

Week 2:
Types of data: Nominal ($=!$): Distinctness. Categorical labels without order. E.g. ID numbers, eye color, zip codes.
Ordinal Scale ($<, >$): Order, no equal intervals (E.g. taste of potato chips on a scale from 1-10), grades, height (tall, medium, short). (Central tendency can be measured by mode or median. The mean and stdev cannot be defined from an ordinal set.)

Interval Scale (+, -): meaningful difference. Order, equal intervals, no true zero. Addition & subtraction. E.g. calendar dates, temperatures in Celsius or Fahrenheit.
Ratio Scale (*, /): meaningful ratios. Order, equal intervals, true zero. E.g. temperature in Kelvin, length, counts, elapsed time multiplication & division (e.g., time to run a race).

Central tendency can be measured by mode, median, or mean.

The mode is the outlier and better represents the central tendency for skewed data. The mode only reflects the most frequently occurring value and does not consider the overall distribution. The mean is sensitive to outliers in the tail of the distribution. Left-Skewed Distribution: The tail extends to the left (negative skew). Relationship: Mean < Median < Mode.

Dispersion can be estimated by the Inter-Quartile Range (IQR), standard deviation. Calculating descriptive statistics: Median is the middle value (average of two middle values). Minimum is the first value. Maximum is the last value. 10th percentile is item at index 0.1*N. 90th percentile is item at index 0.9*N. Range is Maximum minus Minimum. IQR is the difference between the first and third quartile.

How to calculate the IQR:
 i.e.: Q3 - Q1 = 4 - 3 = 1.
 Mean is the sum of values divided by the number of values: $\frac{\sum X_i}{N}$

Variance: $\frac{\sum (X_i - \text{mean})^2}{N-1}$

Standard deviation: $\sqrt{\text{variance}}$
 text data: Requires interpretation, coding or conversion.

Levels of Measurement

	Nominal	Ordinal	Interval	Ratio
Countable	✓	✓	✓	✓
Order defined		✓	✓	✓
Difference defined (addition, subtraction)		✓	✓	✓
Zero defined (multiplication, division)				✓

Measures of Central Tendency

	Nominal	Ordinal	Interval	Ratio
Mode	✓	✓	✓	✓
Median		✓	✓	✓
Mean		✓	✓	✓

Measures of Dispersion

	Nominal	Ordinal	Interval	Ratio
Counts / Distribution	✓	✓	✓	✓
Minimum, Maximum		✓	✓	✓
Range	✓	✓	✓	✓
Percentiles	✓	✓	✓	✓
Standard deviation, Variance		✓	✓	✓

Data acquisition and cleaning:
 Cleaning and Transforming Data
 Why: Real data is often 'dirty'. So: Important to do some data cleaning and transforming first. Typical steps: 1. Type and name conversion. 2. Filtering of missing or inconsistent data. 3. Unifying semantic data representations. 4. Matching entries from different sources. 5. (Later also) Rescaling and optional dimensionality reduction.

Week 3:
Using boxplots
 Mean and stdev are not informative when data is skewed. Box plots summarize data based on 5 numbers:

- Lower inner fence -Q1 - 1.5 * IQR.
- First quartile (Q1) - equivalent to 25th percentile.
- Median (Q2) - equivalent to 50th percentile.
- Third quartile (Q3) - equivalent to 75th percentile.
- Upper outer fence -Q3 + 1.5 * IQR.
- IQR = Q3 - Q1

Values outside fences are outliers. Sometimes include outer fences at Q1-3*IQR and Q3+3*IQR.

Week 4:
Relational databases: In the relational model, a database is a collection of one or more relations. Each relation is a table with rows and columns. Each relation has a schema, which describes the columns, or fields. **Relational Model** defines the structure and relationships between tables using primary and foreign keys.

Definition of relation: A relation is a named, two-dimensional table of data. A relation schema specifies name of relation, and name and data type (domain) of each attribute.

$R = (A_1, A_2, \dots, A_n)$ is a **relation schema**
 • A_1, A_2, \dots, A_n are **attributes**, each having a **data type**

e.g. Student (sid: string, name: string, login: string, add: string, gender: char)
 A relation instance is a set of tuples (a table) for a schema

Some remarks
 Not all tables qualify as a relation:
 - Every relation must have a unique name.
 - Attributes (columns) in tables must have unique names.

• The order of the columns is irrelevant.
 - All tuples in a relation have the same structure

• Constructed from the same set of attributes.
 - Every attribute value is atomic (not multivalued, not composite).

• Every row is unique.
 - Can't have two rows with exactly the same values for all their fields.

• The order of the rows is immaterial.
 RDBMS table extends mathematical relation

- RDBMS allows **duplicate rows**, supports an **order of tuples or attributes**, and allows **null values** for unknown information
SQL – The Structured Query Language
 Supported commands from roughly two categories:
 - **DDL (Data Definition Language)**.
 • Create, drop, or alter the relation schema.
 • Example: `CREATE TABLE name (list_of_columns)`
 - **DML (Data Manipulation Language)**.
 • For retrieval of information also called query language.
 • `INSERT, DELETE, UPDATE`
 • `SELECT ... FROM ... WHERE`

Creating Tables
Creation of tables/relations:

```
CREATE TABLE name ( list-of-columns );
Example: Create the Student table.
CREATE TABLE Student ( sid INTEGER,
                      name VARCHAR(20),
                      login VARCHAR(20),
                      gender CHAR,
                      address VARCHAR(50) );
```

Insertion of new data
Syntax: `INSERT INTO table (list-of-columns) VALUES (list-of-expression)`

Example: `INSERT INTO Students (sid, name) VALUES (53688, "Smith")`

Updating of tuples
Syntax: `UPDATE table SET column = expression WHERE search_condition`

Example: `UPDATE students SET gpa = gpa - 0.1 WHERE gpa >= 3.3`

Deleting of tuples
Syntax: `DELETE FROM table WHERE search_condition`

Example: `DELETE FROM students WHERE name = "Smith"`

Integrity constraints
 A variety of rules to maintain the integrity of data when it is manipulated. The rule must be satisfied for any instance of the database; e.g., domain constraints.

Integrity constraints are declared in the schema. A legal instance of a relation is one that satisfies all specified integrity constraints.

- **Referential integrity** ensures that relationships between tables remain valid and consistent. It enforces rules for foreign keys to maintain data accuracy.

Domain constraints

Domain constraints restrict attributes to valid domains.

- **NULL/NOT NULL**: whether an attribute is allowed to become NULL (unknown).

- **DEFAULT**: to specify a default value.

- **CHECK(condition)**: a Boolean condition that must hold for every tuple in the DB instance.

Example:

```
CREATE TABLE Student (
    sid INTEGER PRIMARY KEY,
    name VARCHAR(20) NOT NULL,
    gender CHAR CHECK (gender IN ('M', 'F', 'T')),
    birthday DATE,
    country VARCHAR(20),
    level INTEGER DEFAULT 1 CHECK (level BETWEEN 1 and 5)
);
```

Primary keys: identifiers of a relation.

Foreign keys: identifiers that enable a dependent relation to refer to its parent relation.

Primary key identifies each tuple of a relation.
Composite Primary Key consisting of more than one attribute.
two primary key as one id

Foreign key is a (set of) attribute(s) in one relation that refers to a tuple in another relation (like a "logical pointer").

Example: Primary & Foreign Keys

```
Student          Enrolled          Unit_of_Study
  sid   name           sid   ucode   semester      ucode   title   credit_pts
  31013  John        31013  I2120  CR          I2120  DB Intro  4
```

```
CREATE TABLE Student (sid INTEGER, CONSTRAINT Student_PK PRIMARY KEY (sid));
CREATE TABLE UoS (ucode CHAR(8), ..., CONSTRAINT UoS_PK PRIMARY KEY (ucode));
CREATE TABLE Enrolled (sid INTEGER, ucode CHAR(8), semester VARCHAR, CONSTRAINT Enrolled_PK FOREIGN KEY (sid) REFERENCES Student, CONSTRAINT Enrolled_PK2 FOREIGN KEY (ucode) REFERENCES UoS, CONSTRAINT Enrolled_PK PRIMARY KEY (sid, ucode));
```

Key constraint: No two distinct tuples can have the same values in all key attributes

Example:

```
CREATE TABLE Enrolled (
    sid INTEGER,
    cid CHAR(8),
    grade CHAR(2),
    PRIMARY KEY (sid),
    PRIMARY KEY (cid, grade) );
    "Students can take only one course and receive a single grade for that course; further, no two students in a course receive the same grade."
```

Foreign Key Constraint
 - For each tuple in the **referring relation** whose foreign key value is 'x', there must be a tuple in the **referred relation** with a candidate key that also has value 'x'.

Keys and Nulls

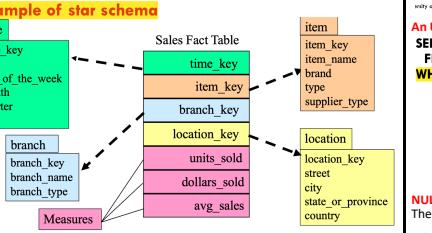
Primary Key: Up to one per table, and must be **unique**. Automatically disallow **NULL values**.

Unique (candidate key): Possibly many candidate keys (specified using **UNIQUE**)

Foreign Key: By default, allows **NULL values**. If there must be a parent tuple, then must combine with **NOT NULL constraint**.

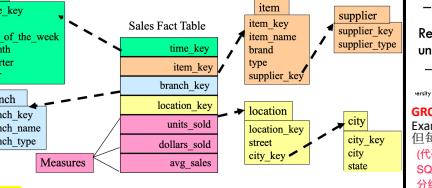
Conceptual modeling of data warehouses
 dimensions & measures instead of relational model. Data warehouse contains a large central table (fact table). Contains the data **without redundancy**. A set of dimension tables.

Star schema



Snowflake schema
 A refinement of a star schema where **some dimensional hierarchy is normalized** into a set of smaller dimension tables, forming a shape similar to a snowflake.

Example of snowflake schema



Week 5:
The SELECT – FROM – WHERE command

Example 1: Which station commence after 1900-1-1?

```
SELECT sitename, commerce, orgcode returned_features
FROM station
WHERE commerce > '1900-1-1';
```

Example 2: How many measurements have we done?

```
SELECT COUNT(*) FROM Measurement;
```

Example 3: List top five measurements ordered by date in descending order.

```
SELECT * FROM Measurement
ORDER BY DATE DESC LIMIT 5;
```

SQL data types

Integers, Floats, Numeric, Date, Timestamp

Strings (CHAR, VARCHAR): CHAR: fixed length; VARCHAR: variable length strings up-to-max length

Comparison operations in SQL

Having clause

Having's作用是将GROUP BY分组之后的每一组数据进行操作

HAVING clause can further filter groups to fulfil a predicate:

```
SELECT uos_code AS unit_of_study, AVG(mark)
FROM Assessment
GROUP BY uos_code;
```

这个条件是对于上述分组结果的过滤。它只保留那些平均成绩大于10的课程单元。

Find the average marks of öcp unit of studies with more than 2 results.

```
SELECT *
FROM TelescopeConfig
WHERE (mindeg BETWEEN -90 AND -50)
AND (maxdeg > -45)
AND (tele_array = 'H168' );
```

Example 1:

```
SELECT *
FROM TelescopeConfig
WHERE tele_array LIKE '%H%';
```

表示查询所有tele_array中包含字母 H 后面为任意多个字符

Comparison operations in SQL

Example 2:

```
SELECT *
FROM TelescopeConfig
WHERE tele_array LIKE '%H%';
```

表示查询tele_array中包含字母 H 后面为任意多个字符

Date and time in SQL

CURRENT_DATE: db system's current date. CURRENT_TIME: db system's current timestamp.

Example:

```
SELECT *
FROM Epoch
WHERE startDate < CURRENT_DATE;
```

Main Operations

- EXTRACT(component FROM date).

• e.g., EXTRACT(year FROM startDate)

- DATE string (Oracle code): TO_DATE(string,template)

• e.g., DATE '2012-03-01'

• Some systems allow templates on how to interpret string.

• Oracle syntax: TO_DATE('01-03-2012', 'DD-Mon-YYYY')

+/- INTERVAL
 - e.g. '2012-04-01' + INTERVAL '36 HOUR'

JOIN: Querying multiple tables

Join 适合两个表中有多个相同列, 而 Neutral join 适合两个表只有一个相同列

Examples:

- Produces the cross-product Station x Organisation:

```
SELECT *
FROM Station, Organisation;
```

- Find the site name, commerce date and organisation name of all stations:

```
SELECT sitename, commerce, organisation
FROM Station, Organisation
WHERE orgcode = code;
```

查询的是选中的三个个别

Controlling study conditions to some degree (active participation of researcher), Strong hypotheses, sample size for desired power and controlled data collection per specified protocols. Establish causality.

Example: randomized control trials.

- 100 subjects. Factor: Average Computer Time. Treatments:

- Control group (computer time: max. 30 minutes). Treatment group (computer time: 2 hours).

- 50 subjects randomly assigned to each treatment. Response: we measure the blood pressure for each group.

Experimental vs observational

Main difference: Most experiments use random assignment, while observational studies do not.

Observational studies typically only establish correlation but not causality

Experimental studies establish causality

Dependent variable: Measure of interest.

Independent variable: Manipulated to observe the effect on dependent variable

Controlled variables: Materials, measurements and methods that don't change.

Research question (Q): Asks whether the independent variable has an effect.

Null hypothesis (H0): The assumption that there is no effect. 零假设是一种假设，指出在研究中没有观察到显著的效果或变化。换句话说，它认为任何变化仅是随机波动而不是受某种因素影响的结果。

Hypothesis testing

We use hypothesis testing to specify whether to accept or reject a claim about a population depending on the evidence provided by a sample of data.

A hypothesis test examines two opposing hypotheses about a population parameter (e.g., the mean):

- The null hypothesis and the alternative hypothesis.

- The null hypothesis represents our initial assumption about the parameter, and we collect evidence to possibly reject the null hypothesis in favour of the alternative hypothesis.

Testing reliability with p-values

Compare to significance level threshold α .

- α is the probability of (wrongly) rejecting H_0 given that it is true (Type I error rate), i.e., false positive. - Commonly use α of 5% or 1%.

values

Accept H_0 | Rejected H_0

H_0 (H_0 is True) | Right Decision

- No difference

H_1 (H_1 is False) | Type I error

- Difference exists

P-value | Indicates

$<\alpha$ | Strong evidence against the null hypothesis

$>\alpha$ | Weak evidence against the null hypothesis

$=\alpha$ | Marginal

Not every test result is correct

$\alpha=0.05$ will erroneously reject H_0 5% of the time

Perform enough tests and you will get a false result (p-hacking): 这是一种数据分析不当调整测试的方法。研究人员可能通过反复试验、选择性报告不同的结果来达到显著性，从而导致得到虚假的结果。

Group by clause

Example: Find company and total amount of sales. 分组后的组可能含有许多行，但每个组中的样本都只有相同的Group by的条件(这里是company)

SQL按照公司名称进行分组

分组。这样的话，查询

结果会为每一家公司生成一行，显示该公司的总销售额。

SELECT Company, SUM(Amount)

FROM Sales

GROUP BY Company

Filtering groups with **HAVING clause**

HAVING 的作用是将GROUP BY分组之后的每一组数据进行操作:

HAVING AVG(mark) > 10 保留那些平均成绩大于10的课程单元。

Tests the null hypothesis that two populations have the same population mean.

方差分析 (ANOVA) 是一种统计方法，用于比较三个或更多组数据的平均值。以确定这些组之间是否存在显著差异。

Assumptions: The samples are independent. - Populations are normally distributed. - Standard deviations are equal (by default).

U test

Nonparametric version of unpaired t-test. 用于比较两个独立样本的中位数差异。它主要适用于数据不是正态分布的情况，或者当样本量较小时，通常被用作替代 t 检验。

- Test the null hypothesis that the distribution underlying sample x is the same as the distribution underlying sample y.

Assumption: The samples are independent.

Analysis of variance (ANOVA)

Tests the null hypothesis two or more groups have the same population mean.

方差分析 (ANOVA) 是一种统计方法，用于比较三个或更多组数据的平均值。以确定这些组之间是否存在显著差异。

Assumptions: The samples are independent. - Populations are normally distributed. - Standard deviations are equal.

Kruskal-Wallis H-test

Nonparametric version of ANOVA.

Test the null hypothesis that the population median of all of the groups are equal. Kruskal-Wallis H 检验是一种无参数统计检验方法，通常用于比较三个或更多独立样本的中位数差异。它是方差分析 (ANOVA) 的一种非参数替代方法。通常用于样本不符合正态分布的情况。

Assumptions: Samples are independent.

Accuracy, precision, recall, f1

Model prediction

Spam (s=1) | Ham (s=0)

Spam (g=1) | FP FN

Ham (g=0) | TP TN

Actual result

Spam (g=1) | TP FN

Ham (g=0) | FP TN

- Accuracy: percentage of correct over all instances.

- Precision: percentage of correct system predictions.

- Recall: percentage of correct gold labels.

- $\bar{P} = \frac{TP}{TP+FN}$

- $\bar{R} = \frac{TP}{TP+FN}$

