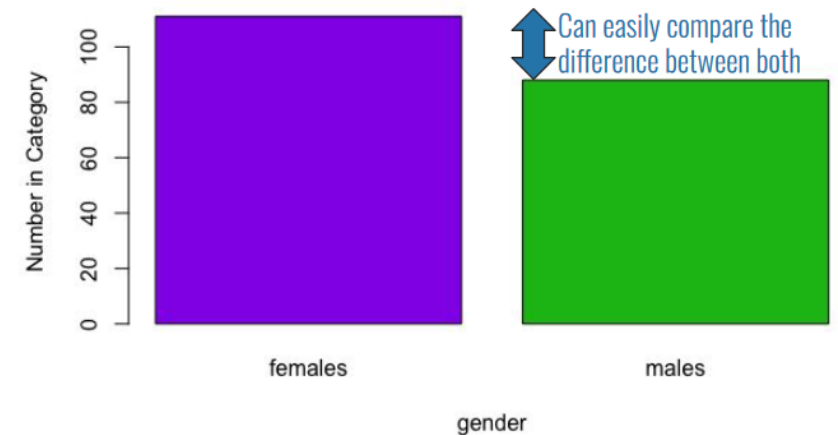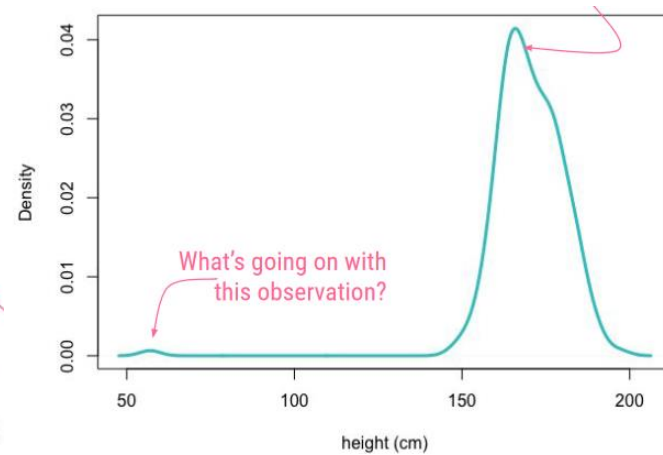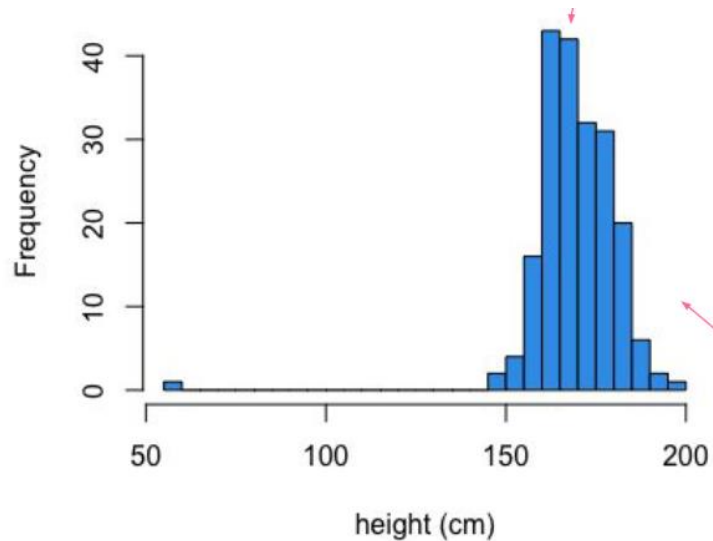# COGS 9 – A05 Discussion

# Deadlines

- Mid Way Team Evaluations E.C. – October 28th (Today)
  - 2 bonus points if you fill the form
- Assignment 2 – October 28th (Today)
  - 40 points
  - Read through and follow the instructions
- Reading Quiz 4 – November 3rd (Thursday)
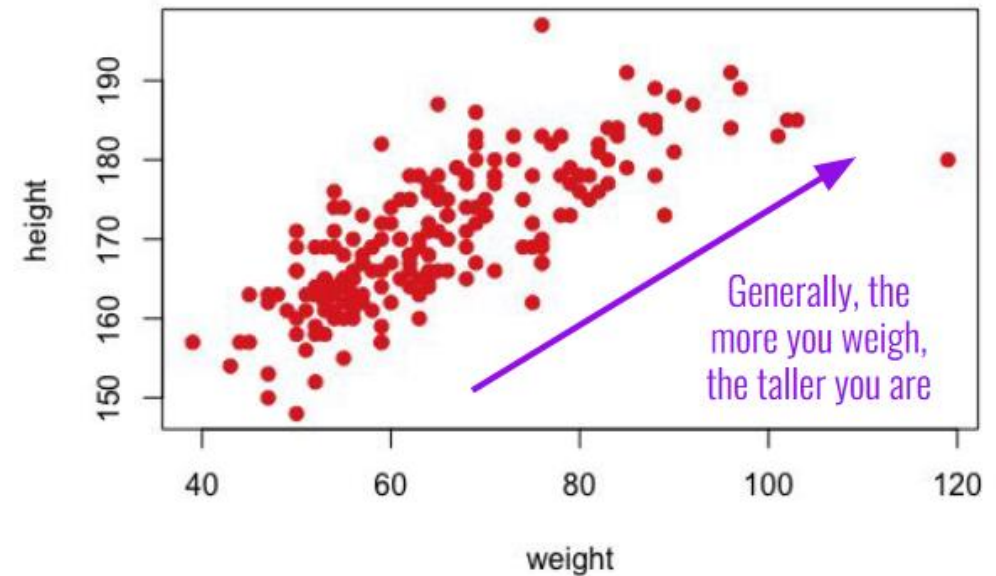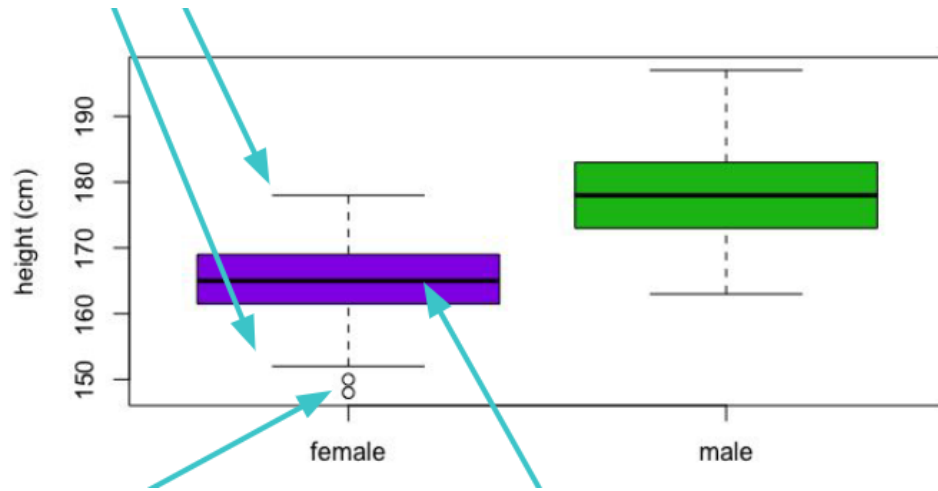  - 10 points

- Any issues with Assignment 1 grades?

# Lecture 9: Data Visualisation

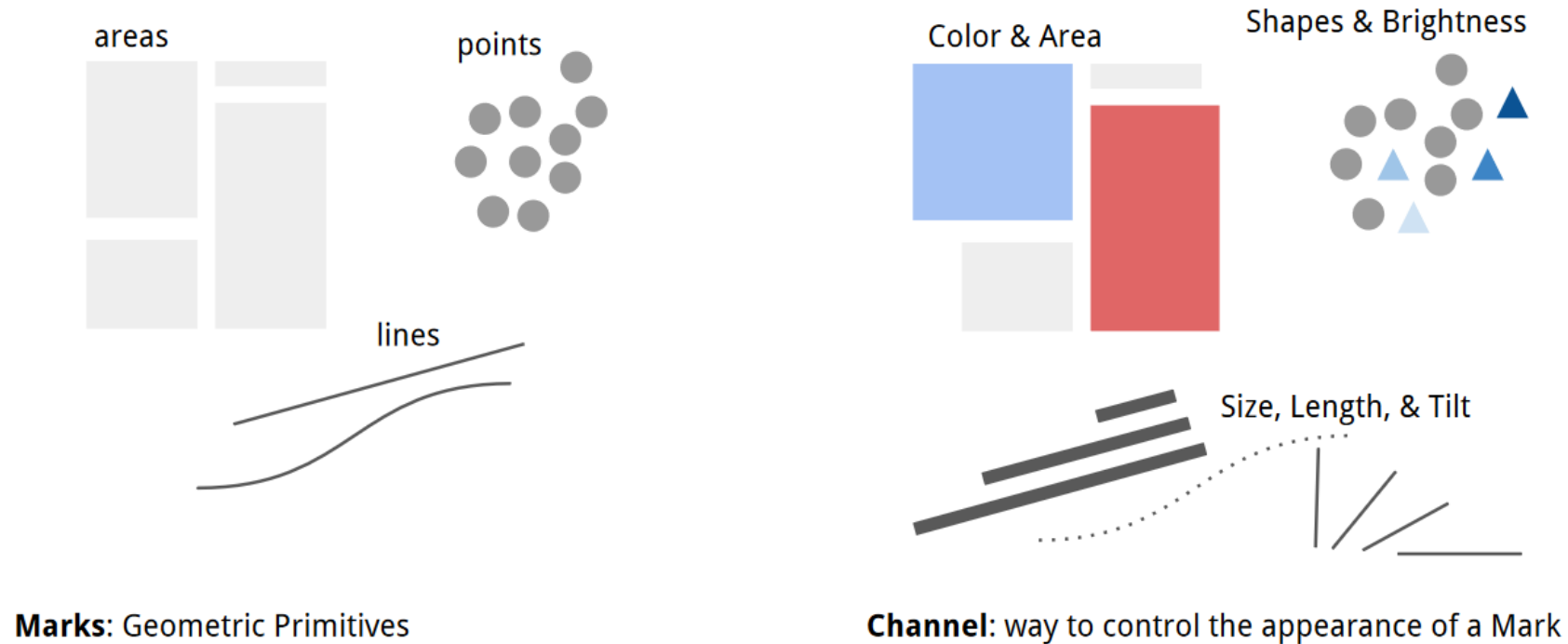- What's the difference between a histogram, densityplot and barplot?

# Lecture 9: Data Visualisation

• When to use a scatterplot and a boxplot?
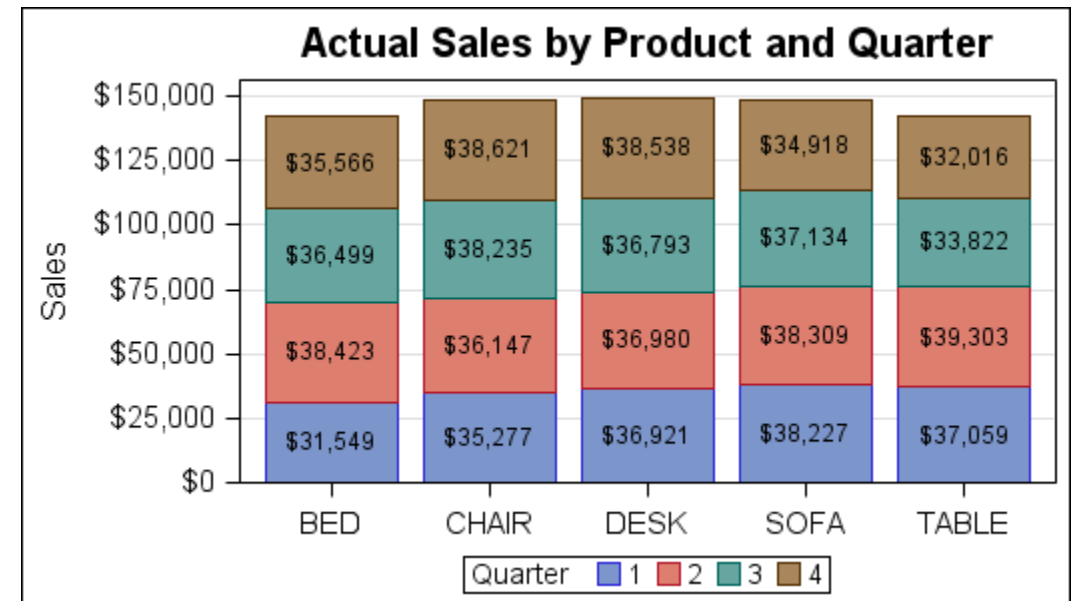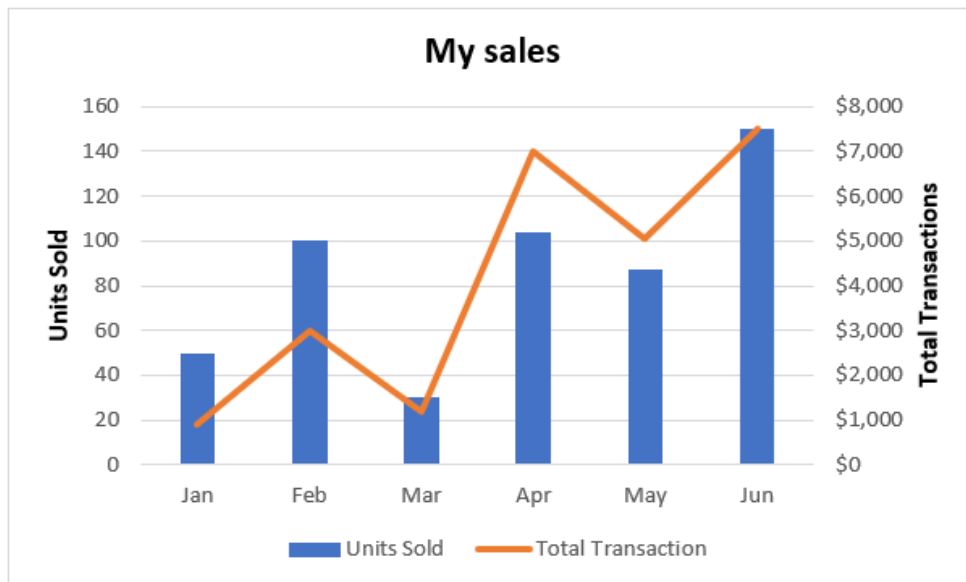
# Lecture 9: Data Visualisation

- Marks and channels



areas    points

lines

**Marks**: Geometric Primitives

Color & Area    Shapes & Brightness

Size, Length, & Tilt

**Channel**: way to control the appearance of a Mark

# Lecture 9: Data Visualisation

- Express: The visual shall express all of, and only, the information provided by the data's attributes. The visual should not add anything extra to the data.

- Effect: How important an attribute is, must match the salience of the channel. Greater importance = greater salience, or more noticeable

# Lecture 9: Data Visualisation

Checklist when creating graphs

[] Consideration for Colorblindness

[] Label the axes

[] Ensure that the data is correct

[] Ensure that the graphic represents the data

[] Make the comparison easy on readers

[] Ensure that the y-axis starts at 0 (What about x-axis?)

[] Choose best visual

[] Keep it Simple Stupid

# Lecture 9: Data Visualisation

Checklist when creating tables

[] Have a top to bottom comparison

[] Logical row ordering

[] Logical column ordering

[] Limit number of rows and columns

[] Informative headers

[] Fix significant digits

[] Format table

# Lecture 9: Descriptive Analysis

- Suppress some of the truth so that humans can understand easily
- Size, shape, missingness, central tendency, variability
- Size: Number of variables and observations
- Shape: Distribution of the variables (Uniform, bimodal, Normal/Gaussian/Bell-shaped, left & right skewed, random)
- Missingness: How much data is missing?
- Central Tendency: Mean, median, mode
- Variability: Variance, Standard Deviation, Range

# Lecture 10: Exploratory Data Analysis

- Data -> Descriptive Analysis -> Exploratory Analysis -> Product
- Exploratory: Inferential, Predictive, Causal, Mechanistic
  - Inferential: Statistics, Frequentist, Bayesian, Text & Geospatial analysis
  - Predictive: Statistical Learning/ML, Deep, Reinforcement Learning
  - Causal: How variable X correlates to Y
  - Mechanistic: How much does variable X affect Y
- Univariate, Bivariate, Multivariate
- Explanatory (Independent) vs Response (Dependent) variables
- Source of data (Zipcode vs hometown), explore missing data
- Don't do EDA to give you the result you want

# Lecture 10: Exploratory Data Analysis

- Checklist of things to do during EDA

[] Investigate missing values

[] Understand outliers

[] Add filters, transform and scale data

[] Calculate numerical summaries

[] Generate plots to explore relationships

[] Handle proportions correctly

[] Use tables to scan data

[] Search for patterns

# Lecture 9 and 10 Demo

# Reading 3: Tidy Data

- Tidy datasets provide a standardized way to link the structure of a dataset (physical layout) with its meaning

- Structure: Rectangular (Rows and columns)

- Semantics: Numbers (Quantitative), Strings (Qualitative)

- Each variable forms a column, each observation forms a row and each type of observational unit forms a table (Codd's 3$^{rd}$ Normal Form)

- How to tidy messy datasets

- Tools for tidying through manipulation, visualization & modelling

- Case study on tidying data using R

# Reading 3: Data organization in spreadsheets

How does cannabis compare to other drugs?

# Reading 4: Attitudes and Perceptions of Data Visualization

- Which visualization do people understand? Which do they pay attention to? Study in a rural area in Pennsylvania
- Interviewed 42 people (diversified in terms of education, age, political affiliation, drug impact) – Asked to rank 10 graphics based on usefulness
- People ranked according to the following
  - Events that occurred in their personal lives
  - Geographical information which impacted them
  - How useful it is to other people
  - Clarity and Novelty
  - Statistical familiarity
  - Source of the graphic (Whether people changed their ranking)
  - Political Identity