

Data types: categorical: nominal(no order), ordinal(order); numerical/quantitative: continuous(Interval: no fixed 0, ratio: fixed 0), discrete(only certain values possible);

Tables: **Reference table**(all data), **Demonstration table**(illustrate a point);

Location statistics: mean, median, mode(may be multiple), percentiles( $L_P = 1 + \frac{P}{100}(n - 1)$ , let  $l$  and  $h$  be the observations at  $[L_P]$  and  $[L_P]$ , then  $p^{th} \text{ percentile} = l + (L_P - [L_P])(h - l)$ );

Scale statistics: **Range**(max – min), **IQR**(Q3– Q1), **Variance**( $\frac{1}{n-1} \sum_i (x_i - \bar{x})^2$ ), **Std**( $\sqrt{\text{variance}}$ ), **MAD**(median of  $|x_i - \text{median}|$ );

**z-score**( $\frac{x - \bar{x}}{\sigma}$ , number of std away from mean, above 2.5 = outlier);

Association statistics: **Covariance**( $\frac{1}{n-1} \sum_i (x_i - \bar{x})(y_i - \bar{y})$ ), **Correlation**( $\frac{\text{covariance}}{\sigma_x \sigma_y}$  or  $\frac{1}{n-1} \sum_i z_{x_i} z_{y_i}$ );

**Strip plot**(data: 1 numerical, task: (cluster, outlier), small datasets); **KDE**(create: choose kernel function + bandwidth => kernel for every datapoint => sum kernels, area: per points  $1/n$ , cannot see maximal and minimal); **Skew**: symmetric, right-skewed(positive) = long right tail, left-skewed(negative) = long left tail;

**Empirical Cumulative Distribution Function**(returns fraction of observations  $\leq x$ )

Search(location, target): lookup(1, 1), browse(1, 0), locate(0, 1), explore(0, 0);

Targets(all data: (trends/patterns, outliers, features), attributes only: (single: (distribution, extremes), many: (dependency, correlation, similarity)));

Marks and Channels(marks: geometric primitive(points, lines, areas, ...), channels: appearance of marks(position, colour(colour-blindness, use blue & orange), size, ...));

Design Strategies(position before colour, natural order: (position, length thickness, area, brightness, saturation), no order: (shape, line style, hue));

Idioms: **Scatterplot**(data: 2 numerical, mark: point, channels: horizontal & vertical position, task: (relations, outliers, trends, clusters)), **Bar chart**(data: (1 categorical, 1 numerical), mark: thick line, channels: (length for value, region per key), task: lookup and compare), **Stacked bar chart**(data: (2 categorical, 1 numerical), mark: vertical stack of lines, channels: (length and hue, region per key), task: (compare lowest category and totals, lookup, part-to-whole, trends)), **Normalized stacked bar chart**(same as previous but every bar = 100%), **Line chart**(data: (2 numerical or 1 numerical and 1 ordinal-categorical), mark: points and connecting lines, channels: (lengths for quantitative values, ordered by key into horizontal regions), task: (trends, relationship)), **Histogram**(data: 1 numerical, task: (distribution, counts), recommended bins =  $\sqrt{n}$ ), **Cumulative Histogram**(data: 1 numerical, task: (distribution, thresholds), height of a bin is sum observations less), **Heatmap**(data: (2 categorical-usually ordinal, 1 numerical) or (3 numerical : 2 binned), mark: (2d matrix, indexed by 2 key), channels: colour by quantity, task: (clusters, outliers, trends)), **Box plot**(data: (1 numerical, ?1 categorical if multiple plots), box: Q1+Q3 + median, outliers: dots for above Q3+1.5IQR & below Q1-1.5IQR, endpoints of whiskers: max and min non-outlier), **Violin plot**(box + kde), **Pie chart**(data: (1 categorical + derived =  $\text{angle} = \frac{\#category}{\#total} 360^\circ$ ) or (1 categorical, 1 numerical), mark: area, channel: angle + colour, layout: radial, scale: at most 12 categories, task: lookup and compare part-to-whole), **Scatterplot matrix**(all possible pairs of axes, by category, scatterplots), **Parallel coordinates**(parallel axes: line representing each item, rectilinear axes: item as point);

Effectiveness(clearly indicates relationship of values, represents quantities accurately, easy comparison, easy observation of ranked order of values, obvious how information should be used);

DMM categories(supervised = predefined target / labelled, global = all the data, local = some data);

Linear regression(kind: (supervised, global) relate  $y$  to  $x$  by  $y = ax + b$ , SSD = Sum of Squared Deviations =  $\sum \text{residuals}^2$ ,  $R^2 = 1 - \frac{SSD}{(n-1)\sigma_y^2}$ );

**clustering**(kind: (unsupervised, global), partition into groups, depends on distance => equal treatment of attributes(same units for similar attributes, relevant distance units, standardize units for dissimilar attributes e.g. z-score));

**k-means clustering**(kind of cluster, pick  $k$  centroids => assign points to nearest centroid => recompute mean of clusters as centroid => repeat till stable);

Distances: **Manhattan**(add up horizontal and vertical distance), **Network**(sum of lengths of edges, good if sparse network and we know all possible movements), **Euclidean**( $\sqrt{x^2 + y^2}$ );

Decision tree(kind: (supervised, global), separate into 'positive' and 'negative' cases);

Evaluation: **Confusion matrix**, **Accuracy** =  $\frac{TP+TN}{TP+TN+FP+FN}$ , **Precision** =  $\frac{TP}{TP+FP}$ , **Recall/Sensitivity** =  $\frac{TP}{TP+FN}$ , **Specificity** =  $\frac{TN}{TN+FP}$

Association rule(kind: (unsupervised, local), find high-confidence associations between subsets / frequently occurring associations,

$\text{support}(X) = \frac{|X|}{\text{total}}$ ,  $\text{confidence}(X \Rightarrow Y) = \frac{|X \cap Y|}{|X|}$ , creation: find frequent itemsets  $\text{length } 1: A, \text{length } 2: A \cap B, \text{length } 3: A \cap B \cap C$  until no sufficient support left then replace  $\cap$  with  $\Rightarrow$ , e.g.  $A \cap B \Rightarrow C$ );

Data Modelling(redundancy: same data too many times, leads to inconsistency, **primary key**: minimal set of attributes in table identifying each row), **Logical Schema**(data modal, logical structure), **Instance**(actual content of database, relation instance = table);

Database Design: **E-R Model**(entity = unique object, entities have attributes, entity set = set of entities that share properties, relationship = association of entities, relationship set = collection of relationship on entity sets, relationship sets can have attributes), **E-R Diagrams syntax**(rectangle = entity set, diamond = relationship set, lines = entity set so relationship set, dotted/dashed line =

rectangle contains attribute set of relationship set, attributes inside rectangle, underline = primary key, simple line = many, arrow = one, double line = all ex: (one to one:  $\leftarrow \text{rel} \rightarrow$ , many to one:  $\leftarrow \text{rel} \dashv$ , many to many:  $\dashv \text{rel} \dashv$ , all to one:  $\dashv \text{rel} =$ );

**SQL**: **Order**(SELECT => FROM => WHERE => GROUP BY => HAVING), **LIKE**(string comparison, %: 0 or more characters), **BOOL**(AND, OR, NOT), **AS**(renames: old-name AS new-name, optional: borrower b = borrower AS b), **Set operators**(UNION, EXCEPT, INTERSET; used to join queries, eliminates duplicates append ALL to retain everything), **Aggregate functions**(COUNT: number of values, MIN: minimum, MAX: maximum, AVG: average, SUM: sum)

**Kinds of data**: **Primary data**(collected by me), **Secondary data**(collected by others), **Features**(measurable properties or characteristics);

**Scientific method**: **Deductive reasoning**(If the premises are true, the conclusion is valid), **Inductive reasoning**(inference from particular to general case, from finite sample to whole population), **Occam's Razor/Parsimony**(prefer the simpler explanation), **Validity**(accurate to the real world, inner: conclusions in the study valid, external: can the conclusions be applied beyond the context of the study), **Reliability**(similar conditions => similar results), **Reproducibility**(ability to replicate findings), **Precision**(errors by the measuring instrument i.e.  $\pm 0.5\text{mm}$ ), **Accuracy**(difference from reality), **Errors**(sources: (measurement process, definition, inadequacy of technology), random: no pattern, systematic: consistent e.g. offset);

**Data sampling**: **Population**(complete set), **Sample**(part of the set, biased if some part is overrepresented), **Convenience sampling**(data that is easier to collect, advantages: (time, effort, money, ...), disadvantages: possible bias threat to external validity), **Random sampling**(everyone equally likely to be included, ignoring knowledge about the population), **Stratified random sampling**(strata: disjoint parts forming the population, proportionate stratified random sampling: random sample from every stratum in equal proportion to the proportion of this stratum in the population, disproportionate stratified random sampling: want to overrepresent a stratum), **Voluntary sampling**(individuals select themselves, self-selection bias is difficult to measure);

**Data Cleaning and Preprocessing**(diagnosing and editing faulty data): **Sources**(equipment or transmission, collection circumstances, manual entry, non-optimal collection protocol), **Handling**(discard: potentially lose a lot of data, introduce bias), impute: (fill in constant, data mining or estimate), do nothing), **Noise reduction in time series data**(e.g. smoothing kernel);

**Filtering time series**(choose window size): **Median**(replace middle by median of values in window), **Mean**(replace middle by mean of values in window, far and near values have same influence), **Gaussian**(choose width  $\omega$ , gaussian kernel  $T$  for  $\omega$ , choose pair  $(t_n, x_n)$ , assign weights  $w_i = T(t_n - t_i, \omega)$  to each value  $x_i$ , further away = lower weight,  $\sum_i^n w_i w_{n-i}$ ,  $\sum w = 1$ ), **Convolution**(example: (gaussian, mean), sum weights = 1), **Differentiation**( $[w_1, w_0, w_1] \Rightarrow [-w_1, -w_0, 0, w_0, w_1]$ );

**Variables/features**: **Independent**(see how it's change impacts dependent value), **Confounding**(variable that's not considered),

**Differentiation**( $p'(t_i) = \frac{p_{i+1} - p_{i-1}}{t_{i+1} - t_{i-1}}$ );

**Distributions**: **Discrete uniform**( $P(X = k) = \frac{1}{n}$ ), **Binomial**( $\binom{n}{k} = \frac{n!}{k!(n-k)!}$ ,  $P(X = k) = \binom{n}{k} p^k (1-p)^{n-k}$ ,  $P(X \leq k) = \sum_k^l P(X = k)$ ),

**Geometric**( $P(X = k) = (1-p)^{k-1} p$ ), **Poisson**(rare events), **Expectation**( $\sum_k k P(X = k)$ , binomial:  $np$ ), **Variance**( $\sum_k (k - E(X))^2 P(X = k)$ , binomial:  $np(1-p)$ ), **Continuous probability**(density function,  $P(X=x)=0$ ,  $P(a \leq X \leq b) = \int_a^b f(x)dx$ , (E, Var, STD replace sum with infinity integral)), **Normal**( $N(\mu, \sigma^2)$ , continuous probability, symmetric around mean  $\mu$ , if  $X \sim N(\mu, \sigma^2)$  then z-scores  $\sim N(0, 1)$ , normal quantile  $z_\alpha$ : value such that  $Z \sim N(0,1)$ ,  $P(Z \leq z_\alpha) = \alpha$ );

**Estimations**: **Central limit theorem**(for large  $n$  the distribution of  $\bar{X}$  can be approximated by  $N(E(X), \frac{\sigma_x^2}{n})$ , for large  $n$  the mean is

normally distributed), **Binomial**(estimate  $\hat{p} = \frac{\#success}{n}$ ,  $E(\hat{p}) = p$ ,  $Var(\hat{p}) = \frac{p(1-p)}{n}$ , central limit theorem applies);

**Confidence intervals**(interval containing true unknown value): significance level =  $\alpha$ , confidence level =  $1-\alpha$ , **Two sided**( $\bar{X} - z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$ ), **One sided**(right:  $\mu \leq \bar{X} + z_{\alpha} \frac{\sigma}{\sqrt{n}}$  left:  $\bar{X} - z_{\alpha} \frac{\sigma}{\sqrt{n}} \leq \mu$ , if sigma unknown use sample std, and studentized  $t_{n-1}$  instead

of z), **Proportion**(same but replace mean with estimate  $\hat{p}$ ), **Difference**( $\mu_1 - \mu_2 \pm z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$ ,  $\hat{p}_1 - \hat{p}_2 \pm z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$ );

**Hypothesis testing**: **Procedure**(choose null hypothesis and alternative hypothesis, if enough evidence => reject null hypothesis) **P-value**(probability of observing a more extreme test statistic in the direction of the alternate hypothesis than observed ( $P(X \geq k)$ ), if p-value <  $\alpha$  => reject), **Critical value**(choose test statistic e.g. mean, compute confidence interval based on  $\alpha$  and hypothesis, if test statistic outside confidence interval =>  $H_0$  rejected);

**Error types**: **Type I**(erroneously rejecting  $H_0$ ), **Type II**(erroneously not rejecting  $H_0$   $\beta$ ), **Power of the test**(correctly rejecting  $H_0$ ,  $1-\beta$ );

**Assumptions-hypothesis tests**: independent observations, come from same distribution, means(normality or "large sample size"), proportions( $np > 5$ ,  $n(1-p) > 5$ ), Normality testing(graphical, goodness-of-fit test e.g. Anderson-Darling(gives p-value, small => not normal), kde, QQplot);

**Assumptions - Linear regression**(diagnostic): expectation is a linear function(scatter plot), additive error term epsilon(scatter plot, residual plot), normality of epsilon(QQplot), independence of observations(QQ plot of standardized residuals, kde plot of standardized residuals, Anderson-Darling test on standardized residuals), equal variance for all observations, no clear pattern in residual plot, studentized residuals within (-2.5, 2.5)