# Data Exploration

# Goal of Data Exploration
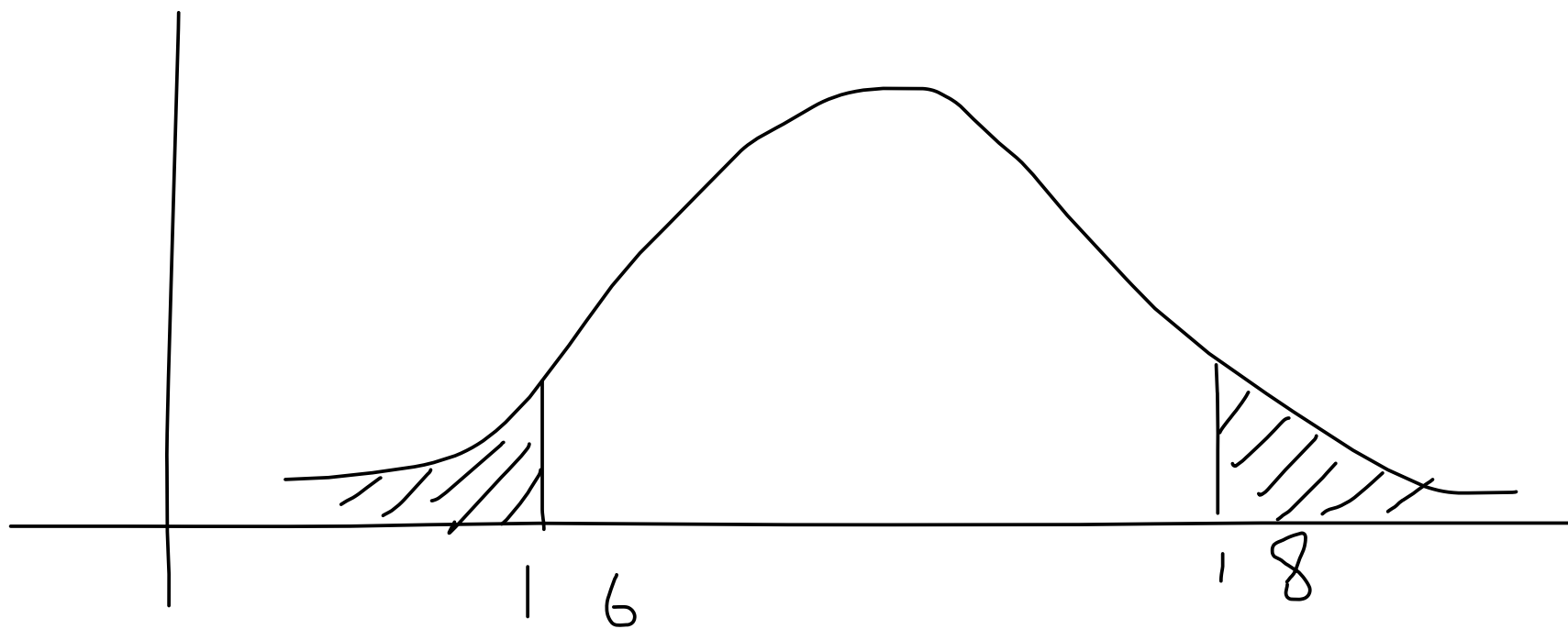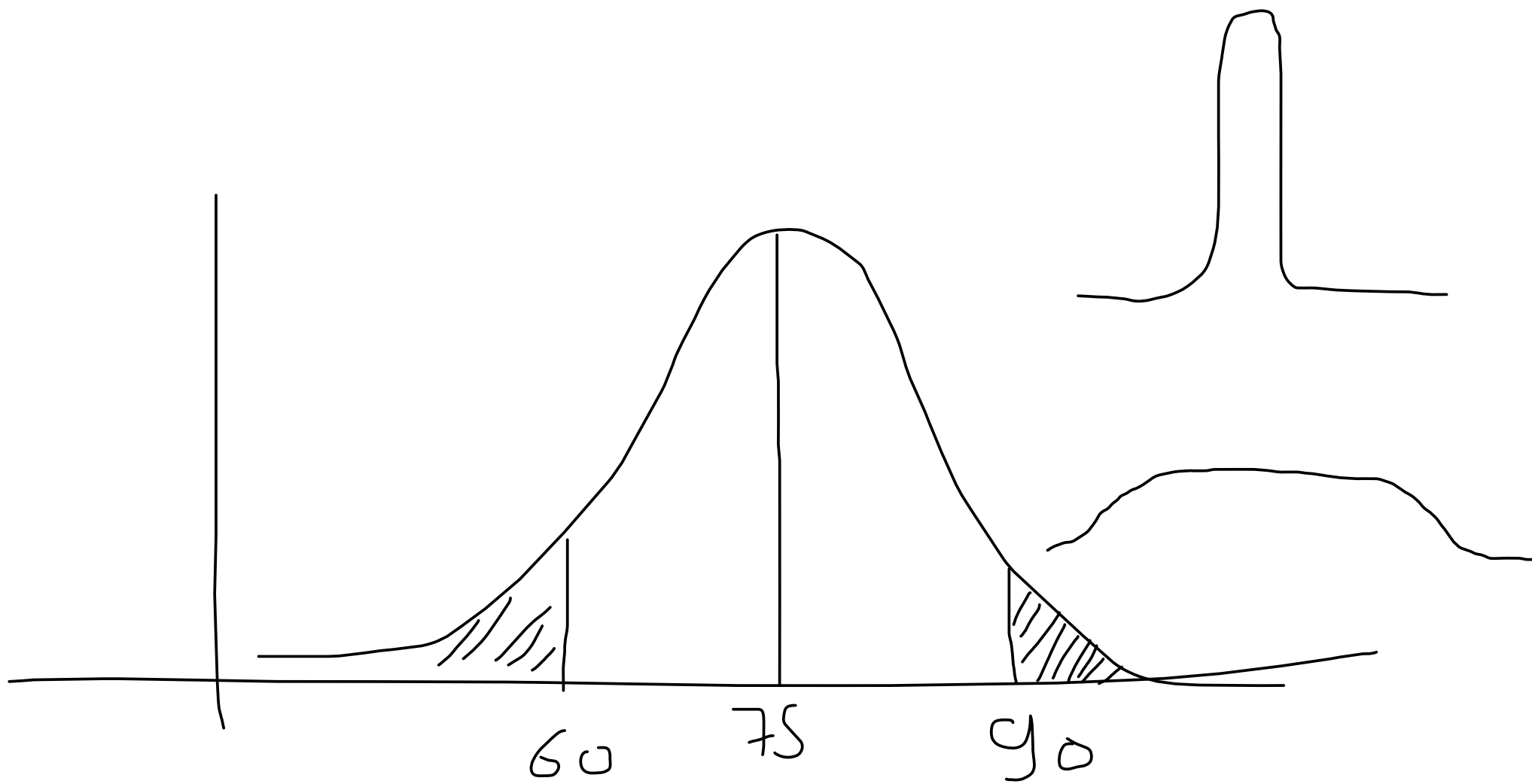
- Goal:
  - Understand the basic characteristics of the data

- Examples for characteristics:
  - Structure
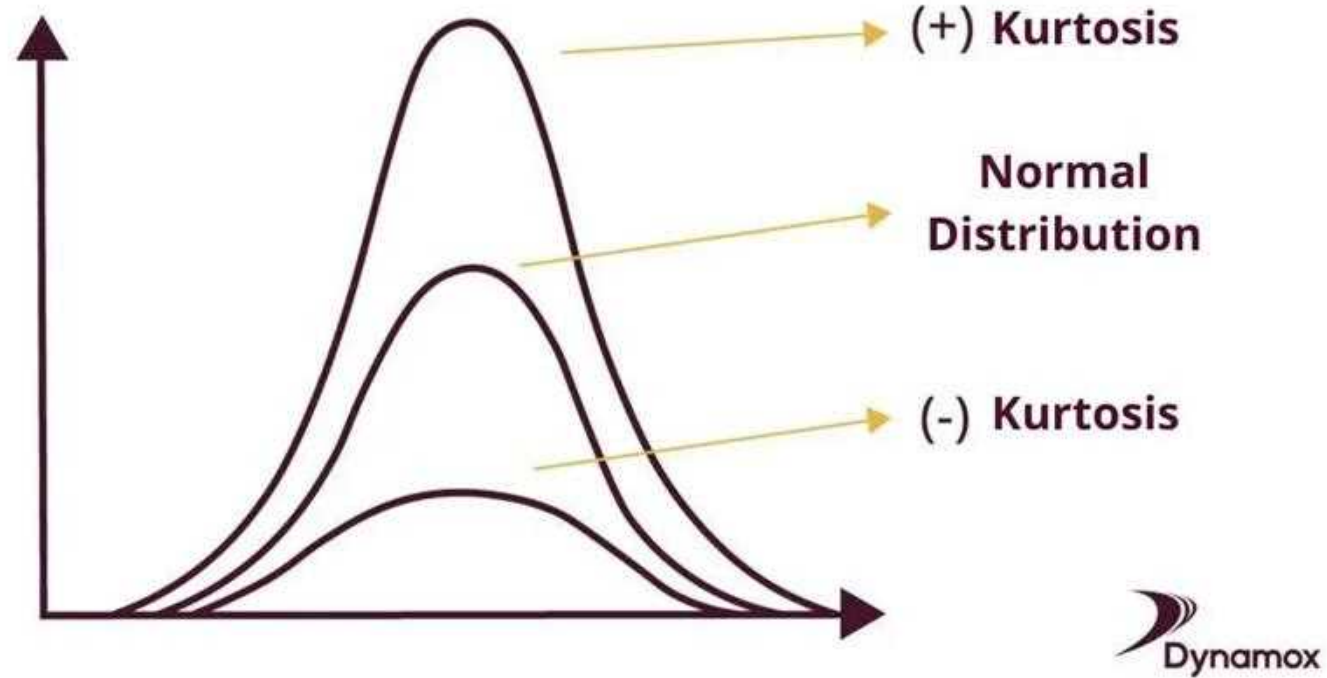  - Size
  - Completeness
  - Relationships

# Descriptive Statistics

- Summarize data through single value

- Common statistics
  - Central tendency (mean/median/mode)
  - Variability (standard deviation,  interquartile range)
  - Range of data (min/max)

- Other important statistics
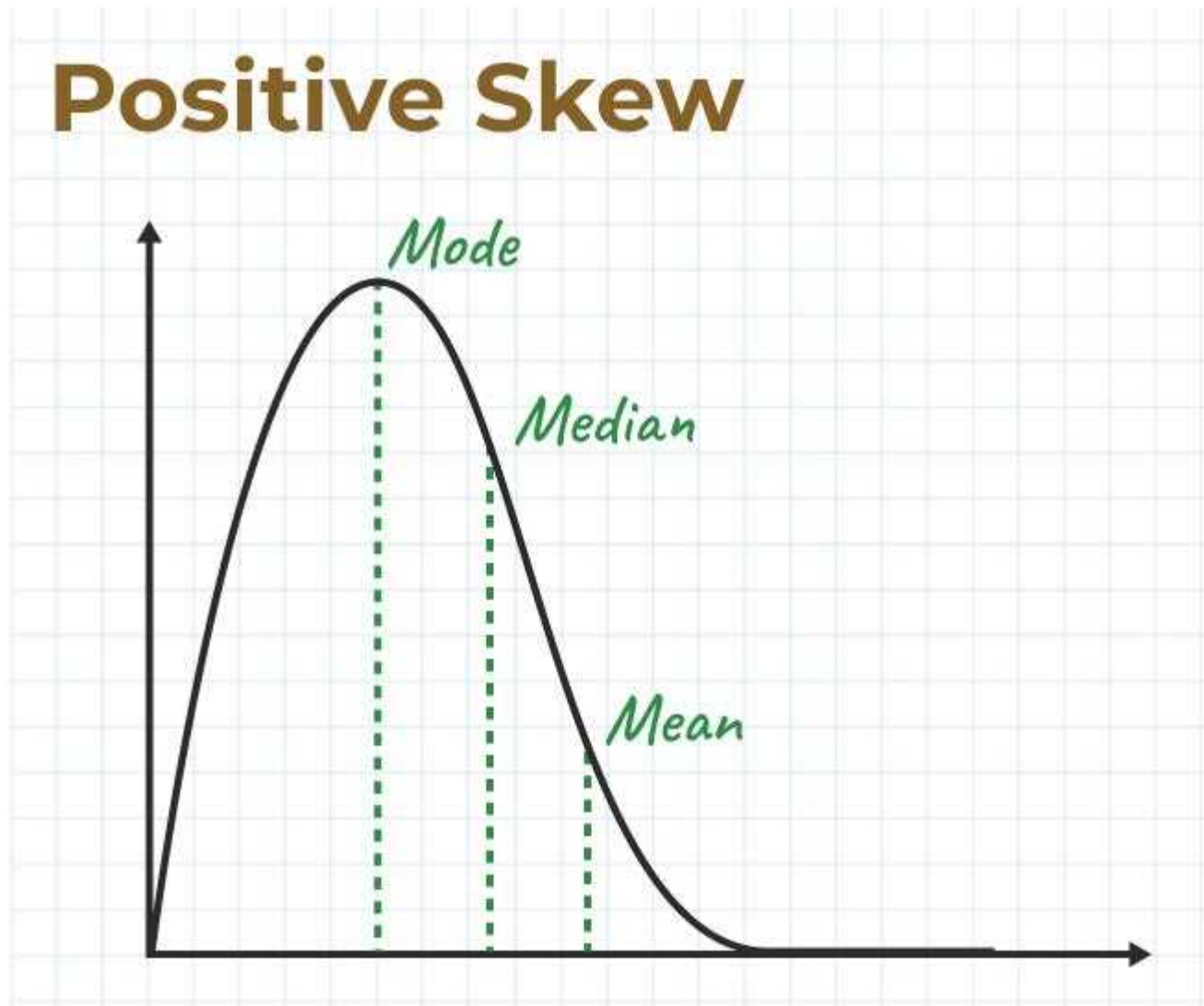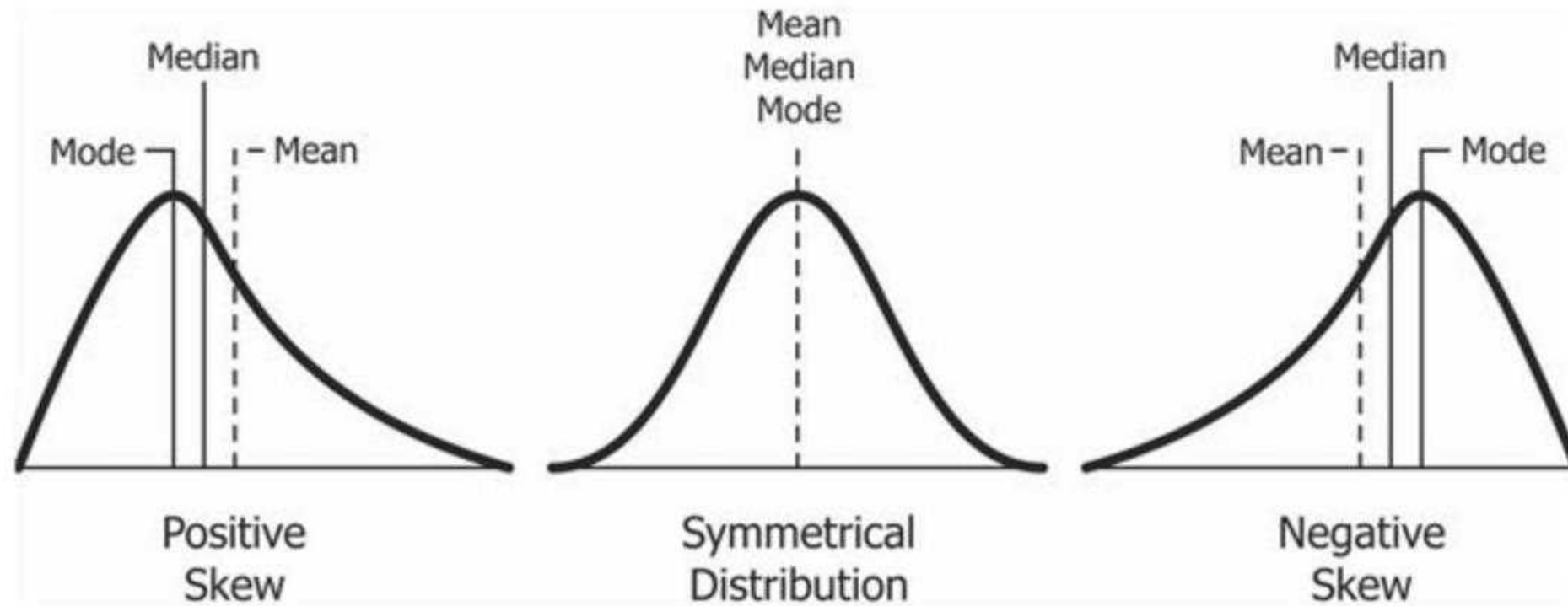  - Kurtosis and skewness for the shape of distributions

16          18

# Kurtosis

# skewness

# skewness



| Positive Skew | Symmetrical Distribution | Negative Skew |

# Central Tendency

- Arithmetic mean
  - $mean(x) = \frac{1}{n}\sum_{i=1}^{n} x_i$ with $x = (x_1, \ldots, x_n) \in \mathbb{R}^n$

- Median
  - The value that separates the higher half from the data of the lower half

- Mode
  - The value that appears most in the data

# Variability

- Measure for the spread of the data

- Standard deviation
    - Measure for the difference of observation to the arithmetic mean
    - $sd(x) = \sqrt{\dfrac{\sum_{i=1}^{n} (x_i - mean(x))^2}{n-1}}$

- Interquartile Range (IQR)
    - Percentile: value below which a given percentage falls
    - Difference between the 75% percentile and the 25% percentile

# percentile

- K-th percentile is x
  - K% of the values are less than x
  - (100 – K) % of the values are larger than x

# percentile

| | 25th PERCENTILE | 50th PERCENTILE | 75th PERCENTILE |
|---|---|---|---|
| EXAMPLE SALARY | $100,000 | $115,000 | $135,000 |
| WHAT IT MEANS | 25% of companies surveyed are paying $100,000 or less for this role (and the other 75% of companies are paying more than $100,000) | 50% of companies surveyed at paying $115,000 or less for this role (and the other 50% of companies are paying more than $115,000) | 75% of companies surveyed at paying $135,000 or less for this role (and the other 25% of companies are paying more than $130,000) |

# Range of data

- Range for which values are observed

- Minimum: Smallest observed value

- Maximum: Largest observed value


- May be strongly distorted by invalid data
  - Makes it also a good tool to discover invalid data
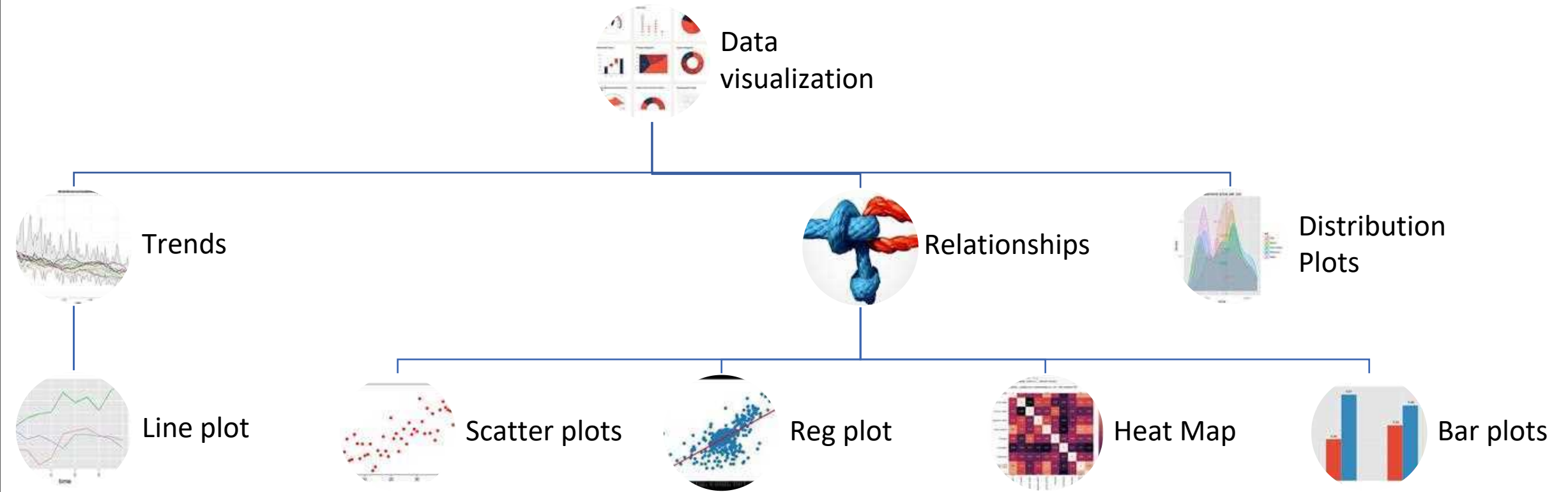
# Example

- Random typing on the keypad
- $x$ = (1, 2, 1, 1, 3, 4, 5, 2, 3, 4, 5, 1, 3, 2, 1, 6, 5, 4, 9, 4, 3, 6, 1, 5, 6, 8, 4, 6, 5, 1, 3, 2, 1, 6, 8, 7, 6, 1, 3, 1, 6, 8, 4, 7, 6, 4, 3, 5, 4, 9, 7, 4, 3, 1, 4, 6, 8, 7, 9, 1, 4, 6, 1, 3, 8, 6, 7, 4, 9, 6, 5, 1, 3, 6, 8, 7)
- central tendency:
  - mean: 4.46052631579
  - median: 4.0
  - mode (count): 1 (14)
- variability
  - sd: 2.41944311488
- range
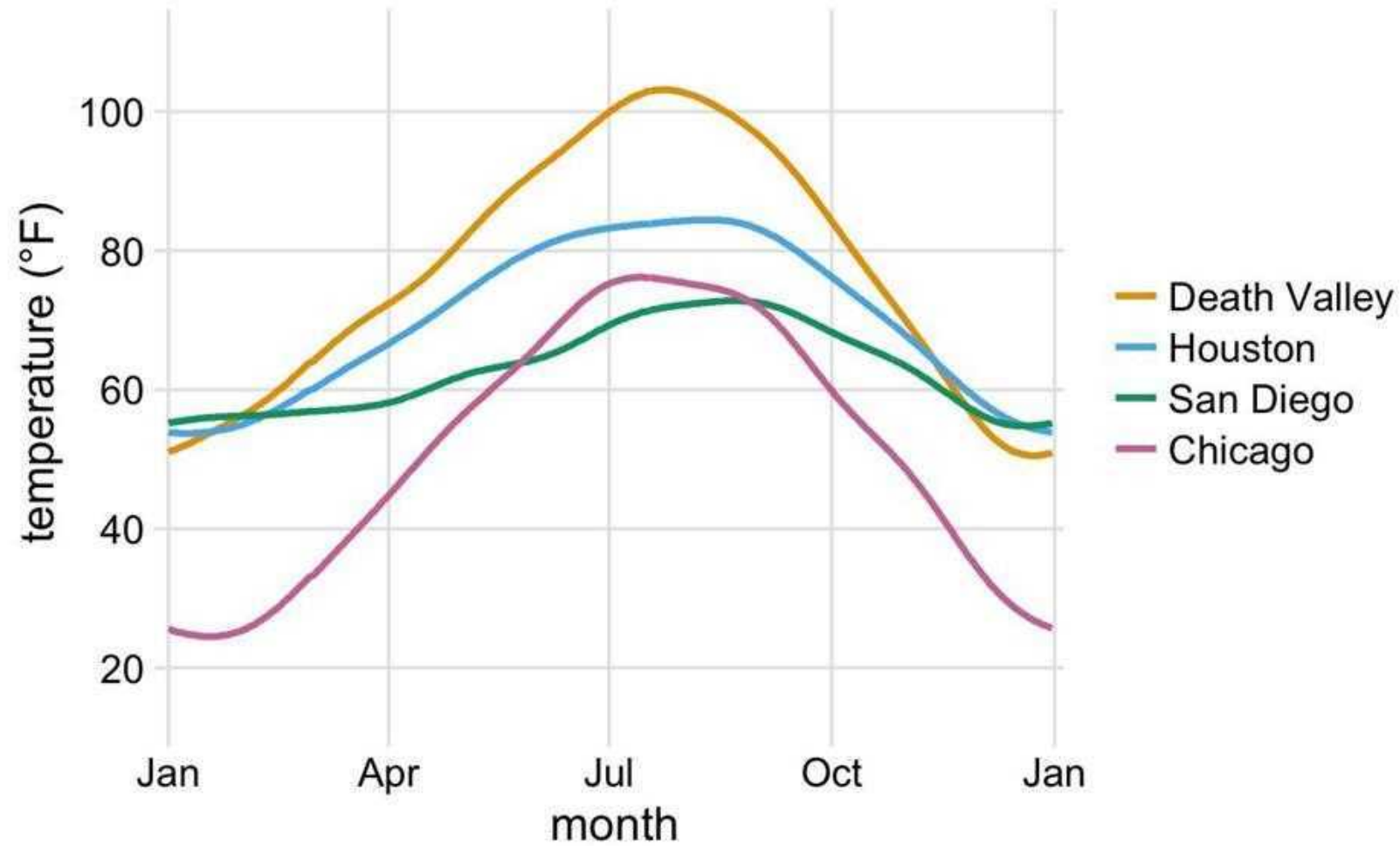  - min: 1
  - max: 9

# Visualization

- A Picture Says More than 1000 Words
- Processed faster than textual information
- Simplify complex information
- Show patterns
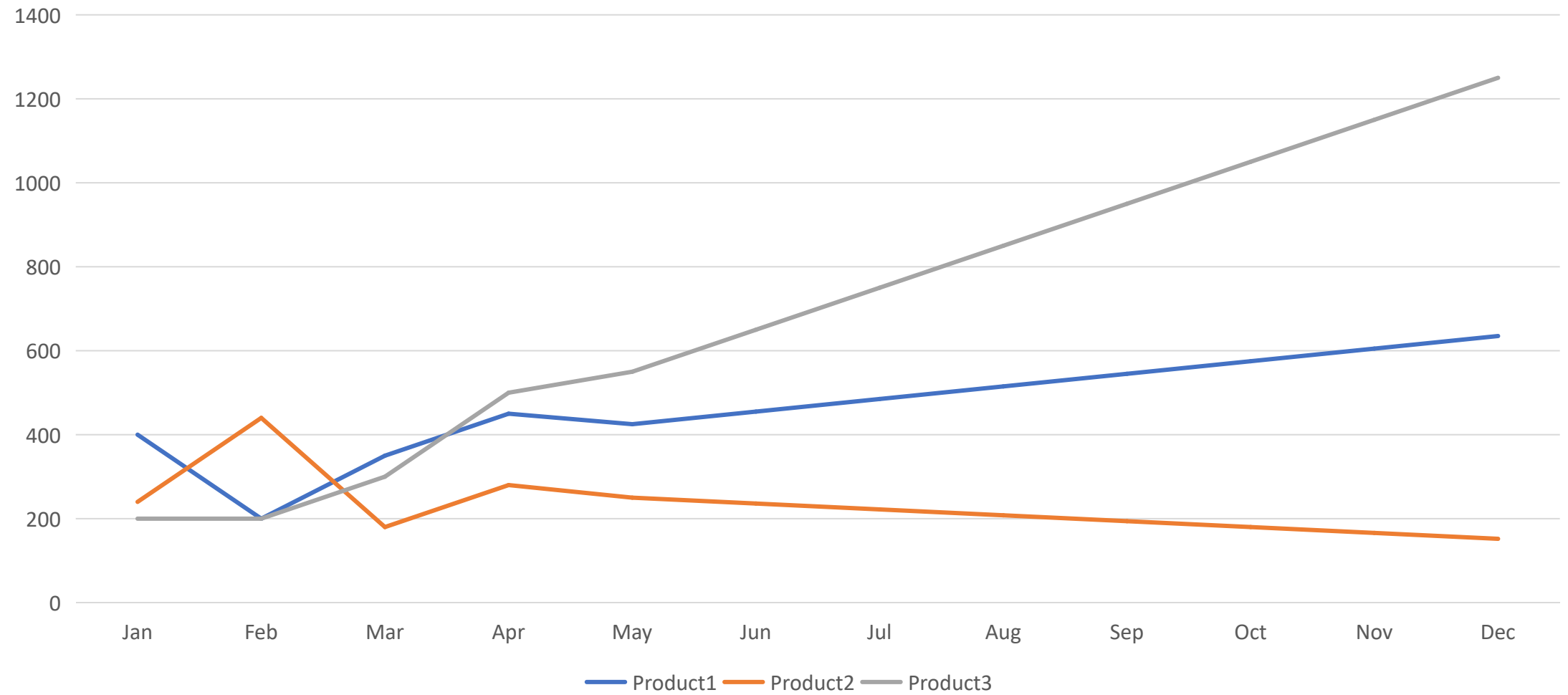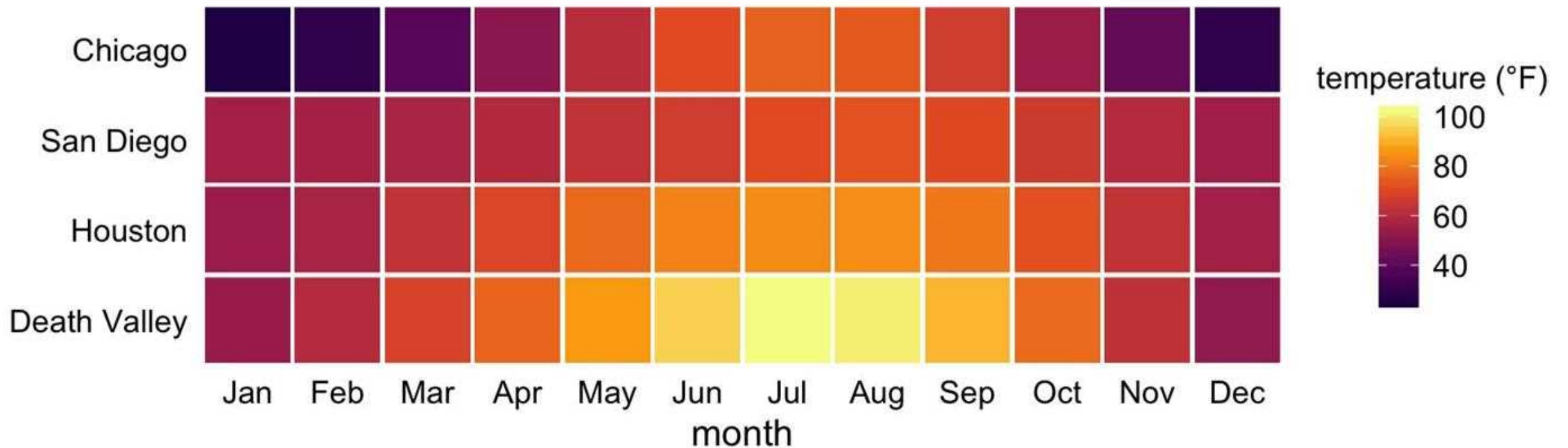- Better long memrization (information retention)
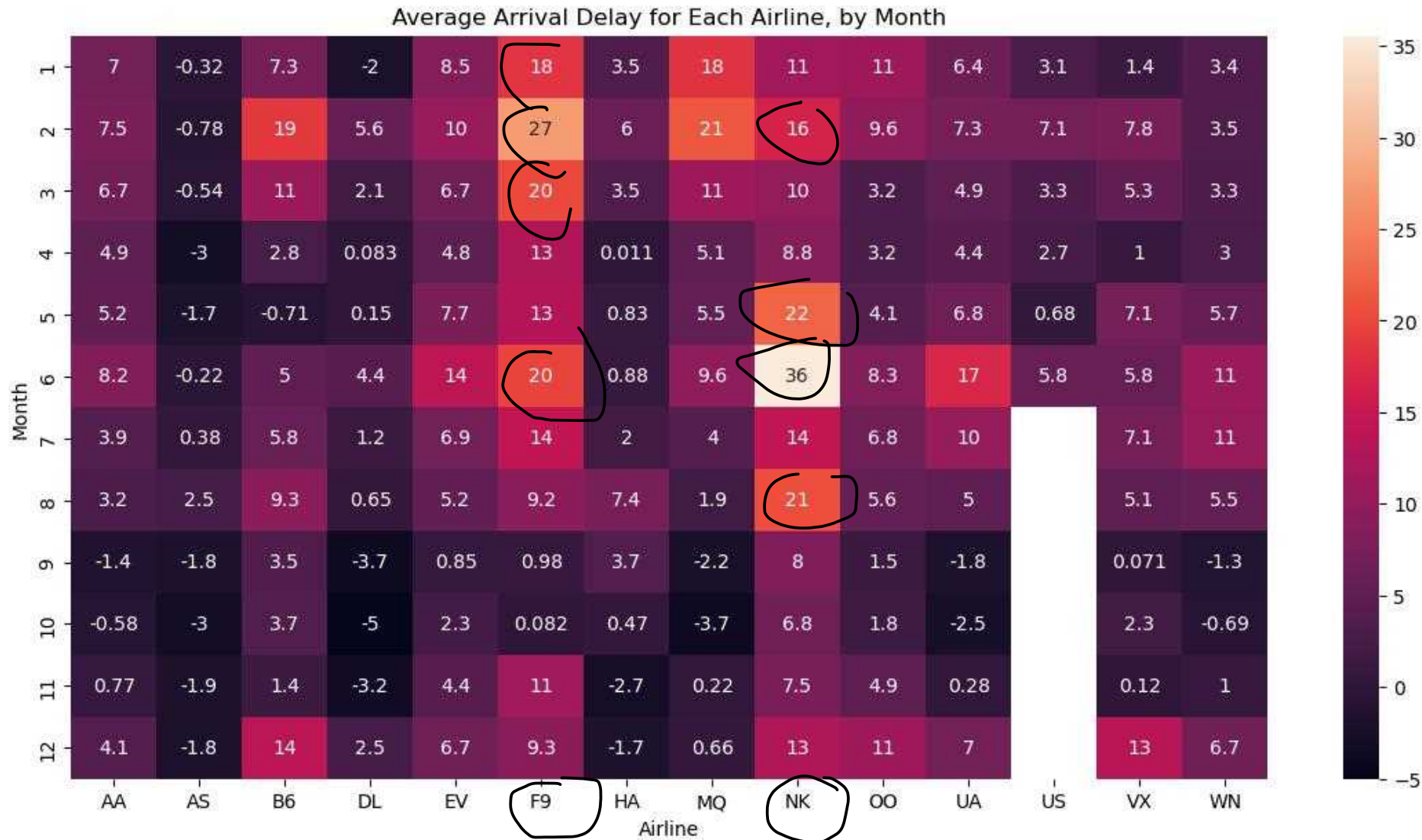
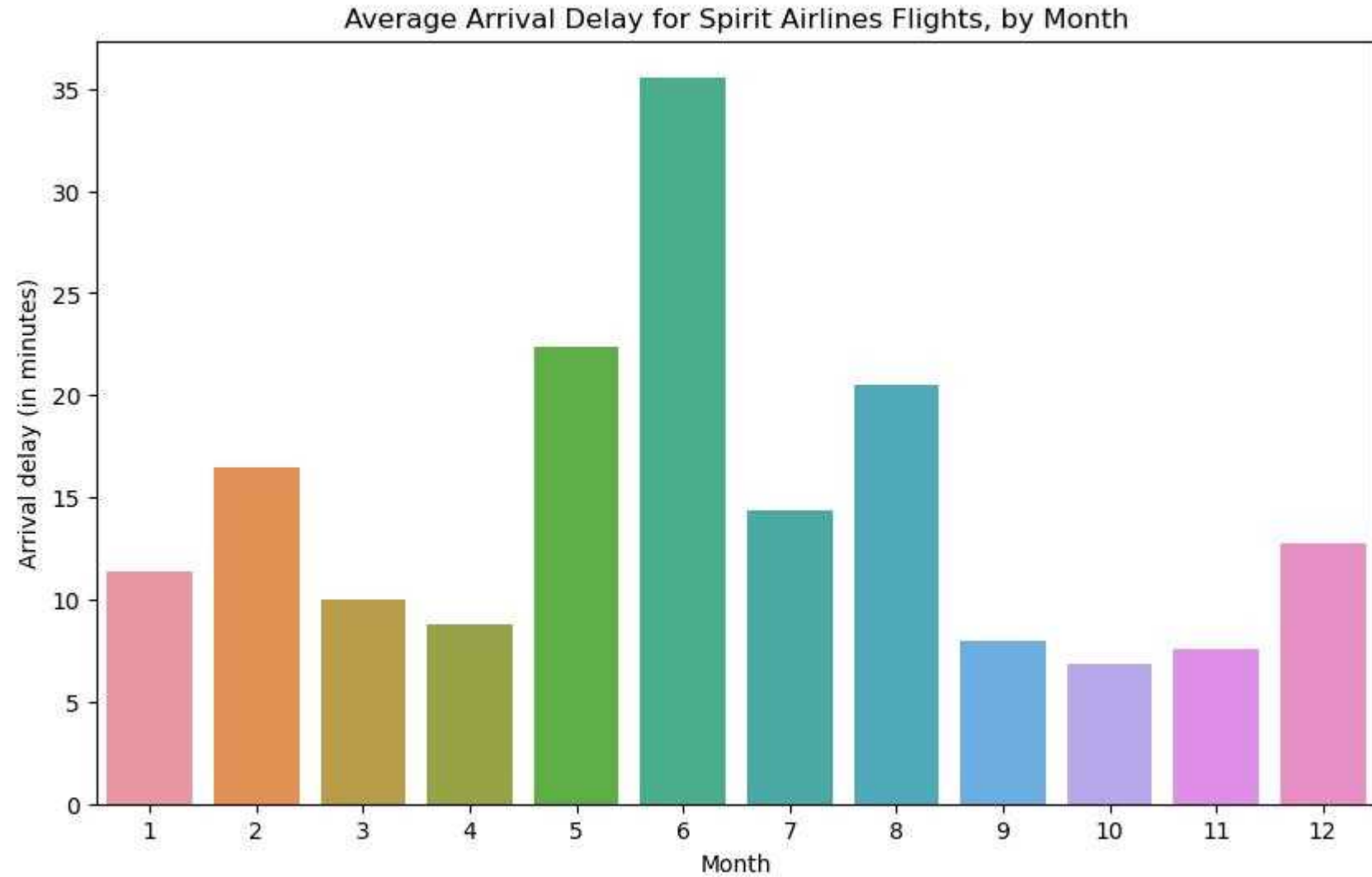# Types of visualizations

# Line plot

# Heat Map

# Arrival time delay Heatmap



Average Arrival Delay for Each Airline, by Month

# Bar plots



Average Arrival Delay for Spirit Airlines Flights, by Month

# Bar plots

# Proportions

# Proportions

# Distributions

# Distribution plot



Distribution of Petal Lengths, by Species

# Scatter Plot

# Scatter Plot

# Map plot

**Product Reviews** — Sentiments Analysis — FACTION A

**Date**
3/21/2011 — 2/6/2018

**Products**
All

**Rating**
All

**Sentiments**
(Blank) | anticipation | fear | negative | sadness | trust
anger | disgust | joy | positive | surprise

**KPIs**

Rating
4.3 ★★★★

Rated Products
8 🛒

Reviews
511 💬

Average Price
$150.73

**Word Frequency by Sentiment**

| Sentiment | Count |
|-----------|-------|
| positive | 407 |
| trust | 305 |
| anticipation | 264 |
| joy | 259 |
| negative | 220 |
| sadness | 139 |
| anger | 115 |
| fear | 107 |
| surprise | 90 |
| disgust | 78 |

**Reviews Classification**

negative 28%
positive 72%

**Reviews by Year/Month/Day**

Sentiment ● negative ● positive

(line chart with axis values 500, 400, 300, 200, 100, 0 and years 2012, 2014, 2016, 2018)

**Total Score by Word**

| Word | Score |
|------|-------|
| perfect | 869 |
| great | 435 |
| top | 414 |
| happy | 372 |
| love | 372 |
| nice | 288 |
| excellent | 198 |
| easy | 165 |
| beautiful | 138 |

(0K – 1K)

**Rating vs Product Score**

(scatter plot, Rating (*) axis 3.5 – 4.5, Score axis up to 1.5)

**Word Cloud**

hard found sturdy front vanity deep company built
perfectly bought damaged
job assembled top hours components received finished missing
buy car plastic arrived shelf
loves wood shelves
fits pictures product bottom looks cheap
furniture close time love piece assembly
daughter received perfect stereo
board particle storage beautiful fit
makes heavy space together
lot size worth
fairly drawers easily price
newer quickly box difficult great
quality excellent desk stand easy
couple follow putting shipping
happy nice bad build instructions
purchase pieces recommend equipment picture
extremely office placed

# Roadmap for data exploration

1. Organize the data set
2. Find the central point for each attribute
3. Understand the spread of the data for each attribute
4. Visualize the distribution of each attribute
5. Watch out for outliers
6. Understanding the relationship between attributes
7. Visualize the relationship between attributes

# Data preparation

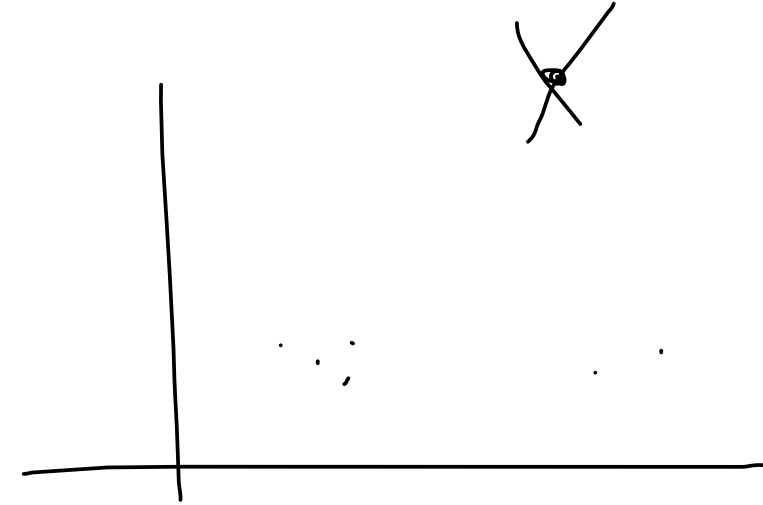# Data pre-processing tasks

| Main Task | Subtasks | Popular Methods |
|-----------|----------|-----------------|
| Data consolidation | Access and collect the data Select and filter the data Integrate and unify the data | SQL queries, software agents, Web services. Domain expertise, SQL queries, statistical tests. SQL queries, domain expertise, ontology-driven data mapping. |
| Data cleaning | Handle missing values in the data | • Fill in missing values (imputations) with most appropriate values (mean, median, min/max, mode, etc.); <br> • recode the missing values with a constant such as "NA"; <br> • remove the record of the missing value; <br> • do nothing. |
| | Identify and reduce noise in the data | Identify the outliers in data with simple statistical techniques (such as averages and standard deviations) or with cluster analysis; once identified, either remove the outliers or smooth them by using binning, regression, or simple averages. |
| | Find and eliminate erroneous data | Identify the erroneous values in data (other than outliers), such as odd values, inconsistent class labels, odd distributions; once identified, use domain expertise to correct the values or remove the records holding the erroneous values. |

# Data pre-processing tasks

| Main Task | Subtasks | Popular Methods |
|---|---|---|
| Data transformation | Normalize the data | Reduce the range of values in each numerically valued variable to a standard range (e.g., 0 to 1 or −1 to +1) by using a variety of normalization or scaling techniques. |
| | Discretize or aggregate the data | If needed, convert the numeric variables into discrete representations using range- or frequency-based binning techniques; for categorical variables, reduce the number of values by applying proper concept hierarchies. |
| | Construct new attributes | Derive new and more informative variables from the existing ones using a wide range of mathematical functions (as simple as addition and multiplication or as complex as a hybrid combination of log transformations). |
| Data reduction | Reduce number of attributes | Use principal component analysis, independent component analysis, chi-square testing, correlation analysis, and decision tree induction. |
| | Reduce number of records | Perform random sampling, stratified sampling, expert-knowledge-driven purposeful sampling. |
| | Balance skewed data | Oversample the less represented or undersample the more represented classes. |

# Summary

- Important to understand the data available

- Summary statistics provide a good overview
  - Can be deceptive!

- Visualization is a powerful way to understand data

- Data prerpration tasks is important because real-world data is not clean and ordered