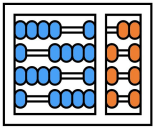# Capacitação profissional em tecnologias de Inteligência Artificial

## Machine Learning Overview

**Prof. Edson Borin**

https://www.ic.unicamp.br/~edson
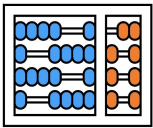
Institute of Computing - UNICAMP
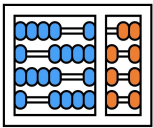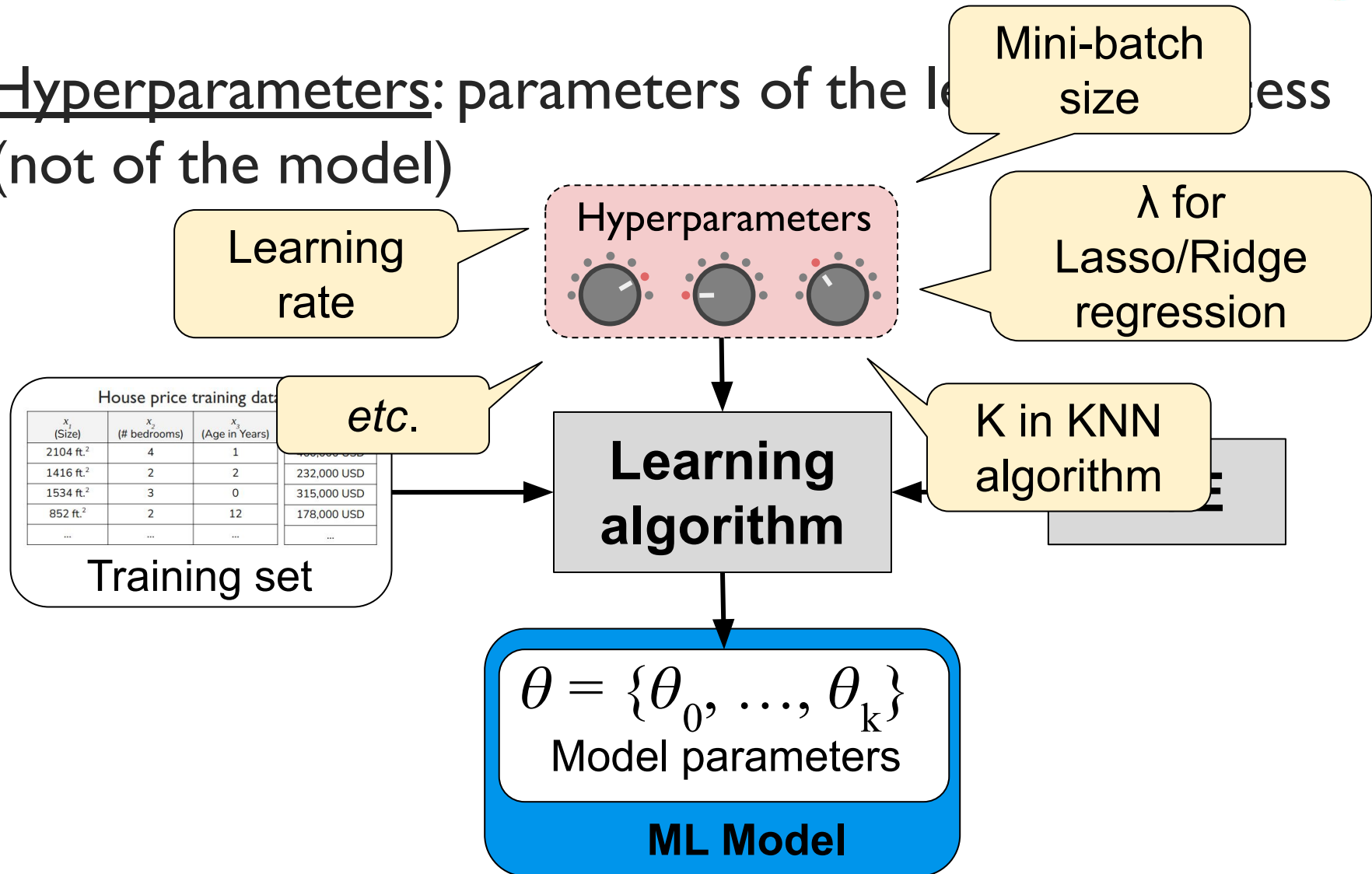
# Hyperparameters tuning

Hyperparameters: parameters of the learning process (not of the model)

<u>Hyperparameters</u>: parameters of the le          ess (not of the model)

Mini-batch size

Learning rate

Hyperparameters

λ for Lasso/Ridge regression

*etc*.

House price training dat

| $x_1$ (Size) | $x_2$ (# bedrooms) | $x_3$ (Age in Years) | |
|---|---|---|---|
| 2104 ft.² | 4 | 1 | |
| 1416 ft.² | 2 | 2 | 232,000 USD |
| 1534 ft.² | 3 | 0 | 315,000 USD |
| 852 ft.² | 2 | 12 | 178,000 USD |
| ... | ... | ... | ... |

Training set

**Learning algorithm**

K in KNN algorithm

$$\theta = \{\theta_0, \dots, \theta_k\}$$

Model parameters

**ML Model**

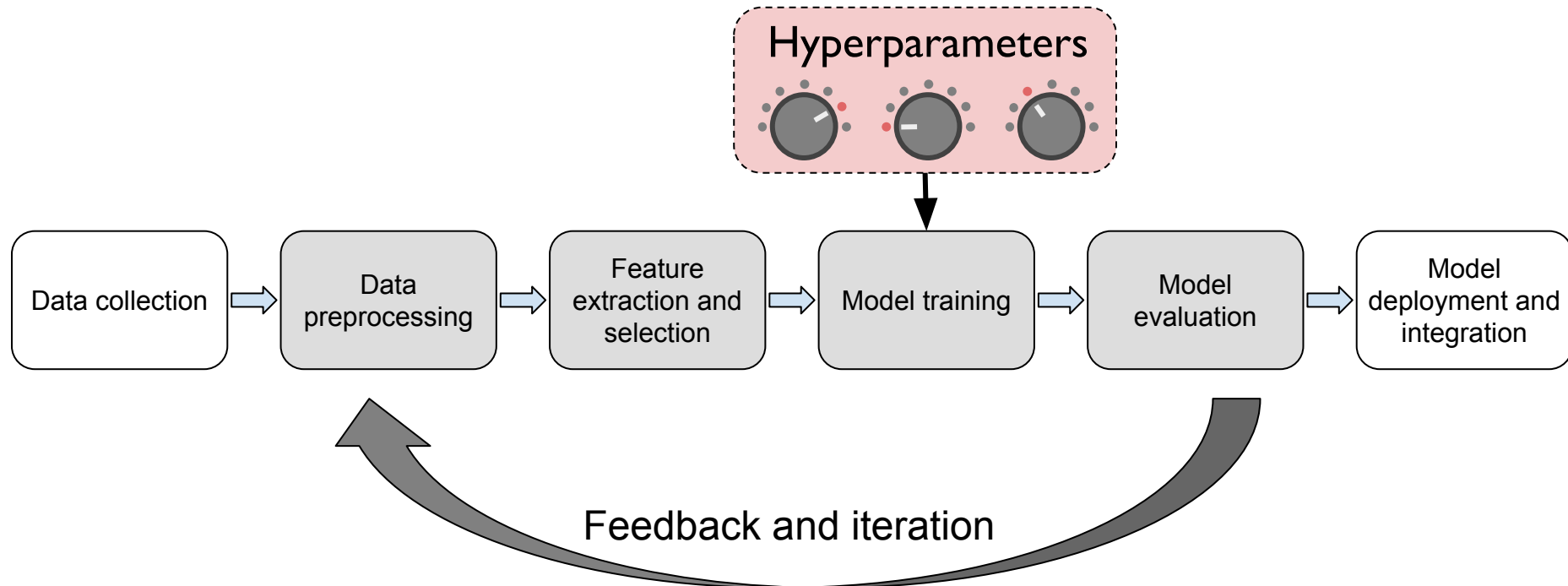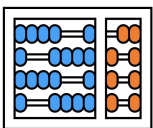# ML Process - Hyperparameters tuning

<u>Hyperparameter tuning</u>: finding the best combination of hyperparameters that causes the learning process to produce the best model!

<u>Hyperparameter tuning</u>: finding the best combination of hyperparameters that causes the learning process to produce the best model!

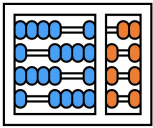- Example: scikit learn SVC models with RBF kernel
  - $C$: regularization parameter
  - $\gamma$: Kernel coefficient
  - Some hyperparameters combinations:
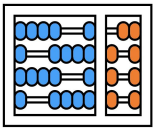    - $(C, \gamma) \in \{ (10, 0.1), (10, 0.2), (100, 0.1), (100, 1.0)\}$

# ML Process - Hyperparameters tuning

<u>Search approach</u>: strategy to evaluate the combinations of hyperparameters

- Several approaches
  - Grid search
  - Random search
  - Bayesian optimization
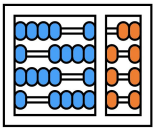  - …

<u>Search approach</u>: **Grid Search**

- Grid search (or parameter sweep) consists on a exhaustive search on a grid defined by the cartesian product of all parameters candidate values

- Example 1:
  - For $C \in \{10, 50, 100\}$, $\gamma = \{0.1, 0.2, 0.5, 1.0\}$, defined by the practitioner
  - $C \times \gamma = \{$ (10, 0.1), (10, 0.2), (10, 0.5), (10, 1.0),
    (50, 0.1), (50, 0.2), (50, 0.5), (50, 1.0),
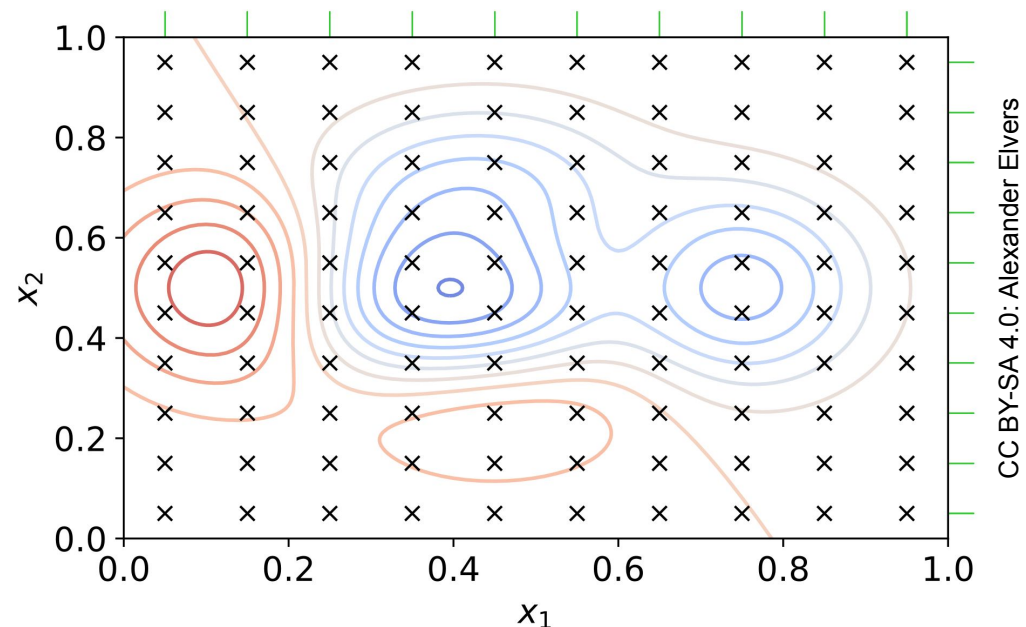    (100, 0.1), (100, 0.2), (100, 0.5), (100, 1.0) $\}$

## Search approach: **Grid Search**

- Grid search (or parameter sweep) consists on a exhaustive search on a grid defined by the cartesian product of all parameters candidate values

- Example:

  - $x_1$ = np.arange(0.05, 1.0, 0.1)
  - $x_2$ = np.arange(0.05, 1.0, 0.1)



CC BY-SA 4.0: Alexander Elvers

# ML Process - Hyperparameters tuning

## Search approach: **Random Search**

- Randomly selects values for hyperparameters
  - Bounds (max, min) values are defined by the user

- Example:
  - $x_1 \in = [0.0, 1.0]$
  - $x_2 \in = [0.0, 1.0]$



CC BY-SA 4.0: Alexander Elvers

<u>Search approach</u>:  **Bayesian optimization**

- Selects next set of hyperparameters to evaluate based on the performance of previous ones
  - Can be adjusted to favor exploring unknown regions or to focus on best regions found so far

- Example:
  - $x_1 \in$ = [0.0, 1.0]
  - $x_2 \in$ = [0.0, 1.0]

CC BY-SA 4.0: Alexander Elvers

# Dataset splitting

# ML Process - Dataset splitting

- On supervised learning tasks, the dataset is usually split into two subsets: training and test
  - <u>Training set</u>: used to train the model (i.e., adjust $\theta$)
  - <u>Test set</u>: check the model generalization
    - Represents new/unseen data

# ML Process - Dataset splitting

- On supervised learning tasks, the dataset is usually split into two subsets: training and test
  - <u>Training set</u>: used to train the model (i.e., adjust $\theta$)
  - <u>Test set</u>: check the model generalization
    - Represents new/unseen data

Dataset

Split

Training set

Test set

**WARNING**: Bad split!

- On supervised learning tasks, the dataset is usually split into two subsets: training and test
  - <u>Training set</u>: used to train the model (i.e., adjust $\theta$)
  - <u>Test set</u>: check the model generalization
    - Represents new/unseen data

Dataset



Split

Training set

Test set



**Random is better!**
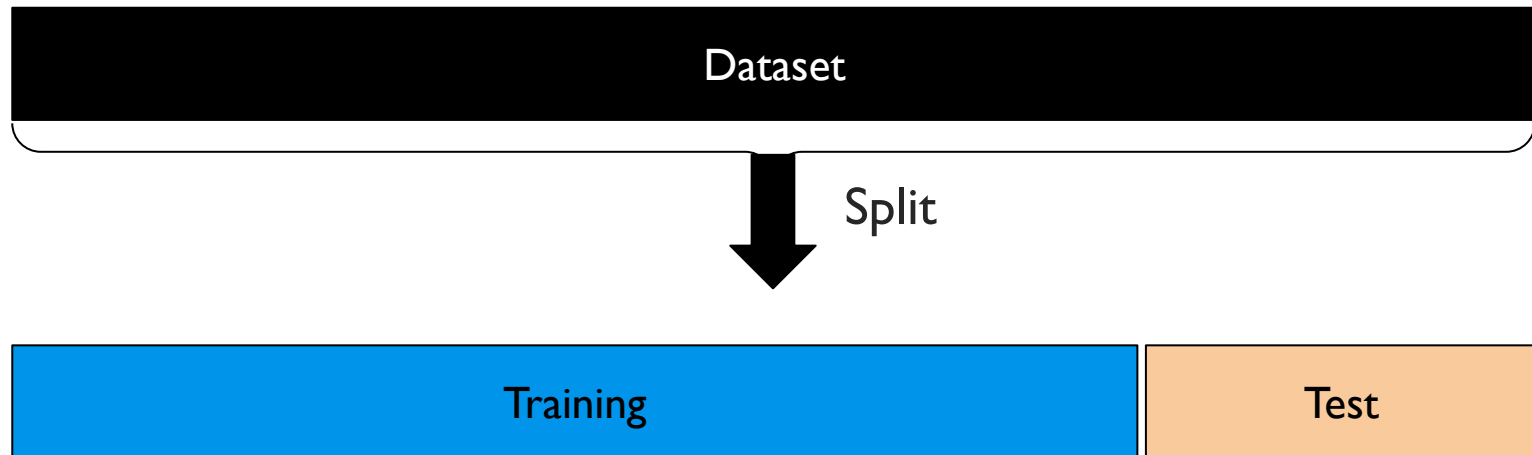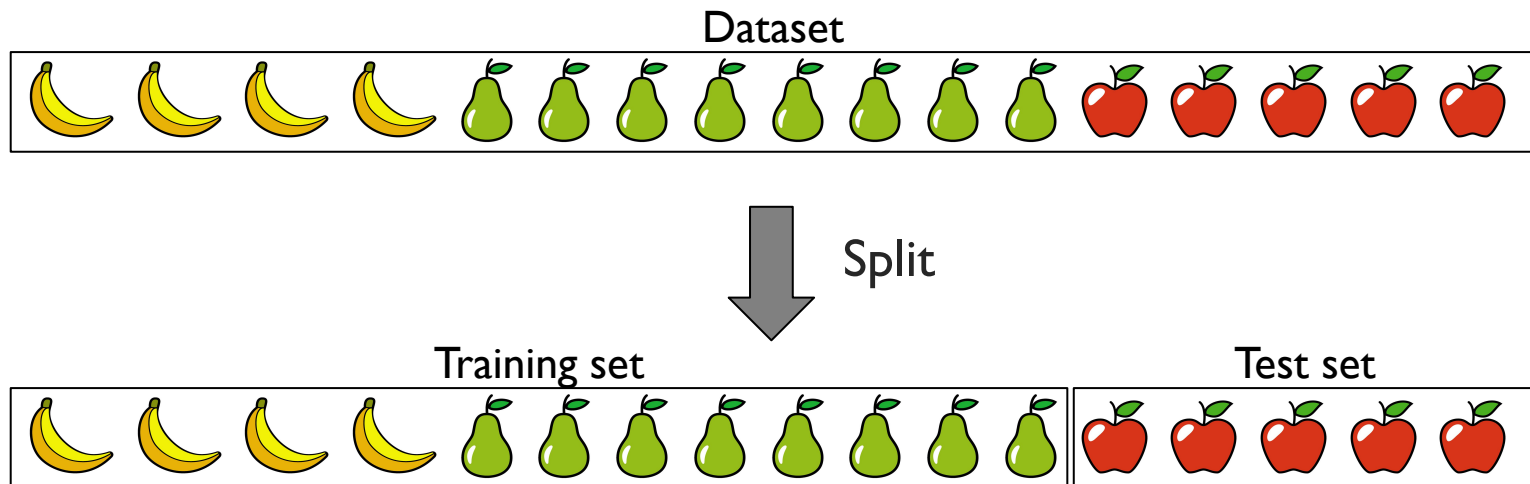
# ML Process - Dataset splitting

- ## On supervised learning tasks, the dataset is usually split into two subsets: training and test
  - ○ <u>Training set</u>: used to train the model (i.e., adjust $\theta$)
  - ○ <u>Test set</u>: check the model generalization
    - ■ Represents new/unseen data

```
from sklearn.model_selection import train_test_split
X = ['y', 'y', 'y', 'r', 'r', 'r', 'r', 'g', 'g', 'g', 'g', 'g', 'g']
y = ['B', 'B', 'B', 'A', 'A', 'A', 'A', 'P', 'P', 'P', 'P', 'P', 'P']

# Split arrays or matrices into random train and test subsets.
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=.2)

print(X_train, X_test)
print(y_train, y_test)
```

```
['r', 'y', 'g', 'r', 'g', 'g', 'y', 'r', 'g', 'g'] ['r', 'y', 'g']
['A', 'B', 'P', 'A', 'P', 'P', 'B', 'A', 'P', 'P'] ['A', 'B', 'P']
```

# ML Process - Dataset splitting

- ● On supervised learning tasks, the dataset is usually split into two subsets: training and test
  - ○ <u>Training set</u>: used to train the model (i.e., adjust $\theta$)
  - ○ <u>Test set</u>: check the model generalization
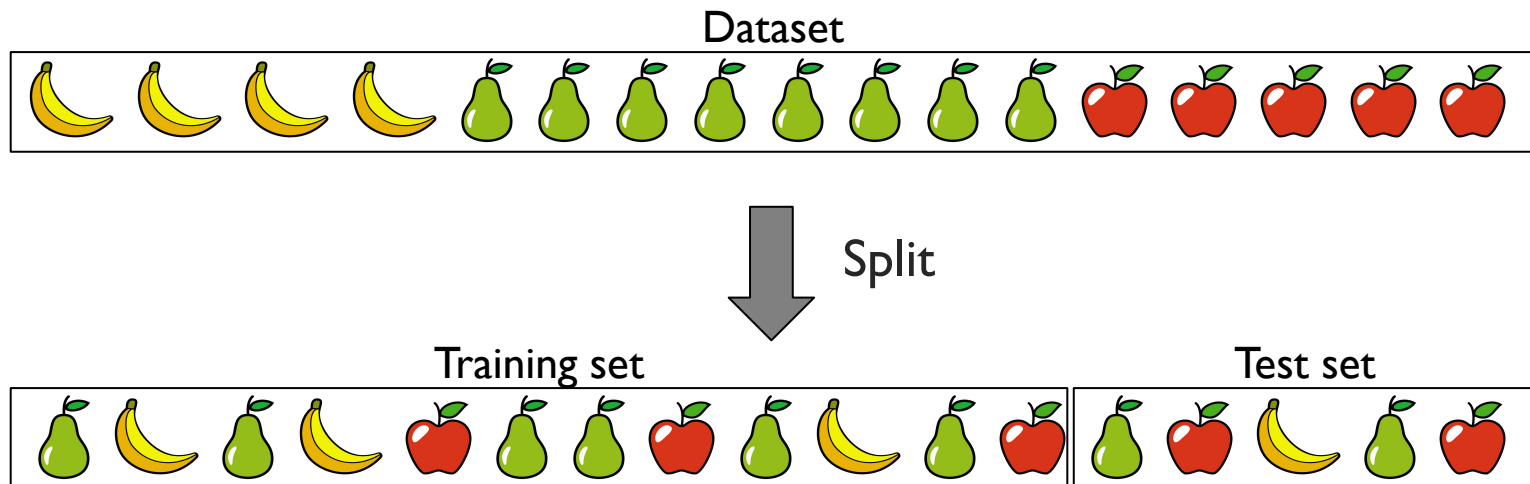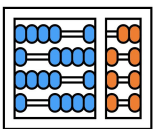    - ■ Represents new/unseen data

```
from sklearn.model_selection import train_test_split
X = ['y', 'y', 'y', 'r', 'r', 'r', 'r', 'g', 'g', 'g', 'g', 'g', 'g']
y = ['B', 'B', 'B', 'A', 'A', 'A', 'A', 'P', 'P', 'P', 'P', 'P', 'P']

# Split arrays or matrices into random train and test subsets.
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=.2)

print(X_train, X_test)
print(y_train, y_test)
```

> Must be chosen carefully!

```
['r', 'y', 'g', 'r', 'g', 'g', 'y', 'r', 'g', 'g
['A', 'B', 'P', 'A', 'P', 'P', 'B', 'A', 'P', 'P
```

# ML Process - Dataset splitting

- On supervised learning tasks, the dataset is usually split into two subsets: training and test
  - <u>Training set</u>: used to train the model (i.e., adjust $\theta$)
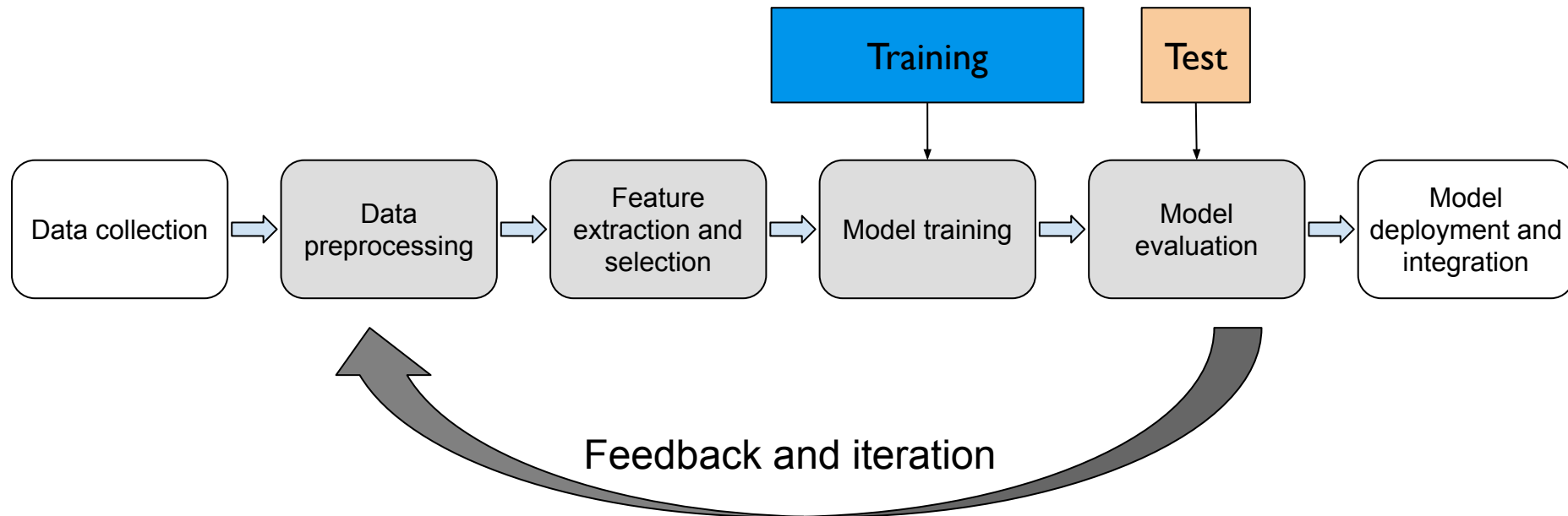  - <u>Test set</u>: check the model generalization
    - Represents new/unseen data

## Never train your model using the test data!

# ML Process - Dataset splitting

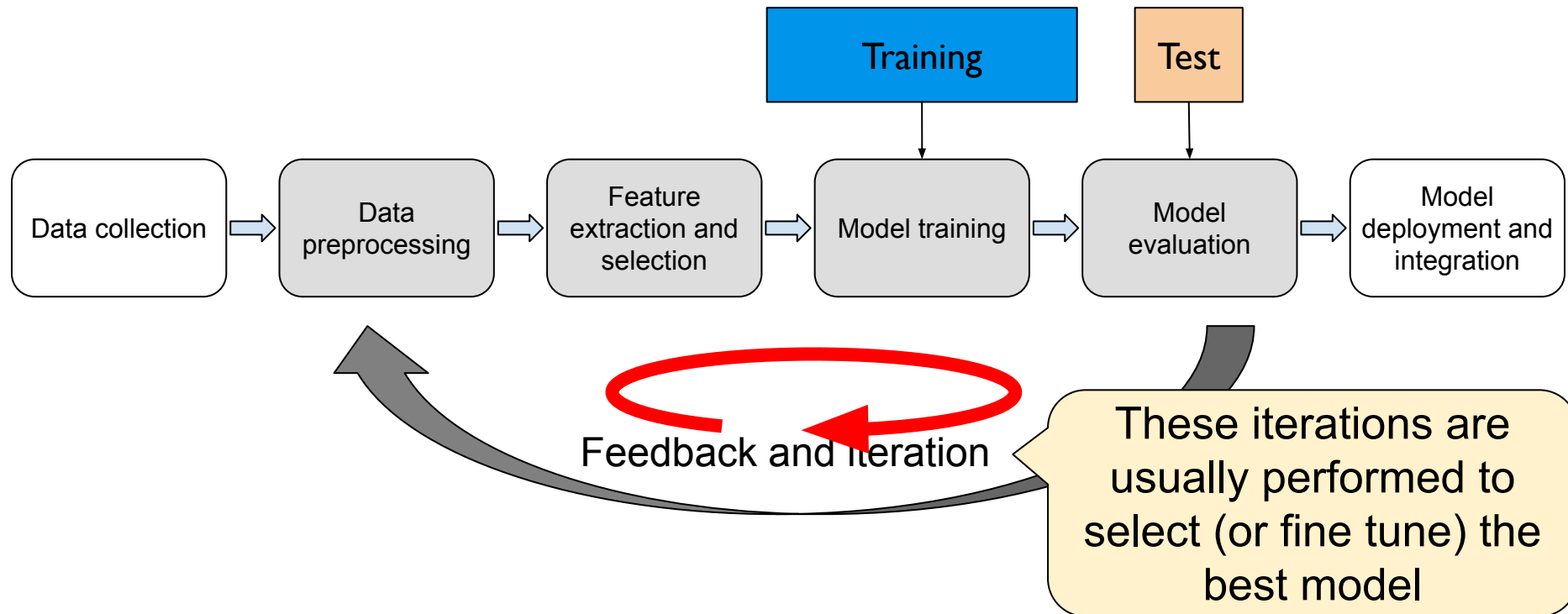## Test set and ML process iterations.

## Previously…

## Test set and ML process iterations.

## Previously…



These iterations are usually performed to select (or fine tune) the best model
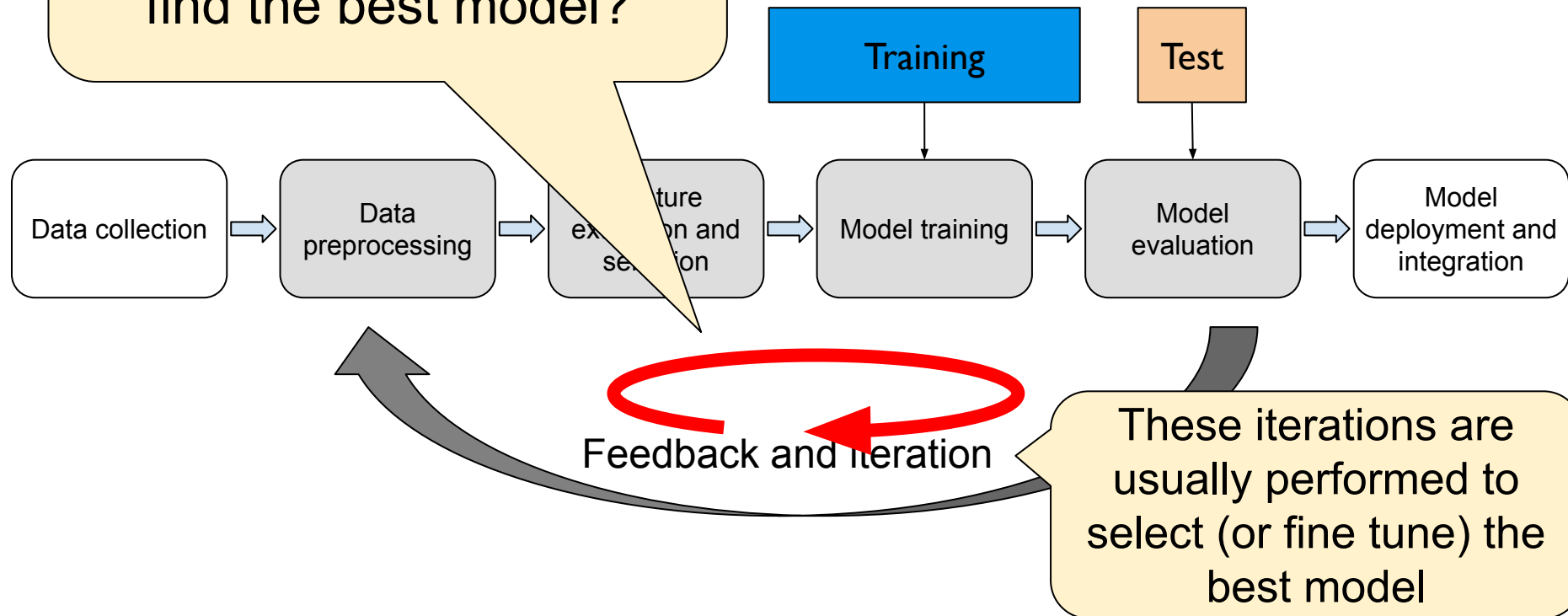
## Test set and ML process iterations.

What happens if we iterate several times to find the best model?

Training

Test

Data collection → Data preprocessing → Feature extraction and selection → Model training → Model evaluation → Model deployment and integration

Feedback and iteration

These iterations are usually performed to select (or fine tune) the best model
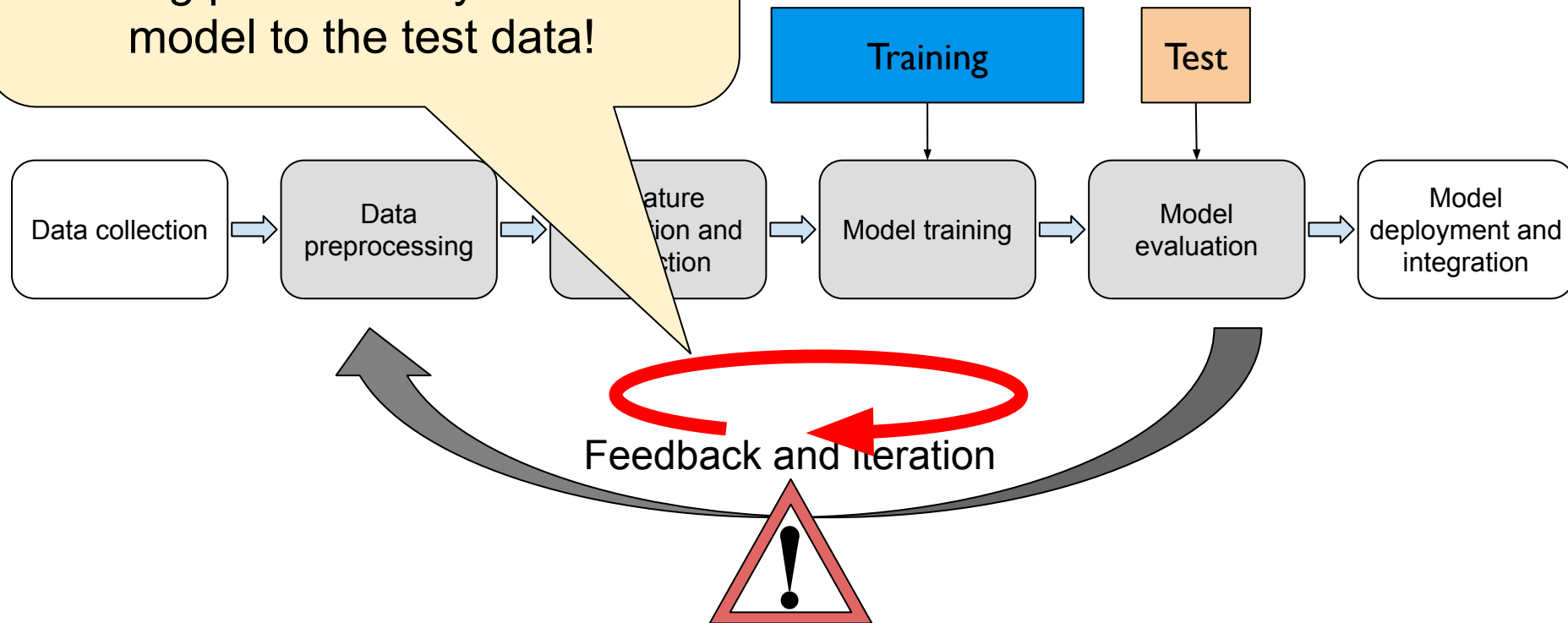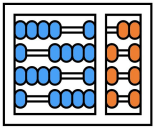
# ML Process - Dataset splitting

## Test set and ML process iterations.

Using the test results to repeatedly change/improve your learning process may overfit the model to the test data!
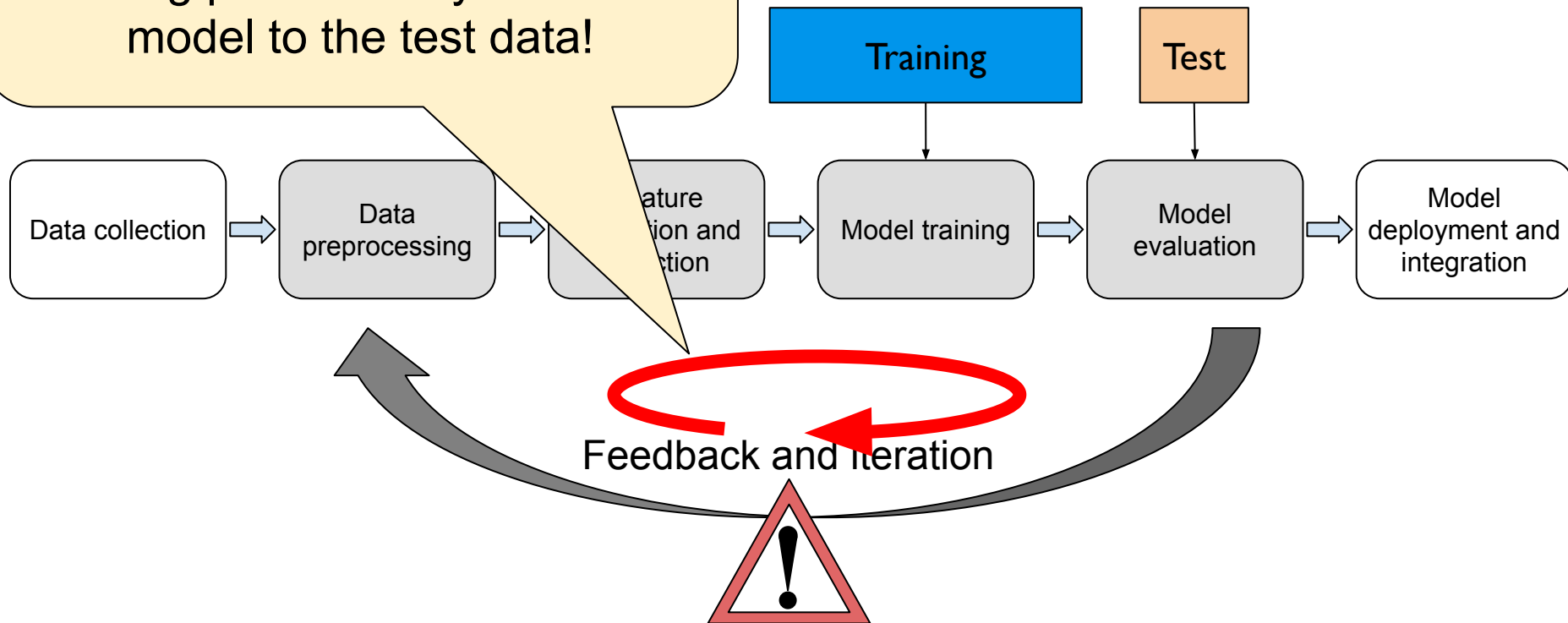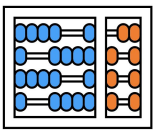
Training

Test

Data collection → Data preprocessing → ...ature ...ion and ...ction → Model training → Model evaluation → Model deployment and integration

Feedback and iteration

# ML Process - Dataset splitting

## Test set and ML process iteration

Test set wears out and cannot be used to evaluate generalization anymore

Using the test results to repeatedly change/improve your learning process may overfit the model to the test data!

| Training | Test |

Data collection → Data preprocessing → Feature selection and extraction → Model training → Model evaluation → Model deployment and integration

Feedback and iteration

# ML Process - Dataset splitting

Cross-validation: use different portions of the training set to train and to evaluate the model



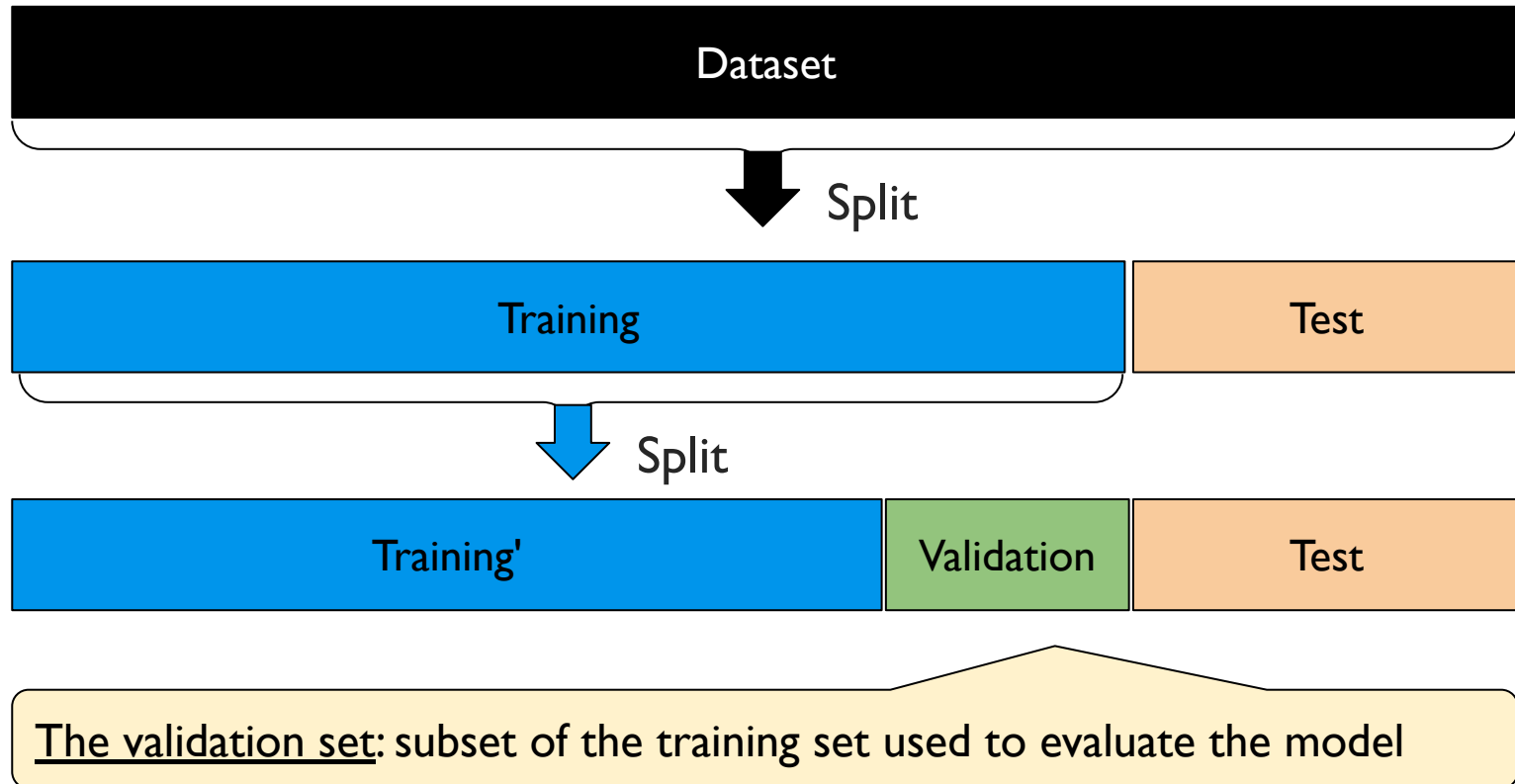The validation set: subset of the training set used to evaluate the model

# ML Process - Dataset splitting

<u>Cross-validation</u>: use different portions of the training set to train and to evaluate the model
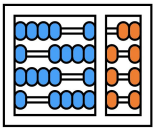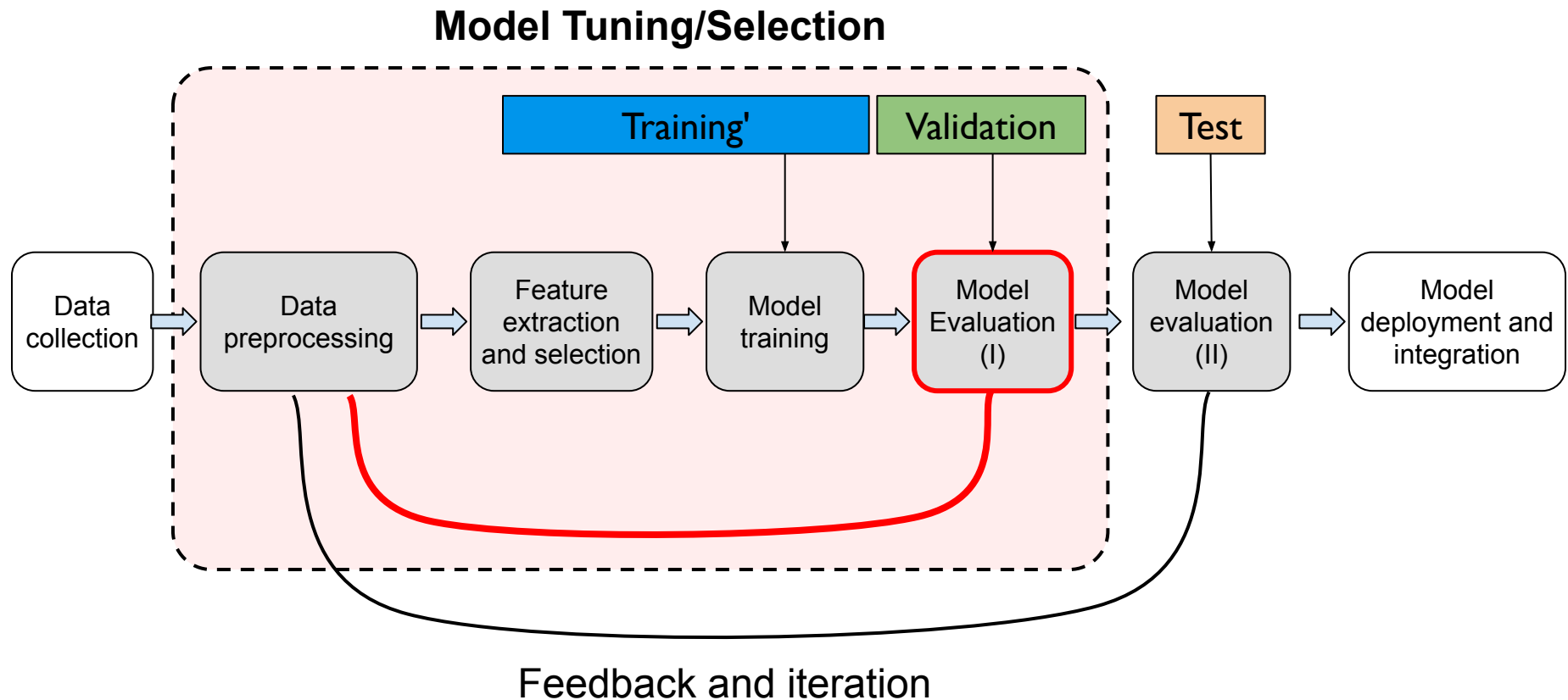
**Model Tuning/Selection**



Feedback and iteration
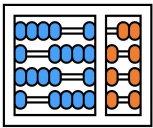
# ML Process - Dataset splitting

## Cross-validation: use different portions of the training set to train and to evaluate the model

**Model Tuning/Selection**



Pick model using validation set results and "double-check" it using the test set.

Cross-validation: use different p[...]
set to train and to evaluate the[...]

> Unfrequently used!
> Ideally, only once!

**Model Tuning/Selection**



> Pick model using validation set results and "double-check" it using the test set.

<u>Cross-validation</u>: use different portions of the training set to train and to evaluate the model

**Model Tuning/Selection**



It is usually a good idea to evaluate the model using the <u>average results of multiple validation sets</u>

Feedback and iteration

<u>Cross-validation</u>: use different portions of the training set to train and to evaluate the model

Several approaches:

- **Holdout method**
- **Leave-one-out cross-validation**
- **k-fold cross-validation**
- Leave-p-out cross-validation
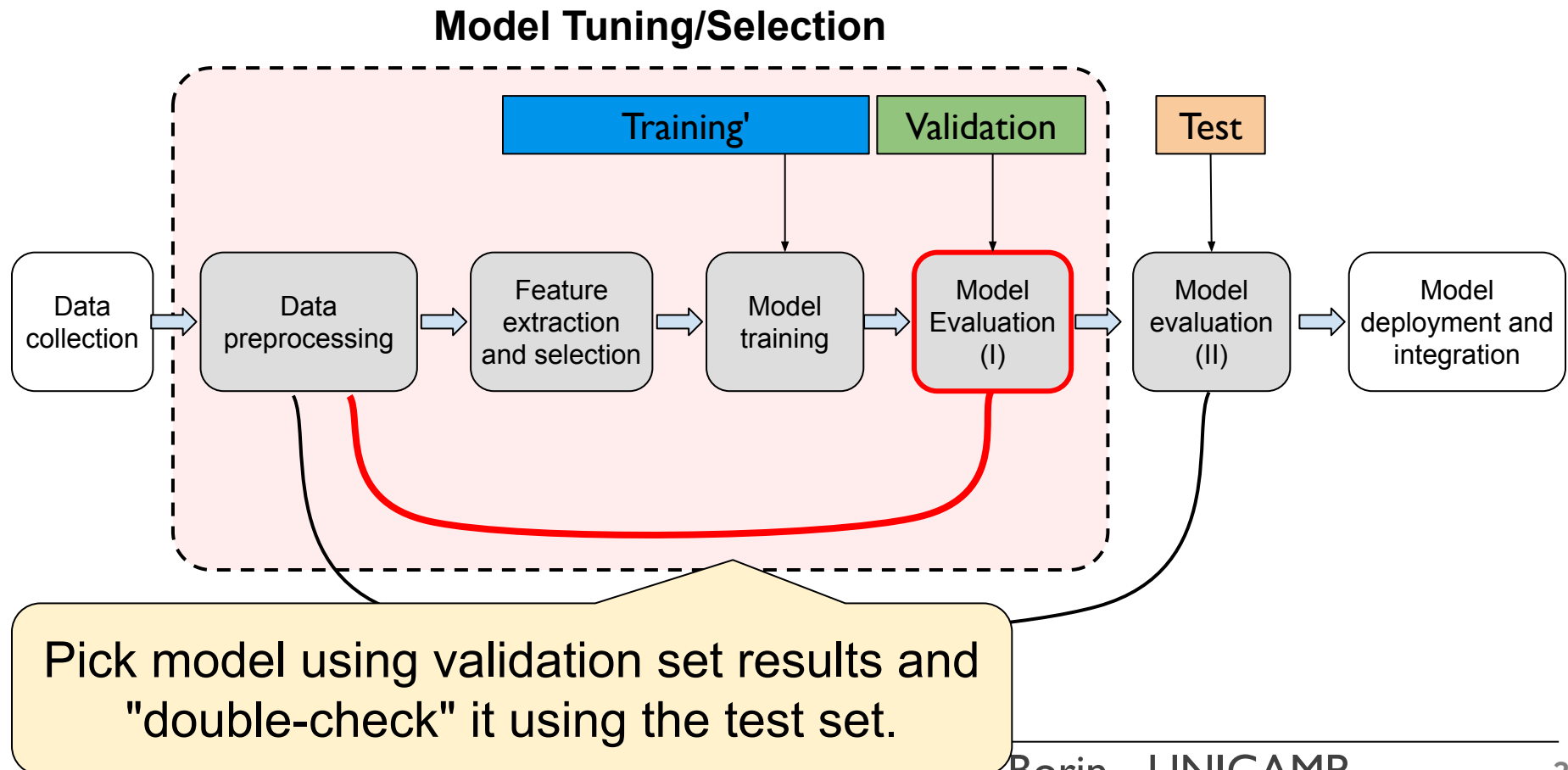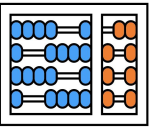- repeated random sub-sampling validation
- k*l-fold cross validation
- …

# ML Process - Dataset splitting

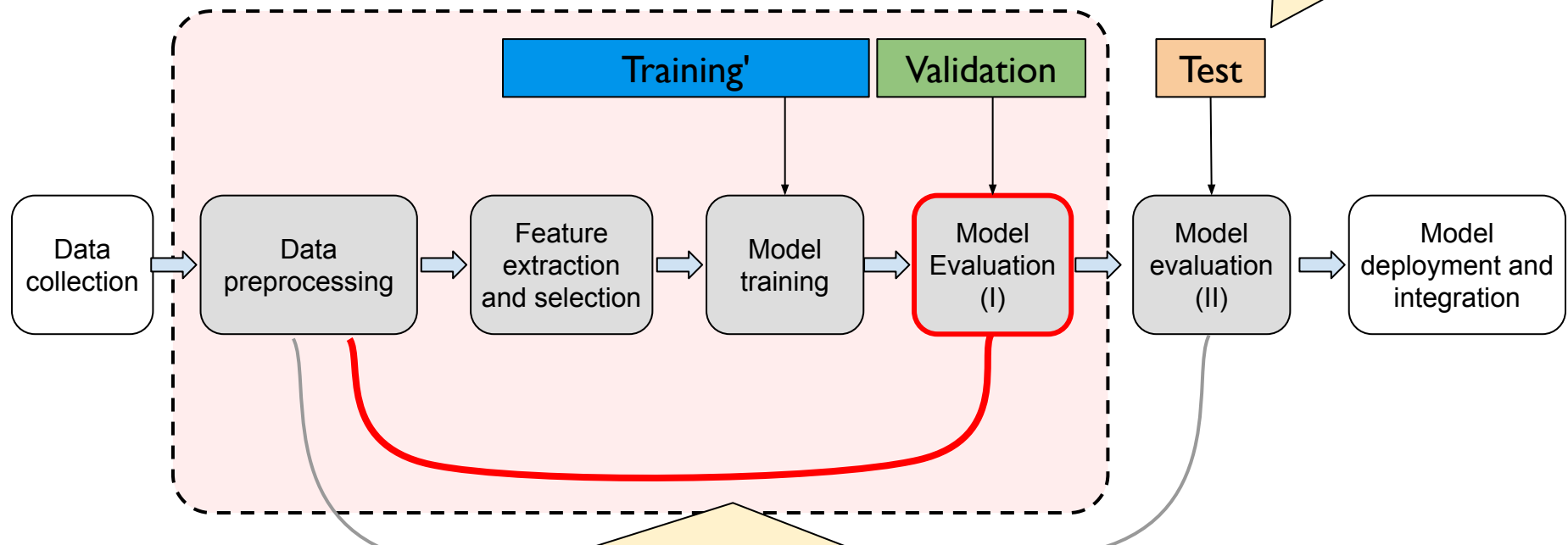<u>Cross-validation</u>: use different portions of the training set to train and to evaluate the model
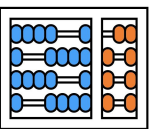
Several approaches:

- **Holdout method**: single train/validation partition randomly selected

Cross-validation: use different portions of the training set to train and to evaluate the model. Several approaches:

- **Holdout method**: single train/validation partition randomly selected

Single partition may cause evaluation bias.

Training

Split

Training'  Validation

Cross-validation: use different portions of the training set to train and to evaluate the model

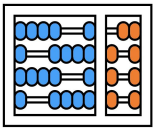Several approaches:

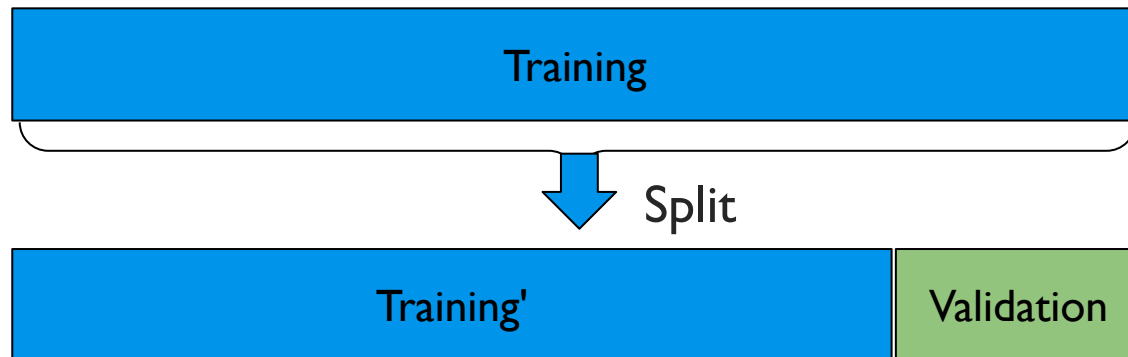- **Leave-one-out cross-validation**: 1 partition per item

# ML Process - Dataset splitting

Cross-validation: use different portions of the training set to train and to evaluate the model

Several approaches:

- **Leave-one-out cross-validation**: 1 part



> Each partition separates one item for validation and the rest for training.

Original training set (shuffled)

| 7 | 1 | 9 | 8 | 2 | 1 | … | 6 | 1 |

Partition 1: | 7 | 1 | 9 | 8 | 2 | 1 | … | 6 | 1 |
Partition 2: | 7 | 1 | 9 | 8 | 2 | 1 | … | 6 | 1 |
Partition 3: | 7 | 1 | 9 | 8 | 2 | 1 | … | 6 | 1 |
…
Partition N: | 7 | 1 | 9 | 8 | 2 | 1 | … | 6 | 1 |

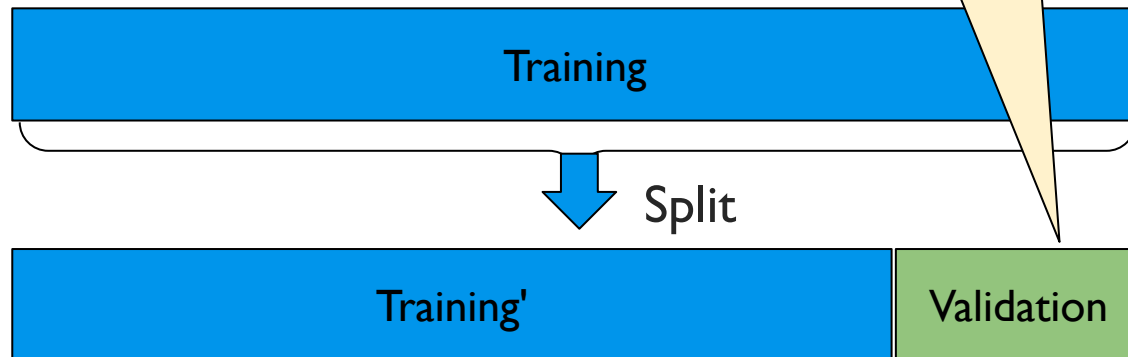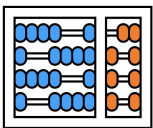**N train/validation partitions**
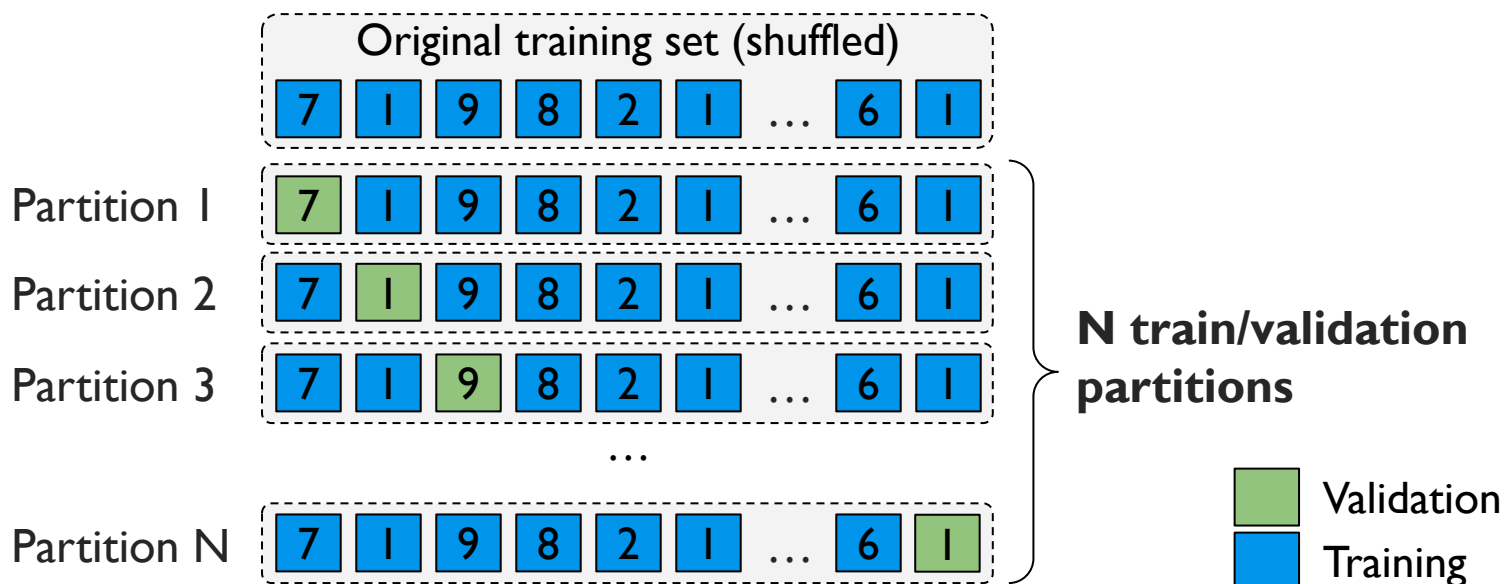
Validation
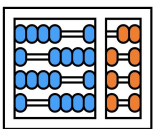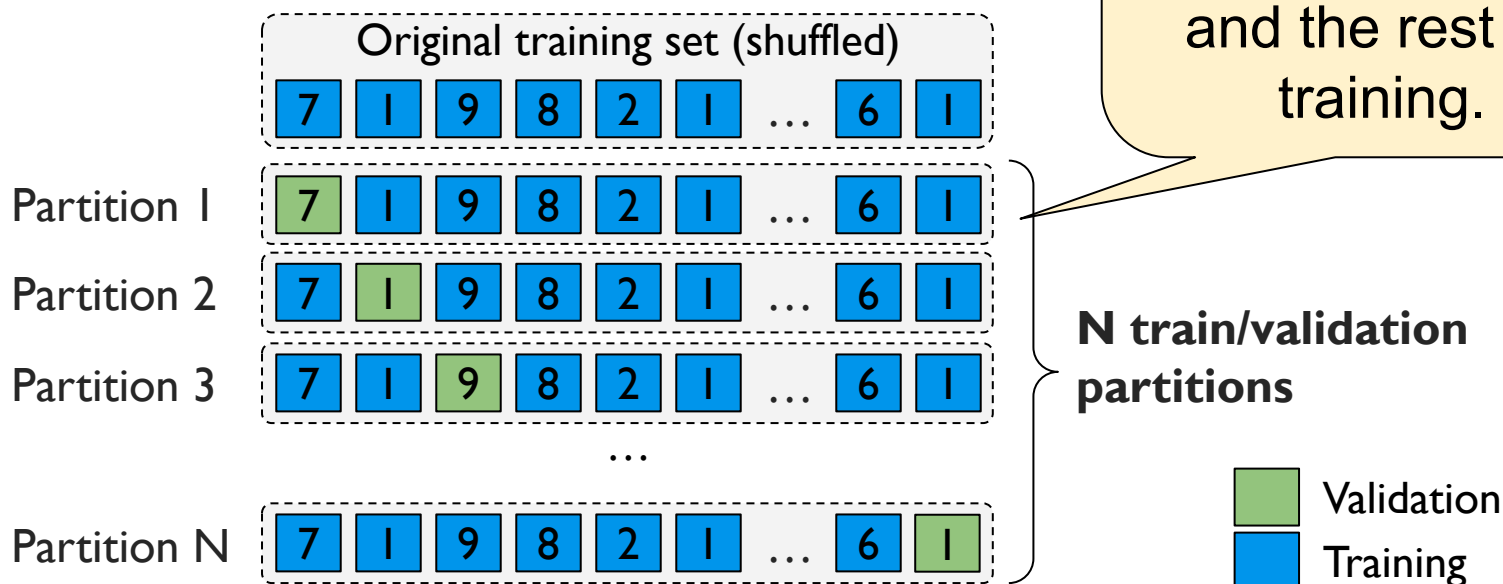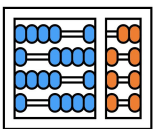Training

# ML Process - Dataset splitting

Cross-validation: use different portions of the training set to train and to evaluate the model

Several approaches:

- **Leave-one-out cross-validation**: 1 partition per item

Cross-validation: use different portions _____ g set to train and to evaluate the model. Several approaches:

- **Leave-one-out cross-validation**: 1 partition per item

> Report average and stdev

> Train and evaluate N times

**Original training set (shuffled)**

| 7 | 1 | 9 | 8 | 2 | 1 | … | 6 | 1 |

Partition 1

| 7 | 1 | 9 | 8 | 2 | 1 | … | 6 | 1 |

Partition 2

| 7 | 1 | 9 | 8 | 2 | 1 | … | 6 | 1 |

Partition 3

| 7 | 1 | 9 | 8 | 2 | 1 | … | 6 | 1 |

…

Partition N

| 7 | 1 | 9 | 8 | 2 | 1 | … | 6 | 1 |

**N train/validation partitions**

- Validation
- Training
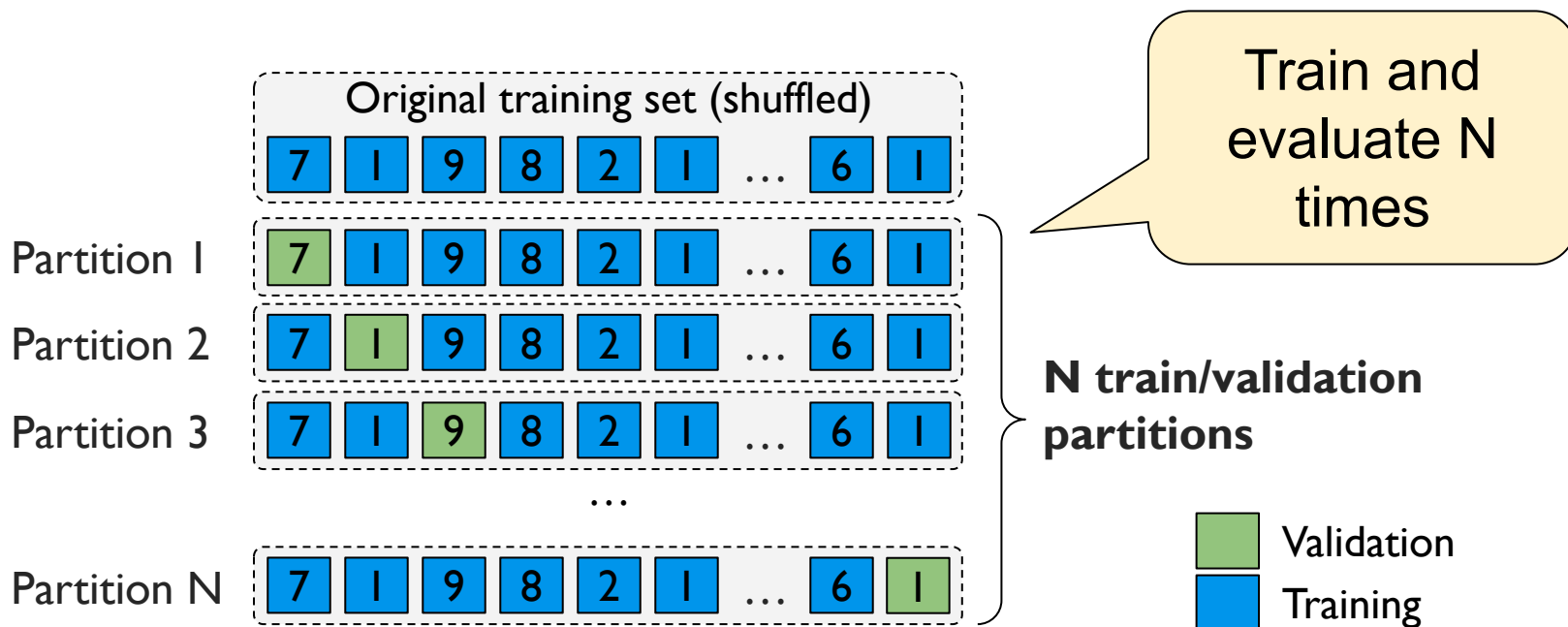
![](abacus logo)

# ML Process - Dataset splitting
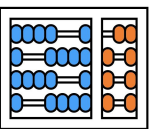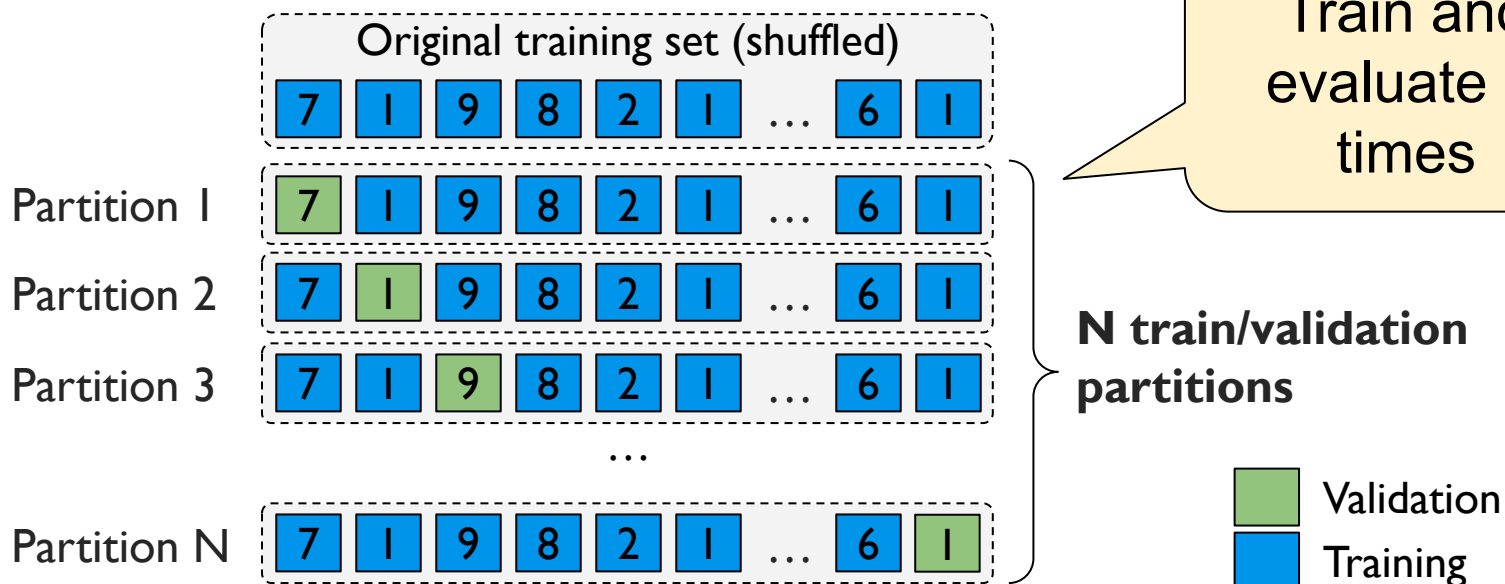
Cross-validation: use different portions ~~of~~ g
set to train and to evaluate the model
Several approaches:

- **Leave-one-out cross-validation**: 1 partition per item

Report average and stdev

Number of partitions grows with the number of items on the training set

⚠️

Train and evaluate N times

Training set (shuffled)

| 8 | 2 | 1 | … | 6 | 1 |

| 8 | 2 | 1 | … | 6 | 1 |

| 8 | 2 | 1 | … | 6 | 1 |

| 8 | 2 | 1 | … | 6 | 1 |

…

Partition N | 7 | 1 | 9 | 8 | 2 | 1 | … | 6 | 1 |

**N train/validation partitions**

🟩 Validation
🟦 Training

<u>Cross-validation</u>: use different portions of the training set to train and to evaluate the model

Several approaches:

- **k-fold cross-validation**: split the data in K folds and generate 1 partition per fold
- Example: 3-fold cross-validation



Original training set (shuffled)

| 7 | 1 | 9 | 8 | 2 | 1 | 6 | 1 | 4 |

Fold 1    Fold 2    Fold 3

# ML Process - Dataset splitting

Cross-validation: use different portions of the training set to train and to evaluate the model
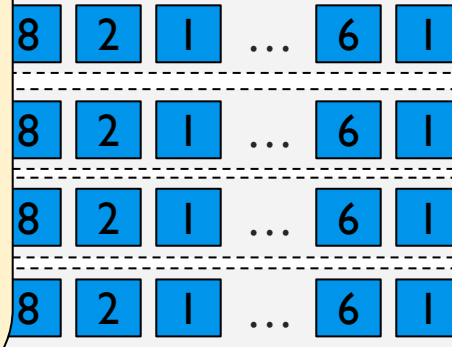
Several approaches:

- **k-fold cross-validation**: split the data in K folds and generate 1 partition per fold
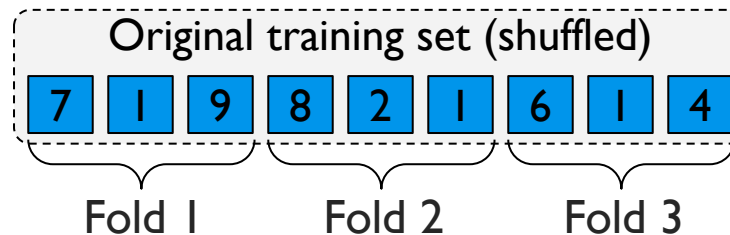- Example: 3-fold cross-validation

# ML Process - Dataset splitting

<u>Cross-validation</u>: use different portions of the training set to train and to evaluate the model

Several approaches:

- **k-fold cross-validation**: split the data in K folds and generate 1 partition per fold
- Example: 3-fold cross-validation

Original training set (shuffled)

| 7 | 1 | 9 | 8 | 2 | 1 | 6 | 1 | 4 |
|---|---|---|---|---|---|---|---|---|

Fold 1        Fold 2        Fold 3

Each partition separates one fold for validation and the rest for training.

Partition 1   | 7 | 1 | 9 | 8 | 2 | 1 | 6 | 1 | 4 |

Partition 2   | 7 | 1 | 9 | 8 | 2 | 1 | 6 | 1 | 4 |

Partition 3   | 7 | 1 | 9 | 8 | 2 | 1 | 6 | 1 | 4 |

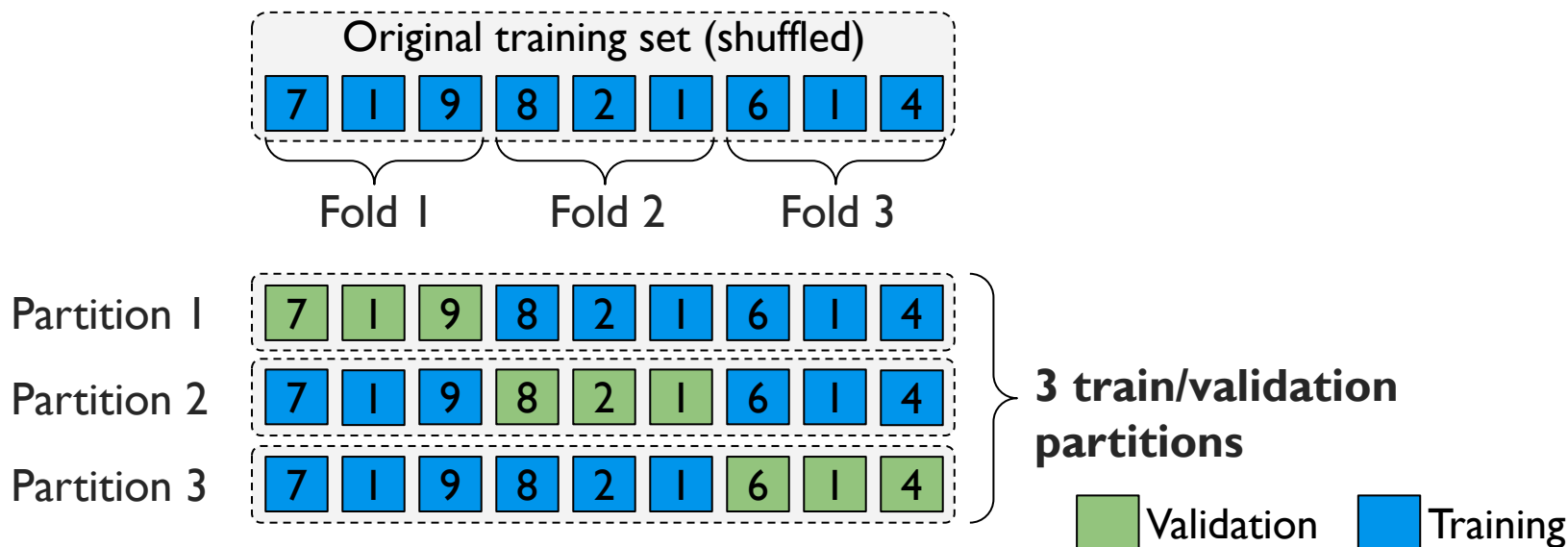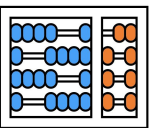**3 train/validation partitions**
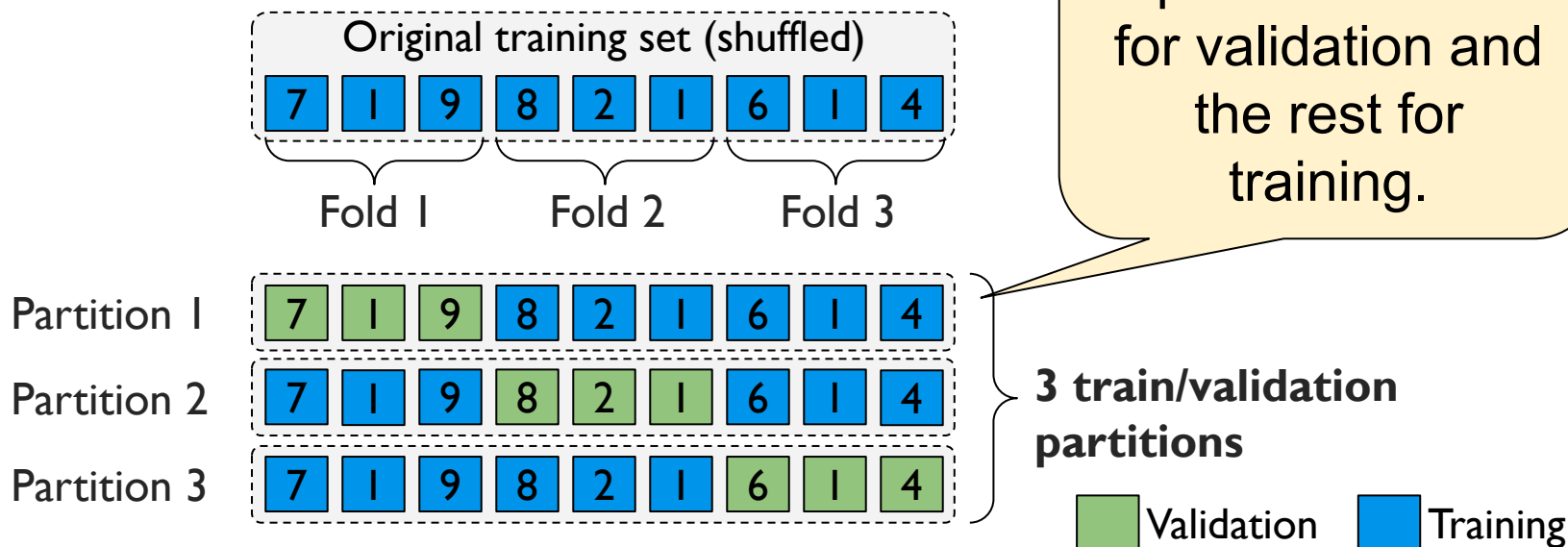
Validation     Training
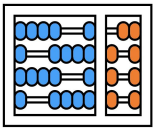
# ML Process - Dataset splitting

Cross-validation: use different portions of the training set to train and to evaluate the model

Several approaches:

- **k-fold cross-validation**: split the data in K folds and generate 1 partition per fold
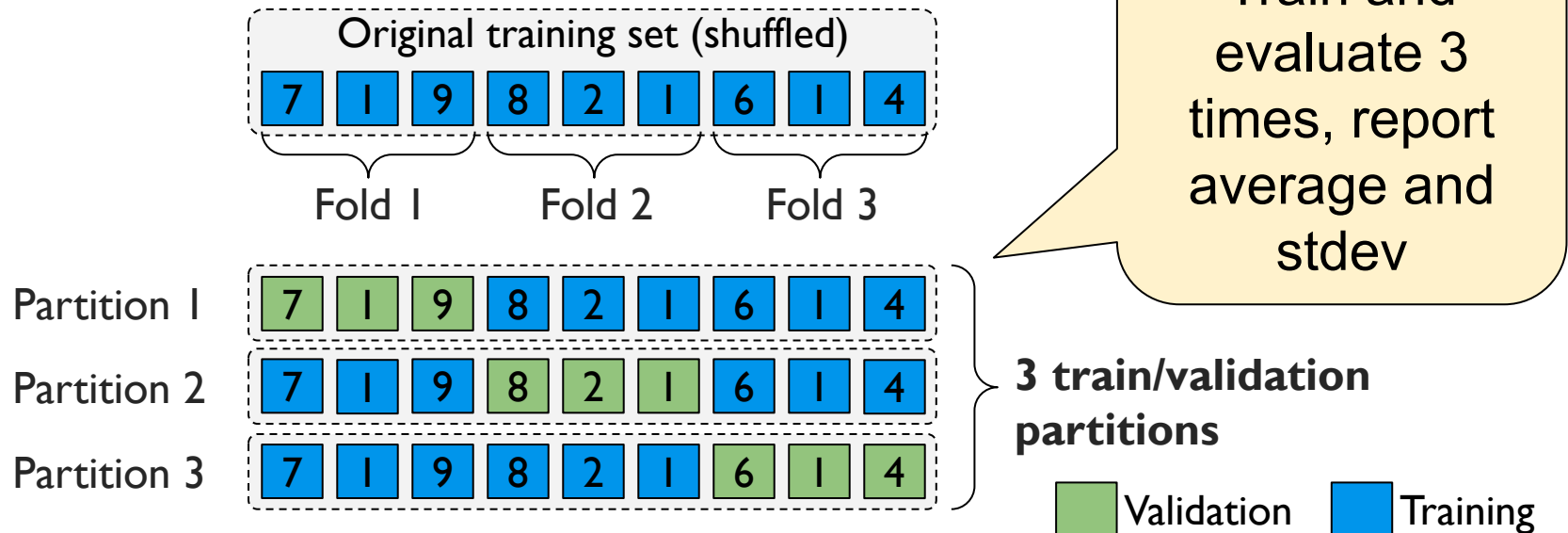- Example: 3-fold cross-validation

Original training set (shuffled)

| 7 | 1 | 9 | 8 | 2 | 1 | 6 | 1 | 4 |

Fold 1    Fold 2    Fold 3

Partition 1    | 7 | 1 | 9 | 8 | 2 | 1 | 6 | 1 | 4 |

Partition 2    | 7 | 1 | 9 | 8 | 2 | 1 | 6 | 1 | 4 |

Partition 3    | 7 | 1 | 9 | 8 | 2 | 1 | 6 | 1 | 4 |

Train and evaluate 3 times, report average and stdev

**3 train/validation partitions**
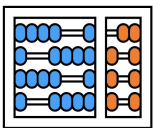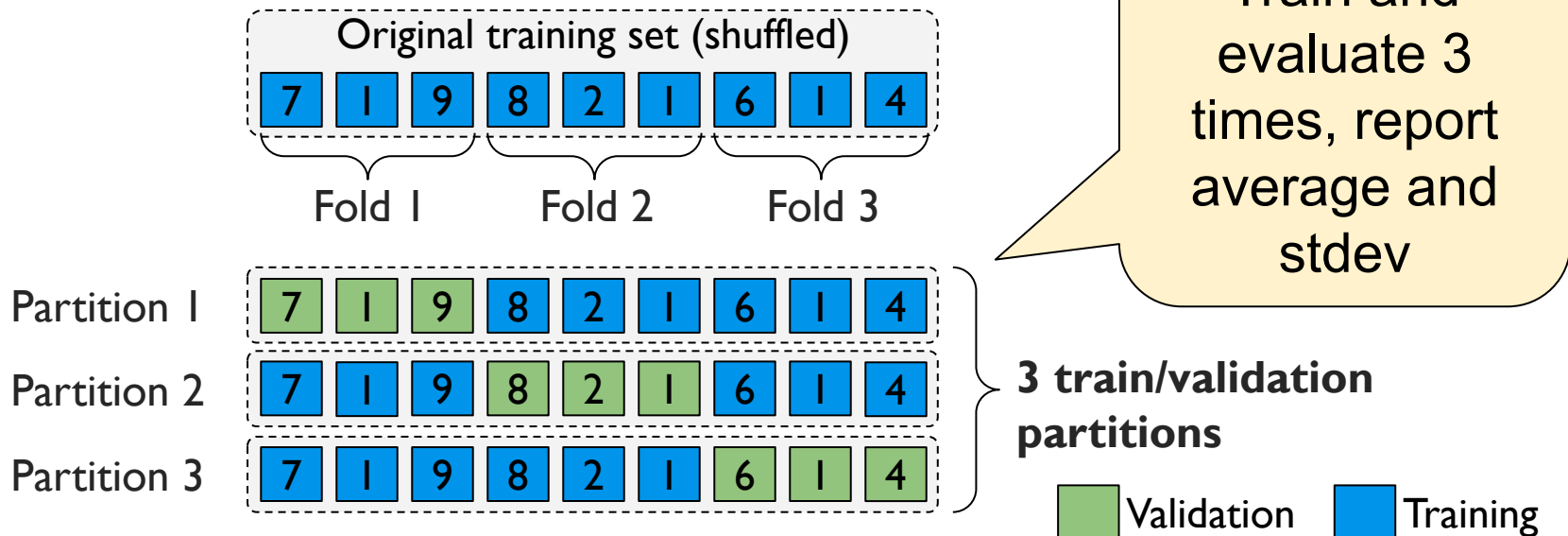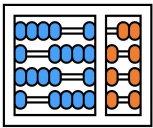
Validation    Training

# ML Process - Dataset splitting

Cross-validation: use different portions of the training set to train and to evaluate the

Several approaches:

- **k-fold cross-validation**: split the ~~~~~~~~ nerate 1 partition per fold
- Example: 3-fold cross-validation

Number of partitions and training/validation operations = K

Train and evaluate 3 times, report average and stdev

Original training set (shuffled)

| 7 | 1 | 9 | 8 | 2 | 1 | 6 | 1 | 4 |

Fold 1       Fold 2       Fold 3

Partition 1   | 7 | 1 | 9 | 8 | 2 | 1 | 6 | 1 | 4 |

Partition 2   | 7 | 1 | 9 | 8 | 2 | 1 | 6 | 1 | 4 |

Partition 3   | 7 | 1 | 9 | 8 | 2 | 1 | 6 | 1 | 4 |

**3 train/validation partitions**
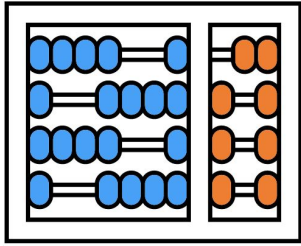
Validation    Training

# ML Process - Dataset splitting

## Key takeaways

- <u>Training Set</u>: part of the dataset used to train the model
- <u>Validation Set</u>: part of the dataset used to evaluate the model when searching for the best model or best set of hyperparameters
- <u>Test set</u>: part of the dataset set aside for final model evaluation. Ideally, should be used only once!
- <u>Cross-validation</u>: resampling method that uses different portions of the training set to train and evaluate models on different iterations
  - <u>k-fold cross-validation</u>: split the data in K folds and generate k partitions - each one using a different fold for validation and the remaining ones for training

# Capacitação profissional em tecnologias de Inteligência Artificial

## Machine Learning Overview

**Prof. Edson Borin**

https://www.ic.unicamp.br/~edson

Institute of Computing - UNICAMP