السورية الخاصة الجامعة
SYRIAN PRIVATE UNIVERSITY

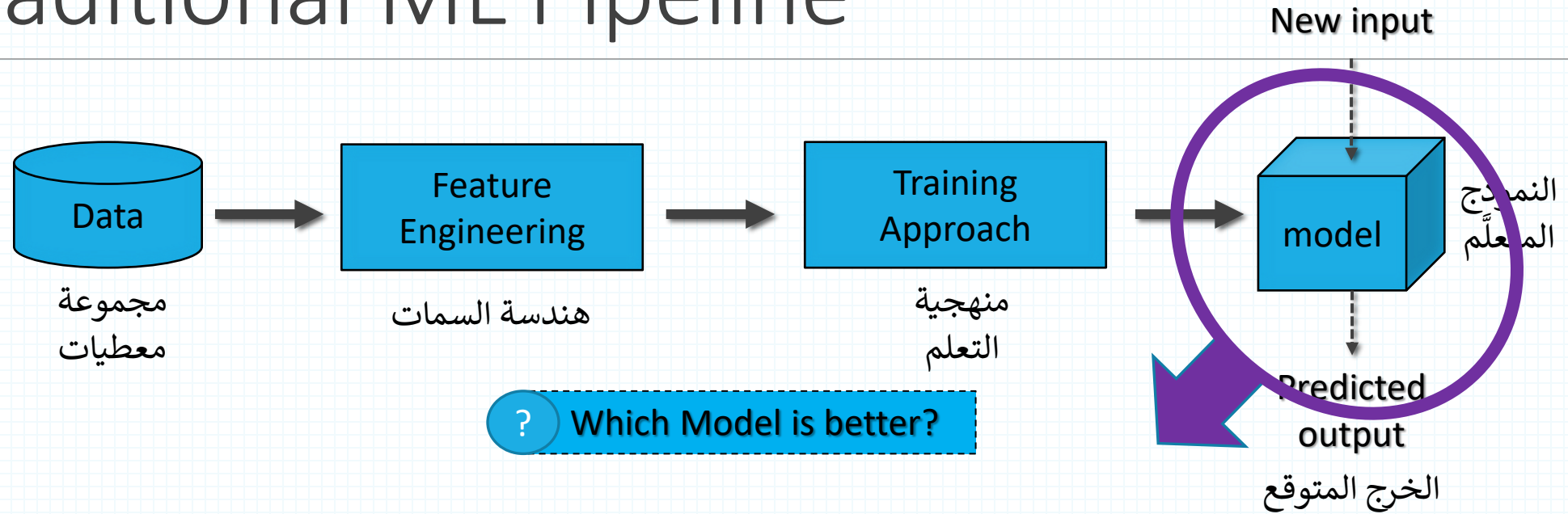| المحاضرة الرابعة | كلية الهندسة المعلوماتية | مقرر تعلم الآلة |
|---|---|---|

# استراتيجيات التدريب والاختبار

د. رياض سنبل

# Traditional ML Pipeline

New input

```
Data → Feature Engineering → Training Approach → model
```

Data
مجموعة معطيات

Feature Engineering
هندسة السمات

Training Approach
منهجية التعلم

model
النموذج المعلَّم

**?  Which Model is better?**

Predicted output
الخرج المتوقع

*"All models are wrong; some are useful."*

⸺George E. P. Box

# Traditional ML Pipeline

New input



| Data | Feature Engineering | Training Approach | model |

مجموعة معطيات     هندسة السمات     منهجية التعلم     النموذج المتعلّم

Predicted output

الخرج المتوقع

**? Which Model is better?**

- Training algorithm ( Example?)
- Parameters (Hyperparameters)?
- Handle unseen cases

Estimation Strategy      Evaluation Metrics

# Hyperparameters

- Hyperparameters are the explicitly specified parameters that control the training process.

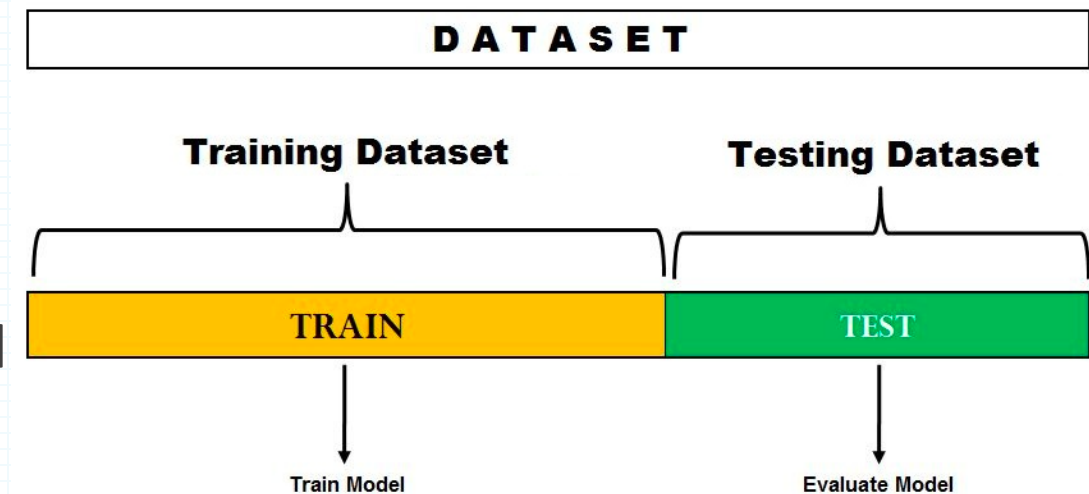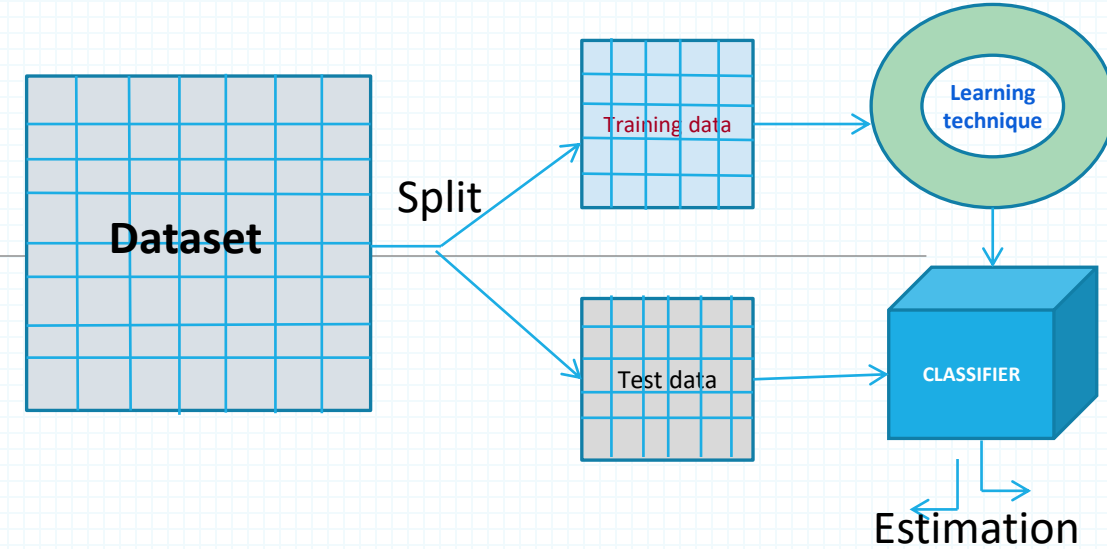- What are the hyperparameters for decision trees?

```
class sklearn.tree.DecisionTreeClassifier(*, criterion='gini', splitter='best', max_depth=None, min_samples_split=2,
min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None, random_state=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, class_weight=None, ccp_alpha=0.0)                                              [source]
```
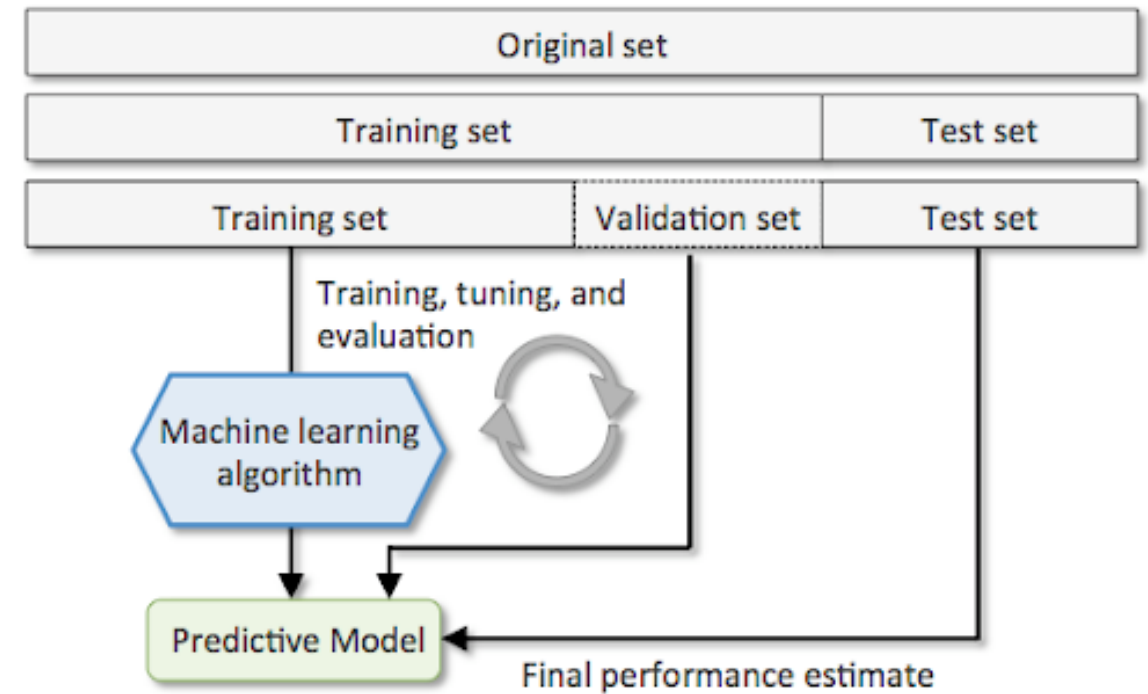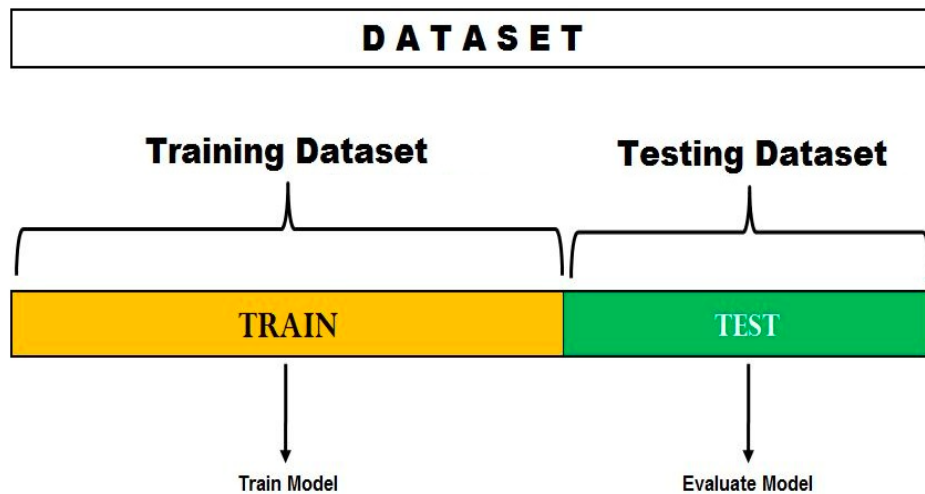
# Estimation Strategies

# The holdout method

- Split dataset into two groups
  - Training set: used to train the classifier
  - Test set: used to estimate the error rate of the trained classifier.

- Ratio of training and testing sets is at the discretion of analyst;
  - Typically **1:1 or 2:1**, and there is a trade-off between these sizes of these two sets.
  - If the training set is too large, then model may be good enough, but estimation may be less reliable due to small testing set and vice-versa.

Dataset → Split → Training data → Learning technique → CLASSIFIER

Test data → CLASSIFIER

Estimation

**DATASET**

**Training Dataset** | **Testing Dataset**

**TRAIN** | **TEST**

Train Model | Evaluate Model

# Holdout method

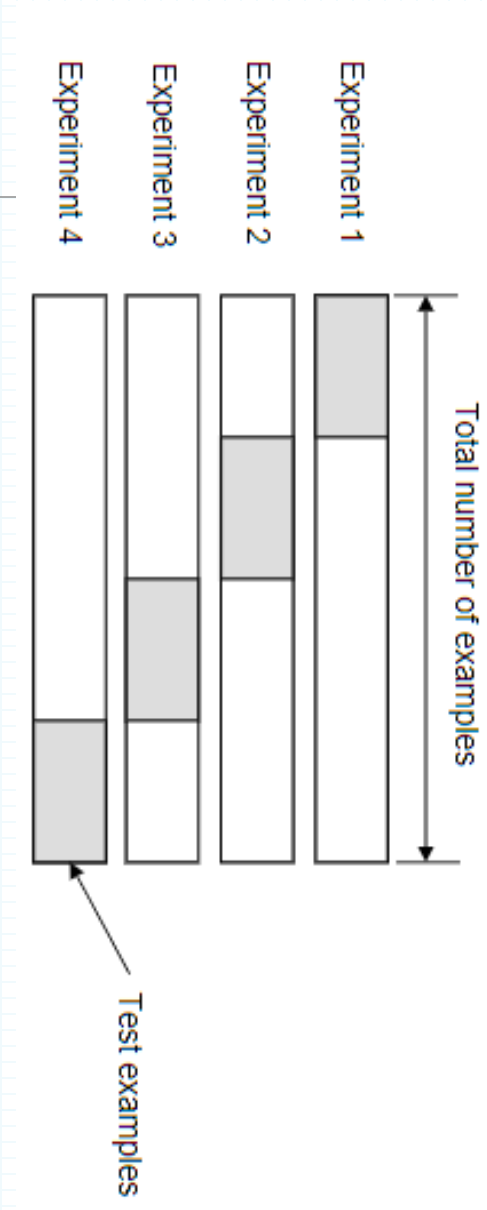But how can we tune hyper-parameters?

# Holdout method - Drawbacks

- The holdout method has **two basic drawbacks**
  - In problems where we have a sparse dataset we may not be able to afford **the "luxury" of setting aside a portion** of the dataset for testing
  - Since it is **a single train-and-test experiment**, the holdout estimate of error rate will be misleading if we happen to get an "unfortunate" split

- The limitations of the holdout can be overcome with a family of re-sampling methods at the expense of higher computational cost
  - Cross Validation
    - Random Subsampling
    - K-Fold Cross-Validation
    - Leave-one-out Cross-Validation
  - Bootstrap

# K-Fold Cross-validation

- Create a K-fold partition of the dataset

- For each of K experiments, use K-1 folds for training and a different fold for testing

- The advantage of K-Fold Cross validation is that **all the examples** in the dataset are eventually used for both training and testing

- The true error is estimated as the average error rate on test examples

# Leave-one-out Cross-Validation

- Leave-one-out is the degenerate case of K-Fold Cross Validation, where K is chosen as the total number of examples
  - For a dataset with N examples, perform N experiments
  - For each experiment use N-1 examples for training and the remaining example for testing