

<b>Business Analytics</b>	<p>Business analytics makes extensive use of statistical analysis, including explanatory and predictive modeling, and fact-based management to drive decision making. Analytics may be used as input for human decisions or may drive fully automated decisions.</p> <ul style="list-style-type: none"> <li>• <i>Descriptive analytics</i> (What has occurred?)</li> <li>• <i>Predictive analytics</i> (What is likely to occur?)</li> <li>• <i>Prescriptive analytics</i></li> </ul>
<b>From data to information</b>	<ol style="list-style-type: none"> <li>1. <i>Data consolidation</i></li> <li>2. <i>Selection and Preprocessing</i></li> <li>3. <i>Predictive analytics</i></li> <li>4. <i>Interpretation and evaluation</i></li> </ol>
<b>Numeric Prediction</b>	Given a collection of data with known numeric outputs, create a function that outputs a predicted value from a new set of outputs
<b>Classification</b>	From data with known labels, create a classifier that determines which label to apply to a new observation
<b>Clustering</b>	Identify “natural” groupings of data (no predefined groups)
<b>Association Rule Analysis</b>	Identify relationships in data from co-occurring terms or items
<b>Machine Learning</b>	<ul style="list-style-type: none"> <li>• <i>Supervised learning</i>: attempt the discovery of the relationships between input attributes and a target attribute. (→ form a description that can be used to predict unseen examples) <ul style="list-style-type: none"> <li>○ Training: estimate the prediction function <math>f</math> of the training data set</li> <li>○ Testing: apply <math>f</math> to a never seen test example</li> </ul> </li> <li>• <i>Unsupervised learning</i>: there is no supervisor and only input data is available. (→ find regularities, irregularities, relationships, similarities and associations in the input)</li> </ul>
<b>Model</b>	<p>Representation of a system that allows for investigation of the properties of the system and, in some cases, prediction of future outcomes.</p> <ul style="list-style-type: none"> <li>• <i>Decision Trees</i></li> <li>• <i>Neural Networks</i></li> </ul>
<b>Causal Inference</b>	<p>Causal inference is the process of deriving cause-and-effect conclusions by reasoning from knowledge and factual evidence.</p> <p>Randomized control trials are the gold standard for causal inference in the social sciences, but often only observational data is available.</p>
<b>Econometrics</b>	Data analysis focusing on causal relationships in economics (Correlation $\neq$ Causation)
<b>Descriptive statistics</b>	Descriptive statistics can be used to summarize the data, either numerically or graphically, to describe the sample (e.g. mean, standard deviation,...)
<b>Inferential statistics</b>	Inferential statistics is used to model patterns in the data, accounting for randomness and drawing inferences about the larger population (estimation, hypothesis testing, forecasting, correlation, regression)
<b>Random Variables</b>	<p><math>X</math> is a random variable if it represents a random draw from some population and is associated with a probability distribution</p> <ul style="list-style-type: none"> <li>• <i>Discrete</i>: random variable can take only selected values</li> <li>• <i>Continuous</i>: random variable can take any value in the real interval</li> </ul> <p><b>Standard Normal</b>: any random variable can be “standardized” by subtracting the mean and dividing by the standard deviation.</p> <p><b>Expected Value</b>: the expected value of a probability weighted average of <math>X</math> is the mean or expected value of the distribution of <math>X</math>.</p>
<b>Random Sample</b>	A random sample is a set of independent, identically distributed random variables (every combination of $n$ sample points has an equal chance of being selected)
<b>Statistical Estimation</b>	<p>Every member of the population has the same chance of being selected in the sample.</p> <p>An estimator is a statistic that is used to estimate an unknown population parameter.</p> <ul style="list-style-type: none"> <li>• <i>Point estimate</i>: sample mean, sample proportion</li> <li>• <i>Interval estimate</i>: confidence interval for mean, confidence interval for proportion</li> </ul> <p>Point estimate is always within the interval estimate</p>
<b>Standard Deviation</b>	The standard deviation of the sample means is equal to the standard deviation of the population divided by the square root of the sample size
<b>Central Limit theorem</b>	The central limit theorem states that the standardized average of any population of random variables is asymptotically $\sim N(0,1)$
<b>Confidence Interval (CI)</b>	<p>Provide us with a range of values that we believe, with a given level of confidence, contains a population parameter CI for the population means</p> <ul style="list-style-type: none"> <li>• Larger sample = smaller interval</li> <li>• Lower confidence interval = smaller interval</li> </ul>

	<ul style="list-style-type: none"> <li>• More variation = larger interval</li> </ul>
<b>Statistical test</b>	<ol style="list-style-type: none"> <li>1. <i>Formulate hypothesis</i></li> <li>2. <i>Collect data to test hypothesis</i> (systematic error can be controlled by statistical significance or by confidence interval)</li> <li>3. <i>Accept or reject hypothesis</i></li> </ol>
<b>Hypothesis testing</b>	<ol style="list-style-type: none"> <li>1. <i>State null and alternative hypothesis</i> (<math>H_0</math> usually a statement of no effect)</li> <li>2. <i>Choose <math>\alpha</math> level</i> (related to confidence interval) - probability of falsely rejecting the <math>H_0</math></li> <li>3. <i>Calculate test statistic and find p-value</i> (how far data are from null hypothesis)</li> <li>4. <i>State conclusion</i> <ol style="list-style-type: none"> <li>a. <math>P \leq \alpha</math>, reject null hypothesis</li> <li>b. <math>P &gt; \alpha</math>, insufficient evidence to reject null hypothesis</li> </ol> </li> </ol> <p><b>Type I error:</b> reject null hypothesis when the null hypothesis is actually true  <b>Type II error:</b> fail to reject the null hypothesis when the alternative hypothesis is true</p>
<b>Student t-Distribution</b>	When the population is normally distributed, the statistic t is Student t distributed (bell-shaped and symmetric around zero)
<b>T-Test</b>	<ul style="list-style-type: none"> <li>• <i>Single sample:</i> tests whether a sample mean is significantly different from a pre-existing value</li> <li>• <i>Paired samples:</i> tests the relationship between 2 linked samples</li> <li>• <i>Independent samples:</i> tests the relationship between 2 independent populations</li> </ul>
<b>Selected Statistical Tests</b>	<ul style="list-style-type: none"> <li>• <i>Parametric tests</i> <ul style="list-style-type: none"> <li>○ T-Tests: compares two sample means or tests a single sample mean</li> <li>○ F-Test: compares the equivalence of variances of two samples</li> </ul> </li> <li>• <i>Non-parametric tests</i> <ul style="list-style-type: none"> <li>○ Wilcoxon signed-rank test: independence of two means for 2 paired samples when normality is not assumed</li> <li>○ Mann-Whitney-U test: used for 2 independent samples</li> <li>○ Kruskal-Wallis-Test: equivalence of multiple means in case of several non-normally distributed samples</li> </ul> </li> <li>• <i>Tests of the Probability Distribution</i> <ul style="list-style-type: none"> <li>○ Kolmogorov-Smirnov and Chi-square test: used to determine whether two underlying probability distributions differ, or whether an underlying probability distribution differs from a hypothesized distribution</li> </ul> </li> </ul>
<b>Linear Regression</b>	<p>Regressions identify relationships between dependent and independent variables. The linear regression is a statistical tool for numerical predictions</p> $Y = \beta_0 + \beta_1 X + \varepsilon$ (first order linear model) $\beta_0$ and $\beta_1$ are unknown, therefore, are estimated from the data
<b>Ordinary Least Squares (OLS)</b>	<p>Minimize the sum of squared residuals.</p> <p>OLS requires choosing values of the estimated coefficients, such that Residual Sum of Squares (RSS) is as small as possible for the sample.</p>
<b>Residual Sum of Squares (RSS)</b>	Sum of squared differences between the points and the regression line (how well line fits the data)
<b>Total Deviation</b>	The Total Sum of Squares (TSS) is the sum of the Explained Sum of Squares (ESS) and the RSS
<b>Coefficient of Determination</b>	$R^2$ measure the proportion of the variation in y that is explained by the variation in x ( $R^2 = 1 \rightarrow$ perfect match between the line and data points, $R^2 = 0 \rightarrow$ no linear relationship between x and y)
<b>Multiple Linear Regression</b>	<p>A p-variable regression model can be expressed as a series of equations. Beta coefficients are known as partial regression coefficients</p> $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 = \mathbf{X}\hat{\beta}$ <p><b>Adjusted <math>R^2</math>:</b> It represents the proportion of variability of y explained by X. <math>R^2</math> is adjusted so that models with a different number of variables can be compared</p> <p>F-Test: significant F indicates a linear relationship between y and at least one of the <math>X_s</math></p> <p>T-Test: significant t indicates that the variable in question influences the response variable while controlling for other explanatory variables</p>
<b>Model Specification</b>	<p>In regression analysis the specification is the process of developing a regression model (selecting an appropriate functional form and choosing which variables to include). The model might include irrelevant variables or omit relevant variables.</p> <p>Non-linear models are challenging, but some nonlinear regression problems can be linearized.</p>
<b>Subset Selection</b>	<p>Fit a parsimonious model that explains variation in Y with a small set of predictors</p> <ul style="list-style-type: none"> <li>• <i>Best subset</i> (computationally expensive)</li> <li>• <i>Backward elimination</i> (top-down approach)</li> <li>• <i>Forward selection</i> (bottom-up approach)</li> <li>• <i>Stepwise regression</i> (combines forward and backward)</li> </ul>

<b>Gauss-Markov Theorem</b>	<p>The Gauss-Markov theorem states that in a linear regression model in which errors</p> <ul style="list-style-type: none"> <li>• Have expectation zero</li> <li>• Are uncorrelated</li> <li>• Have equal variances</li> </ul> <p>The best (lowest variance) linear unbiased estimator (BLUE) of the coefficients is given by the ordinary least squares (OLS) estimator.</p>
<b>Unbiased</b>	Expected value for estimator "is true" ( $E(\beta) = \beta$ )
<b>Consistent</b>	$\text{Var}(\beta)$ decreases with increasing sample size $n$
<b>Efficient</b>	Estimator $\beta$ has lower variance than any other estimator
<b>OLS assumptions</b>	<ul style="list-style-type: none"> <li>• <i>Linearity</i>: linear relationship in parameters <math>\beta</math> (when linearity does not hold, try to reformulate)</li> <li>• <i>No multicollinearity</i>: no linear dependency between predictors</li> <li>• <i>Homoscedasticity</i>: residuals exhibit constant variance (the spread of the data points does not change much)</li> <li>• <i>No autocorrelation</i>: there is no correlation between the <math>i</math> and <math>j</math> residual terms</li> <li>• <i>Exogeneity</i>: expected value of the residual vector, given <math>X</math>, is 0</li> </ul>
<b>Outlier</b>	<p>An outlier is an observation that is unusually large or small. <b>Possibilities:</b></p> <ul style="list-style-type: none"> <li>• There was an error in recording the value</li> <li>• The point does not belong in the sample</li> <li>• The observation is valid</li> </ul>
<b>Multicollinearity check</b>	<ol style="list-style-type: none"> <li>1. <i>Calculate the correlation coefficient for each pair of predictor variables</i> → Large correlations (greater than the correlations between predictor and response) indicate problems.</li> <li>2. <i>Variance Inflation Factor (VIF)</i>: <math>VIF = \frac{1}{1-R_k^2}</math> where the <math>R_k^2</math> is the value when the predictor in question (<math>k</math>) is set as the dependent variable (how much of the variance can be explained) → remove variables with VIF scores greater than 10</li> </ol> <p>If the variable has a non-significant t-value, then either</p> <ul style="list-style-type: none"> <li>• The variable is not related to the response</li> <li>• The variable is not related to the response, but it is not required in the regression because it is strongly related to a third variable that is in the regression</li> </ul> <p>→ The usual remedy is to drop one or more variables from the model</p>
<b>Heteroscedasticity</b>	<p>When the requirement of a constant variance is violated.</p> <p>Breusch-Pagan test or White test are used to check for heteroscedasticity. If there is heteroscedasticity, the estimated <math>\text{Var}(\beta)</math> is biased and OLS might not be efficient anymore</p>
<b>Autocorrelation</b>	<p>By examining the residuals over time, no pattern should be observed. Reasons of autocorrelation:</p> <ul style="list-style-type: none"> <li>• Omission of an important variable</li> <li>• Functional misfit</li> <li>• Measurement error in the independent variable</li> </ul> <p>Durbin-Watson (DW) statistic to test for first order autocorrelation.</p>
<b>Modeling seasonality</b>	<p>A regression can estimate both the trend and additive seasonal indexes:</p> <ul style="list-style-type: none"> <li>• Create dummy variables which indicate the season</li> <li>• Regress on time and the seasonal variables</li> <li>• Use the multiple regression model to forecast</li> </ul> <p><math>y = \beta_0 + \beta_1 t + \beta_2 Q_1 + \beta_3 Q_2 + \beta_4 Q_3</math> (if all 4 <math>Q</math>s are modelled → multicollinearity!)</p>
<b>Exogeneity</b>	Other factors, which are not explicitly accounted for in the model but are contained in the error term, are not correlated with $X$
<b>Endogeneity</b>	Endogeneity is given when an independent variable is correlated with the error term and the covariance is not null (e.g. omitted variable)
<b>Cross-section data</b>	Refers to data observing many subjects at the same point in time, or without regard to differences in time (there might be omitted variables describing important characteristics of individuals)
<b>Panel Data</b>	<p>A panel data set is one where there are repeated observations on the same units</p> <ul style="list-style-type: none"> <li>• <i>Balanced panel</i>: every unit is surveyed in every time period</li> <li>• <i>Unbalanced panel</i>: some individuals have not been recorded in some time period</li> </ul> <p>Individual effects:</p> <ul style="list-style-type: none"> <li>• <i>Fixed effects</i>: individual-specific effects are correlated to other covariates (endogeneity)</li> <li>• <i>Random effects</i>: individual-specific effects are uncorrelated to other covariates</li> </ul> <p>The Hausman test can help decide on one or the other (the test takes into account the covariance matrix of the FE and RE estimators as well as the estimates and follows a chi-square distribution)</p>
<b>Fixed effect model</b>	<p>Treat <math>\lambda</math> (the individual-specific heterogeneity) as a constant for each individual</p> <p><math>y_{it} = (\beta_0 + \lambda_i) + \beta_1 x_{1it} + \beta_2 x_{2it} + \dots + \beta_p x_{pit} + \varepsilon_{it}</math></p>

	<b>Estimators:</b> first differences, within, between, least squares dummy variable estimator
<b>Random effect model</b>	The larger the variance of individual effects $\lambda$ , the more random effects accounts for it. $Y_{it} = \beta_0 + \beta_1 X_{1it} + \beta_2 X_{2it} + \beta_3 X_{3i} + \beta_4 X_{4i} + \varepsilon_{it}$
<b>Logistic Regression</b>	The logistic function $\Pr[Y X]$ constraints the estimated probabilities to lie between 0 and 1. $\Pr[X Y]$ is the estimated probability that the $i^{\text{th}}$ case is in a category and $\beta_0 + \beta_1 X$ is the regular linear regression equation $\rightarrow$ Probability of success ( $Y=1$ ) given the predictor variable ( $X$ ) is a logistic function (non-linear) <ul style="list-style-type: none"> <li><math>\beta_0</math> is the regression constant</li> <li><math>\beta_1 X</math> is the regression slope (if <math>\beta_1 &lt; 0</math> and <math>X</math> increases <math>\rightarrow</math> the odds go down, opposite by <math>\beta_1 &gt; 0</math>)</li> </ul> <b>Odds:</b> by algebraic manipulation, the logistic regression equation can be written in terms of an odd of success (range from 0 to INF). <b>Logit:</b> taking the natural log of both sides, we can write the equation in terms of logits (log-odds). The presence of multicollinearity will not lead to biased coefficients, but it will have an effect on the standard errors. The inclusion of irrelevant variables can result in a poor model fit (consult Wald statistics). Multiple Logistic Regression: more than one independent variable
<b>Multinomial Logit model</b>	The dependent variable, $Y$ , is a discrete variable that represents a choice, or category, from a set of mutually exclusive choices or categories (more than 2). $\rightarrow$ <b>conditional logit</b> Ordered logit models have ordinal dependent variables.
<b>Generalized Linear Models (GLM)</b>	GLMs are a general class of linear models that are made up for three components: <ul style="list-style-type: none"> <li><i>Random component:</i> identifies dependent variable <math>\mu</math> and its probability distribution</li> <li><i>Systematic component:</i> identifies the set of explanatory variables (<math>X_1, \dots, X_k</math>)</li> <li><i>Link function:</i> identifies a function of <math>\mu</math> that is a linear function of the explanatory variables <ul style="list-style-type: none"> <li>Identity link: form used in normal linear regression models</li> <li>Log link: used when <math>\mu</math> cannot be negative as when data are Poisson counts</li> <li>Logit link: used when <math>\mu</math> is bounded between 0 and 1 as when data are binary</li> </ul> </li> </ul>
<b>Maximum Likelihood Estimation (MLE)</b>	Maximum Likelihood estimation is a statistical method for estimating the coefficients of a model. (MLE involves finding the coefficients that make the log of the likelihood function (LL<0) as large as possible). Likelihood Function (L): measures the probability of observing the particular set of dependent variable values that occur in the sample
<b>Goodness of Fit</b>	<ul style="list-style-type: none"> <li><i>Null model:</i> assumes one parameter (the intercept) for all of the data points, which means you only estimate 1 parameter</li> <li><i>Fitted model:</i> assumes you can explain your data points with <math>p</math> parameters and an intercept term, so you have <math>p+1</math> parameters</li> <li><i>Null deviance</i> (<math>-2\ln(L(\text{null}))</math>): how much is explained by a model with only the intercept</li> <li><i>Residual deviance</i> (<math>-2\ln(L(\text{fitted}))</math>): small values mean that the fitted model explains the data well</li> <li><i>Likelihood ratio test:</i> Non-significant <math>\chi^2</math> values indicate that a significant amount of the variance is unexplained</li> <li><i>Wald test:</i> used to test the statistical significance of each coefficient in the model hypothesis that <math>\beta_i = 0</math></li> </ul>
<b>Count Variables</b>	Count variables are non-negative integers. (OLS bad regression model as count variables cannot be negative and are often highly skewed).
<b>Count models</b>	<ul style="list-style-type: none"> <li><i>Poisson regression model</i></li> <li><i>Negative binomial regression model</i></li> </ul>
<b>Poisson Regression</b>	In Poisson regression, $y$ is typically conceptualized as a rate (positive coefficients indicate higher rate). Poisson models are non-linear (coefficients don't have a simple linear interpretation) Poisson Model has a log form; exponentiation aids interpretation <b>Assumptions:</b> <ul style="list-style-type: none"> <li>The mean and variance are the same (often not met in real data <math>\rightarrow</math> variance is often greater than <math>\mu</math>: Overdispersion)</li> <li><b>Overdispersion</b> <math>\rightarrow</math> standard errors will be underestimated; potential overconfidence in results</li> </ul> $\rightarrow$ Negative binomial regression as alternative to Poisson regression
<b>Zero-Inflation</b>	If outcome variable has many zero values it tends to be highly skewed. <b>Zero-inflated model:</b> assume two types of groups in your sample: <ol style="list-style-type: none"> <li>Type A: Always zero (no probability of non-zero values)</li> <li>Type ~A: non-zero chance of positive count value</li> </ol> $\rightarrow$ Use logit to model group membership (A or ~A) $\rightarrow$ Use Poisson or NB regression to model counts for those in group ~A $\rightarrow$ Compute probabilities based on those results

<b>Classification</b>	<p>Given a database D of tuples and a set of classes C, the classification problem is to define a mapping <math>D \rightarrow C</math> where each x is assigned to one class. A class contains precisely those tuples mapped to it. Prediction is similar, but usually implies a mapping to numeric values instead of a class C.</p> <p>Algorithms:</p> <ul style="list-style-type: none"> <li>• <i>Logistic Regression</i></li> <li>• <i>Statistical Modeling</i> (e.g. Naïve Bayes)</li> <li>• <i>Decision Trees</i>: divide and conquer</li> <li>• <i>Classification Rules</i> (e.g. PRISM)</li> <li>• <i>Instance-Based Learning</i> (e.g. kNN)</li> <li>• <i>Support Vector Machines</i></li> <li>• ...</li> </ul>
<b>Naïve Bayes Classifier</b>	<p>Naïve Bayes classifier takes all attributes into account. <b>Assumptions:</b></p> <ul style="list-style-type: none"> <li>• All attributes are equally important</li> <li>• All attributes independent (value of one value tells nothing about value of another attribute)</li> </ul> <p>Adding to many redundant attributes will cause problems</p>
<b>Bayesian (Belief) Networks</b>	<p>Bayesian Belief Network describe conditional independence among subsets of attributes: combining prior knowledge about dependencies among variables with observed training data (Graphical representation: directed acyclic graph)</p> <p>Ex. <math>P(A,B,C,D,E) = P(A) P(B A) P(C A,B) P(D A,B,C) P(E A,B,C,D)</math></p> <p><b>Inference:</b> other than the joint probability of specific events, we may want to infer the probability of an event, given observations about a subset of other variables</p>
<b>Conditional probability</b>	<p>Probability of events change when we know something of the world.  <math>P(A B) \rightarrow</math> Probability of A given that we know B is true <math>= P(A \cap B) / P(B)</math>          If events A and B do not influence each other <math>\rightarrow P(A B) = P(A)P(B)</math></p> $Pr[h e] = \frac{Pr[e h]Pr[h]}{Pr[e]}$ <p><b>Normalization:</b></p> $Pr[h e] = \frac{Pr[e_1 h]Pr[e_2 h] \dots Pr[e_n h]Pr[h]}{Pr[e]}$ <p>Dealing with numeric attributes <math>\rightarrow</math> <b>Probability density function</b></p> $f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (\text{normal or gaussian probability})$ <p>For unknown distributions <math>\rightarrow</math> <b>Kernel Density Estimate</b></p>
<b>Zero Frequency Problem</b>	<p>Attribute value does not occur with every class value <math>\rightarrow</math> add 1 to the numerator for every attribute-class combination, and the probability can never be 0</p> <p><b>Modified Probability Estimates:</b> in some cases adding a constant different from 1 might be more appropriate</p>
<b>Probability Laws</b>	<ul style="list-style-type: none"> <li>• <i>Chain rule</i>: the joint distribution is independent of the ordering</li> <li>• <i>Conditional independence</i>: event A causes something in event B (rain <math>\rightarrow</math> umbrella)</li> </ul>
<b>Decision Trees</b>	<ul style="list-style-type: none"> <li>• <i>Internal node</i>: test on an attribute</li> <li>• <i>Branch</i>: outcome of the test</li> <li>• <i>Leaf node</i>: class label or class label distribution</li> </ul>
<b>Building Decision Tree</b>	<ul style="list-style-type: none"> <li>• <i>Top-down tree construction</i>: all training examples are at the root. Partition the examples recursively by choosing one attribute each time</li> <li>• <i>Bottom-up tree pruning</i>: remove subtrees or branches, in a bottom-up manner, to improve the estimated accuracy on new cases (on the basis of an evaluation)</li> </ul> <p>The output of decision trees can be used for descriptive as well as predictive purposes</p>
<b>Attribute selection</b>	<ul style="list-style-type: none"> <li>• Choose the attribute which will result in the smallest tree (not good)</li> <li>• Choose the attribute that produces the "purest" nodes (heuristics)</li> <li>• Choose attribute that results in the greatest information gain</li> </ul>
<b>Entropy</b>	<p>Gives the information required in bits. The distribution's entropy represents the info required to predict an event</p> $\text{entropy}(p_1, p_2, \dots, p_n) = -p_1 \log_2 p_1 - p_2 \log_2 p_2 \dots - p_n \log_2 p_n$ <p>(entropy = 0 <math>\rightarrow</math> perfectly ordered system)</p>
<b>Expected Information Gain</b>	<p>Gain(S,a) is the information gained adding a sub-tree. Problems:</p> <ul style="list-style-type: none"> <li>• Attributes with a large number of values (e.g. ID codes)</li> <li>• Subsets are more likely to be pure if there is a large number of values <math>\rightarrow</math> Overfitting</li> </ul> <p><b>Gain ratio:</b> modification of the information gain that reduces its bias on high-branch attributes. Gain ratio takes number and size of branches into account when choosing an attribute</p> <p>(Problem: gain ratio may overcompensate <math>\rightarrow</math> choose an attribute just because intrinsic information is very low)</p>




	Intrinsic Information: how much info do we need to tell which branch an instance belongs to
<b>Splitting Criterion</b>	<ul style="list-style-type: none"> <li>• <i>Gini Index</i>: frequency of positive and negative classes (select the split that decreases the Gini index most)</li> <li>• <i>C4.5 and CART</i></li> </ul>
<b>Industrial-Strength Algorithm</b>	<p>For an algorithm to be useful in a wide range of real-world applications it must: (e.g. C4.5)</p> <ul style="list-style-type: none"> <li>• <i>Permit numeric attributes</i>: every attribute has many possible split points (evaluate info gain) Numeric attributes may be tested several times along a path in the tree</li> <li>• <i>Allow missing values</i></li> <li>• <i>Be robust in the presence of noise</i></li> </ul>
<b>Handling Missing Values (Tree)</b>	<ul style="list-style-type: none"> <li>• <i>Ignore instances with missing values</i></li> <li>• <i>Ignore attributes with missing values</i></li> <li>• <i>Treat missing value as another nominal value</i></li> <li>• <i>Estimate missing value</i></li> <li>• <i>Follow the leader</i>: an instance with a missing value is sent down the branch with the most instances</li> <li>• <i>Partition the instance</i>: send down part of the instance proportional to the number of training instances</li> </ul>
<b>Overfitting</b>	<p>Chasing every abnormality (noise, outliers) causes overfitting:</p> <ul style="list-style-type: none"> <li>• Decision tree gets too large and complex</li> <li>• Good accuracy on training set, poor accuracy on test set</li> <li>• Does not generalize the data any more</li> </ul> <p>→ <b>Prune the tree</b></p>
<b>Propositional rules</b>	Rules comparing attributes to constants are called propositional rules
<b>Pruning</b>	<ul style="list-style-type: none"> <li>• <i>Prepruning</i>: tries to decide a priori when to stop creating subtrees (halt tree construction)</li> <li>• <i>Postpruning</i>: simplifies an existing decision tree (construct complete tree) <ul style="list-style-type: none"> <li>◦ Subtree replacement: replace a subtree with a single leaf node</li> </ul> </li> </ul> <p>To determine if a node should be replaced, compare the error rate estimate for the node with the combined error rates of its children → replace the node if combined error rate is higher</p>
<b>Data understanding</b>	<ul style="list-style-type: none"> <li>• <i>Quantity</i> <ul style="list-style-type: none"> <li>◦ Number of instances (&gt;5000)</li> <li>◦ Number of attributes (&lt;50)</li> <li>◦ Number of targets (&gt;100 for each class)</li> </ul> </li> <li>• <i>Visualization</i></li> <li>• <i>Data summaries</i></li> </ul>
<b>Data preparation</b>	<ul style="list-style-type: none"> <li>• <i>Data cleaning</i>: <ul style="list-style-type: none"> <li>◦ Missing values</li> <li>◦ Discretization: reduces the number of values for a continuous attribute (binning) <ul style="list-style-type: none"> <li>▪ Equal Frequency: bins with equal number of instances</li> <li>▪ Class Dependent (supervised discretization)</li> </ul> </li> </ul> </li> <li>• <i>Conversion</i>: ordered attributes to numeric attributes</li> <li>• <i>Building balanced train sets</i> → <b>Subset Selection</b> <ul style="list-style-type: none"> <li>◦ Best subset (computationally expensive)</li> <li>◦ Backward elimination (top-down approach)</li> <li>◦ Forward selection (bottom-up approach)</li> <li>◦ Stepwise regression (combines forward/backward)</li> </ul> </li> </ul>
<b>Validity</b>	<ul style="list-style-type: none"> <li>• <i>External validity</i>: a statistical study has external validity if its inferences and conclusions can be generalized from the population and setting studied to other populations and settings</li> <li>• <i>Internal validity</i>: a statistical analysis has internal validity if the statistical inferences about causal effects are valid for the population being studied. Problems: <ul style="list-style-type: none"> <li>◦ Misspecification of the functional form of the model</li> <li>◦ Measurements errors in the independent variable</li> <li>◦ Simultaneous causality</li> <li>◦ Omitted variable bias</li> <li>◦ Sample selection bias (if data is not collected via a randomized controlled trial)</li> </ul> </li> </ul>
<b>Data Collection</b>	<ul style="list-style-type: none"> <li>• <i>Randomized controlled trials (RCTs)</i>: randomized experiments, where each subject is randomly assigned to a treated group or a control group in order to control for extraneous factors</li> <li>• <i>Quasi-experiments</i>: compare natural groups and measure effects without randomization of the subjects. The independent variable is controlled, but the assignment of subjects is not random</li> </ul>

	<ul style="list-style-type: none"><li>• <i>Observational studies</i>: draw inferences from a sample to a population where the independent variable is not under control of the researcher<ul style="list-style-type: none"><li>○ <u>Cross-sectional studies</u>: data collection at one specific point in time</li><li>○ <u>Longitudinal study</u>: repeated observations of the same variables over long periods of time</li><li>○ <u>Panel study</u>: group of subjects is closely monitored over a span of time</li><li>○ <u>Case-control study</u>: two existing groups differing in outcome are identified and compared on the basis of some supposed causal attribute</li></ul></li></ul>									
Confounding variable	<p>Is an extraneous variable in a statistical model that correlates (directly or inversely) with both the dependent variable and the independent variable, in a way that “explains away” some or all of the correlation between these two variables.</p> <p><b>Identification strategies:</b></p> <ul style="list-style-type: none"><li>• <i>Randomized controlled trials</i></li><li>• <i>Fixed effects models for panel data</i>: eliminates alternative explanations that are “fixed” across units</li><li>• <i>Propensity score matching</i>: only difference in similar subjects is treatment<ul style="list-style-type: none"><li>a. Estimate propensity score: individual being selected in the treatment</li><li>b. Match subjects with similar propensity score (balance the pretreatment covariates)</li><li>c. Evaluate quality of matching: check if the treatment and comparison group are similar</li><li>d. Evaluate outcomes</li></ul></li><li>• <i>Instrument variables</i>: determines observed versus unobserved explanations for taking treatment, and only uses observed portion</li><li>• <i>Regression discontinuity analysis</i></li><li>• <i>Difference-in-differences for quasi-experimental data</i>: if the treatment in quasi-experiment is as if subjects were randomly assigned, we can use the differences regression</li><li>• ...</li></ul>									
Common trend assumption	Treatment and control group have the same overall trend									
Experiments	<ul style="list-style-type: none"><li>• <i>Lab experiment</i>: create a situation with desired conditions, manipulate some variables while controlling others, examine the dependent variable</li><li>• <i>Field experiment</i>: research study in a natural setting, manipulate some variables, examine the dependent variable</li></ul> <table><tr><td></td><td>Randomized experiment</td><td>Quasi-experiment</td></tr><tr><td>Field</td><td>High internal validity/ High external validity</td><td>Low internal validity / High external validity</td></tr><tr><td>Lab</td><td>High internal validity/ Low external validity</td><td>Low internal validity/ Low external validity</td></tr></table>		Randomized experiment	Quasi-experiment	Field	High internal validity/ High external validity	Low internal validity / High external validity	Lab	High internal validity/ Low external validity	Low internal validity/ Low external validity
	Randomized experiment	Quasi-experiment								
Field	High internal validity/ High external validity	Low internal validity / High external validity								
Lab	High internal validity/ Low external validity	Low internal validity/ Low external validity								
CRISP Data Mining Process	<ol style="list-style-type: none"><li>1. <i>Business understanding</i></li><li>2. <i>Data evaluation</i></li><li>3. <i>Data preparation</i></li><li>4. <i>Modeling</i></li><li>5. <i>Evaluation</i></li><li>6. <i>Deployment</i></li></ol>									
Bias-Variance Tradeoff	<p>The bias-variance tradeoff provides a conceptual framework for determining a good model</p> <ul style="list-style-type: none"><li>• Models with too few parameters are inaccurate because of a large bias (not enough flexibility) → <b>Underfitting</b>: model is too simple to represent all relevant characteristics (high bias, low variance)</li><li>• models with too many parameters are inaccurate because of a large variance (too much sensitivity to the sample) → <b>Overfitting</b>: model is too complex and fits irrelevant characteristics/noise (low bias, high variance)</li></ul>									
Model selection	<p>Estimating performances of different models to choose the best one (minimum of the test error)</p> <ul style="list-style-type: none"><li>• <i>Akaike Information Criterion (AIC)</i></li><li>• <i>Minimum description length</i></li><li>• <i>Resampling methods</i> (cross validation, jackknife, bootstrap, etc...<ul style="list-style-type: none"><li>○ Holdout procedure: reserve some data for testing</li><li>○ Stratified holdout: guarantee that class are proportionally represented in the test and training set</li><li>○ Repeated holdout: randomly select holdout set several times and average the error rate estimates</li></ul></li></ul>									

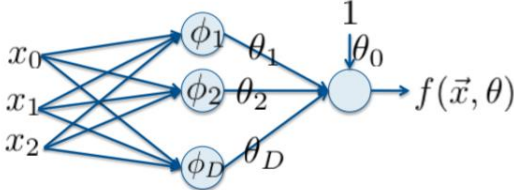
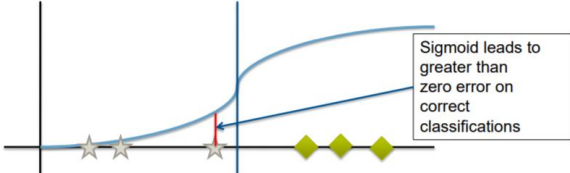

<b>Resampling methods</b>	<ul style="list-style-type: none"> <li>• <i>Cross validation</i>: Train most of the data and use remaining for testing (use all parts of the data once for testing) → k-fold cross validation: divide data in k partition and use 1 for testing (estimate error rates)</li> <li>• <i>Jackknife</i> (Leave-One Out Holdout): use all but one instance for training. Each iteration is evaluated by predicting the omitted instance</li> <li>• <i>Bootstrap</i>: sampling several times with replacement from training set to form a “bootstrap” data set. Some observations are considered more than once and others not at all. Prediction are made for original training set (process repeated many times)</li> </ul>
<b>Model assessment</b>	Having chosen a model, estimating the prediction error on new data
<b>Measuring errors</b>	<ul style="list-style-type: none"> <li>• <math>Error\ rate = (False\ negative\ (FN) + false\ positive\ (FP)) / N\ (Instances)</math></li> <li>• <math>Recall\ (hit\ rate) = TP / (TP+FN)</math></li> <li>• <math>Precision = TP / (TP+FP)</math></li> <li>• <math>Specificity = TN / (TN+FP)</math></li> <li>• <math>False\ alarm\ rate = FP / (FP+TN)</math></li> </ul>
<b>Cost-sensitive learning</b>	Weighting of instances according to costs (ex. Increase the “no” instances in training --> when testing on the original test data set, there will be fewer false positives)
<b>Gain Curve</b>	Instances are sorted according to their predicted probability (x axis = sample size, y axis = number of positives) <ul style="list-style-type: none"> <li>• <i>Random List</i> (ex. 5% of random list have 5% of targets)</li> <li>• <i>Model-Ranked List</i> (ex. 5% of model ranked list have 20% of targets <math>Gain(5\%) = 20\%</math>)</li> </ul>
<b>Lift Curve</b>	How much times better the model-ranked is in comparison to the random list
<b>ROC Curves</b>	(Receiver Operating Characteristics) Recall vs. False alarm rate: y axis shows percentage of true positives in sample, x axis shows percentage of false positives in sample <ul style="list-style-type: none"> <li>• <i>Jagged curve</i>: one set of test data</li> <li>• <i>Smooth curve</i>: use cross-validation and average</li> </ul>
<b>Real-world comparison studies</b>	<ul style="list-style-type: none"> <li>• <i>Logistic regression</i> (discriminant analysis) - widely used</li> <li>• <i>Decision trees</i> – widely used</li> <li>• <i>K-nearest neighbor</i></li> <li>• <i>Non-parametric statistical methods</i></li> <li>• <i>Neural networks</i></li> </ul>
<b>Algorithmic Information Theory</b>	<ul style="list-style-type: none"> <li>• <i>Kolmogorov complexity</i></li> <li>• <i>Minimum description length principle</i></li> </ul> <p>A good model is a simple model that achieves high accuracy on the given data</p> <p><b>Theory 1:</b> very simple, elegant theory that explains the data almost perfectly (preferable)</p> <p><b>Theory 2:</b> significantly more complex theory that produces the data without mistakes</p>
<b>Kolmogorov complexity</b>	The Kolmogorov complexity (K) of a binary object is the length of the shortest program that generates this object on a universal Turing machine (random strings are not compressible) <u>Kolmogorov complexity is not computable!</u> (it needs approximation)
<b>Minimum description length principle</b>	MDL restricts the set of allowed codes in such a way that it becomes possible (computable) to find the shortest codelength of the data, relative to the allowed codes (model selection criterion) <b>DL</b> = space required to describe a theory + space required to describe the theory's mistakes
<b>Computational learning theory</b>	<ul style="list-style-type: none"> <li>• <i>Probably approximately correct (PAC) learning</i>: only reasonable expectation of a learner is that with high probability it learns a close approximation to the target concept</li> <li>• <i>Vapnik-Chervonenkis (VC) theory</i>: provides a measure of the expressiveness of infinite hypothesis spaces</li> <li>• <i>Bayesian inference</i></li> </ul> <p><b>Concepts:</b></p> <ul style="list-style-type: none"> <li>• <i>Sample complexity</i>: how many training examples are needed for a learner to converge to a successful hypothesis</li> <li>• <i>Computational complexity</i>: how much computational effort is needed for a learner to converge to a successful hypothesis</li> <li>• <i>Mistake bound</i>: how many trainings examples will the learner misclassify before converging to a successful hypothesis</li> </ul>
<b>Ensembles</b>	Combining multiple models <ul style="list-style-type: none"> <li>• <b>(+)</b> often improves predictive performance</li> <li>• <b>(-)</b> usually produces output that is very hard to analyze</li> </ul> <p>Methods:</p> <ul style="list-style-type: none"> <li>• <i>Bagging</i></li> <li>• <i>Random Forests</i></li> </ul>

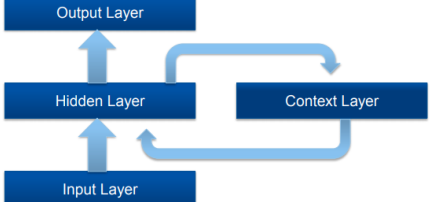
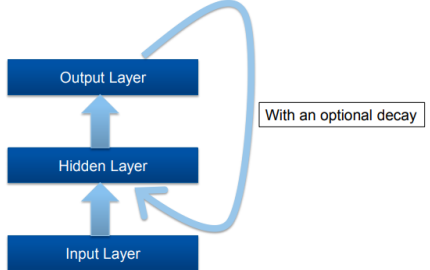
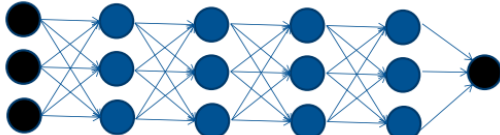


	<ul style="list-style-type: none"> <li>• <i>Boosting</i></li> <li>• <i>Stacking</i></li> </ul>
<b>Bagging</b>	<p>Combining predictions by voting/averaging</p> <ol style="list-style-type: none"> <li>1. Sample several training sets of size <math>n</math> from the population</li> <li>2. Build a classifier for each training set</li> <li>3. Combine the classifiers' predictions <ul style="list-style-type: none"> <li>- If the learning scheme is unstable (small change in the data causes big change in the model) bagging almost always improves the performance</li> </ul> </li> </ol> <ul style="list-style-type: none"> <li>• (+) Can be applied to numeric prediction and classification</li> <li>• (+) Can help a lot if the data is noisy</li> <li>• (+) Can easily be parallelized because ensemble members are created independently</li> </ul>
<b>Bias-Variance decomposition</b>	<p>The bias-variance decomposition is used to analyze how much restriction to a single training set affects performance.</p> <ul style="list-style-type: none"> <li>• <i>Bias</i>: expected error of the ensemble classifier on new data</li> <li>• <i>Variance</i>: component of the expected error due to the particular training set being used to build our classifier</li> <li>• <i>Total expected error</i> = bias + variance</li> </ul> <p>→ Combining multiple classifiers generally decreases the expected error by reducing variance</p>
<b>Random Forests</b>	<p>Random forests randomize data and features (more generally applicable than bagging).</p> <ol style="list-style-type: none"> <li>1. Draw a bootstrap sample from the data</li> <li>2. Grow a "random" tree, where at each node, the best split is chosen among <math>m</math> randomly selected variables. The tree is grown to maximum size and not pruned back</li> <li>3. Store the resulting decision tree</li> <li>4. For each of the decision tree predict class of instance</li> </ol>
<b>Boosting</b>	<p>Boosting tries to minimize the bias in terms of training performance of simple learners (ex. AdaBoost)</p> <p>Boosting needs weight, but you can apply boosting without weights (resample data with probability determined by weights)</p> <ul style="list-style-type: none"> <li>- Boosting implements forward stagewise additive modelling (well-known statistical technique)</li> </ul>
<b>Forward Stagewise additive model</b>	<ol style="list-style-type: none"> <li>1. Build simple regression model</li> <li>2. Gather residuals, learn model predicting residuals, and repeat</li> <li>3. Sum up individual predictions from all regression models</li> </ol>
<b>Additive Regression</b>	<p>Additive regression greedily minimizes squared error of ensemble if base learner minimizes squared error</p>
<b>Stacking</b>	<p>In stacking, the predictions from heterogeneous classifiers are used as input into a meta-learner, which attempts to combine the predictions to create a final best predicted classification</p> <ul style="list-style-type: none"> <li>• <i>Level-0 Models</i>: <ul style="list-style-type: none"> <li>○ Decision tree</li> <li>○ Naïve Bayes</li> <li>○ Instance-based</li> </ul> </li> <li>• <i>Level-1 Model</i> <ul style="list-style-type: none"> <li>○ Meta learner</li> </ul> </li> </ul>
<b>Meta Learning</b>	<ol style="list-style-type: none"> <li>1. Holdout part of the training set</li> <li>2. Use remaining part for training level-0 methods</li> <li>3. Use holdout data to train level-1 learning</li> <li>4. Retrain level-0 algorithms with all the data</li> </ol> <p>If the base learner can output class probabilities, use those as input to meta learner instead of plain classifications</p>
<b>Clustering</b>	<p>Find a natural partitioning of the data set into a number of clusters such that:</p> <ul style="list-style-type: none"> <li>- <i>Intra-cluster similarity is maximized</i> (items in same cluster are similar)</li> <li>- <i>Inter-cluster similarity is minimized</i> (items in different clusters are different)</li> </ul> <p>Clusters are not known a priori.</p> <p>Many different methods and algorithms:</p> <ul style="list-style-type: none"> <li>- <i>Numeric and/or nominal data</i></li> <li>- <i>Deterministic vs. probabilistic</i></li> <li>- <i>Partitional vs. overlapping</i></li> <li>- <i>Hierarchical vs. flat</i></li> </ul>
<b>Hierarchical Clustering</b>	<p><b>Bottom up:</b></p> <ol style="list-style-type: none"> <li>1- Start with single instance clusters</li> <li>2- At each step, join the two closest clusters</li> </ol> <p><b>Top down:</b></p> <ol style="list-style-type: none"> <li>1- Start with one universal cluster</li> </ol>

	2- Find two clusters 3- Proceed recursively on each subset <b>(MST Algorithm: compute the minimal spanning tree of the graph)</b>	
<b>K-Means clustering</b>	K-Means is an example of a partial clustering algorithm 1- Pick a number (k) of cluster centers (at random) → number of clusters determined a priori 2- Assign every item to its nearest cluster center (e.g. using Euclidean distance) 3- Move each cluster center to the mean of its assigned items 4- Repeat steps 2 and 3 until convergence Results can vary significantly depending on the initial choice of seeds. To increase chance of finding global optimum restart with different random seeds	
<b>K-Means Pros and Cons</b>	<b>(+)</b> Simple and understandable <b>(+)</b> Items automatically assigned to clusters	<b>(-)</b> Must pick number of clusters before hand <b>(-)</b> All items forced into clusters (sensitive to outliers)
<b>Probability-based clustering</b>	Model each cluster with a probability distribution (mixture). Each probability distribution gives the probability of an instance being in a given cluster. 1- Start with initial guesses for the parameters 2- Calculate cluster probabilities for each instance (expectation) 3- Re-estimate the distribution parameters from probabilities (maximization) 4- Repeat The maximum found by EM could be a local optimum, so repeat several times with different initial values. How do we know the parameters for the mixture? Use an iterative approach similar in spirit to the k-means algorithm. <b>Extending the model</b> <ul style="list-style-type: none"> <li>- <i>Multiple clusters</i>: use k normal distributions</li> <li>- <i>Multiple attributes</i>: multiply the probabilities of all attributes to get the probability of an instance (in case of correlation among attributes use multivariate normal distribution)</li> <li>- <i>For nominal attributes</i>: create probability distributions for the values</li> </ul>	
<b>Dimensionality problem</b>	Dimensionality reduction is used when the number of variables exceeds the number of observations	
<b>Principal Component Analysis (PCA)</b>	Principal component analysis (PCA) converts a set of possibly correlated variables into a (possibly smaller) set of values of linearly uncorrelated variables called principal components. The principal components are orthogonal (they are the Eigenvectors of the symmetric covariance matrix) → Rotation of the axes → Eigenvalues explain the proportion of variance explained by PC → The PCA score for any of the x is just its coefficient in each of the y PCA Stages <ol style="list-style-type: none"> <li>1. <i>Calculate Zero mean data</i> (calculate the mean of each column and subtract the mean from each value)</li> <li>2. <i>Calculate the variance matrix</i>, which summarizes the relationship between variables (if variables are measured in different units use the correlation matrix)</li> <li>3. <i>Calculate the Eigenvectors and Eigenvalues of the covariance matrix</i></li> <li>4. <i>Order the Eigenvectors by Eigenvalues</i> (from highest to lowest)</li> <li>5. <i>Rotate to get the Eigenvectors as axes</i></li> </ol> You can decide to ignore the components of lesser significance (lose some information) → How many dimensions should remain? Take enough Eigenvalues to cover 80-90% of the total variance (the first k components display as much as possible of the variation in the data) <ul style="list-style-type: none"> <li>- The first PC is direction of maximum variance from the origin, subsequent PCs are orthogonal to first PC and describe maximum residual variance</li> </ul> PCA can be used for example for image compression	
<b>PCA Assumptions</b>	<ul style="list-style-type: none"> <li>- PCA assumes relationships among variables are linear (if the structure of the data is nonlinear, the principal axes will not be an efficient and informative summary of the data)</li> <li>- PCA uses the Euclidean distance among points assuming continuous variables. With discrete variables special techniques are in order)</li> </ul>	
<b>Singular Value Decomposition</b>	SVD of the data matrix can be seen as an alternative technique to compute the same Eigenvectors. $A = USV^T$  <ul style="list-style-type: none"> <li>- The diagonal values of S are called the <i>singular values</i> (it is accustomed to sort them by size)</li> <li>- The Columns of U are called the <i>left singular vectors</i></li> </ul>	

	<ul style="list-style-type: none"> <li>- The columns of <math>V</math> are called the <i>right singular vectors</i> (principal axes)</li> <li>- The columns of <math>US</math> are the <i>principal components</i> (scores)</li> <li>- Singular values of the SVD decomposition of the matrix <math>A</math> is the square root of the Eigenvalues of the matrix <math>(AA^t)</math> or <math>(A^tA)</math></li> </ul>
<b>Principal Component Regression</b>	<p>PCR will tend to do well in cases when the first few principal components are sufficient to capture most of the variation in the predictors as well as the relationship with the response.</p> <ol style="list-style-type: none"> <li>1. PCA to compress the data</li> <li>2. PC regression</li> </ol> <p>In PCR, the number of principal components is typically chosen</p>
<b>Partial Least Squares</b>	<p>Partial Least Squares (PLS) is just like PC Regression except in how the components are computed. PLS makes use of the response <math>y</math> in order to identify new features that do not only approximate the old features well, but also that are related to the response.</p> <p>→ Weights are calculated from the covariance matrix of the predictors</p> <p>→ Weights reflect the covariance structure between predictors and response <math>y</math></p>
<b>Regularization</b>	<p>By regularizing the estimator in some way, its variance will be reduced. If the corresponding increase in bias is small, this will be worthwhile.</p> <ul style="list-style-type: none"> <li>• <i>Subset selection</i> (forward, backward, all subsets)</li> <li>• <i>Ridge regression</i></li> <li>• <i>The lasso</i></li> </ul>
<b>Ridge regression</b>	<p>Ridge coefficient minimize a penalized RSS. This is a biased estimator that for some value of <math>\lambda &gt; 0</math> may have smaller mean squared error than the least squares estimator.</p> <p>Ridge regression estimates will be more biased than the OLS ones but have lower variance</p>
<b>The Lasso</b>	<p>The lasso coefficients minimize the quantity. The lasso has a major advantage over the ridge regression, in that it produces simpler and more interpretable models that involved only a subset of predictors.</p> <p>The lasso can generate more accurate predictions compared to ridge regression.</p> <p>→ Cross-validation can be used in order to determine which approach is better on a particular data set</p>
<b>Ridge-PCA-PLS-Lasso</b>	<ul style="list-style-type: none"> <li>• Ridge regression and PCR outperform PLS in prediction</li> <li>• Lasso outperforms ridge when there are a moderate number of sizable effects, rather than many small effects. It also produces more interpretable models.</li> </ul>
<b>Association Rule Discovery</b>	<p>Aims to discover interesting correlations or other relationships in large databases.</p> <ul style="list-style-type: none"> <li>- <i>Market Basket Analysis</i>: analyze customer buying habits by finding associations and correlations between the different items that customers place in their “shopping basket”</li> </ul> <p><b>Steps:</b></p> <ol style="list-style-type: none"> <li>1- Find all (frequent) itemsets that meet minimum support</li> <li>2- Find all rules that meet minimum confidence</li> <li>3- Prune</li> </ol> <p><b>Subset Property:</b> every subset of a frequent set is frequent</p> <p><b>Apriori Algorithm:</b> use one-item sets to generate two-item sets, two-item sets to generate three-item sets, ...</p>
<b>Support of itemset</b>	<p><math>\text{Supp}(I)</math> is the proportion of transactions that support (contain) <math>I</math></p>
<b>Frequent Itemset</b>	<p>A frequent itemset <math>I</math> is one with at least the minimum support (<math>\text{supp}(i) &gt; \text{minsupp}</math>). Association rules with maximum support and confidence are sometimes called “strong” rules.</p> <p>Frequent itemsets represent sets of items which are positively correlated.</p>
<b>Recommender Systems</b>	<p>Systems for recommending items to users based on examples of their preferences</p> <ul style="list-style-type: none"> <li>- <i>Collaborative filtering</i>: based on similarities among users tastes. For a given user, find other users whose ratings strongly correlate with the current user → recommend items rated highly by these similar users <ul style="list-style-type: none"> <li>o Product associations (ex. 90% of users who like A and B also like C)</li> <li>o User associations (ex. 90% of products liked by A and B are also liked by C)</li> <li>o Combination of product and user associations</li> </ul> </li> <li>(-) Cold start: needs to be enough other users already in the system to find a match</li> <li>(-) Sparsity: hard to find users that have rated the same items</li> <li>(-) First rater: cannot recommend an item that has not been previously rated</li> <li>(-) Popularity bias: tends to recommend popular items</li> <li>- <i>Content-Based filtering</i>: recommendations are based on information on the content of items rather than on other users’ opinions <ul style="list-style-type: none"> <li>(+) Able to recommend to users with unique tastes</li> <li>(+) Able to recommend new and unpopular items</li> </ul> </li> </ul> <p>→ <b>SVD</b>: doesn’t work on sparse matrices, but one can estimate the vectors by minimizing the Euclidian distance between <math>r_{ui}</math> and the dot product for the relevant vectors of <math>U</math> and <math>V</math></p>

<b>Similarity Weigjting</b>	Typically use Pearson correlation coefficient between ratings for active user, a, and another user, u
<b>Significance Weigjting</b>	Include significance weights, $s_{a,u}$ , based on number of co-rated items, n → Important not to trust correlations based on very few co-rated items
<b>Rating prediction</b>	Predict a rating for each item (for active user A) by using the k selected neighbor users
<b>Neural Network</b>	 <ul style="list-style-type: none"> <li>- Edges multiply the signal by a weight</li> <li>- Nodes apply a function, <math>\phi_d</math></li> <li>- Each neuron has its own bias and weights</li> </ul> <p>Combining (Activation) Function – <b>Sigmoid function classification error</b></p>  <p><b>Perceptron Error</b></p>  <p><b>Problems with Neural networks:</b></p> <ul style="list-style-type: none"> <li>• <i>Interpretation of hidden layers</i> (what are the hidden layers doing?) <ul style="list-style-type: none"> <li>◦ Feature extraction (the non-linearities in the feature extraction can make interpretation of hidden layers very difficult --&gt; NN treated as black boxes)</li> </ul> </li> <li>• <i>Overfitting</i></li> </ul>
<b>Multi-Layer Feed-Forward Networks</b>	Multi-layer networks can represent arbitrary functions, but an effective an effective learning algorithm for such networks was thought to be difficult. A typical multi-layer network consists of an input, hidden and output layer, each fully connected to the next, with activation feeding forward. The weights determine the overall function computed <b>Feed-Forward:</b> predictions are fed forward through the network to classify --> The output from one layer is the input to the next (each layer has its own sets of weights)
<b>Cost Function (NN)</b>	Use a cost function to compute the average (misclassification) cost for all training data (for a particular set of weights and biases it computes a single number) --> <b>Minimize the cost function!</b> (-) There is no “closed-form” solution (+) Convex (for any pair of points within a region, the line is in the region)
<b>Backpropagation</b>	Backpropagation is a method used in artificial neural networks to calculate a gradient that is needed in the calculation of the weights to be used in the network. It describes how a single training example, starting from the output neurons, determines the goal for the neurons on the next layer and steps backwards recursively <ul style="list-style-type: none"> <li>• <i>Epoch:</i> one pass through the training set, with an adjustment to the network weights for each training example. (perform as many epochs as needed to reduce the classification error)</li> <li>• <i>Loss function:</i> function that maps values of one or more variables onto a real number intuitively representing the cost associated with the event</li> <li>• <i>Risk function R:</i> expectation of the loss function</li> </ul>
<b>Gradient descent</b>	<ol style="list-style-type: none"> <li>1. Gradient points in the direction of fastest increase</li> <li>2. To minimize R, move in the opposite direction</li> <li>3. Nearly guaranteed to converge to the minimum</li> </ol> <p>→ Can only find a local minimum unless the function is convex → Can Oscillate if the steps are too large (stall if derivate is ever 0 not at the minimum)</p>
<b>Linear Regression Neural Networks</b>	The product of two linear transformations is itself a linear transformation (nothing special) → Non-linearities to identify complex regions in space
<b>Linearly separability</b>	Two classes of points are linearly separable, if there exists a line such that all the points of the class fall on one side of the line, and all the points of the other class fall on the other side of the line

	<p>If <math>f(x)</math> is linear, the NN can only draw straight decision boundaries  → Use the non-linear, differentiable sigmoidal “logistic” function <math>f(x)</math>, which can draw complex boundaries</p>
<b>Hidden nodes</b>	<p>The number of nodes in the hidden layer affects generality and convergence</p> <ul style="list-style-type: none"> <li>• Too few hidden nodes → convergence may fail</li> <li>• Few but not too few nodes → possibly slow convergence but good generalization</li> <li>• Too many nodes → rapid convergence, but “overfitting” happens</li> </ul>
<b>Error Backpropagation</b>	<p>Error backpropagation unravels the multivariate chain rule and solves the gradient for each partial component separately. The target values for each layer come from the next layer  → Feeds the errors back along the network</p>
<b>Recurrent Neural Networks</b>	<p>Output or hidden layer information is stored in a context or memory layer  → Learn patterns in time series and sequential learning tasks</p>  <p><b>Time Delayed Recurrent Neural Networks (TDRNN):</b> Output layer from time <math>t</math> are used as inputs to the hidden layer at time <math>t+1</math></p> 
<b>Deep Neural Networks</b>	<p>Many-layer neural network architectures capable of learning the true underlying features and “feature logic”, and therefore generalize very well. Regularization is used to combat overfitting.</p>  <p>There is no universally agreed threshold of depth dividing shallow learning from deep learning (most researchers agree that deep learning has multiple nonlinear layers)</p>
<b>Problems with Neural Networks</b>	<ul style="list-style-type: none"> <li>• <i>Lack of transparency</i> (where is the knowledge?)</li> <li>• <i>Difficulty in predicting convergence</i></li> <li>• <i>Difficulty in scaling up</i> (NNs are often useful subsystems, but highly complex systems must be carefully structured into separately trainable subsystems)</li> </ul>



Point Estimate	$\bar{X} = \frac{\sum X}{n}$
Estimate of variability in population	$s = \sqrt{\frac{1}{n-1} \sum_i (X_i - \bar{X})^2}$
True Standard deviation in sample mean	$SD = \sigma / \sqrt{n}$
Standard error of sample mean	$SE = s / \sqrt{n}$
95% Confidence interval	Depends on degrees of freedom
Welch Test (2 Samples, independent)	$t_0 = \frac{\bar{x}_2 - \bar{x}_1 - \mu_0}{s_{\bar{x}-\bar{y}}}$
Paired t-test (2 Samples, dependent)	$t_0 = \frac{\bar{d} - \mu_0}{s_d} \sqrt{n}$
Z-test / Gauss test (1 Sample, $\sigma$ known)	$Z = \frac{\bar{X} - \mu_0}{\sigma / \sqrt{n}}$ Rejection region: $\mu \neq \mu_0 \quad  z  \geq z_{\alpha/2}$ $\mu > \mu_0 \quad z \geq z_{\alpha}$ $\mu < \mu_0 \quad z \leq -z_{\alpha}$
T-test (1 Sample, $\sigma$ unknown)	$t = \frac{\bar{d} - \Delta_0}{s / \sqrt{n}}$ Rejection region: $\mu_d \neq \Delta_0 \quad  t  \geq t_{\alpha/2, n-1}$ $\mu_d > \Delta_0 \quad t \geq t_{\alpha, n-1}$ $\mu_d < \Delta_0 \quad t \leq -t_{\alpha, n-1}$
Formulas for coefficients	$\beta_1 = \frac{\sum_i^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_i^n (x_i - \bar{x})^2} = \frac{\frac{1}{n} \sum_i^n x_i y_i - \bar{x} \bar{y}}{\frac{1}{n} \sum_i^n x_i^2 - \bar{x}^2}$
Residual Sum of Squares (RSS)	$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$
Total Deviation	$\sum (y - \bar{y})^2 = \sum (\hat{y} - \bar{y})^2 + \sum (y - \hat{y})^2$ Total deviation (TSS) = explained deviation (ESS) + unexplained deviation (RSS)
R <sup>2</sup> (Coefficient of Determination)	$R^2 = \frac{TSS - RSS}{TSS} = 1 - \frac{RSS}{TSS} = \frac{ESS}{TSS}$ $R^2 = 1 \rightarrow$ Perfect match between the line and data points $R^2 = 0 \rightarrow$ There is no linear relationship between x and y
Multiple linear regression model	$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 = \mathbf{X} \hat{\beta}$
Variance Inflation Factor (VIF)	$VIF = \frac{1}{1 - R_k^2}$
Durbin-Watson statistic (DW)	$DW = \frac{\sum_{i=2}^n (e_i - e_{i-1})^2}{\sum_{i=1}^n e_i^2}$ <ul style="list-style-type: none"> <li><math>DW = 2</math> – no autocorrelation</li> <li><math>DW = 0</math> – perfect positive autocorrelation</li> <li><math>DW = 4</math> – perfect negative autocorrelation</li> </ul>
Logit (log odds)	$\ln \left( \frac{p(x)}{1-p(x)} \right)$ $p = 0.50$ , then logit = 0 $p = 0.70$ , then logit = 0.84 $p = 0.30$ , then logit = -0.84

Odds of success	$\frac{p(X)}{1 - p(X)} = e^{\beta_0 + \beta_1 X}$
Likelihood function for the logit model	$\Pr(Y_i = 1) = F(\beta_0 + \beta_1 X_{1i}) = \frac{e^{(\beta_0 + \beta_1 X_{1i})}}{1 + e^{(\beta_0 + \beta_1 X_{1i})}} = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_{1i})}}$
Likelihood Ratio test	$D = -2 \ln \left( \frac{L(\text{null})}{L(\text{fitted})} \right) = -2 (LL(\text{null}) - LL(\text{fitted}))$
McFadden $R^2$	$R_{McFadden}^2 = 1 - \frac{LL(\text{fitted})}{LL(\text{null})}$
Conditional Probability (Bayes Rule)	$\Pr[h e] = \frac{\Pr[e h]\Pr[h]}{\Pr[e]}$
Normalization	$\Pr[h e] = \frac{\Pr[e_1 h]\Pr[e_2 h] \dots \Pr[e_n h]\Pr[h]}{\Pr[e]}$
Density function $f(x)$	$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$
Entropy( $p_1, p_2, \dots, p_n$ )	$\text{entropy}(p_1, p_2, \dots, p_n) = -p_1 \log_2 p_1 - p_2 \log_2 p_2 \dots - p_n \log_2 p_n$
Gini Index	$Gini(S) = 1 - P^2 - N^2 \in [0, 0.5]$ with $P = p / (p + n)$ $N = n / (p + n)$
C4.5 Method	$e = p = \left( f + \frac{z^2}{2n} + z^* \sqrt{\frac{f}{n} - \frac{f^2}{n} + \frac{z^2}{4n^2}} \right) / \left( 1 + \frac{z^2}{n} \right)$ $f$ is the error on the training data $n$ is the number of instances covered by the node
Euclidian distance	$d_2(x, y) = \left( \sum_{j=1}^p  x_j - y_j ^2 \right)^{1/2}$
Manhattan distance	$d(x, y) = \left( \sum_{j=1}^p  x_j - y_j  \right)$
Correlation Matrix	$\rho = \text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$
Reconstruction of Original data with one Eigenvector	$X \approx \text{PCA Scores} * \text{Eigenvectors} + \text{original mean}$
Similarity Weighting	$C_{a,u} = \frac{\text{cov}(r_a, r_u)}{\sigma_{r_a} \sigma_{r_u}}$
Significance Weighting	$W_{a,u} = S_{a,u} C_{a,u}$ $S_{a,u} = \begin{cases} 1 & \text{if } n > 50 \\ \frac{n}{50} & \text{if } n \leq 50 \end{cases}$
Rating prediction	$p_{a,i} = \bar{r}_a + \sum_{u=1}^k \frac{w_{a,u}(r_{u,i} - \bar{r}_u)}{\sum_{u=1}^k  w_{a,u} }$
Neural network	$f(\vec{x}, \theta) = \sum_{n=0}^{N-1} \sum_{d=1}^D \theta_d \phi_d(x_n) + \theta_0$

**Logistic Neuron Optimization  
(gradient descent)**

$$R(\theta) = \frac{1}{2N} \sum_{i=0}^{N-1} (t_i - g(\theta^T x_i))^2$$

derivative of  $f(z)^2 \Rightarrow 2f(z)f'(z)$  (chain rule)

$$\nabla_{\theta} R = \frac{1}{2N} \sum_{i=0}^{N-1} 2(t_i - g(\theta^T x_i))(-1)g'(\theta^T x_i)x_i = 0$$