



CS37300: Data Mining and Machine Learning

Exploratory Data Analysis

Sep 13 2023

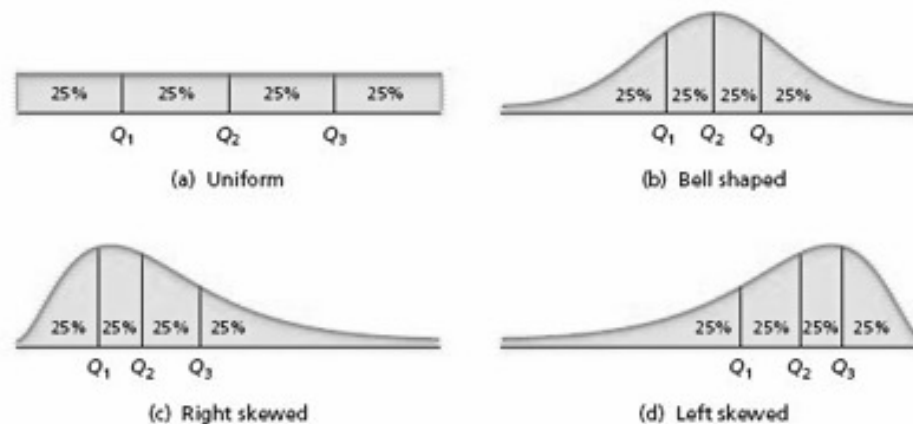


Exploratory data analysis

- Data analysis approach that employs a number of (mostly graphical) techniques to:
 - Maximize insight into data
 - Uncover underlying structure
 - Identify important variables
 - Detect outliers and anomalies
 - Test underlying modeling assumptions
 - Develop parsimonious models
 - **Generate hypotheses from data**

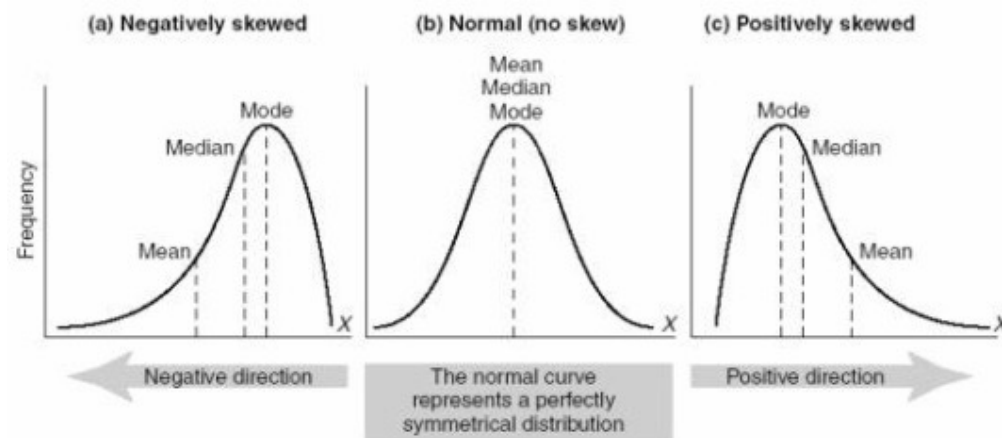
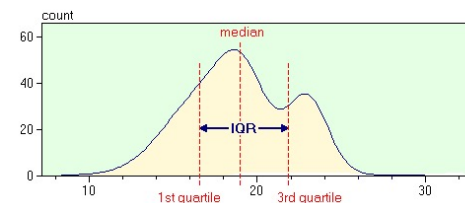
Data summarization

- Measures of location
 - Mean: $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x(i)$
 - Median: value with 50% of points above and below
 - Quartile: value with 25% (75%) points above and below
 - Mode: most common value



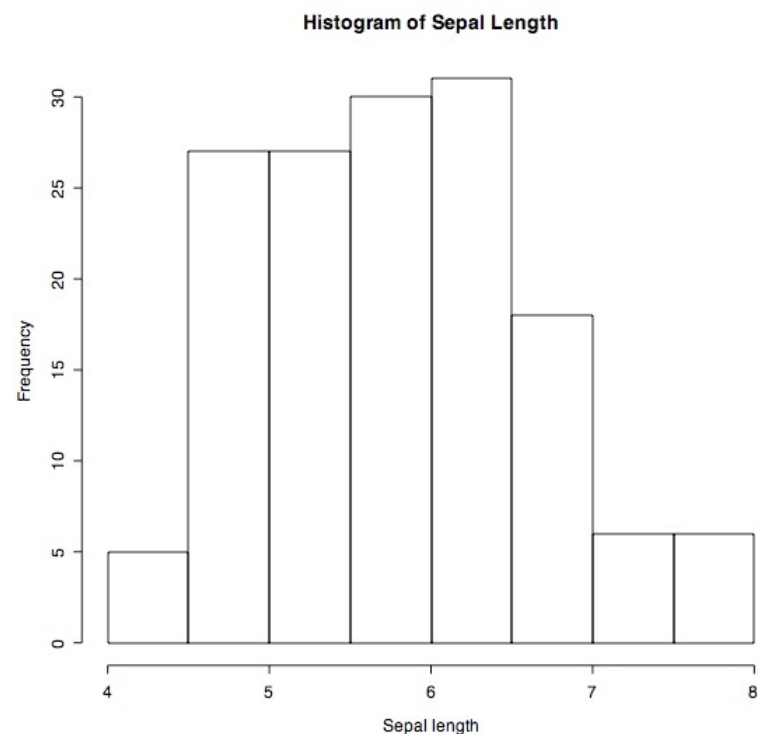
Data summarization

- Measures of dispersion or variability
 - Variance: $\hat{\sigma}_k^2 = \frac{1}{n} \sum_{i=1}^n (x(i) - \mu)^2$
 - Standard deviation: $\hat{\sigma}_k = \sqrt{\frac{1}{n} \sum_{i=1}^n (x(i) - \mu)^2}$
 - Range: difference between max and min point
 - Interquartile range: difference between 1st and 3rd Q
 - Skew: $E \left[\left(\frac{X - \mu}{\sigma} \right)^3 \right] \Leftrightarrow \frac{\sum_{i=1}^n (x(i) - \hat{\mu})^3}{(\sum_{i=1}^n (x(i) - \hat{\mu})^2)^{\frac{3}{2}}}$



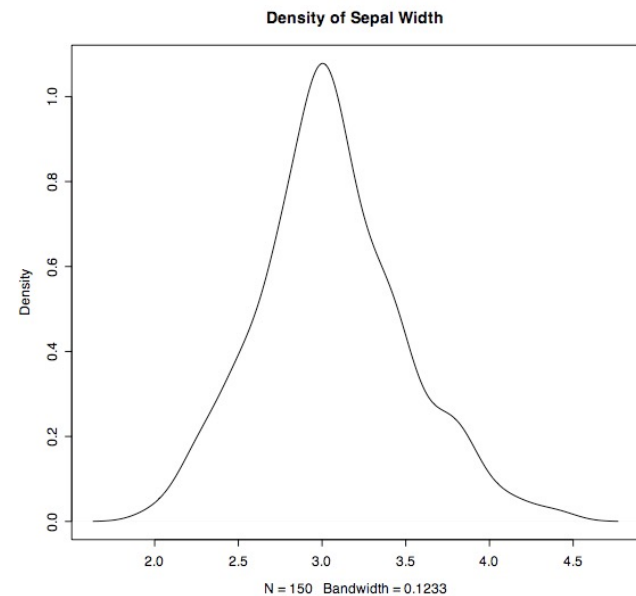
Histograms (1D)

- Most common plot for univariate data
- Split data range into equal-sized bins, count number of data points that fall into each bin
- Graphically shows:
 - Center (location)
 - Spread (scale)
 - Skew
 - Outliers
 - Multiple modes



Histogram limitations

- Histograms can be misleading for small datasets
 - Slight changes in the data or binning approach can result in different histograms
- Solution: smoothed density plots
 - Use kernel function to estimate density at each point x , pools information from neighboring points

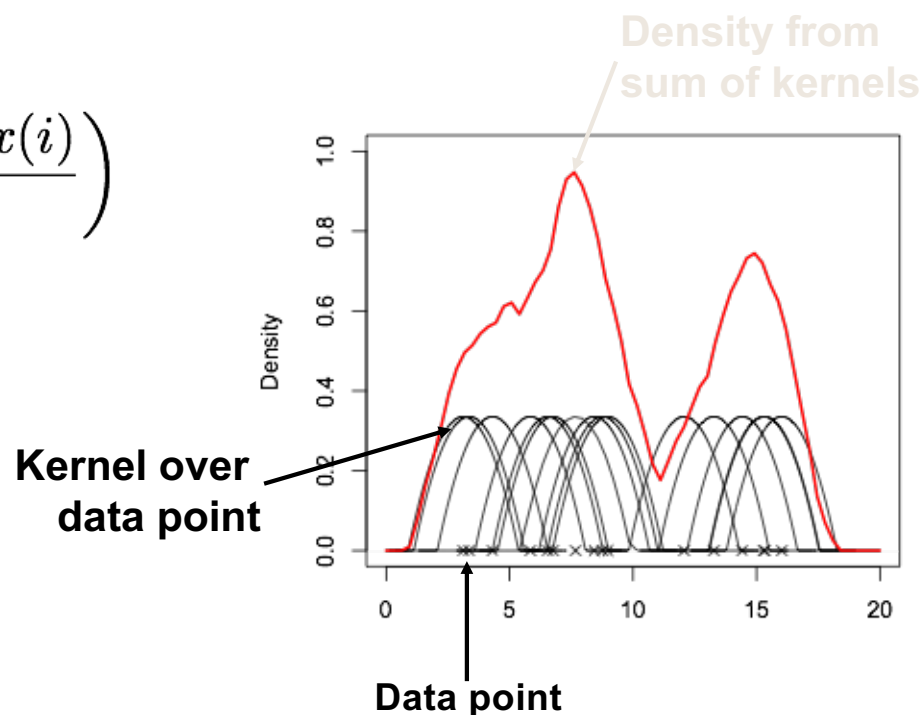


Density plots

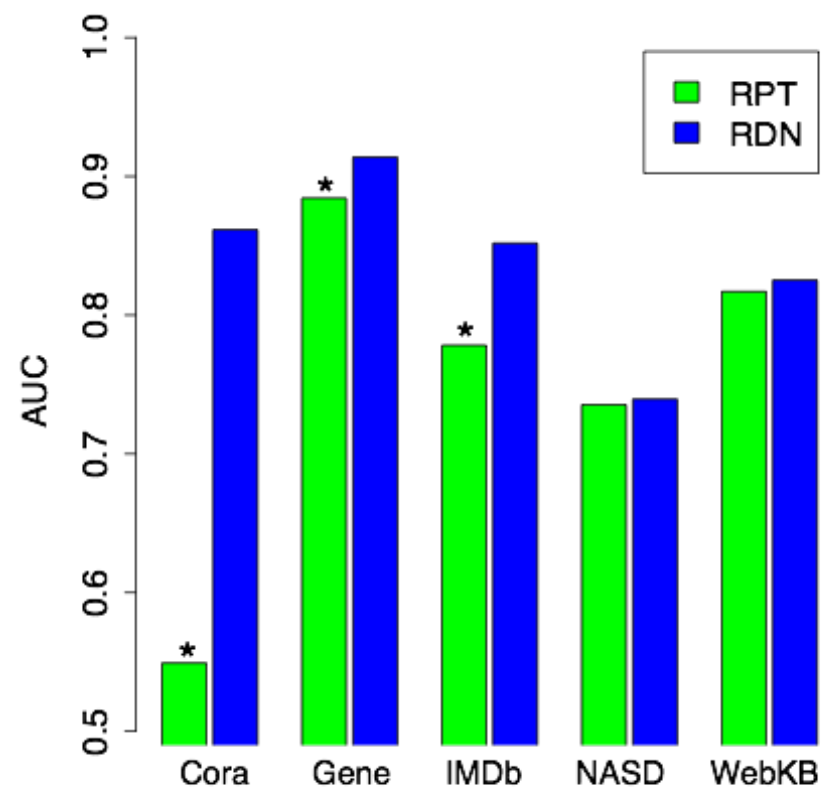
- Estimated density is:

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^n K \left(\frac{x - x(i)}{h} \right)$$

- Two parameters:
 - **Kernel function** K
(e.g., Gaussian, Epanechnikov)
 - Bandwidth h

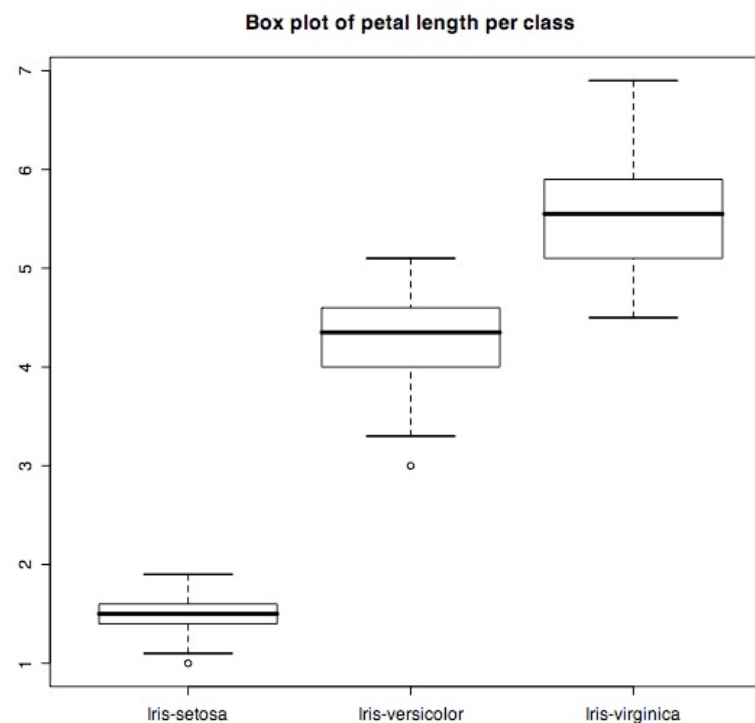


Bar plots



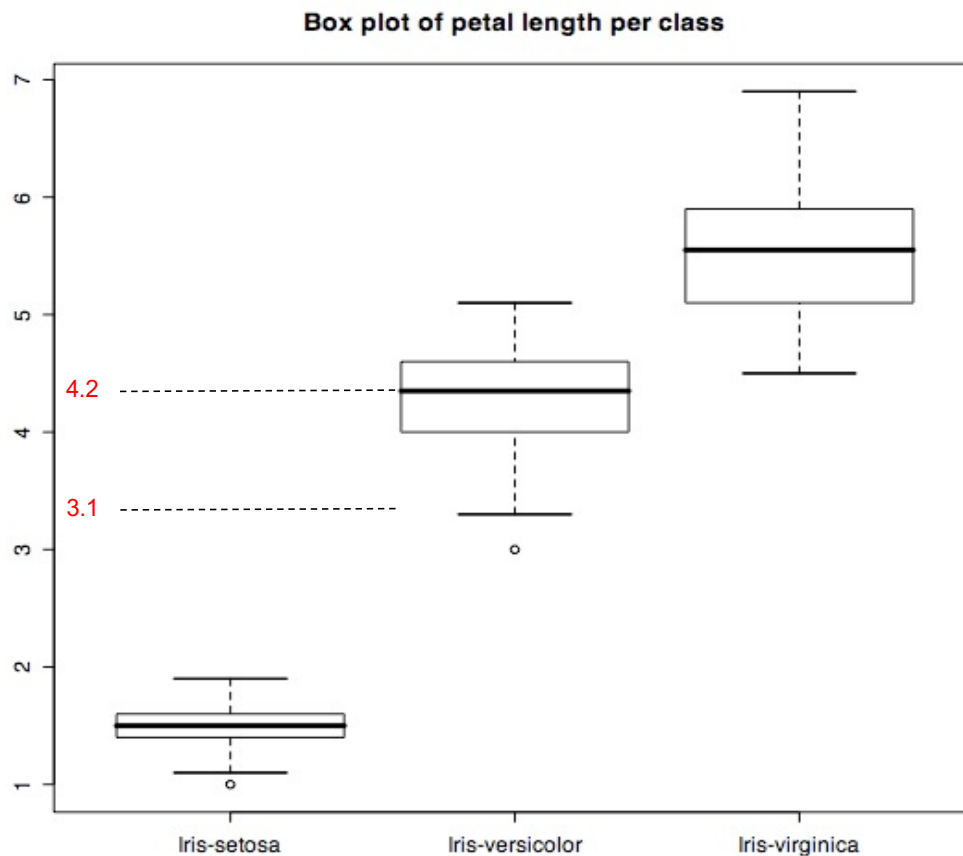
Box plot (2D)

- For each discrete value X , calculate quartiles and range of associated Y values
- Data summary for:
minimum, first quartile, median, third quartile, and maximum
- Can also plot outliers separately



Interpreting Box Plots

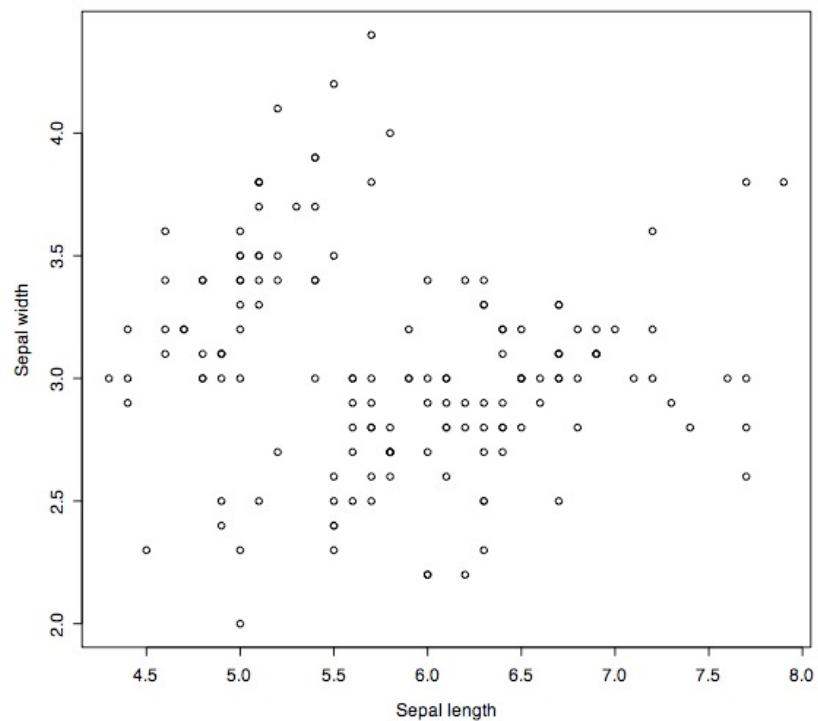
- Petals of Iris-Versicolor are:
 - Always longer than 3.1?
 - At least 50% of the petals are longer than 4



Scatter plot (2D)

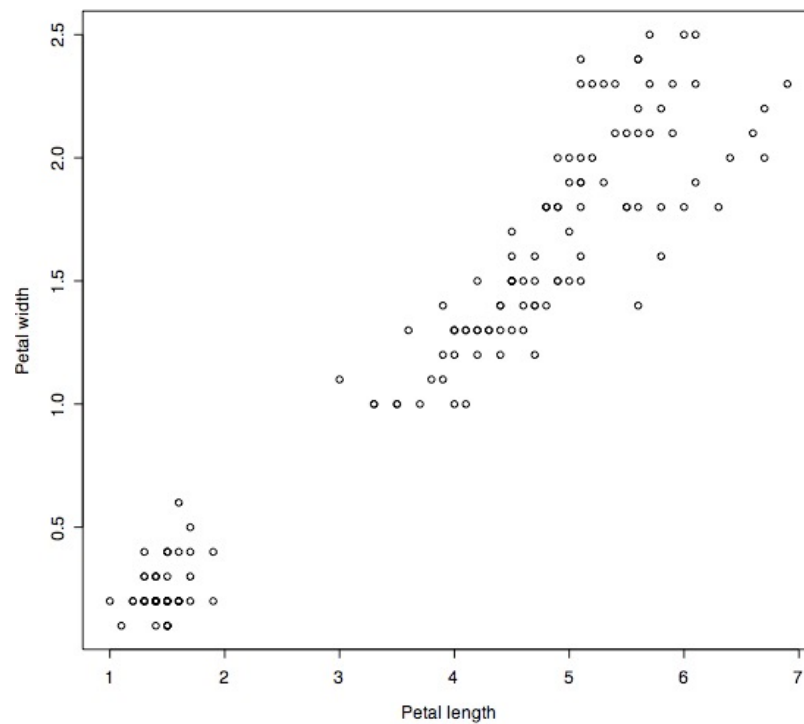
- Most common plot for bivariate data
 - Horizontal X axis: the suspected **independent** variable
 - Vertical Y axis: the suspected **dependent** variable
- Graphically shows:
 - If X and Y are related
 - Linear or non-linear relationship
 - If the variation in Y depends on X
 - Outliers

No relationship

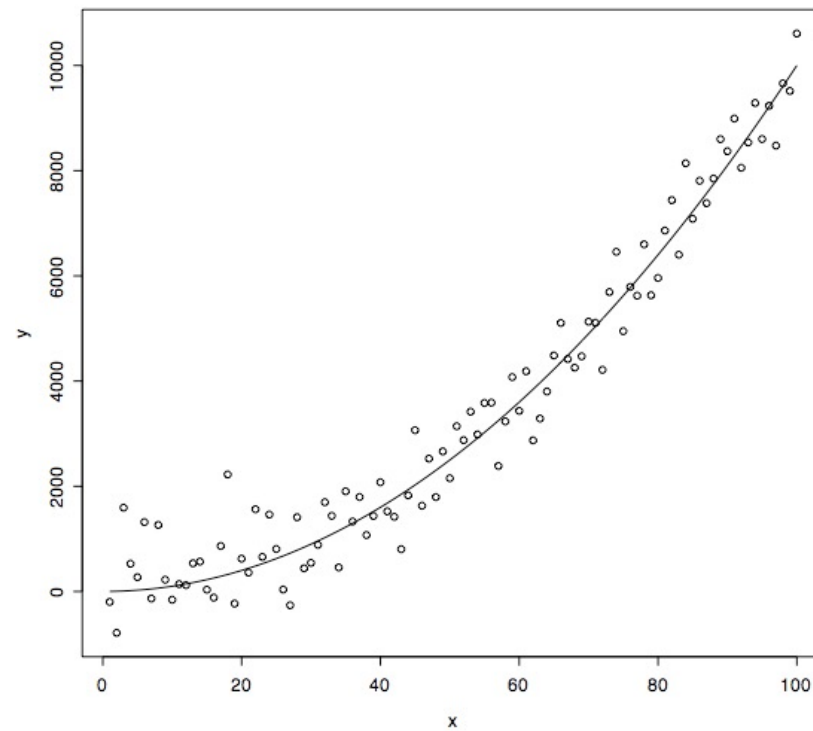


Should this variable be excluded or “pre-pruned” when building a decision tree? Why or why not?

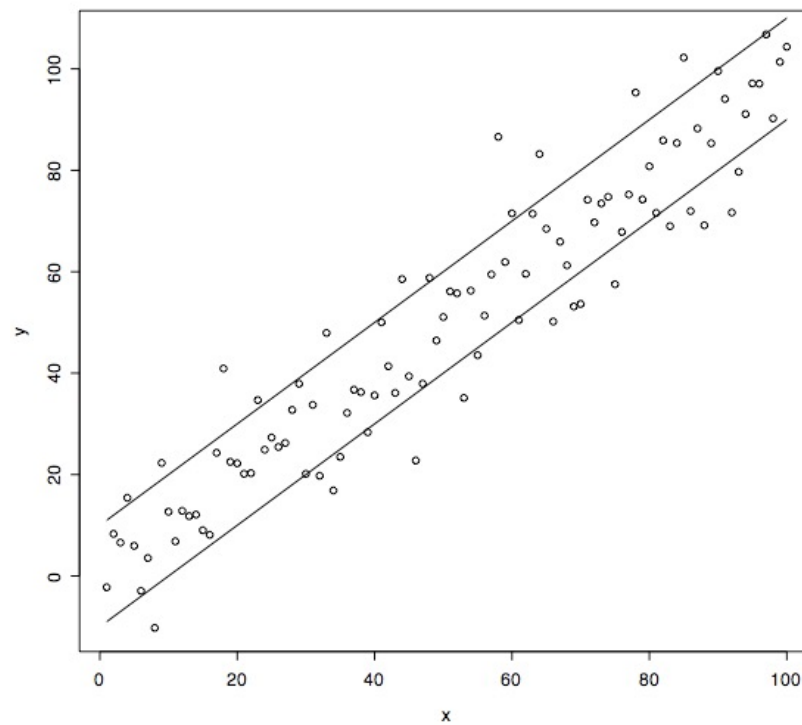
Linear relationship



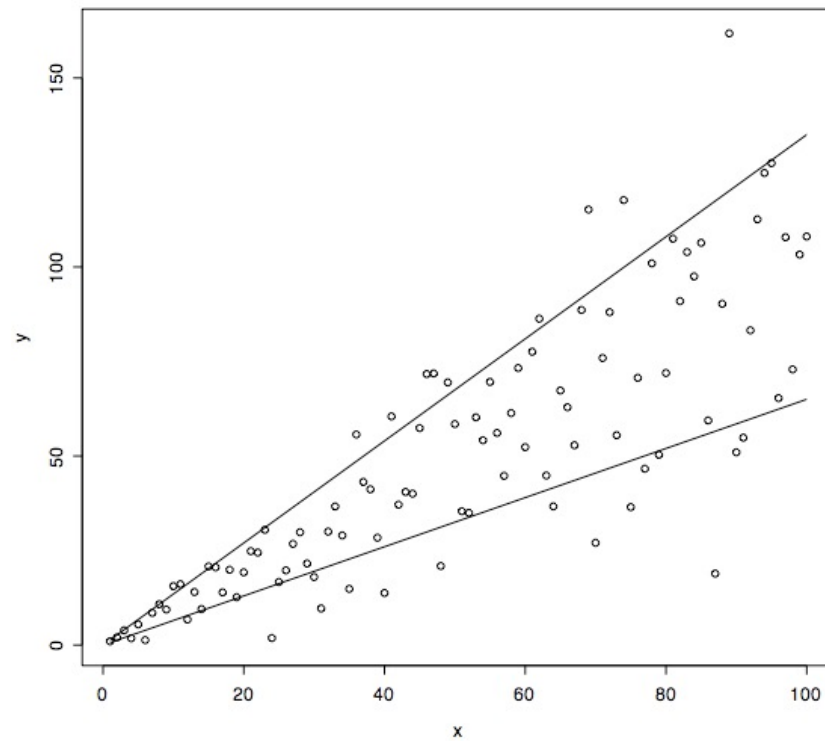
Non-linear relationship



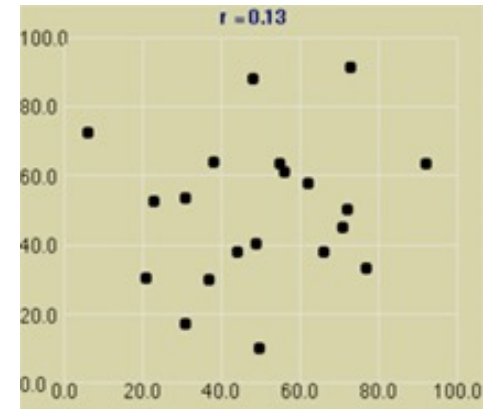
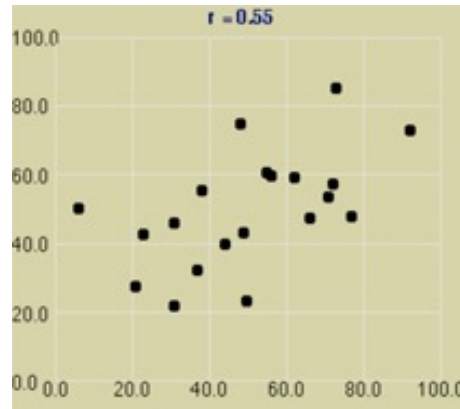
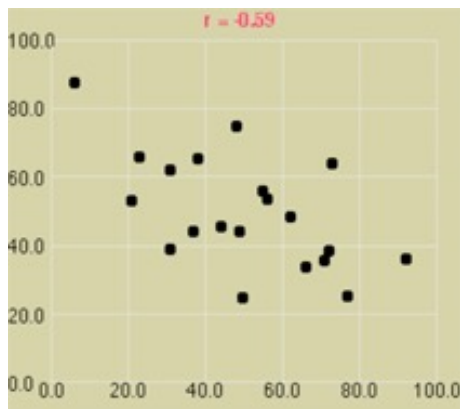
Homoskedastic



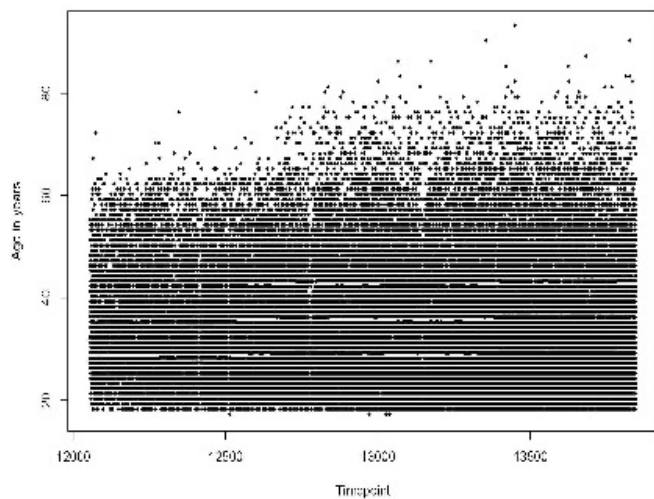
Heteroskedastic



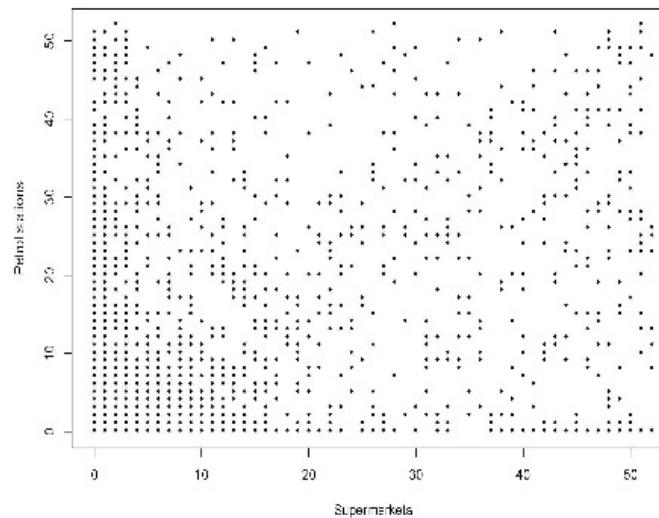
Which one of the plots describes a positive correlation?



Scatterplot limitations

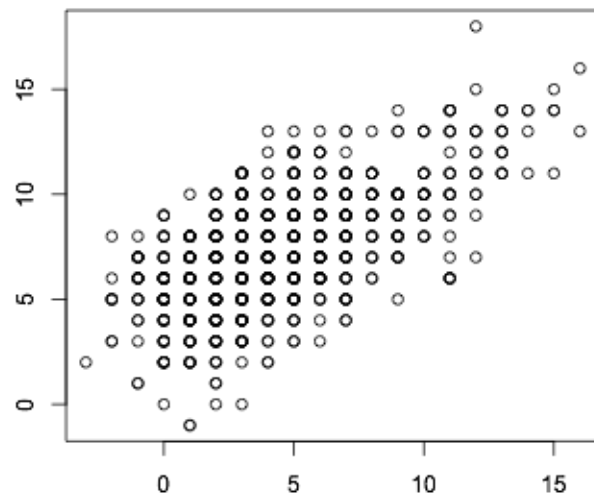


Too much data

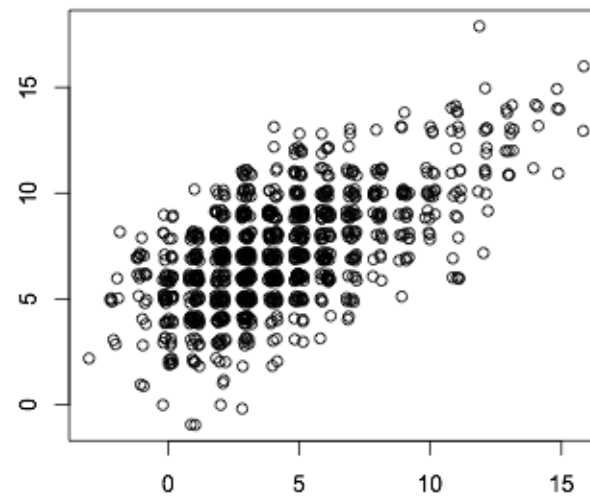


Overprinting

Scatterplot limitations



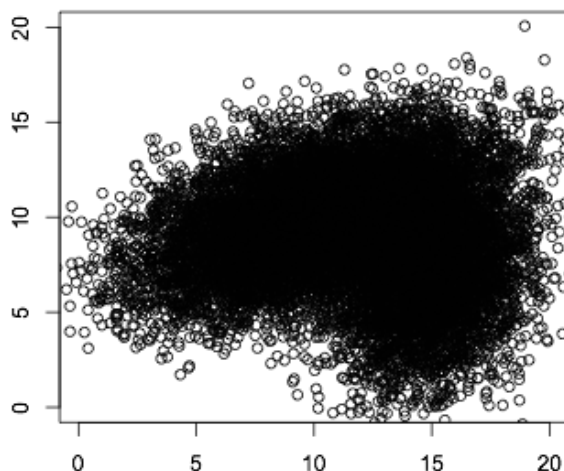
Overprinting



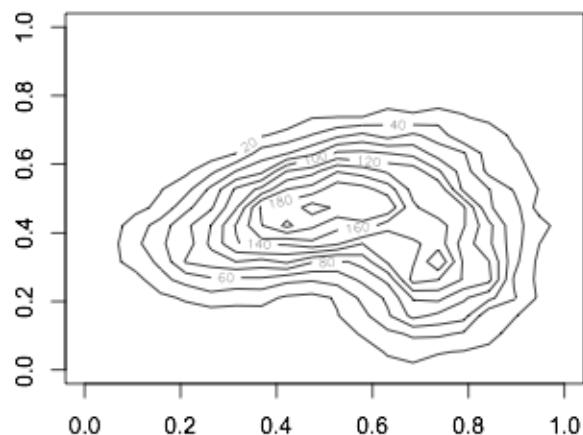
Solution: Jitter points

Contour plot (3D)

- Limitations of 2D scatterplot (e.g., when there is too much data to discern relationship)



- Solution: represent a 3D surface by plotting constant z slices (contours) in a 2D format

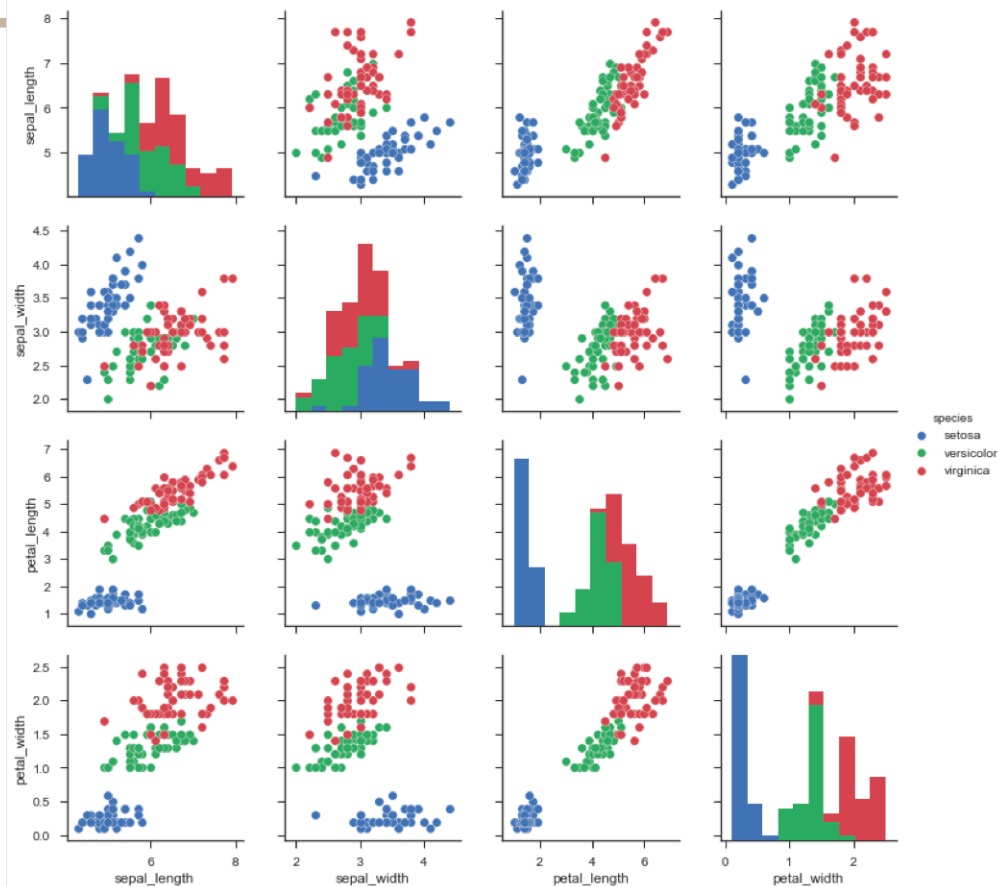


Scatterplot matrix

Good to check for
linear relationships
in multivariate datasets

```
import seaborn as sns
import matplotlib.pyplot as plt
sns.set(style="ticks")

df = sns.load_dataset("iris")
sns.pairplot(df, hue="species")
plt.show()
```



Summary – Data Exploration

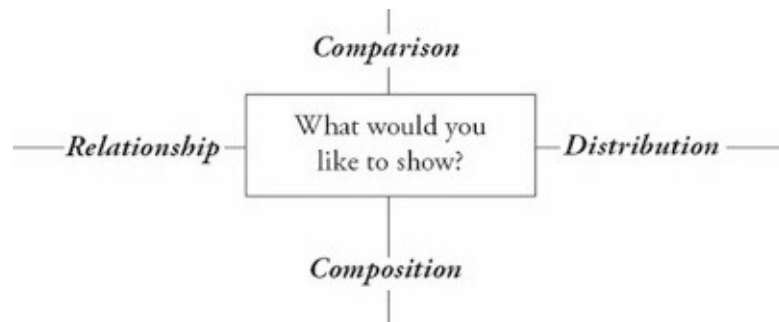


Chart Suggestions—A Thought-Starter

