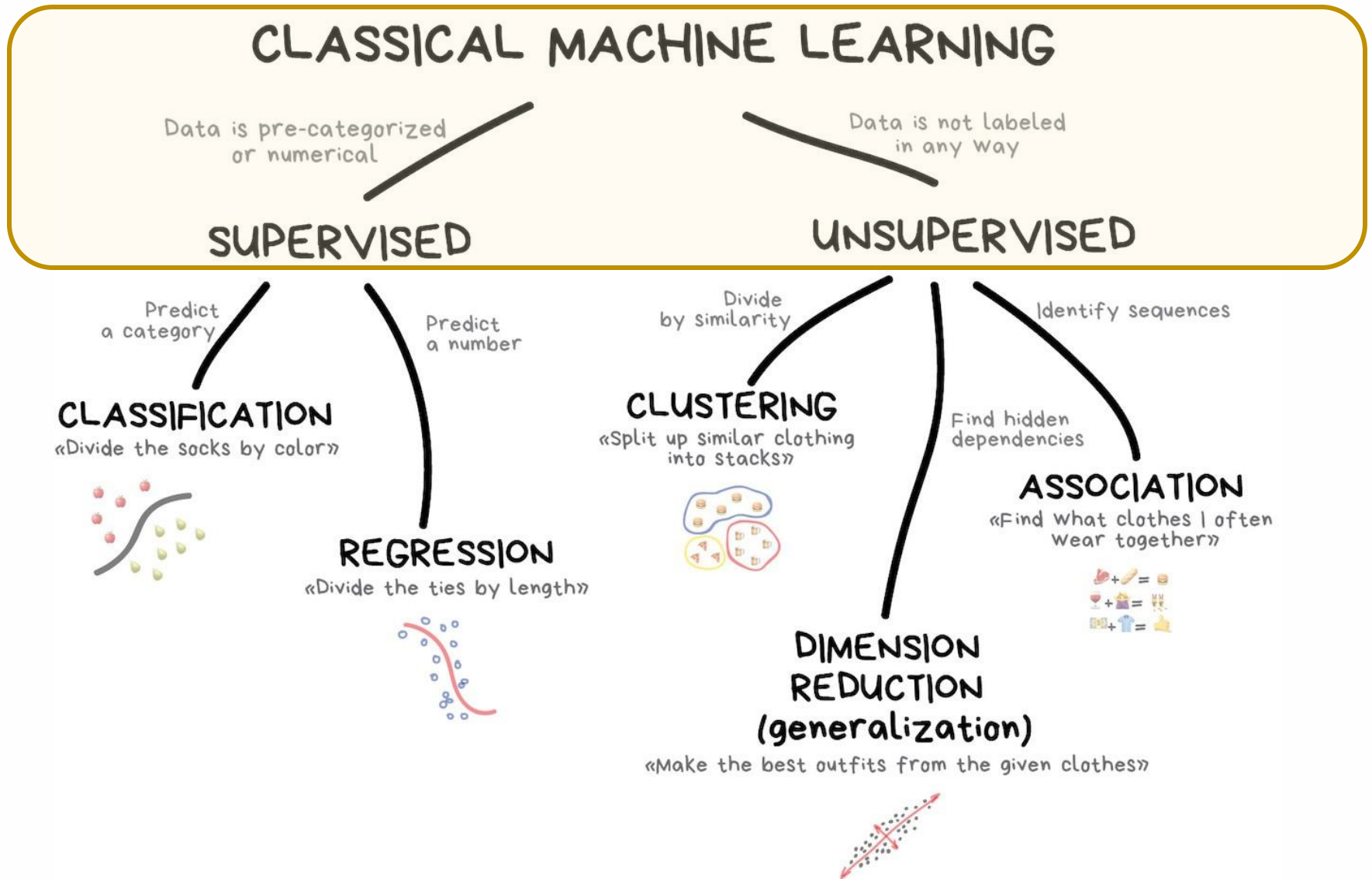
MODULE 5 – LECTURE OUTLINE

- Intro to Machine Learning (ML)
- Classification of ML algorithms
- *Supervised* ML examples
 - Logistic Regression
 - Decision trees
- *Unsupervised* ML examples
 - K-means clustering
 - PCA
 - PLS-DA

Supervised or Unsupervised ML algorithms?

....another conceptual framework

A fundamental distinction: supervised and unsupervised ML



Source: Image from https://vas3k.com/blog/machine_learning/index.html

A fundamental distinction: supervised and unsupervised ML

- ML includes many different algorithms that can be used for understanding data. These algorithms can be classified as:
 - **Supervised Learning Algorithms:**
 - building a model to estimate or predict an output based on one or more inputs
 - **Regression:** Modeling a relationship, the typical output variable is continuous (e.g. weight, height, time, etc.) or dichotomous.
 - **Classification:** Splits objects based on one of the attributes known beforehand. The the typical output variable is categorical (e.g. male or female, pass or fail, benign or malignant, etc.)
 - **Unsupervised Learning Algorithms:**
 - finding structure and relationships among inputs. There is no “supervising” output
 - **Clustering:** Finding “clusters” of observations in a dataset that are similar to each other (*based on unknown features*).
 - **Association:** Finding “rules” that can be used to draw associations. For example, if a patient has a high biomarker X, he will have a low biomarker Y.
 - **Dimension reduction:** Assembling specific features into more high-level ones (e.g. PCA)

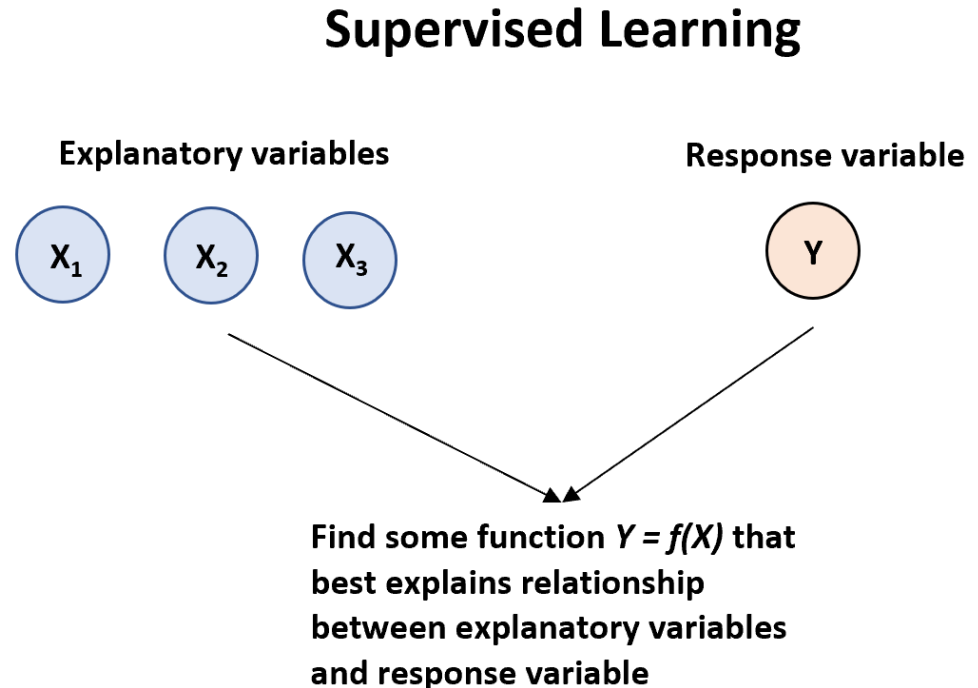
MODULE 5 – LECTURE OUTLINE

- Intro to Machine Learning (ML)
- Classification of ML algorithms
- *Supervised* ML examples
 - Logistic Regression
 - Decision trees
- *Unsupervised* ML examples
 - K-means clustering
 - PCA
 - PLS-DA

Supervised ML algorithms

Supervised Learning Algorithms mechanics

- A supervised learning algorithm can be used when we have **one or more explanatory variables** ($X_1, X_2, X_3, \dots, X_p$) and a **response variable** (Y) and we would like to find some function that describes the relationship between the explanatory variables and the response variable:
- $Y = f(X) + \varepsilon$
- where
 - $f()$ represents **systematic information that X provides about Y** and where
 - ε is a random error term independent of X with a mean of zero.



Source: <https://www.statology.org/supervised-vs-unsupervised-learning/>

Supervised Learning Algorithms **purpose**

There are two main reasons to use supervised learning algorithms:

- 1. Prediction:** We often use a set of explanatory variables to predict the value of some response variable (e.g. using square footage and number of bedrooms to predict home price)
 - 2. Inference:** We may be interested in understanding the way that a response variable is affected as the value of the explanatory variables change (e.g. how much does home price increase, on average, when the number of bedrooms increases by one?)
- Depending on whether our goal is inference or prediction (or a mix of both), we may use different methods for estimating the function f . For example, linear models offer easier interpretation but non-linear models that are difficult to interpret may offer more accurate prediction.*

Supervised Learning: commonly used algorithms

Most commonly used supervised learning algorithms:

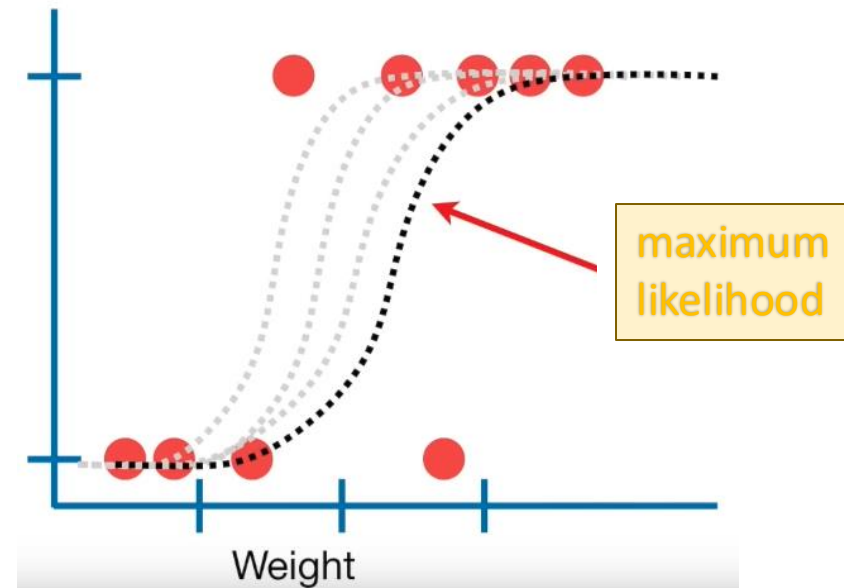
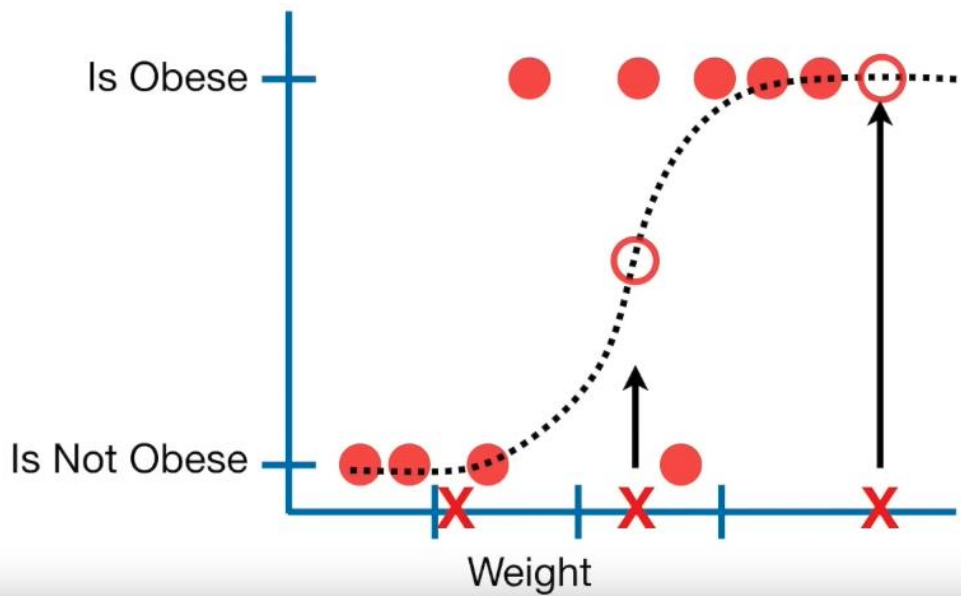
- Linear regression
- Logistic regression
- Linear discriminant analysis
- Quadratic discriminant analysis
- Decision trees
- Naive bayes
- Support vector machines
- Neural networks

Logistic Regression for classification

An example of **supervised** ML algorithm

Purpose of logistic regression

- **Logistic regression** tells the **PROBABILITY** of some phenomenon (e.g. being obese) and is normally used for **binary classification**:
 - e.g. if the probability a mouse is obese > 50%, we will classify it as obese
- Unlike linear regression (that uses **least square**) the logistic regression line is fit using the **maximum likelihood** criterion

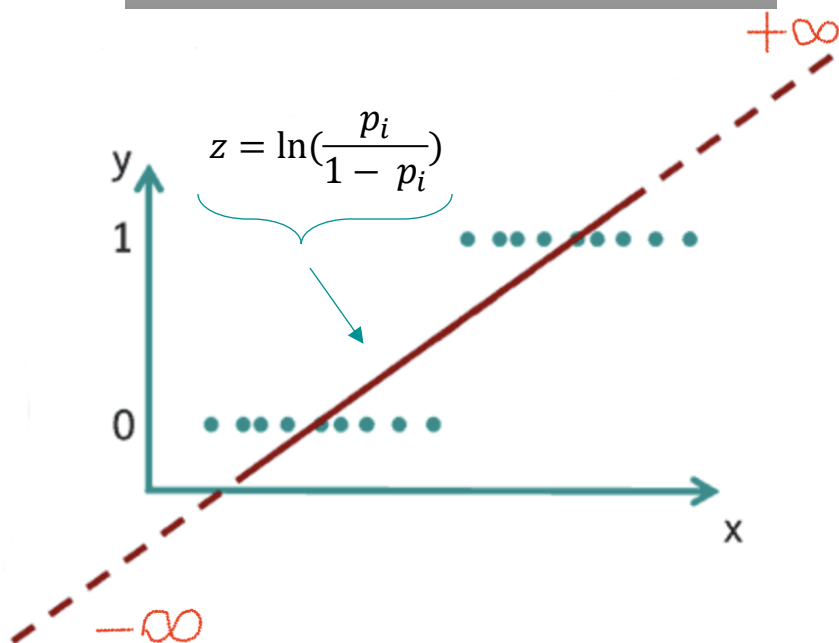


Source: <https://www.youtube.com/watch?v=yYKR4sgzI8&t=204s>

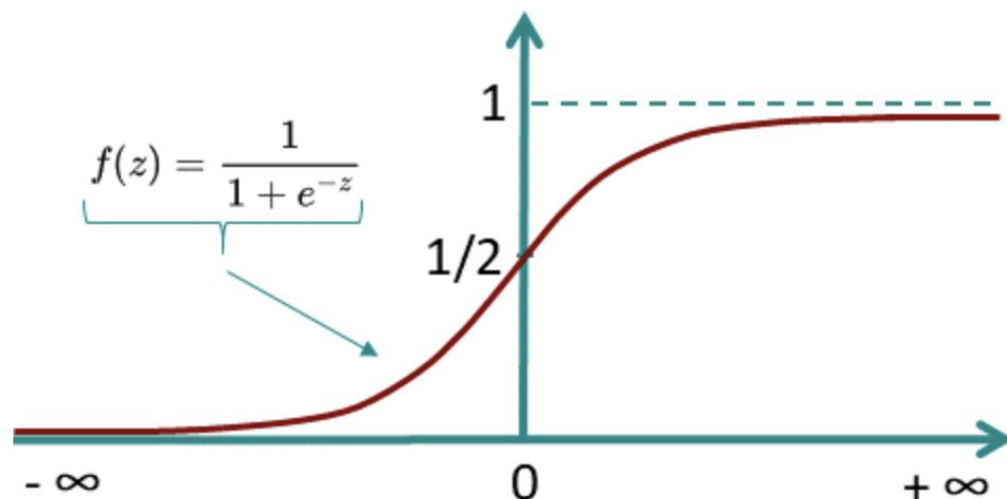
The “*logic*” of logistic regression

- **Logistic regression** is similar to **linear regression**, but... now the **response variable y** can only take **values 0 and 1**
- So we need *some* function to restrict the value range for the prediction between 0 and 1: the **logistic function** (or **sigmoid** function).

Logit = Log(odds) function



Logistic function (*sigmoid*)



Source: <https://datatab.net/tutorial/logistic-regression>

The logit and inverse-logit functions

Given $p_i = P(y = 1 | x_1 + x_2 + \dots + x_k)$, i.e. is the probability that $y = 1$ given a vector of independent variables x_k , linear regression is not appropriate. Instead, we can use **logistic regression**:

- 1) The **logit** function classifies the outcomes as *log(odds)*, but these has range $[-\infty, +\infty]$

logit $logit(p_i) = \ln\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \dots + \beta_k x_{k,i}$
(or “link” function)

- 2) To obtain a *meaningful* probability value for $y = 1$, the **sigmoid** function maps values back to the range $[0, 1]$:

inverse-logit $p_i = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \dots + \beta_k x_{k,i})}}$
(aka **logistic** function)

- 3) Finally, the result from the **sigmoid** function ($[0,1]$) is then passed through a decision rule (i.e. a **threshold**) to divide the outcome into classes as required

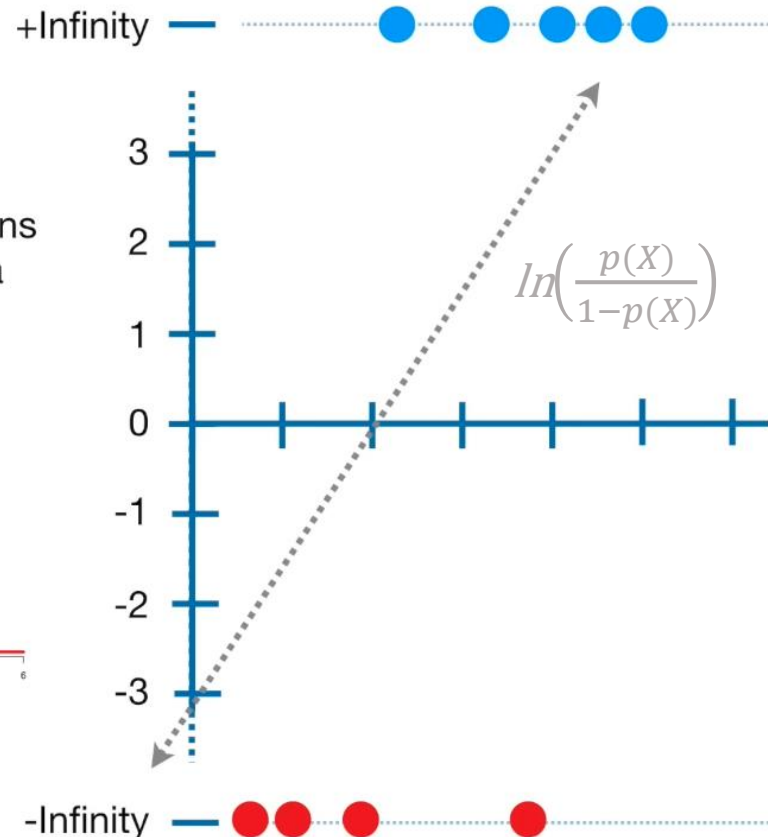
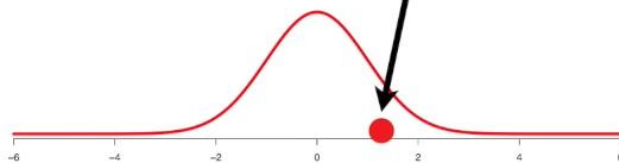
Interpreting a coefficient of Y (disease) for continuous X (weight)

- In logistic regression, **coefficients** are in terms of the *log(odds)*
 - INTERCEPT**: when *weight = 0*, *log(odds)* of disease = -3.476
 - SLOPE**: for every *additional 1 unit* of weigh, *log(odds of disease)* increases by 1.825
- Sometimes **coefficients** are converted to odds, simply by $odds = e^{\beta_j}$

$$y = -3.48 + 1.83 \times \text{weight}$$

Again, the z-value is the number of standard deviations the estimate is from 0 on a standard normal curve...

Coefficients:			
	Estimate	Std. Error	z value
(Intercept)	-3.476	2.364	-1.471
weight	1.825	1.088	1.678



Interpreting a coefficient of Y (obesity) for discrete X (YES mutated gene)

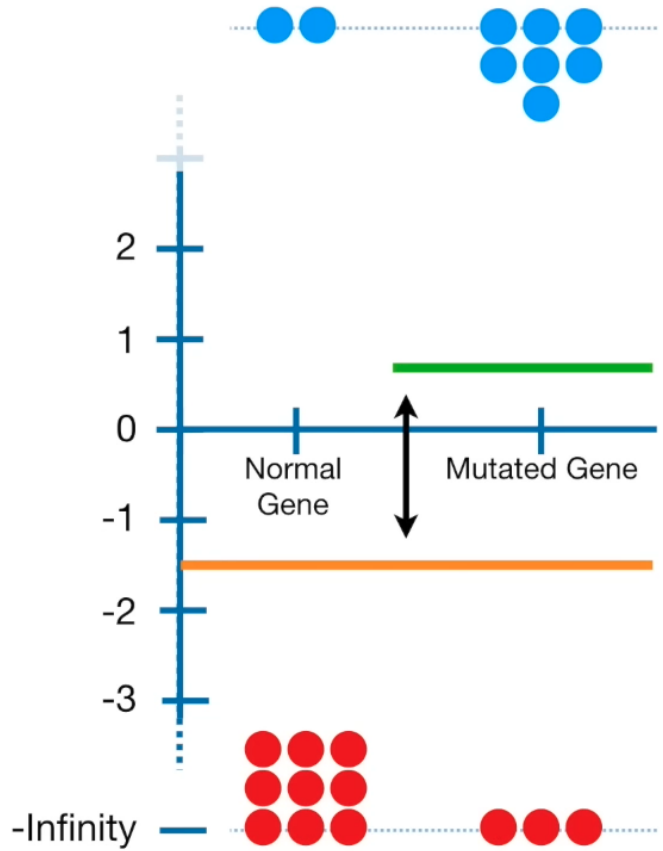
Coefficients:				
	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-1.5041	0.7817	-1.924	0.0544
geneMutant	2.3514	1.0427	2.255	0.0241

log (odds of Y)

$$= \log(\text{odds}_{\text{gene}_{\text{normal}}}) \times B_1 + \log\left(\frac{\text{odds}_{\text{gene}_{\text{mutated}}}}{\text{odds}_{\text{gene}_{\text{normal}}}}\right) \times B_2$$

↑
 ...can be converted into division, this term is a **log(odds ratio)**.

It tells us, on a log scale, how much having the mutated gene increases (or decreases) the odds of a mouse being obese.



$$\log(\text{odds of Y}) = \log(2/9) \times \beta_1 + \log\left(\frac{7/3}{2/9}\right) \beta_2 = -1.5 \times \beta_1 + 2.35 \times \beta_2$$

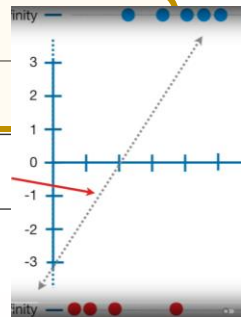
Fitting a logistic regression model (on biopsy data)

- If p_i is the probability that $y_i = 1$, i.e. the tumor is malignant, I can fit a **null model** with all the covariates x_k (continuous vars that take values from 1 to 10)

$$\text{logit}(p_i) = \ln\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \dots + \beta_k x_{k,i}$$

Logistic regression model results

term	estimate	std.error	statistic	p.value	signif. lev.
(Intercept)	-9.4169	1.1637	-8.0921	0.0000	***
clump_thickness	0.4984	0.1434	3.4758	0.0005	***
uniform_size	0.0992	0.2356	0.4211	0.6737	
uniform_shape	0.2809	0.2655	1.0580	0.2901	
marginal_adhesion	0.2688	0.1285	2.0917	0.0365	*
single_epith_size	0.0800	0.1679	0.4765	0.6337	
bare_nuclei	0.5446	0.0983	3.5054	0.0005	***
bland_chromatin	0.4083	0.1859	2.1958	0.0281	*
normal_nuclei	0.2196	0.1283	1.7119	0.0869	
mitosis	0.4356	0.3513	1.2397	0.2151	



Classification and Regression Trees (CART)

An example of **supervised** ML algorithm

Random forest classification

An example of **supervised** ML algorithm

MODULE 5 – LECTURE OUTLINE

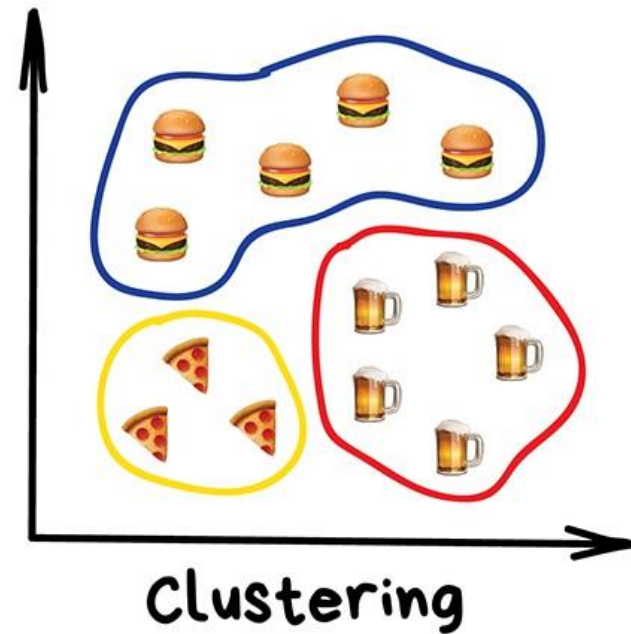
- Intro to Machine Learning (ML)
- Classification of ML algorithms
- *Supervised* ML examples
 - Logistic Regression
 - Decision trees
- *Unsupervised* ML examples
 - K-means clustering
 - PCA
 - PLS-DA

K-means Clustering

An example of **Unsupervised** ML algorithm

Clustering is a classification with no predefined classes

- *"It divides objects based on unknown features. Machine chooses the best way"*



Source: Image from https://vas3k.com/blog/machine_learning/index.html

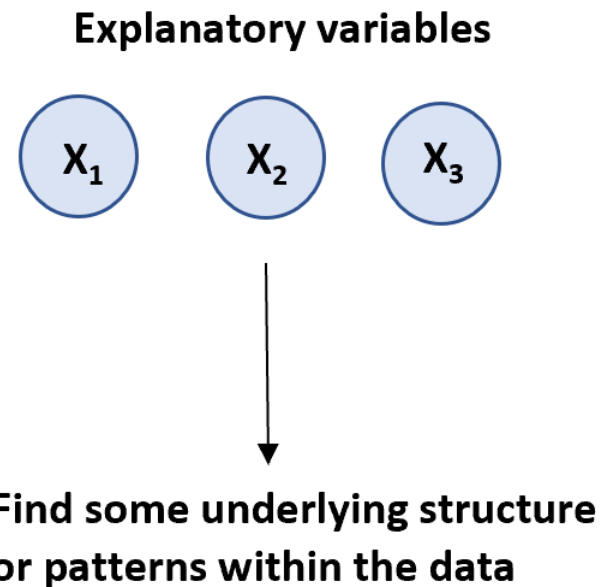
PCA for dimension reduction

An example of **Unsupervised** ML algorithm

Unsupervised Learning Algorithms mechanics

An unsupervised learning algorithm can be used when we have a list of variables ($X_1, X_2, X_3, \dots, X_p$) and we would simply like to find underlying structure or patterns within the data.

Unsupervised Learning



Source: <https://www.statology.org/supervised-vs-unsupervised-learning/>

Unsupervised Learning Algorithms typical purpose

There are two main types of unsupervised learning algorithms:

- 1. Clustering:** Using these types of algorithms, we attempt to find “clusters” of observations in a dataset that are similar to each other. This is often used in retail when a company would like to identify clusters of customers who have similar shopping habits so that they can create specific marketing strategies that target certain clusters of customers.
- 2. Association:** Using these types of algorithms, we attempt to find “rules” that can be used to draw associations. For example, retailers may develop an association algorithm that says “if a customer buys product X they are highly likely to also buy product Y.”

Unsupervised Learning: commonly used algorithms

- Most commonly used unsupervised learning algorithms:
 - Principal component analysis
 - K-means clustering
 - K-medoids clustering
 - Hierarchical clustering
 - Apriori algorithm

Summary: Supervised vs. Unsupervised Learning

- Here are the key differences between supervised and unsupervised learning algorithms:

	Supervised Learning	Unsupervised Learning
Description	Involves building a model to estimate or predict an output based on one or more inputs.	Involves finding structure and relationships from inputs. There is no “supervising” output.
Variables	Explanatory and Response variables	Explanatory variables only
End goal	Develop model to (1) predict new values or (2) understand existing relationship between explanatory and response variables	Develop model to (1) place observations from a dataset into a specific cluster or to (2) create rules to identify associations between variables.
Types of algorithms	(1) Regression and (2) Classification	(1) Clustering and (2) Association

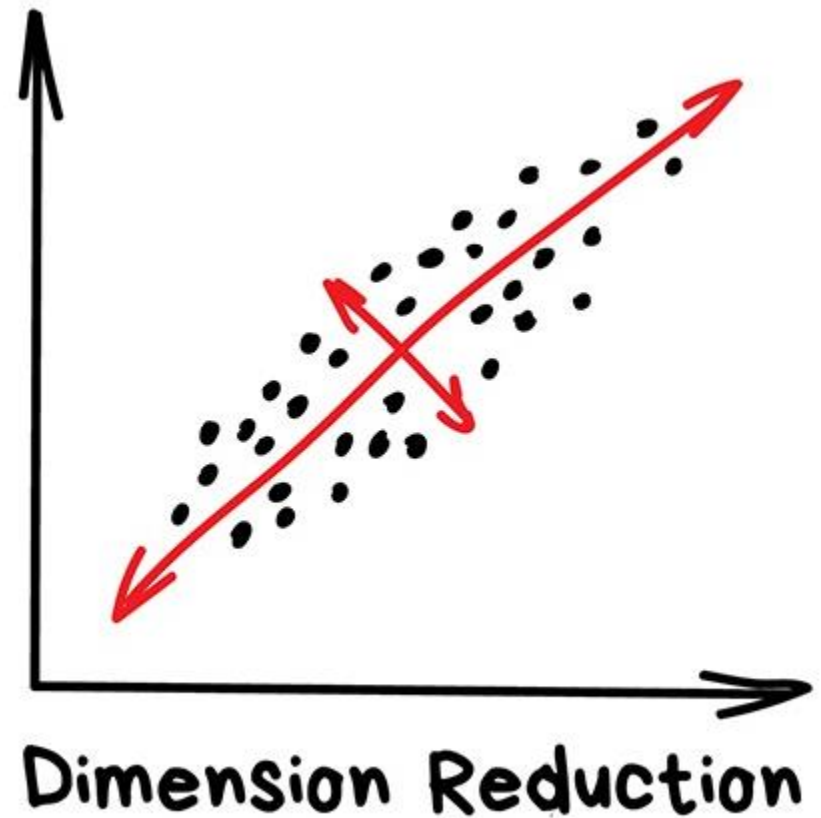
Source: <https://www.statology.org/supervised-vs-unsupervised-learning/>

Principal Component Analysis (PCA)

A type of *unsupervised* learning algorithm for
dimensionality reduction

Purpose of PCA

- The goal of PCA is to transform a high-dimensional dataset into a lower-dimensional dataset while retaining as much of the variance in the data as possible.
- Common use cases of PCA:
 1. to reduce the dimensionality of high-dimensional datasets
 2. to visualize the structure of the data
 3. to remove noise and redundant information from the data
 4. as a preprocessing step for other machine learning algorithms



Source: Image from https://vas3k.com/blog/machine_learning/index.html

Covariance

Population mean is unknown

$$\text{var}(x) = \frac{\sum_i^n (x_i - \bar{x})^2}{N - 1}$$

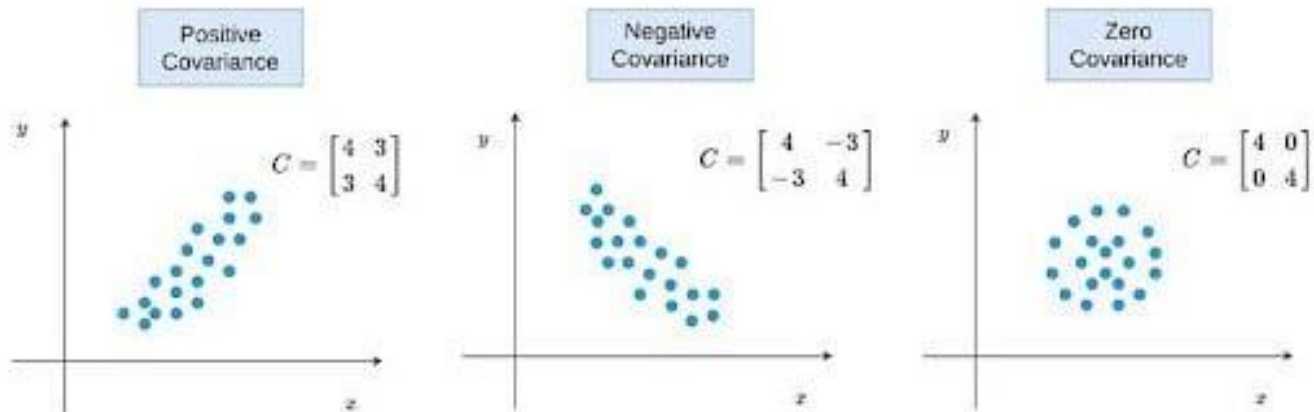
Population mean is unknown

$$\text{cov}(x, y) = \frac{\sum_i^n (x_i - \bar{x}) \cdot (y_i - \bar{y})}{N - 1}$$

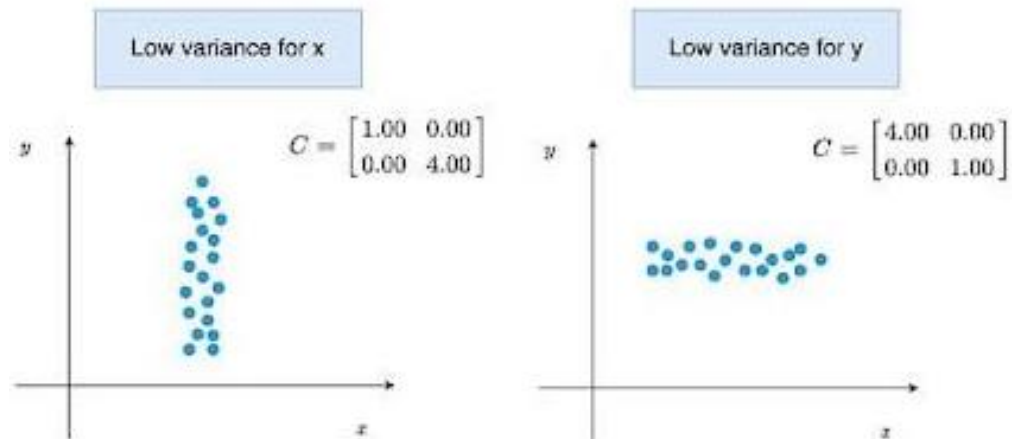
$$\begin{array}{c} \begin{array}{ccc} & x & y & z \\ x & \left[\begin{array}{ccc} \text{var}(x) & \text{cov}(x, y) & \text{cov}(x, z) \end{array} \right] \\ y & \left[\begin{array}{ccc} \text{cov}(x, y) & \text{var}(y) & \text{cov}(y, z) \end{array} \right] \\ z & \left[\begin{array}{ccc} \text{cov}(x, z) & \text{cov}(y, z) & \text{var}(z) \end{array} \right] \end{array} \end{array}$$

Variance measures how the values vary in a variable.
Covariance measures how changes in one variable are associated with changes in a second variable.

Covariance



Positive, negative and zero covariance.



Different variances and zero covariance.

Source: <https://builtin.com/data-science/covariance-matrix>

PCA

PCA originally is a [linear algebra](#) operation.

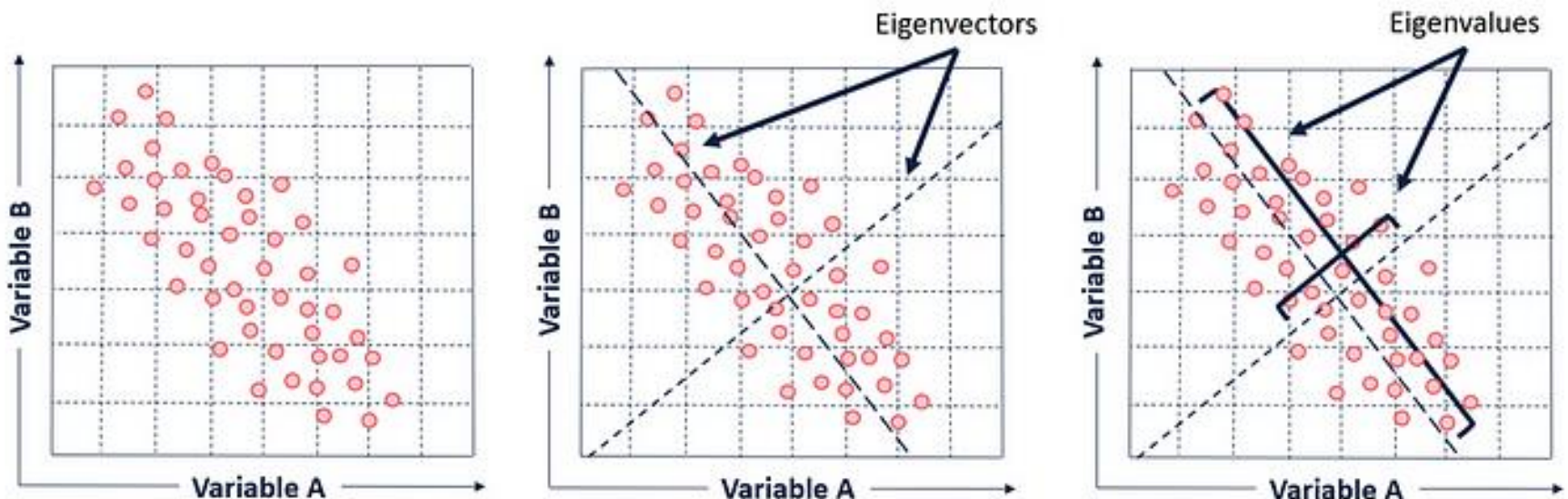
It is a transformation method that creates (weighted [linear](#)) combinations of the original variables in a data set, with the intent that the new combinations will capture as much [variance](#) in the dataset as possible while eliminating correlations (i.e., redundancy).

PCA creates the new variables using the eigenvectors and eigenvalues calculated from the [covariance matrix](#) of your original variables.

Eigenvectors & Eigenvalues

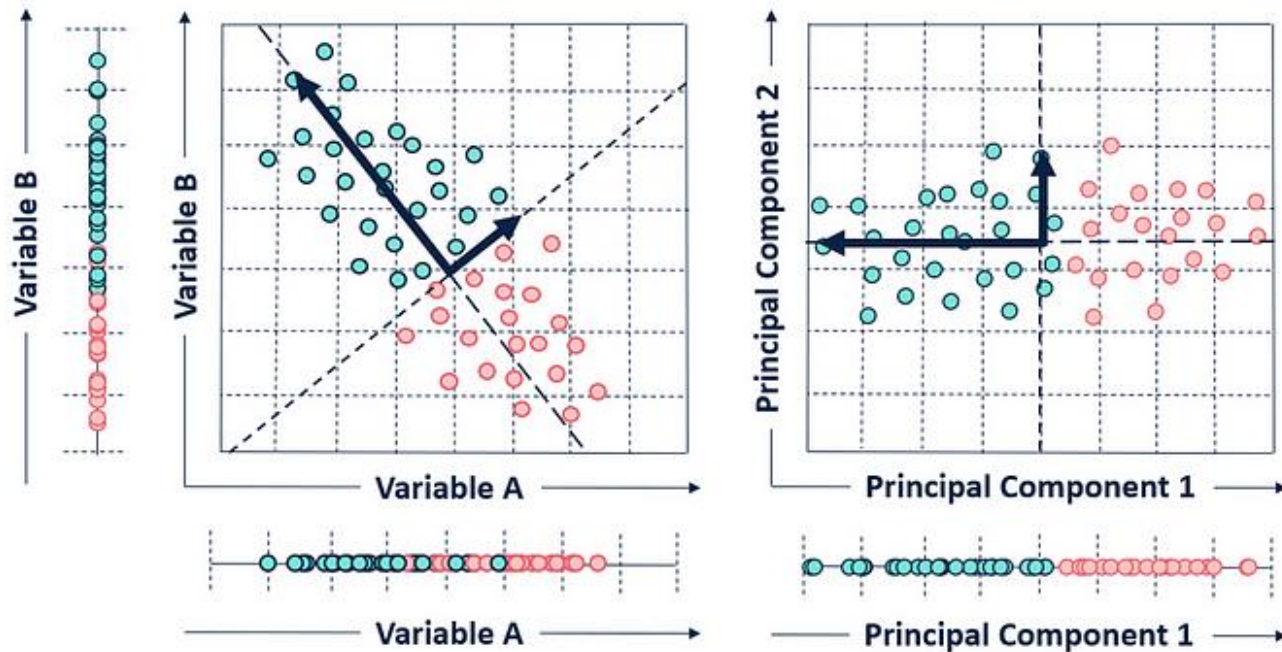
In the context of PCA

- The **eigenvectors** of the covariance matrix define the directions of the principal components calculated by PCA.
- The **eigenvalues** associated with the eigenvectors describe the variance along the new axis.



Source: <https://towardsdatascience.com/tidying-up-with-pca-an-introduction-to-principal-components-analysis-f876599af383>

Principal components

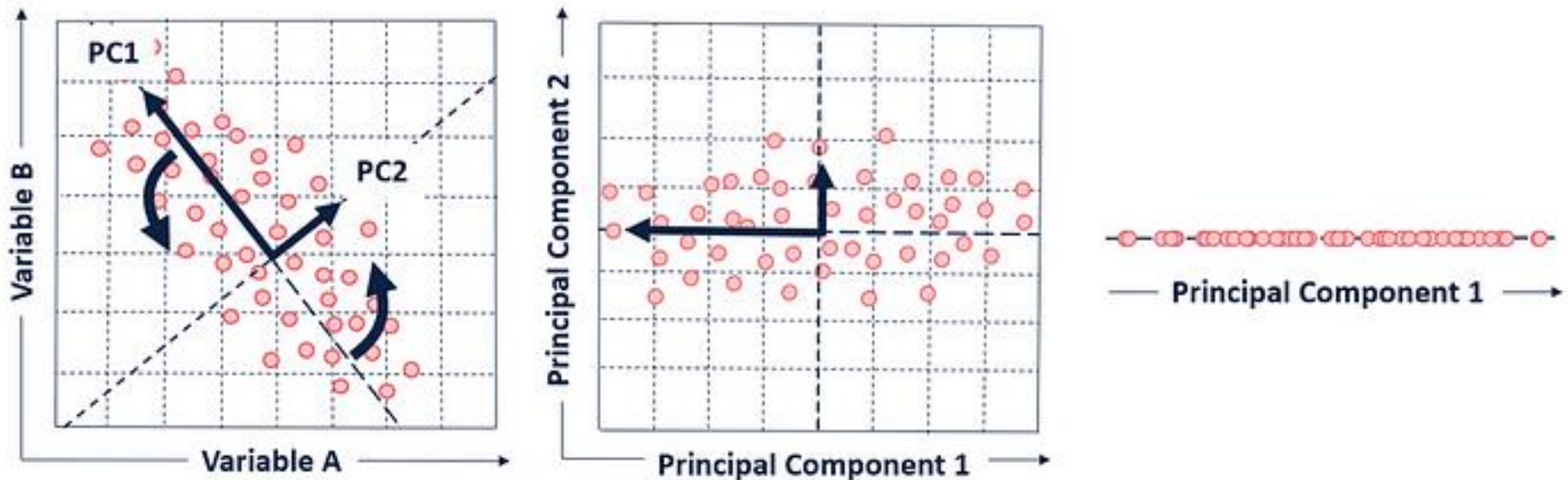


Principal Component 1 accounts for variance from both variables A and B. (dimension reduction)

The principal components (eigenvectors) are sorted by descending eigenvalue. The principal components with the highest eigenvalues are “picked first” as principal components because they account for the most variance in the data.

Source: <https://towardsdatascience.com/tidying-up-with-pca-an-introduction-to-principal-components-analysis-f876599af383>

Principal components



To convert our original points, we create a projection matrix.

This projection matrix is just the selected eigenvectors concatenated to a matrix. We can then multiply the matrix of our original observations and variables by our projection matrix.

The output of this process is a transformed data set, projected into our new data space — made up of our principal components!

Source: <https://towardsdatascience.com/tidying-up-with-pca-an-introduction-to-principal-components-analysis-f876599af383>

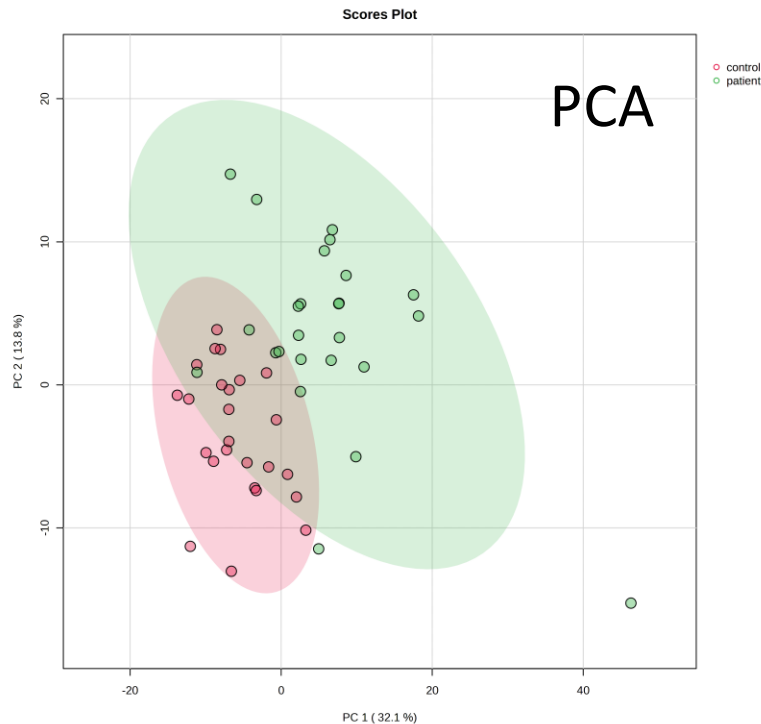
PLS Discriminant Analysis (PLS-DA)

A *supervised* alternative to PCA
performing simultaneous dimensionality
reduction and classification

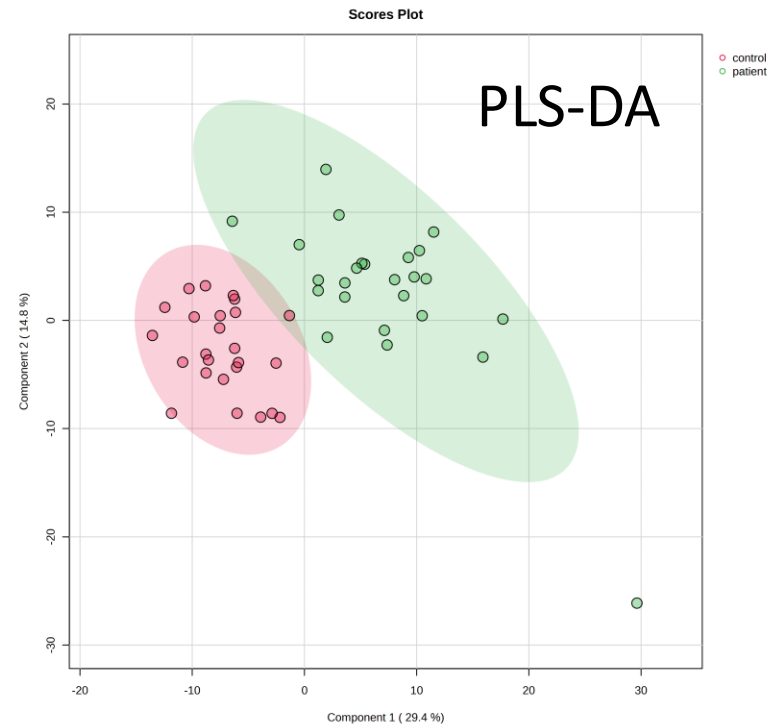
Purpose: PLS-DA vs PCA

- PCA is completely unsupervised (i.e. you don't know in advance if there are classes in your dataset)
- In PLS-DA you know how your dataset is divided in classes from the response vector Y . The goal here is then to project the predictors into a space, while maximizing the ratio $= \frac{\textit{Between group distances}}{\textit{Within group distances}}$
- Common scenarios for using PLS-DA: omics sciences.

Scores plot: PCA vs PLS-DA



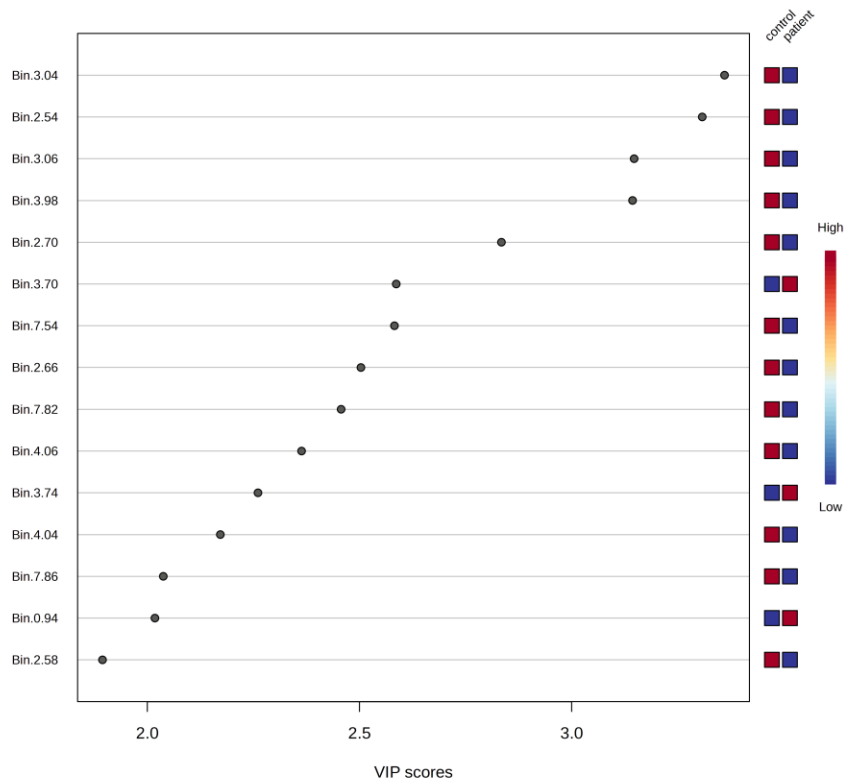
Samples projected in the space of Principal Components



Samples projected in the space of latent variables (components) that maximize the separation between groups

Source: Test data ([NMR spectral bins](https://www.metaboanalyst.ca)) provided by METABOANALYST platform: <https://www.metaboanalyst.ca>

Feature Importance in PLS-DA

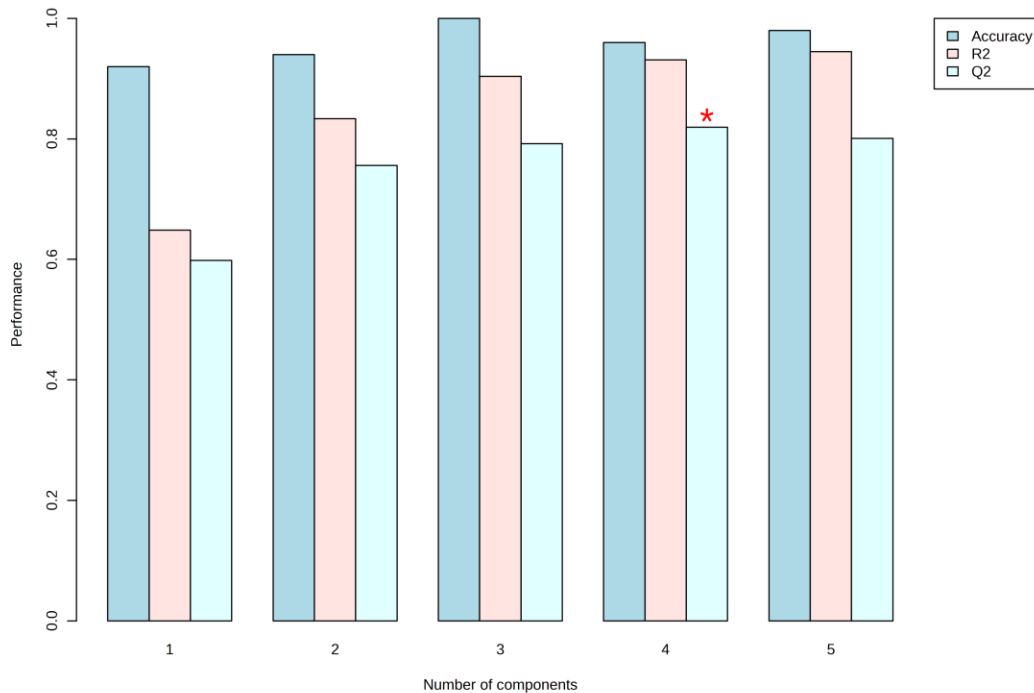


VIP (Variable Importance in Projection) scores, ranking the variables based on their significance in the PLS-DA **model of classification**.

...very useful to select potential biomarkers!

Source: Test data ([NMR spectral bins](https://www.metaboanalyst.ca)) provided by METABOANALYST platform: <https://www.metaboanalyst.ca>

Cross validation in PLS-DA



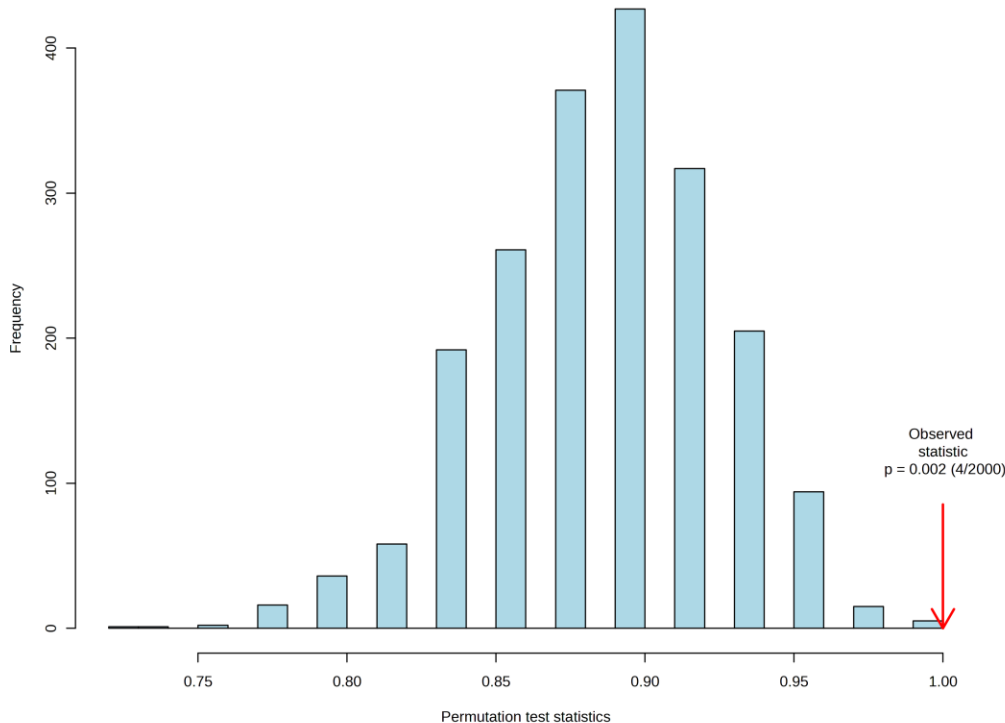
PLS-DA generate a model of classification.

By partitioning the dataset and iteratively testing the model, cross validation estimate the predictive ability of the model.

Q^2 is an analogous of R^2 in regression: the higher the better!

Source: Test data ([NMR spectral bins](https://www.metaboanalyst.ca)) provided by METABOANALYST platform: <https://www.metaboanalyst.ca>

Permutation in PLS-DA



Permutation testing is a non-parametric approach to assess the significance of a model's results.

In the context of PLS-DA, this test helps verify whether the observed classification accuracy is better than what would be expected by chance.

Test data ([NMR spectral bins](https://www.metaboanalyst.ca)) provided by METABOANALYST platform: <https://www.metaboanalyst.ca>