

Applied Machine Learning (AML)

Exploratory Data Analysis

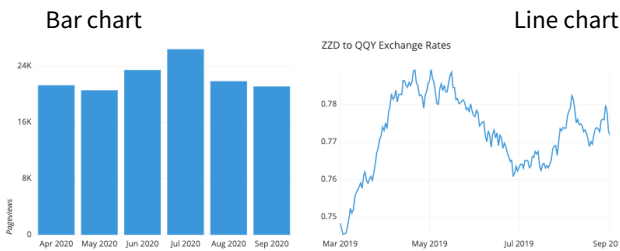
Data Visualisation

Oisin Mac Aodha • Siddharth N.

Plotting Data

Plot Types

- temporal change
- part-to-whole composition
- distribution
- group comparison
- inter-variable relations

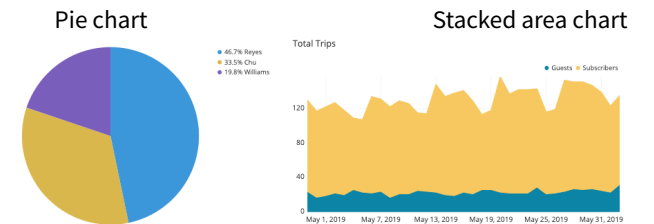


Figures: chartio.com

Plotting Data

Plot Types

- temporal change
- part-to-whole composition
- distribution
- group comparison
- inter-variable relations



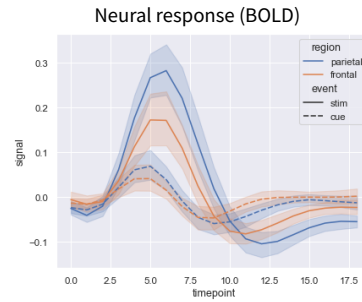
Figures: chartio.com

Features of a good plot

- title
- labelled axes
- axes ranges and ticks
- clarity (colour/thickness)
- legend

informative:
convey as much as necessary

clean:
avoid overfilling & redundancy



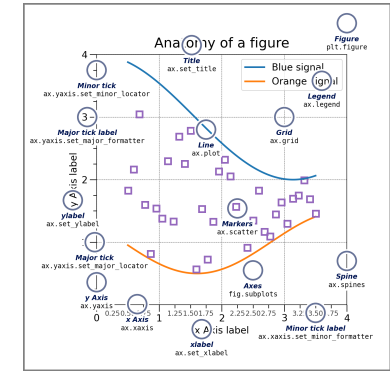
Figures: Matplotlib—atomy of a figure

Features of a good plot

- title
- labelled axes
- axes ranges and ticks
- clarity (colour/thickness)
- legend

informative:
convey as much as necessary

clean:
avoid overfilling & redundancy



Figures: Matplotlib—atomy of a figure

Features of a good plot

- title
- labelled axes
- axes ranges and ticks
- clarity (colour/thickness)
- legend

informative:
convey as much as necessary

clean:
avoid overfilling & redundancy

Relatively easy to think about
when data is low dimensional

What do we do when data
is high dimensional?

Dimensionality Reduction

Figures: Matplotlib—atomy of a figure

Curse of Dimensionality

Manifold Hypothesis

High-dimensional data in the real world really lies on low-dimensional manifolds within that high-dimensional space.

- Data is typically high dimensional
vision: 10^4 pixels, text: 10^6 words
- Example: handwritten digits (MNIST)
 - 28×28 pixels $\rightarrow \{0, 1\}^{784}$ possible “images”
 - only a very small number of these images are actually real
 - true dimensionality: actual variation of pen strokes!



Dealing with high dimensionality

Statistics

- ML involves some form of “counting” observations and features
 - count within some regions
e.g. constructing histograms
 - use counts to construct predictors
e.g. decision trees
- As dimensionality grows, fewer observations per region

Mitigation

- domain knowledge / feature engineering
- modelling assumptions about features independence, smoothness, symmetry
- reduce data dimensionality
construct a new set of dimensions / variables

Dimensionality Reduction

Goal: Represent data using a “few” variables

- compression: preserve as much information/structure as possible
- discrimination: only keep information that enables task (e.g. classification)

Selection

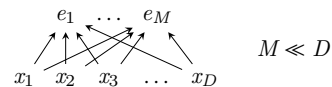
- subset of all features

$$x_1, x_2, x_3, \dots, x_{D-1}, x_D$$

- relevant to task
e.g. ‘credit history’ \rightarrow loan?

Transformation

- construct a new set of dimensions



- transformation of original
e.g. linear $F \Rightarrow e = Fx$

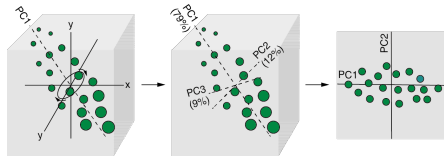
Dimensionality Reduction

PCA

Principal Components Analysis (PCA)

Define principal components (PCs)

- 1st PC: direction of *greatest* variation in the data
- 2st PC: \perp 1st PC; greatest *remaining* variation
- ...and so on until D , for $x \in \mathbb{R}^D$.
- First $M \ll D$ components \rightarrow new basis dimensions
- ...transform coordinates of each data point to new basis



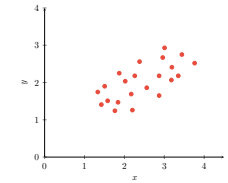
Rationale

- variation along direction = *information*
- transform basis \rightarrow fit maximum information into M dimensions

PCA: Basics

$$X = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix} \quad X \in \mathbb{R}^{N \times D}, \mathbf{x}_i \in \mathbb{R}^D \quad (\text{data})$$

$$S = \frac{1}{N} X^T X \quad S \in \mathbb{R}^{D \times D} \quad (\text{covariance, assuming 0-mean})$$

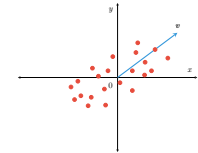


Intuition

Repeated transformation using the covariance (S) turns towards direction of maximum variance (example)

$$Sv = \begin{bmatrix} 2.0 & 0.8 \\ 0.8 & 0.6 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \begin{bmatrix} -1.2 \\ 0.2 \end{bmatrix} \stackrel{S}{=} \dots \stackrel{S}{=} \begin{bmatrix} -14.1 \\ -6.4 \end{bmatrix} \stackrel{S}{=} \begin{bmatrix} -33.3 \\ -15.1 \end{bmatrix}$$

where the *slope* converges to 0.454



Goal: Find v such that $Sv = \lambda v$

PCA: Maximising Variance

Recall Xv projects X onto v

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

$\max v^T S v$, s.t. $v^T v = 1$
solved using *Lagrange multipliers* as

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

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^T

$$v^T S v = v^T \lambda v$$

$$= \lambda v^T v$$

$$= \lambda \quad \square$$

$\lambda \rightarrow$ max variance

PCA: Finding Principal Components

More generally, solve for $Sv = \lambda v$ using Eigen decomposition

$$V = [v_1, \dots, v_D], \Lambda = \begin{bmatrix} \lambda_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_D \end{bmatrix} \quad v_i \in \mathbb{R}^D, V \in \mathbb{R}^{D \times D}, \Lambda^{D \times D}$$

Eigenvalues

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

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

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

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

Eigenvectors

Find i^{th} eigenvector by solving $Sv_i = \lambda_i v_i$

$$\begin{bmatrix} 2.0 & 0.8 \\ 0.8 & 0.6 \end{bmatrix} \begin{bmatrix} v_{1,1} \\ v_{1,2} \end{bmatrix} = 2.36 \begin{bmatrix} v_{1,1} \\ v_{1,2} \end{bmatrix} \Rightarrow v_1 = \begin{bmatrix} 2.2 \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} 2.0 & 0.8 \\ 0.8 & 0.6 \end{bmatrix} \begin{bmatrix} v_{2,1} \\ v_{2,2} \end{bmatrix} = 0.23 \begin{bmatrix} v_{2,1} \\ v_{2,2} \end{bmatrix} \Rightarrow v_2 = \begin{bmatrix} -0.41 \\ 0.91 \end{bmatrix}$$

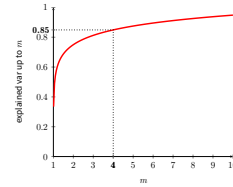
PCA: Picking number of dimensions

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

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

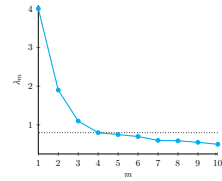
Explained variance

- 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

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

More generally, projected data E

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

Reconstruction

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

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

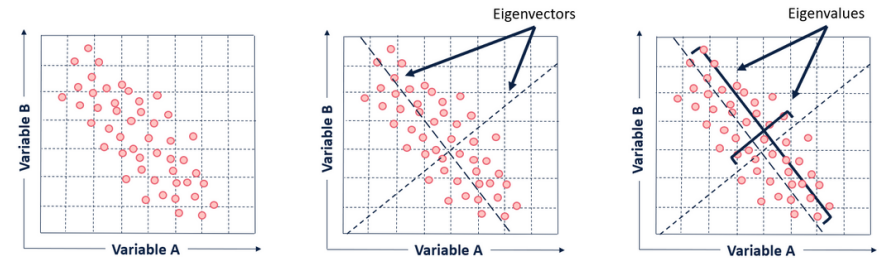
More generally, reconstructed data \hat{X}

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

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

PCA: Overview and Use

Characteristics

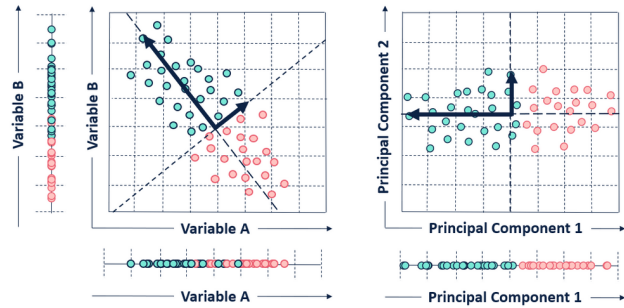


Dimensionality Reduction

PCA: Examples

PCA: Overview and Use

Use: Classification



Figures: Sydney Firmin @ towardsdatascience.com

PCA Example 1: UK Food Consumption

$$X \in \mathbb{R}^{4 \times 17}$$

	England	N Ireland	Scotland	Wales
Alcoholic drinks	375	135	458	475
Beverages	57	47	53	73
Carcase meat	245	267	242	227
Cereals	1472	1494	1462	1582
Cheese	105	66	103	103
Confectionery	54	41	62	64
Fats and oils	193	209	184	235
Fish	147	93	122	160
Fresh fruit	102	674	957	1137
Fresh potatoes	720	1033	566	874
Fresh Veg	253	143	171	265
Other meat	685	586	750	803
Other Veg	488	355	418	570
Processed potatoes	198	187	220	203
Processed Veg	360	334	337	365
Soft drinks	1374	1508	1572	1256
Sugars	156	139	147	175

Projecting to 1 component (V_1)



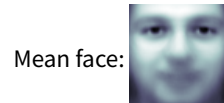
Projecting to 2 components (V_2)



Figures: setosa.io Data: Mark Richardson

PCA Example 2: Eigenfaces

Data $X \in \mathbb{R}^{300 \times 4096}$
 Image $x \in \mathbb{R}^{64 \times 64}$ is flattened to \mathbb{R}^{4096}



Principal Component Faces:



PCA Example 2: Eigenfaces

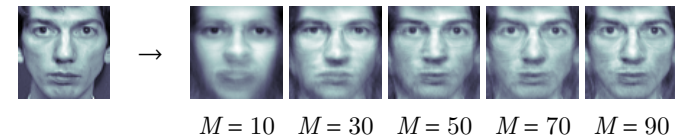
Projection

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

$$x_i = \text{Mean Face} + e_1 + e_2 + e_3 + \dots$$

Reconstruction

Reconstructing face \hat{x}_i using M components



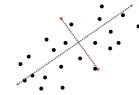
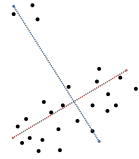
(90 \ll 4096!)

PCA: Limitations

Sensitivity

- outliers or scaling dimensions
 - changes variance along dimension
 - changes principal components
- **fix:** normalise—zero mean unit variance

$$x' = \frac{x - \mu}{\sigma}$$
- find outliers using interquartile range (IQR)
 - 'spread' of middle 50% of values
 - median(upper quartile) - median(lower quartile)
 - define 'outlier' as values $> 1.5 * IQR$

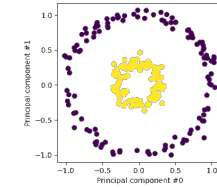
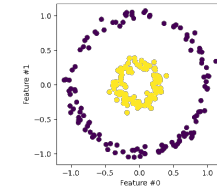
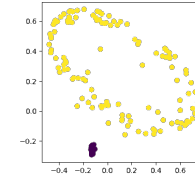


Removing outliers

PCA: Limitations

Linearity

- 1D: line; 2D: plane
- transform to handle non-linearity



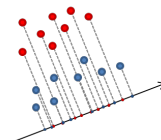
PCA: Limitations

Unsupervised

- maximises data variance along few directions
- ignorant of class labels
- could be hard to separate classes



Data



Projection on PC1

EDA: Summary

- Broad range of visualisation types
- Need to think about what information goes into a visualisation
- Actual data dimensionality \ll observed dimensionality
- For high-dimensional data
 - domain knowledge / feature engineering
 - modelling assumption: independence / smoothness / symmetry etc.
 - dimensionality reduction: selection / transformation
- Principal Components Analysis (PCA)
 - choose directions that maximise variation (eigenvectors)
 - for smaller number of components M , pack information
 - examples: UK food consumption, Eigenfaces