# INTRO TO DATA SCIENCE
## CLASS 5: MODEL EVALUATION

# LAST TIME:

# – INTRO TO MACHINE LEARNING
# – OVERVIEW OF K NEAREST NEIGHBORS

# QUESTIONS?

# QUESTIONS?

## WHAT WAS THE MOST INTERESTING THING YOU LEARNED?

## WHAT WAS THE HARDEST TO GRASP?

# I. QUICK REVIEW OF CLASSIFICATION PROBLEMS
# II. ERRORS, UNDERFITTING & OVERFITTING
# III. CROSS VALIDATION

# IV. LAB: CROSS VALIDATION IN SCIKIT-LEARN

# I. QUICK REVIEW OF CLASSIFICATION PROBLEMS

| | continuous | categorical |
|---|---|---|
| supervised | ??? | ??? |
| unsupervised | ??? | ??? |

|  | **continuous** | **categorical** |
|---|---|---|
| **supervised** | regression | classification |
| **unsupervised** | dimension reduction | clustering |

# Here's (part of) an example dataset:

**Fisher's *Iris* Data**

| Sepal length ⇕ | Sepal width ⇕ | Petal length ⇕ | Petal width ⇕ | Species ⇕ |
|---|---|---|---|---|
| 5.1 | 3.5 | 1.4 | 0.2 | *I. setosa* |
| 4.9 | 3.0 | 1.4 | 0.2 | *I. setosa* |
| 4.7 | 3.2 | 1.3 | 0.2 | *I. setosa* |
| 4.6 | 3.1 | 1.5 | 0.2 | *I. setosa* |
| 5.0 | 3.6 | 1.4 | 0.2 | *I. setosa* |
| 5.4 | 3.9 | 1.7 | 0.4 | *I. setosa* |
| 4.6 | 3.4 | 1.4 | 0.3 | *I. setosa* |
| 5.0 | 3.4 | 1.5 | 0.2 | *I. setosa* |

# Here's (part of) an example dataset:

**features**

**Fisher's *Iris* Data**

| Sepal length ⇕ | Sepal width ⇕ | Petal length ⇕ | Petal width ⇕ | Species ⇕ |
|---|---|---|---|---|
| 5.1 | 3.5 | 1.4 | 0.2 | *I. setosa* |
| 4.9 | 3.0 | 1.4 | 0.2 | *I. setosa* |
| 4.7 | 3.2 | 1.3 | 0.2 | *I. setosa* |
| 4.6 | 3.1 | 1.5 | 0.2 | *I. setosa* |
| 5.0 | 3.6 | 1.4 | 0.2 | *I. setosa* |
| 5.4 | 3.9 | 1.7 | 0.4 | *I. setosa* |
| 4.6 | 3.4 | 1.4 | 0.3 | *I. setosa* |
| 5.0 | 3.4 | 1.5 | 0.2 | *I. setosa* |

*Here's (part of) an example dataset:*

**Fisher's *Iris* Data**

| Sepal length ⇕ | Sepal width ⇕ | Petal length ⇕ | Petal width ⇕ | Species ⇕ |
|---|---|---|---|---|
| 5.1 | 3.5 | 1.4 | 0.2 | *I. setosa* |
| 4.9 | 3.0 | 1.4 | 0.2 | *I. setosa* |
| 4.7 | 3.2 | 1.3 | 0.2 | *I. setosa* |
| 4.6 | 3.1 | 1.5 | 0.2 | *I. setosa* |
| 5.0 | 3.6 | 1.4 | 0.2 | *I. setosa* |
| 5.4 | 3.9 | 1.7 | 0.4 | *I. setosa* |
| 4.6 | 3.4 | 1.4 | 0.3 | *I. setosa* |
| 5.0 | 3.4 | 1.5 | 0.2 | *I. setosa* |

*features*

*class labels*
*(qualitative)*
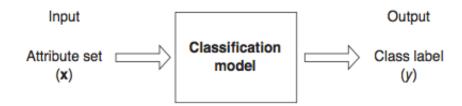
*Q: What does "supervised" mean?*

*Q: What does "supervised" mean?*

*A: We know the labels.*

**Fisher's *Iris* Data**

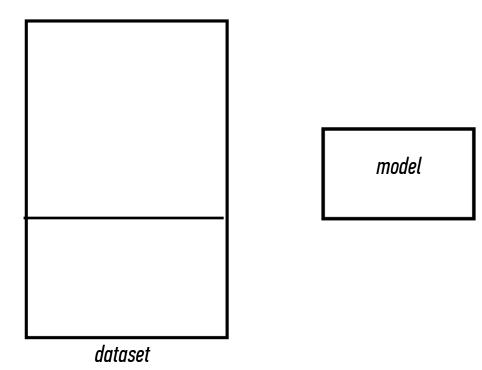| Sepal length ⇕ | Sepal width ⇕ | Petal length ⇕ | Petal width ⇕ | Species ⇕ |
| --- | --- | --- | --- | --- |
| 5.1 | 3.5 | 1.4 | 0.2 | *I. setosa* |
| 4.9 | 3.0 | 1.4 | 0.2 | *I. setosa* |
| 4.7 | 3.2 | 1.3 | 0.2 | *I. setosa* |
| 4.6 | 3.1 | 1.5 | 0.2 | *I. setosa* |
| 5.0 | 3.6 | 1.4 | 0.2 | *I. setosa* |
| 5.4 | 3.9 | 1.7 | 0.4 | *I. setosa* |
| 4.6 | 3.4 | 1.4 | 0.3 | *I. setosa* |
| 5.0 | 3.4 | 1.5 | 0.2 | *I. setosa* |

*class labels*

*(qualitative)*

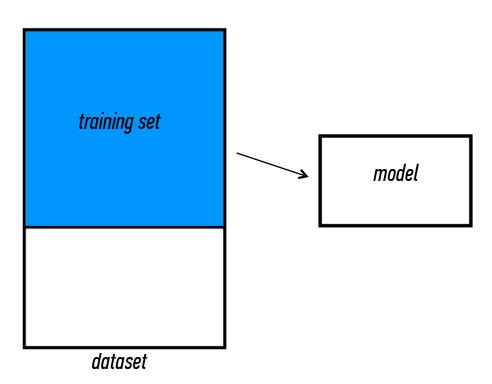*Q: How does a classification problem work?*
*A: Data in, predicted labels out.*



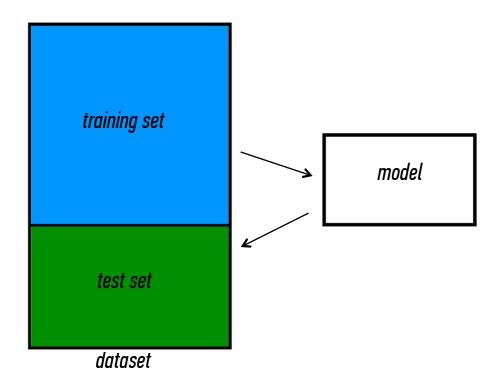**Figure 4.2.** Classification as the task of mapping an input attribute set $x$ into its class label $y$.

*Q: What steps does a classification problem require?*

dataset

model

# Q: What steps does a classification problem require?

1) split dataset

dataset

model

*Q: What steps does a classification problem require?*

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



training set

model

dataset

# Q: What steps does a classification problem require?

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

## Q: What steps does a classification problem require?

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



training set

test set

dataset

new data

model

predictions

# Q: What steps does a classification problem require?

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



training set

test set

dataset

model

new data

predictions

**NOTE**
This new data is called *out of sample* data.

We don't know the labels for these OOS records!

# II. ERRORS, UNDERFITTING & OVERFITTING

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

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

1) training error

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

1) training error
2) generalization error

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

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

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

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



training set

test set

dataset

model

new data

**NOTE**

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

predictions

*Q: Why should we use training & test sets?*

*Q: Why should we use training & test sets?*

*Thought experiment:*
*Suppose instead, we train our model using the entire dataset.*

# 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?*

*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).*

*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!*

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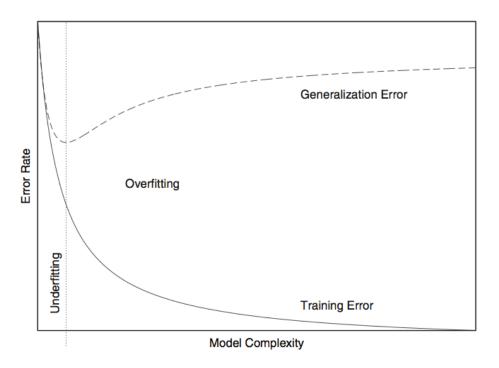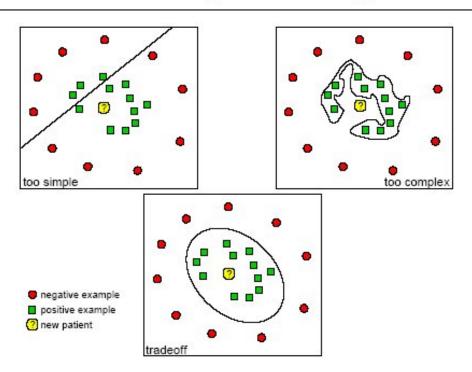
**FIGURE 18-1.** *Overfitting: as a model becomes more complex, it becomes increasingly able to represent the training data. However, such a model is overfitted and will not generalize well to data that was not used during training.*

source: *Data Analysis with Open Source Tools*, by Philipp K. Janert. O'Reilly Media, 2011.

source: http://www.dtreg.com

*source: http://www.dtreg.com*

# 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!
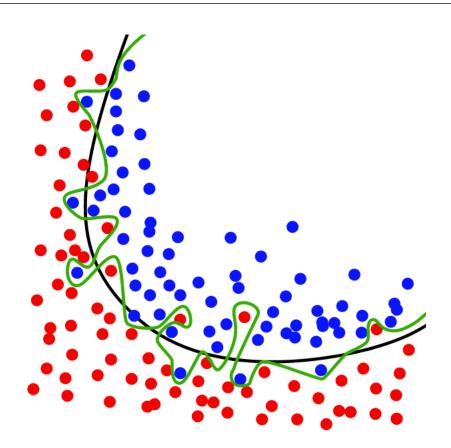
**NOTE**

This phenomenon is called *overfitting*.

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

*Suppose we do the train/test split.*

*Suppose we do the train/test split.*

*Q: How well does generalization error predict OOS accuracy?*

*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.*

*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?*

*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!*

*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.*

*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.

# BIAS-VARIANCE

*Something is still missing!*

*Something is still missing!*

*Q: How can we do better?*

*Something is still missing!*

*Q: How can we do better?*

*Thought experiment:*

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

*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?*

# 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!

*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.*

# III. CROSS VALIDATION

*Steps for n-fold cross-validation:*

# Steps for n-fold cross-validation:

1) Randomly split the dataset into n equal partitions.

## 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.

*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.*

*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.*

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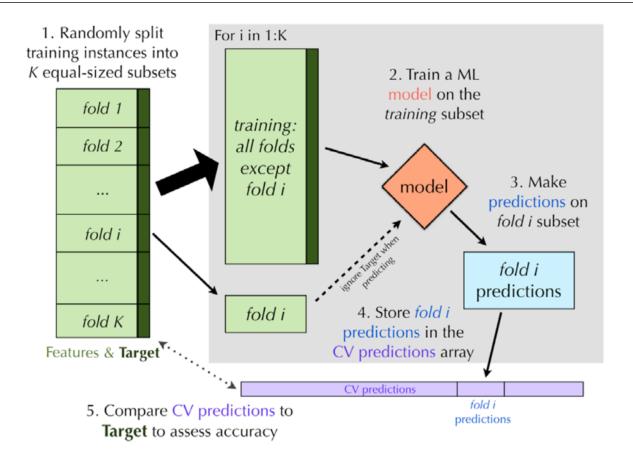
1) *Randomly split the dataset into n equal partitions.*
2) *Use partition 1 as test set & union of other partitions as training set.*
3) *Find generalization error.*
4) *Repeat steps 2-3 using a different partition as the test set at each iteration.*
5) *Take the average generalization error as the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

**CROSS-VALIDATION: 5-FOLD EXAMPLE**

| Dataset | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 | Accuracy |
|---|---|---|---|---|---|---|
| 1 | Test | Train | Train | Train | Train | $k_1$ % |
| 2 | Train | Test | Train | Train | Train | $k_2$ % |
| 3 | Train | Train | Test | Train | Train | $k_3$ % |
| 4 | Train | Train | Train | Test | Train | $k_4$ % |
| 5 | Train | Train | Train | Train | Test | $k_5$ % |

*5-Fold Generalization Error = $(k_1 + k_2 + k_3 + k_4 + k_5) / 5$*

1. Randomly split training instances into K equal-sized subsets

fold 1
fold 2
...
fold i
...
fold K

Features & **Target**

For i in 1:K

training: all folds except fold i

fold i

ignore Target when predicting

2. Train a ML model on the training subset

model

3. Make predictions on fold i subset

fold i predictions

4. Store fold i predictions in the CV predictions array

CV predictions

fold i predictions

5. Compare CV predictions to **Target** to assess accuracy

*Features of n-fold cross-validation:*

# Features of n-fold cross-validation:

1) More accurate estimate of OOS prediction error.

*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.*

# 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

*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.*

*Last time:*

- *Types of machine learning problems / algorithms*
- *Generalization*

*This time:*

- *Train / Test Split*
- *Errors, Overfitting and underfitting*
- *Cross validation*

# LAB: CROSS VALIDATION WITH KIN