



THE UNIVERSITY *of* EDINBURGH
informatics

Applied Machine Learning (AML)

Class Starting at 4:10pm

Oisin Mac Aodha • Siddharth N.

Applied Machine Learning

Week 5: Representing Data and Exploratory Data Analysis

*This slides will be made available on the project website after the class.
This session will be recorded.*

Overview

- 1) Outline your tasks this for week
- 2) Discussion of Week 4's topics

Week 5: Your tasks for this week

- 1) Attend and complete Lab 2 - solutions will be online later this week
- 2) Watch videos for week 5 - **Optimisation** and **Generalisation**
- 3) Ask questions on Piazza if stuck
- 4) Continue working on the coursework
- 5) Start Tutorial 2 which takes places next week - link in week 6

Reality of Applied Machine Learning

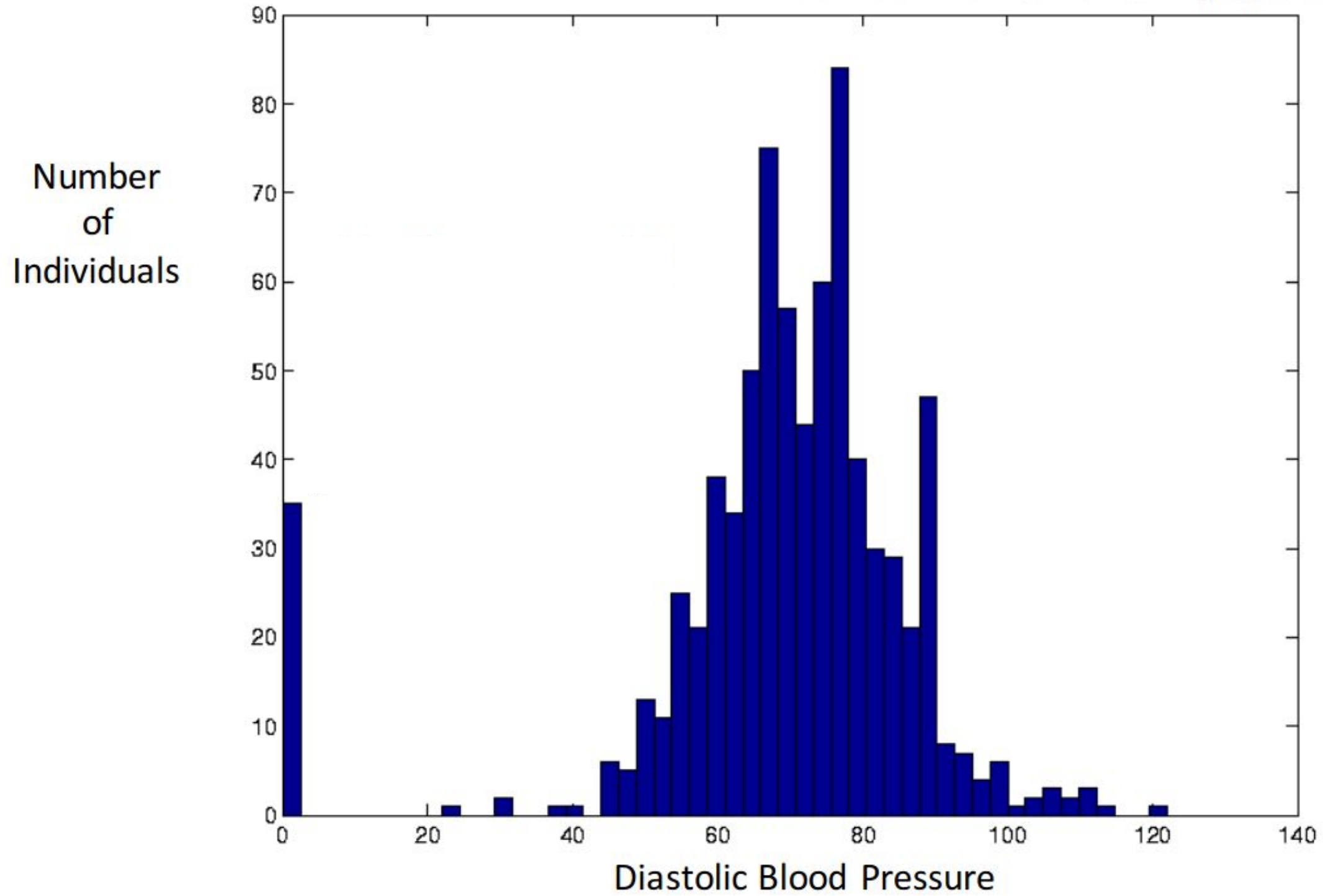
- Majority of time is not spent designing new ML algorithms
- It will be spent:
 - understanding the input data
 - defining the ML task
 - cleaning the data
 - designing features/attributes
 - visualising the data
 - comparing models (evaluating performance)
 - repeat

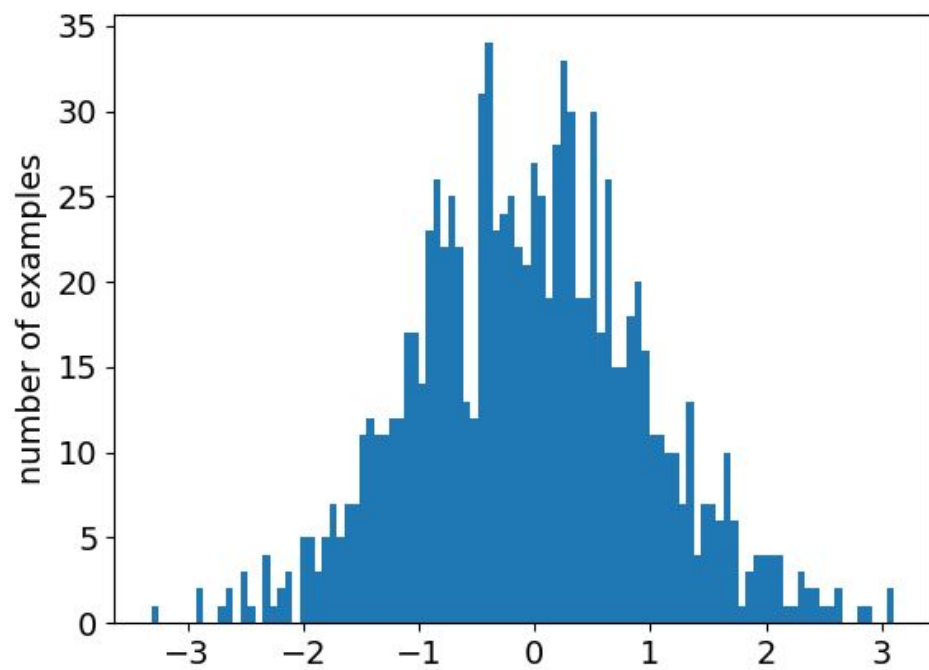
Visualisation

1) Crucial part of exploratory data analysis

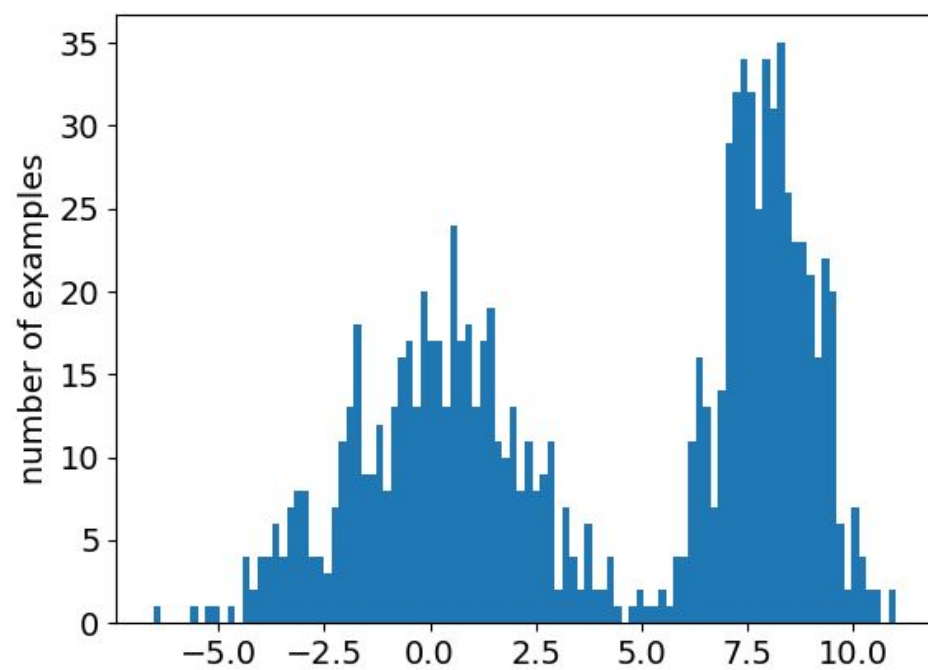
Identify problems/anomalies in your data

Pima Indians Diabetes Data,
From UC Irvine Machine Learning Repository

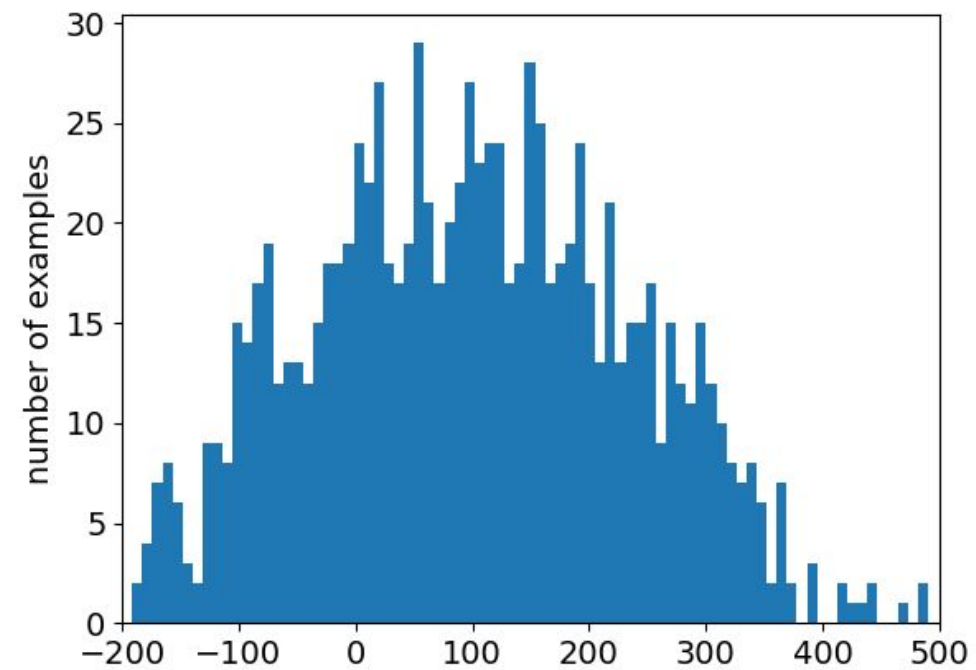




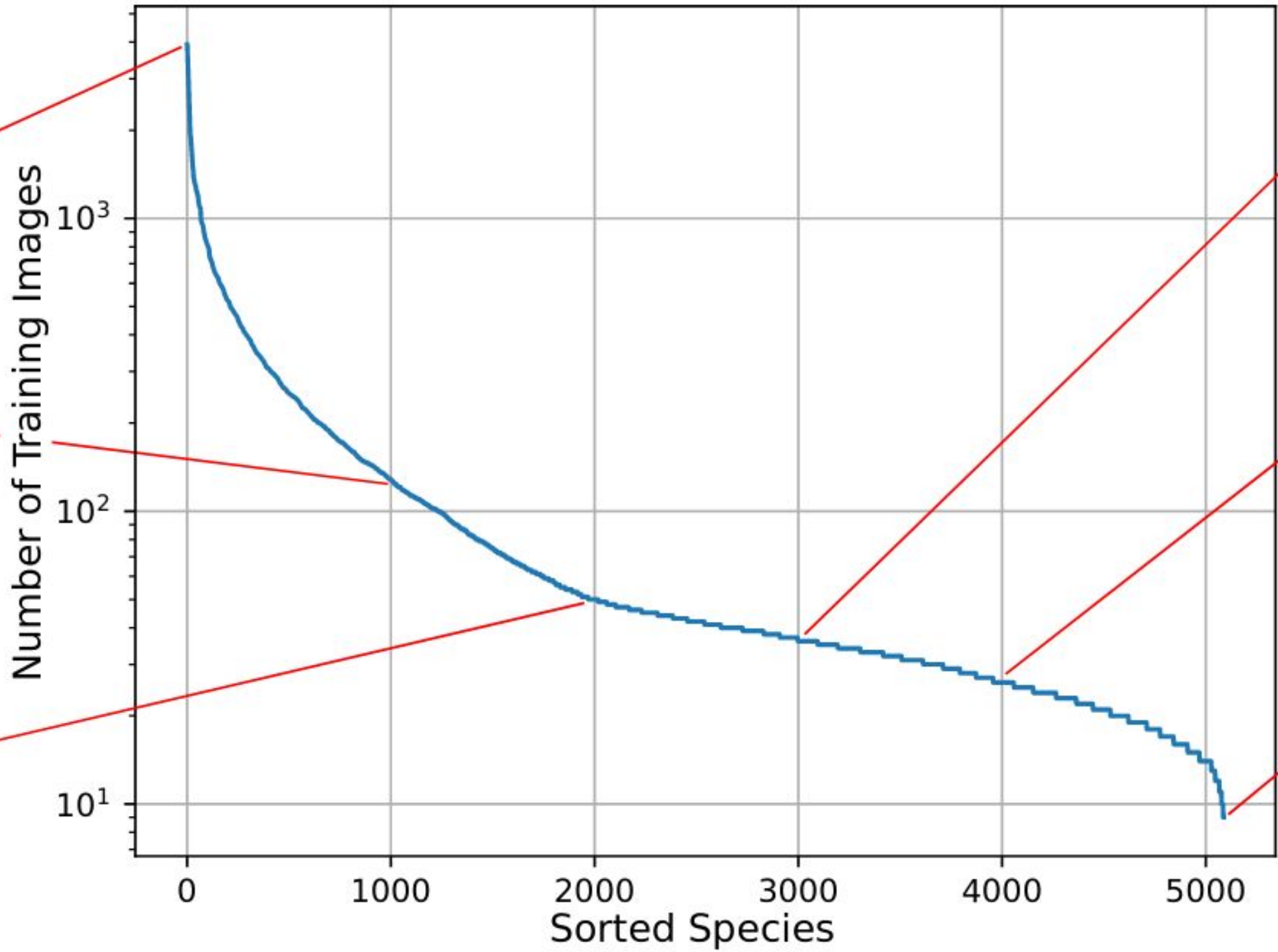
X_1



X_2



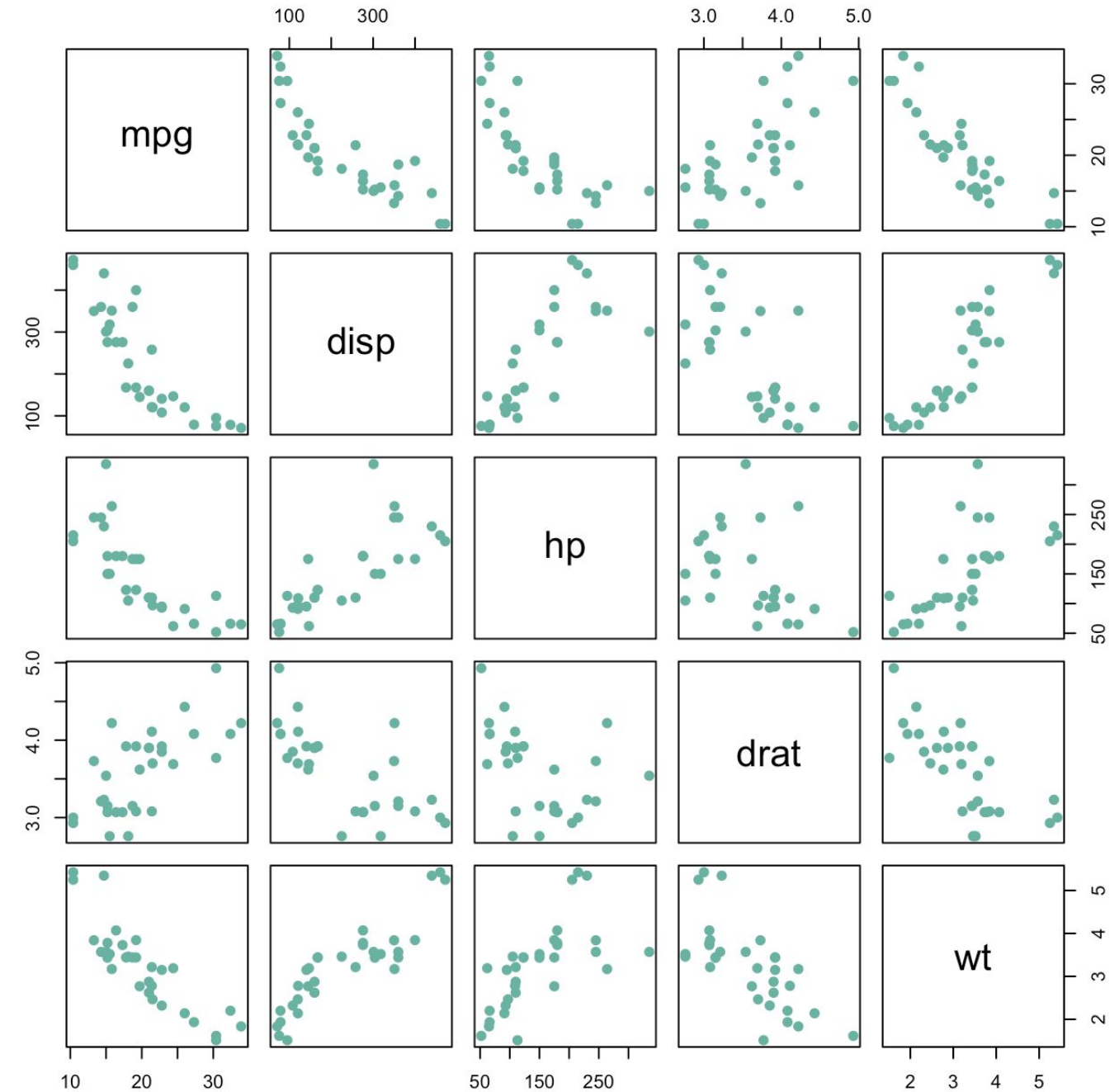
X_3



Higher Dimensions are Challenging

Visualization is essential but not scalable as the number of dimensions increase

Scatter plot matrix



Visualisation

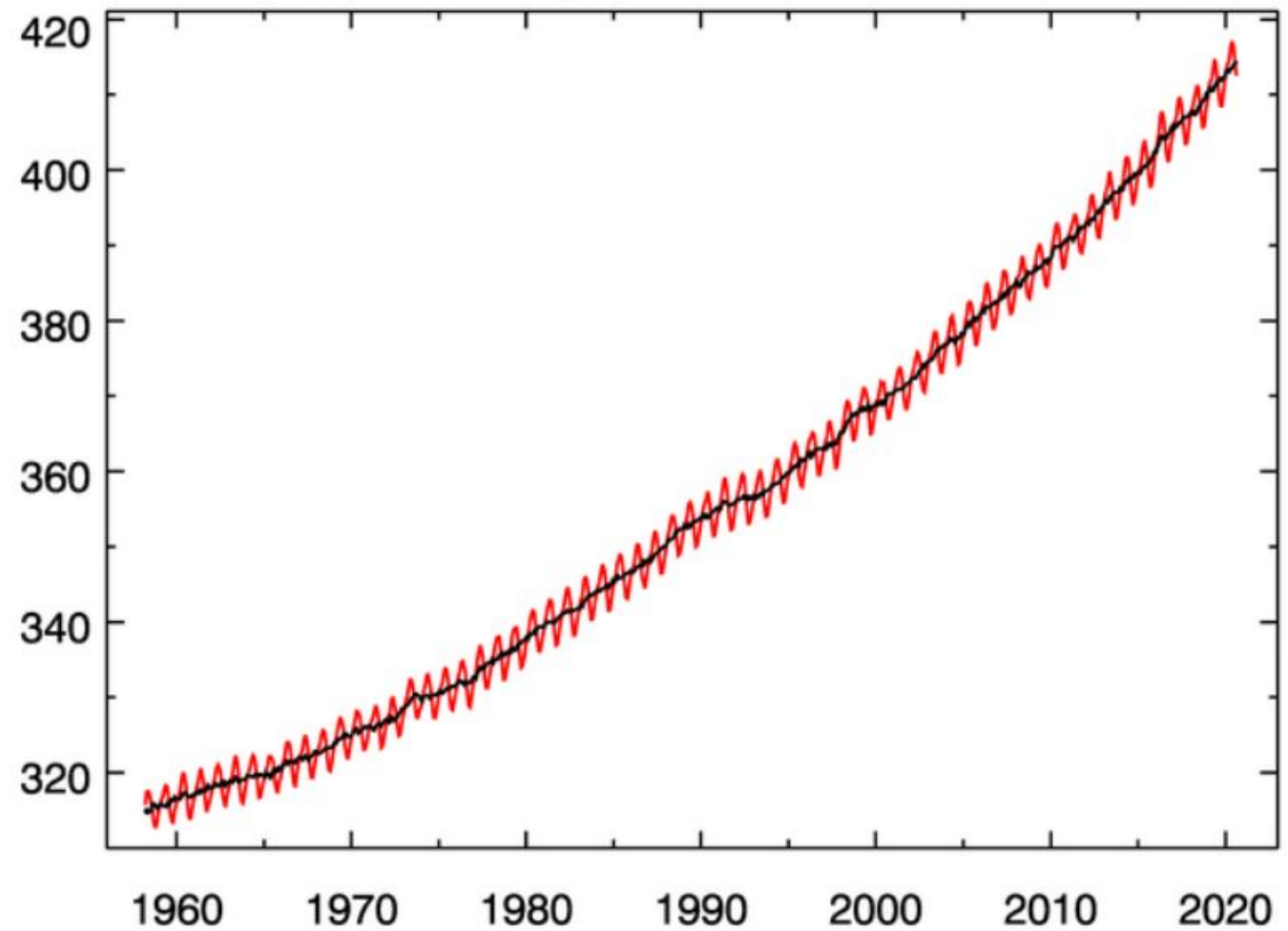
1) Crucial part of exploratory data analysis

Identify problems/anomalies in your data

1) Visualisation to present results

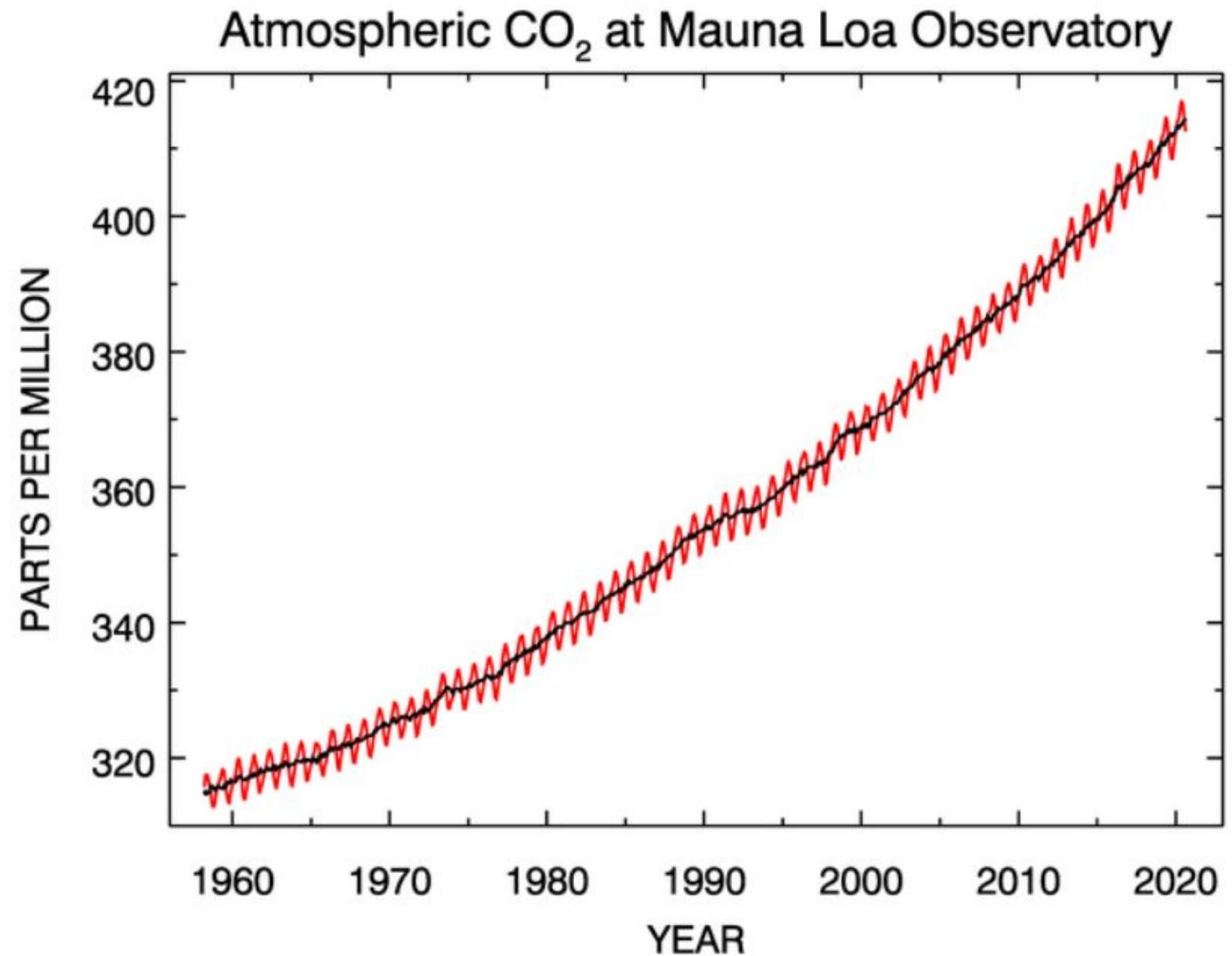
Convey message / summarise results

Common Visualization Mistakes

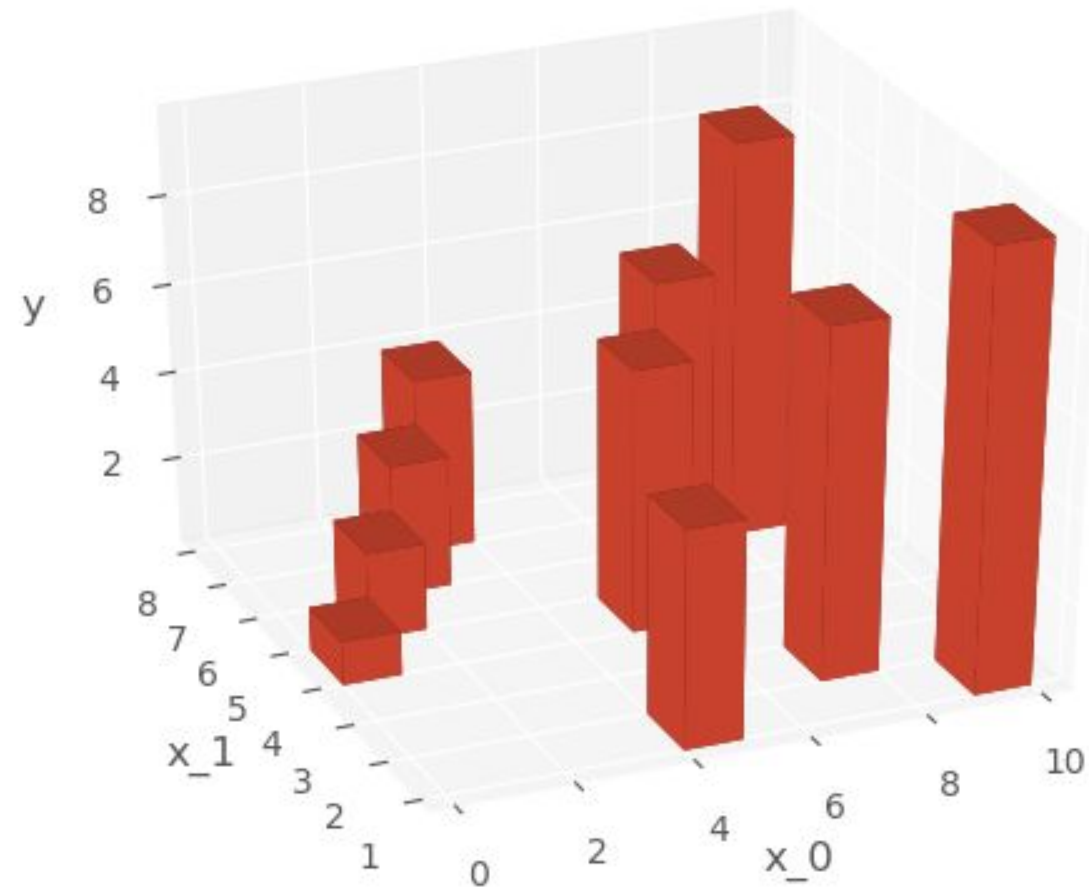


Common Visualization Mistakes

Label your axes!

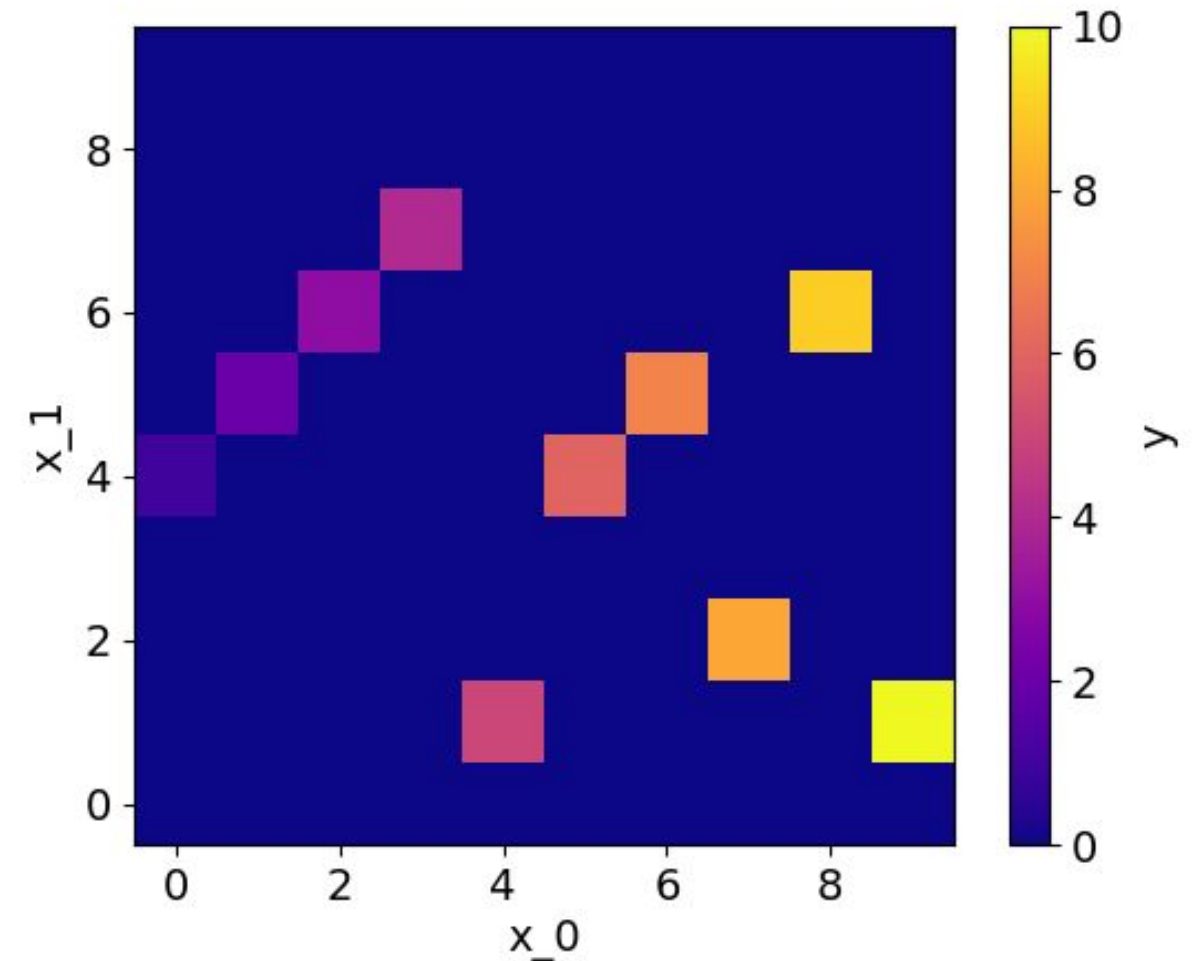
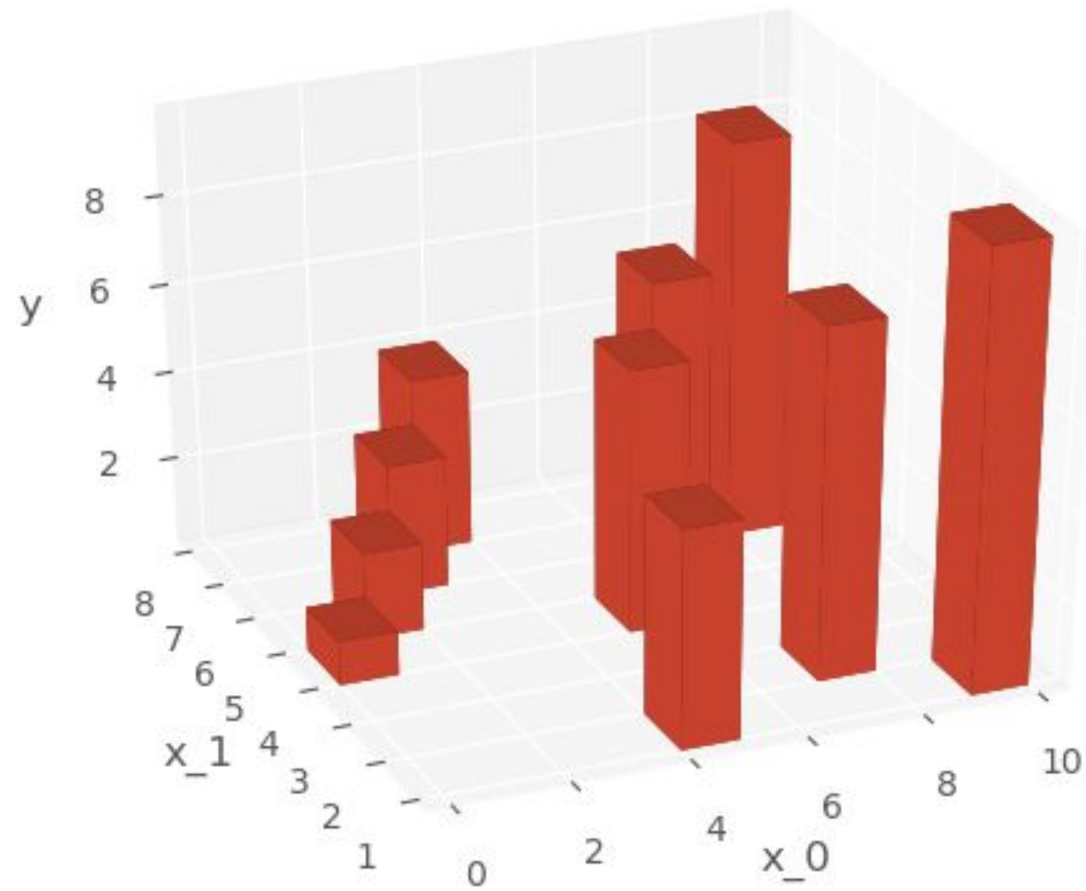


Common Visualization Mistakes



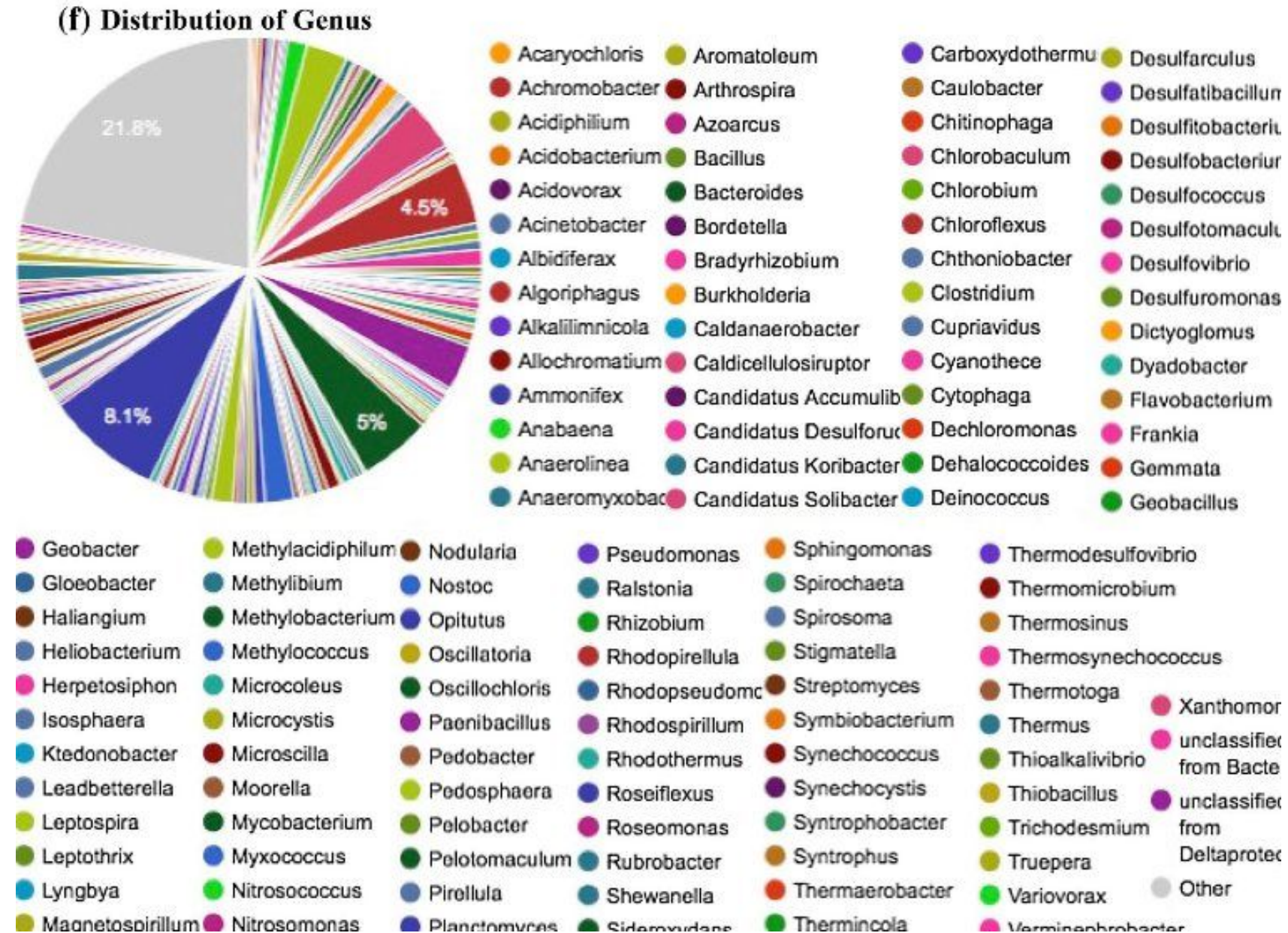
Common Visualization Mistakes

Don't overcomplicate things e.g. 3D is rarely a good idea - use colour

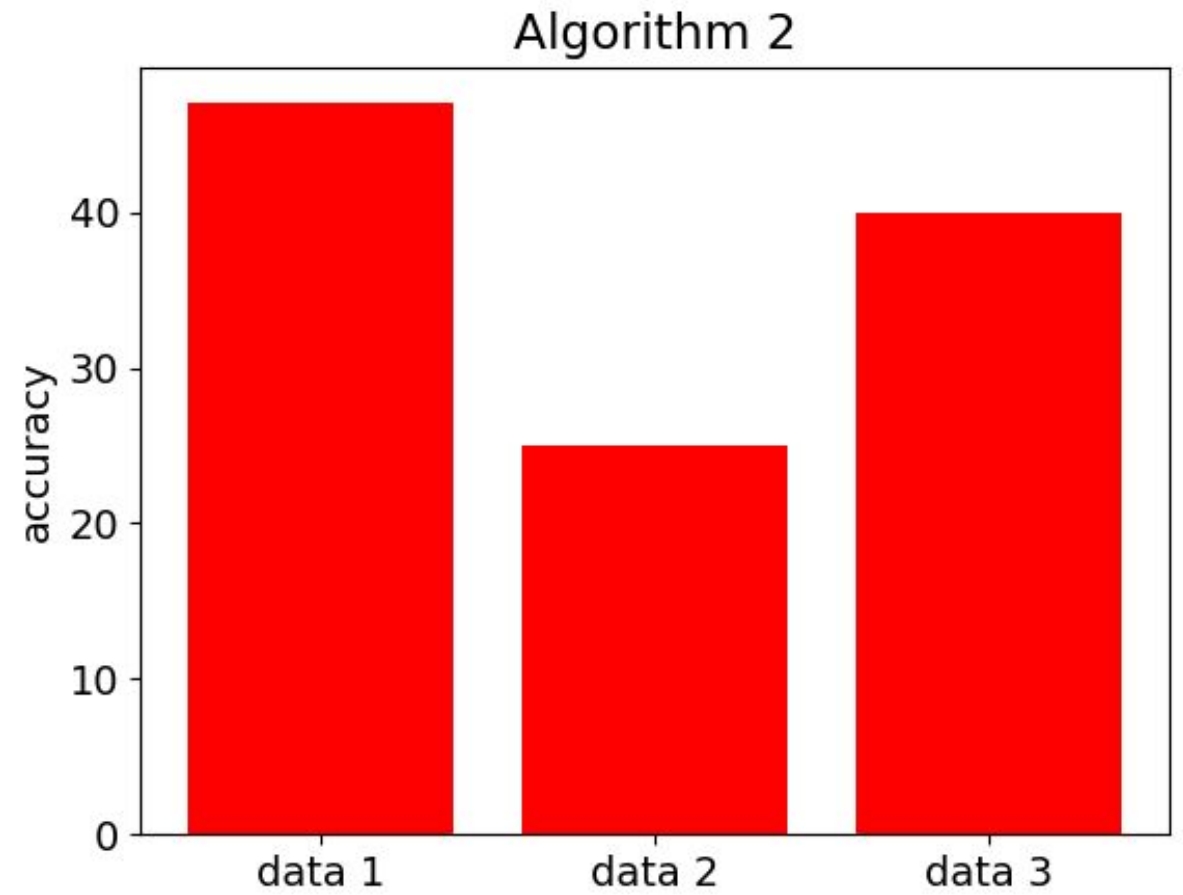
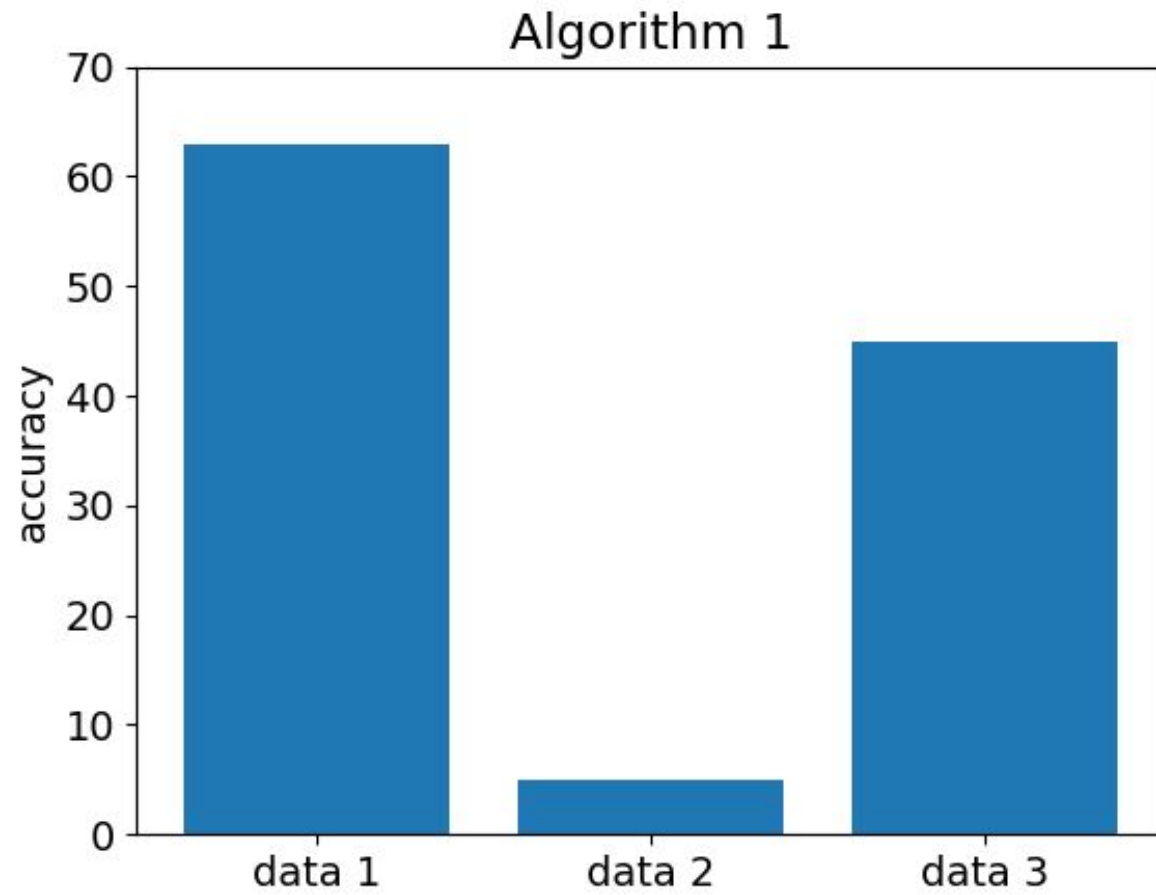


Common Visualization Mistakes

Keep it "clean"

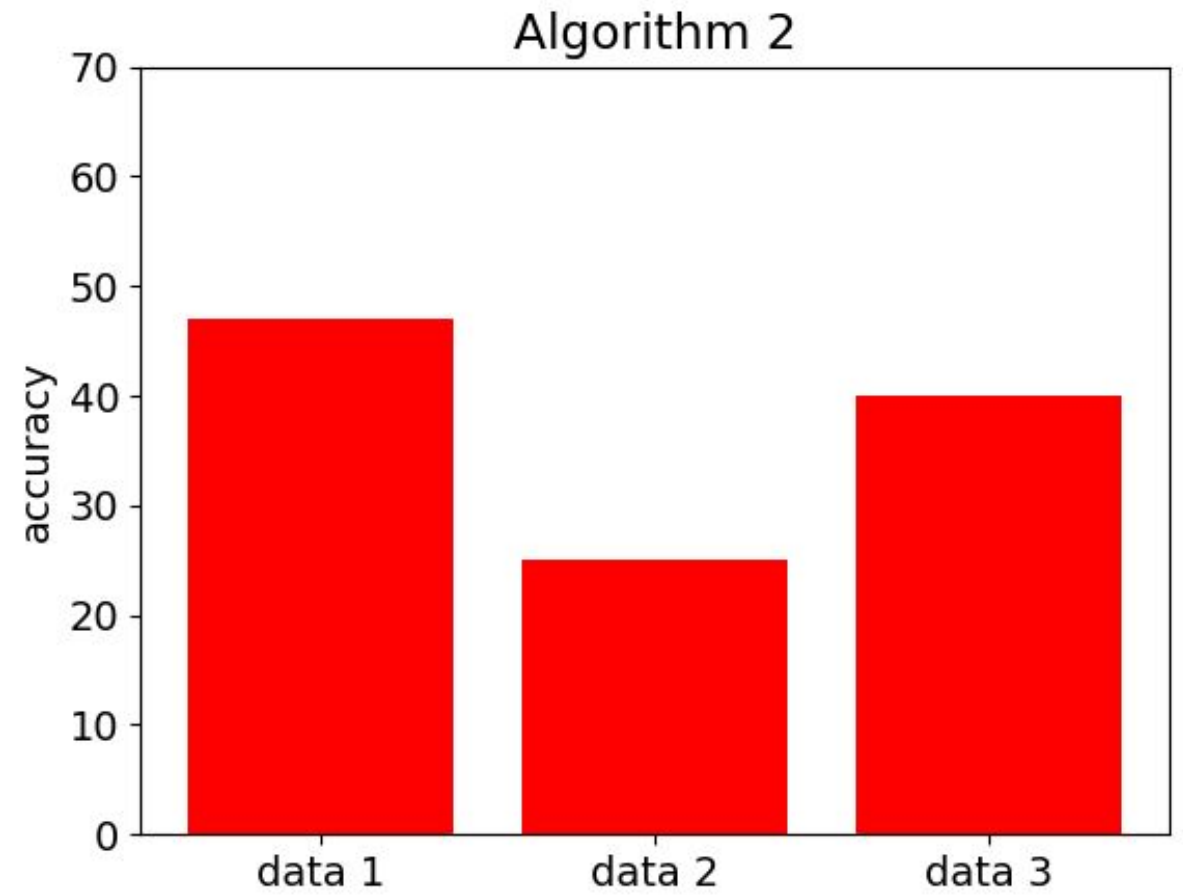
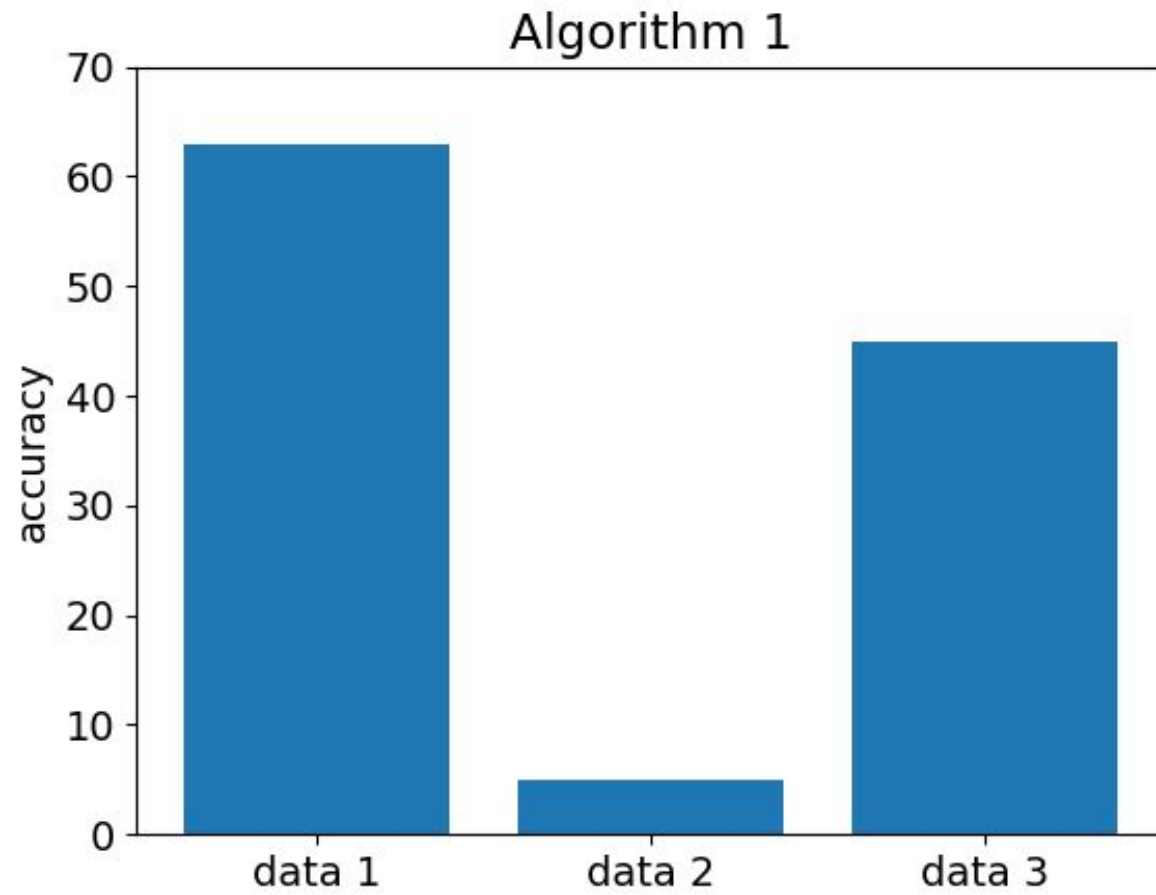


Common Visualization Mistakes



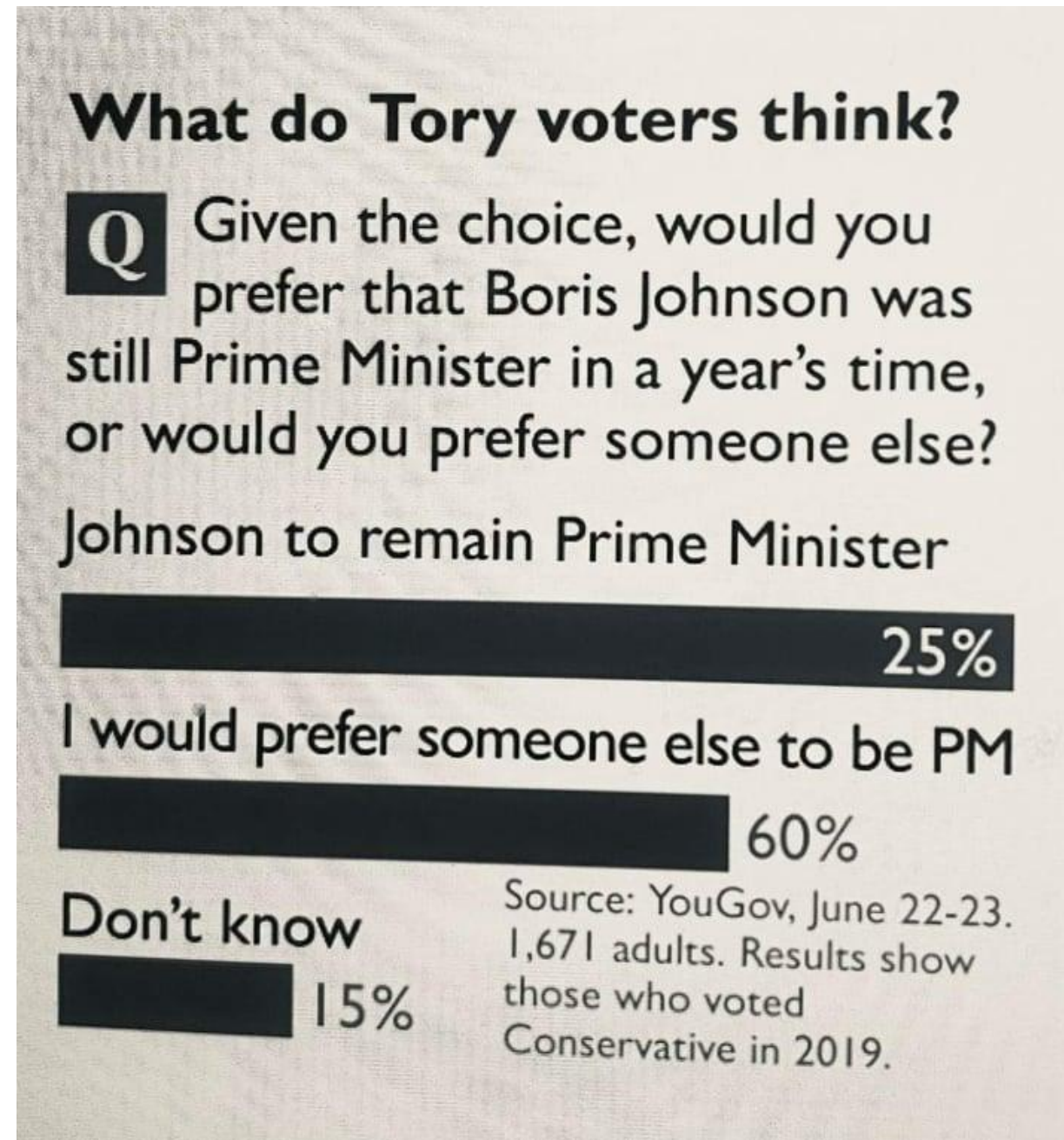
Common Visualization Mistakes

Don't mislead



Common Visualization Mistakes

Don't mislead



Extra Resources

GOV.UK:

<https://www.gov.uk/government/publications/a-bite-sized-guide-to-visualising-data-a-dstl-biscuit-book>



EU Publication Office:

<https://data.europa.eu/apps/data-visualisation-guide/>



Additional Links

Chartjunk - “decorations in a graphic that have no purpose”

https://www.edwardtufte.com/bboard/q-and-a-fetch-msg?msg_id=00040Z



Examples of bad plots used in media

<https://junkcharts.typepad.com>



Plots to avoid

https://genomicsclass.github.io/book/pages/plots_to_avoid.html



Classifying New Data - Ham Likelihood

- Given an *new* email we would like to be able to classify it
- For example, given the test email:
“review us now”
- $\mathbf{x}_t = [0, 1, 0, 1, 0, 0]^\top$

- Class priors:

$$p(\text{spam}) = 4/6 \quad p(\text{ham}) = 2/6$$

- Per-class likelihoods:

$p(x_d \text{spam})$	$p(x_d \text{ham})$	x_d
2/4	0/2	password
1/4	2/2	review
3/4	1/2	send
3/4	1/2	us
3/4	1/2	your
1/4	1/2	account

$$\begin{aligned} p(\mathbf{x}_t|\text{ham}) &= p(0, 1, 0, 1, 0, 0|\text{ham}) \\ &= \left(1 - \frac{0}{2}\right)\left(\frac{2}{2}\right)\left(1 - \frac{1}{2}\right)\left(\frac{1}{2}\right)\left(1 - \frac{1}{2}\right)\left(1 - \frac{1}{2}\right) = 0.0625 \end{aligned}$$

PCA: Maximising Variance

Recall Xv projects X onto v

$$\begin{aligned}\text{Var}[Xv] &= \frac{1}{N}(Xv)^\top(Xv) \\ &= \frac{1}{N}v^\top X^\top Xv \\ &= v^\top \frac{X^\top X}{N}v \\ &= v^\top Sv\end{aligned}$$

$$\max v^\top Sv, \text{ s.t. } v^\top v = 1$$

solved using *Lagrange multipliers* as

$$\max \underbrace{v^\top Sv - \lambda(v^\top v - 1)}_{\mathcal{L}}$$

computing derivative w.r.t v and setting = 0

$$\frac{d\mathcal{L}}{dv} = 2Sv - 2\lambda v = 0$$

$$Sv = \lambda v \quad \square$$

$v \rightarrow$ direction of max variance

$$Sv = \lambda v$$

left multiply by v^\top

$$\begin{aligned}v^\top Sv &= v^\top \lambda v \\ &= \lambda v^\top v \\ &= \lambda \quad \square\end{aligned}$$

$\lambda \rightarrow$ max variance

PCA: Finding Principal Components

More generally, solve for $SV = \Lambda V$ using Eigen decomposition

$$V = [\mathbf{v}_1, \dots, \mathbf{v}_D], \Lambda = \begin{bmatrix} \lambda_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_D \end{bmatrix} \quad \mathbf{v}_i \in \mathbb{R}^D, V \in \mathbb{R}^{D \times D}, \Lambda^{D \times D}$$

Eigenvalues

Solve $|S - \lambda I| = 0$

$$\begin{vmatrix} 2.0 - \lambda & 0.8 \\ 0.8 & 0.6 - \lambda \end{vmatrix} = 0$$

$$\lambda^2 - 2.6\lambda + 0.56 = 0$$

$$\implies \{\lambda_1, \lambda_2\} = \{2.36, 0.23\}$$

Eigenvectors

Find i^{th} eigenvector by solving $S\mathbf{v}_i = \lambda_i\mathbf{v}_i$

$$\begin{bmatrix} 2.0 & 0.8 \\ 0.8 & 0.6 \end{bmatrix} \begin{bmatrix} v_{1,1} \\ v_{1,2} \end{bmatrix} = 2.36 \begin{bmatrix} v_{1,1} \\ v_{1,2} \end{bmatrix} \implies \mathbf{v}_1 = \begin{bmatrix} 2.2 \\ 1 \end{bmatrix}$$
$$\begin{bmatrix} 2.0 & 0.8 \\ 0.8 & 0.6 \end{bmatrix} \begin{bmatrix} v_{2,1} \\ v_{2,2} \end{bmatrix} = 0.23 \begin{bmatrix} v_{2,1} \\ v_{2,2} \end{bmatrix} \implies \mathbf{v}_2 = \begin{bmatrix} -0.41 \\ 0.91 \end{bmatrix}$$

PCA: Picking number of dimensions

Given: eigenvectors $V = [v_1, \dots, v_D]$; Require: $M \ll D$

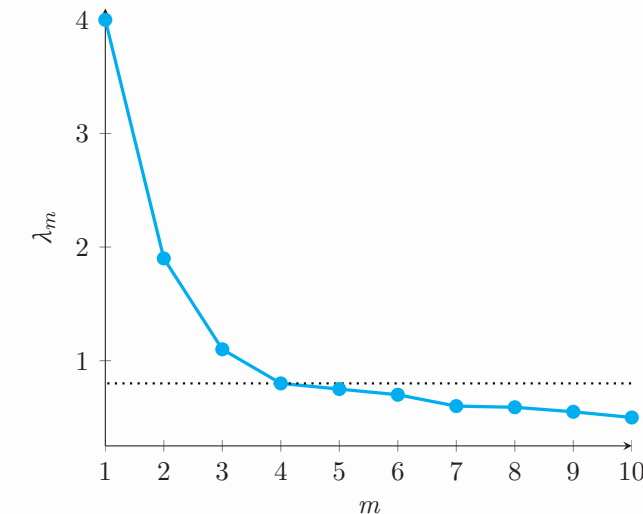
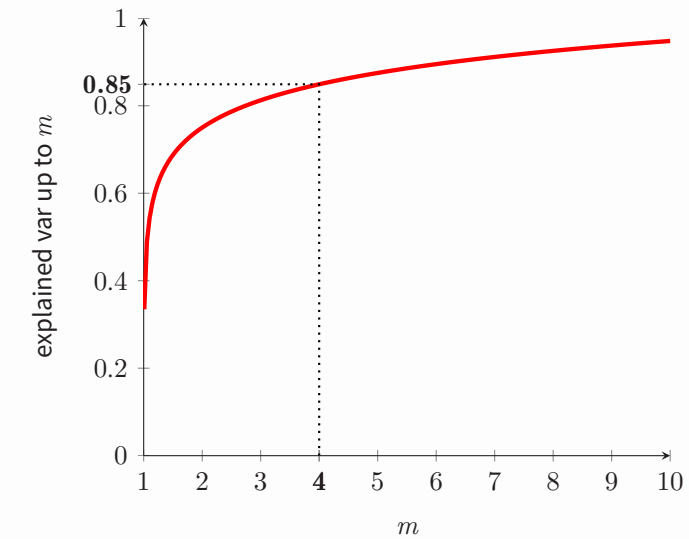
Known: eigenvalue $\lambda_i = \text{variance along } v_i$

Explained variance

- sort eigenvectors s.t. $\lambda_1 \geq \dots \geq \lambda_D$
- choose top M eigenvectors that explain “most” variance (typically 85%, 90%, or 95%)

Elbow plot

- plot eigenvalues in descending order $\lambda_1 \geq \dots \geq \lambda_D$
- choose point at which curve “bends” most (i.e. elbow)



PCA: Dimensionality Reduction

Let $V_M = [v_1, \dots, v_M] \in \mathbb{R}^{D \times M}$ denote the *truncated* eigenvector matrix for $M \ll D$

Reduction

Dimensionality reduction on data x_i

$$e_i^\top = x_i^\top V_M \in \mathbb{R}^M$$

More generally, projected data E

$$\begin{aligned} E &= [e_1^\top, \dots, e_N^\top] \\ &= [x_1^\top V_M, \dots, x_N^\top V_M] \\ &= X V_M \in \mathbb{R}^{N \times M} \end{aligned}$$

Reconstruction

Recover data \hat{x}_i from e_i using V_M^\top

$$\hat{x}_i^\top = e_i^\top V_M^\top = (x_i^\top V_M) V_M^\top \in \mathbb{R}^D$$

More generally, reconstructed data \hat{X}

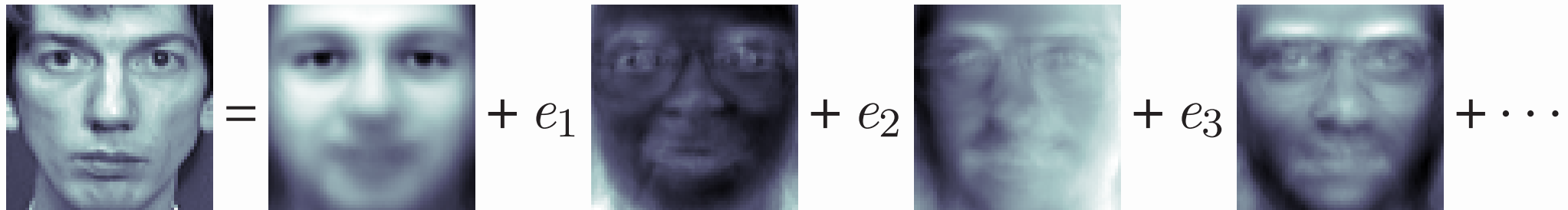
$$\begin{aligned} \hat{X} &= [\hat{x}_1^\top, \dots, \hat{x}_N^\top] \\ &= X V_M V_M^\top \in \mathbb{R}^{N \times D} \end{aligned}$$

$V_M V_M^\top \in \mathbb{R}^{D \times D}$ is the data *projection* matrix

PCA Example 2: Eigenfaces

Projection

Projecting face x_i onto $e_i = [e_{i1}, \dots, e_{iM}]$



Reconstruction

Reconstructing face \hat{x}_i using M components



(90 \ll 4096!)

$M = 10$ $M = 30$ $M = 50$ $M = 70$ $M = 90$

The MovieFlix Corporation wants to improve the algorithm it uses to recommend new movies to its customers. It has decided to release a portion of its data to the public in the hope that researchers will be able to find patterns in users' preferences for the different movies. The data is represented as a matrix X where rows $i = 1, \dots, n$ correspond to customers, and columns $j = 1, \dots, c$ correspond to movies, and X_{ij} is a numeric value that reflects the degree to which customer i enjoyed movie j . Assume the data is complete; i.e., we know the rating of every user for every movie.

MovieFlix decides to carry out PCA on X treating each row (i.e. each customer) as a data vector. How will this help them understand the patterns in the preference data?

PCA relationship between X and X^T

Instead of thinking about the data as rows of customer responses, we can think of a movie as the response it elicits from N customers.



PCA relationship between X and X^T

Instead of thinking about the data as rows of customer responses, we can think of a movie as the response it elicits from N customers.

recall $Sv = \lambda v$

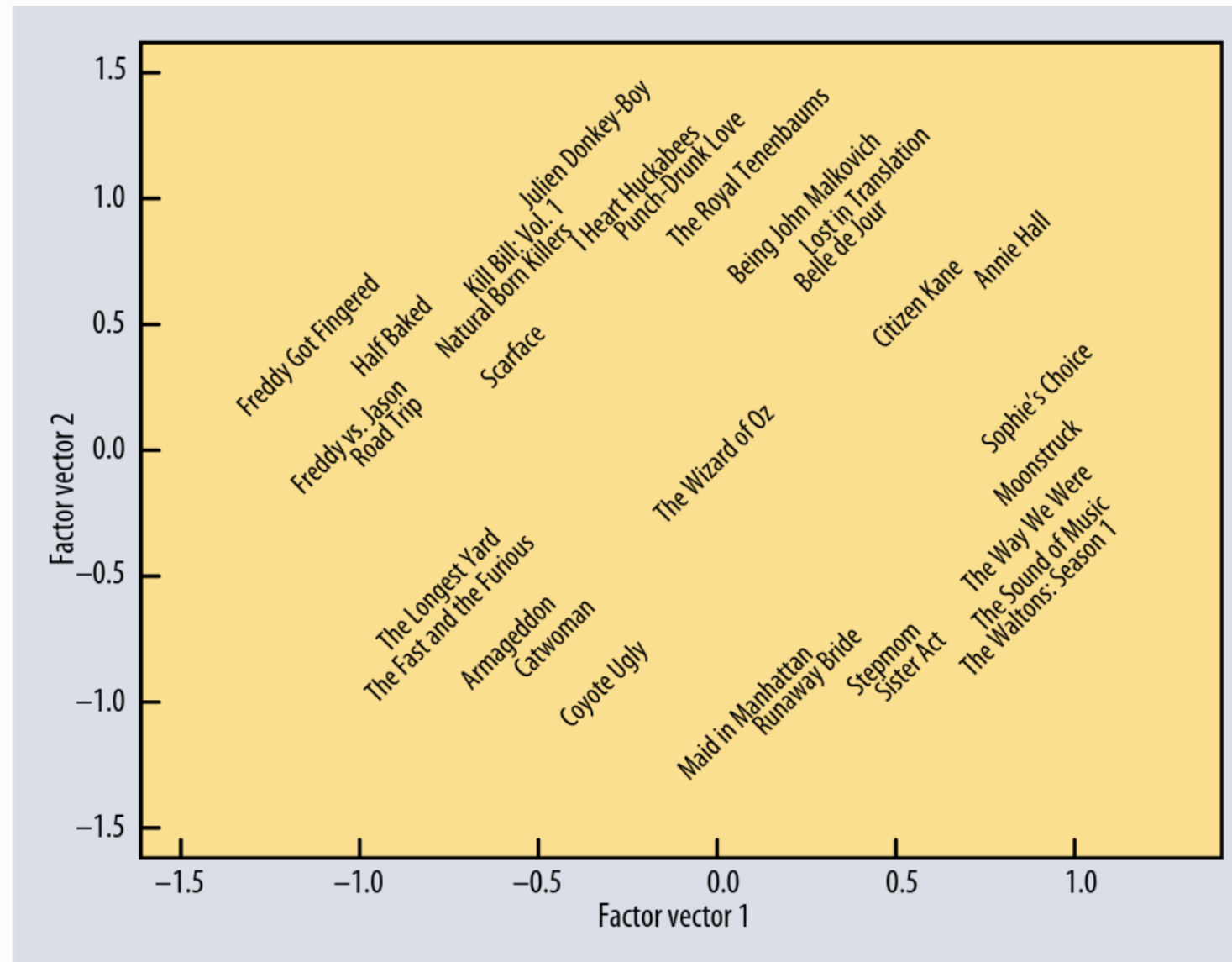
$$\left(\frac{1}{N}x^T x\right)v = \lambda v \quad (\because S = \frac{1}{N}x^T x)$$

pre-multiply by x

$$\left(\frac{1}{N}x x^T\right)(xv) = \lambda(xv) \Rightarrow S'v' = \lambda v'$$

$S' = \frac{1}{N}(x^T)^T x^T = \frac{1}{N}x x^T$ v' new pc direction \hookrightarrow same eigenvalues!

In the famous Netflix challenge, a similar low-rank idea can be used as in PCA, but almost all of the entries in X are missing—not every person has seen every film!



credit: Yehuda Koren, Robert Bell, & Chris Volinsky, Matrix Factorization Techniques for Recommender Systems