# Tree Classifiers PU Machine learning mit Java, Weka und UIMA

Nils Reiter, nils.reiter@uni-koeln.de

> December 1, 2020 Winter term 2020/21

# ▲ Organisatorisches

### 📽 Übungen

- Studienleistung besteht in der Abgabe aller\* Übungen
- Bitte auch halbfertige Lösungen pushen, zeigen Rede- und Erklärbedarf

\*...mit den üblichen Ausnahmen für Krankheit, Coronablues, Liebeskummer etc.

# ▲ Organisatorisches

📽 Übungen

- Studienleistung besteht in der Abgabe aller\* Übungen
- Bitte auch halbfertige Lösungen pushen, zeigen Rede- und Erklärbedarf
- P Anmeldung zu Modulprüfungen: in den nächsten Tagen via Klips

\*...mit den üblichen Ausnahmen für Krankheit, Coronablues, Liebeskummer etc. Reiter Tree Classifiers

# A Organ<u>isatorisches</u>



#### Tree Classifiers

<sup>\*...</sup>mit den üblichen Ausnahmen für Krankheit, Coronablues, Liebeskummer etc.

Exercise 4

# Section 1

# Exercise 4

```
20
     public static void main(String[] args) throws Exception {
       // Parse command line options. Not in exercise. but useful
21
       // This is done with the libarary JewelCLT
22
23
       // http://jewelcli.lexicalscope.com
24
       Options options = CliFactory.parseArguments(Options.class, args);
25
26
       // load data set
27
       File inputFile = new File(options.getInput()):
28
       ArffLoader loader = new ArffLoader():
29
       loader.setFile(inputFile):
30
       Instances instances = loader.getDataSet():
31
       instances.setClassIndex(instances.numAttributes() - 1);
32
33
       // Inititalize and use first filter
34
       StringToNominal filter0 = new StringToNominal();
35
       filter0.setAttributeRange("first-last"):
36
       filter0.setInputFormat(instances);
37
       instances = Filter (useFilter (instances, filter0);
38
39
       // Inititalize and use second filter
       MergeInfrequentNominalValues filter1 = new MergeInfrequentNominalValues();
40
41
       filter1.setAttributeIndices("first-last"):
42
       filter1.setMinimumFrequency(10):
43
       filter1.setInputFormat(instances):
44
       instances = Filter.useFilter(instances. filter1);
45
46
       // Initialise the classifier. potentially with parameters
47
       NaiveBayes nb = new NaiveBayes();
48
49
       // Evaluation with the built-in handling of cross validation and evaluation
50
       Evaluation evaluation = new Evaluation(instances);
51
       evaluation.crossValidateModel(nb, instances, options.getNumberOfFolds(), new Random(options.getRandomSeed()));
52
       System.out.println(evaluation.toClassDetailsString());
53
```

4 5	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class	-
6	0,733	0,122	0,983	0,733	0,840	0,379	0,893	0,986	0	
7	0,661	0,051	0,253	0,661	0,366	0,385	0,946	0,515	B-EntityPER	
8	0,568	0,086	0,200	0,568	0,296	0,297	0,842	0,360	I-EntityPER	
9	0,720	0,053	0,131	0,720	0,222	0,292	0,938	0,256	B-EntityLOC	
.0	0,700	0,073	0,150	0,700	0,247	0,300	0,893	0,420	I-EntityLOC	
1	0,000	0,002	0,000	0,000	0,000	-0,001	0,475	0,000	B-EntityWRK	
.2	0,000	0,005	0,000	0,000	0,000	-0,002	0,690	0,002	I-EntityWRK	
3 Weighted Avg.	0,723	0,117	0,910	0,723	0,789	0,373	0,892	0,931		

1 Nov 30, 2020 10:56:39 AM com.github.fommil.jni.JniLoader liberalLoad s

### **Potential Issues**

Filters need to know their input format first

filter.setInputFormat(instances)

Order of filter application matters

► Do feature removal first, because performance

### **Potential Issues**

- Filters need to know their input format first
  - filter.setInputFormat(instances)
- Order of filter application matters
- Do feature removal first, because performance
- Evaluation
  - weka.classifiers.Evaluation is easy to use and handles everything
  - Sometimes you need more control (e.g., debugging)  $\rightarrow$  next slide

Exercise 4

### Reference solution: doCrossValidation(...)

```
65
     public static void doCrossValidation(Instances instances, Classifier classifier, Options options) throws Exception {
66
67
       // for each fold
68
      for (int f = 0: f < options.getNumberOfFolds(): f++) {</pre>
69
         System.err.print("Fold " + f + ": "):
70
71
         // split the data into train and test
         // (we don't have control over the random generator here)
72
         Instances train = instances.trainCV(options.getNumberOfFolds(), f); 7
73
         Instances test = instances.testCV(options.getNumberOfFolds(), f);
74
75
76
         // train classifier on training data
         classifier_buildClassifier(train)
77
78
79
         // Use Evaluation class to get predictions
80
         Evaluation eval = new Evaluation(test):
81
         double[] predictions = eval.evaluateModel(classifier.test):
82
83
         // print all or the first 100 (whichever is smaller) predictions
         // from this fold
84
85
         for (int i = 0: i < (Math.min(test.numInstances(), 100)); i++) {</pre>
           System.err.print(test.get(i).classValue() + " ");
86
87
           System.err.print(predictions[i] + "|"):
88
         3
89
         System.err.println():
90
91
```

#### Exercise 4

# Reference solution: doCrossValidation(...)



 15
 Fold
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0.0
 0

- Computers don't have random numbers (because they are deterministic machines)
- Real randomness can only achieved with specialised hardware



Computer & Accessories > Accessories > Cables & Accessories > Cables > USB Cables





#### Truerng V3 – USB Hardware

Zufallszahlengenerator

★★★★☆ × 25 ratings

#### Our Price: €97.44 & FREE Delivery

Prices for items sold by Amazon include temporarily reduced VAT. Depending on your delivery address, VAT may vary at Checkout. For other items, please see details. Information on the reduced VAT in Germany.

#### New (4) from €97.44 + FREE Shipping

- High performance Speed: > 400 kilobit/second
- Improved Internal Whitening
- Native support for Windows and Linux

#### €97.44

& FREE Delivery

Arrives: **Dec 18 - 31** Details Fastest delivery: **Dec 3 - 5** Details

#### In stock.

Quantity: 1 \$

 Add to Basket

 Buy Now

#### Secure transaction

Dispatched from and sold by Boss Bross. For Returns, please check the seller link.

#### Reiter

#### Tree Classifiers

- Computers don't have random numbers (because they are deterministic machines)
- ▶ Real randomness can only achieved with specialised hardware
- Computers: Pseudo-random numbers
  - $\blacktriangleright$  Start with a seed value  $x_0$
  - Function to generate  $x_{n+1}$  from  $x_n$ 
    - This is a research area on its own

- Computers don't have random numbers (because they are deterministic machines)
- Real randomness can only achieved with specialised hardware
- Computers: Pseudo-random numbers
  - $\blacktriangleright$  Start with a seed value  $x_0$
  - Function to generate  $x_{n+1}$  from  $x_n$ 
    - This is a research area on its own
- Possible seed value sources
  - Mouse movement (not really random ...)
  - Memory allocation
  - Fixed number (0, 23, 42, ...)
- ► Why?
  - Predictability of lab conditions (= reproducibility of results)
  - Not used when models are applied productively

- Computers don't have random numbers (because they are deterministic machines)
- Real randomness can only achieved with specialised hardware
- Computers: Pseudo-random numbers
  - $\blacktriangleright$  Start with a seed value  $x_0$
  - Function to generate  $x_{n+1}$  from  $x_n$ 
    - This is a research area on its own
- Possible seed value sources
  - Mouse movement (not really random ...)
  - Memory allocation
  - Fixed number (0, 23, 42, ...)

► Why?

- Predictability of lab conditions (= reproducibility of results)
- Not used when models are applied productively

_	Random random = new java.util.Random(5);
$\leq$	<pre>for (int i = 0; i &lt; 10; i++) {</pre>
	<pre>System.err.println(random.nextInt());</pre>
	}

- Computers don't have random numbers (because they are deterministic machines)
- ▶ Real randomness can only achieved with specialised hardware
- Computers: Pseudo-random numbers
  - Start with a seed value  $x_0$
  - Function to generate  $x_{n+1}$  from  $x_n$ 
    - This is a research area on its own
- Possible seed value sources
  - Mouse movement (not really random ...)
  - Memory allocation
  - Fixed number (0, 23, 42, ...)

► Why?

- Predictability of lab conditions (= reproducibility of results)
- Not used when models are applied productively

```
Random random = new java.util.Random (5);
for (int i = 0; i < 10; i++) {</pre>
  System.err.println(random.nextInt());
                                         1
                                   -1157408321
                                    758500184
                                    379066948
                                   -1667228448
                                   2099829013
                                    -236332086
                                    1983575708
                                    -745003013
                                    1926715444
                                    1836354642
```

### Git

#### Some files don't belong in the git repository

- target/, because it's compiled code
- Eclipse meta data files: .classpath, .project, .settings
- IntelliJ meta data files: .idea
- Netbeans files: ...

### Git

#### Some files don't belong in the git repository

- target/, because it's compiled code
- Eclipse meta data files: .classpath, .project, .settings
- IntelliJ meta data files: .idea
- Netbeans files: ...
- gitignore A file in a git repository
  - Lists file patterns that are to be ignored by git
  - List of patterns for various file types: https://github.com/github/gitignore

# Section 2

### **Decision Trees**

# Recap: The Machine-Learning Recipe

- 1. Your problem
  - What are its structural properties?
  - What do you believe/know is relevant information to make a prediction?
  - To what extent can humans solve it?
  - Do you have non-formalisable requirements?

# Recap: The Machine-Learning Recipe

- 1. Your problem
  - What are its structural properties?
  - What do you believe/know is relevant information to make a prediction?
  - To what extent can humans solve it?
  - Do you have non-formalisable requirements?
- 2. Pick an appropriate algorithm
  - There are many out there
  - Appropriate: It should fit to your problem!
  - Gather training/testing data

# Recap: The Machine-Learning Recipe

- 1. Your problem
  - What are its structural properties?
  - What do you believe/know is relevant information to make a prediction?
  - To what extent can humans solve it?
  - Do you have non-formalisable requirements?

#### 2. Pick an appropriate algorithm

- There are many out there
- Appropriate: It should fit to your problem!
- Gather training/testing data
- 3. Evaluate fairly
  - What's the performance of a realistic and optimised baseline?
  - What's an appropriate evaluation metric?
  - Report evaluation results including all settings and constraints

# **Decision Trees**

#### Prediction Model – Toy Example





#### Reiter

Prediction Model



- Each non-leaf node in the tree represents one feature
- Each leaf node represents a class label
- Each branch at this node represents one possible feature value
  - Number of branches = number of possible values for this feature

Prediction Model



- Each non-leaf node in the tree represents one feature
- Each leaf node represents a class label
- Each branch at this node represents one possible feature value
  - Number of branches = number of possible values for this feature
- Make a prediction for x:
  - 1. Start at root node
  - 2. If it's a leaf node
    - assign the class label
  - 3. Else
    - Check node which feature is to be tested
    - Extract feature value
    - Follow corresponding branch
    - Go to 2

### Weka



### ID3

Quinlan (1986)

- Core idea: The tree represents splits of the training data
  - 1. Start with the full data set  $D_{\text{train}}$  as D
  - 72. If D only contains members of a single class:
    - Done, the node is labeled with the class
    - 3. Else if there are no more features to select:
      - Done, the node is labeled with the majority class
    - 4. Else if D is empty  $\sim$ 
      - Done, the node is labeled with the majority class of its parent
    - 5. Else:
      - Select a feature  $f_i$
      - Extract feature values of all instances in D
      - Split the data set according to  $f_i: D = D_v \cup D_w \cup D_u \dots$
      - For each subset as D, go back to 2

Remaining question: How to select features?

- What is a good feature?
  - One that maximizes homogeneity in the split data set

- What is a good feature?
  - One that maximizes homogeneity in the split data set

Justile que Falice 20

- "Homogeneity"
  - Increase ${ <math> \mathbf{A} \mathbf{A} \mathbf{A} \mathbf{O}$  } = {  $\mathbf{O}$  }  $\cup { \mathbf{A} \mathbf{A} \mathbf{A} }$

 $\{ \spadesuit \spadesuit \spadesuit \heartsuit \} = \{ \blacklozenge \} \cup \{ \spadesuit \spadesuit \heartsuit \}$ 

No increase

1

6

- What is a good feature?
  - One that maximizes homogeneity in the split data set
- "Homogeneity"
  - Increase {♠♠♠♡} = {♡} ∪ {♠♠♠} ← better split!
    No increase

$$\{ \bigstar \bigstar \bigstar \heartsuit \} = \{ \bigstar \} \cup \{ \bigstar \bigstar \heartsuit \}$$

Homogeneity: Entropy/information

Shannon (1948)

- What is a good feature?
  - One that maximizes homogeneity in the split data set
- "Homogeneity"
  - ▶ Increase  $\{ \spadesuit \spadesuit \spadesuit \heartsuit \} = \{ \heartsuit \} \cup \{ \clubsuit \spadesuit \spadesuit \} \leftarrow \text{better split!}$
  - ► No increase
    - $\{ \bigstar \spadesuit \clubsuit \heartsuit \} = \{ \clubsuit \} \cup \{ \clubsuit \clubsuit \heartsuit \}$
- Homogeneity: Entropy/information

Shannon (1948)

- ▶ Rule: Always select the feature with the highest *information gain* (IG)
  - (= the highest reduction in entropy = the highest increase in homogeneity)

Feature Selection Entropy (Shannon, 1948)

number of classes present in Xrelative frequency of the class  $H(X) = -\sum_{i=1}^{n} p(x_i) \log p(x_i)$ 



# Feature Selection (2)



# Feature Selection (3)

#### Example



20 / 30



C4.5 Quinlan (1993)

ID3, but with modifications:

- Continuous attributes
- Data with missing values
- Pruning the tree

### **Continuous Attributes**

The training cases T are first sorted on the values of the attribute A being considered. There are only a finite number of these values, so let us denote them in order as  $\{v_1, v_2, \ldots, v_m\}$ . Any threshold value lying between  $v_i$  and  $v_{i+1}$  will have the same effect of dividing the cases into those whose value of the attribute A lies in  $\{v_1, v_2, \ldots, v_i\}$  and those whose value is in  $\{v_{i+1}, v_{i+2}, \ldots, v_m\}$ . There are thus only m-1 possible splits on A, all of which are examined. (Quinlan, 1993, p. 25)

### **Continuous Attributes**

The training cases T are first sorted on the values of the attribute A being considered. There are only a finite number of these values, so let us denote them in order as  $\{v_i, v_2, \ldots, v_m\}$ . Any threshold value lying between  $v_i$  and  $v_{i+1}$  will have the same effect of dividing the cases into those whose value of the attribute A lies in  $\{v_1, v_2, \ldots, v_i\}$  and those whose value is in  $\{v_{i+1}, v_{i+2}, \ldots, v_m\}$ . There are thus only m-1 possible splits on A, all of which are examined. (Quinlan, 1993, p. 25)

I.e., we use all possible threshold values to split continuous features in two sub sets, and select the one with the highest information gain as usual.

Missing Values Quinlan (1993, pp. 28 f.)

- Calculate information gain as before, but
  - Take into account only instances whose attribute values are known
  - Weight information gain by the portion of known values

In my experience, almost all decision trees can benefit from simplification. (Quinlan, 1993, p. 36)

- Decision trees tend to overfit (e.g., with an 'id feature')
- Pruning the tree prevents that
- C4.5: Post-training pruning
  - Instead of stopping to expand the tree during training

In my experience, almost all decision trees can benefit from simplification. (Quinlan, 1993, p. 36)

- Decision trees tend to overfit (e.g., with an 'id feature')
- Pruning the tree prevents that
- C4.5: Post-training pruning
  - Instead of stopping to expand the tree during training
- General idea:
  - ► For each non-leaf node: Check if pruning reduces error rate
  - Error rate can't accuracy measured on training data, because all tree changes decrease accuracy

Pessimistic Pruning

- Pruning decision based on confidence intervals
  - For each node, we can calculate the observed error rate f
  - We are interested in the real error rate
  - Given a level of confidence, calculate the possible worst case (= upper bound) error rate
  - If the pruned tree has better worst case error rate, prune it

Pessimistic Pruning

- Pruning decision based on confidence intervals
  - $\blacktriangleright$  For each node, we can calculate the observed error rate f
  - $\blacktriangleright$  We are interested in the real error rate p
  - ▶ Given a level of confidence, calculate the possible worst case (= upper bound) error rate
  - If the pruned tree has better worst case error rate, prune it
- > Two pruning operations. Selection based on conf. intervals
  - Replace subtree with leaf
  - Replace subtree with one of its branches

Pessimistic Pruning

- Pruning decision based on confidence intervals
  - $\blacktriangleright$  For each node, we can calculate the observed error rate f
  - $\blacktriangleright$  We are interested in the real error rate p
  - ▶ Given a level of confidence, calculate the possible worst case (= upper bound) error rate
  - If the pruned tree has better worst case error rate, prune it
- > Two pruning operations. Selection based on conf. intervals
  - Replace subtree with leaf
  - Replace subtree with one of its branches

#### Reduced Error Pruning

- ► An alternative pruning strategy, not used by C4.5
- Prune by estimated errors on held-out data set

#### Options



#### Options

- C4.5 pruning options
  - -C <confidence level>: Pruning confidence level
- Reduced error pruning
  - ► -R: Reduced error pruning
  - -N <number of folds>

#### Options

- C4.5 pruning options
  - -C <confidence level>: Pruning confidence level
- Reduced error pruning
  - ► -R: Reduced error pruning
  - -N <number of folds>
- Other pruning options
  - ► -ण: Turn (all) pruning off
  - ▶ ③ When pruning, don't replace subtrees with branches, only with leafs
  - M<number of instances>: pre-pruning. Only split it at least <number of instances> remain in a subset

#### Options

- C4.5 pruning options
  - -C <confidence level>: Pruning confidence level
- Reduced error pruning
  - ► -R: Reduced error pruning
  - -N <number of folds>
- Other pruning options
  - ► -v: Turn (all) pruning off
  - ▶ -s: When pruning, don't replace subtrees with branches, only with leafs
  - -M <number of instances>: pre-pruning. Only split it at least <number of instances> remain in a subset
- Other options
  - ► -B: Only binary splits
  - -doNotMakeSplitPointActualValue: Do not enforce split points of numeric features to be actual values

Source: https://weka.8497.n7.nabble.com/Details-of-J48-pruning-parameters-td42456.html

Reiter

#### Tree Classifiers

# J48 Pruning Options

#### Example

- Data set: German credit
- Setup: 10-fold cross validation

# J48 Pruning Options

#### Example

- Data set: German credit
- Setup: 10-fold cross validation

Setting	Weighted F1	Tree depth
No pruning	0.718	10
C4.5 pruning	0.751	) 70
Reduced error pruning	0.719	8)

# Random Forest

Breiman (2001)

- ► A tree-based classifier that performs well on many tasks
  - Well suited for limited/small data set sizes

Random Forest Breiman (2001)



- A tree-based classifier that performs well on many tasks
  - Well suited for limited/small data set sizes

Bagging

- Generate variants of the training set by sampling with replacement from it
- Train a decision tree on each variant
- Create a prediction with all trees, and return the average / majority

# Random Forest

Breiman (2001)

- A tree-based classifier that performs well on many tasks
  - Well suited for limited/small data set sizes

Bagging

- Generate variants of the training set by sampling with replacement from it
- Train a decision tree on each variant
- Create a prediction with all trees, and return the average / majority

#### Feature bagging

- At each split, choose between a random subset of features
- Avoids very similar trees, if one feature is a strong predictor

Weka

weka.classifiers.trees.RandomForest

Summary

- Classification algorithm
- Built around trees, recursive learning and prediction

Pros

- Highly transparent
- Reasonably fast
- Dependencies between features can be incorporated into the model

Cons

- No pairwise dependencies
- May lead to overfitting but pruning and bagging help
- Variants exist (e.g., CART)
- Random forest

Exercise 5

# Exercise 5