



CMPS 460 – Spring 2022

# MACHINE LEARNING

Tamer Elsayed

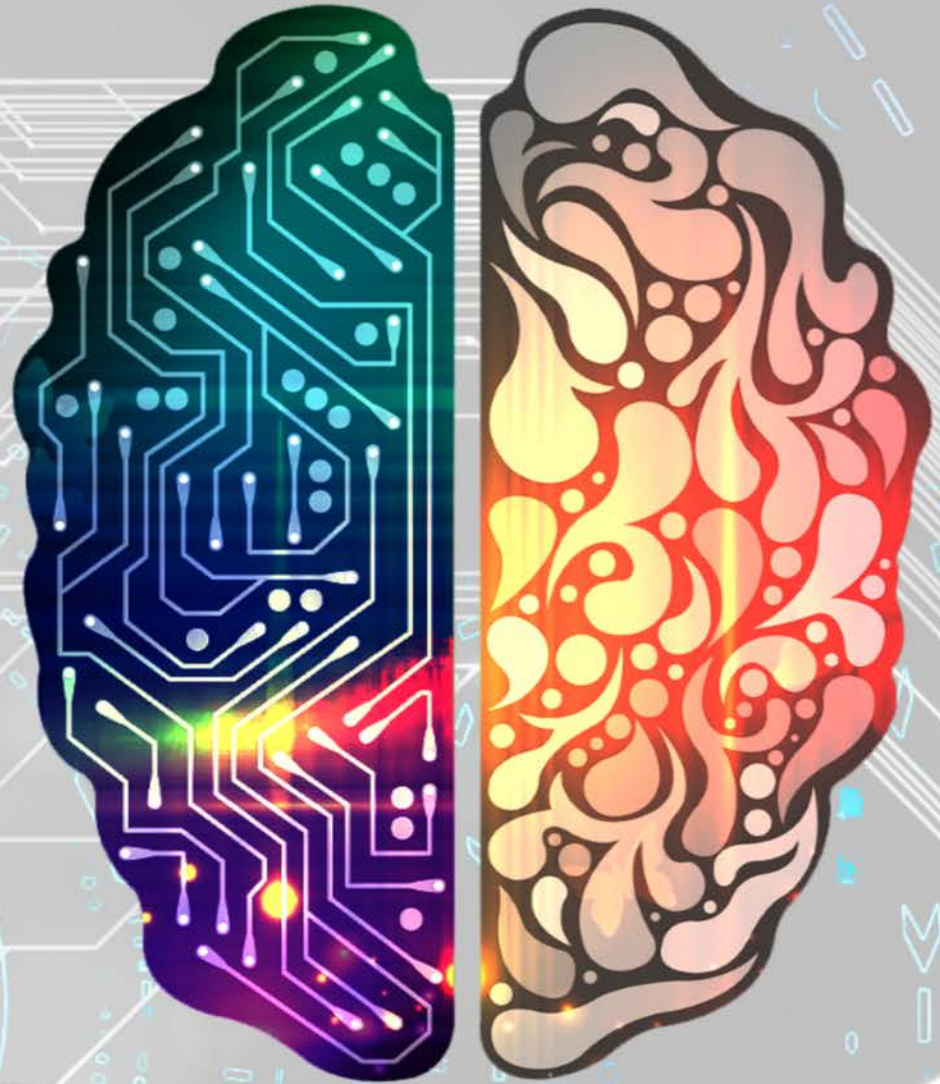


Image hosted by: WittySparks.com | Image source: Pixabay.com

5.c

## Practical Issues: Cross Validation & Debugging



Sec 5.6,  
5.8



# Cross-Validation

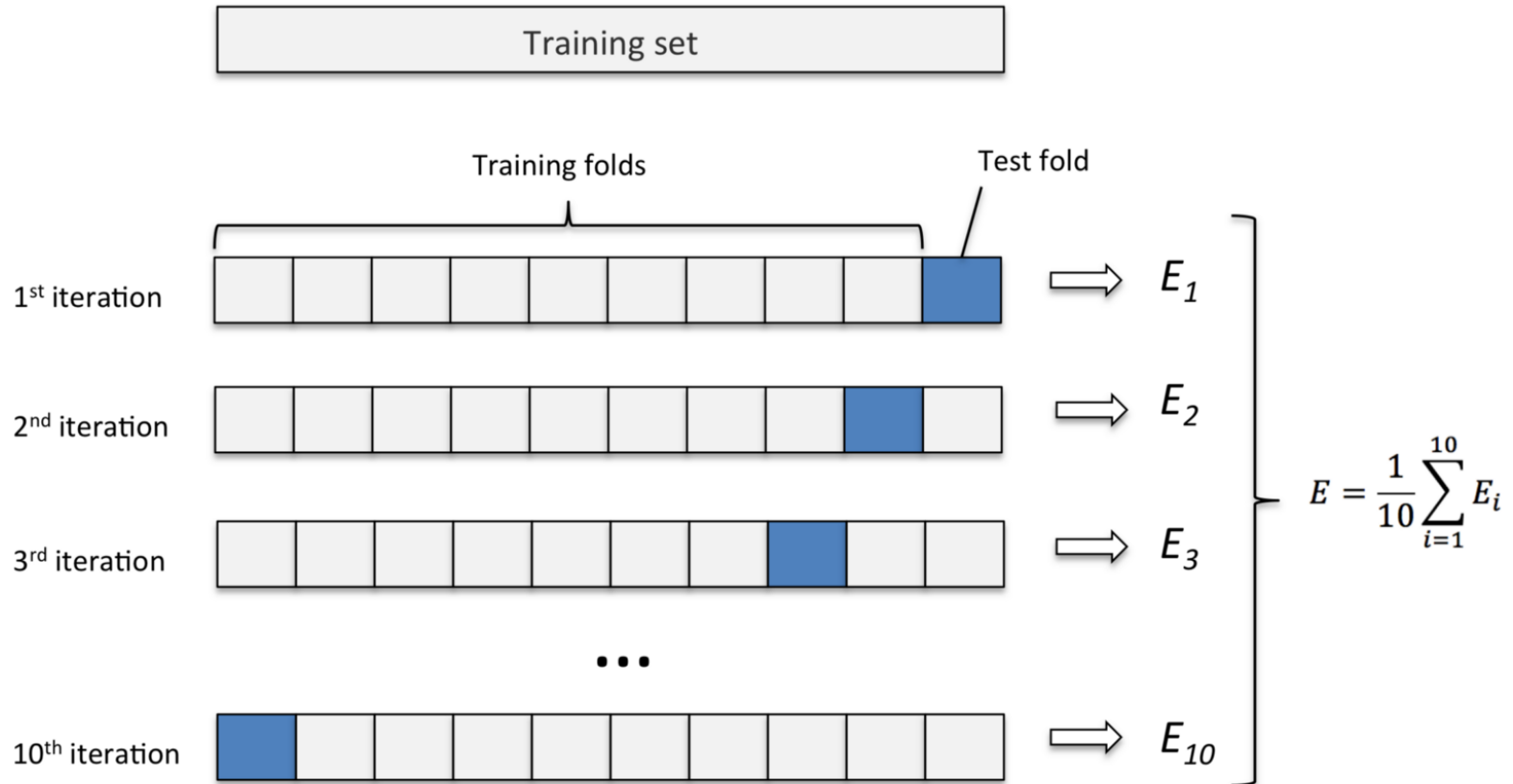
# Disadvantage of Train-Test split?

- Random split ==> Data might not be representative!
- Test only on part of the available data!

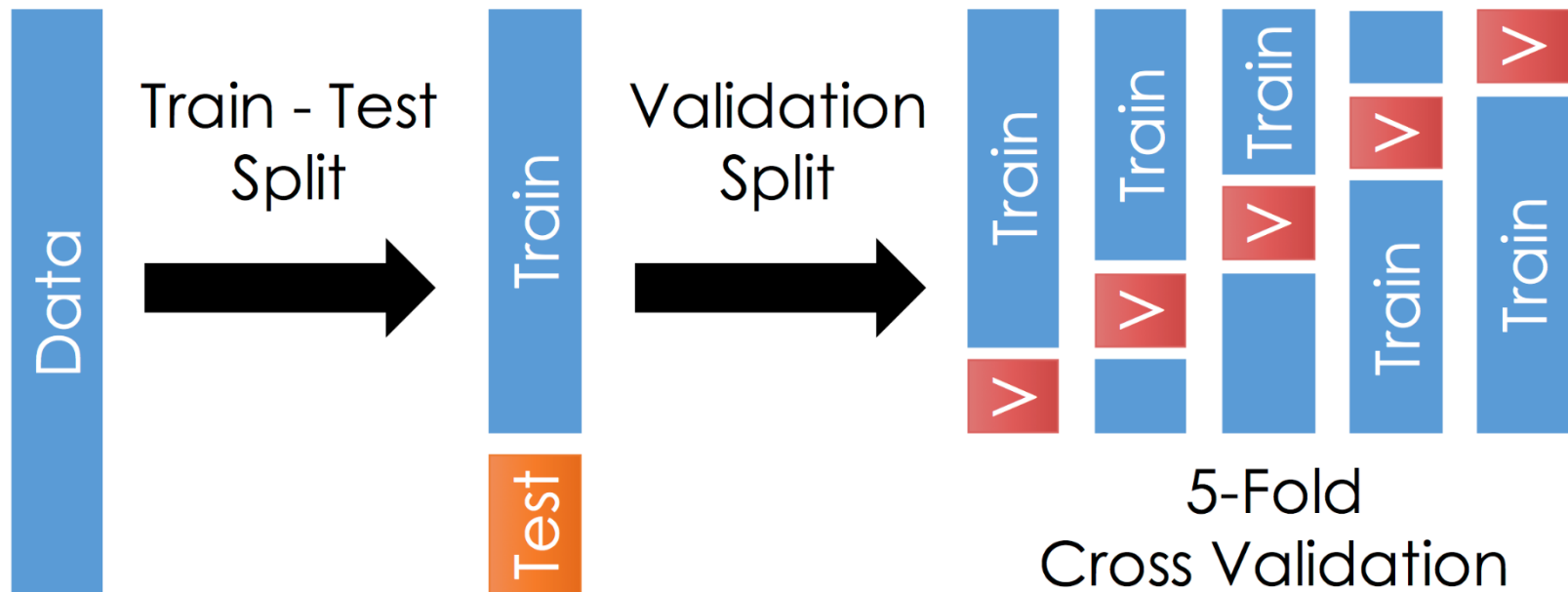
**Solution: Cross-Validation**

# K-fold Cross-Validation

K=10 here



# With Hyper-Parameters?





## Algorithm 8 **CROSSVALIDATE**(*LearningAlgorithm*, *Data*, *K*)

```
1:  $\hat{\epsilon} \leftarrow \infty$  // store lowest error encountered so far
2:  $\hat{\alpha} \leftarrow \text{unknown}$  // store the hyperparameter setting that yielded it
3: for all hyperparameter settings  $\alpha$  do
4:    $err \leftarrow [ ]$  // keep track of the  $K$ -many error estimates
5:   for  $k = 1$  to  $K$  do
6:      $train \leftarrow \{(x_n, y_n) \in Data : n \bmod K \neq k - 1\}$ 
7:      $test \leftarrow \{(x_n, y_n) \in Data : n \bmod K = k - 1\}$  // test every  $K$ th example
8:      $model \leftarrow \text{Run } LearningAlgorithm \text{ on } train$ 
9:      $err \leftarrow err \oplus \text{error of } model \text{ on } test$  // add current error to list of errors
10:  end for
11:   $avgErr \leftarrow \text{mean of set } err$ 
12:  if  $avgErr < \hat{\epsilon}$  then
13:     $\hat{\epsilon} \leftarrow avgErr$  // remember these settings
14:     $\hat{\alpha} \leftarrow \alpha$  // because they're the best so far
15:  end if
16: end for
```

# Leave-one-out (LOO) C.V.

- When  $K=N$

Why would we do that?

Any drawbacks?

- Very natural for kNN.





# Debugging Learning Algorithms



THIS IS YOUR MACHINE LEARNING SYSTEM?

YUP! YOU POUR THE DATA INTO THIS BIG PILE OF LINEAR ALGEBRA, THEN COLLECT THE ANSWERS ON THE OTHER SIDE.

WHAT IF THE ANSWERS ARE WRONG?

JUST STIR THE PILE UNTIL THEY START LOOKING RIGHT.



# Debugging!

- You've implemented a learning algorithm ..
- You try it on some train/dev/test data ..
- You get really bad performance ..

**What's going on?!**

- Is the **data** too noisy?
- Is the **learning problem** too hard?
- Is the **implementation** of the learning algorithm buggy?
- ...

# Strategies for Isolating Causes of Errors

- Is the problem with generalization to test data?
  - Can learner fit the training data?
    - Yes: problem is in generalization to test data
      - too complicated model family
      - not enough data
    - No: problem is in representation
      - need better features
      - better data
- Train/test mismatch?
  - Try reselecting train/test by shuffling training and test data together.

# Strategies for Isolating Causes of Errors

- **Is learning algorithm implementation correct?**
  - Measure loss rather than accuracy (make sure it is minimized)
  - Hand-craft a toy dataset
  - (if possible) Compare against reference implementation
- **Is representation adequate?**
  - If you can't fit training data, you might need better features.
  - Can you overfit if you add a cheating feature that perfectly correlates with correct class?
    - If not (near) 0% error: too many noisy features OR a bug!
    - If (near) 0% error: work on better features, or change learning model
- **Do you have enough data?**
  - Try training on less training set, how much does it hurt performance?