

Changelog: 14 Oct 2016, 30 Oct 2106

## Basic Concepts

*Induction, Generalization, Evaluation*

Lecture 2: Part 1

Marina Santini

Acknowledgements  
Pictures retrieved from the web

2016 Lecture 2: Basic Concepts - Part 1 1

## Lecture 2: Required Reading

- Handout
- Daumé III (2015: 8-10; 19-24)
- Witten et al. (2011): Ch 2

2016 Lecture 2: Basic Concepts - Part 1 2

## Outline

- Induction
- Generalization
- Splitting the Data
- Evaluation Measures
- Modelling
- Parameters
- Inductive Bias

2016 Lecture 2: Basic Concepts - Part 1 3

## Supervised Learning

- also called Inductive Learning
- Inductive learning is the general theory behind supervised learning

2016 Lecture 2: Basic Concepts - Part 1 4

## What is inductive learning?

- input (x)
- output (y)
- estimate the function (f)

More specifically:  
 "the problem is to generalize from the samples and the mapping to be useful to estimate the output for new samples in the future. In practice it is almost always too hard to estimate the function, so we are looking for very good approximations of the function."

Jason Browlee  
ML Mastery Blog

2016 Lecture 2: Basic Concepts - Part 1 5

## Induction

**Induction** is the process of reaching a general conclusion from specific examples

```

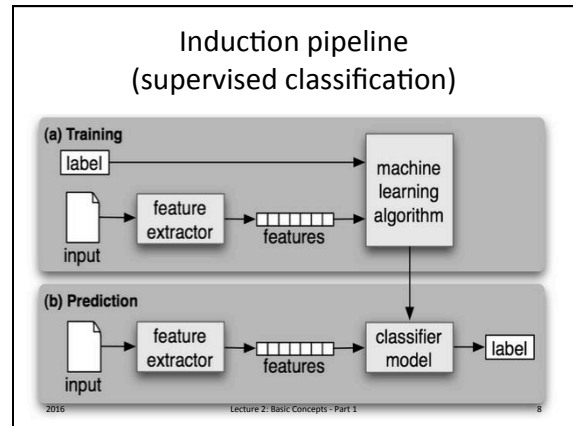
graph LR
    A[Examples] -- Generalize --> B[Model]
    B -- "Instantiate for another case" --> C((Prediction))
    
```

2016 Lecture 2: Basic Concepts - Part 1 6

### Inductive Machine Learning (Simplified)

- The goal of inductive machine learning is to take some training data and use it to induce a model.
- This model will be evaluated on the test data
- The machine learning algorithm has succeeded if its performance on the test data is high.

2016 Lecture 2: Basic Concepts - Part 1 7



### Example: Predict the class of an unseen iris flower

Sepal length – Sepal width – Petal length – Petal width – Type  
5.2      3.7      1.7      0.3      ???

Sepal length (cm)	Sepal width (cm)	Petal length (cm)	Petal width (cm)	Type
5.1	3.5	1.4	0.2	<i>Iris setosa</i>
4.9	3.0	1.4	0.2	<i>Iris setosa</i>
4.7	3.2	1.3	0.2	<i>Iris setosa</i>
4.6	3.1	1.5	0.2	<i>Iris setosa</i>
5.0	3.6	1.4	0.2	<i>Iris setosa</i>
7.0	3.2	4.7	1.4	<i>Iris versicolor</i>
6.4	3.2	4.5	1.5	<i>Iris versicolor</i>
6.9	3.1	4.9	1.5	<i>Iris versicolor</i>
5.5	2.3	4.0	1.3	<i>Iris versicolor</i>
6.5	2.8	4.6	1.5	<i>Iris versicolor</i>
6.3	3.3	6.0	2.5	<i>Iris virginica</i>
5.8	2.7	5.1	1.9	<i>Iris virginica</i>
7.1	3.0	5.9	2.1	<i>Iris virginica</i>
6.3	2.9	5.6	1.8	<i>Iris virginica</i>
6.5	3.0	5.8	2.2	<i>Iris virginica</i>

Require us to generalize from the training data

2016 Lecture 2: Basic Concepts - Part 1 9

### Data: The iris dataset

Sepal length (cm)	Sepal width (cm)	Petal length (cm)	Petal width (cm)	Type
5.1	3.5	1.4	0.2	<i>Iris setosa</i>
4.9	3.0	1.4	0.2	<i>Iris setosa</i>
4.7	3.2	1.3	0.2	<i>Iris setosa</i>
4.6	3.1	1.5	0.2	<i>Iris setosa</i>
5.0	3.6	1.4	0.2	<i>Iris setosa</i>
7.0	3.2	4.7	1.4	<i>Iris versicolor</i>
6.4	3.2	4.5	1.5	<i>Iris versicolor</i>
6.9	3.1	4.9	1.5	<i>Iris versicolor</i>
5.5	2.3	4.0	1.3	<i>Iris versicolor</i>
6.5	2.8	4.6	1.5	<i>Iris versicolor</i>
6.3	3.3	6.0	2.5	<i>Iris virginica</i>
5.8	2.7	5.1	1.9	<i>Iris virginica</i>
7.1	3.0	5.9	2.1	<i>Iris virginica</i>
6.3	2.9	5.6	1.8	<i>Iris virginica</i>
6.5	3.0	5.8	2.2	<i>Iris virginica</i>

**Three components:**

1. Class label (aka "label", denoted y)
2. Features (aka "attributes")
3. Feature values (aka "attribute values", denoted x) ⇒ Features can be binary, nominal or continuous

• A labeled dataset is a collection of (x,y) pairs

2016 Lecture 2: Basic Concepts - Part 1 10

### Task

- Predict the class for this "test" example:

Sepal length – Sepal width – Petal length – Petal width – Type  
5.2      3.7      1.7      0.3      ???

Sepal length (cm)	Sepal width (cm)	Petal length (cm)	Petal width (cm)	Type
5.1	3.5	1.4	0.2	<i>Iris setosa</i>
4.9	3.0	1.4	0.2	<i>Iris setosa</i>
4.7	3.2	1.3	0.2	<i>Iris setosa</i>
4.6	3.1	1.5	0.2	<i>Iris setosa</i>
5.0	3.6	1.4	0.2	<i>Iris setosa</i>
7.0	3.2	4.7	1.4	<i>Iris versicolor</i>
6.4	3.2	4.5	1.5	<i>Iris versicolor</i>
6.9	3.1	4.9	1.5	<i>Iris versicolor</i>
5.5	2.3	4.0	1.3	<i>Iris versicolor</i>
6.5	2.8	4.6	1.5	<i>Iris versicolor</i>
6.3	3.3	6.0	2.5	<i>Iris virginica</i>
5.8	2.7	5.1	1.9	<i>Iris virginica</i>
7.1	3.0	5.9	2.1	<i>Iris virginica</i>
6.3	2.9	5.6	1.8	<i>Iris virginica</i>
6.5	3.0	5.8	2.2	<i>Iris virginica</i>

Require us to generalize from the training data

2016 Lecture 2: Basic Concepts - Part 1 11

### Generalization

- Machine learning is all about finding patterns in data.
- The most central concept in machine learning is **generalization**: how to generalize beyond the examples that have been provided at "training time." to new examples that you see at "test time."

2016 Lecture 2: Basic Concepts - Part 1 12

## Generalization

- Learning is generalization not memorization
- Generalization in the context of ML is the ability of a learning system to perform accurately on new/unseen data after having learned from existing examples.

2016 Lecture 2: Basic Concepts - Part 1 13

## Ex:

2016 Lecture 2: Basic Concepts - Part 1 14

## Splitting the data (1)

- Training data & Test Data:
  - Common splits in LT: 70/30; 80/20; 90/10
- **TEST DATA MUST BELONG TO THE SAME STATISTICAL DISTRIBUTION AS THE TRAINING DATA**

2016 Lecture 2: Basic Concepts - Part 1 15

## Face Recognition

Training examples of a person

Test images

ORI dataset,  
AT&T Laboratories, Cambridge UK

Lecture 1: Definitions and Examples 2016

## Example: Letter vs non-letter classification

Training set

Test set

2016

## Cross-Validation: 10-fold

- In 10-fold cross-validation you break you training data up into 10 equally-sized partitions.
- You train a learning algorithm on 9 of them and test it on the remaining 1.
- You do this 10 times, each holding out a different partition as the test data.
- Typical choices for n-fold are 2, 5, 10.
- 10-fold cross validation is the most common.

2016 Lecture 2: Basic Concepts - Part 1 18

### Leave One Out

- Leave One Out (or LOO) is a simple cross-validation. Each learning set is created by taking all the samples except one, the test set being the sample left out.

2016

Lecture 2: Basic Concepts - Part 1

19

### Stratification

- Proportion of each class in the training set and test sets is the same as the proportion in the original sample.

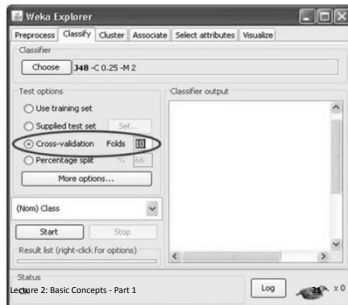
2016

Lecture 2: Basic Concepts - Part 1

20

### Weka Evaluation Methods

- Ex: 10-fold cross-validation



2016

Lecture 2: Basic Concepts - Part 1

21

### True and False Positives and Negatives

- **True positives** are correct labels that are correctly identified.
- **True negatives** are incorrect labels that are correctly identified as incorrect.
- **False positives (or Type I errors)** are incorrect labels that are incorrectly identified as correct.
- **False negatives (or Type II errors)** are correct labels that are incorrectly identified as incorrect.

2016

Lecture 2: Basic Concepts - Part 1

22

### Confusion Matrix

- Usually, the rows are the observed/actual class labels and the columns the predicted class labels.
- Each cell contains the number of predictions made by the classifier that fall into that cell.

		Predicted	
		Positive	Negative
Actual	Positive	True Positive	False Positive
	Negative	False Negative	True Negative

Truth Table Confusion Matrix

2016

Lecture 2: Basic Concepts - Part 1

23

### Multi-Class Confusion Matrix

If a classification system has been trained to distinguish between cats, dogs and rabbits, a confusion matrix will summarize the results:

		Predicted		
		Cat	Dog	Rabbit
Actual class	Cat	5	3	0
	Dog	2	3	1
	Rabbit	0	2	11

2016

Lecture 2: Basic Concepts - Part 1

24

### True and False Positives and Negatives & evaluation measures

Given these four numbers , we can define the following measures:

- Accuracy
- Precision
- Recall
- F-Measure

2016 Lecture 2: Basic Concepts - Part 1 25

### Accuracy

Accuracy measures the percentage of correct results that a classifier has achieved.

$$\text{Accuracy} = (TP + TN) / (TP + TN + FP + FN)$$

Watch out: accuracy is a tricky metric:  
it can be biased towards the most frequent class

2016 Lecture 2: Basic Concepts - Part 1 26

### Precision, Recall, F-Measure in Classification

- **Precision** the number of correctly classified positive examples divided by the number of examples labeled by the system as positive  $P = TP / (TP + FP)$
- **Recall** the number of correctly classified positive examples divided by the number of positive examples in the data.  $R = TP / (TP + FN)$
- The **F-Measure** (or **F-Score**) a combination of the above.  
Sokolova and Lapalme (2009)

2016 Lecture 2: Basic Concepts - Part 1 27

### Confusion matrix: Sokolova and Lapalme (2009)

“A systematic analysis of performance measures for classification tasks”

Data class	Classified as <i>pos</i>	Classified as <i>neg</i>	$\begin{bmatrix} tp & fn \\ fp & tn \end{bmatrix}$
<i>pos</i>	true positive ( <i>tp</i> )	false negative ( <i>fn</i> )	
<i>neg</i>	false positive ( <i>fp</i> )	true negative ( <i>tn</i> )	

2016 Lecture 2: Basic Concepts - Part 1 28

### Summary: Accuracy, Precision, Recall, F-measure

- Accuracy =  $(TP + TN) / (TP + TN + FP + FN)$
- Precision =  $TP / (TP + FP)$
- Recall =  $TP / (TP + FN)$
- F-measure (harmonic) =  $2 * ((precision * recall) / (precision + recall))$

2016 Lecture 2: Basic Concepts - Part 1 29


### Weka: Output

Example of a classifier output

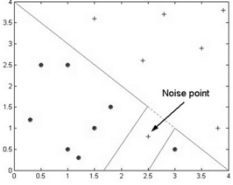
2016 Lecture 2: Basic Concepts - Part 1

### Disturbing factors: Noise

- Unexplained or random variation in the data
- Anomaly



#### Overfitting due to Noise



Lecture 2: Basic Concepts - Part 1 31

### Underfitting & Overfitting

- Underfitting: the model has not learned enough from the data and is unable to generalize.
- Overfitting: the model has learned too many idiosyncrasies (noise) and is unable to generalize.

2016 Lecture 2: Basic Concepts - Part 1 32

### Generalization: Overfitting & Underfitting

- Overfitting occurs when the model fits the training data too well and does not generalize so it performs badly on the test data. Overfitting is often a result of an excessively complicated model.
- Underfitting occurs when the model does **not** fit the data well enough. Underfitting is often a result of an excessively simple model.
- Both overfitting and underfitting lead to **poor predictions** on unseen examples.
- A model that overfits or underfits does not generalize well.

2016 Lecture 2: Basic Concepts - Part 1 33

### Parameters

Models can have many parameters and finding the best combination of parameters is not trivial.

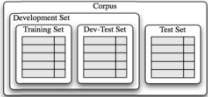
2016 Lecture 2: Basic Concepts - Part 1 34

### Hyperparameters

A hyperparameter is a parameter that controls other parameters of the model.

2016 Lecture 2: Basic Concepts - Part 1 35

### Ex: Testing parameters and hyperparameters via development set



- Split your data into 70% training data, 10% development data and 20% test data.
- For each possible setting of the hyperparameters:
  - Train a model using that setting on the training data
  - Compute the model error rate on the development data
  - From the above collection of models, choose the one that achieve the lowest error rate on development data.
  - Evaluate that model on the test data to estimate future **test performance**.

2016 Lecture 2: Basic Concepts - Part 1 36

## Inductive Bias: Definition

"The inductive bias of a learning algorithm is the set of assumptions that the learner uses to predict outputs given inputs that it has not encountered".

Tom Mitchell, 1980

2016

Lecture 2: Basic Concepts - Part 1

37

## Inductive Bias

Each learning method has a different inductive bias.

The inductive bias affects the performance of a classifier.

2016

Lecture 2: Basic Concepts - Part 1

38

## Not everything is learnable

- Noise at feature level
- Noise at class label level
- Features are insufficient
- Class labels are controversial
- Inductive bias not appropriate for the kind of problem we try to learn

2016

Lecture 2: Basic Concepts - Part 1

39

## Key Point

- Our goal when we choose a machine learning model is the following:
  - the model should perform well on future, unseen data.

2016

Lecture 2: Basic Concepts - Part 1

40

## Quiz 1: Stratification

What does it mean "stratified" cross-validation?

1. The examples of a class are all in the training set, and the rest of the classes are in the test set.
2. The proportion of each class in the sets are the same as the proportion in the original sample
3. None of the above.

2016

Lecture 2: Basic Concepts - Part 1

41

## Quiz 2: Data Splits

Which are recommended splits between training and test data?

1. 80/20
2. 50/50
3. 10/90

2016

Lecture 2: Basic Concepts - Part 1

42

### Quiz 3: Overfitting

What does "overfitting" mean?

1. The model has not learned enough from the data and is unable to generalize
2. The proportion of each class in the sets is the same as the proportion in the original sample
3. None of the above.

2016

Lecture 2: Basic Concepts - Part 1

43

### Quiz 4: Accuracy

Why is accuracy unreliable?

1. Because it can be biased towards the most frequent class.
2. Because it always guesses wrong.
3. None of the above

2016

Lecture 2: Basic Concepts - Part 1

44

The end

2016

Lecture 2: Basic Concepts - Part 1

45