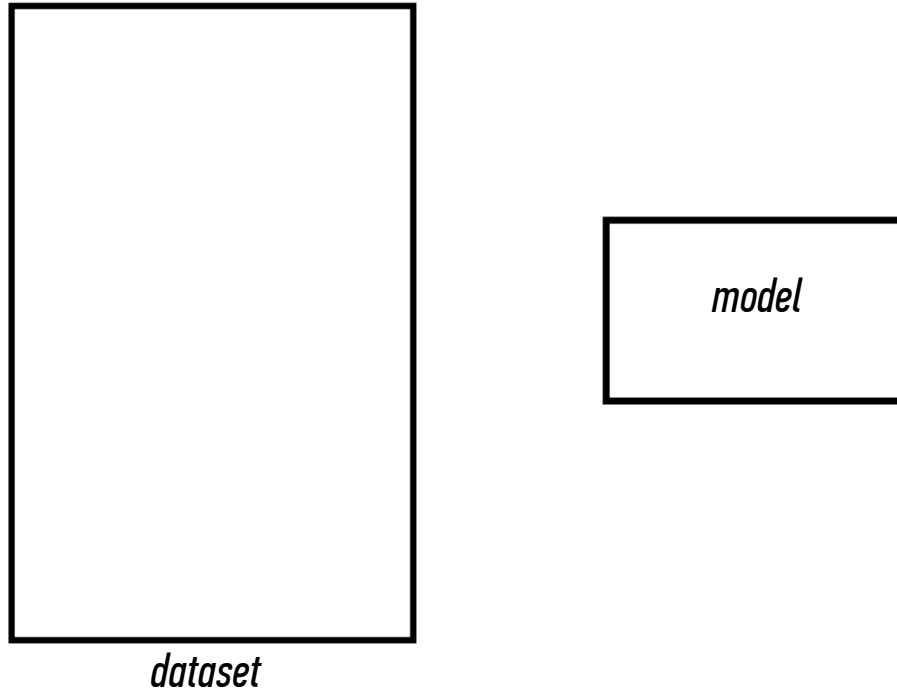


INTRO to DATA SCIENCE

CROSS-VALIDATION

CLASSIFICATION PROBLEMS

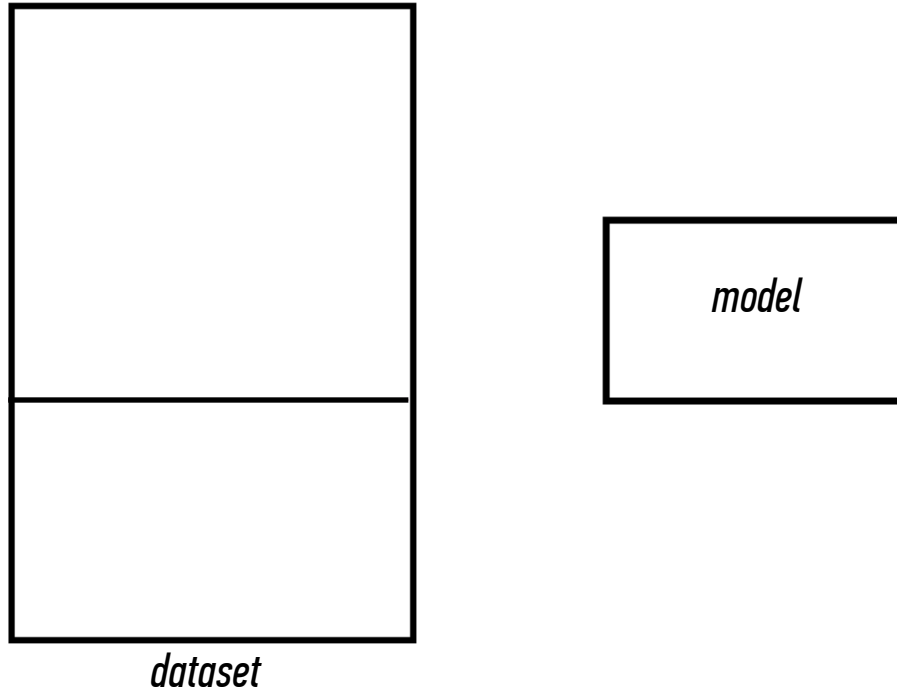
Q: What steps does a classification problem require?



CLASSIFICATION PROBLEMS

Q: What steps does a classification problem require?

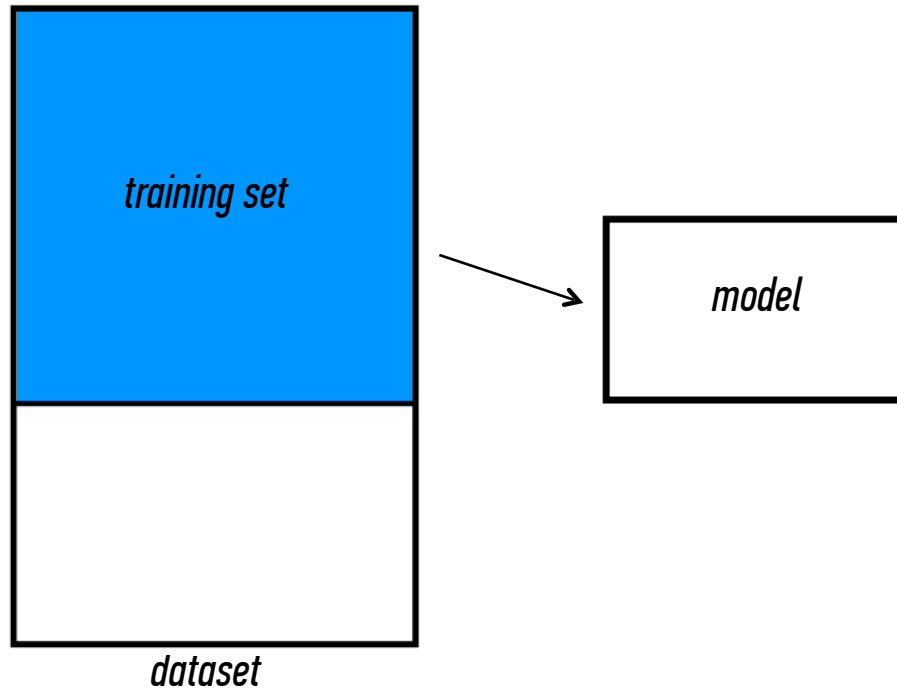
1) split dataset



CLASSIFICATION PROBLEMS

Q: What steps does a classification problem require?

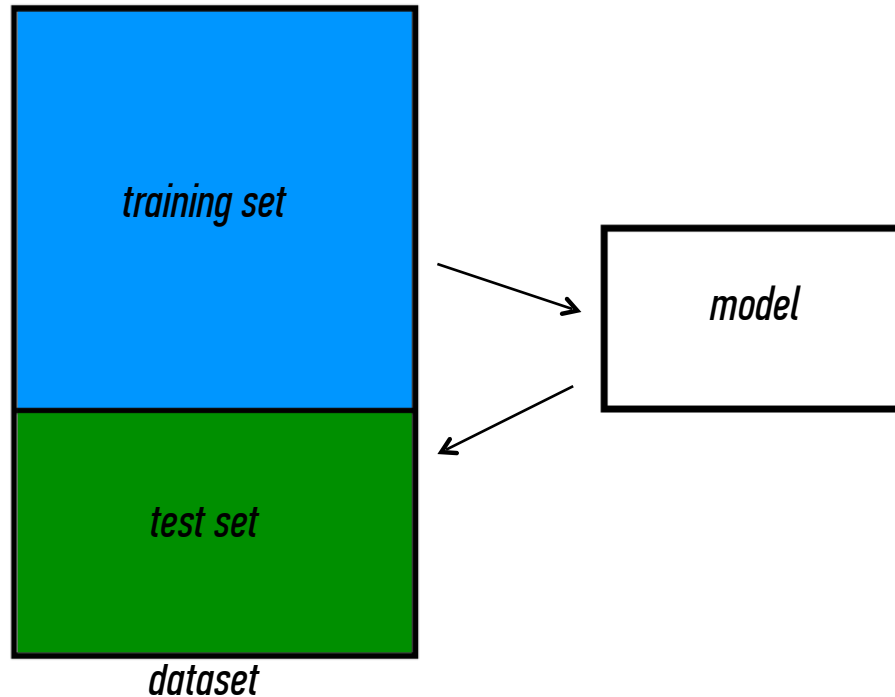
- 1) split dataset*
- 2) train model*



CLASSIFICATION PROBLEMS

Q: What steps does a classification problem require?

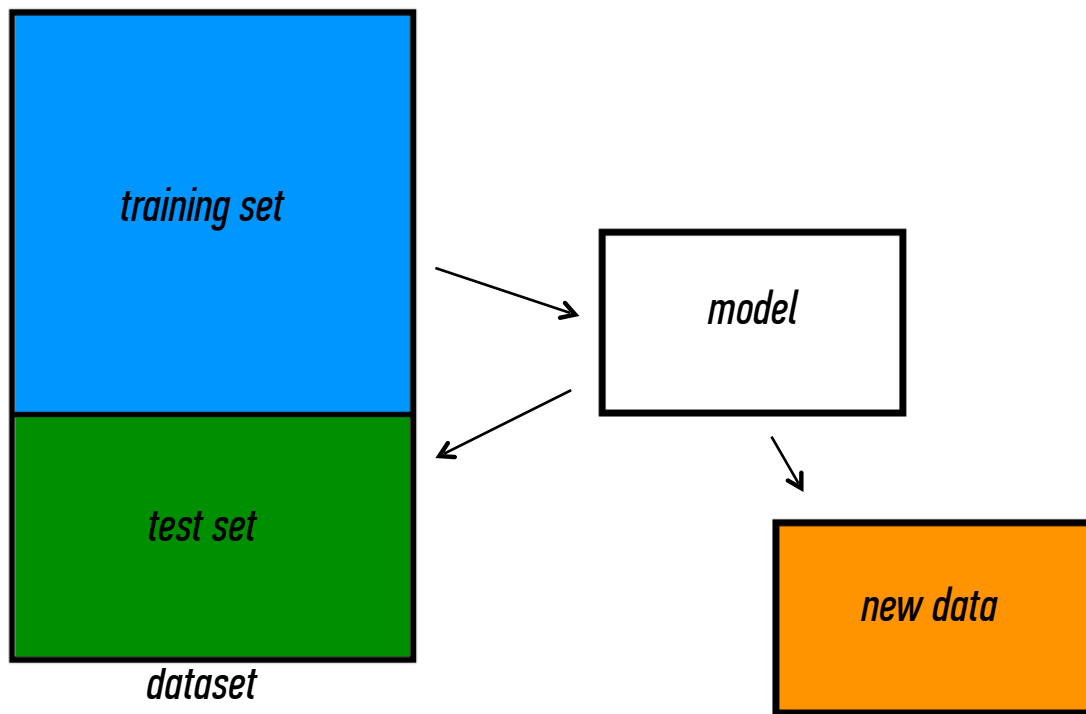
- 1) split dataset*
- 2) train model*
- 3) test model*



CLASSIFICATION PROBLEMS

Q: What steps does a classification problem require?

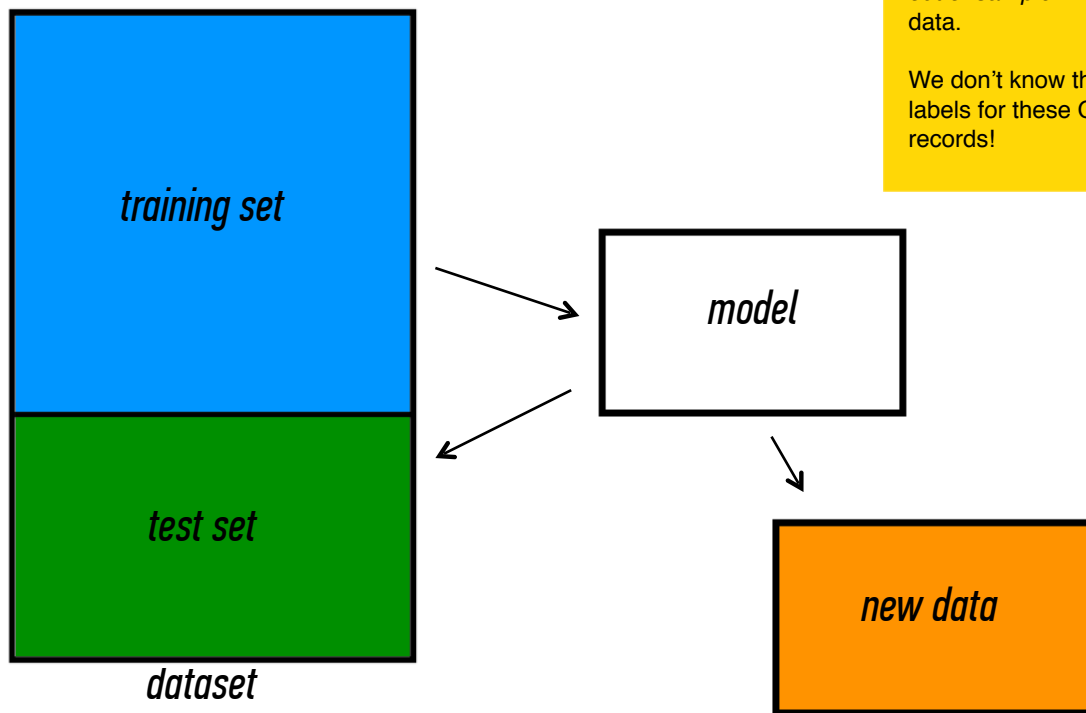
- 1) split dataset*
- 2) train model*
- 3) test model*
- 4) make predictions*



CLASSIFICATION PROBLEMS

Q: What steps does a classification problem require?

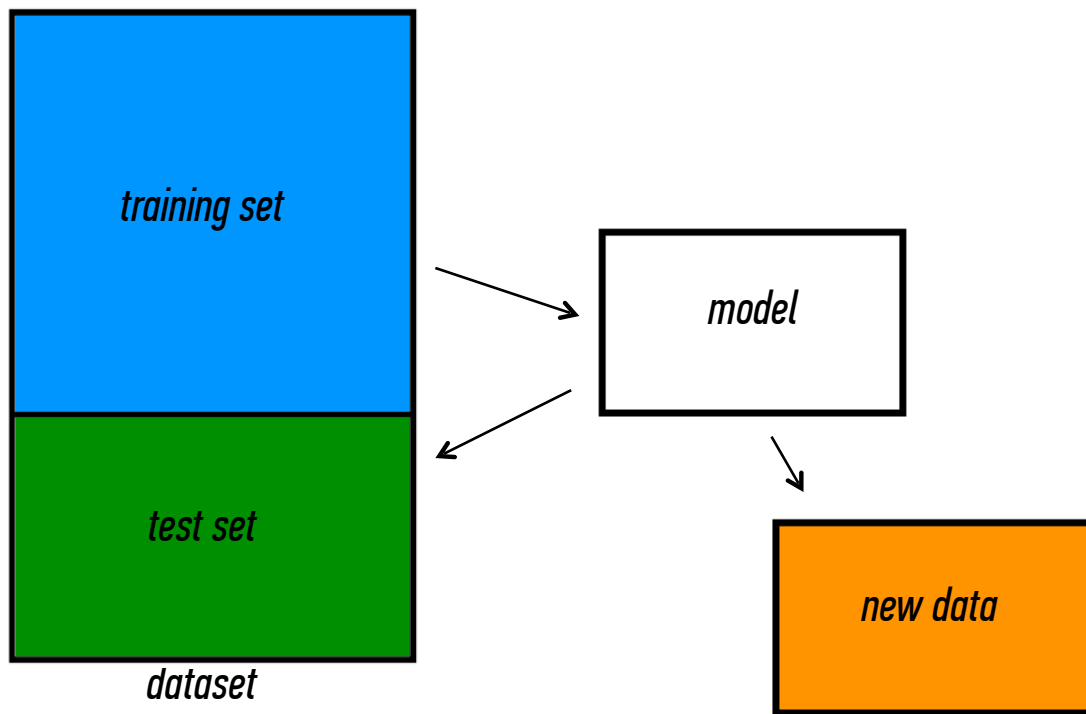
- 1) split dataset*
- 2) train model*
- 3) test model*
- 4) make predictions*



III. BUILDING EFFECTIVE CLASSIFIERS

BUILDING EFFECTIVE CLASSIFIERS

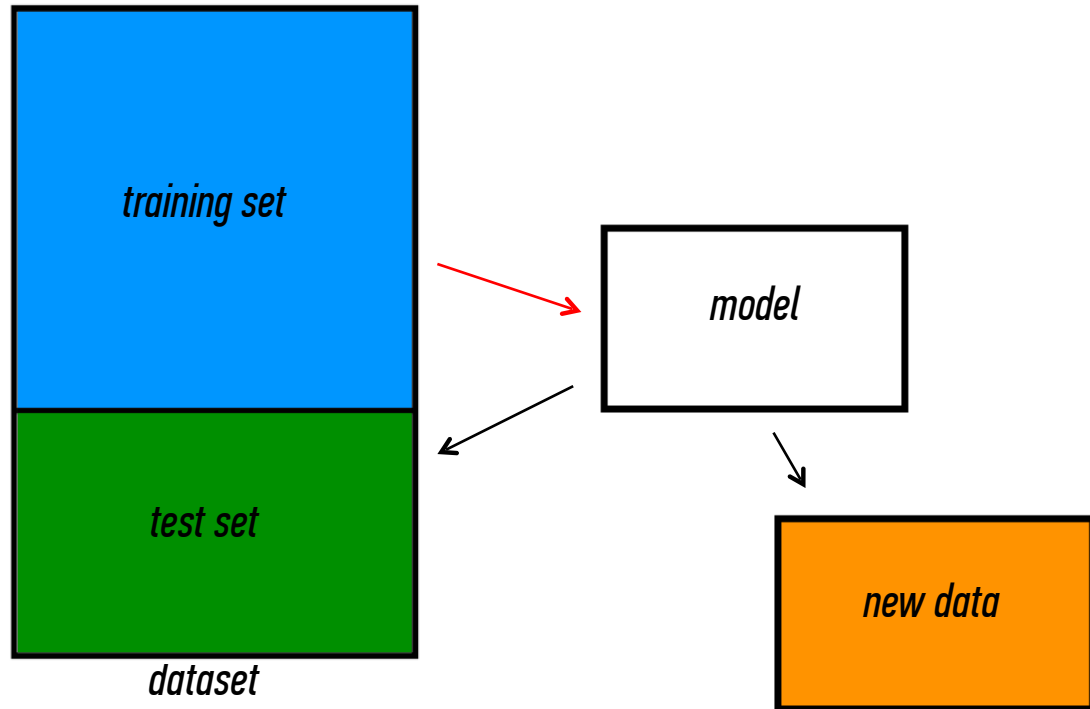
Q: What types of prediction error will we run into?



BUILDING EFFECTIVE CLASSIFIERS

Q: What types of prediction error will we run into?

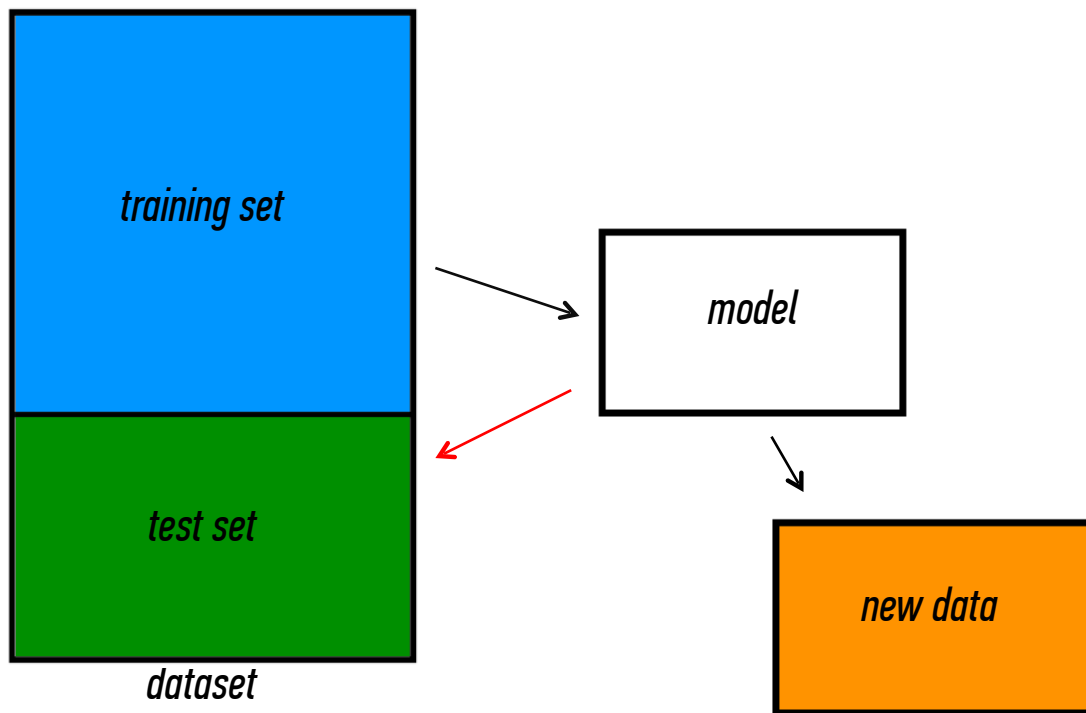
1) training error



BUILDING EFFECTIVE CLASSIFIERS

Q: What types of prediction error will we run into?

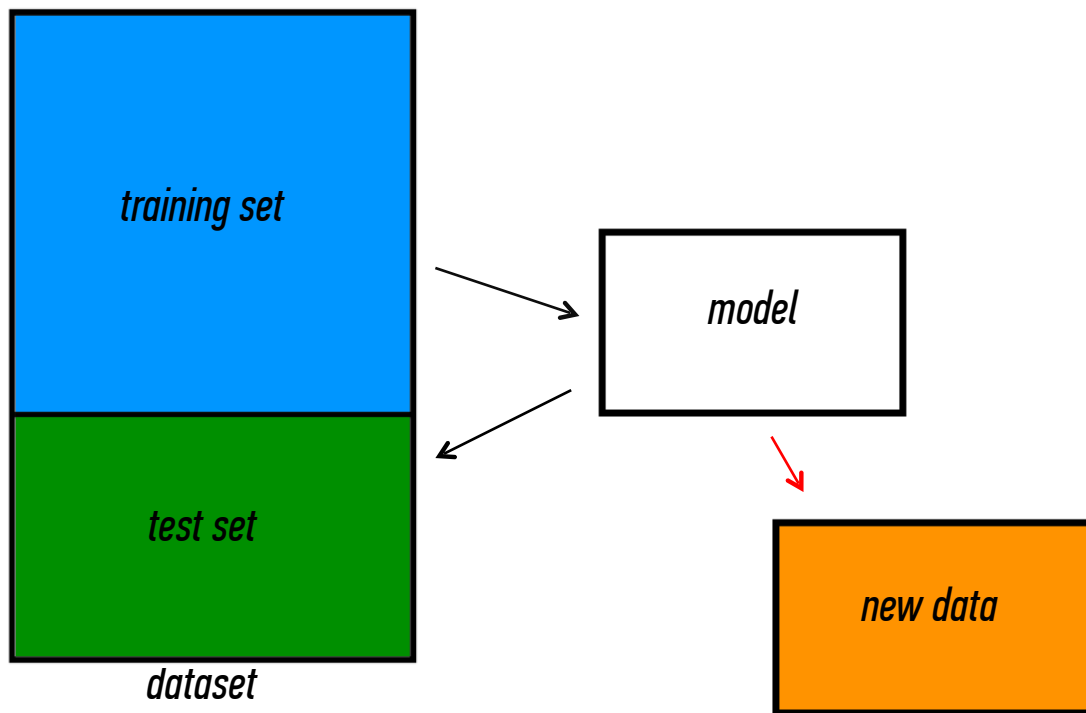
- 1) training error*
- 2) generalization error*



BUILDING EFFECTIVE CLASSIFIERS

Q: What types of prediction error will we run into?

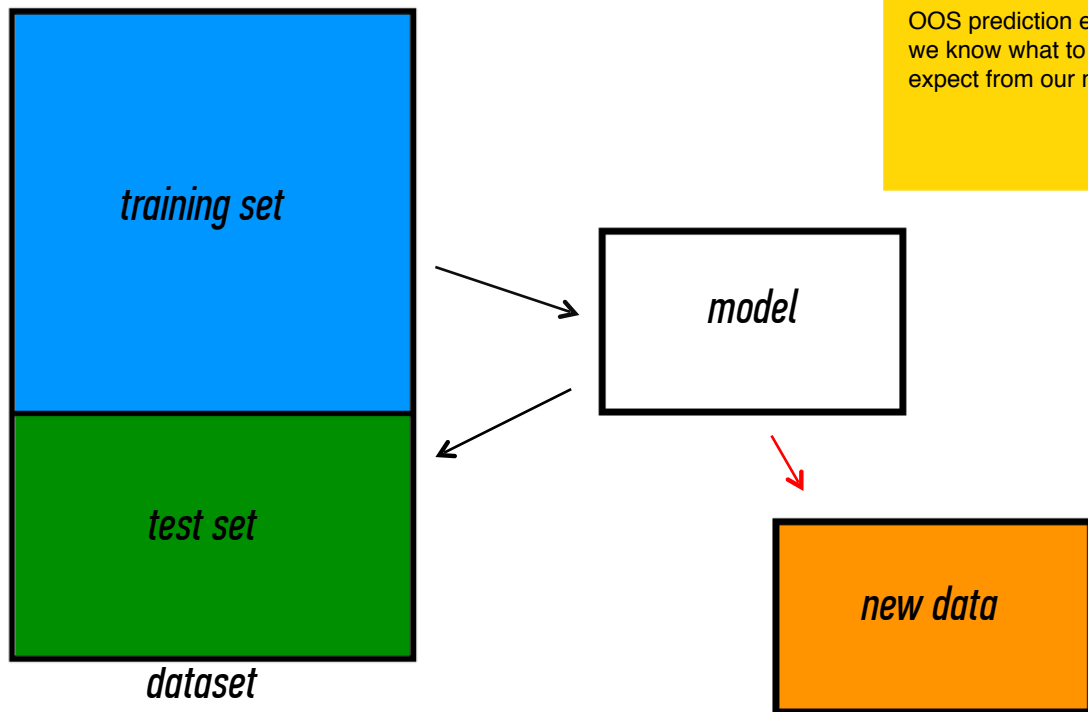
- 1) training error*
- 2) generalization error*
- 3) OOS error*



BUILDING EFFECTIVE CLASSIFIERS

Q: What types of prediction error will we run into?

- 1) training error*
- 2) generalization error*
- 3) OOS error*



NOTE

We want to estimate OOS prediction error so we know what to expect from our model.

TRAINING ERROR

Q: Why should we use training & test sets?

TRAINING ERROR

Q: Why should we use training & test sets?

Thought experiment:

Suppose instead, we train our model using the entire dataset.

TRAINING ERROR

Q: Why should we use training & test sets?

Thought experiment:

Suppose instead, we train our model using the entire dataset.

Q: How low can we push the training error?

TRAINING ERROR

Q: Why should we use training & test sets?

Thought experiment:

Suppose instead, we train our model using the entire dataset.

Q: How low can we push the training error?

- We can make the model arbitrarily complex (effectively “memorizing” the entire training set).*

TRAINING ERROR

Q: Why should we use training & test sets?

Thought experiment:

Suppose instead, we train our model using the entire dataset.

Q: How low can we push the training error?

- We can make the model arbitrarily complex (effectively “memorizing” the entire training set).*

A: Down to zero!

TRAINING ERROR

Q: Why should we use training & test sets?

Thought experiment:

Suppose instead, we train our model using the entire dataset.

Q: How low can we push the training error?

- We can make the model arbitrarily complex (effectively “memorizing” the entire training set).*

A: Down to zero!

NOTE

This phenomenon is called *overfitting*.

OVERFITTING

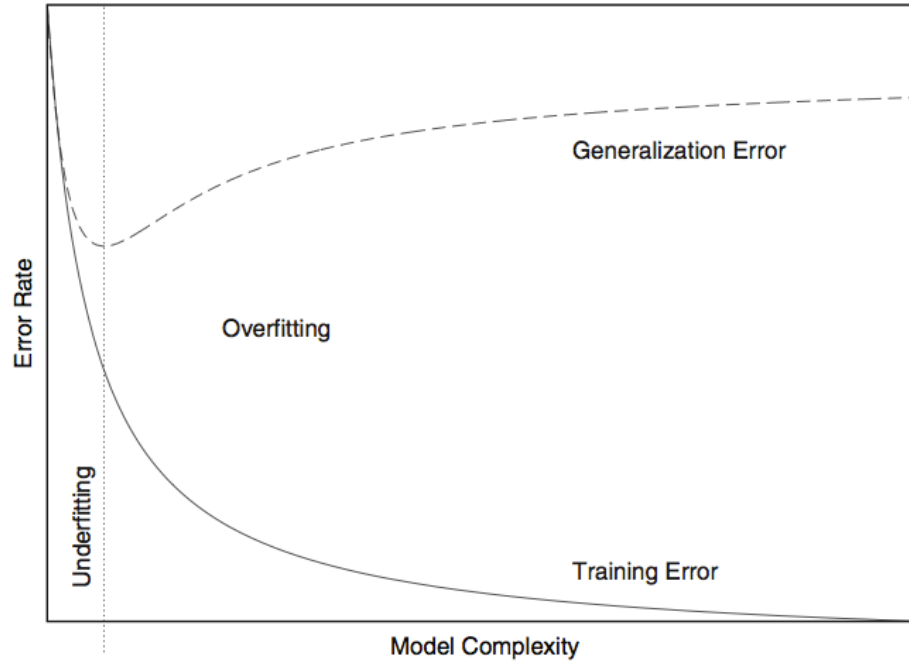
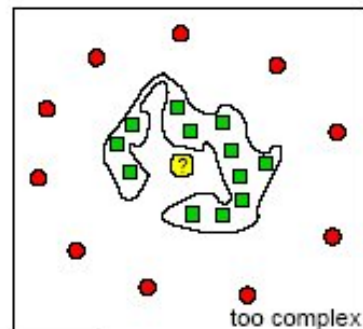
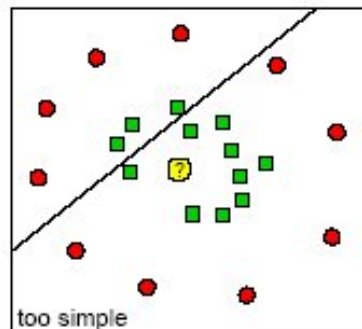


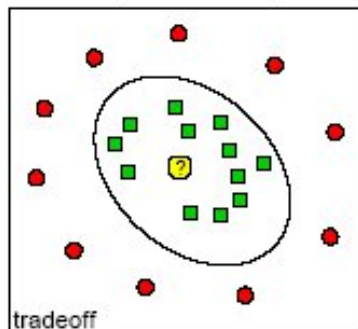
FIGURE 18-1. *Overfitting: as a model becomes more complex, it becomes increasingly able to represent the training data. However, such a model is overfitted and will not generalize well to data that was not used during training.*

OVERFITTING - EXAMPLE

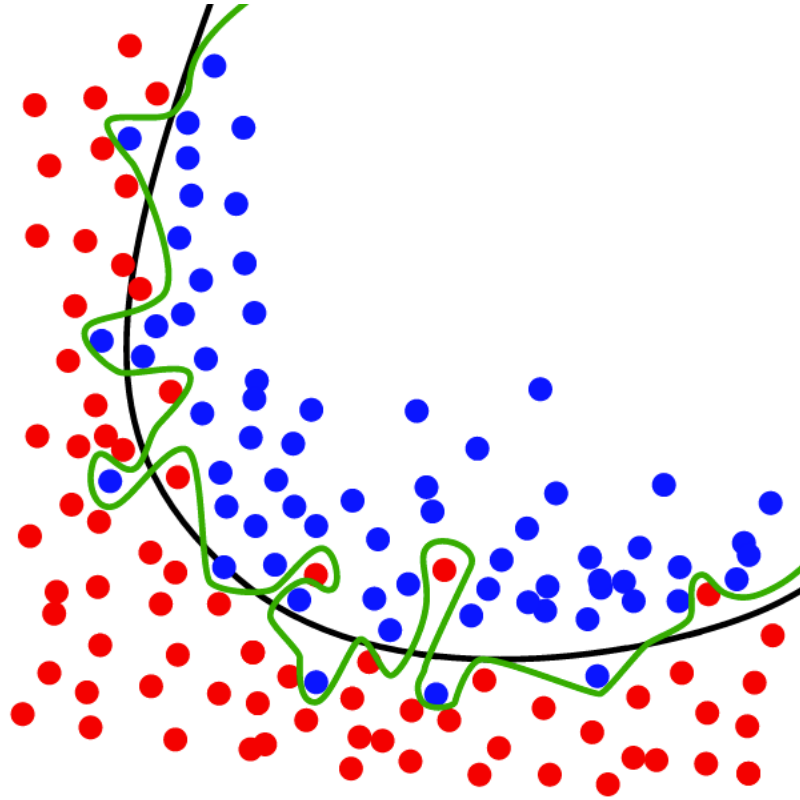
Underfitting and Overfitting



- negative example
- positive example
- new patient



OVERFITTING - EXAMPLE



TRAINING ERROR

Q: Why should we use training & test sets?

Thought experiment:

Suppose instead, we train our model using the entire dataset.

Q: How low can we push the training error?

- We can make the model arbitrarily complex (effectively “memorizing” the entire training set).*

A: Down to zero!

A: Training error is not a good estimate of OOS accuracy.

NOTE

This phenomenon is called *overfitting*.

GENERALIZATION ERROR

Suppose we do the train/test split.

GENERALIZATION ERROR

Suppose we do the train/test split.

Q: How well does generalization error predict OOS accuracy?

GENERALIZATION ERROR

Suppose we do the train/test split.

Q: How well does generalization error predict OOS accuracy?

Thought experiment:

Suppose we had done a different train/test split.

GENERALIZATION ERROR

Suppose we do the train/test split.

Q: How well does generalization error predict OOS accuracy?

Thought experiment:

Suppose we had done a different train/test split.

Q: Would the generalization error remain the same?

GENERALIZATION ERROR

Suppose we do the train/test split.

Q: How well does generalization error predict OOS accuracy?

Thought experiment:

Suppose we had done a different train/test split.

Q: Would the generalization error remain the same?

A: Of course not!

GENERALIZATION ERROR

Suppose we do the train/test split.

Q: How well does generalization error predict OOS accuracy?

Thought experiment:

Suppose we had done a different train/test split.

Q: Would the generalization error remain the same?

A: Of course not!

A: On its own, not very well.

GENERALIZATION ERROR

Suppose we do the train/test split.

Q: How well does generalization error predict OOS accuracy?

Thought experiment:

Suppose we had done a different train/test split.

Q: Would the generalization error remain the same?

A: Of course not!

A: On its own, not very well.

NOTE

The generalization error gives a *high-variance estimate* of OOS accuracy.

GENERALIZATION ERROR

Something is still missing!

GENERALIZATION ERROR

Something is still missing!

Q: How can we do better?

GENERALIZATION ERROR

Something is still missing!

Q: How can we do better?

Thought experiment:

Different train/test splits will give us different generalization errors.

GENERALIZATION ERROR

Something is still missing!

Q: How can we do better?

Thought experiment:

Different train/test splits will give us different generalization errors.

Q: What if we did a bunch of these and took the average?

GENERALIZATION ERROR

Something is still missing!

Q: How can we do better?

Thought experiment:

Different train/test splits will give us different generalization errors.

Q: What if we did a bunch of these and took the average?

A: Now you're talking!

GENERALIZATION ERROR

Something is still missing!

Q: How can we do better?

Thought experiment:

Different train/test splits will give us different generalization errors.

Q: What if we did a bunch of these and took the average?

A: Now you're talking!

A: Cross-validation.

CROSS-VALIDATION

Steps for n -fold cross-validation:

CROSS-VALIDATION

Steps for n -fold cross-validation:

1) Randomly split the dataset into n equal partitions.

CROSS-VALIDATION

Steps for n -fold cross-validation:

- 1) Randomly split the dataset into n equal partitions.*
- 2) Use partition 1 as test set & union of other partitions as training set.*

CROSS-VALIDATION

Steps for n -fold cross-validation:

- 1) Randomly split the dataset into n equal partitions.*
- 2) Use partition 1 as test set & union of other partitions as training set.*
- 3) Find generalization error.*

CROSS-VALIDATION

Steps for n -fold cross-validation:

- 1) Randomly split the dataset into n equal partitions.*
- 2) Use partition 1 as test set & union of other partitions as training set.*
- 3) Find generalization error.*
- 4) Repeat steps 2-3 using a different partition as the test set at each iteration.*

CROSS-VALIDATION

Steps for n -fold cross-validation:

- 1) Randomly split the dataset into n equal partitions.*
- 2) Use partition 1 as test set & union of other partitions as training set.*
- 3) Find generalization error.*
- 4) Repeat steps 2-3 using a different partition as the test set at each iteration.*
- 5) Take the average generalization error as the estimate of OOS accuracy.*

CROSS-VALIDATION

Features of n -fold cross-validation:

CROSS-VALIDATION

Features of n -fold cross-validation:

1) More accurate estimate of OOS prediction error.

CROSS-VALIDATION

Features of n -fold cross-validation:

- 1) More accurate estimate of OOS prediction error.*
- 2) More efficient use of data than single train/test split.*
 - Each record in our dataset is used for both training and testing.*

CROSS-VALIDATION

Features of n -fold cross-validation:

- 1) More accurate estimate of OOS prediction error.*
- 2) More efficient use of data than single train/test split.*
 - Each record in our dataset is used for both training and testing.*
- 3) Presents tradeoff between efficiency and computational expense.*
 - 10-fold CV is 10x more expensive than a single train/test split*

CROSS-VALIDATION

Features of n-fold cross-validation:

- 1) More accurate estimate of OOS prediction error.*
- 2) More efficient use of data than single train/test split.*
 - Each record in our dataset is used for both training and testing.*
- 3) Presents tradeoff between efficiency and computational expense.*
 - 10-fold CV is 10x more expensive than a single train/test split*
- 4) Can be used for model selection.*