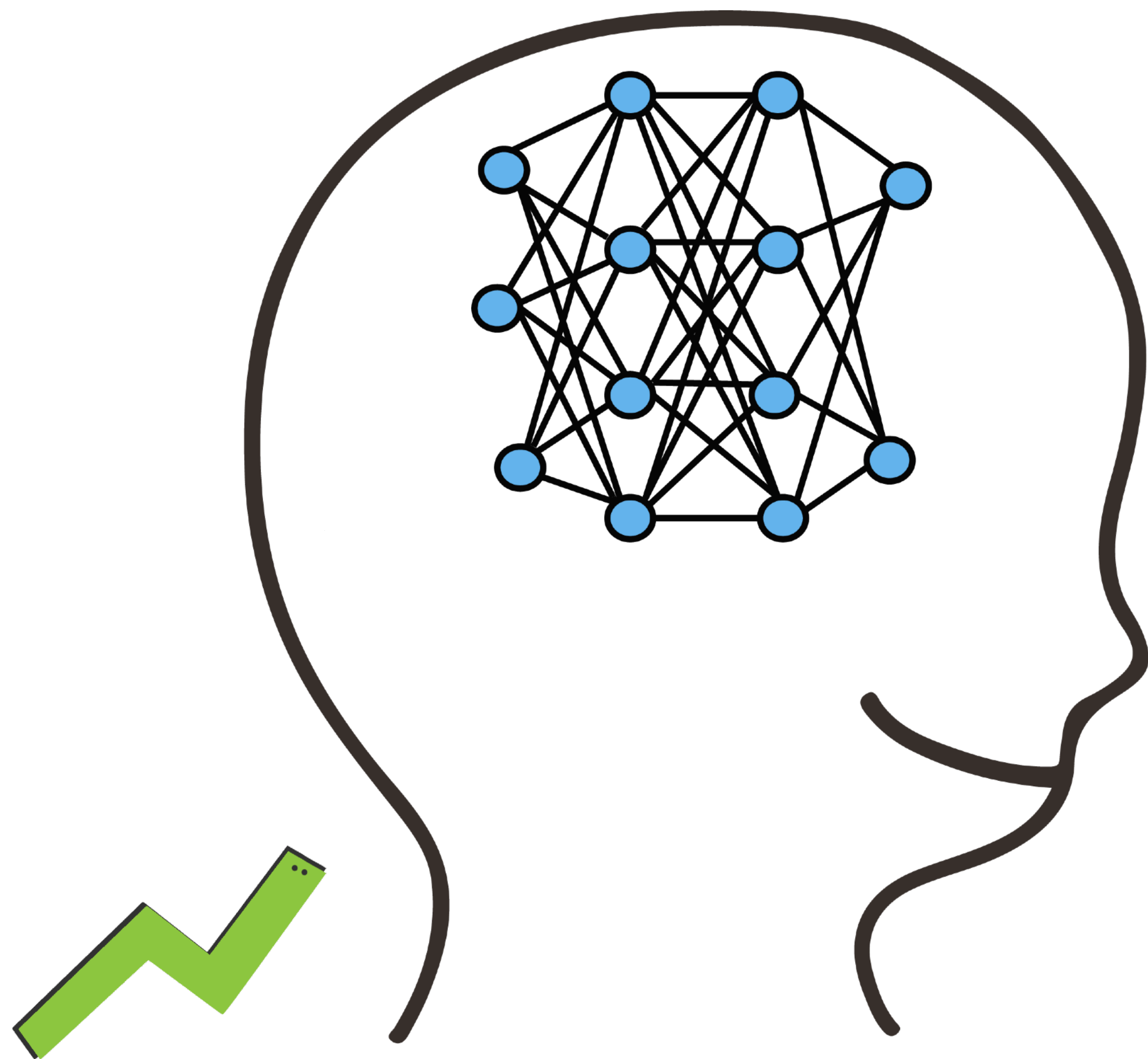


資料分析方法(一)

Data Analytics Methods - Part1



統計分析

Statistical Analysis



描述性統計 (Descriptive Statistics)

- 或稱「敘述性統計」
- 集中量數：呈現資料集中的情形，如：(算術)平均、中位數、眾數等
- 變異量數：呈現資料分散的情形，如：全距（最大值 - 最小值）、標準差、四分位數等

‣ DataFrame.describe()

| | A | B | C | D |
|-------|----------|----------|----------|----------|
| count | 6.000000 | 6.000000 | 6.000000 | 6.000000 |
| mean | 0.473862 | 0.615370 | 0.568419 | 0.622193 |
| std | 0.252262 | 0.312380 | 0.164988 | 0.329959 |
| min | 0.080301 | 0.202910 | 0.317583 | 0.047279 |
| 25% | 0.361322 | 0.365285 | 0.462811 | 0.507737 |
| 50% | 0.474192 | 0.685850 | 0.660573 | 0.725291 |
| 75% | 0.684415 | 0.851957 | 0.677213 | 0.809227 |
| max | 0.736302 | 0.951857 | 0.692137 | 0.962875 |



幾何平均數 (Geometric Mean)

- 適用於計算比率數據的變化率

$$G = \sqrt[n]{\prod_{i=1}^n x_i} = \sqrt[n]{x_1 \cdot x_2 \cdots x_n}$$

- e.g. 營業額成長：12%, 15%, -4%, -10%, 6%
- `scipy.stats.gmean([1.12, 1.15, 0.96, 0.9, 1.06]) => 1.04 (4%)`



調和平均數 (Harmonic Mean)

- 數值倒數的算術平均數的倒數，又稱為「倒數平均數」。

$$H = \frac{n}{\frac{1}{x_1} + \frac{1}{x_2} + \dots + \frac{1}{x_n}}$$

- e.g. 台北到高雄坐高鐵平均時速300公里、高雄到台北坐台鐵普悠瑪號平均時速150公里，全程平均時速是多少？
- `scipy.stats.hmean([300, 150]) => 200`



截尾平均數 (Trimmed Mean)

- 平均數容易受到極端值影響
- 截尾平均數會將極端值去除後再取算術平均
 - 自訂上下限： `scipy.stats.tmean(array-like data, (lower limit, upper limit))`
 - 截尾後的標準差 (tstd) 、變異數 (tvar) 、最大值 (tmax) 、最小值 (tmin)
 - 依比例去除： `scipy.stats.trim_mean(array-like data, proportiontocut)`

四分位數 (Quartile)

- 將數據從小到大排列
 - 第一四分位數 (Q_1) : 在 $1/4$ 位置的數，又稱「較小四分位數」
 - 第二四分位數 (Q_2) : 在 $1/2$ 位置的數，又稱「中位數」
 - 第三四分位數 (Q_3) : 在 $3/4$ 位置的數，又稱「較大四分位數」
 - 四分位距 (IQR) = $Q_3 - Q_1$
 - e.g. 1, 2, 3, 4, 5, 6, 7, 8
 - 內插法 : $Q_1 = 2.75$ 、 $Q_2 = 4.5$ 、 $Q_3 = 6.25$

Notes

► 四分位數的計算方法爭議

四分位數確切的數值計算方法仍具爭議
Scipy和Pandas計算出的值有少許誤差。



基本統計函式表

| 統計函式 | | Pandas DataFrame | Scipy.stats | Numpy |
|------|-----|---------------------------------------|--|---|
| 敘述統計 | | DataFrame.describe() | describe(<i>data</i>) | x |
| 算術 | 平均數 | DataFrame.mean() | x | mean(<i>data</i>) |
| 幾何 | | x | gmean(<i>data</i>) | x |
| 調和 | | x | hmean(<i>data</i>) | x |
| 截尾 | | x | trim_mean(<i>data</i> , <i>proportiontocut</i>) tmean(<i>data</i> , (<i>lower limit</i> , <i>upper limit</i>)) | x |
| 加權 | | x | x | average(<i>data</i> , <i>weights</i>) |
| 截尾 | 最大值 | x | tmax(<i>data</i> , (<i>lower limit</i> , <i>upper limit</i>)) | x |
| | 最小值 | x | tmin(<i>data</i> , (<i>lower limit</i> , <i>upper limit</i>)) | x |
| | 標準差 | x | tstd(<i>data</i> , (<i>lower limit</i> , <i>upper limit</i>)) | x |
| | 變異數 | x | tvar(<i>data</i> , (<i>lower limit</i> , <i>upper limit</i>)) | x |
| 最大值 | | DataFrame.max() | x | max(<i>data</i>) |
| 最小值 | | DataFrame.min() | x | min(<i>data</i>) |
| 中位數 | | DataFrame.median() | x | median(<i>data</i>) |
| 標準差 | | DataFrame.std() | x | std(<i>data</i>) |
| 變異數 | | DataFrame.var() | x | var(<i>data</i>) |
| 四分位數 | | DataFrame.quantile(<i>quantile</i>) | mstats.mquantiles(<i>data</i>) | x |
| 眾數 | | DataFrame.mode() (0.19.1版) | mode(<i>data</i>) 8 (0.18.1版) | x (1.11版) |

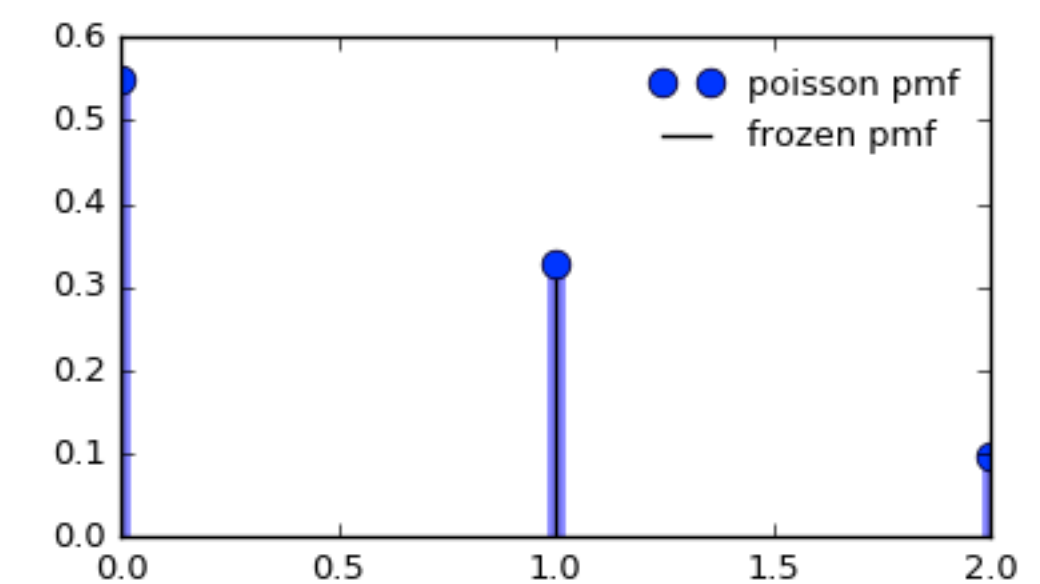
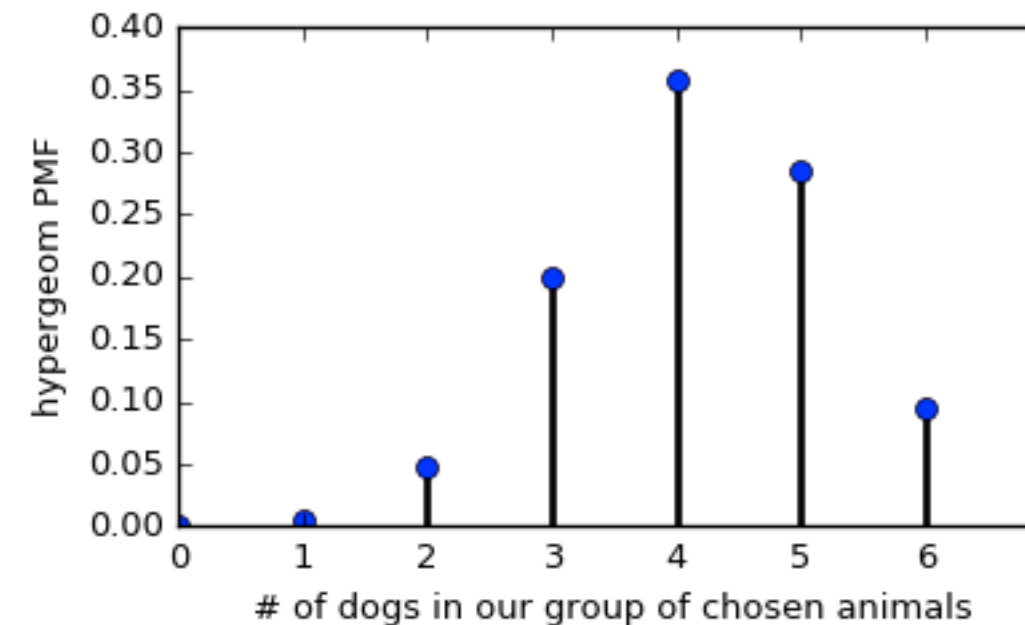
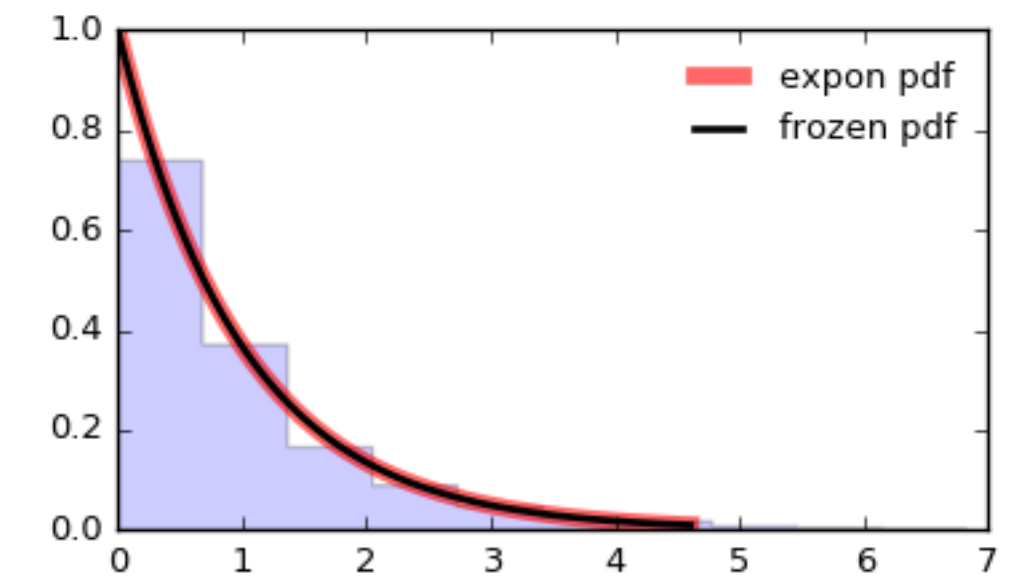
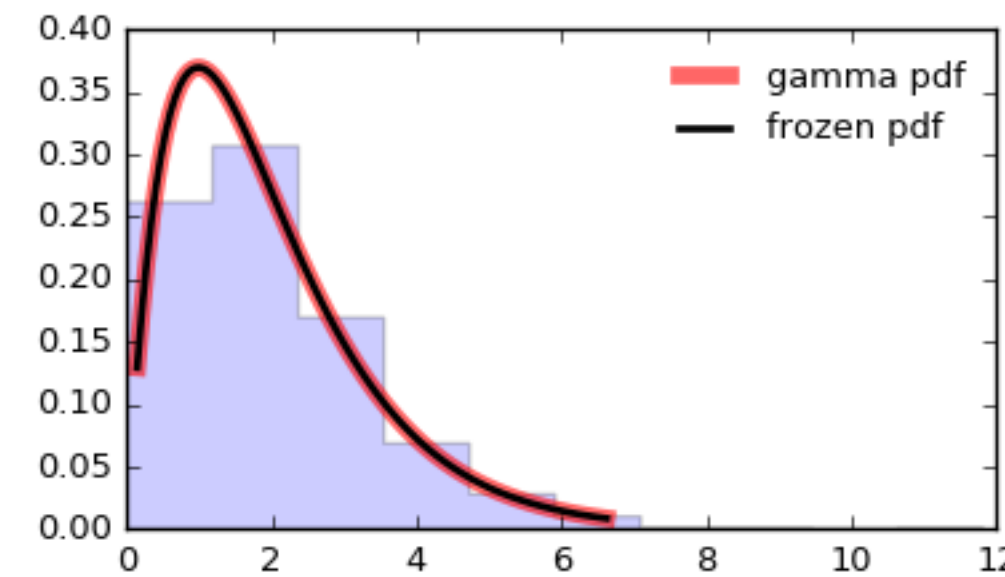
機率分佈

- 連續機率分佈 (Continuous Distributions) :

- 伽瑪分佈 (Gamma Distribution)
- 指數分佈 (Exponential Distribution)
- 常態分佈 (Normal Distribution)
- 均勻分佈 (Uniform Distribution)
- 卡方分佈 (Chi-square Distribution)

- 間斷機率分佈 (Discrete Distributions) :

- 白努力分佈 (Bernoulli Distribution)
- 二項式分佈 (Binomial Distribution)
- 負二項式分佈 (Negative Binomial Distribution)
- 波式分佈 (Poisson Distribution)
- 超幾何分佈 (Hypergeometric Distribution)





Scipy.stats

- Scipy 統計函式：<https://docs.scipy.org/doc/scipy/reference/stats.html>

Continuous distributions

| | |
|---------------------------|---|
| alpha | An alpha continuous random variable. |
| anglit | An anglit continuous random variable. |
| arcsine | An arcsine continuous random variable. |
| beta | A beta continuous random variable. |
| betaprime | A beta prime continuous random variable. |
| bradford | A Bradford continuous random variable. |
| burr | A Burr (Type III) continuous random variable. |
| burr12 | A Burr (Type XII) continuous random variable. |
| cauchy | A Cauchy continuous random variable. |
| chi | A chi continuous random variable. |
| chi2 | A chi-squared continuous random variable. |
| cosine | A cosine continuous random variable. |
| dgamma | A double gamma continuous random variable. |
| dweibull | A double Weibull continuous random variable. |
| erlang | An Erlang continuous random variable. |
| expon | An exponential continuous random variable. |

Discrete distributions

| | |
|---------------------------|---|
| bernoulli | A Bernoulli discrete random variable. |
| binom | A binomial discrete random variable. |
| boltzmann | A Boltzmann (Truncated Discrete Exponential) random variable. |
| dlaplace | A Laplacian discrete random variable. |
| geom | A geometric discrete random variable. |
| hypergeom | A hypergeometric discrete random variable. |
| logser | A Logarithmic (Log-Series, Series) discrete random variable. |
| nbinom | A negative binomial discrete random variable. |
| planck | A Planck discrete exponential random variable. |
| poisson | A Poisson discrete random variable. |
| randint | A uniform discrete random variable. |
| skellam | A Skellam discrete random variable. |
| zipf | A Zipf discrete random variable. |



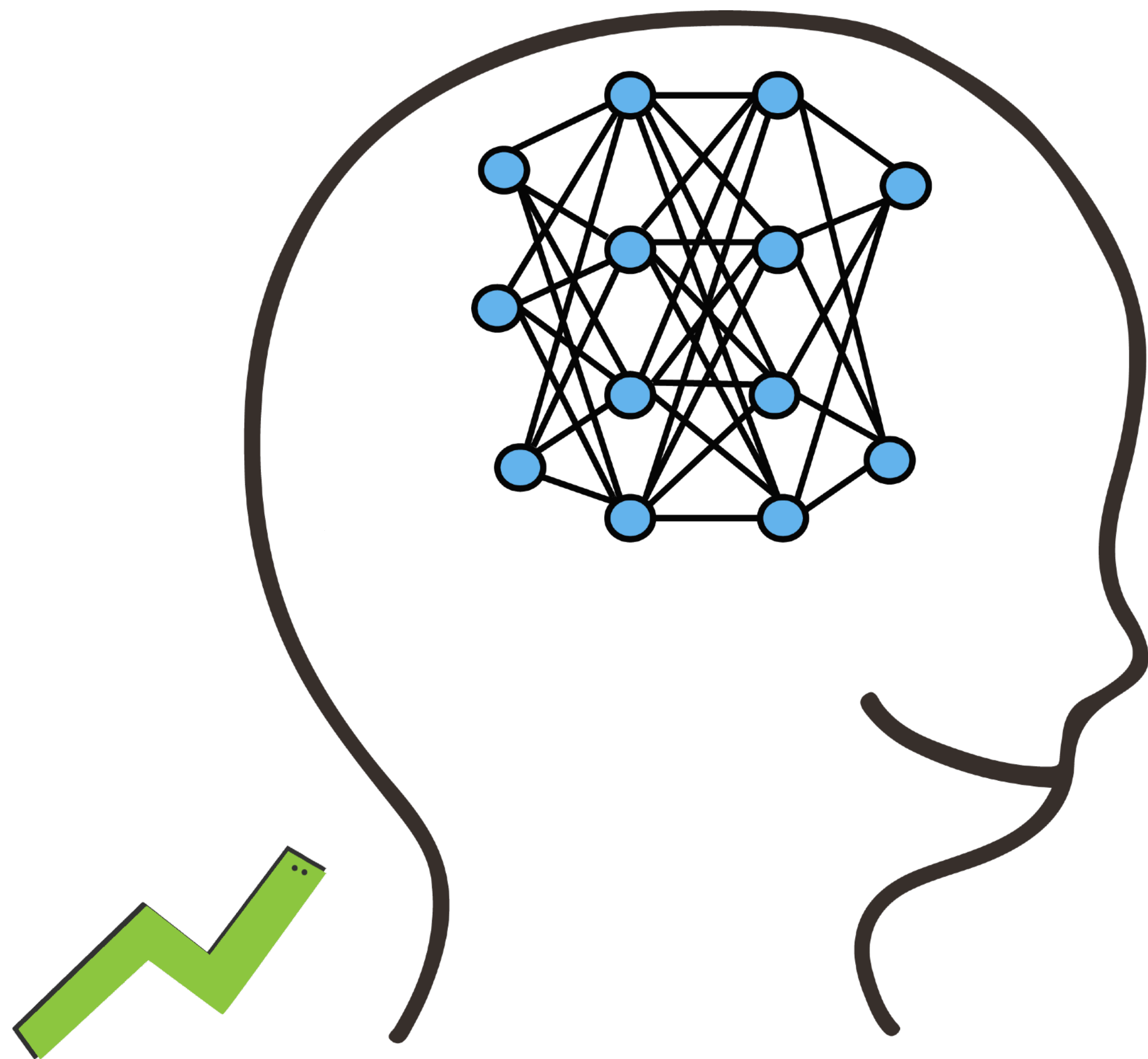
其他常用統計函式

- One-way ANOVA
- F value
- T-test
- 相關性 (Correlation)
- 峰度 (kurtosis)
- 偏態 (skewness)
- 共變數 (Covariance)
- ...



References

- Numpy 統計関式 : <https://docs.scipy.org/doc/numpy/reference/routines.statistics.html>
- Scipy 統計関式 : <https://docs.scipy.org/doc/scipy/reference/stats.html>
- Pandas DataFrame 統計関式 : <http://pandas.pydata.org/pandas-docs/stable/api.html#api-dataframe-stats>

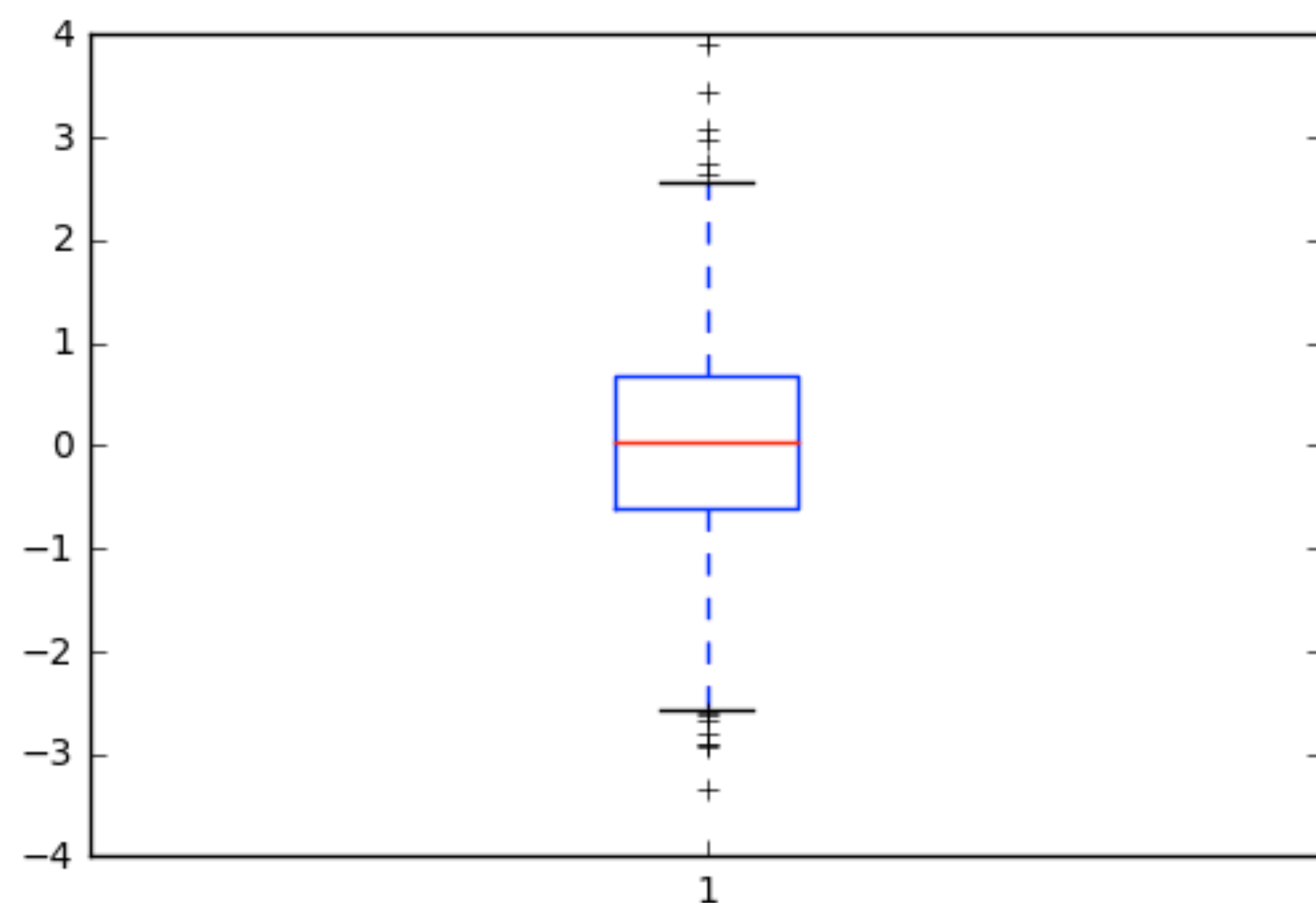


異常値偵測

Anomaly Detection



異常值偵測(1) - 四分位數與箱形圖

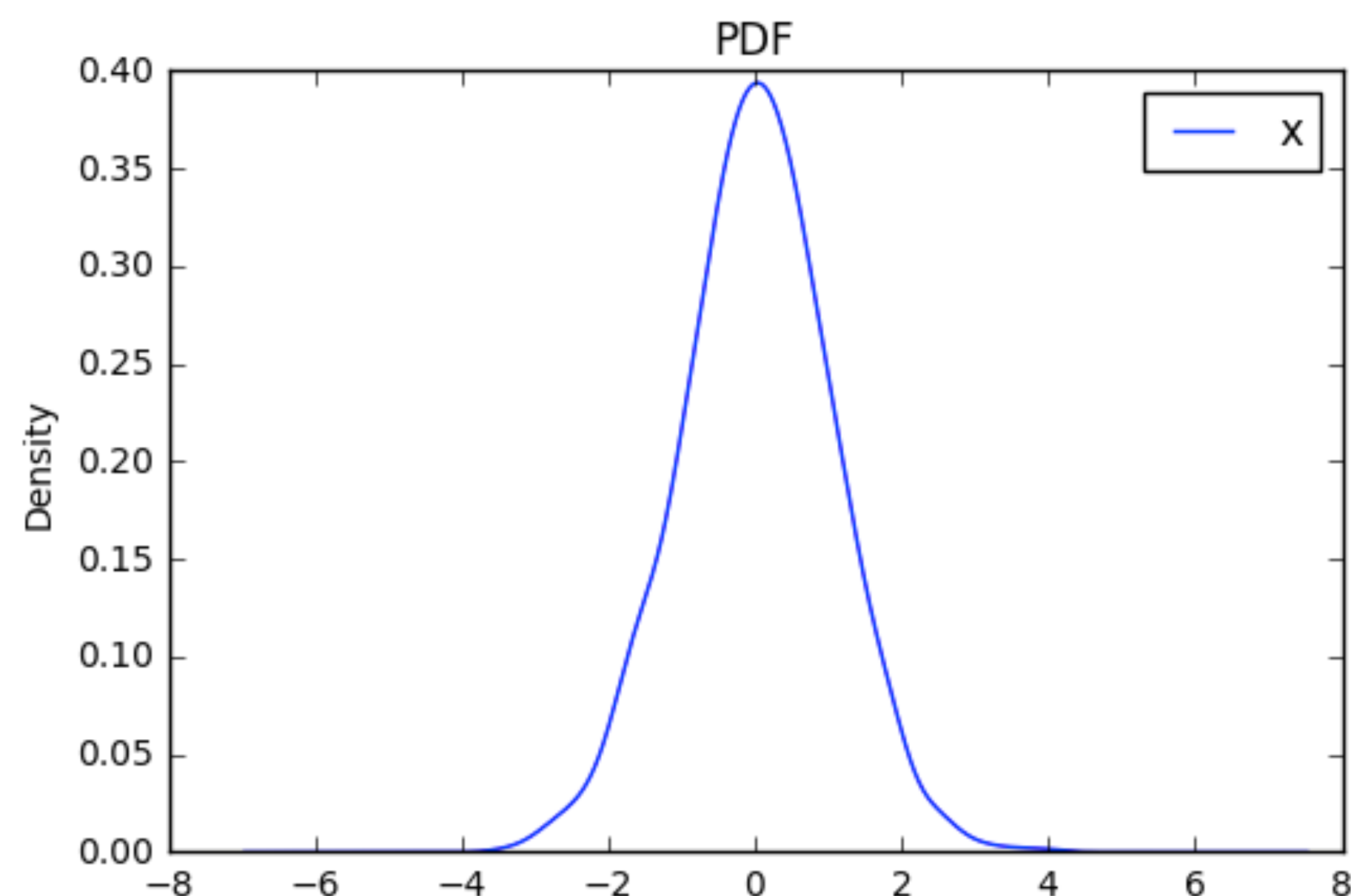


`plt.boxplot(x, showfliers=True)`

- 四分位間距(InterQuartile Range, IQR)
 - $Q_3 - Q_1$
- 最大值： $Q_3 + 1.5 * IQR$
- 最小值： $Q_1 - 1.5 * IQR$
- 異常值：高於最大值、低於最小值



異常值偵測(2) - 常態分佈與標準差



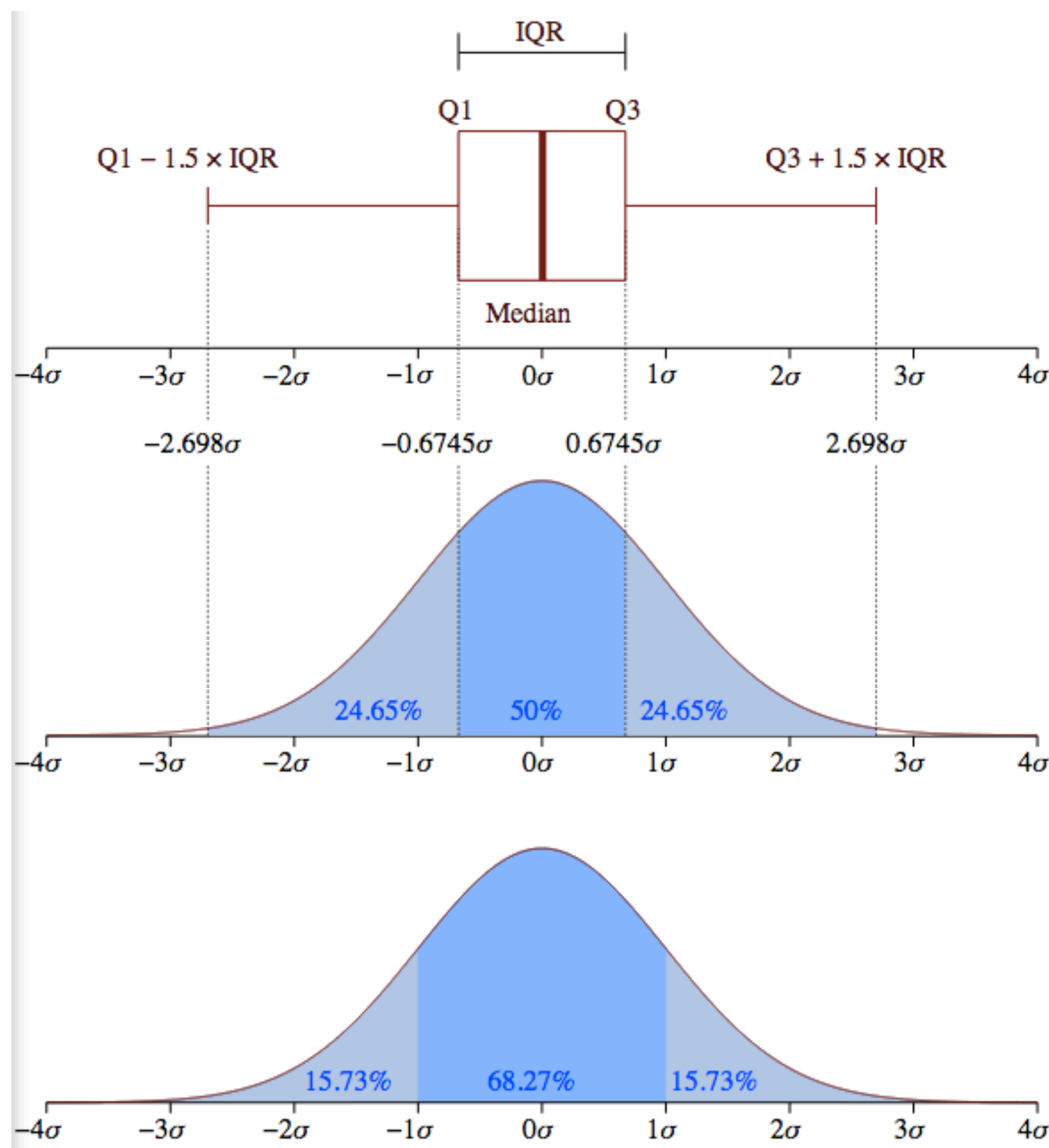
`df.plot(kind='kde')`

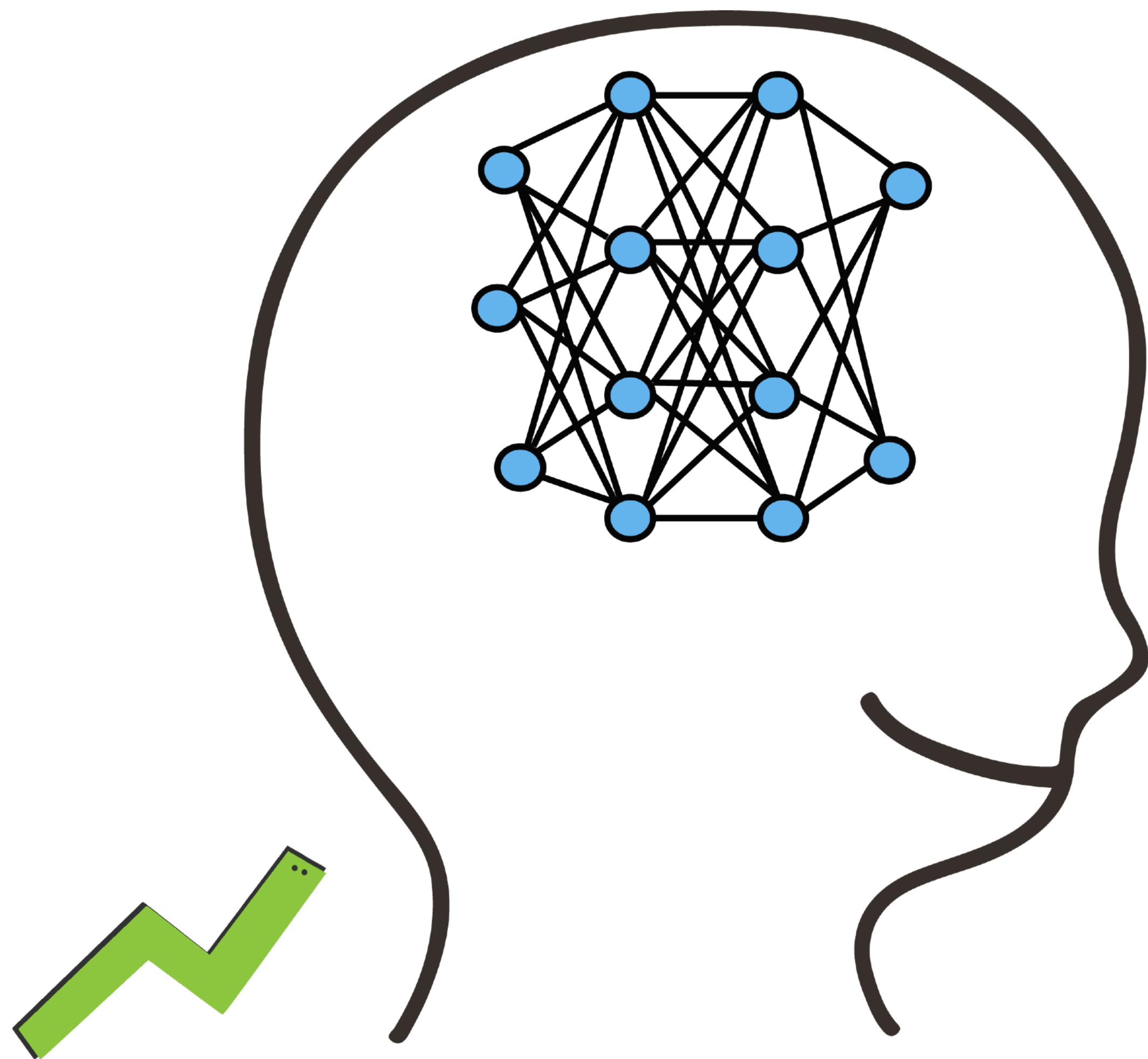
- 平均 (mean) 、標準差 (std, σ)
- 上限： $\text{mean} + 3 * \sigma$ (sigma)
- 下限： $\text{mean} - 3 * \sigma$ (sigma)
- 異常值：高於上限、低於下限



比較

- 若資料型態傾向常態分佈（例如：身高、體重、成績），適合使用標準差的判定方式
- 若異常值過大或過小，容易過度影響標準差，則建議使用四分位數和箱型圖，因為大於 Q_3 和小於 Q_1 的值不論離多遠都不會影響四分位數的值，所以在判定異常值效果好





相關性分析

Correlation Analysis

Pearson 相關係數

- 最常用的相關係數 (-1~+1)

$$\rho_{X,Y} = \frac{\overset{\substack{\text{共變數} \\ \downarrow}}{\text{cov}(X,Y)}}{\underset{\substack{\uparrow \\ \text{標準差}}}{\sigma_X \sigma_Y}} = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y}$$

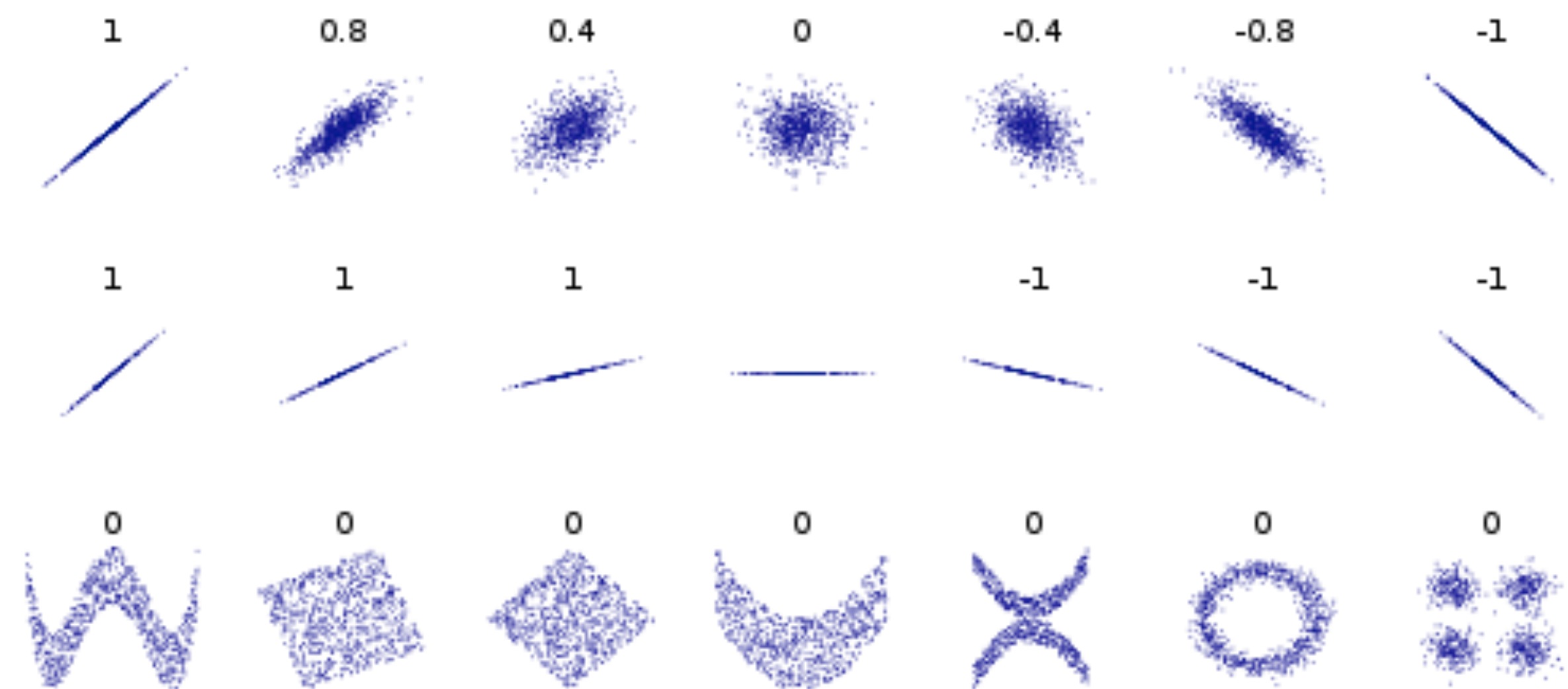
$$r = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

(Pearson, 1917)

相關程度

- 相關程度

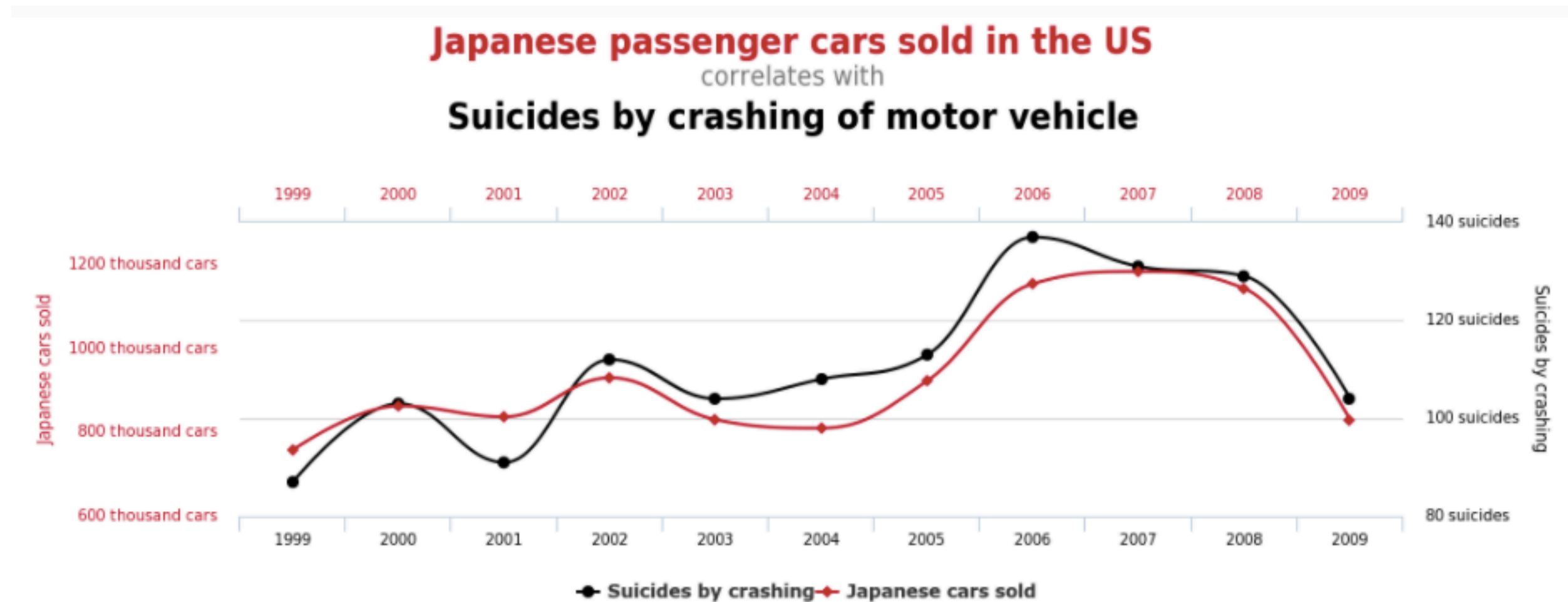
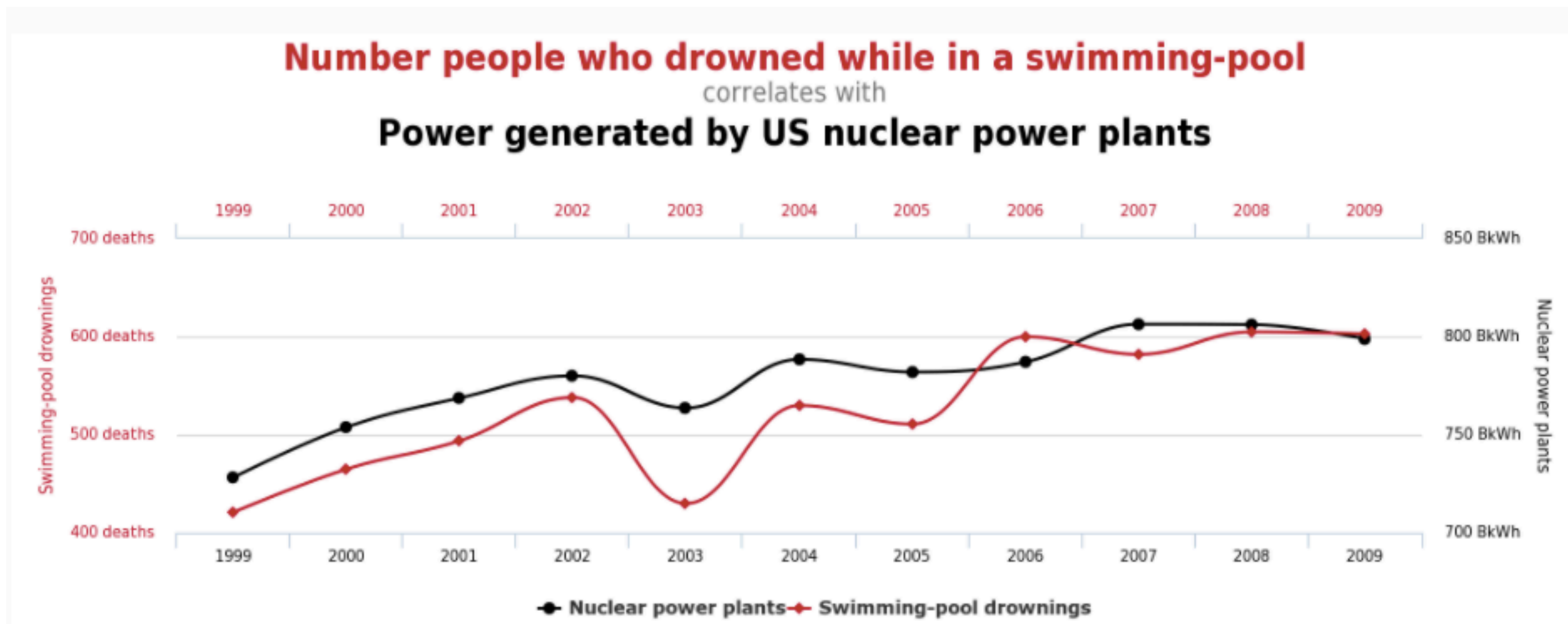
- ▶ $+1$: 完全正相關
- ▶ -1 : 完全負相關
- ▶ 0.3 至 -0.3 : 低度 (正/負)相關
- ▶ $(+/-) 0.3$ 至 0.6 : 中度 (正/負)相關
- ▶ $(+/-) 0.6$ 至 0.9 : 高度 (正/負)相關
 - ▶ `DataFrame.corr()`



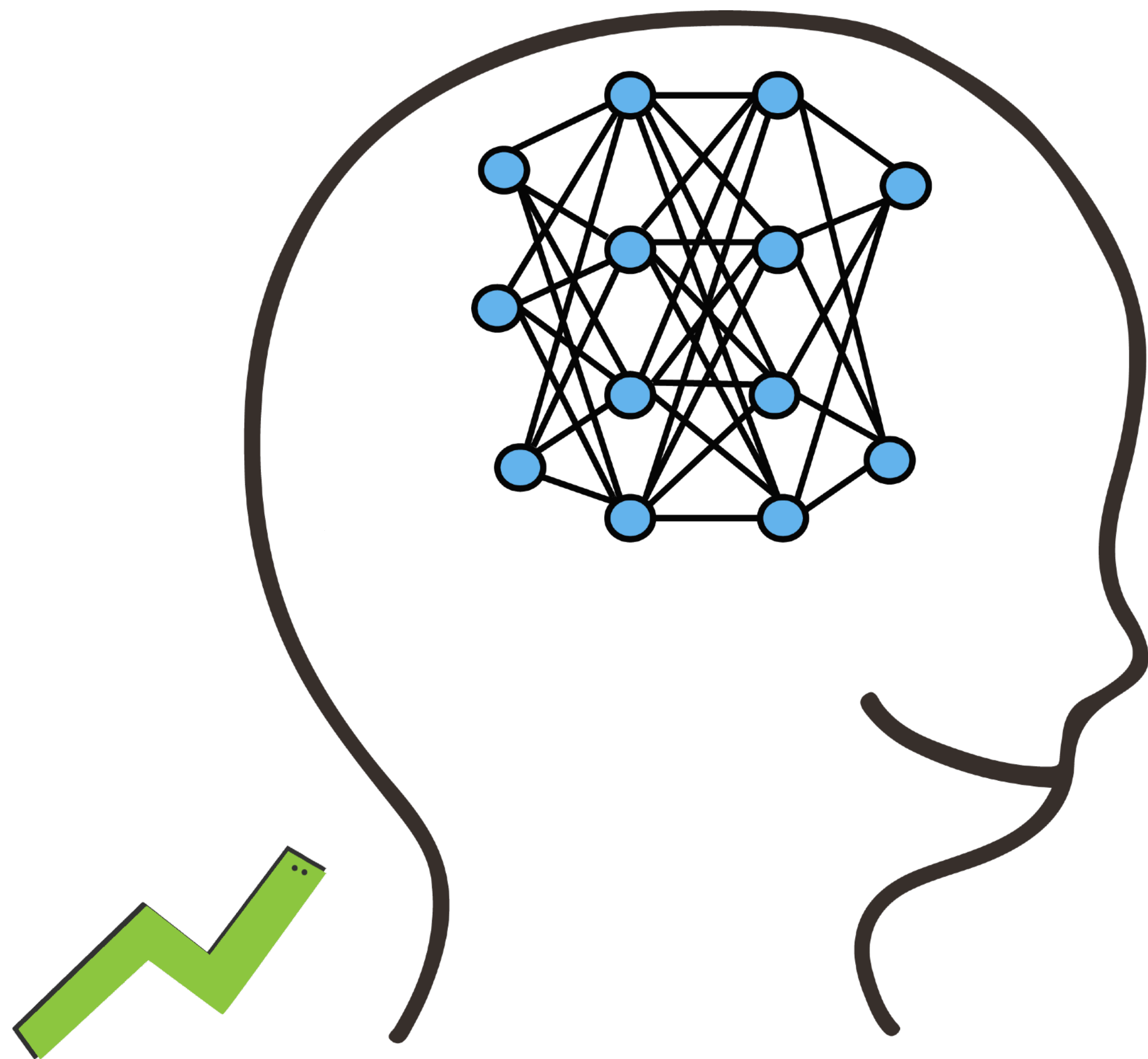
(from wikipedia)



相關程度不等於有因果



(<http://tylervigen.com/>)



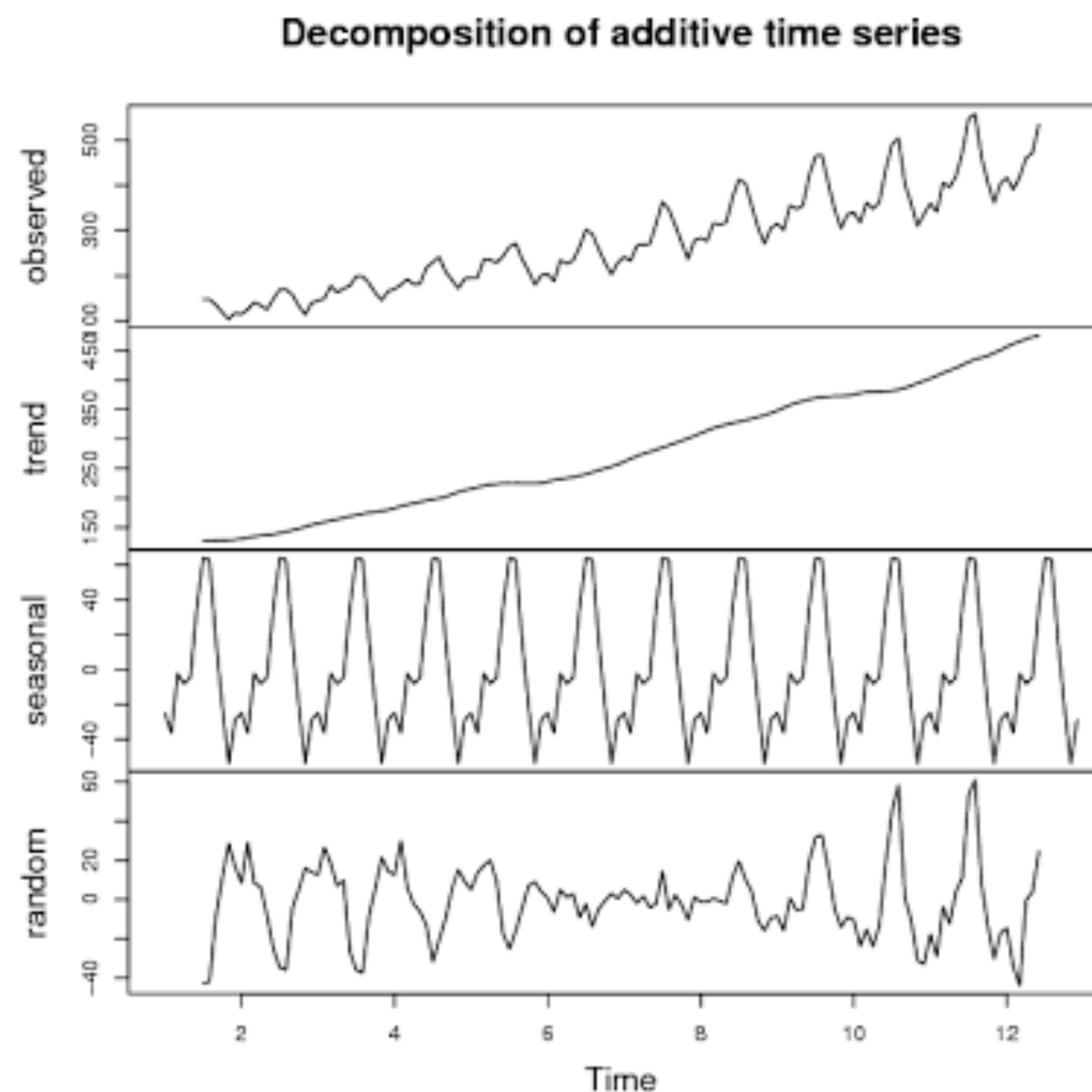
時間序列處理與週期性分析

Time Series Processing & Periodic Analysis



時間序列分析 (Time Series Analysis)

- 主要目的：
 - 分析現象
 - 預測未來
- 時間序列分解：
 - 季節性 (Seasonality)
 - 不規則 (Irregular / Random)
 - 趨勢 (Trend)



Additive Decomposition

Observed series =
Trend + Seasonal + Irregular

$$O_t = T_t + S_t + I_t$$

時間序列分析有更多複雜的模型和
分析方法...

(Australian Bureau of Statistics, 2005; Yanchang Z., 2015)



Python 時間資料型態

- time, calendar
- datetime
 - datetime.date : (year, month, day)
 - datetime.time : (hour, minute, second, microsecond)
 - datetime.datetime: (year, month, day, hour, minute, second, microsecond)
 - datetime.timedelta: (days, seconds, microseconds)
- Documents : <https://docs.python.org/3/library/datetime.html>



String & Datetime 轉換

- string to datetime: `datetime.strptime(str, format)`
- datetime to string: `datetime.strftime(datetime, format)`

| Directive | Meaning | Example |
|-----------|---|--|
| %a | Weekday as locale's abbreviated name. | So, Mo, ..., Sa (de_DE) |
| %A | Weekday as locale's full name. | Sonntag, Montag, ..., Samstag (de_DE) |
| %w | Weekday as a decimal number, where 0 is Sunday and 6 is Saturday. | 0, 1, ..., 6 |
| %d | Day of the month as a zero-padded decimal number. | 01, 02, ..., 31 |
| %b | Month as locale's abbreviated name. | Jan, Feb, ..., Dez (de_DE) |
| %B | Month as locale's full name. | Januar, Februar, ..., Dezember (de_DE) |
| %m | Month as a zero-padded decimal number. | 01, 02, ..., 12 |
| %Y | Year with century as a decimal number. | 0001, 0002, ..., 2013, 2014, ..., 9999 |
| %H | Hour (24-hour clock) as a zero-padded decimal number. | 00, 01, ..., 23 |
| %I | Hour (12-hour clock) as a zero-padded decimal number. | 01, 02, ..., 12 |
| %p | Locale's equivalent of either AM or PM. | am, pm (de_DE) |
| %M | Minute as a zero-padded decimal number. | 00, 01, ..., 59 |
| %S | Second as a zero-padded decimal number. | 00, 01, ..., 59 |
| %f | Microsecond as a decimal number, zero-padded on the left. | 000000, 000001, ..., 999999 |

- 節錄自：<https://docs.python.org/3/library/datetime.html>



Pandas DatetimeIndex

- 把DataFrame的index轉為DatetimeIndex型態
 - e.g `df.index = pd.to_datetime(df.index,format='%Y-%m-%d')`
- DatetimeIndex 時間分割/聚合
 - e.g. `DataFrame.groupby([columns]).agg_func()`
 - e.g. `DataFrame.resample('M').agg_func()`
 - e.g. `DataFrame.resample('Q-NOV').agg_func()` #Q-EndMonth
- Document (DatetimeIndex的attributes、methods) : <http://pandas.pydata.org/pandas-docs/stable/generated/pandas.DatetimeIndex.html>

| | | | | | | | | | | | | |
|-------|----|---|----|----|---|----|----|---|----|----|----|----|
| M | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
| Q-DEC | Q1 | | | Q2 | | | Q3 | | | Q4 | | |
| Q-NOV | Q1 | | Q2 | | | Q3 | | | Q4 | | Q1 | |