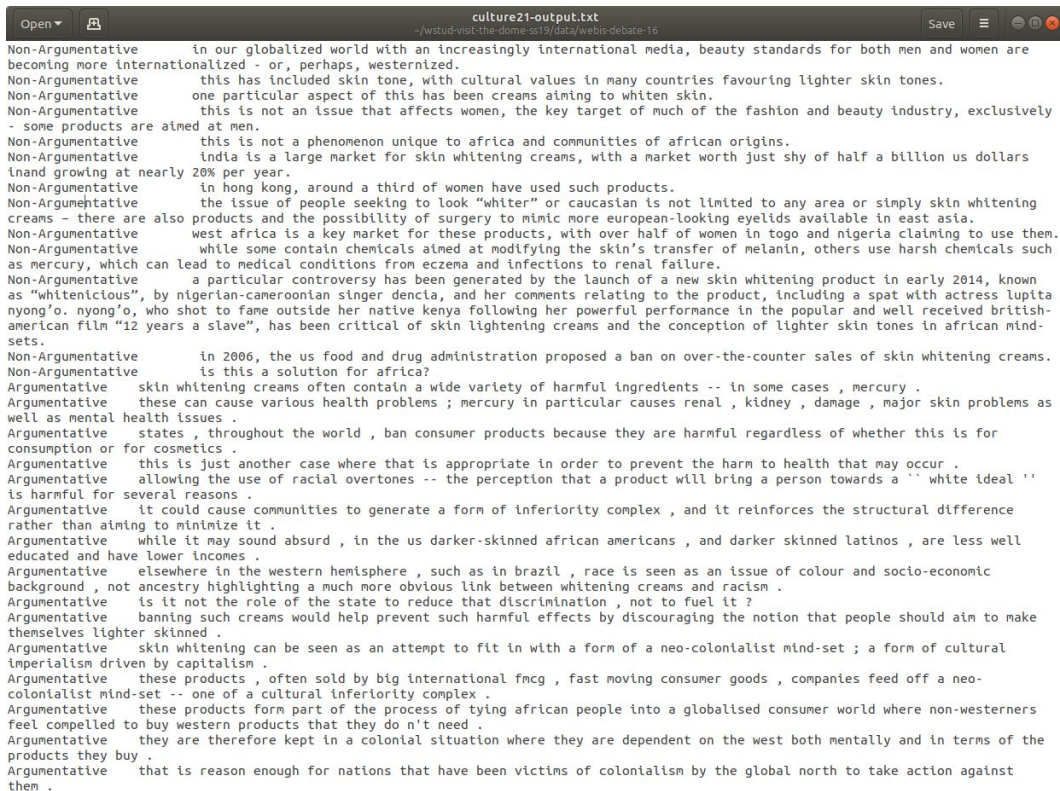


Program with UIMA Framework

Annie Lukas - 2019.05.06

Text to XMI files

- Each line contains one statement, with Argumentative Discourse Unit Type (ADU)
- To convert the dataset in UIMA format (XMI) with annotated ADU



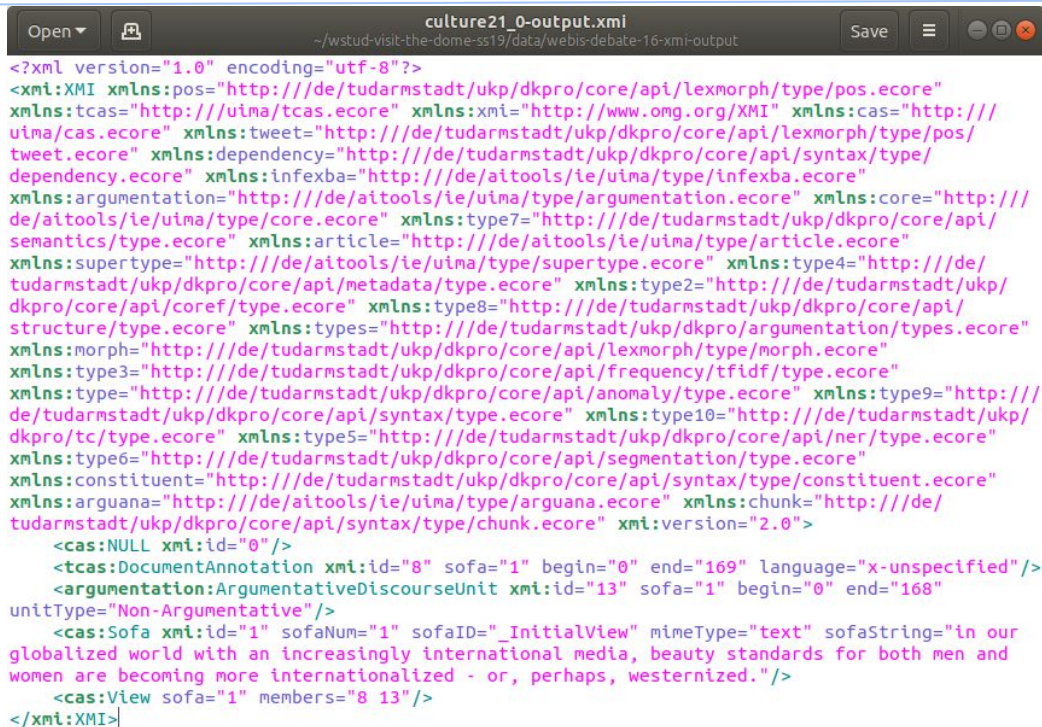
Text to XMI Converter

- Created a Text to XMI Converter (used example), where defined:
- Input directory
- Output directory
- Collection Reader: Plain Text Reader
- AE: Dummy

```
public class Text2XMIConverter {  
  
    // -----  
    // PARAMETERS  
    // -----  
  
    /**  
     * The path of root directory of the files to be processed.  
     */  
    private static final String INPUT_COLLECTION_DIR =  
        "data/webis-debate-16/";  
  
    /**  
     * The path of the XMI file of the collection reader to be used to iterate  
     * over all files to be processed.  
     */  
    private static final String COLLECTION_READER_PATH =  
        "../aitools4-ie-uima/conf/uima-descriptors/collection-readers/UIMAPlainTextReader.xml";  
  
    /**  
     * The path of the XMI file of the analysis engine to be used to process  
     * the files.  
     */  
    private static final String ANALYSIS_ENGINE_PATH =  
        "../aitools4-ie-uima/conf/uima-descriptors/primitive-AEs/template/"  
        + "DummyAnalysisEngine.xml";  
  
    /**  
     * The path of the directory where the XMI files shall be written to  
     */  
    private static final String OUTPUT_COLLECTION_DIR =  
        "data/webis-debate-16-xmi-output/";  
}
```

XMI Output

- Split each line and extract the content to be document text
- Annotated the statement with corresponding ADU type



```
Open  [icon] culture21_0-output.xml Save [menu] [window controls]
~/wstud-visit-the-dome-ss19/data/webis-debate-16-xmi-output

<?xml version="1.0" encoding="utf-8"?>
<xmi:XMI xmlns:pos="http://de/tudarmstadt/ukp/dkpro/core/api/lexmorph/type/pos.ecore"
xmlns:tcas="http://uima/tcas.ecore" xmlns:xmi="http://www.omg.org/XMI" xmlns:cas="http://
uima/cas.ecore" xmlns:tweet="http://de/tudarmstadt/ukp/dkpro/core/api/lexmorph/type/pos/
tweet.ecore" xmlns:dependency="http://de/tudarmstadt/ukp/dkpro/core/api/syntax/type/
dependency.ecore" xmlns:infexba="http://de/aitools/ie/uima/type/infexba.ecore"
xmlns:argumentation="http://de/aitools/ie/uima/type/argumentation.ecore" xmlns:core="http://
de/aitools/ie/uima/type/core.ecore" xmlns:type7="http://de/tudarmstadt/ukp/dkpro/core/api/
semantics/type.ecore" xmlns:article="http://de/aitools/ie/uima/type/article.ecore"
xmlns:supertype="http://de/aitools/ie/uima/type/supertype.ecore" xmlns:type4="http://de/
tudarmstadt/ukp/dkpro/core/api/metadata/type.ecore" xmlns:type2="http://de/tudarmstadt/ukp/
dkpro/core/api/coref/type.ecore" xmlns:type8="http://de/tudarmstadt/ukp/dkpro/core/api/
structure/type.ecore" xmlns:types="http://de/tudarmstadt/ukp/dkpro/argumentation/types.ecore"
xmlns:morph="http://de/tudarmstadt/ukp/dkpro/core/api/lexmorph/type/morph.ecore"
xmlns:type3="http://de/tudarmstadt/ukp/dkpro/core/api/frequency/tfidf/type.ecore"
xmlns:type="http://de/tudarmstadt/ukp/dkpro/core/api/anomaly/type.ecore" xmlns:type9="http://
de/tudarmstadt/ukp/dkpro/core/api/syntax/type.ecore" xmlns:type10="http://de/tudarmstadt/ukp/
dkpro/tc/type.ecore" xmlns:type5="http://de/tudarmstadt/ukp/dkpro/core/api/ner/type.ecore"
xmlns:type6="http://de/tudarmstadt/ukp/dkpro/core/api/segmentation/type.ecore"
xmlns:constituent="http://de/tudarmstadt/ukp/dkpro/core/api/syntax/type/constituent.ecore"
xmlns:arguana="http://de/aitools/ie/uima/type/arguana.ecore" xmlns:chunk="http://de/
tudarmstadt/ukp/dkpro/core/api/syntax/type/chunk.ecore" xmi:version="2.0">
  <cas:NULL xmi:id="0"/>
  <tcas:DocumentAnnotation xmi:id="8" sofa="1" begin="0" end="169" language="x-unspecified"/>
  <argumentation:ArgumentativeDiscourseUnit xmi:id="13" sofa="1" begin="0" end="168"
unitType="Non-Argumentative"/>
  <cas:Sofa xmi:id="1" sofaNum="1" sofaID="_InitialView" mimeType="text" sofaString="in our
globalized world with an increasingly international media, beauty standards for both men and
women are becoming more internationalized - or, perhaps, westernized."/>
  <cas:View sofa="1" members="8 13"/>
</xmi:XMI>
```

Simple Pipeline

- Input: XMI Files
- Output: Statement with annotated token and part of speech tag
- Collection Reader: UIMA Annotation File Reader

```
public class CollectionProcessor {  
  
    // -----  
    // PARAMETERS  
    // -----  
  
    /**  
     * The path of root directory of the training files to be processed.  
     */  
    private static final String COLLECTION_ROOT_DIR =  
        "data/webis-debate-16-xmi-output/";  
  
    /**  
     * The path of the XMI file of the collection reader to be used to iterate  
     * over all files to be processed.  
     */  
    private static final String COLLECTION_READER_PATH =  
        "../aitools4-ie-uima/conf/uima-descriptors/collection-readers/"  
        + "UIMAAnnotationFileReader.xml";  
  
    /**  
     * The path of the XMI file of the analysis engine to be used to process  
     * the files.  
     */  
    private static final String ANALYSIS_ENGINE_PATH =  
        "conf/uima-descriptors/aggregate-AEs/"  
        + "AggregateAE.xml";  
  
    /**  
     * The path of the directory where the XMI files shall be written to  
     */  
    private static final String OUTPUT_COLLECTION_DIR =  
        "data/webis-debate-16-xmi-output-token-partofspeech/";  
}
```


Aggregate Analysis Engine

- Contains
 - Sentence Splitter
 - Tokenizer
 - Lemma and Part of Speech Tagger

```
<delegateAnalysisEngine key="InfexBASentenceSplitter">  
  <import location="../../../aitools4-ie-uima/conf/uima-descriptors/  
</delegateAnalysisEngine>  
<delegateAnalysisEngine key="InfexBATokenizer">  
  <import location="../../../aitools4-ie-uima/conf/uima-descriptors/  
</delegateAnalysisEngine>  
<delegateAnalysisEngine key="TT4jLemmaAndPartOfSpeechTagger">  
  <import location="../../../aitools4-ie-uima/conf/uima-descriptors/
```

Infex BA Sentence Splitter

- Developed in the Infex BA project and improve in the Argument Analysis project later on
- Does not require any input annotations and produces sentence annotations

Infex BA Tokenizer

- Rather effective but still efficient rule-based tokenizer that was developed in the Infex BA project and improved in the Argument Analysis project later on
- Requires sentence annotations and produces token annotations

Part of Speech Tag

- POS tagging is process of marking up a word in a text as a corresponding to a part of speech, based on its definition and its context; e.g identification of words as nouns, verbs, adjectives, adverbs

TAG	DESCRIPTION	EXAMPLE
CC	conjunction, coordinating	<i>and, or, but</i>
CD	cardinal number	<i>five, three, 13%</i>
DT	determiner	<i>the, a, these</i>
EX	existential there	<i>there were six boys</i>
FW	foreign word	<i>mais</i>
IN	conjunction, subordinating or preposition	<i>of, on, before, unless</i>
JJ	adjective	<i>nice, easy</i>
JJR	adjective, comparative	<i>nicer, easier</i>
JJS	adjective, superlative	<i>nicest, easiest</i>
LS	list item marker	
MD	verb, modal auxiliary	<i>may, should</i>
NN	noun, singular or mass	<i>tiger, chair, laughter</i>
NNS	noun, plural	<i>tigers, chairs, insects</i>
NNP	noun, proper singular	<i>Germany, God, Alice</i>
NNPS	noun, proper plural	<i>my, your, our</i>
PDT	predeterminer	<i>extremely, loudly, hard</i>
POS	possessive ending	<i>better</i>
PRP	pronoun, personal	<i>best</i>
PRPS	pronoun, possessive	<i>about, off, up</i>
RB	adverb	<i>%</i>
RBR	adverb, comparative	<i>what to do?</i>
RBS	adverb, superlative	<i>oh, oops, gosh</i>
RP	adverb, particle	<i>think</i>
SYM	symbol	<i>she thinks</i>
TO	infinitival to	<i>I think</i>
UH	interjection	<i>they thought</i>
VB	verb, base form	<i>a sunken ship</i>
VBZ	verb, 3rd person singular present	<i>thinking is fun</i>
VBP	verb, non-3rd person singular present	<i>which, whatever, whichever</i>
VBD	verb, past tense	<i>what, who, whom</i>
VBN	verb, past participle	<i>whose, whosever</i>
VBG	verb, gerund or present participle	<i>where, when</i>
WDT	wh-determiner	<i>..?*</i>
WP	wh-pronoun, personal	<i>,</i>
WPS	wh-pronoun, possessive	<i>:</i>
WRB	wh-adverb	<i>(</i>
.	punctuation mark, sentence closer	<i>)</i>
,	punctuation mark, comma	
:	punctuation mark, colon	
(contextual separator, left paren	
)	contextual separator, right paren	

Output XML Files

1.	CC	Coordinating conjunction
2.	CD	Cardinal number
3.	DT	Determiner
4.	EX	Existential <i>there</i>
5.	FW	Foreign word
6.	IN	Preposition or subordinating conjunction
7.	JJ	Adjective
8.	JJR	Adjective, comparative
9.	JJS	Adjective, superlative
10.	LS	List item marker
11.	MD	Modal
12.	NN	Noun, singular or mass
13.	NNS	Noun, plural
14.	NNP	Proper noun, singular
15.	NNPS	Proper noun, plural
16.	PDT	Predeterminer
17.	POS	Possessive ending
18.	PRP	Personal pronoun
19.	PRP\$	Possessive pronoun
20.	RB	Adverb
21.	RBR	Adverb, comparative
22.	RBS	Adverb, superlative
23.	RP	Particle
24.	SYM	Symbol
25.	TO	to
26.	UH	Interjection
27.	VB	Verb, base form
28.	VBD	Verb, past tense
29.	VBG	Verb, gerund or present participle
30.	VBN	Verb, past participle
31.	VBP	Verb, non-3rd person singular present
32.	VBZ	Verb, 3rd person singular present
33.	WDT	Wh-determiner
34.	WP	Wh-pronoun
35.	WP\$	Possessive wh-pronoun
36.	WRB	Wh-adverb

```
Open  [icon] culture21_0-output.xml  Save [icon] [icon] [icon]
~/wstudies-the-dome-study-web-site/16-xml-output-token-partofspeech

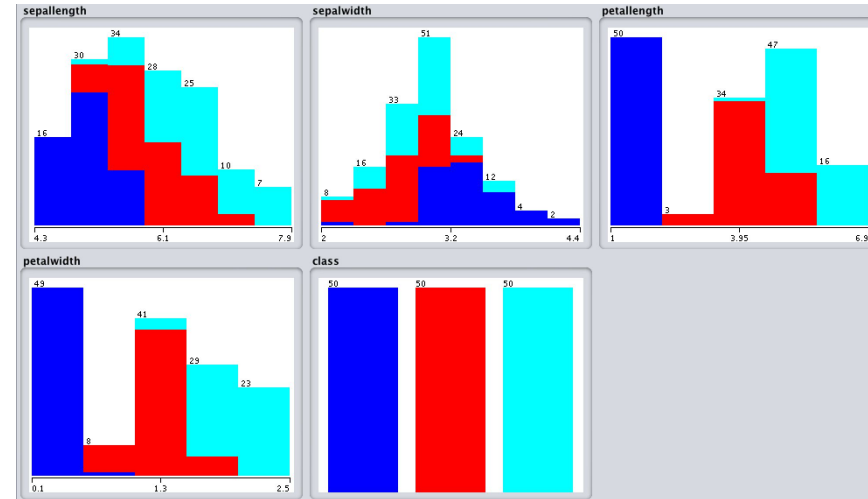
<?xml version="1.0" encoding="utf-8"?>
<xml:XML xmlns:pos="http://de/tudarnstadt/ukp/dkpro/core/api/lexmorph/type/pos.ecore" xmlns:tcas="http://uina/tcas.ecore"
xmlns:xmi="http://www.omg.org/XMI" xmlns:cas="http://uina/cas.ecore" xmlns:tweet="http://de/tudarnstadt/ukp/dkpro/core/api/lexmorph/
type/pos/tweet.ecore" xmlns:dependency="http://de/tudarnstadt/ukp/dkpro/core/api/syntax/type/dependency.ecore" xmlns:infexba="http://
de/atools/le/uina/type/infexba.ecore" xmlns:argumentation="http://de/atools/le/uina/type/argumentation.ecore" xmlns:core="http://de/
atools/le/uina/type/core.ecore" xmlns:type7="http://de/tudarnstadt/ukp/dkpro/core/api/semantics/type.ecore" xmlns:article="http://de/
atools/le/uina/type/article.ecore" xmlns:supertype="http://de/atools/le/uina/type/supertype.ecore" xmlns:type4="http://de/
tudarnstadt/ukp/dkpro/core/api/metadata/type.ecore" xmlns:type2="http://de/tudarnstadt/ukp/dkpro/core/api/coref/type.ecore"
xmlns:type8="http://de/tudarnstadt/ukp/dkpro/core/api/structure/type.ecore" xmlns:types="http://de/tudarnstadt/ukp/dkpro/argumentation/
types.ecore" xmlns:morph="http://de/tudarnstadt/ukp/dkpro/core/api/lexmorph/type/morph.ecore" xmlns:type3="http://de/tudarnstadt/ukp/
dkpro/core/api/frequency/trifid/type.ecore" xmlns:type="http://de/atools/le/uina/type/argumentation.ecore" xmlns:anomaly/type.ecore"
xmlns:type9="http://de/tudarnstadt/ukp/dkpro/core/api/syntax/type.ecore" xmlns:type18="http://de/tudarnstadt/ukp/dkpro/tc/type.ecore"
xmlns:type5="http://de/tudarnstadt/ukp/dkpro/core/api/ner/type.ecore" xmlns:type6="http://de/tudarnstadt/ukp/dkpro/core/api/
segmentation/type.ecore" xmlns:constituent="http://de/tudarnstadt/ukp/dkpro/core/api/syntax/type/constituent.ecore"
xmlns:arguana="http://de/atools/le/uina/type/arguana.ecore" xmlns:chunk="http://de/tudarnstadt/ukp/dkpro/core/api/syntax/type/
chunk.ecore" xmi:version="2.0">
  <cas:NULL xmi:id="0"/>
  <tcas:DocumentAnnotation xmi:id="1" sofa="11" begin="0" end="169" language="x-unspecified"/>
  <argumentation:ArgumentativeDiscourseUnit xmi:id="6" sofa="11" begin="0" end="168" unitType="Non-Argumentative"/>
  <core:SourceDocumentInformation xmi:id="18" sofa="11" begin="0" end="0" url="http://home/ciso478/wstudies-the-dome-ss19/data/websi-
debate-16-xml-output/culture21_0-output.xml" offsetInSource="0" documentSize="169" lastSegment="false"/>
  <core:Sentence xmi:id="26" sofa="11" begin="0" end="169"/>
  <core:Token xmi:id="38" sofa="11" begin="0" end="2" lemma="in" pos="IN"/>
  <core:Token xmi:id="41" sofa="11" begin="3" end="6" lemma="our" pos="PP$"/>
  <core:Token xmi:id="52" sofa="11" begin="7" end="17" lemma="globalize" pos="VVN"/>
  <core:Token xmi:id="63" sofa="11" begin="18" end="23" lemma="world" pos="NN"/>
  <core:Token xmi:id="74" sofa="11" begin="24" end="28" lemma="with" pos="IN"/>
  <core:Token xmi:id="85" sofa="11" begin="29" end="31" lemma="an" pos="DT"/>
  <core:Token xmi:id="96" sofa="11" begin="32" end="44" lemma="increasingly" pos="RB"/>
  <core:Token xmi:id="107" sofa="11" begin="45" end="58" lemma="international" pos="JJ"/>
  <core:Token xmi:id="118" sofa="11" begin="59" end="64" lemma="medium" pos="NN$"/>
  <core:Token xmi:id="129" sofa="11" begin="64" end="65" lemma="," pos=","/>
  <core:Token xmi:id="140" sofa="11" begin="66" end="72" lemma="beauty" pos="NN"/>
  <core:Token xmi:id="151" sofa="11" begin="73" end="82" lemma="standard" pos="NNS"/>
  <core:Token xmi:id="162" sofa="11" begin="83" end="86" lemma="for" pos="IN"/>
  <core:Token xmi:id="173" sofa="11" begin="87" end="91" lemma="both" pos="DT"/>
  <core:Token xmi:id="184" sofa="11" begin="92" end="95" lemma="man" pos="NNS"/>
  <core:Token xmi:id="195" sofa="11" begin="96" end="99" lemma="and" pos="CC"/>
  <core:Token xmi:id="206" sofa="11" begin="100" end="105" lemma="woman" pos="NNS"/>
  <core:Token xmi:id="217" sofa="11" begin="106" end="109" lemma="be" pos="VP"/>
  <core:Token xmi:id="228" sofa="11" begin="110" end="118" lemma="become" pos="VVC"/>
  <core:Token xmi:id="239" sofa="11" begin="119" end="123" lemma="more" pos="RBR"/>
  <core:Token xmi:id="250" sofa="11" begin="124" end="141" lemma="internationalize" pos="VVN"/>
  <core:Token xmi:id="261" sofa="11" begin="142" end="143" lemma="," pos=","/>
  <core:Token xmi:id="272" sofa="11" begin="144" end="146" lemma="or" pos="CC"/>
  <core:Token xmi:id="283" sofa="11" begin="146" end="147" lemma="," pos=","/>
  <core:Token xmi:id="294" sofa="11" begin="148" end="155" lemma="perhaps" pos="RB"/>
  <core:Token xmi:id="305" sofa="11" begin="155" end="156" lemma="," pos=","/>
  <core:Token xmi:id="316" sofa="11" begin="157" end="168" lemma="westernize" pos="VVN"/>
  <core:Token xmi:id="327" sofa="11" begin="168" end="169" lemma="." pos="SENT"/>
  <cas:Sofa xmi:id="11" sofaNum="1" sofaId="InitialView" mimeType="text" sofaString="In our globalized world with an increasingly
international media, beauty standards for both men and women are becoming more internationalized - or, perhaps, westernized."/>
  <cas:View sofa="11" members="1 6 18 26 30 41 52 63 74 85 96 107 118 129 140 151 162 173 184 195 206 217 228 239 250 261 272 283 294
305 316 327"/>
</xml:XML>
```

Weka

- An open source Java software that has a collection of machine learning algorithms for data mining tasks
- A powerful tool for understanding and visualizing machine learning algorithms on your local machine
- Contains tools for data preparation, classification, regression, clustering, and visualization
- Provides an easy way to apply many different algorithms to your data and see which one will give the best results

Visualise class distribution vs features

- Attribute information: sepal length, sepal width, petal length, petal width
- Class: Iris- Setosa, Versicolour, Virginica
- Petal length and petal width has high class correlation (0.95)
- Even class distribution (50 each)



Classify with Random Forest Classifier

- Random Forest is an ensemble learning algorithm that can be used for classification regression and other tasks. It works by constructing a multitude of decision trees at training time and outputting the predicted class.
- 10 fold cross validation: Form 10 test sets by splitting dataset into disjoint sets of similar size. Each time, train the classifier by 9 sets and test by the remaining set.

Summary of run information

=== Run information ===

```
Scheme:      weka.classifiers.trees.RandomForest -P 100 -I 100 -num-slots 1 -K 0 -M 1.0 -V 0.001 -S 1
Relation:    iris
Instances:   150
Attributes:  5
              sepallength
              sepalwidth
              petallength
              petalwidth
              class
Test mode:   10-fold cross-validation
```

Result

- Accuracy:
95.3%

=== Classifier model (full training set) ===

RandomForest

Bagging with 100 iterations and base learner

weka.classifiers.trees.RandomTree -K 0 -M 1.0 -V 0.001 -S 1 -do-not-check-capabilities

Time taken to build model: 0.14 seconds

=== Stratified cross-validation ===

=== Summary ===

Correctly Classified Instances	143	95.3333 %
Incorrectly Classified Instances	7	4.6667 %
Kappa statistic	0.93	
Mean absolute error	0.0408	
Root mean squared error	0.1621	
Relative absolute error	9.19 %	
Root relative squared error	34.3846 %	
Total Number of Instances	150	

Result (continue)

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	1.000	0.000	1.000	1.000	1.000	1.000	1.000	1.000	Iris-setosa
	0.940	0.040	0.922	0.940	0.931	0.896	0.991	0.984	Iris-versicolor
	0.920	0.030	0.939	0.920	0.929	0.895	0.991	0.982	Iris-virginica
Weighted Avg.	0.953	0.023	0.953	0.953	0.953	0.930	0.994	0.989	

=== Confusion Matrix ===

```

a  b  c  <-- classified as
50  0  0 | a = Iris-setosa
 0 47  3 | b = Iris-versicolor
 0  4 46 | c = Iris-virginica

```

		Truth	
		P	N
Hypothesis	P	TP (a)	FP (b)
	N	FN (c)	TN (d)

Precision:

$$\frac{a}{a+b}$$

Recall:

$$\frac{a}{a+c}$$

F-measure:

$$F_{\alpha} = \frac{1 + \alpha}{\frac{1}{precision} + \frac{\alpha}{recall}}$$

$\alpha = 1$
 $\alpha \in (0; 1)$
 $\alpha > 1$

harmonic mean
 favor precision over recall
 favor recall over precision