

Lecture 10:

Review & Advanced Regression Techniques

Essential to pass the course

Bonus to get higher scores

Road Map of the Statistics Part

	Lecture 5	Lecture 6	Lecture 7	Lecture 8	Lecture 9
Quantification Technique	Mean, variance, skewness, & kurtosis	Pearson's Correlation (Linear relationship)	Linear regression (OLS)	Model Selection	TLS / PCA / EOF
Uncertainty & Significance	Gaussian distribution Chi-2 distribution	<code>r, p = scipy.stats.pearsonr(x, y)</code>	<code>results.summary()</code>	Training error vs. prediction error	
Assumptions	Data is Gaussian or follows specific types of distribution Independent Sampling	Data is Gaussian Independent Sampling	x is noise free Error is Gaussian Independent Sampling Equal err variance		
Test assumptions	K-S test		Auto-correlation (Effective Sample Size)	Split datasets Cross Validation	Regression Dilution Total Least Square
Treatment		Bootstrapping	Block Bootstrapping	BIC	

Explore a single dataset - a crucial building block for later analysis (Problem 1)

Lecture/Practical 5

Remove mean /
seasonal cycle...

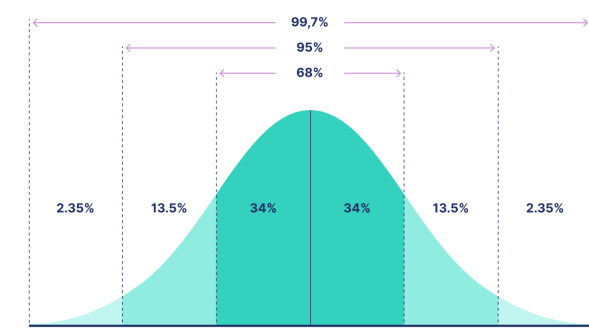
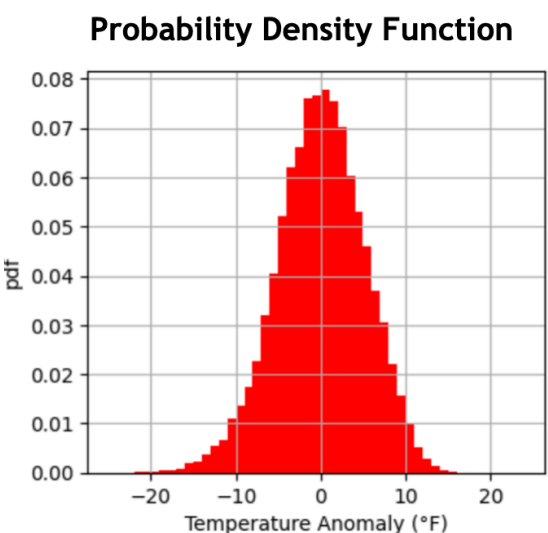
Plot and visualise
your data

Calculate Summary
Statistics

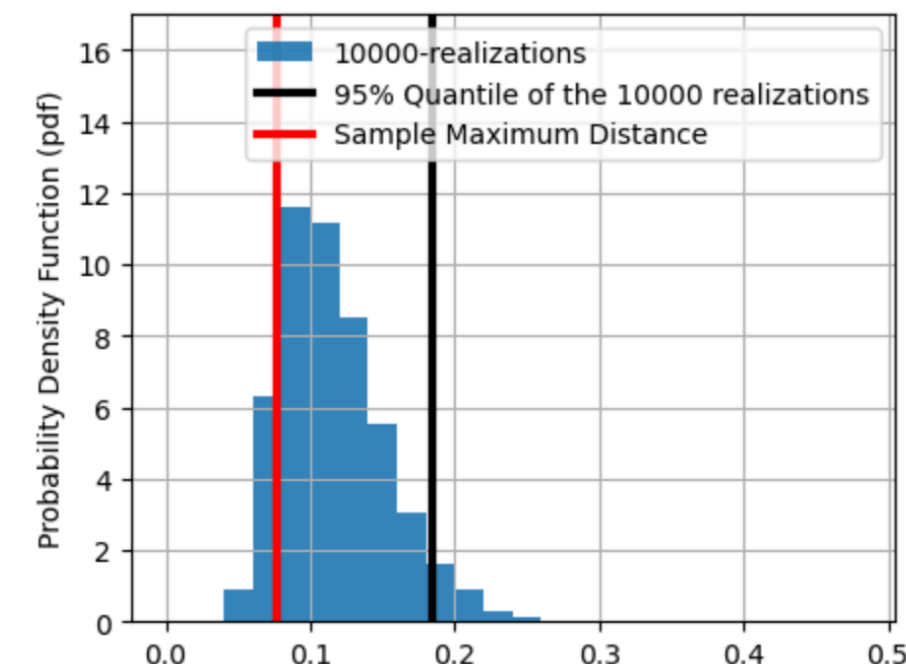
Is data Gaussian or
Normally distributed?

K-S test

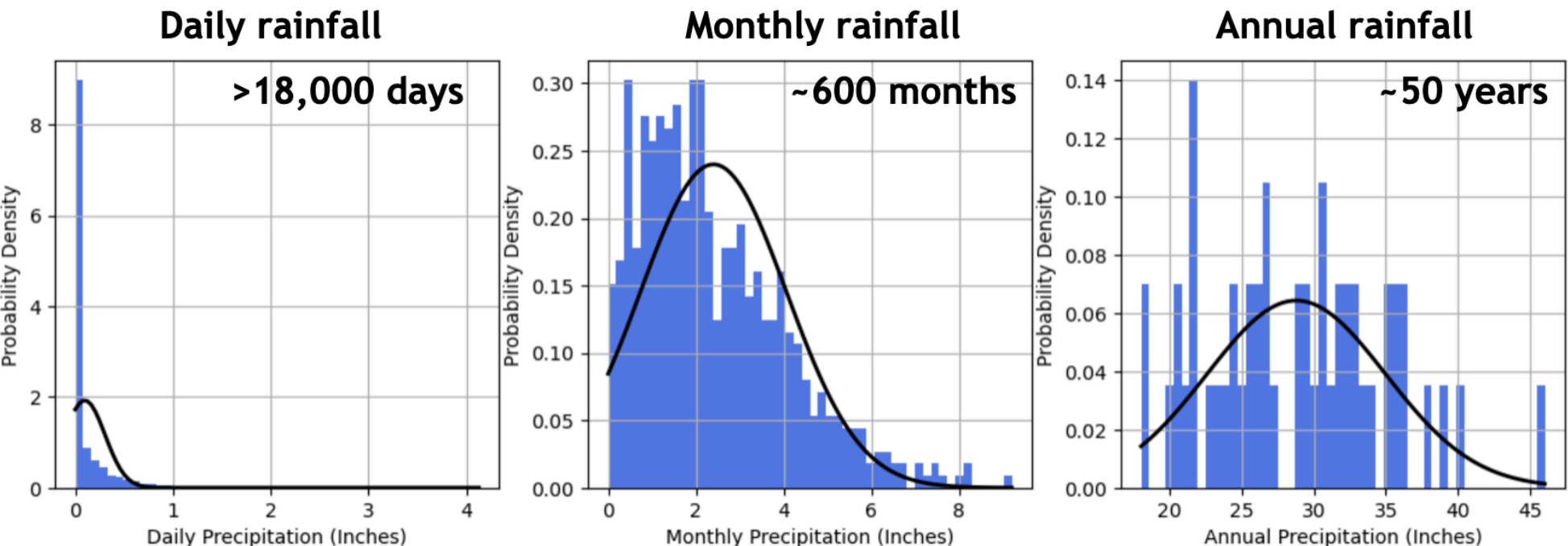
p-value?



Skewness of 0
Kurtosis of 3



Central Limit Theorem:
Likely Gaussian when average or sum over **>30 independent** samples (such as monthly or annual mean values).



Terminology	Meaning
Null hypothesis (H_0)	This is our starting assumption that the effect being studied does not exist
Alternative hypothesis (H_1)	This is what we might believe to be true if we find sufficient evidence against the null hypothesis.
Test statistics	This is a calculated value from our data that we use to test our hypothesis.
Null distribution	This represents what we would expect to see from our test statistic purely by chance if the null hypothesis were true.
significance level (α)	This is a threshold we set to decide when to reject the null hypothesis.
p-value	The probability of obtaining our data, or something more extreme, if the null hypothesis is true.

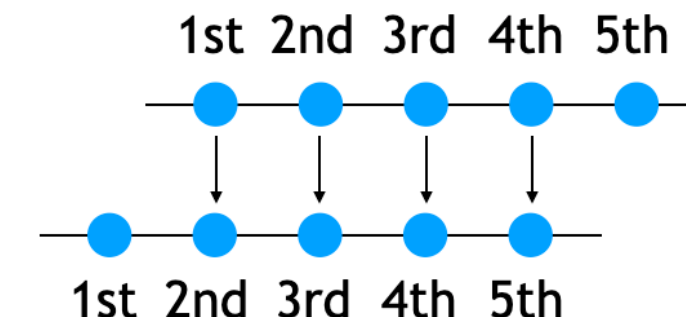
Lecture/Practical 7

Are data independent?

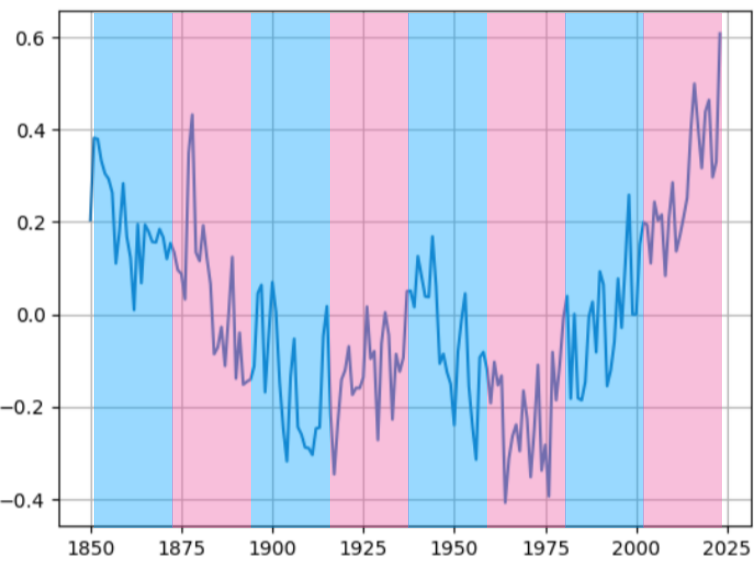
Auto-correlation test

No

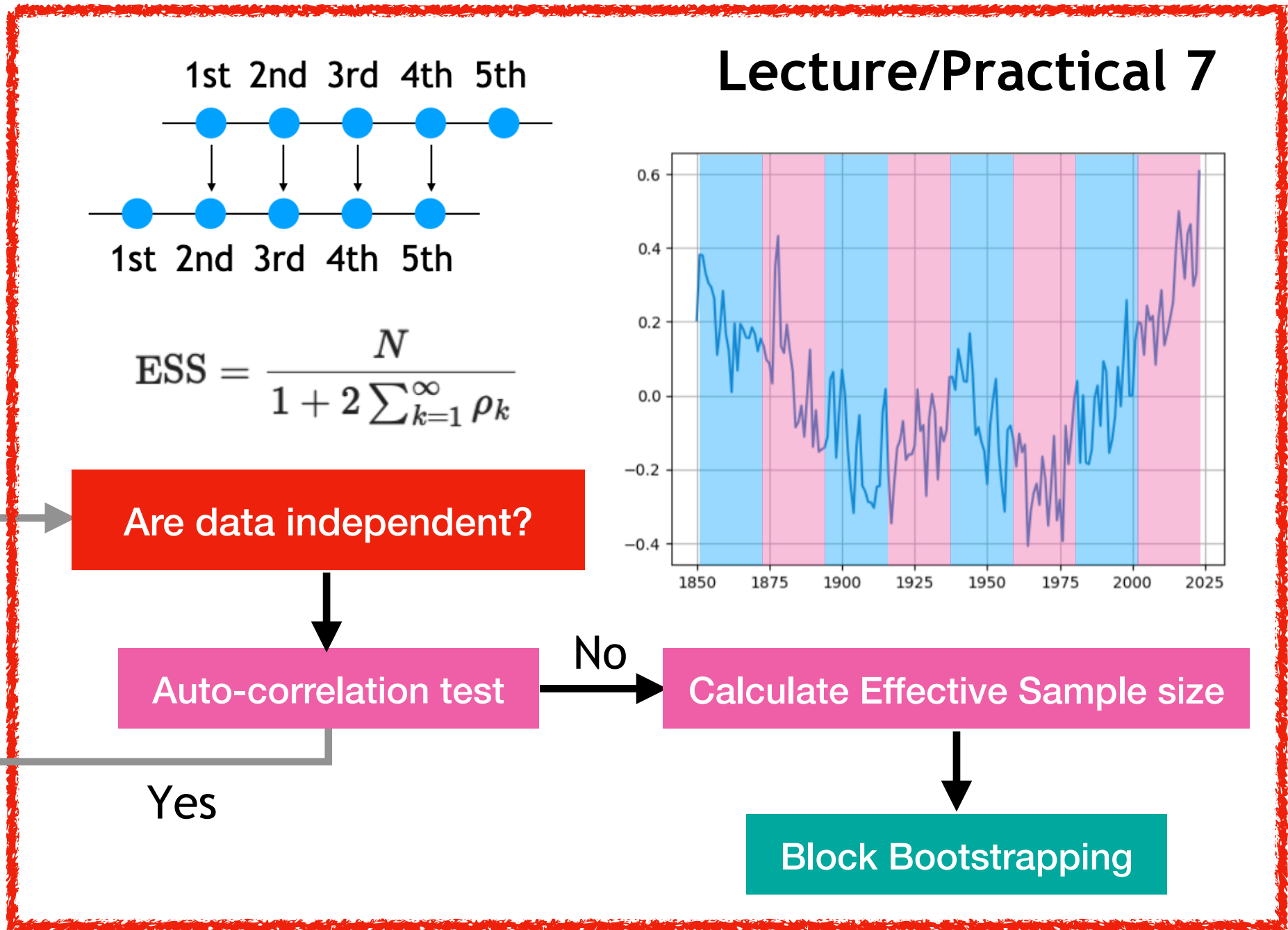
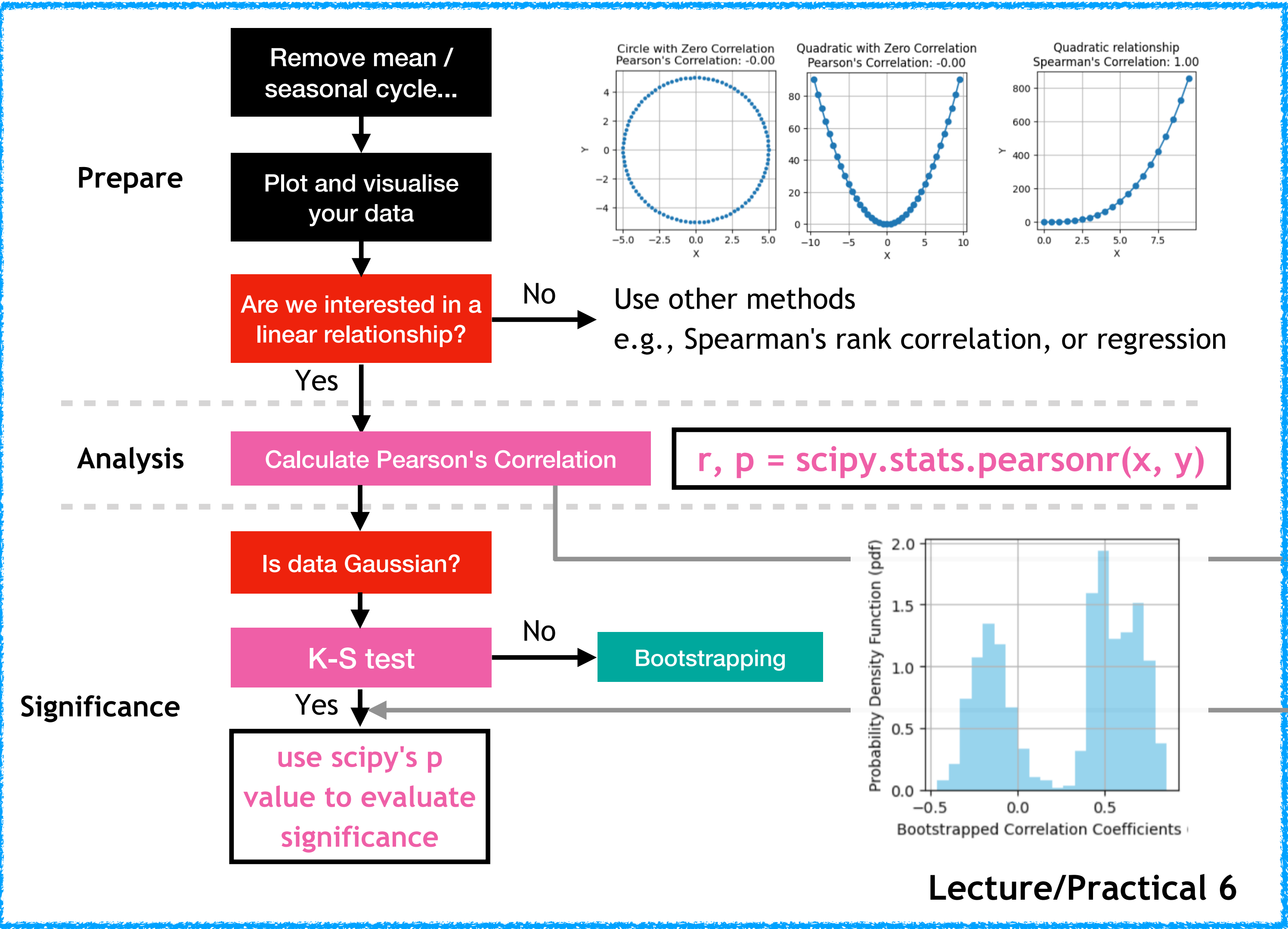
Calculate Effective
Sample size



$$ESS = \frac{N}{1 + 2 \sum_{k=1}^{\infty} \rho_k}$$



Steps for evaluating correlations between two variables (Problem 2)

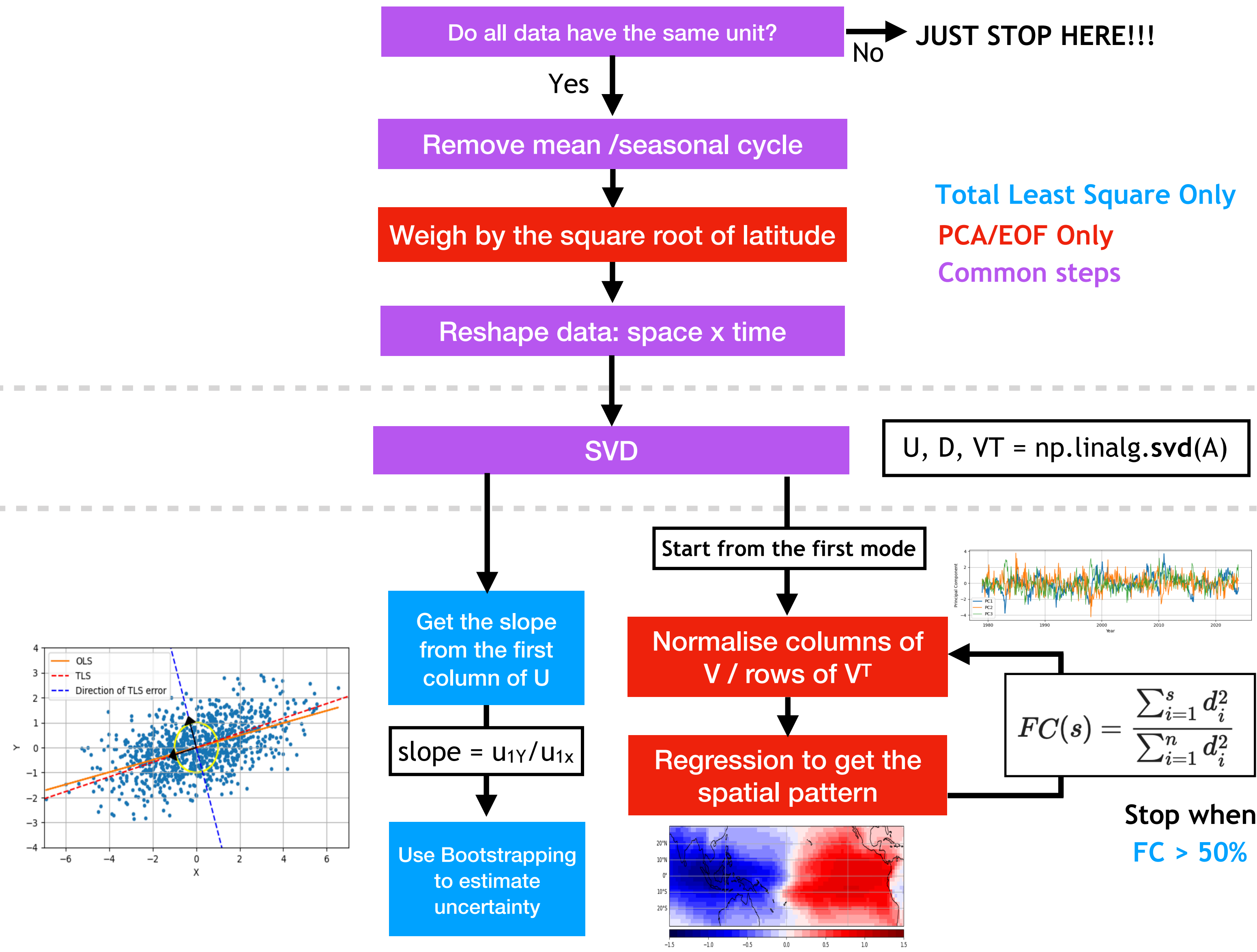


Steps for evaluating correlations between two variables (Problem 3)

Prepare

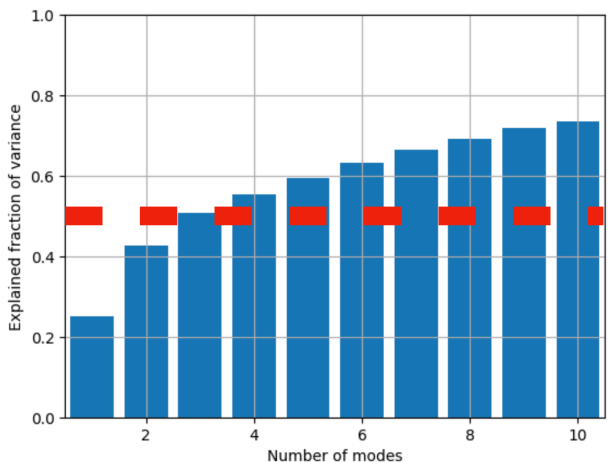
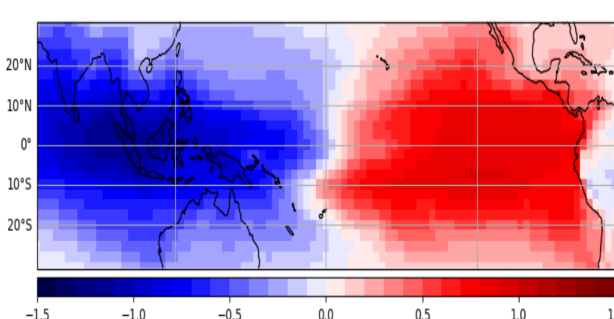
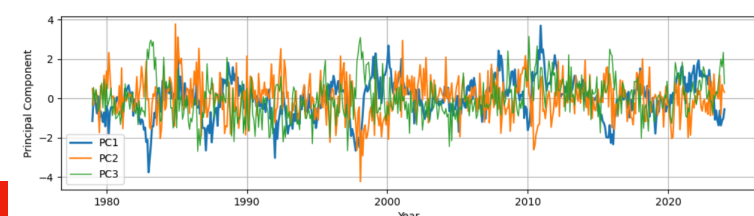
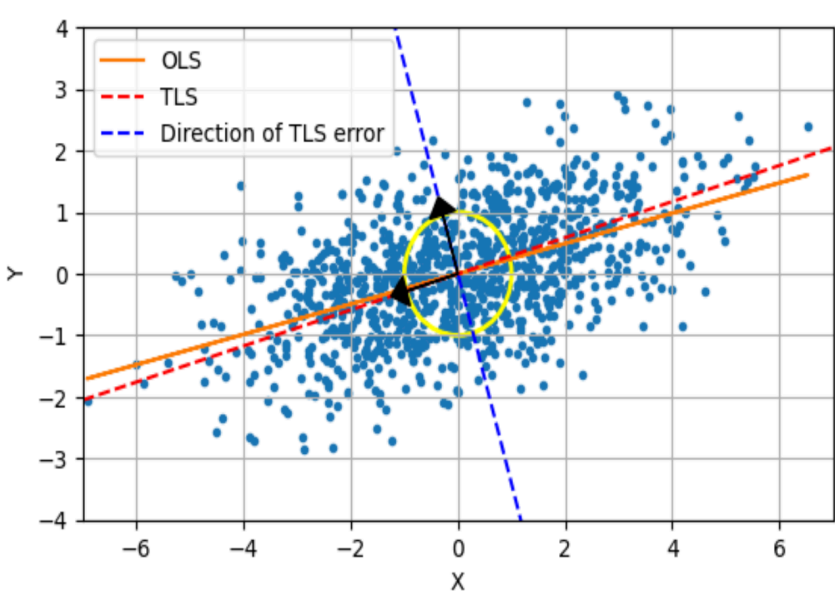
Analysis

Post-processs

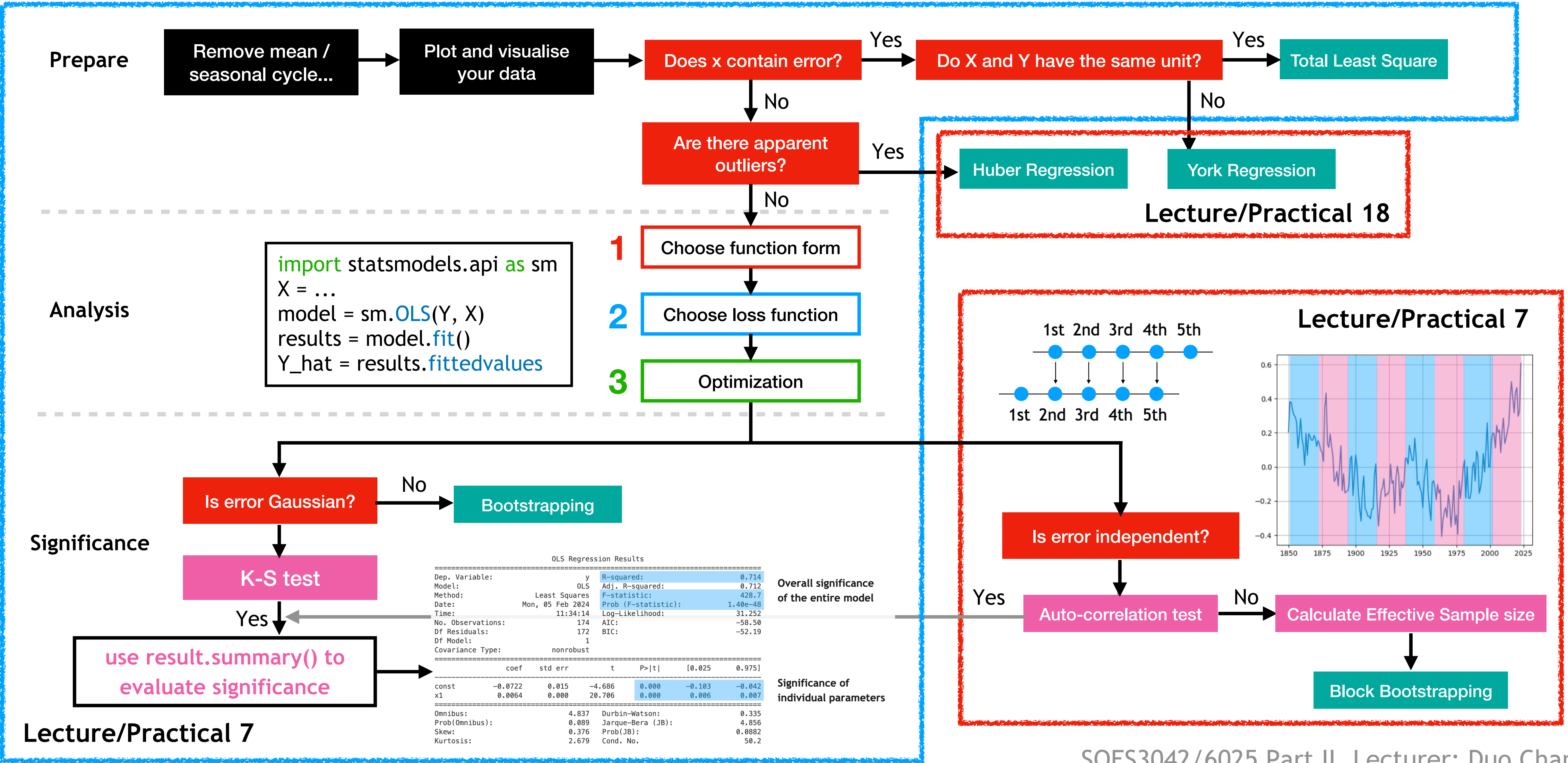


Total Least Square Only
PCA/EOF Only
Common steps

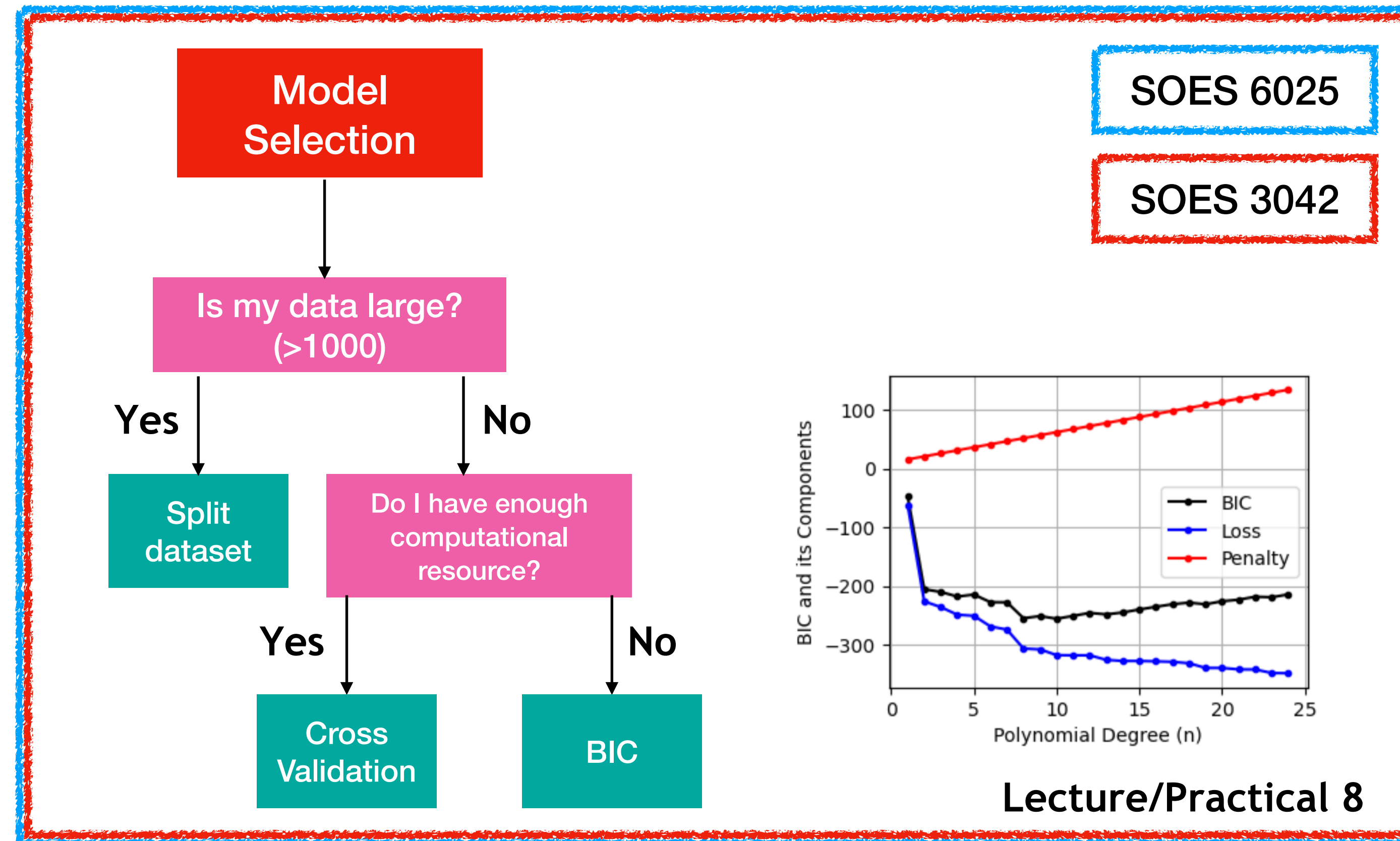
- (1) Individual columns of **U** are **orthogonal**.
The new directions are perpendicular to each other.
- (2) Individual columns of **V** are **orthogonal**.
Pearson's correlations of locations in the new coordinate is zero.
- (3) D_i is ranked in a **descending order**.



Steps for linear regression using ordinary least square (Problem 4)



Steps for Model Selection (Problem 4)



Accounting for Regression Dilution: York Regression

When the error of individual data points is **known**

York Regression is a way to account for regression dilution.

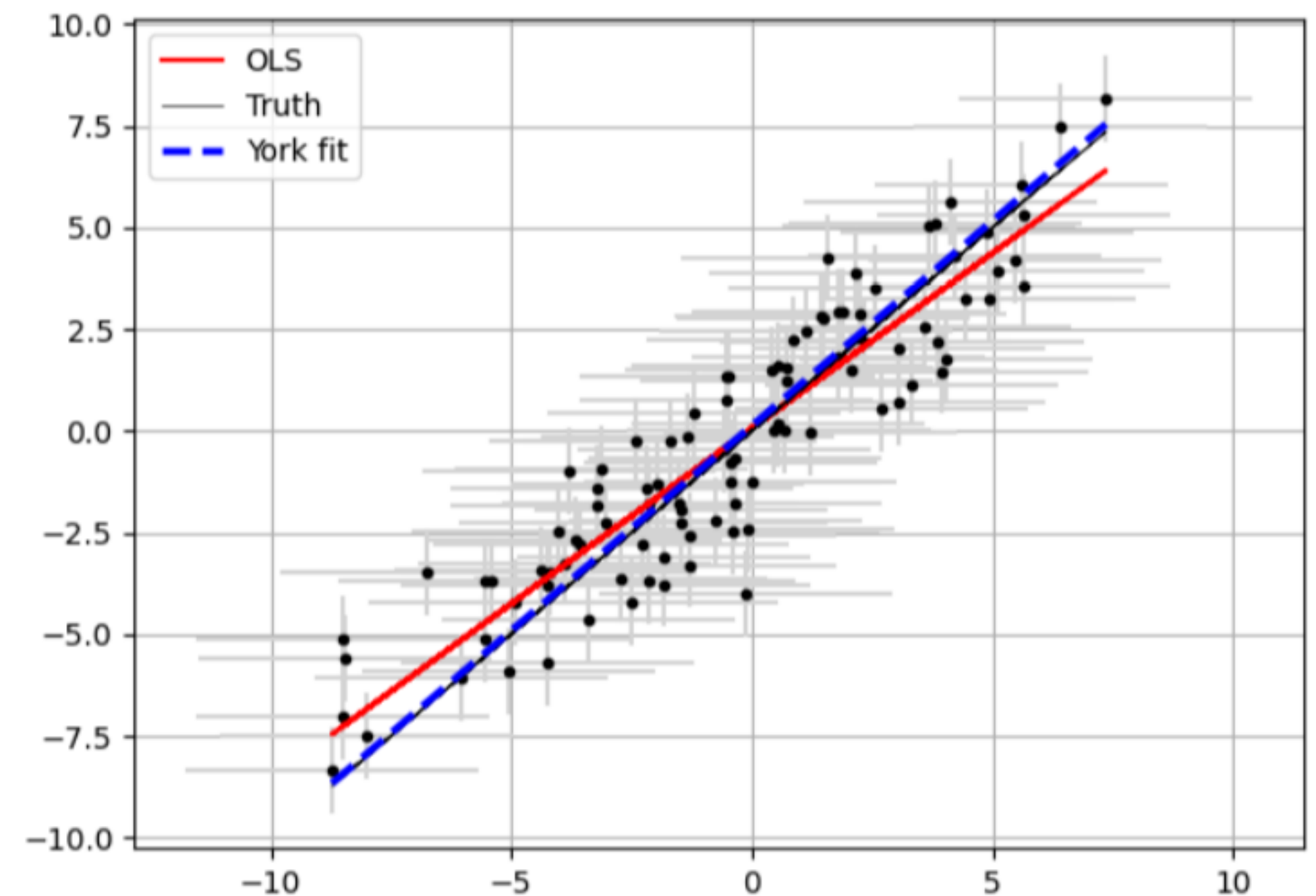
$$P(x|\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{(x - \mu)^2}{2\sigma^2}\right]$$

$$\sigma^2 = \sigma_y^2 + \alpha^2 \sigma_x^2 - 2\alpha\rho\sigma_x\sigma_y$$

Weight is a function
of fitted slope

Input: x , y , σ_x , σ_y , ρ (optional)

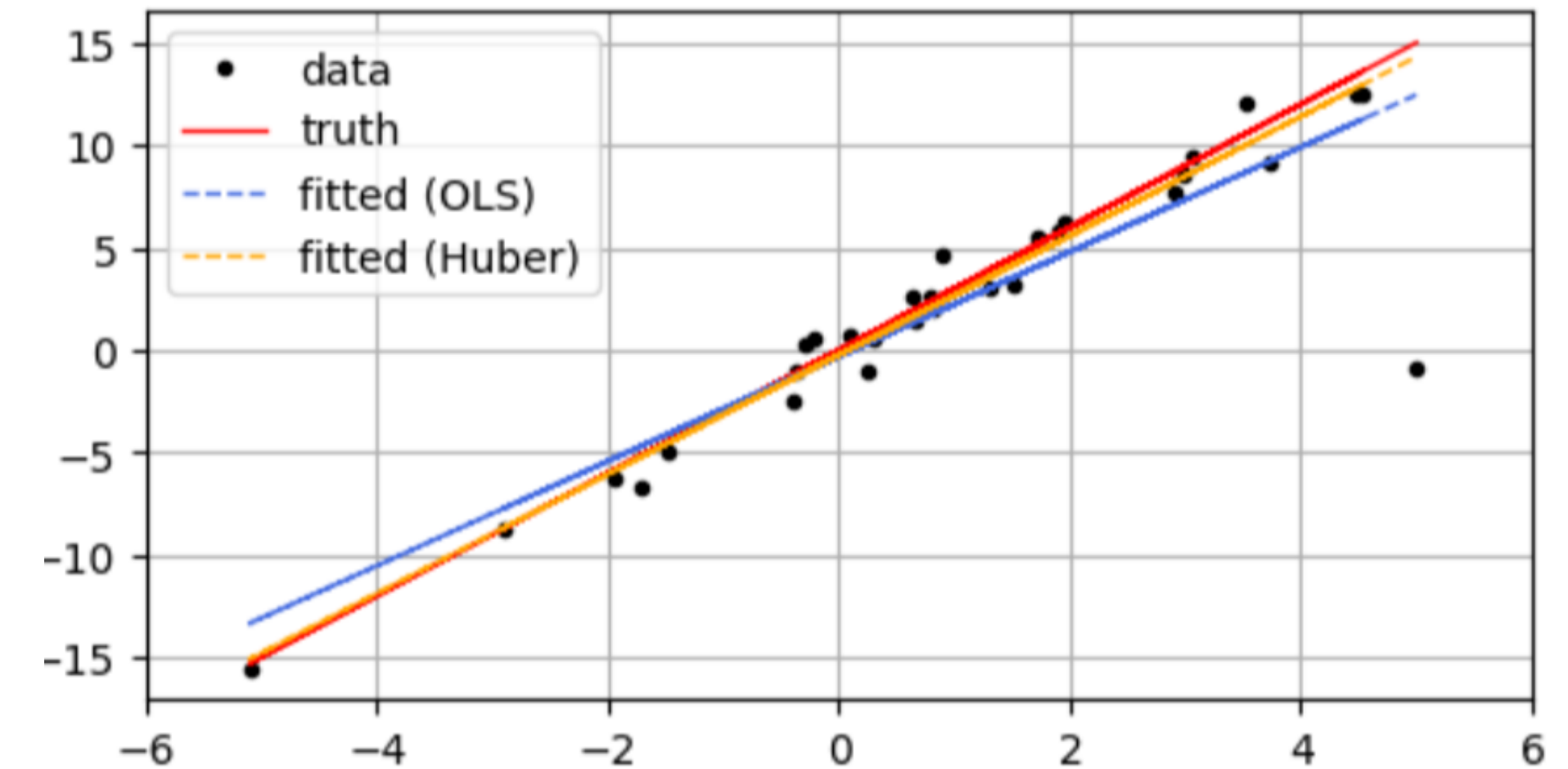
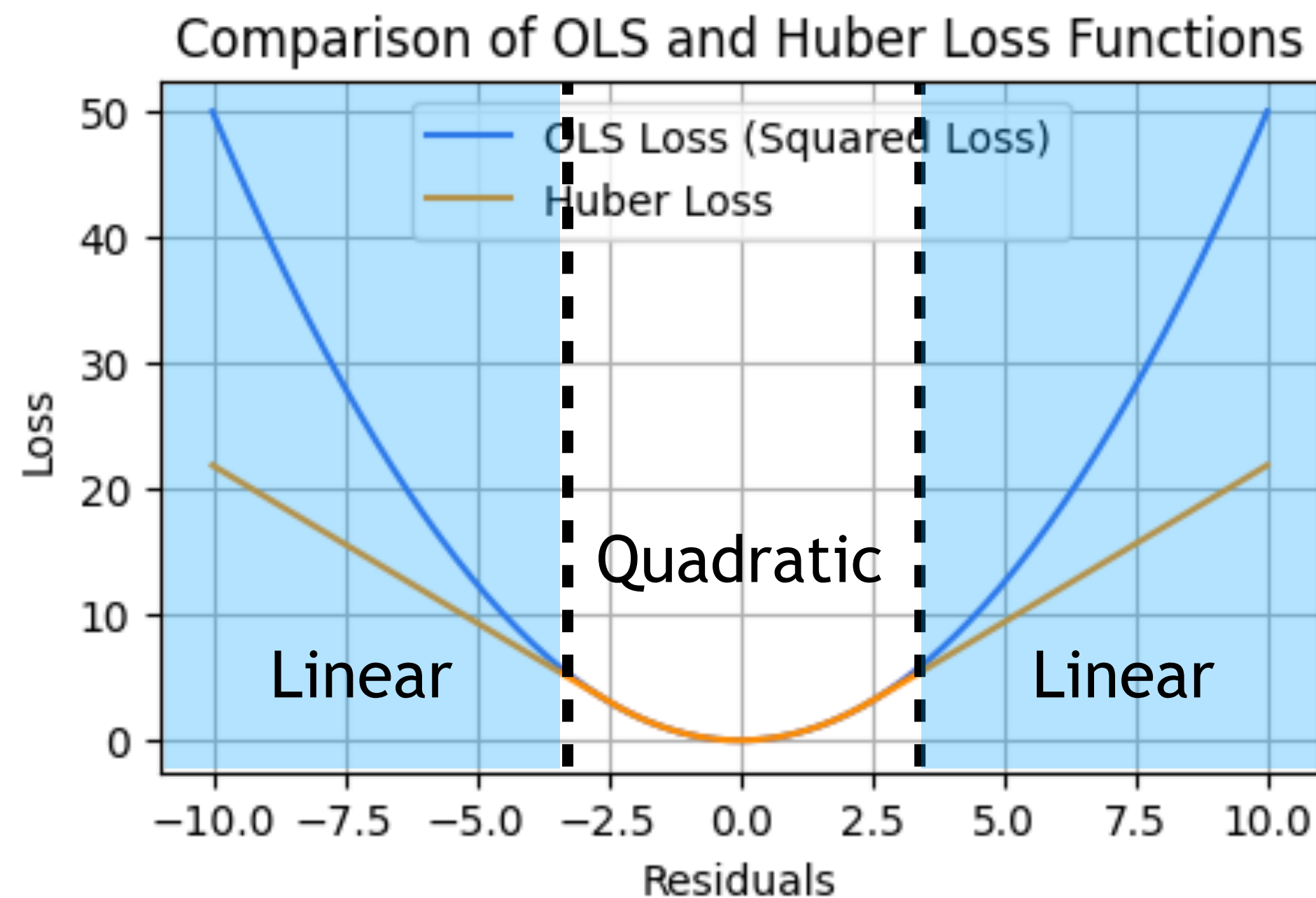
The algorithm uses an iterative approach that converges to the solution.



Account for Outliers: Huber Regression

When the error of individual data points is **unknown**

We need to estimate which data points are more reliable



```
import statsmodels.api as sm  
  
X = sm.add_constant(years)  
  
model = sm.RLM(y, X,  
               M=sm.robust.norms.HuberT())  
  
results = model.fit()  
  
GMST_hat = results.fittedvalues
```

Generalised Regression and Machine Learning

1

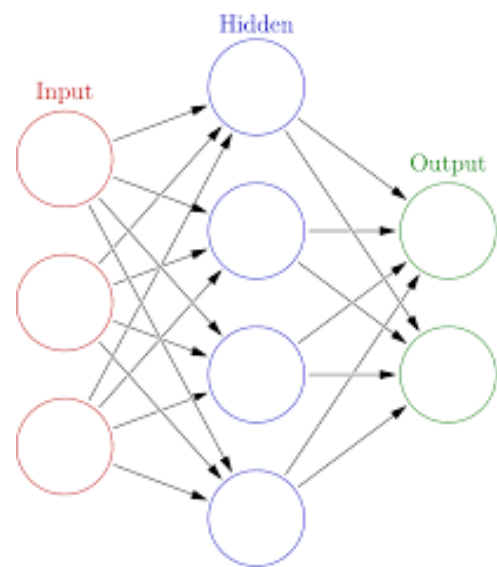
Choose a functional form

2

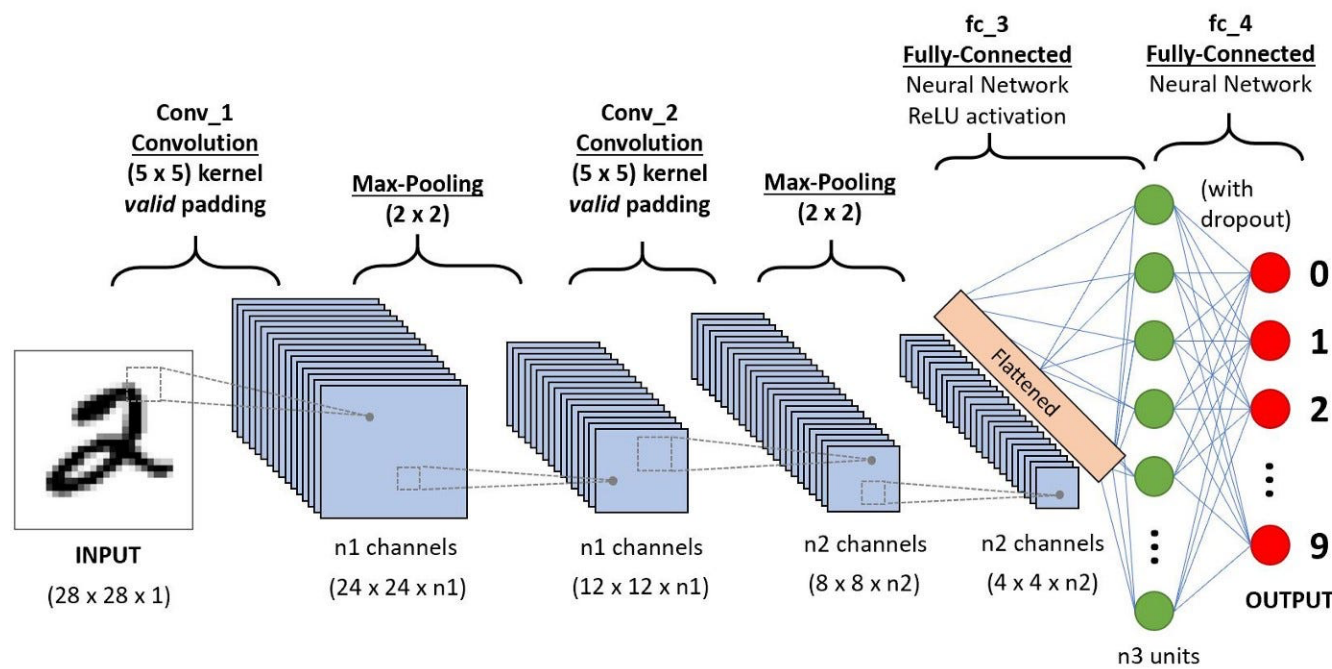
Define what it means by fit

3

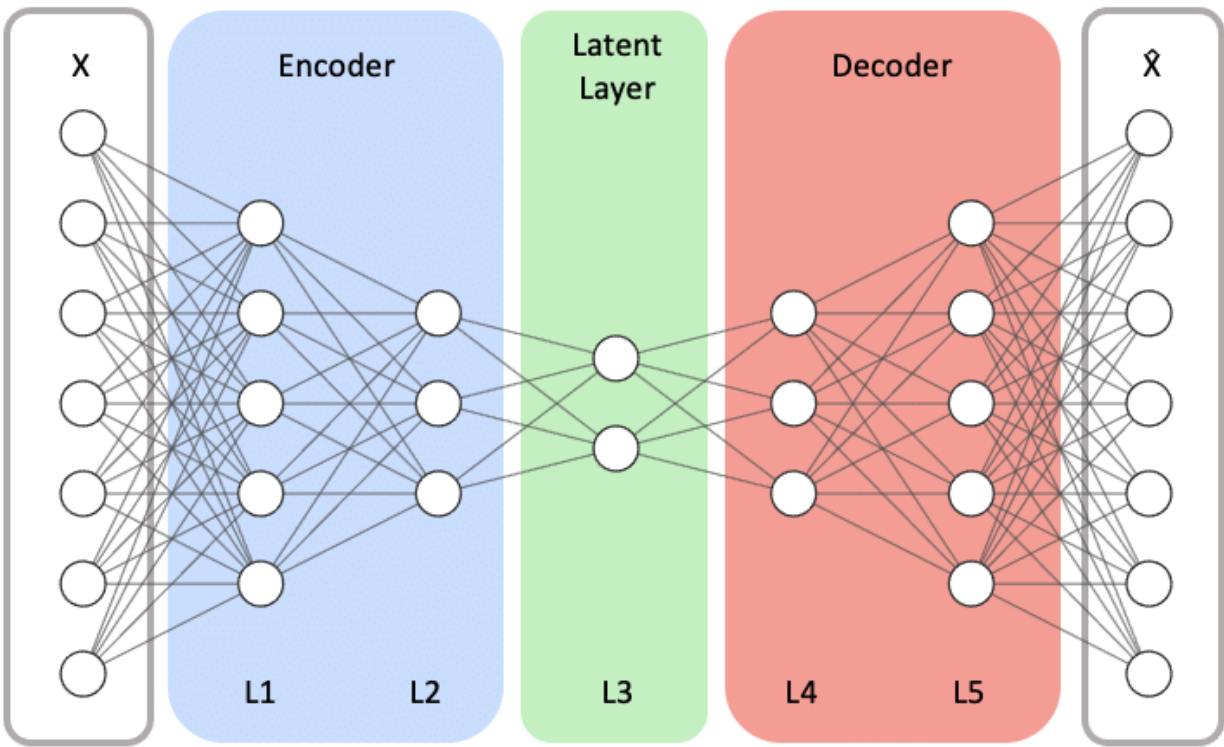
Find that "most"



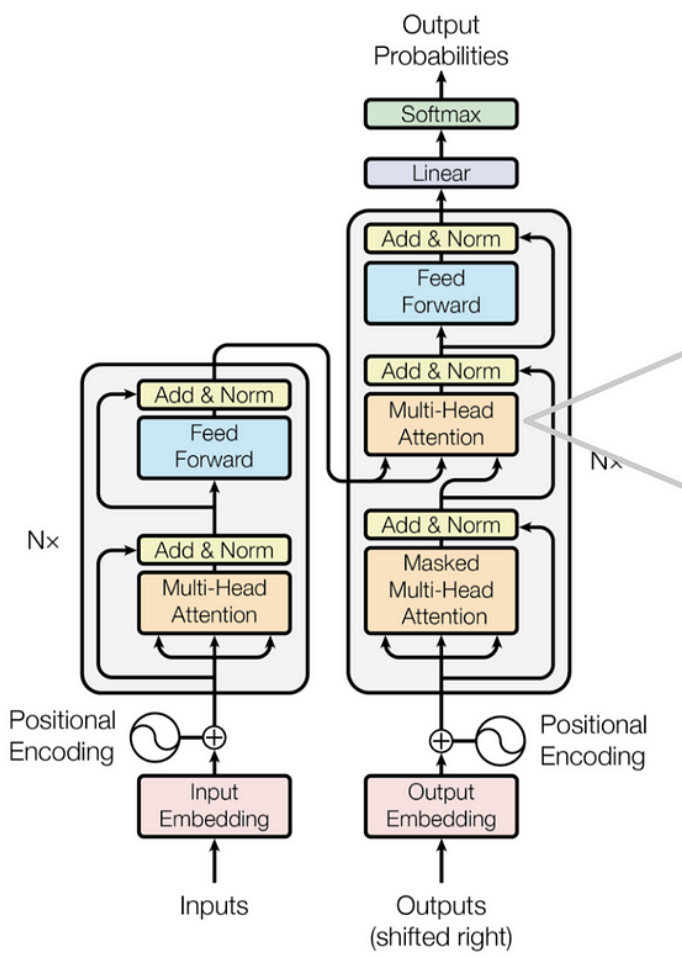
Neural Network



Convolutional Neural Network (Image recognition; weather forecast)

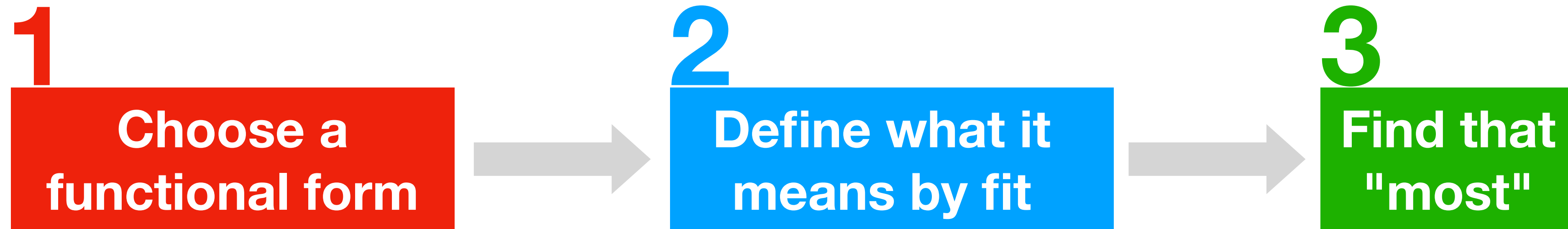


Auto encoder - decoder (Dimension reduction; non-linear EOF)



Transformer (GPT)

Generalised Regression and Machine Learning



Other loss functions

Absolute Loss (L1 Loss) - LASSO regression

Kullback-Leibler (KL) Divergence - For machine learning categorical data

Poisson Loss - For count data.

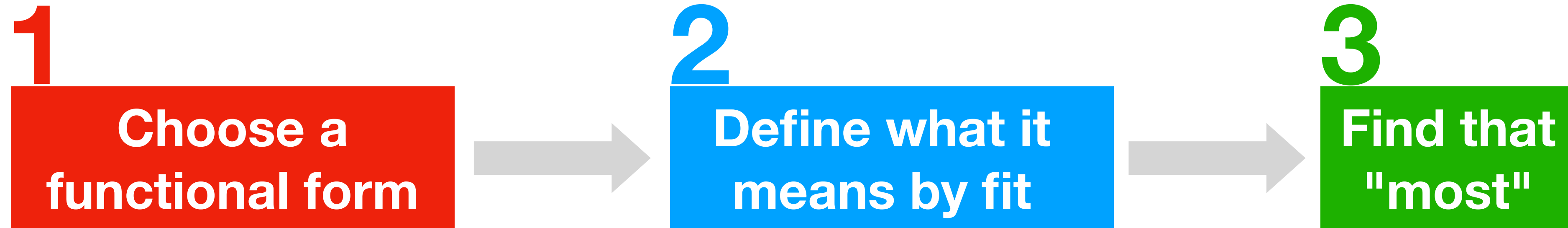
Binomial Loss - For binary outcomes or proportions.

Multinomial Loss - For multi-class categorical data.

Exponential Loss - For time-to-event or survival data.

Gamma Loss - For positively skewed continuous data.

Generalised Regression and Machine Learning



Optimisation

```
from scipy.optimize import minimize
result = minimize(neg_log_likelihood, initial_params,
                  args=(y, omega), bounds=bounds)
```

Classic Optimisation Methods:

- Gradient Descent** - Iterative gradient minimisation.
- Newton's Method** - Uses second derivatives.
- Quasi-Newton (BFGS)** - Approximates second derivatives.
- Conjugate Gradient** - Efficient for large-scale problems.
- Simplex Method** - Linear programming.
- Lagrange Multipliers** - Constrained optimisation.

Advanced/ML Optimisation Methods:

- Stochastic Gradient Descent (SGD)** - Mini-batch updates.
- Adam** - Adaptive learning rates.
- Bayesian Optimization** - Probabilistic global optimisation.
- Genetic Algorithms** - Evolution-inspired optimisation.
- Simulated Annealing** - Randomised global search.