

CSC 561: Neural Networks and Deep Learning

Loss, Overfitting, Model Selection

Marco Alvarez

Department of Computer Science and Statistics
University of Rhode Island

Spring 2024



Loss functions

The “learning” problem

- Finding a hypothesis (classifier/regressor) that best approximates the target function

for $h \in \mathcal{H}$ and $\forall (x_i, y_i) \sim P$, we want $h(x_i) \approx f(x_i)$

ML uses **search** and **optimization** (to **minimize expected loss**)

$$\mathbb{E} \left[l(h, x_i, y_i) \right]_{(x_i, y_i) \sim P}$$

3

Approximating the expected loss

$$\mathbb{E} \left[l(h, x_i, y_i) \right]_{(x_i, y_i) \sim P} \text{ expected loss}$$

$$\approx L = \frac{1}{n} \sum_{i=1}^n l(h, x_i, y_i) \text{ empirical loss}$$

the **law of large numbers** states that the arithmetic mean of the values almost surely converges to the expected value as the number

4

0/1 loss

$$l_{0/1}(h, x_i, y_i) = I(h(x_i) \neq y_i)$$

indicator
function

Prediction	Target
5	5
1	9
2	2
7	7
8	0
0	0
0	8
3	3
6	6
4	4

Empirical loss?

5

Practice

X0	X1	X2	Y
1	0	0	-1
1	1	0	+1
1	1	1	+1
1	0	1	+1

$$h_w(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$$

$$\sigma(z) = \begin{cases} +1 & \text{if } z > 0 \\ -1 & \text{if } z \leq 0 \end{cases}$$

zero-one loss for $\mathbf{w}_a = [0,0,0]^T$?

zero-one loss for $\mathbf{w}_b = [0,1,0]^T$?

zero one loss for $\mathbf{w}_c = [-1,2,2]^T$?

6

Squared loss

$$l_{sq}(h, x_i, y_i) = (h(x_i) - y_i)^2$$

penalizes big
mistakes

Prediction	Target
1.2	1.4
2.3	2.3
1.1	1.2
3.4	4.1
2.3	2.5
1.1	1.1
2.5	2.6
3.1	3.2
1.7	1.8
2.3	2.3

Empirical loss?

7

Absolute loss

$$l_{abs}(h, x_i, y_i) = |h(x_i) - y_i|$$

Prediction	Target
1.2	1.4
2.3	2.3
1.1	1.2
3.4	4.1
2.3	2.5
1.1	1.1
2.5	2.6
3.1	3.2
1.7	1.8
2.3	2.3

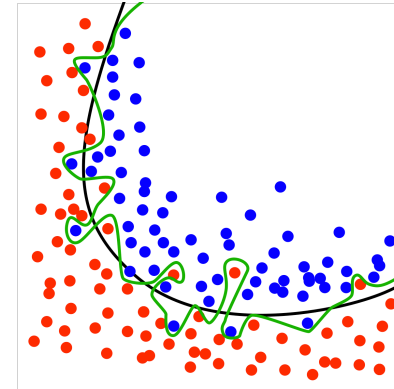
Empirical loss?

8

Overfitting

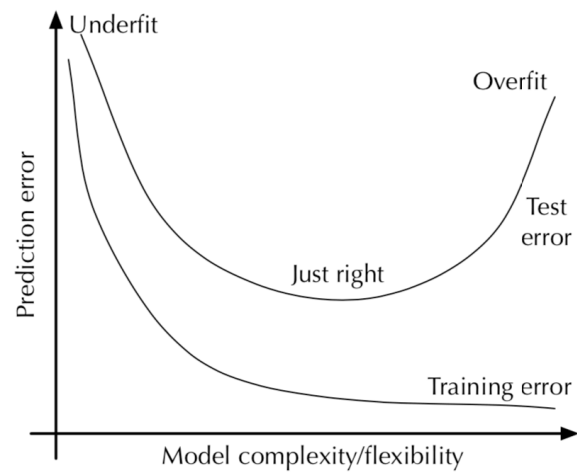
Overfitting

- Learning a model that “knows” the training data very well but does not generalize



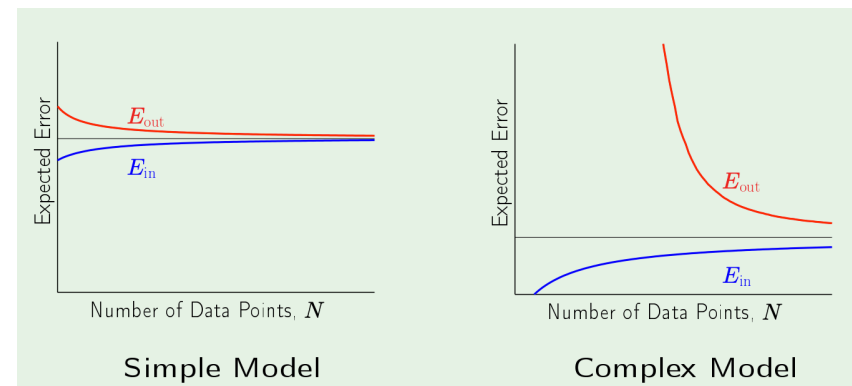
10

Model complexity



11

Number of data instances



<https://work.caltech.edu/lectures.html>

12

Overfitting

- Reasons
 - ✓ model is too complex
 - ✓ model is fitting noise present in the training data
 - ✓ training data is not a representative sample of the distribution
- How to prevent?
 - ✓ use more training data
 - ✓ use fewer features
 - ✓ regularize your model

13

Generalization

- We can use a ML method to calculate:

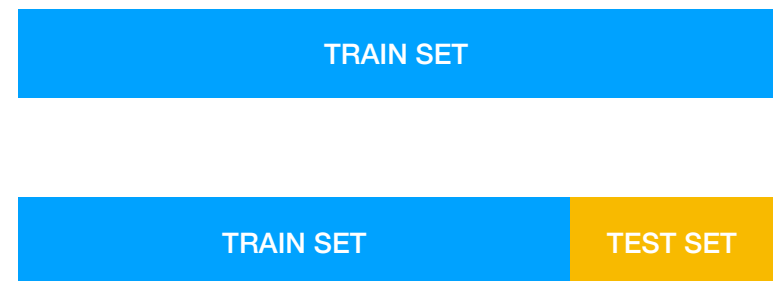
$$g = \arg \min_{h \in \mathcal{H}} L(h, \mathcal{D})$$

- Problem: it may overfit the training data
 - ✓ we want better generalization
- Solution: split your data in train, validation, test
 - ✓ use train and validation to select the best hypothesis
 - ✓ use test for final evaluation and report

14

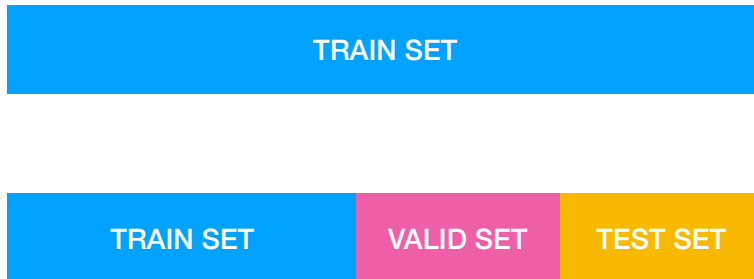
Model selection

Train and test

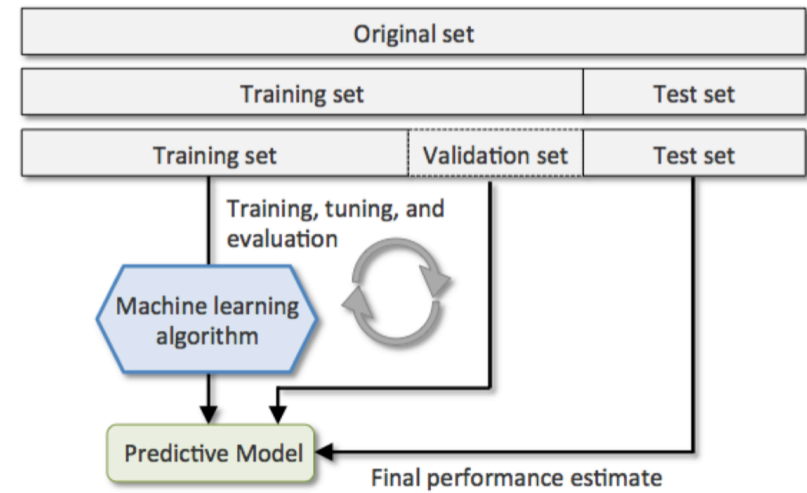


16

Train, validation, and test



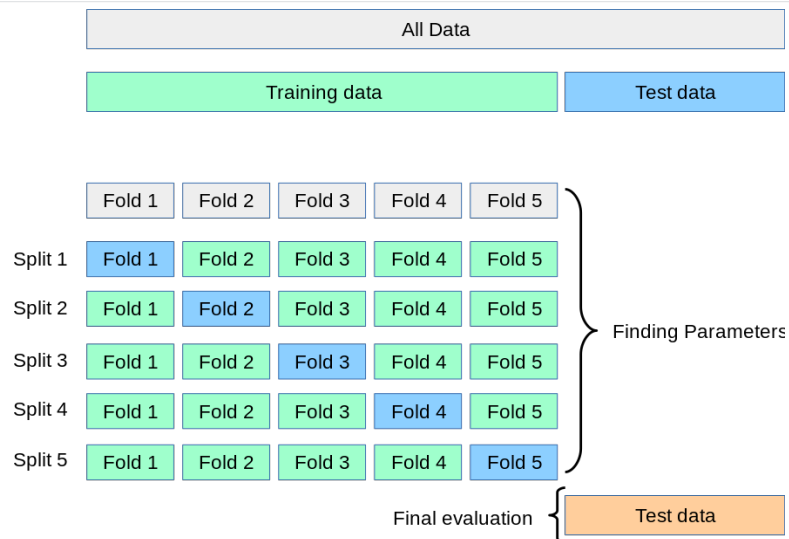
17



from INFO-4604: Applied Machine Learning, Fall 2017, Michael Paul, Univ. of Colorado

18

k-fold cross validation



https://scikit-learn.org/stable/modules/cross_validation.html

19

Stratified cross validation



Stratified cross validation aims at having the same class distribution within each fold

<https://towardsdatascience.com/cross-validation-explained-evaluating-estimator-performance-e51e5430ff85>

20

Evaluation

Confusion matrix (2 classes)

		Predicted condition	
		Positive (PP)	Negative (PN)
Actual condition	Positive (P)	True positive (TP)	False negative (FN)
	Negative (N)	False positive (FP)	True negative (TN)

https://en.wikipedia.org/wiki/Confusion_matrix

22

Confusion matrix (example)

		Predicted condition	
		Cancer	Non-cancer
Actual condition	Cancer	6	2
	Non-cancer	1	3

https://en.wikipedia.org/wiki/Confusion_matrix

23

Evaluation metrics (2 classes)

accuracy (ACC)

$$ACC = \frac{TP + TN}{P + N} = \frac{TP + TN}{TP + TN + FP + FN}$$

F1 score

is the harmonic mean of precision and sensitivity

$$F_1 = 2 \cdot \frac{PPV \cdot TPR}{PPV + TPR} = \frac{2TP}{2TP + FP + FN}$$

Matthews correlation coefficient (MCC)

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

https://en.wikipedia.org/wiki/Confusion_matrix

24

Confusion matrix (example >2 classes)

