

# Prediction



UNIVERSITÄT  
HEIDELBERG  
ZUKUNFT  
SEIT 1386



Community health workers in a school, Rabai county, Kenya, September 2021

**WASHA, Takwimu, UKZN, 30 August 2023**

**Till Bärnighausen**, Heidelberg Institute of Global Health, University Hospital and Medical Faculty, Heidelberg University



# Quantitative analyses serve many important functions in health systems research

## 4 FUNCTIONS

1. Description
2. Discovery – unsupervised machine learning
  - Dimension reduction: PCA
  - Cluster analysis: k means
3. **Prediction – supervised machine learning**
  - **Penalized regression**
  - **kNN**
4. Causation
  - IV
  - RDD
  - FE

# Causal and predictive analysis require very different approaches

## COMPARISON

PSM = propensity score matching, IV = instrumental variable analysis, RDD = regression discontinuity design, kNN = k nearest neighbors

	Causation	Prediction
<b>Disciplines</b>	<ul style="list-style-type: none"><li>• Epidemiology</li><li>• Economics</li></ul>	<ul style="list-style-type: none"><li>• Machine learning</li><li>• Computer science</li></ul>
<b>Foundation</b>	<ul style="list-style-type: none"><li>• Theory-based</li><li>• Hypothesis testing</li></ul>	<ul style="list-style-type: none"><li>• Atheoretical</li><li>• Data-driven insight</li></ul>
<b>Purposes</b>	<ul style="list-style-type: none"><li>• Understanding</li><li>• Policy guidance</li><li>• Regulatory approval</li></ul>	<ul style="list-style-type: none"><li>• Intervention targeting</li><li>• Intervention tailoring</li><li>• Now- and forecasting</li></ul>
<b>Goal of approach</b>	<ul style="list-style-type: none"><li>• Minimize bias</li></ul>	<ul style="list-style-type: none"><li>• Optimize bias-variance trade-off</li></ul>
<b>Approach</b>	<ul style="list-style-type: none"><li>• Estimation</li></ul>	<ul style="list-style-type: none"><li>• Training-(validation)-testing</li><li>• Complexity reduction</li></ul>
<b>Example</b>	<ul style="list-style-type: none"><li>• Ordinary multiple regression</li><li>• PSM</li><li>• IV and RDD</li></ul>	<ul style="list-style-type: none"><li>• Regularized multiple regression</li><li>• kNN</li><li>• Neural networks</li></ul>

# What is the generic objective in prediction?

## MATHEMATICAL INTUITION

- **Data:**  $(Y_i, x_{i1}, \dots, x_{ip})$  for  $i = 1, \dots, n$
- **Objective:** predict  $Y$  for a given **new input**  $x_{new} = (x_1, \dots, x_p)$
- Two major **categories** are (machine learning language)
  - **Regression:** continuous data
  - **Classification:** discrete values representing classes

	Cluster membership (reference: cluster 1)			
	Cluster 2	Cluster 3	Cluster 4	Cluster 5
<b>Age</b> (years)	1.055*** (0.0004)	1.038*** (0.005)	1.165*** (0.001)	1.084*** (0.0004)
<b>Education</b> (years)	0.986*** (0.001)	1.063*** (0.010)	1.007*** (0.001)	1.031*** (0.001)
<b>Wealth index</b> (standard deviation)	0.910*** (0.004)	1.412*** (0.051)	1.443*** (0.006)	1.628*** (0.005)
<b>Rural Residence</b> (reference: urban)	1.000*** (0.009)	0.951*** (0.091)	0.990*** (0.012)	1.116*** (0.009)
<b>Religion</b> (reference: Hindu)				
Muslim	1.124*** (0.011)	1.373*** (0.116)	1.558*** (0.015)	1.522*** (0.011)
Christian	1.289*** (0.013)	0.986*** (0.165)	1.570*** (0.018)	1.290*** (0.014)
Sikh	0.919*** (0.030)	1.170*** (0.252)	2.189*** (0.030)	1.768*** (0.025)
Buddhist	1.260*** (0.033)	0.915** (0.410)	1.810*** (0.040)	1.228*** (0.035)
Jain	1.290*** (0.112)	0.018*** (0.00005)	1.199*** (0.128)	1.013*** (0.105)
Jewish	0.003*** (0.00000)	0.526*** (0.00001)	0.029*** (0.00000)	1.334*** (0.00001)
Zoroastrian	4.175*** (0.257)	0.580*** (0.0002)	3.354*** (0.278)	3.313*** (0.282)
No religion	2.492*** (0.168)	0.112*** (0.0002)	5.235*** (0.193)	2.173*** (0.186)
Other religions	1.860*** (0.031)	0.190*** (0.001)	2.425*** (0.041)	1.505*** (0.036)
Constant	0.244*** (0.018)	0.001 (0.196)	0.003 (0.027)	0.056*** (0.019)

**Penalized regression takes us away from the unbiased coefficient estimators – to achieve higher predictive accuracy**

**MATHEMATICAL INTUITION**

- **OLS**  $\min RSS$
- **LASSO**  $\min(RSS + \delta_1 \sum_{j=1}^p |\beta_j|)$
- **Ridge**  $\min(RSS + \delta_2 \sum_{j=1}^p \beta_j^2)$
- **Elastic net**  $\min(RSS + \delta_1 \sum_{j=1}^p |\beta_j| + \delta_2 \sum_{j=1}^p \beta_j^2)$

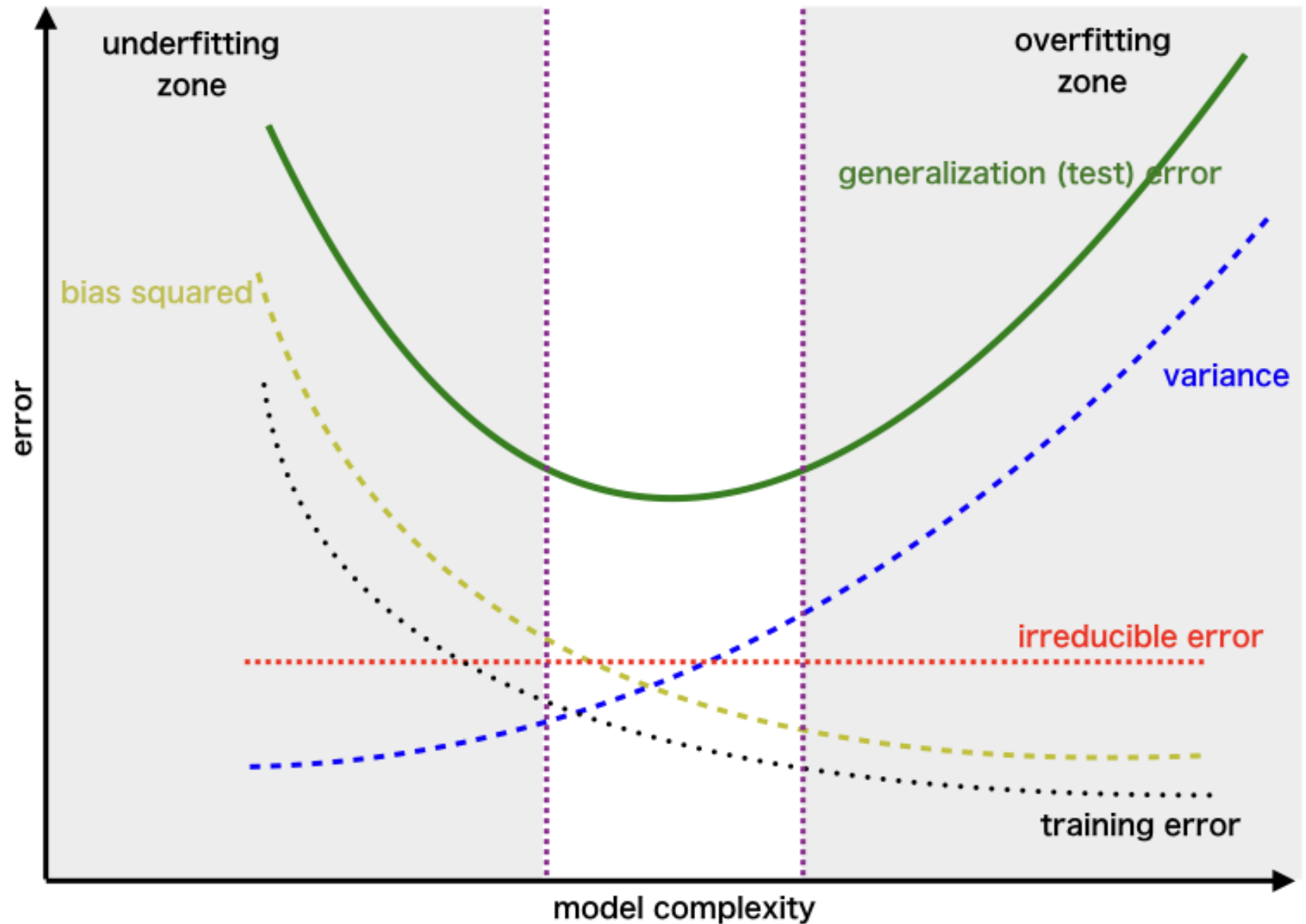


# For best prediction, we trade-off bias and variance

## CONCEPTS

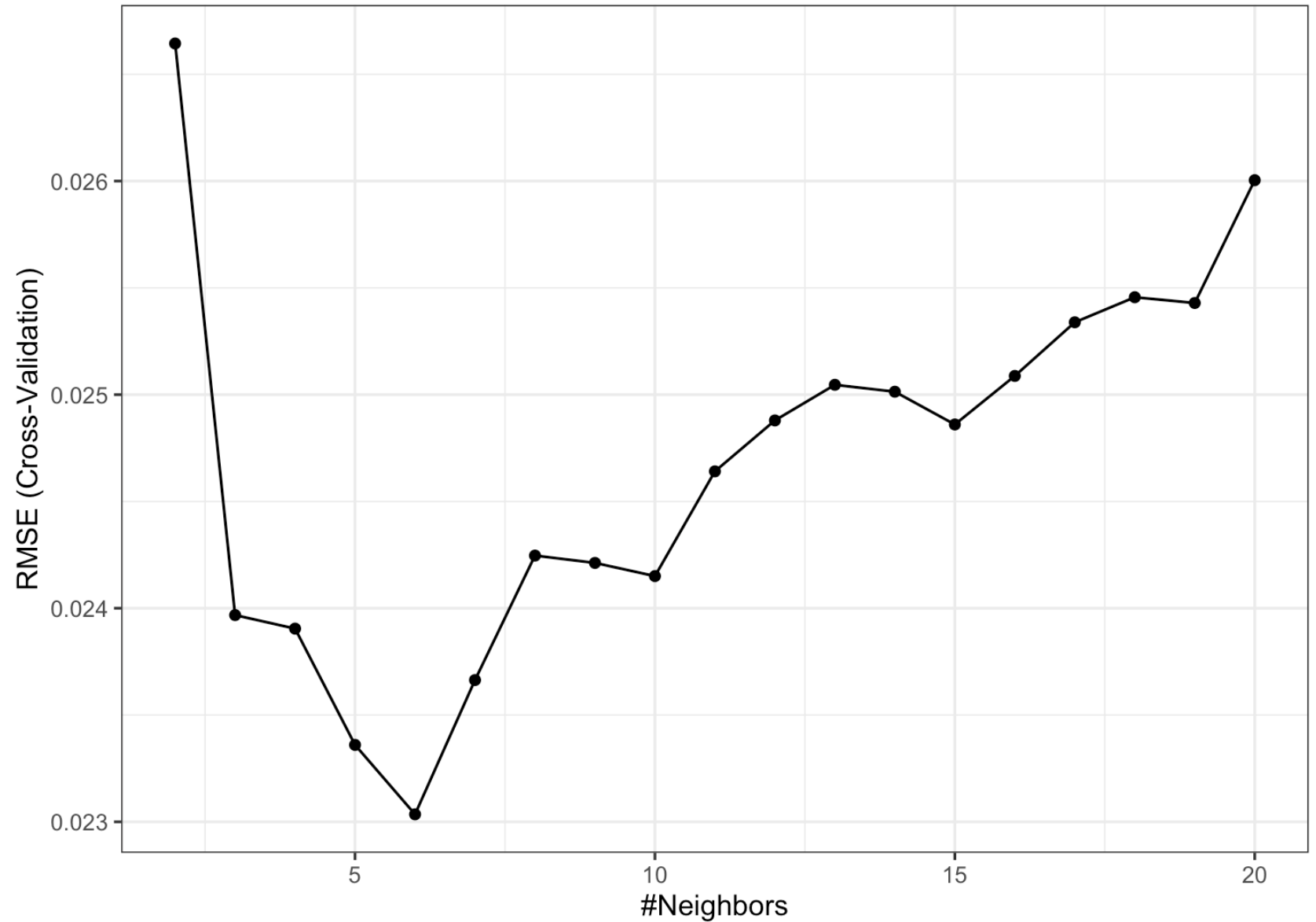
Source:

<https://www.geeksforgeeks.org/ml-bias-variance-trade-off/>



**We  
expect/hope  
to find a  
minimum in  
a loss  
function  
value**

**EXAMPLE**



RMSE = root mean squared error



# **k-nearest neighbors (kNN) is a simple non-parametric supervised learning method**

## **OVERVIEW**

- Fix & Hodges 1951
- Non-parametric
- Memoryless
- Choice of  $k$  defines locality
- As a local method: strong for low-dimensional large data
- Weak for understanding

# How do we define near?

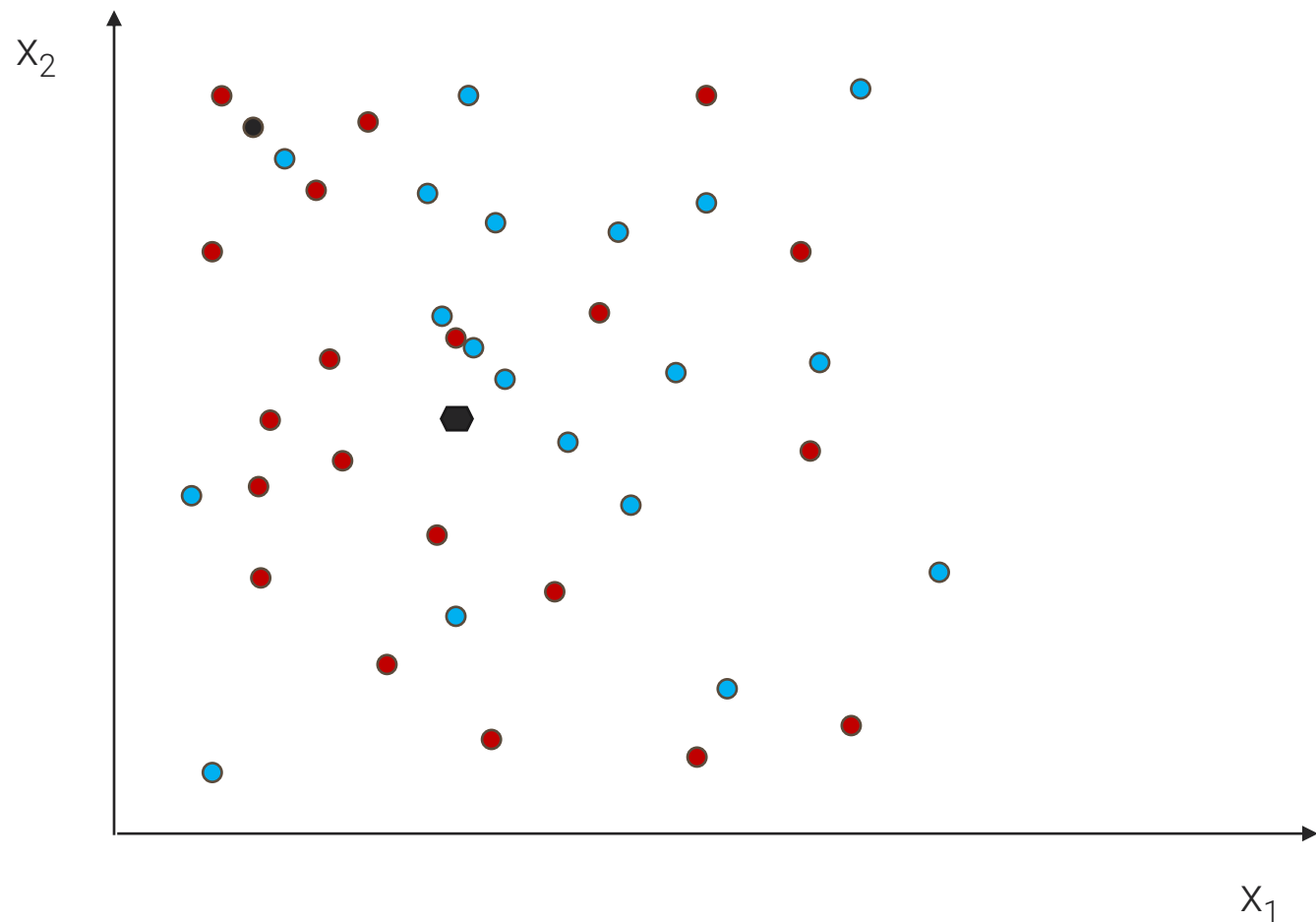
## MINKOWSKI – DISTANCE METRICS

In  $\mathbb{R}^q$

$$\|\mathbf{x}' - \mathbf{x}_j\|^p = \left( \sum_{i=1}^q |(\mathbf{x}')_i - (\mathbf{x}_j)_i|^p \right)^{1/p}$$

# For the binary case, the intuition is the majority vote

## BASIC IDEA



# kNN works for binary classification ...

## FORMULA

$$\mathcal{Y} = \{1, -1\}$$

$$f_{\text{KNN}}(\mathbf{x}') = \begin{cases} 1 & \text{if } \sum_{i \in \mathcal{N}_K(\mathbf{x}')} y_i \geq 0 \\ -1 & \text{if } \sum_{i \in \mathcal{N}_K(\mathbf{x}')} y_i < 0 \end{cases}$$

... and for multiple-class classification ...

FORMULA

$$f_{\text{KNN}}(\mathbf{x}') = \arg \max_{y \in \mathcal{Y}} \sum_{i \in \mathcal{N}_K(\mathbf{x}')} \mathcal{I}(y_i = y)$$

... as well as for regression

FORMULA

$$\mathbf{f}_{KNN}(\mathbf{x}') = \frac{1}{K} \sum_{i \in \mathcal{N}_K(\mathbf{x}')} \mathbf{y}_i$$

# Weighted knn may boost performance

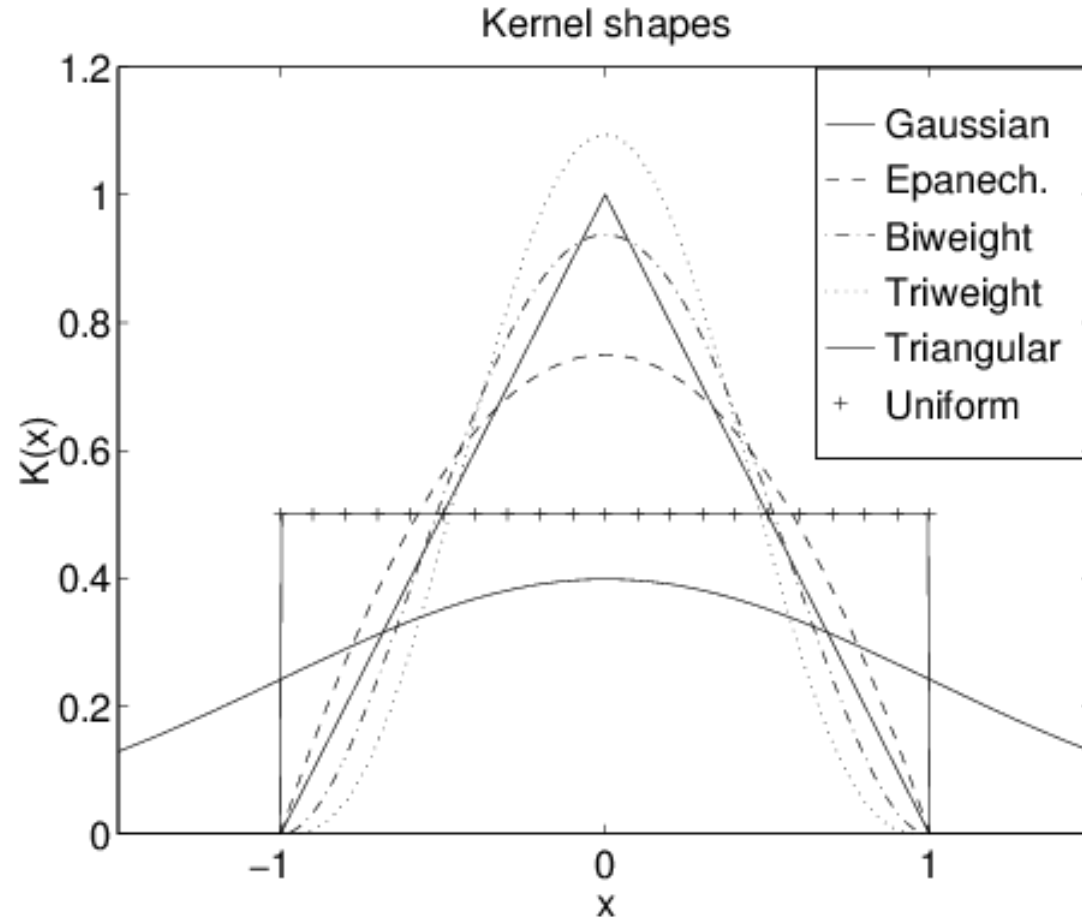
WEIGHTED KNN AND KERNEL TRICK

$$K(x, y) = \langle \varphi(x), \varphi(y) \rangle$$



# We can use kernel functions for weighting distance

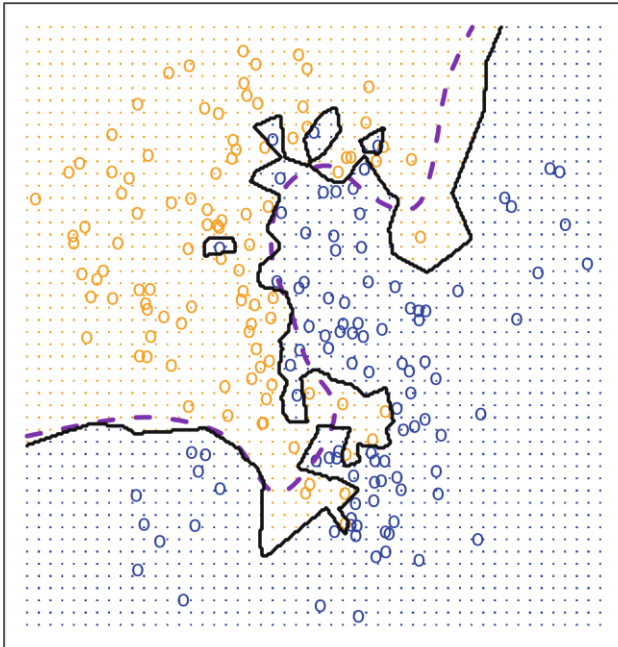
## KERNEL TYPES



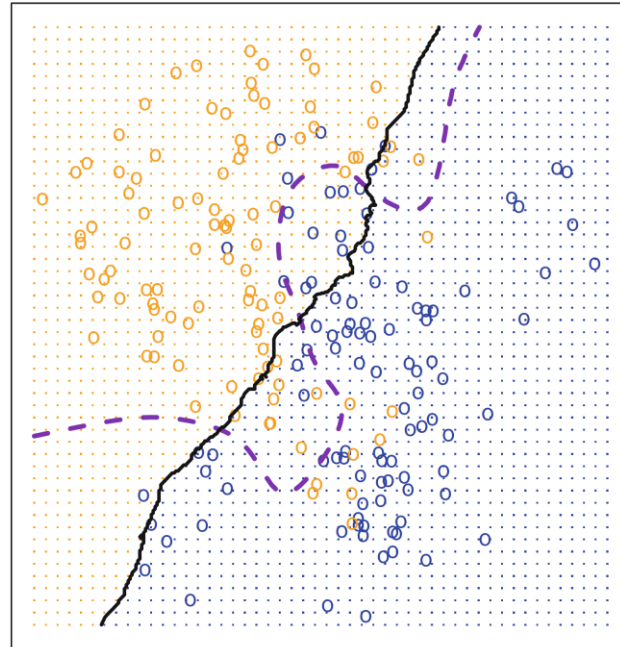
# kNN needs to be 'tuned'

## EXAMPLE

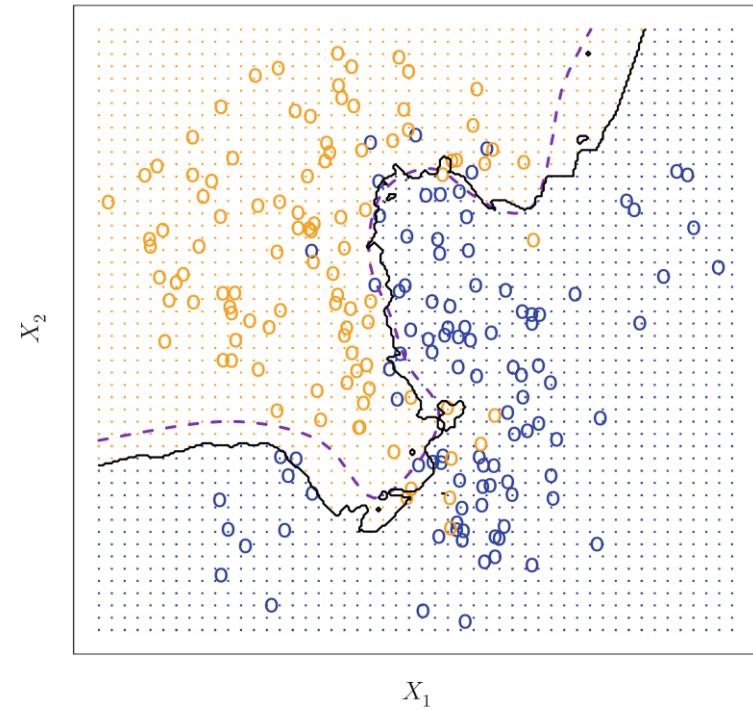
KNN: K=1



KNN: K=100



KNN: K=10



# Hyperparameters are used during the learning process – parameters are the result

## EXAMPLES

- **Hyperparameters**

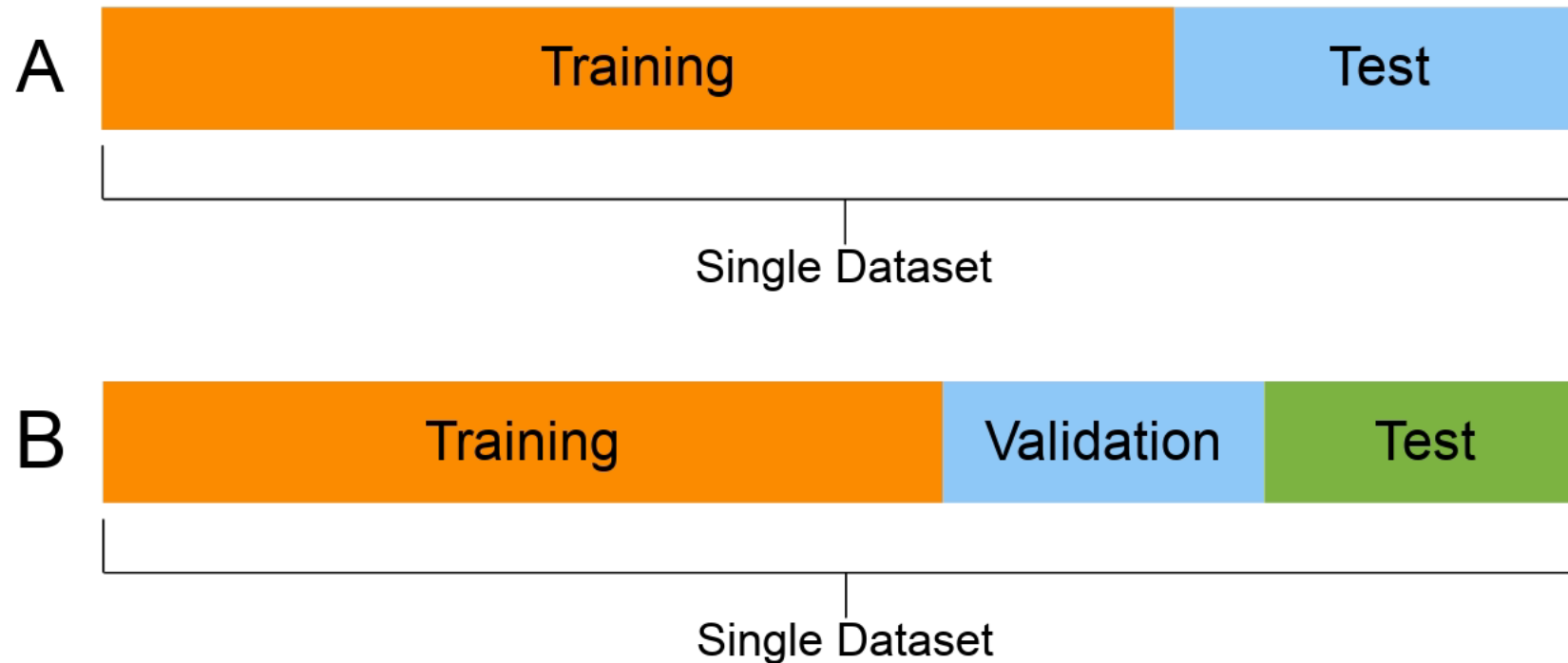
- Train-test split ratio
- Choice of optimization algorithm
- Number of principal components in PCA
- Kernel size
- $K$  in  $k$  means cluster analysis
- $K$  in knn analysis
- Penalty term weights in penalized regression

- **Parameters**

- Weights that generate principal components as linear combinations of the original data
- Cluster centroids in cluster analysis
- The actual nearest neighbors in knn
- Coefficients in penalized regression

# Predictive practice is fundamentally different from causal practice

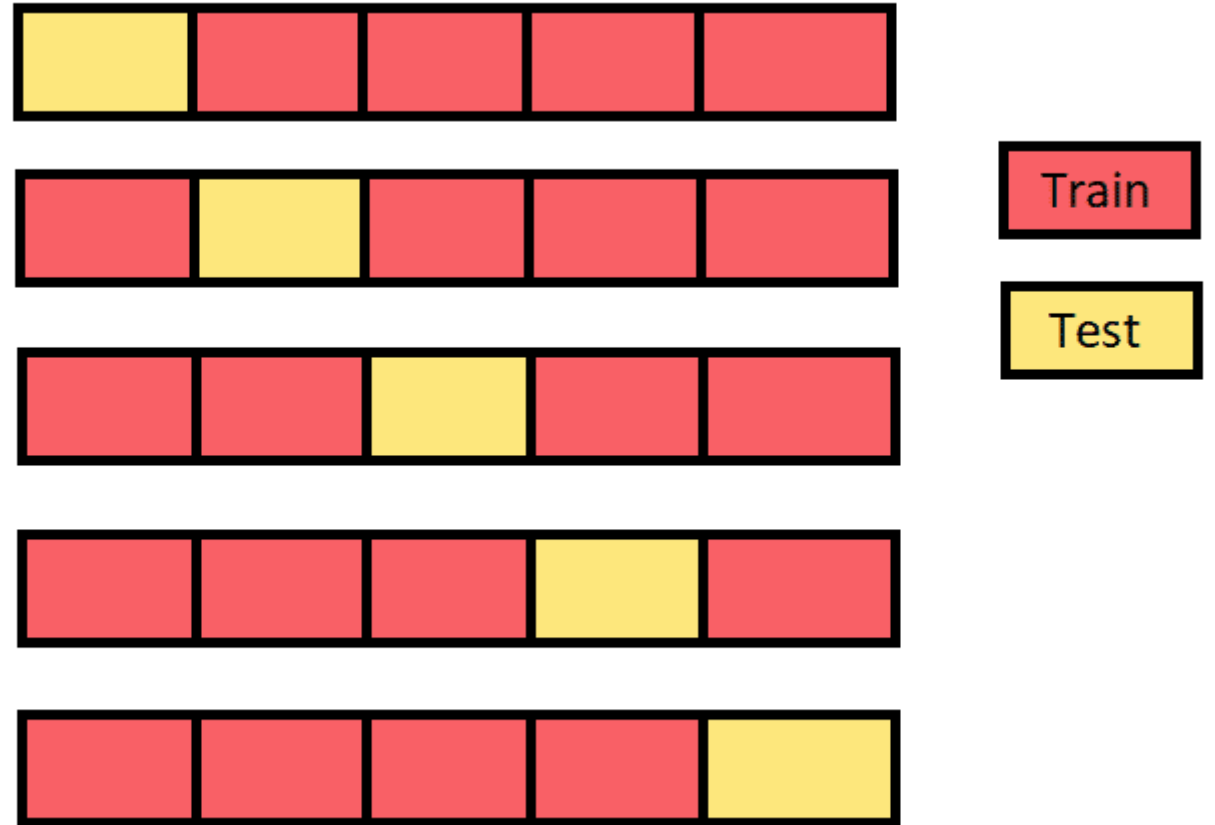
## VALIDATION DATASET FOR HYPERPARAMETER TUNING



[https://commons.wikimedia.org/wiki/File:ML\\_dataset\\_training\\_validation\\_test\\_sets.png](https://commons.wikimedia.org/wiki/File:ML_dataset_training_validation_test_sets.png)

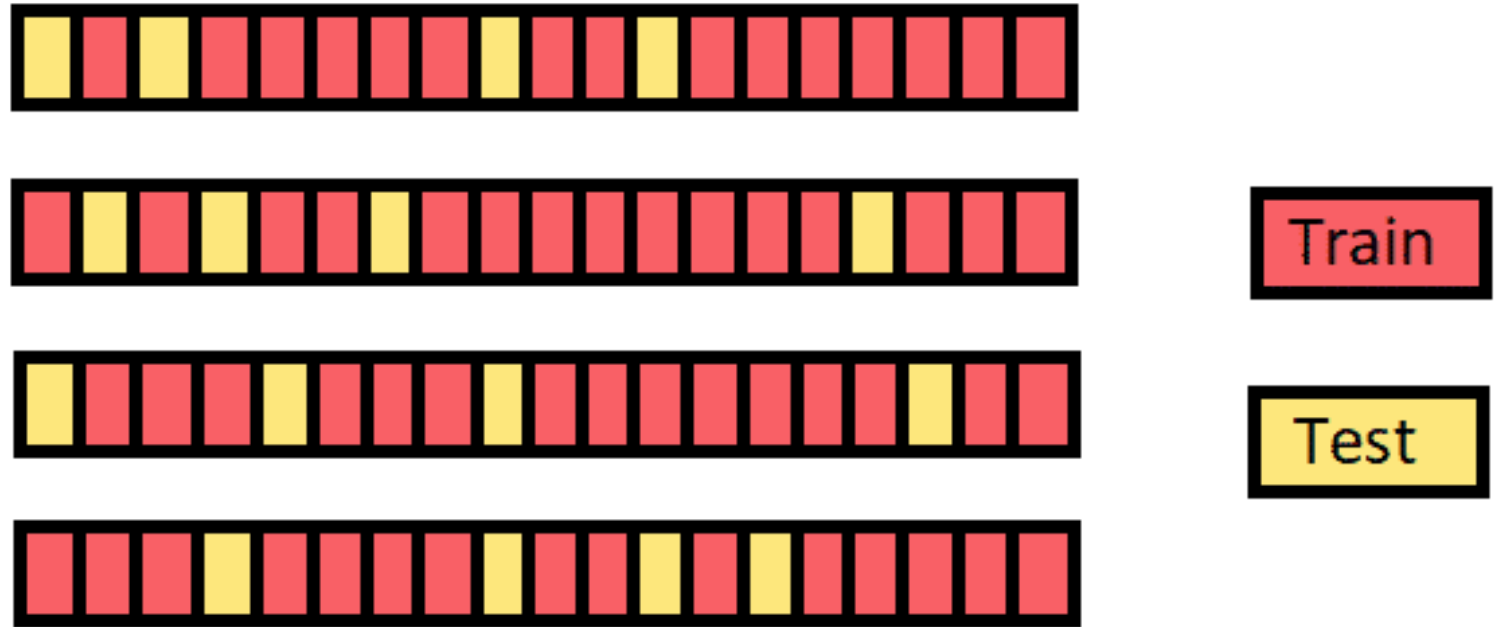
# Cross-validation increases the efficiency of hyperparameter tuning

K-FOLD



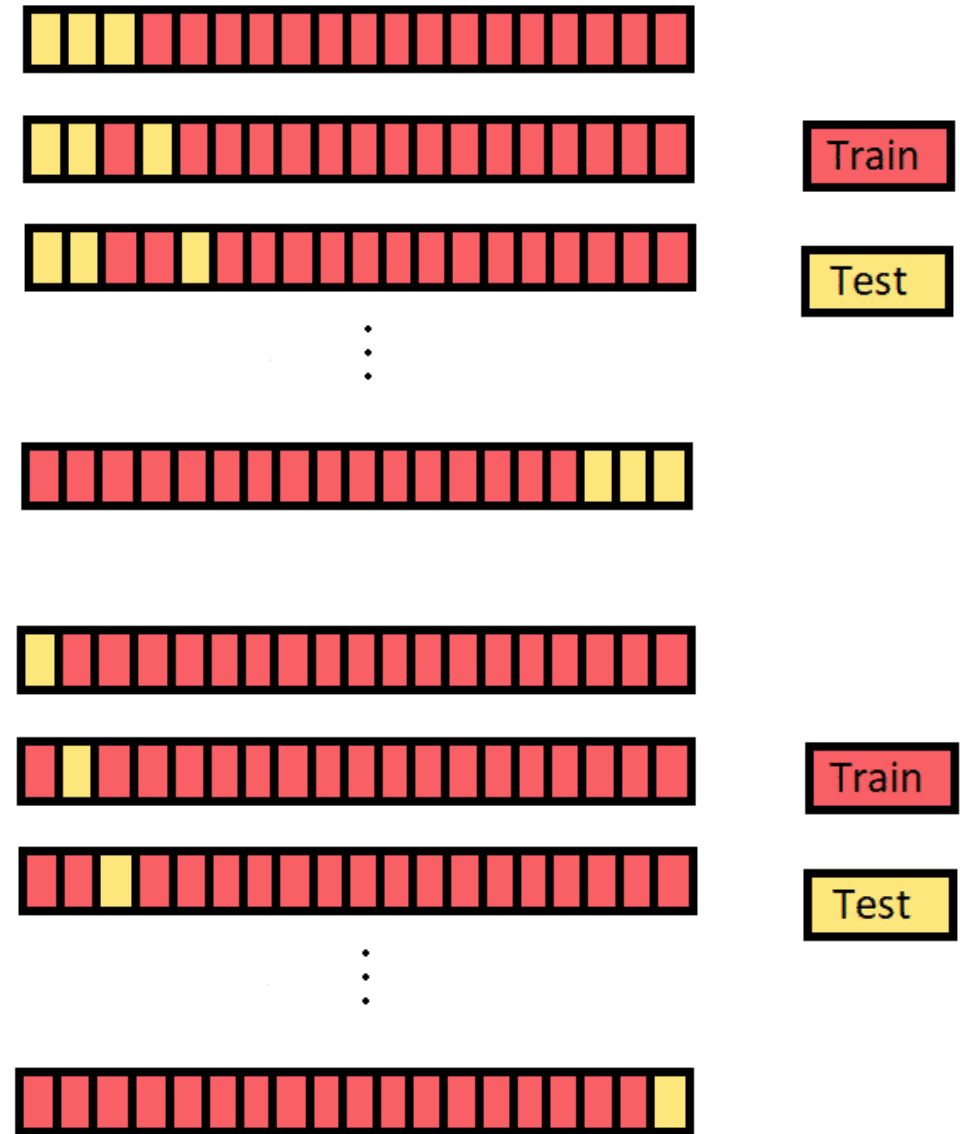
# Cross-validation increases the efficiency of hyperparameter tuning

MONTE CARLO



# Cross-validation increases the efficiency of hyperparameter tuning

LEAVE P OUT, LEAVE ONE OUT

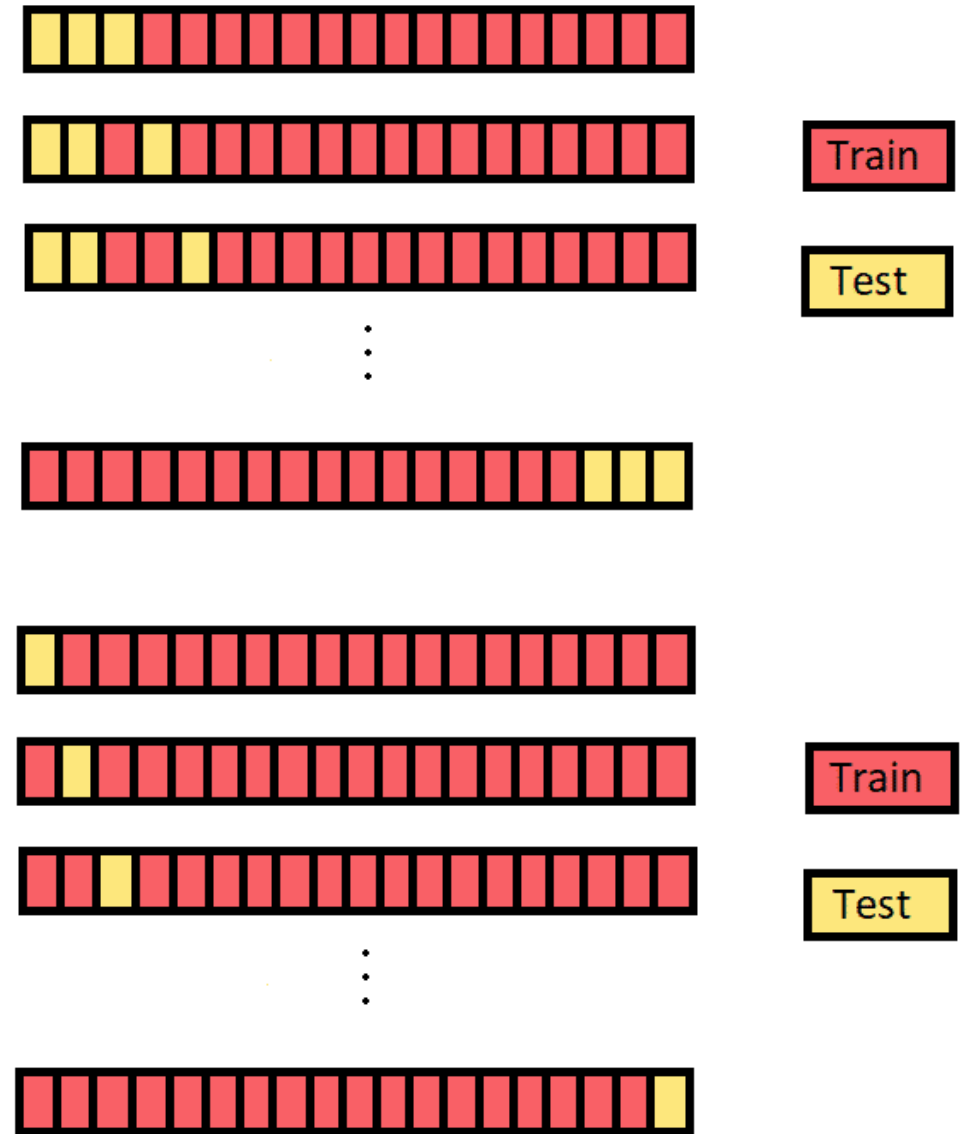


<https://aiaspirant.com/cross-validation/>



# Cross-validation increases the efficiency of hyperparameter tuning

LEAVE P OUT, LEAVE ONE OUT



<https://aiaspirant.com/cross-validation/>

# The curse of dimensionality is particularly problematic for similarity-based algorithms

BELLMAN 1961

- General: Sample size needed to estimate a function with a given level of accuracy grows exponentially with the dimensionality of the data
- Specifically for our topic: For similarity-based algorithms (k means, knn), the number of instances that need to be accessed for precise estimation grows exponentially with data dimensionality

# Dimension reduction is typically an important data preprocessing step for kNN

## EXAMPLES

- ***Linear*** –
  - Unsupervised: PCA, FA, SVM, ...
  - Supervised: LDA, PLS ...
- ***Non-linear*** – kernel PCA, FAMD, t-SNE, ...

PCA = principal component analysis, SVM = support vector machine, LDA = linear discriminant analysis, PLS = partial least squares, FAMD = factor analysis for mixed data, t-SNE = t-distributed stochastic neighbor embedding