

Introduction to Machine Learning and Deep Learning

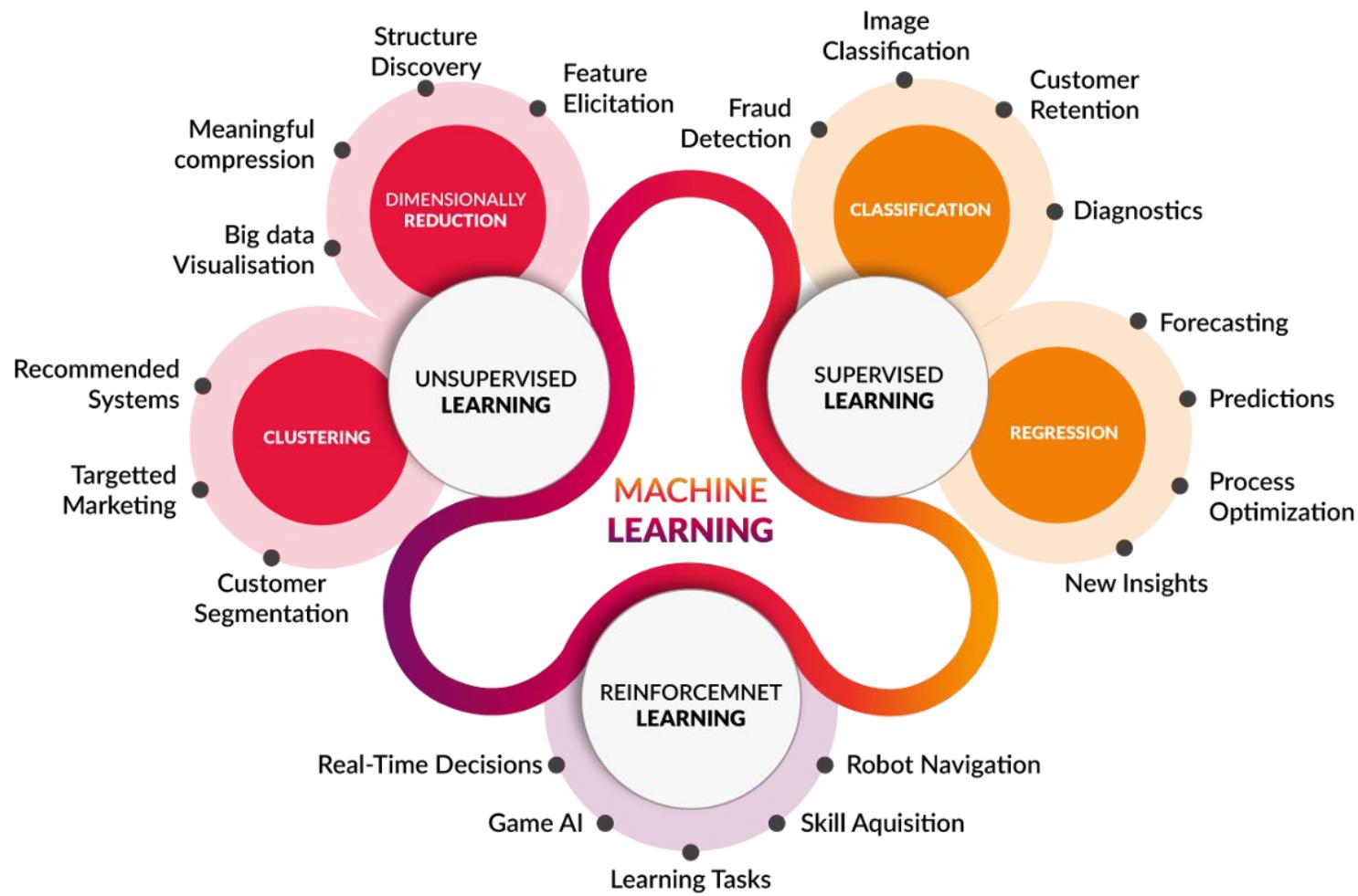
Dr Mahmoud El-Haj
SCC, Lancaster University



InfoLab21

Office C28, InfoLab21
m.el-haj@Lancaster.ac.uk

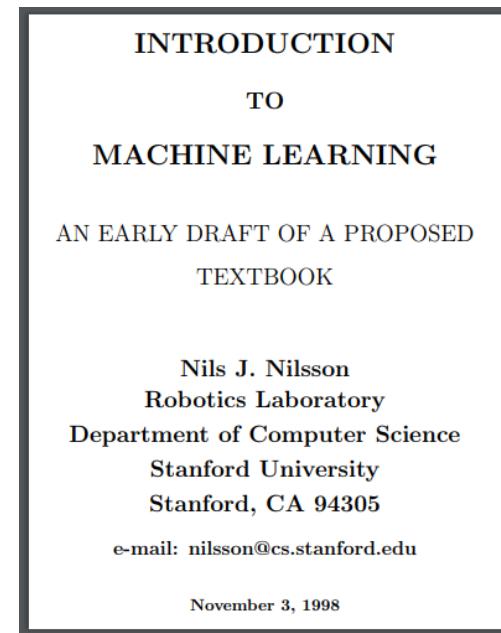
MACHINE LEARNING



RELEVANT BOOK FOR MACHINE LEARNING

- Introduction to Machine Learning: An Early Draft of a Proposed Textbook. N. J. Nilsson, Stanford.

<http://ai.stanford.edu/people/nilsson/mlbook.html>

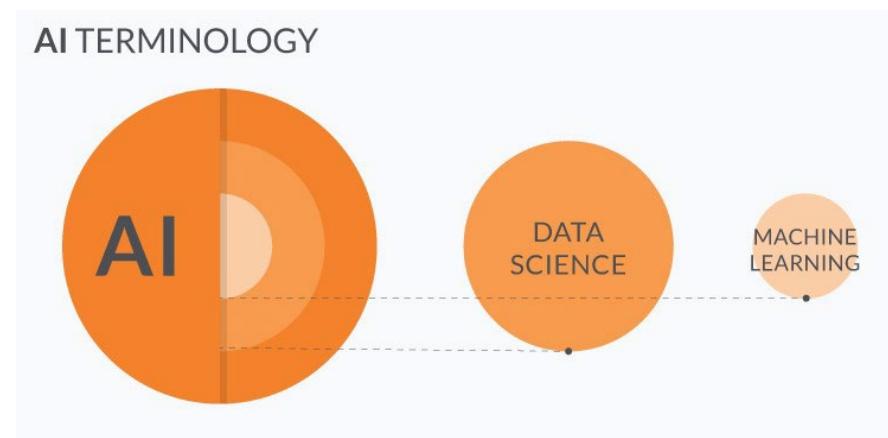


YOUR LEARNING OUTCOMES

- Understand fundamental Machine Learning concepts and current issues.
- Understand the difference between learning methods, algorithms and how to apply that on Automatic Text Classification.

MACHINE LEARNING – FIRST THINGS FIRST

- Machine:
Semi or fully automated device that magnifies human physical and/or mental capabilities in performing one or more operations.
- Where does NLP fit into AI?
Machine Learning is a current application of AI





Part 0: Introduction to Machine Learning

History
Definitions
Applications
Timeline
Philosophy



WHAT IS LEARNING?

- Learning:

“to gain knowledge,
or understanding of,
or skill in,
by study, instruction, or experience,”
and “modification of a behavioural tendency by experience.”

(Nils Nilsson)
- Definition: “improving performance in some task with experience”

(Tom Mitchell).

WHAT IS MACHINE LEARNING?

- broadly: a machine learns whenever it changes its structure, program, or data (based on its inputs or in response to external information) in such a manner that its expected future performance improves.
- *Arthur Samuel (1959): “Field of study that gives computers the ability to learn without being explicitly programmed”.*
- *He wrote a program that learnt to play checkers (BE: draughts)*

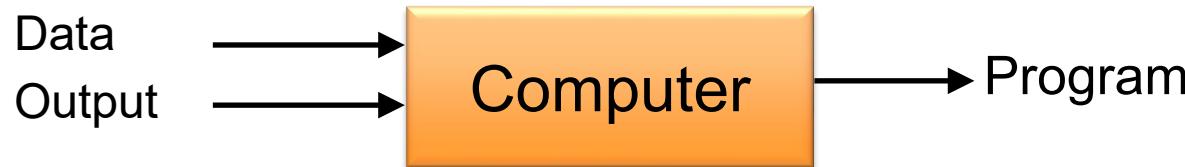


MACHINE LEARNING VS PROGRAMMING

Traditional Programming

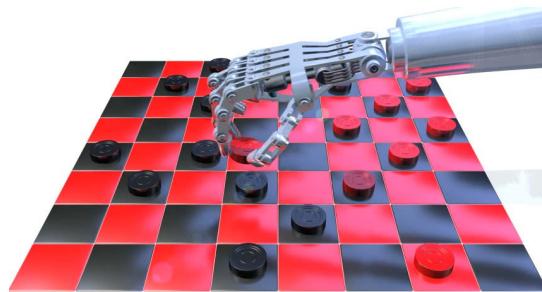


Machine Learning



WHAT IS MACHINE LEARNING?

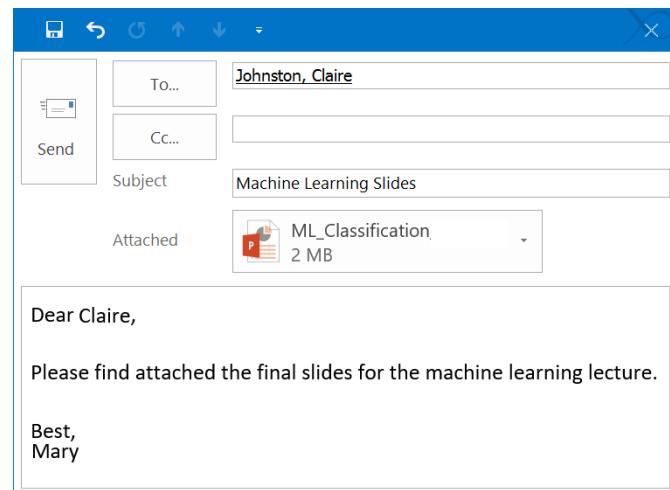
- Tom Mitchel (1998): a computer program is said to learn from experience E with respect to some task T and some performance measure P , if its performance on T , as measured by P , improves with experience E ;
- T : Task of playing checkers.
- E : Experience of having the program to play 10s of thousands of games against itself.
- P : Probability of winning the next game of checkers against a new opponent.



WHAT IS MACHINE LEARNING?

- Suppose your email application watches which emails you attach or don't attach files to, and based on that it learns how to better remind you if you forgot to attach a file. Can you identify **T**, **E** and **P** in this example?

- T** Classify whether an email needs or does not need an attachment
E Watching you attaching or not attaching files to emails
P The number of emails correctly classified as needs/does not need attachment
 None of the above- it's not a machine learning case.



WHY NOW?

- Huge amount of data
- Computational power

2018 *This Is What Happens In An Internet Minute*



TEXTUAL DATA

TEXTUAL DATA



APPLICATIONS OF MACHINE LEARNING

Machine learning is preferred approach to:

Speech recognition – Alexa, Siri, Google ...etc

Face Recognition – Facebook tag a friend suggestion

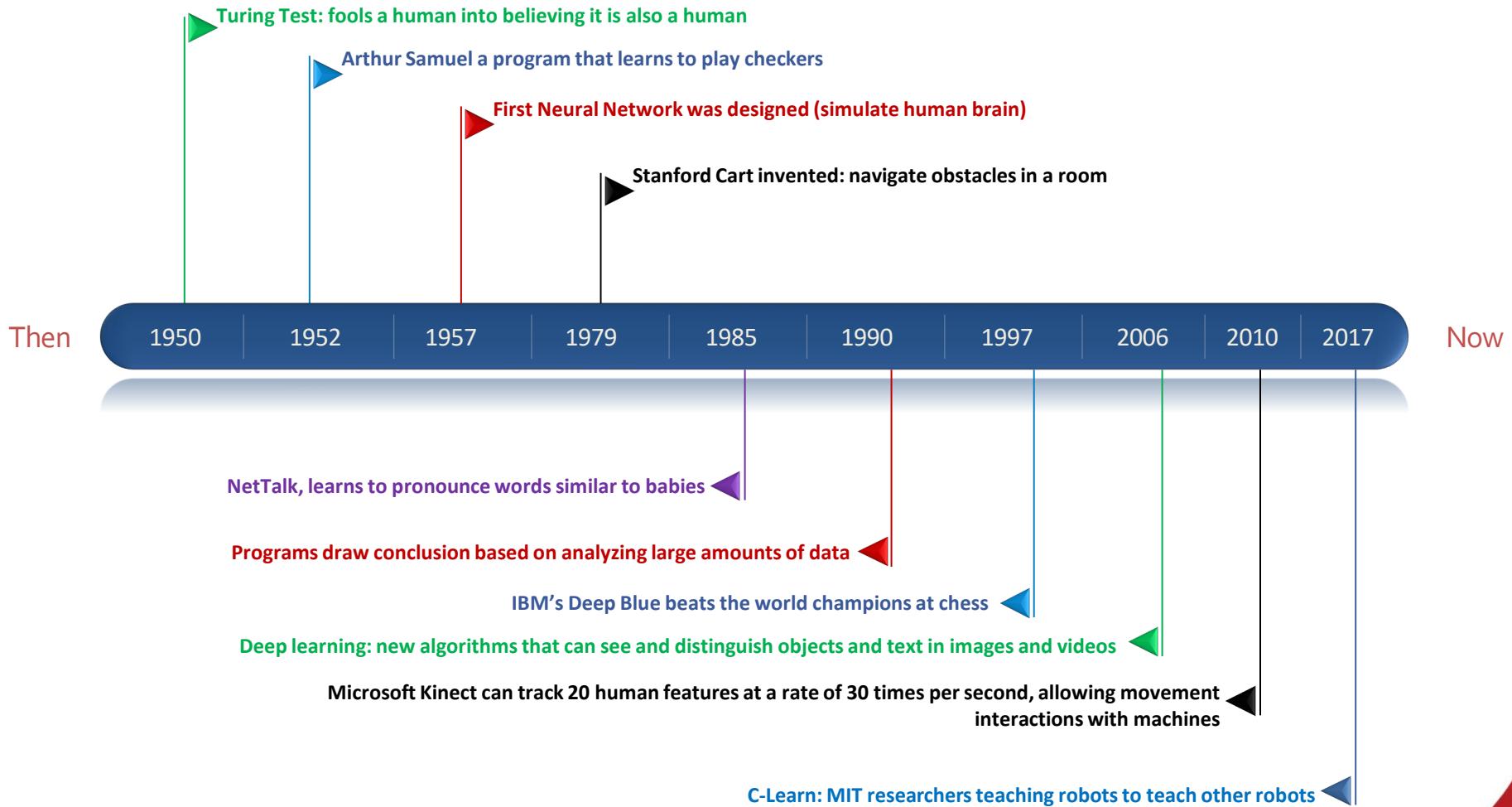
Fraud Detection - PayPal

Travel Shortest path - Google Maps, Uber

Recommendations – Netflix, Spotify

Etc.....

MACHINE LEARNING TIMELINE

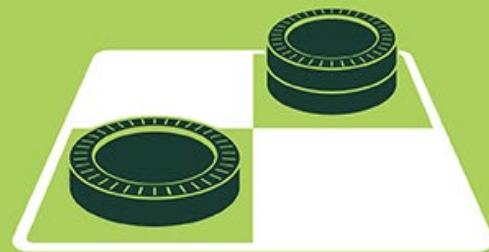




MACHINE LEARNING TIMELINE

ARTIFICIAL INTELLIGENCE

Early artificial intelligence stirs excitement.



1950's

1960's

1970's

1980's

1990's

2000's

2010's

MACHINE LEARNING

Machine learning begins to flourish.



DEEP LEARNING

Deep learning breakthroughs drive AI boom.





DEMO

NetTALK

<https://www.youtube.com/watch?v=gakJlr3GecE>

Stanford Cart

<https://www.youtube.com/watch?v=ypE64ZLwC5w>

IBM Deep Blue Short Documentary

<https://www.youtube.com/watch?v=NJarxpYyoFI>

WILL MACHINE LEARNING REPLACE Us?

ML will unlikely be able to replace every human job.

But! Machines have been slowly replacing human workers since the industrial revolution- 18th century!



MAJOR JOBS REPLACED WITH MACHINE LEARNING

- Assembly-line and factory workers.
- Bus, taxi, and truck drivers. Tesla / Uber
- Phone operators, telemarketers and receptionists.
- Cashiers. Self-checkout lines
- Bank tellers and clerks. ATMs, digital currency i.e. Bitcoin
- Packing, stockroom and warehouse moving. Amazon
- Prescription. Pharmacy robots
- Auto Pilots, auto bartenders, postal machines, surgeon robots, soldiers, travel agents,etc

MACHINE LEARNING PHILOSOPHY

- Limitations (relevant data, interpretation of results)
- Ethics: learning bias (high paid job ads for males)
- Discrimination: automatic decisions (credit scoring)

MACHINE LEARNING COMPETITIONS

Kaggle Competitions: Challenge yourself with real-world machine learning problems.

<https://www.kaggle.com/competitions>

A cartoon robot with a large head, wearing a white shirt and red tie, stands next to a list of Kaggle competitions. The competitions are listed in descending order of prize money:

Prize Money	Number of Teams
\$1,500,000	424 teams
\$1,200,000	3,780 teams
\$100,000	83 teams
\$50,000	1,349 teams

Passenger Screening Algorithm Challenge
Improve the accuracy of the Department of Homeland Security's threat recognition algorithms
Featured · 24 days to go · terrorism, image, object detection

Zillow Prize: Zillow's Home Value Prediction (Zestimate)
Can you improve the algorithm that changed the world of real estate?
Featured · 2 months to go · housing, real estate

Mercari Price Suggestion Challenge
Can you automatically suggest product prices to online sellers?
Featured · 3 months to go

Statoil/C-CORE Iceberg Classifier Challenge
Ship or iceberg, can you decide from space?
Featured · 2 months to go · weather, shipping, image, binary classification

Part 1: Machine Learning Fundamentals

Learning Process

Aim

Training and Testing

Cross-validation and Overfitting

Learning Algorithms

Introduction to NB, SVM and RF

Features Extraction

LEARNING PROCESS

- The process of learning begins with :
 1. Observations / data (examples, direct experience, instructions)
 2. Finding patterns (**e.g.**, those who buy cereal buy milk)
 3. Make better decisions (**e.g.**, place milk near cereal, but away from seafood)

AIM

- The primary aim is to allow the computers to learn (**training**) automatically and adjust actions accordingly to make decisions on future events (**testing**).

TRAINING

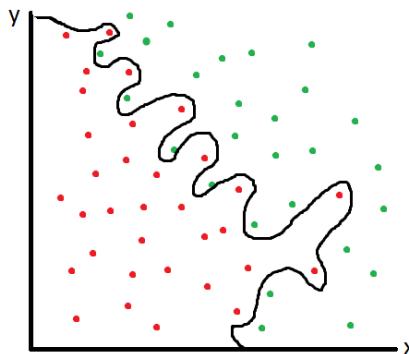
- Machine Learning Algorithms learn from data (**training data**).
- To find relationships, develop understanding, make decision and evaluate their confidence from the data given.
- The quality and quantity of the training data has as much to do with the success as the algorithms themselves.
- Training ends up with building a model.

TESTING

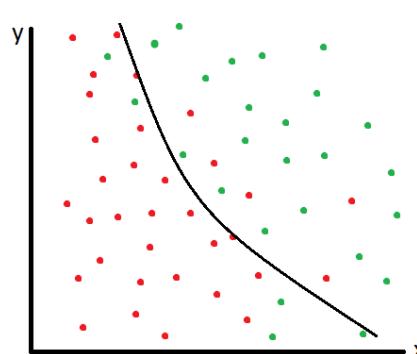
- To evaluate the trained model we use data that the model has never seen before (**testing data**)
- Usually the dataset is split into 90% training and 10% testing (or 70%, 30%) on different random folds (this is called cross validation).
- Cross validation helps avoiding overfitting.
- Best practice is to test your model on a testing dataset similar to the training one as otherwise you risk misclassification errors.

OVERFITTING

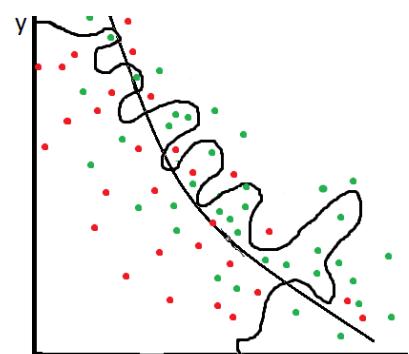
- Overfitting: your model fits with a particular training dataset but fails to fit additional data or predict future observations reliably.
- Avoid creating a model that is tailored only to your training data (e.g., assuming all fruit are round where training dataset includes only apples and oranges).
- ML works on predicting **y** instances given **x** instances:



Overfitting



Balanced



Unseen data

MACHINE LEARNING ALGORITHMS

Algorithms simply put, they learn from and make predictions on data.

Such algorithms overcome weaknesses of static (hard coded) program instructions by making data-driven predictions or decisions through building a model from sample inputs.

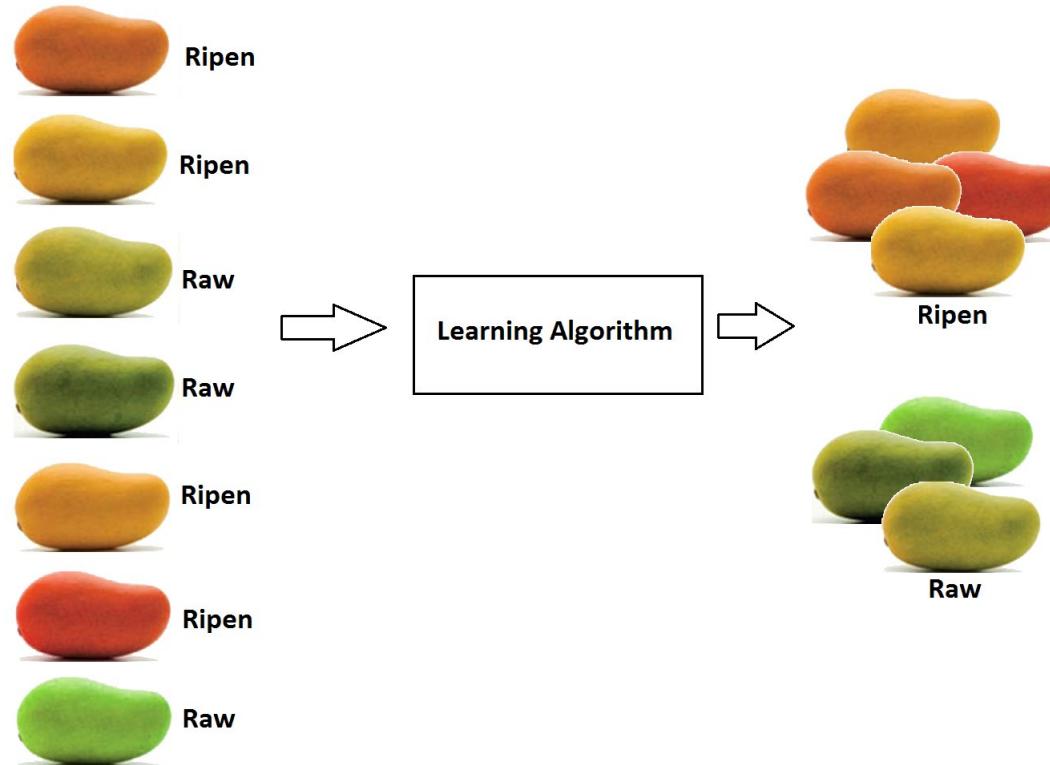
MACHINE LEARNING ALGORITHMS

- **Supervised**
 - Naïve Bayes, Decision Trees (Random Forest),...
- **Unsupervised**
 - Clustering: K-means.
- **Semi-supervised**
 - A combination of both

Note that both theory and empirical evidence provide little systematic guidance about which method is likely to dominate; performance depends on particular setting and dataset so NLP researchers often try multiple methods and then pick the best performer in the particular circumstances

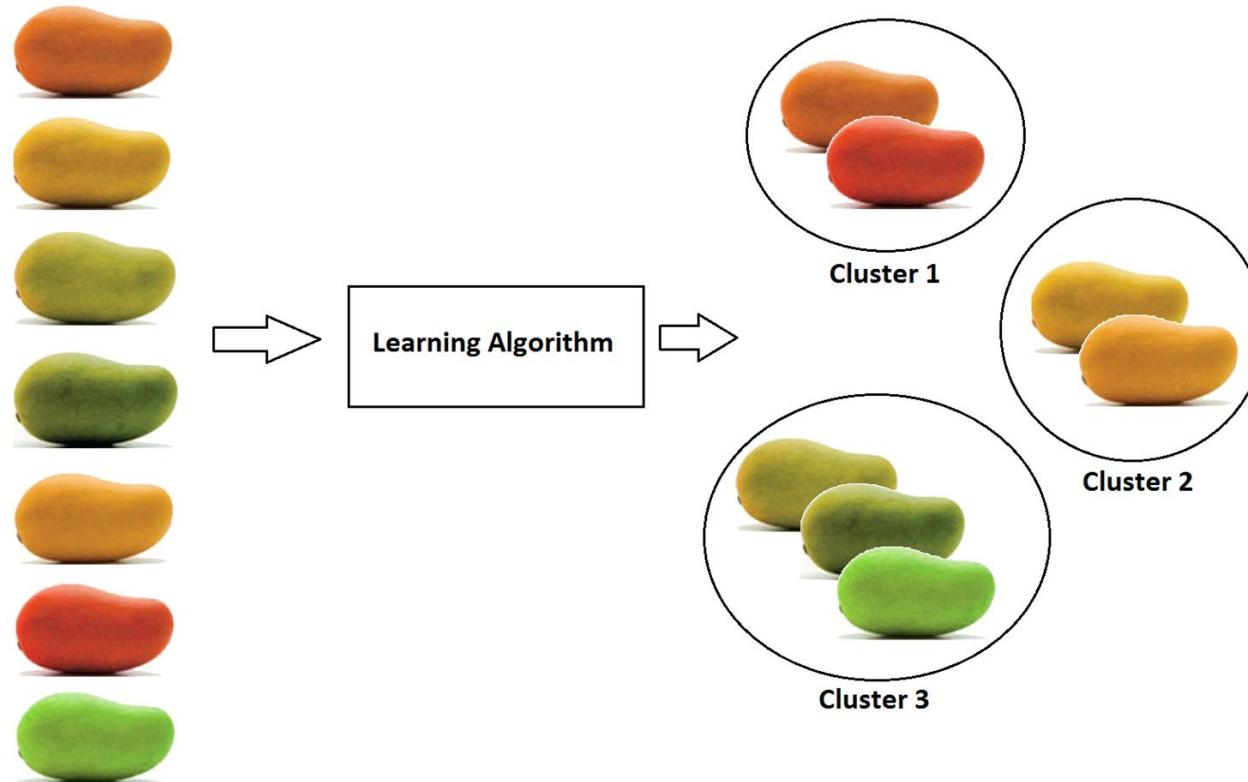
SUPERVISED MACHINE LEARNING

- Apply what learnt in the past to new data
- Using labelled examples to predict future events



UNSUPERVISED MACHINE LEARNING

- Training data is unlabelled.
- Describe a hidden structure from unlabelled data.



SEMI-SUPERVISED MACHINE LEARNING

- Falls somewhere between supervised and unsupervised
- Using both labelled and unlabelled data
- Why semi?

Acquiring labelled data is costly.

FEATURES AND FEATURES-EXTRACTION

- Feature: an individual measurable property or characteristic of a phenomenon being observed.



- Shape, colour, texture, taste, price, origin...

FEATURES REDUCTION

- Having irrelevant features in your data can decrease the accuracy of many models.
- **Reduces Overfitting:** Less redundant data means less opportunity to make decisions based on noise.
- **Improves Accuracy:** Less misleading data means modelling accuracy improves.
- **Reduces Training Time:** Less data means that algorithms train faster.



NAÏVE BAYES (NB)

Set of algorithms based on Bayes Theorem (**probability** of an event based on knowledge of conditions that might be related to the event). $P(\text{Play} | \text{Sunny, clear sky})$

Why naïve? Because it treats features independently.

A fruit may be considered to be an apple if it is **red, round**, and about 3" in **diameter**.

NB classifier considers each of those “features” to contribute independently to the probability that the fruit is an apple, regardless of any correlations between features. Which is not always the case.

$P(\text{Apple} | \text{red, round, 3"}\text{diam})$

NAÏVE BAYES (NB)

Advantages

- It's relatively simple to understand and build
- It's easily trained, even with a small dataset
- It's fast!
- It's not sensitive to irrelevant features (handles missing values)

Disadvantages

- It assumes every feature is independent, which isn't always the case (rain vs clouds).

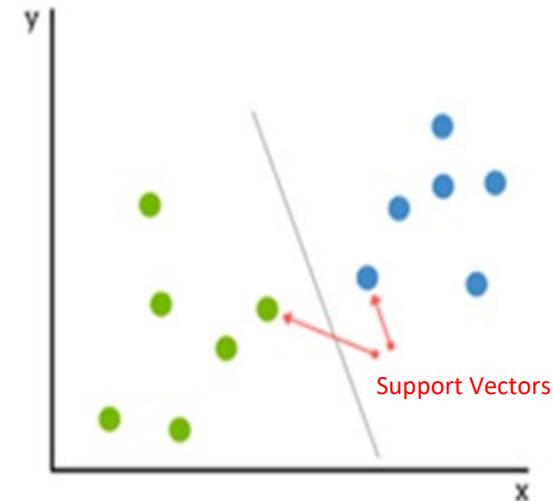
SUPPORT VECTOR MACHINE (SVM)

Supervised classifier. Finds a **hyperplane** that best divides a dataset into two classes.

Think of a hyperplane as a line that linearly separates and classifies a set of data.

Support vectors are the data points nearest to the hyperplane.

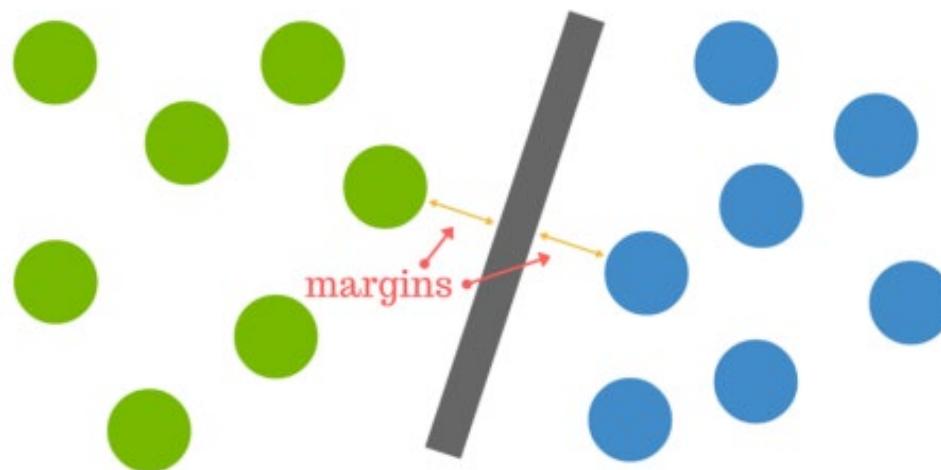
Points that are close to the opposing class.



SUPPORT VECTOR MACHINE (SVM)

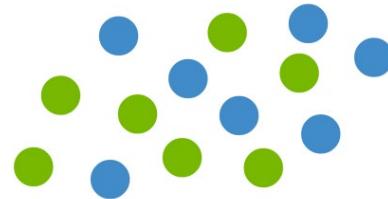
The distance between the hyperplane and the nearest data point from either set is known as the **margin**.

The goal is to choose a hyperplane with the greatest possible margin between the hyperplane and any point within the training set.

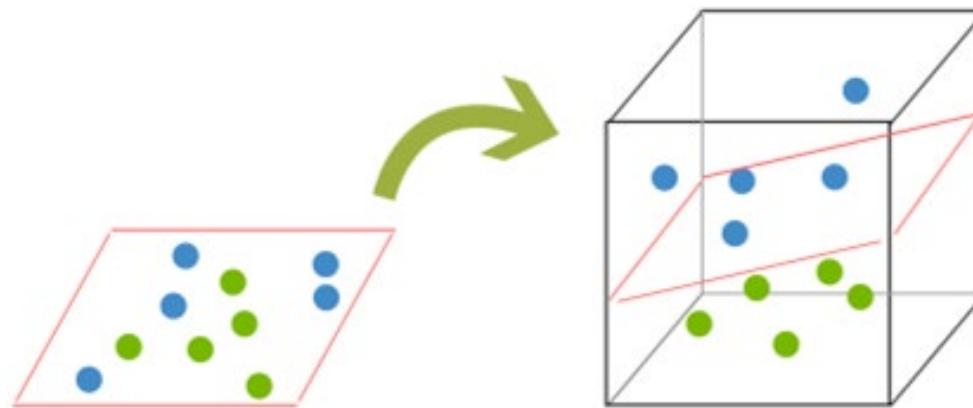


SUPPORT VECTOR MACHINE (SVM)

But what happens when there is no clear hyperplane?



In order to classify a dataset like the one above it's necessary to move away from a 2d view of the data to a 3d view



SUPPORT VECTOR MACHINE (SVM)

Advantages

- Accuracy
- Works well on smaller cleaner datasets

Disadvantages

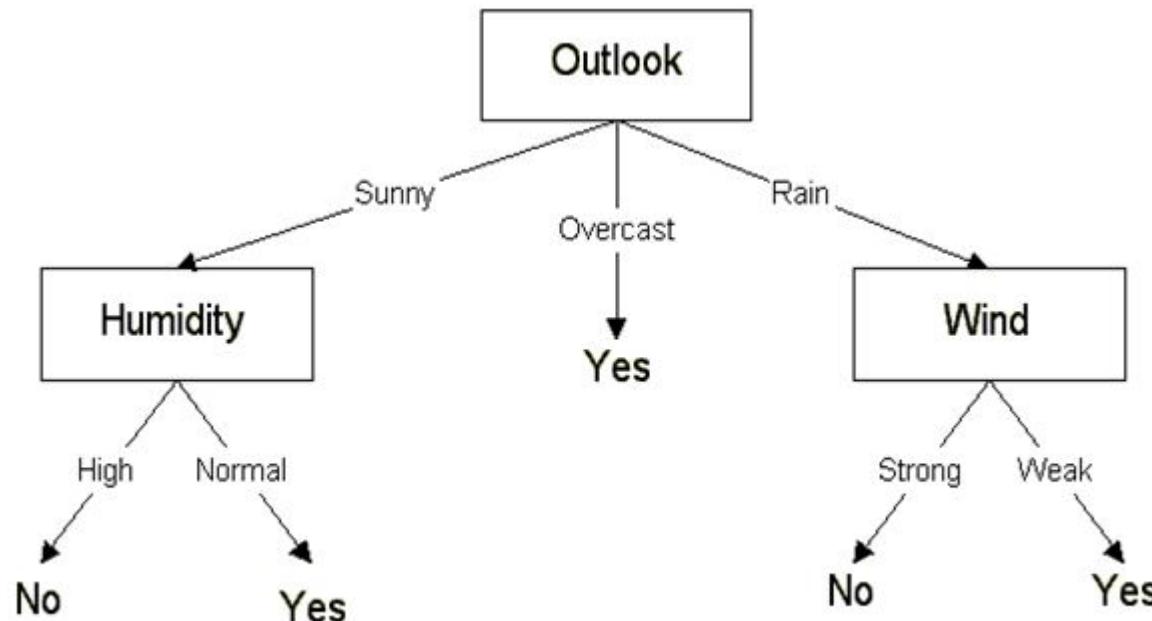
- Not suited to larger datasets as the training time with SVMs can be high
- Less effective on noisier datasets with overlapping classes
- Requires you to handle missing data.

RANDOM FOREST (RF)

- Very popular.
- Also known as random decision forest, it is a ensemble method.
- Ensemble means it uses multiple learning models to gain better predictive results.
- Therefore, creating an entire ‘forest’ of random uncorrelated decision trees to arrive at the best possible answer.

RANDOM FOREST (RF)

- Top-down approach where the root nodes creates binary splits until a certain criteria is met.
- The more trees in the forest the more robust is the prediction and thus high accuracy.



RANDOM FOREST (RF)

Advantages

- Robust with high accuracy
- Handles missing values.
- Handles small and large datasets

Disadvantages

- Overfitting in noisy training data (poor performance on unseen data).
- Very vulnerable to changes in testing and training data

How MUCH TRAINING DATA IS REQUIRED

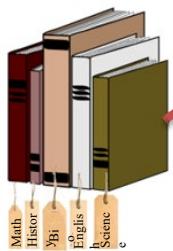
- The quality and amount of training data is often the single most dominant factor that determines the performance of a model.
- But exactly how much training data do you need? The correct answer is: **it depends on:**
 - task
 - performance to achieve
 - input features
 - noise in data
 - noise in features
 - and complexity of the model (number of classes).

SUMMARY OF PART 2

Unsupervised



Supervised

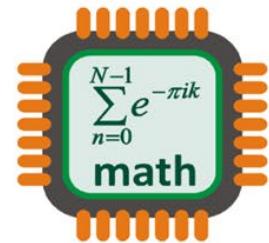


Feature Extraction



Unseen Data

ML Algorithm



Clustering

Predictive Model



Classification



Part 2: Text Classification

Introduction to Text Classification

Textual Features

Practical Text Classification Example

WHAT IS TEXT CLASSIFICATION

- Text Classification assigns one or more classes to a sample of text (e.g., sentence, document, etc.) according to its content.
- Examples:
 - Detecting spam emails (Spam / Not Spam)
 - Detecting news articles topics (Politics, Sports, Science,)
 - Assist search engines (Relevant results on Google/Bing search engines)
- Examples in Acc&Fin:
 - Determining the sentiment of a sentence or media article
 - Identifying potential fraud cases

FEATURES FOR TEXTUAL DATA

- Verbs, Nouns, Pronouns
- Adverbs, Function words
- Word length
- Number of syllables
- Topic
- Words frequency
- Different words
- Syntax, Semantic
- Readability
-etc

FEATURES FOR TEXTUAL DATA

Chairman's Statement

Chairman's statement: This is **my** first report to you as **Chairman** of your board of directors, following the successful admission of the **Group**'s shares to trading on AIM on 10 January 2008. 1700 **Group** Plc was incorporated on 9 May 2007 but did not trade until the acquisition of its first trading subsidiary, Hamblin Selection Ltd ("Hamblin") on 14 November 2007. The subsequent acquisition of the **Group**'s second trading subsidiary, Inspired Selection Ltd ("Inspired"), was completed on 10 January 2008. I would like to express **my** thanks to the **Group**'s employees

Governance Report

Corporate **governance** statement: Statement of compliance with the **Combined Code**. The **Group** recognises the value of the Principles of Good **Governance** and The **Combined Code** on Corporate **Governance** published in 2006 by the Financial Reporting Council ('the **Combined Code**'). The **Group** intends to comply with the **Combined Code** so far as is practical and appropriate for a public **group** of its size and nature. The **Group** supports the recommendations on corporate **governance** of the Quoted Companies Alliance (QCA) and has implemented steps to reach compliance.

Part 3: Deep Learning (Neural Network)

Introduction to Deep Learning



DEEP LEARNING



“Now! ... That should clear up
a few things around here!”

DEEP LEARNING

- is a subfield of machine learning concerned with algorithms inspired by the structure and function of the brain called **artificial neural networks**.
- Generally speaking, deep learning is a machine learning method that takes in an input **X**, and uses it to predict an output of **Y**.

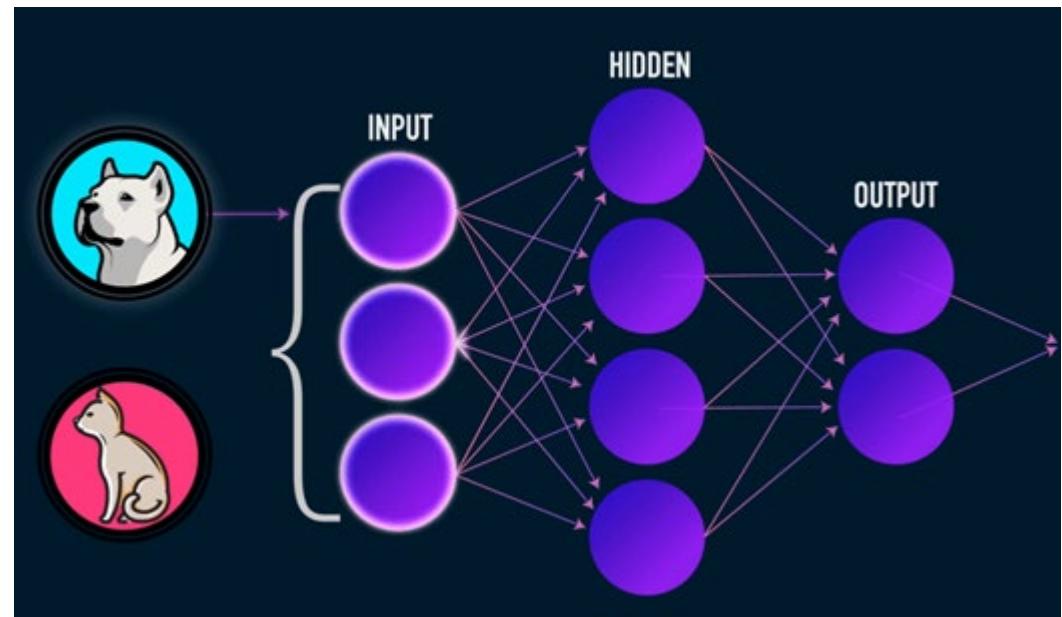
How DOES IT WORK?

- Given a large dataset of input and output pairs,
- a deep learning algorithm will try to minimize the difference between its prediction and expected output.
- By doing this, it tries to learn the association/pattern between given inputs and outputs
- this in turn allows a deep learning model to generalize to inputs that it hasn't seen before.



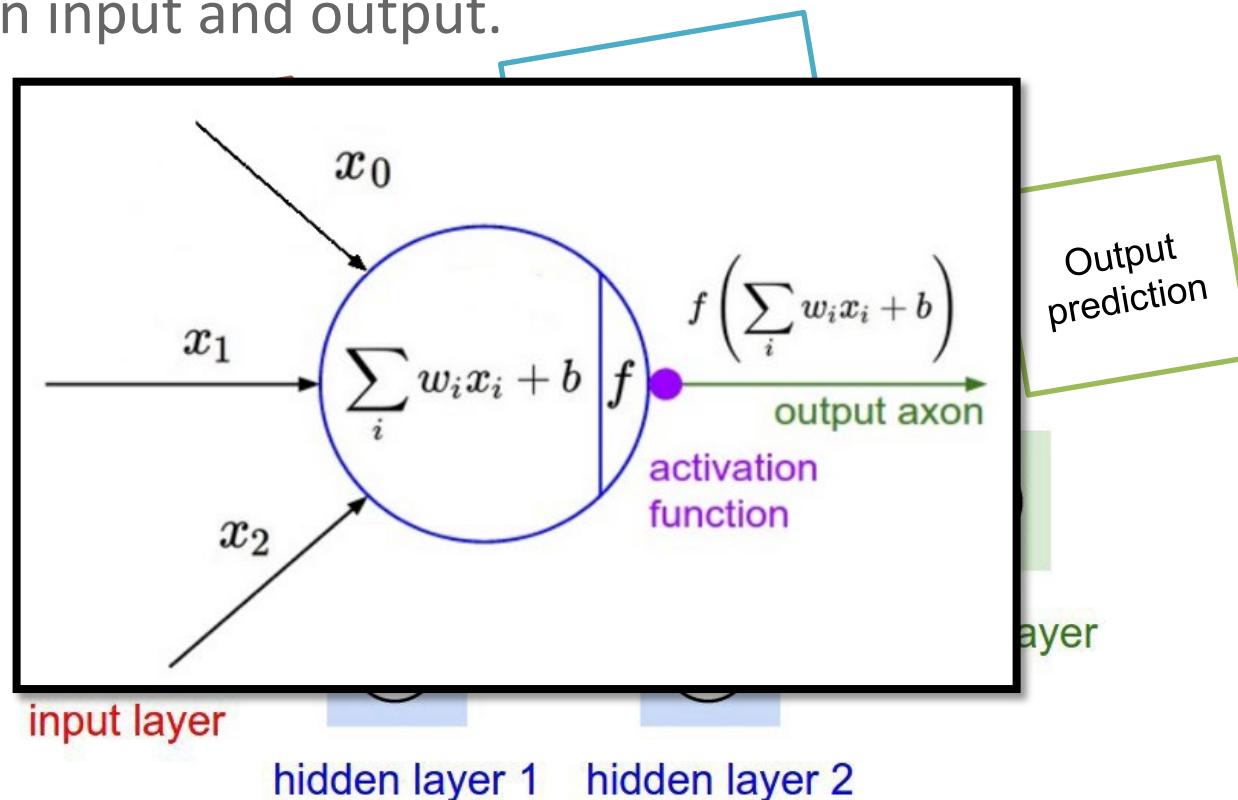
EXAMPLE

If image is labelled as dog but DL predicts cat, then DL learns the features of the given image associated with dog.



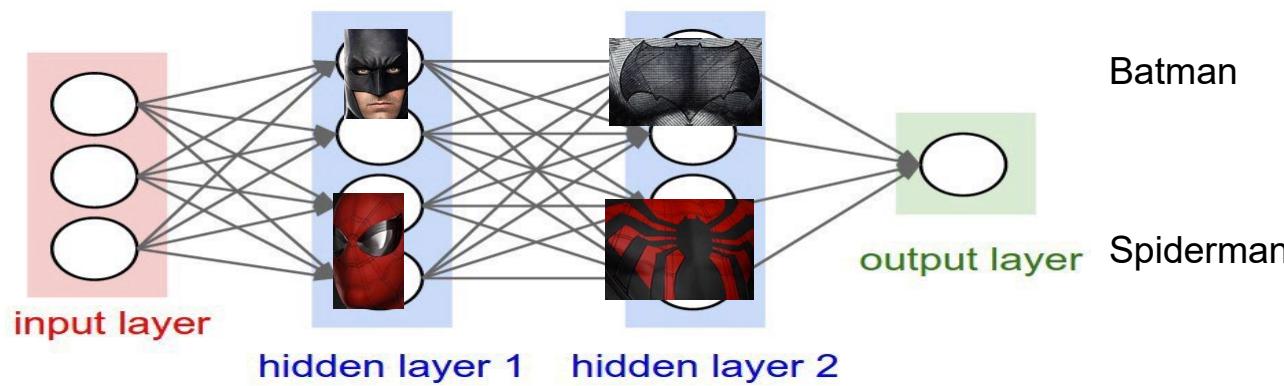
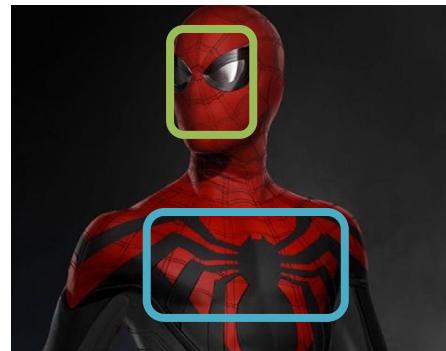
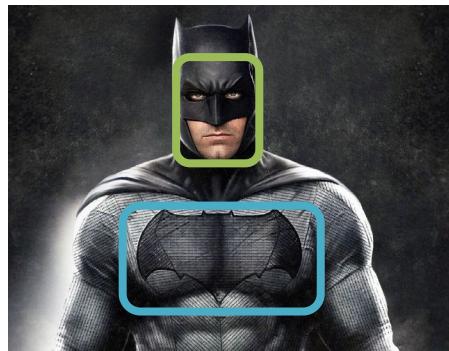
How Do DEEP LEARNING ALGORITHMS “LEARN”?

- DL algorithms use Neural Network to find associations between input and output.



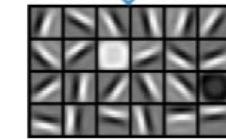
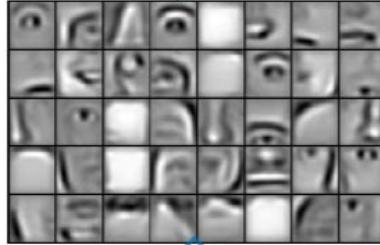
WHAT IS A HIDDEN LAYER?

- artificial neurons take in a set of weighted inputs and produce an output through an activation function

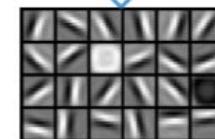


ANOTHER EXAMPLE

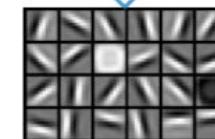
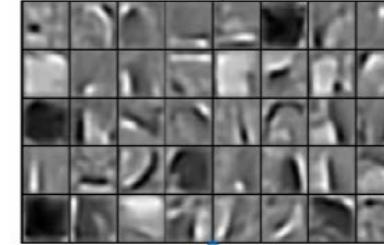
Faces



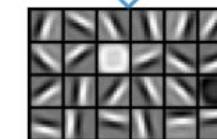
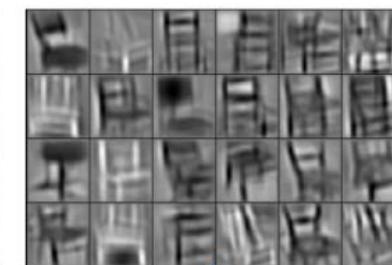
Cars



Elephants



Chairs



LOSS FUNCTION

- After the DL NN passes its inputs all the way to its outputs,
- the network evaluates how good its prediction was
- by comparing prediction to output.

The diagram illustrates the Mean Square Error loss function. On the left, there is a small image of Batman labeled "Expected Output". On the right, there is a small image of Spider-Man labeled "Prediction". In the center, the formula for the loss is shown:

$$\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

"Mean Square Error" los function

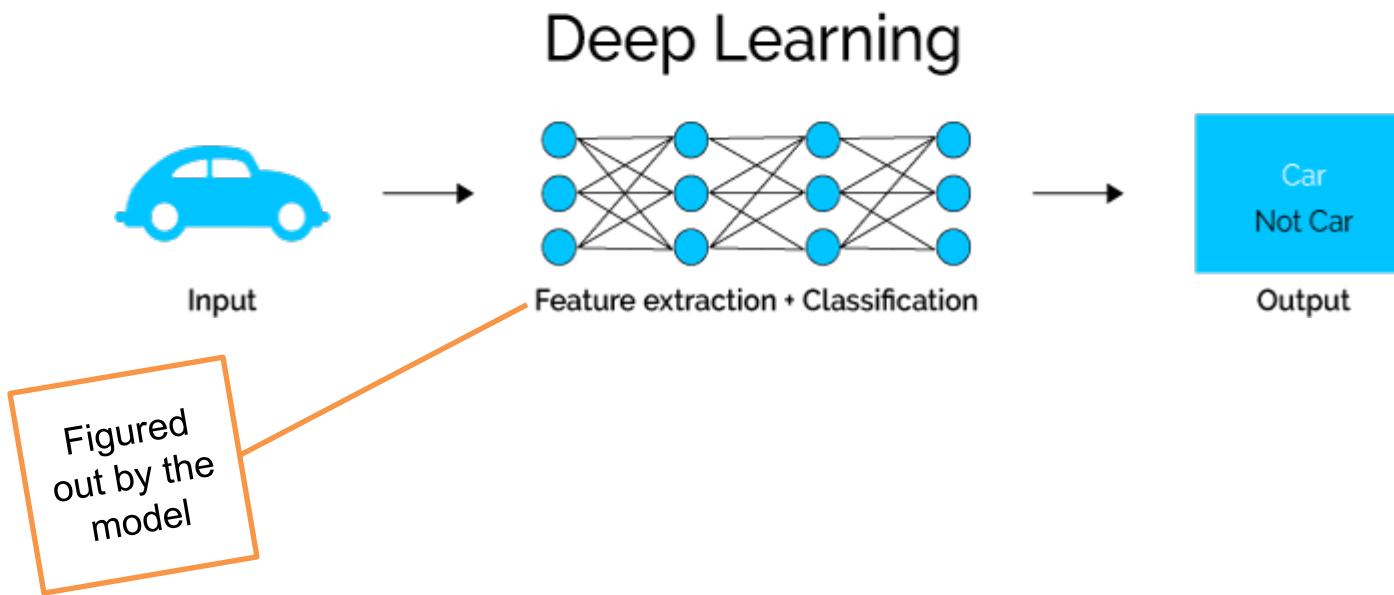
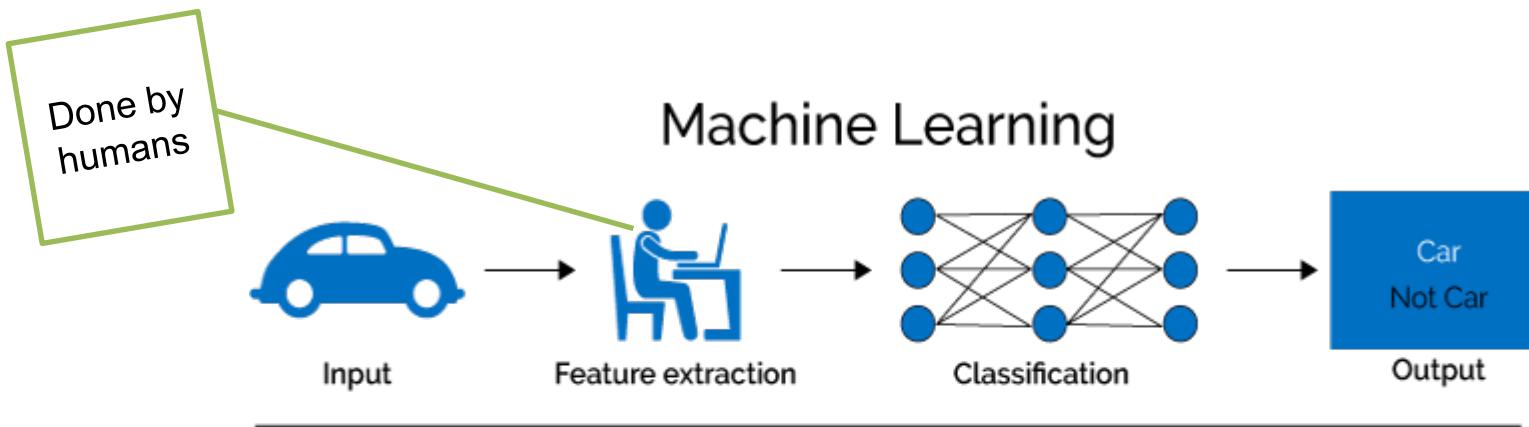
DEEP LEARNING GOAL

- Minimize Loss by adjusting weights of the network.
- Using something called “Back propagation”
- the network backtracks through all its layers to update the weights of every node in the opposite direction of the loss function
- in other words, every iteration of back propagation should result in a smaller loss function than before.
- the continuous updates of the weights of the network ultimately turns it into a precise function approximator — one that models the relationship between inputs and expected outputs.

SO WHY IS IT CALLED “DEEP” LEARNING?

- The “deep” part of deep learning refers to creating deep neural networks.
- This refers a neural network with a large amount of layers — with the addition of more weights, the neural network improves its ability to approximate more complex functions.
-

DEEP LEARNING VS TRADITIONAL MACHINE LEARNING



DEEP LEARNING VS TRADITIONAL MACHINE LEARNING

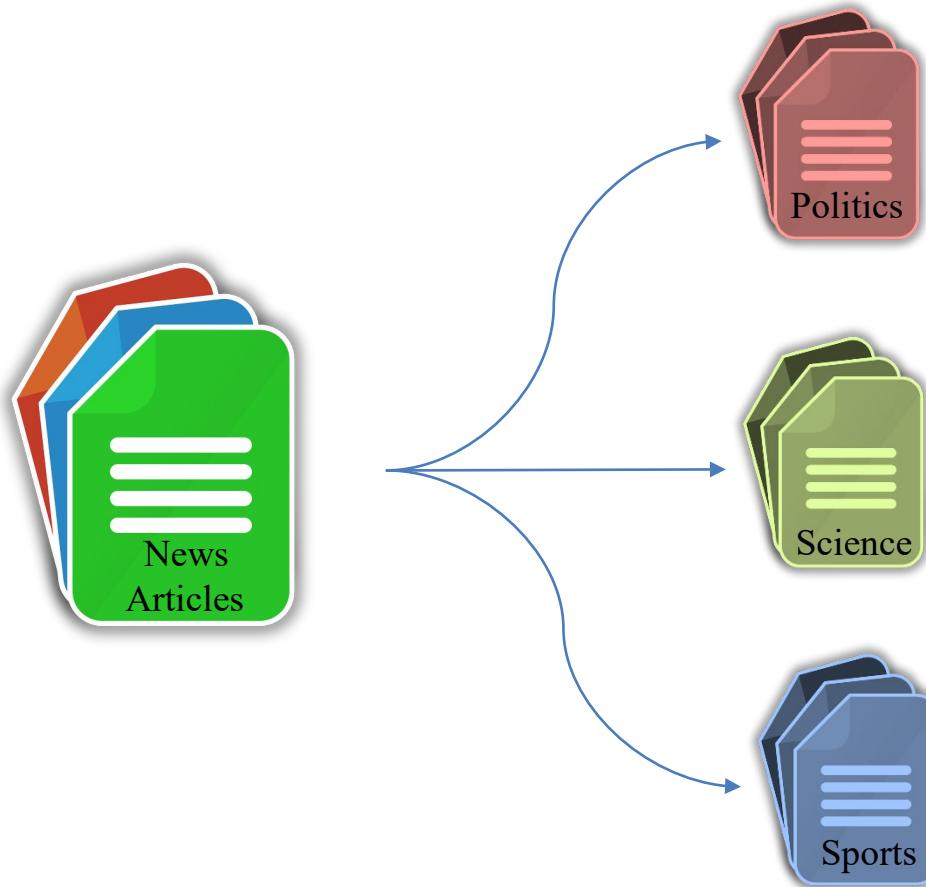
- Deep learning works better with more data (the larger the better).
- With Small data traditional machine learning algorithms are preferred.
- Deep learning requires high computational time and high end infrastructure (GPU).
- When there is lack of domain understanding for features selection Deep Learning outshine (worry less about features).
- Deep learning works well with complex problems (image classification, NLP, speech recognition ...etc).

CONCLUSION

- Deep learning is ultimately an expansive field, and is far more complex than I've described it to be.
- Various types of neural networks exist for different tasks (e.g. Convolutional NN for computer vision, Recurrent NN for NLP),
- and go far and beyond the basic neural network that I've covered.



PRACTICAL: CLASSIFYING DOCUMENTS



THE TASK

- Building a model using **WEKA** to classify Wikipedia Articles into either Spanish or English articles. So basically detecting the language of the articles.



DATASET

- 70 Spanish Wikipedia articles
- 70 English Wikipedia articles
- Articles saved in plain text file formats ending with .txt.

MACHINE ALGORITHM TOOL

- For the purpose of this workshop we will be using WEKA machine learning algorithms to classify our data.



WHAT IS WEKA?



- Weka is a collection of machine learning algorithms for data mining tasks.
- Weka contains tools for:
 - data pre-processing,
 - **classification**,
 - regression,
 - clustering,
 - and visualization

WEKA JAVA

- WEKA is originally written in Java but in this workshop we will be using the Graphical User Interface (GUI) application.
- So no coding required.

USING WEKA TO CLASSIFY TEXT

- - Classify text into two classes
- - Compare NB vs RF vs SVM
- - Reduce attributes (features)

READING DATASET INTO ARFF

- ARFF:
Attribute-
Relation File
Format
- Used for
WEKA
- Structured
format

```
% This is a toy example, the UCI weather dataset.  

% Any relation to real weather is purely coincidental  

@relation weather  

@attribute outlook {sunny, overcast, rainy}  

@attribute temperature real  

@attribute humidity real  

@attribute windy {TRUE, FALSE}  

@attribute play {yes, no}  

@data  

sunny,85,85,TRUE,no  

sunny,80,90,FALSE,no  

overcast,83,86,TRUE,yes  

rainy,70,96,FALSE,yes  

rainy,68,80,FALSE,yes  

rainy,65,70,TRUE,no  

overcast,64,65,TRUE,yes  

sunny,72,95,FALSE,no  

sunny,69,70,FALSE,yes  

rainy,75,80,FALSE,yes  

sunny,75,70,TRUE,yes  

overcast,72,90,TRUE,yes  

overcast,81,75,FALSE,yes  

rainy,71,91,TRUE,no
```

Comment ← **Dataset name**

Attributes ← **Target / Class variable**

Data Values ←

CONFUSION MATRIX (FOR CLASSIFICATIONS)

Confusion Matrix		True Condition	
		Actual +ve	Actual -ve
Predicted Condition	Pred. +ve	TP	FP (Type I)
	Pred. -ve	FN (Type II)	TN

- TP (or TN) Rate = Percentage of True Positives (Negatives)
 - Aka recall (next slide), TPR = probability of detection
 - TPR = sensitivity, TNR = specificity
 - $TPR = TP/(TP+FN)$; $TNR = TN/(TN+FP)$
- FP Rate = Percentage of False Positives (FN similarly for Negatives)
 - Aka fallout, FPR = probability of a false alarm
 - $FPR = FP/(FP+TN)$

NLP MEASURES¹

¹: POWERS, D.M.W; EVALUATION: FROM PRECISION, RECALL, AND F-MEASURE TO ROC, INFORMEDNESS, MARKEDNESS & CORRELATION, JMLT, 2011

$$\text{recall} = \frac{\text{Number of document retrieved that are relevant}}{\text{Total number of documents that are relevant}}$$

i.e. what proportion of actual positives were identified correctly?
aka – sensitivity, same measure as True Positive Rate

$$\text{precision} = \frac{\text{Number of document retrieved that are relevant}}{\text{Total number of documents that are retrieved}}$$

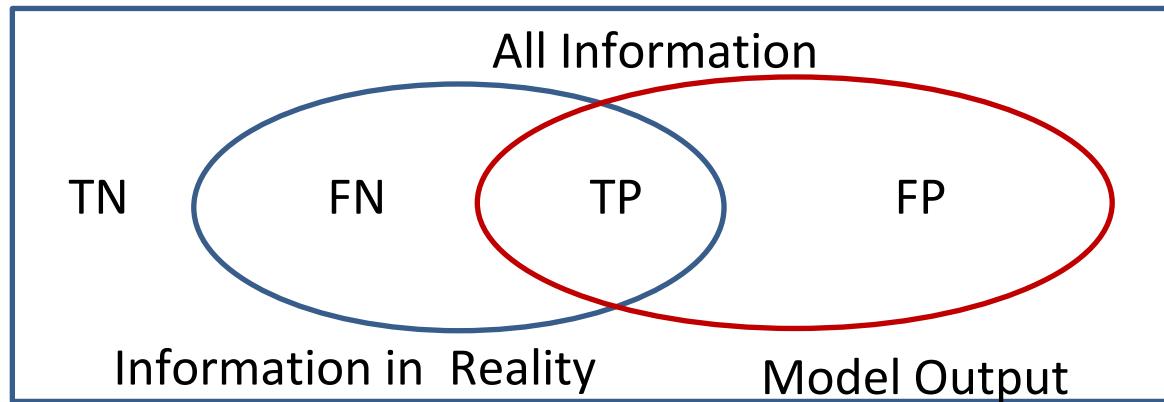
i.e. what proportion of positive identifications were actually correct?

$$\text{accuracy} = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}}$$

i.e. what proportion of predictions were actually correct?

$$F \text{ measure} = \frac{2 \times \text{recall} \times \text{precision}}{\text{recall} + \text{precision}}$$

NLP FORMULA



- **TP = True Positive**
- **FP = False Positive**
- $\text{Accuracy} = (|\text{TP}| + |\text{TN}|) / (|\text{TP}| + |\text{TN}| + |\text{FP}| + |\text{FN}|)$ used in Lift charts
- $\text{Recall} = |\text{TP}| / (|\text{TP}| + |\text{FN}|)$ used in ROC curves²
- $\text{Precision} = |\text{TP}| / (|\text{TP}| + |\text{FP}|)$

²: Receiver Operating Characteristic Curves, plots TPR vs FPR



Text Classification Practical

WEKA Classifier



NEEDED SOFTWARE

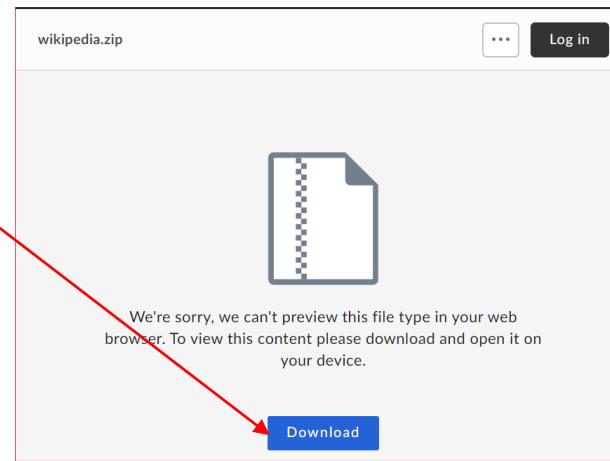
1 – WEKA

2 – NOTEPAD ++



NEEDED MATERIALS

- Create a new folder called **it classification** on your desktop:
C:\Users\username\Desktop\classification
- Go to <http://bit.ly/2XPZ0RI> and download **Wikipedia.zip** (click Download):



- Downloaded as .zip files. Right click and extract all contents.

HANDS-ON TASK STEPS

1. Convert dataset into ARFF file format.
2. Prepare dataset and extract features
3. Train a classifier using training data
4. Create a classification model
5. Test the classification model

CONVERT DATASET INTO ARFF FILE FORMAT

- Our data comes in as text files (folder1: EN, folder2: ES) with each containing 70 text (.txt) files
- We want to automatically convert the dataset into ARFF using a single command line.
- To do this start by running WEKA GUI Chooser



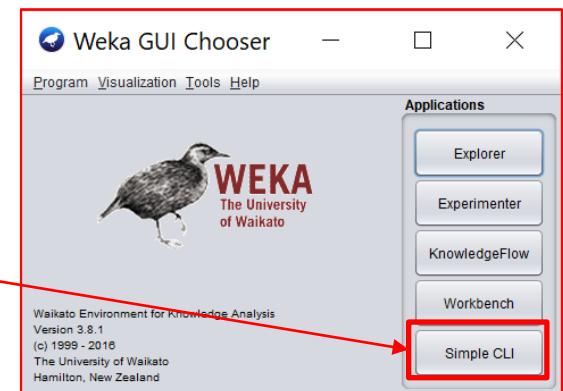
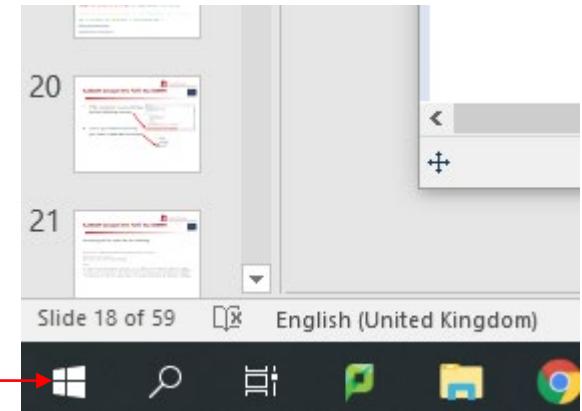
EN-madrid-ES-EN-1.txt
EN-madrid-ES-EN-2.txt
EN-madrid-ES-EN-3.txt
EN-madrid-ES-EN-4.txt



ES-madrid-ES-EN-1.txt
ES-madrid-ES-EN-2.txt
ES-madrid-ES-EN-3.txt
ES-madrid-ES-EN-4.txt

CONVERT DATASET INTO ARFF FILE FORMAT

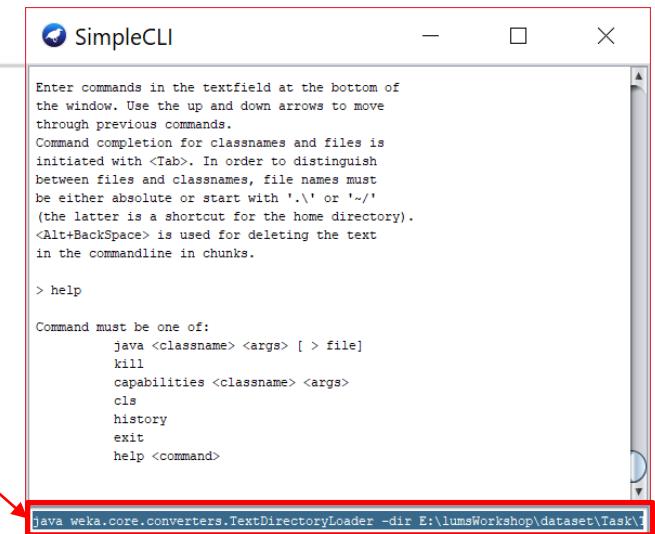
- To do this start by running WEKA GUI Chooser
- If you can't find it, click Windows and type WEKA
- Click on Simple CLI





CONVERT DATASET INTO ARFF FILE FORMAT

- Type following command in the command box below and hit enter:



```
SimpleCLI

Enter commands in the textfield at the bottom of the window. Use the up and down arrows to move through previous commands.
Command completion for classnames and files is initiated with <Tab>. In order to distinguish between files and classnames, file names must be either absolute or start with '.' or '/' (the latter is a shortcut for the home directory).
<Alt+BackSpace> is used for deleting the text in the commandline in chunks.

> help

Command must be one of:
    java <classname> <args> [ > file]
    kill
    capabilities <classname> <args>
    cls
    history
    exit
    help <command>

java weka.core.converters.TextDirectoryLoader -dir E:\lumsWorkshop\dataset\Task\1
```

java weka.core.converters.TextDirectoryLoader -dir Training_Directory > output_Arff_File

**java weka.core.converters.TextDirectoryLoader -dir
C:\Users\username\Desktop\classification\wikipedia>
C:\Users\username\Desktop\classification\wiki.arff**

Note: You may need to “skip” the backslashes “\” by adding another one in “\\”.

Note: if your path has spaces you may need to surround it by double quotations e.g. “C:/my path/”

E.g.:

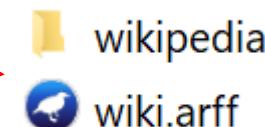
\\lancs\\homes\\24\\elhaj\\Downloads\\classification\\wikipedia
becomes
\\\\lancs\\\\homes\\\\24\\\\elhaj\\\\Downloads\\\\classification\\\\wikipedia

CONVERT DATASET INTO ARFF FILE FORMAT

- If the command is successful you get the following message

```
SimpleCLI
> cls
>
> java weka.core.converters.TextDirectoryLoader -dir
"E:\classification\wikipedia" > "E:\classification\wiki.arff"
Finished redirecting output to 'E:\classification\wiki.arff'.
```

- And in your dataset directory you'll find a **wiki.arff** file created





CONVERT DATASET INTO ARFF FILE FORMAT

Resulting arff file looks like the following

```
1 @relation E_classification_wikipedia
2
3 @attribute text string
4 @attribute @@class@@ {EN,ES}
5
6 @data
7
8 'THE AMERICA MUSEUM\r\nThe America museum is a national Spanish museum ',EN
9 'MUSEO DE AMÉRICA\r\nEl Museo de América es un Museo Nacional español con sede en Madrid',ES
10
```

CONVERT DATASET INTO ARFF FILE FORMAT

But what about **testing** dataset?

WEKA is very strict with arff files created for training and testing and requires both to have the same format, structure and attributes.

For example if your training dataset has features such as colour, shape and size but your testing has only colour and shape, WEKA will consider them mismatching.

CONVERT DATASET INTO ARFF FILE FORMAT

- Best and easy way to create an arff file for your testing (unseen) dataset is by taking that out of the training arff file before you do any further work on it.
- For example our arff file has 70 ES and 70 EN instances.
- We'll take a copy of the file keeping only 5 instances for each class.
- We'll then replace each ES and EN with “?”
- User notepad++ to open and edit files.

CONVERT DATASET INTO ARFF FILE FORMAT

```
@relation E_classification_wikipedia
@attribute text string
@attribute @@class@@ {EN,ES}

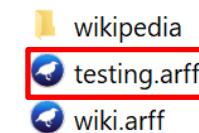
@data
'space reserved for the museological purposes for which it was designed.\r\n',EN
' espacio para los fines museológicos para los que había sido concebido.\r\r',ES
'd his business with parks in other continents, as it was Six Flags\'s case.',EN
'su negocio con parques en otros continentes, como era el caso de Six Flags.',ES
```

Convert known labels
to unknown (?)

```
@relation E_classification_wikipedia
@attribute text string
@attribute @@class@@ {EN,ES}

@data
'space reserved for the museological purposes for which it was designed.\r\n',?
' espacio para los fines museológicos para los que había sido concebido.\r\r',?
'd his business with parks in other continents, as it was Six Flags\'s case.',?
'su negocio con parques en otros continentes, como era el caso de Six Flags.',?
```

- Save the arff copy file as **testing.arff** in the dataset directory.
- We'll come back to this file later in the session.

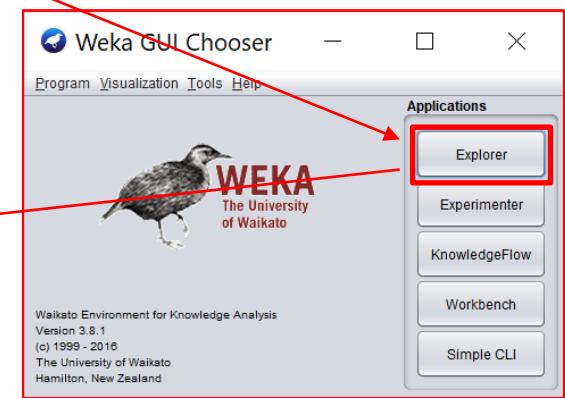
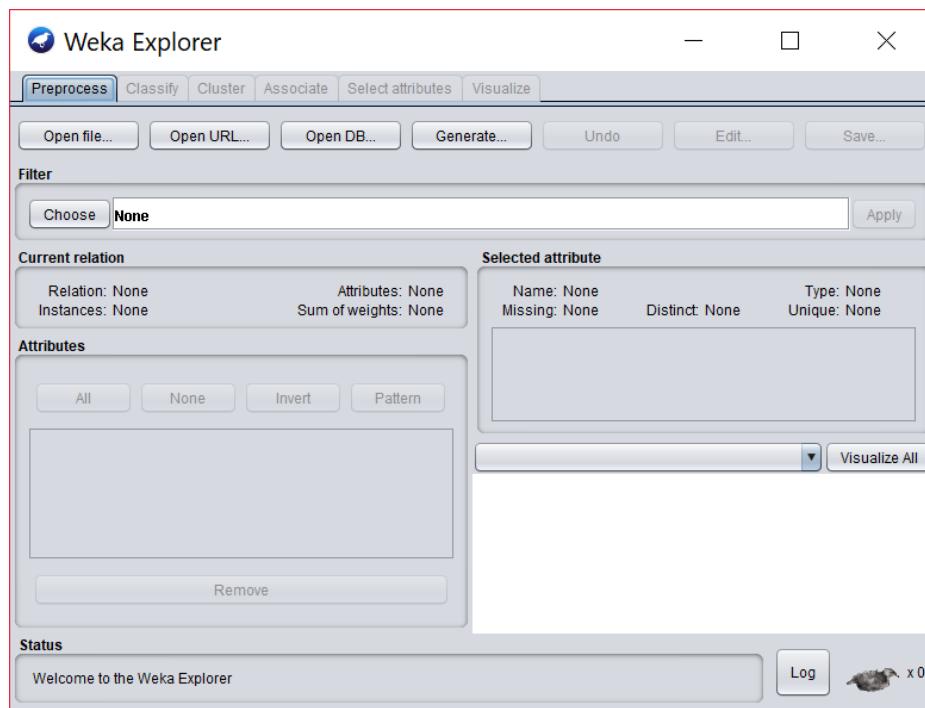


STEPS

1. Convert dataset into ARFF file format.
2. **Prepare dataset and extract features**
3. Train a classifier using training data
4. Create a classification model
5. Test the classification model

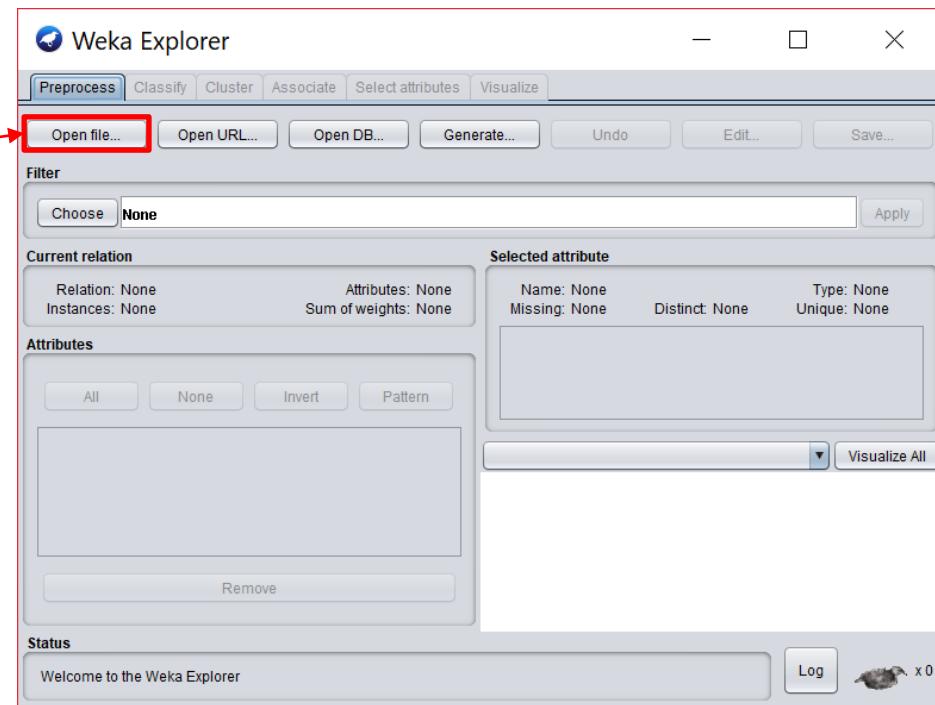
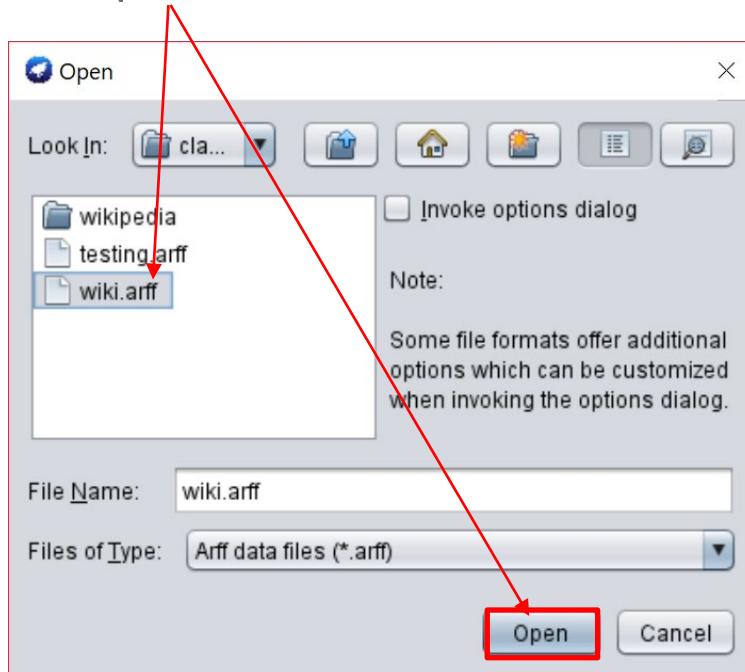
PREPARE DATASET AND EXTRACT FEATURES

- In WEKA GUI Chooser click on “Explorer”



PREPARE DATASET AND EXTRACT FEATURES

- We start by loading the arff we created into WEKA Explorer to view the contents.
- Navigate to your arff file and click Open



PREPARE DATASET AND EXTRACT FEATURES

Weka Explorer

Preprocess (circled in red)

Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... Undo Edit... Save...

Filter

Choose None Apply Stop

Current relation

Relation: E_classification_wi... Attributes: 2
Instances: 140 Sum of weights: 140

Selected attribute

No.	Label	Count	Weight
1	EN	70	70.0
2	ES	70	70.0

Name: @@class@@ Type: Nominal
Missing: 0 (0%) Distinct: 2 Unique: 0 (0%)

Attributes

All	None	Invert	Pattern
No.	Name		
1	text		
2	@@class@@		

Your labels(attributes)

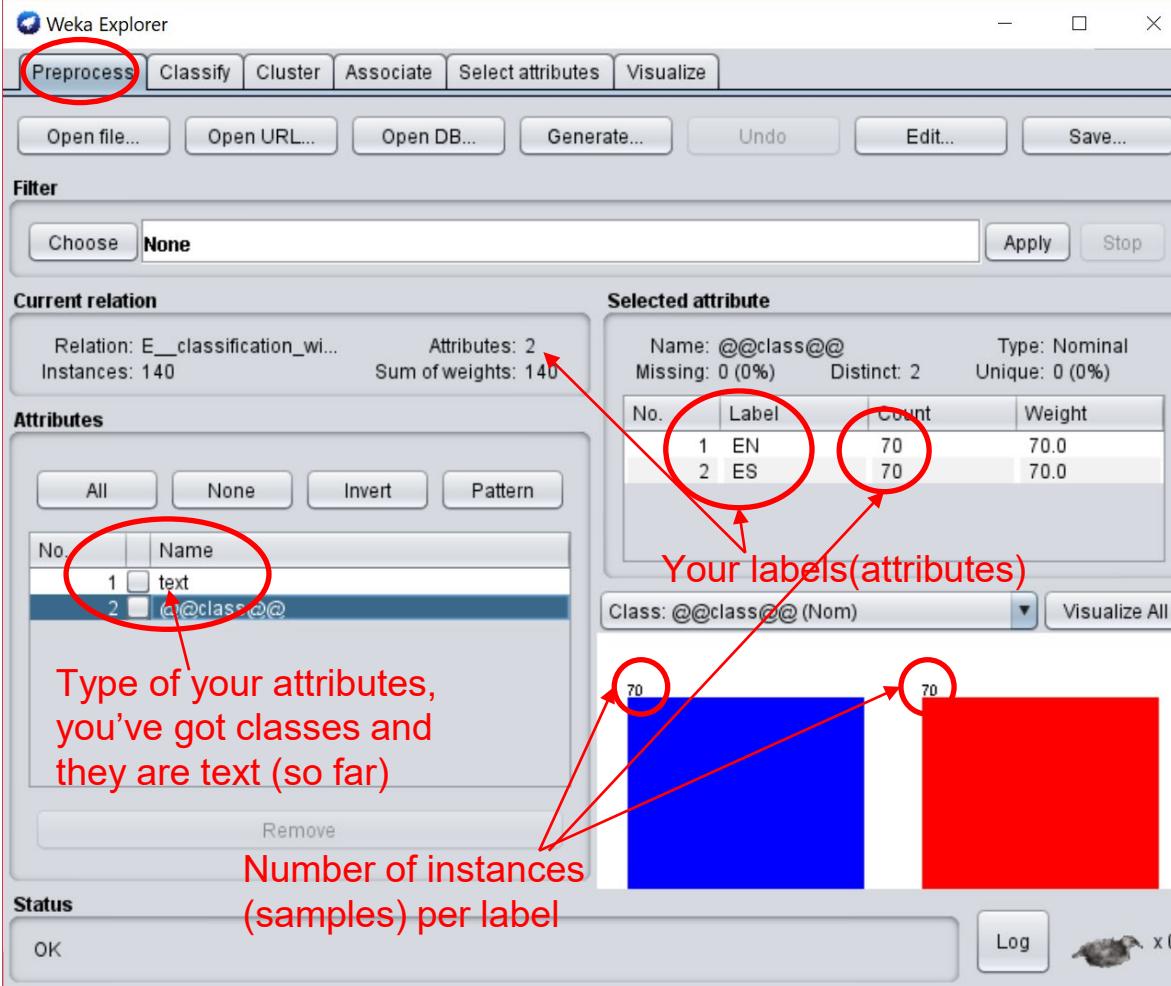
Class: @@class@@ (Nom) Visualize All

70 70

Type of your attributes, you've got classes and they are text (so far)

Number of instances (samples) per label

Status OK Log x 0



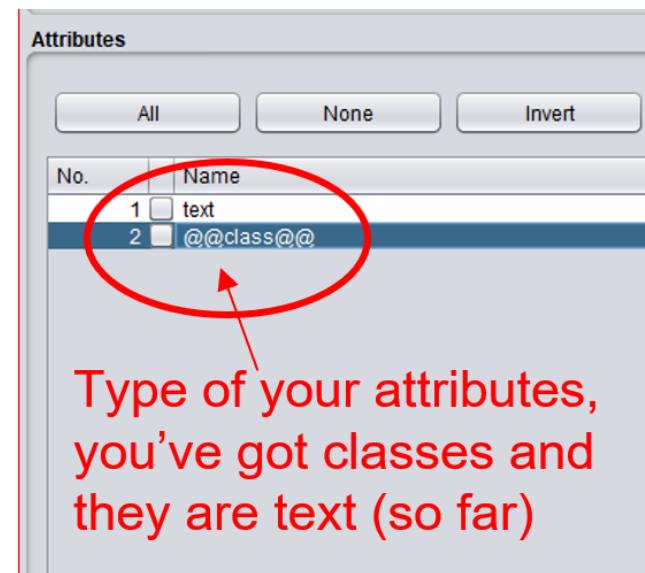
PREPARE DATASET AND EXTRACT FEATURES

As shown when loading the arff the data is made of “text”.

Algorithms deal with numbers not textual data.

Therefore we need to convert our text (string) attributes into set of numbers.

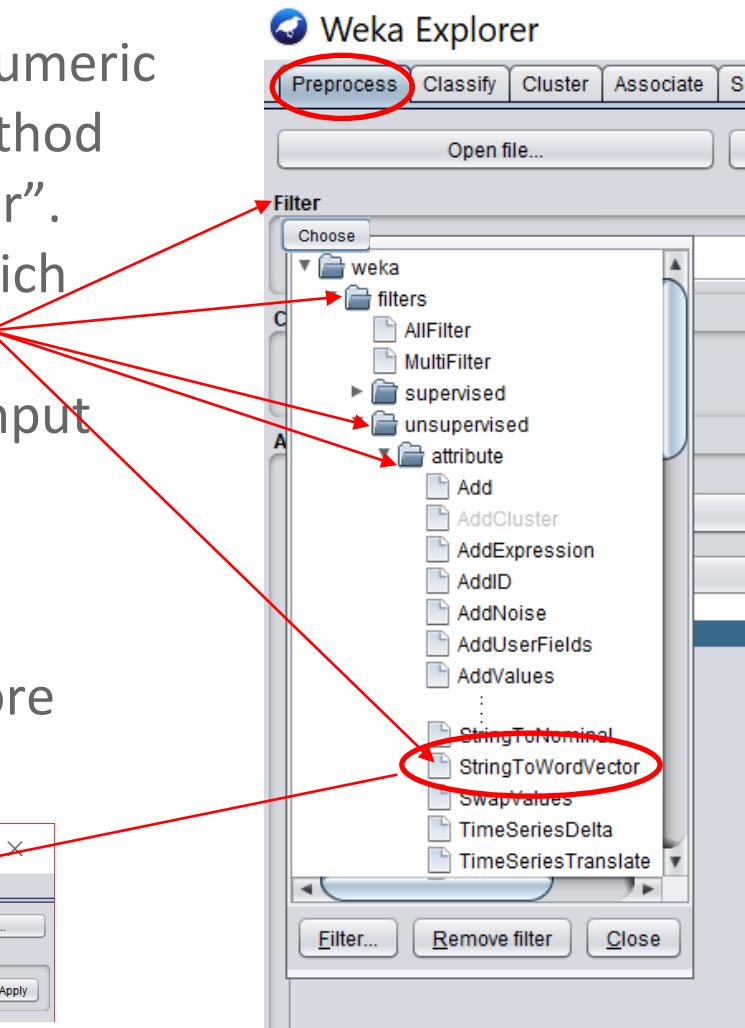
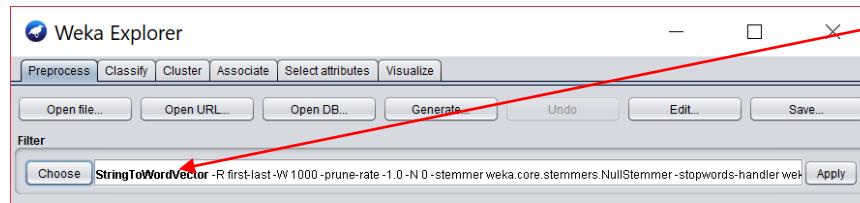
The numbers represent word occurrence information from the text.



PREPARE DATASET AND EXTRACT FEATURES

To convert the textual data into numeric attributes we will use a WEKA method (filter) called “StringToWordVector”. The method is “unsupervised” which means it extract word occurrence information without needing an input from you.

Select and click on “StringToWordVectors” to see more properties

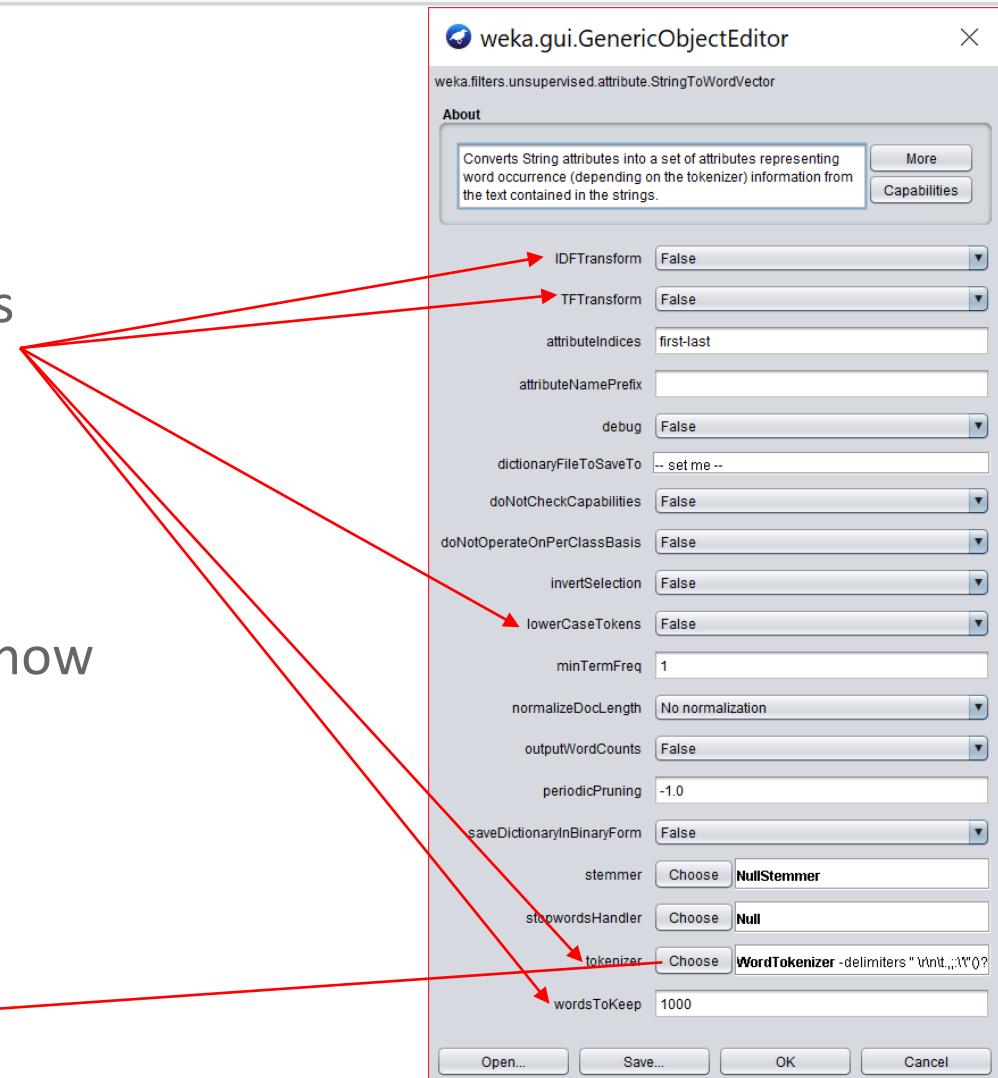
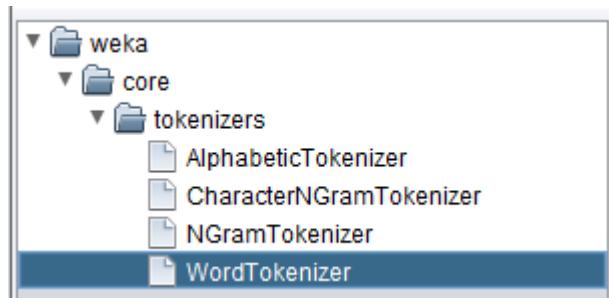


PREPARE DATASET AND EXTRACT FEATURES

StringToWordVector extract features from the text mainly related to frequencies.

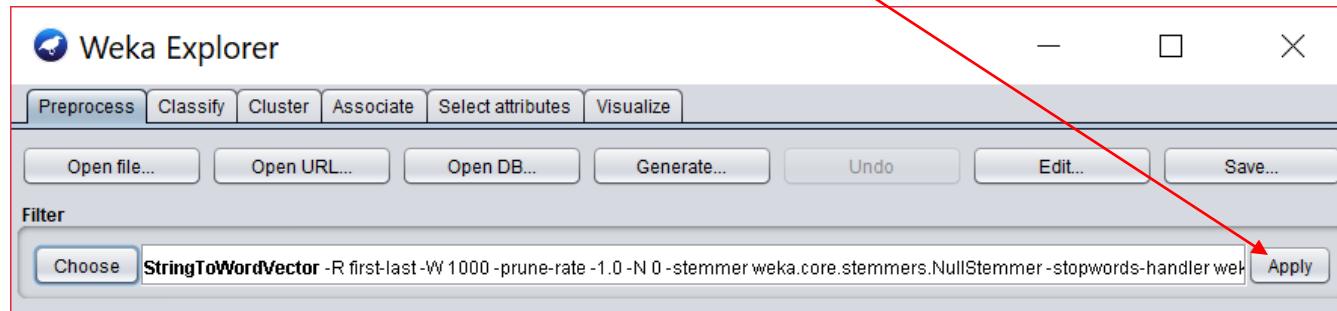
The options on the right helps customise what features to extract.

For example you can control how to split words



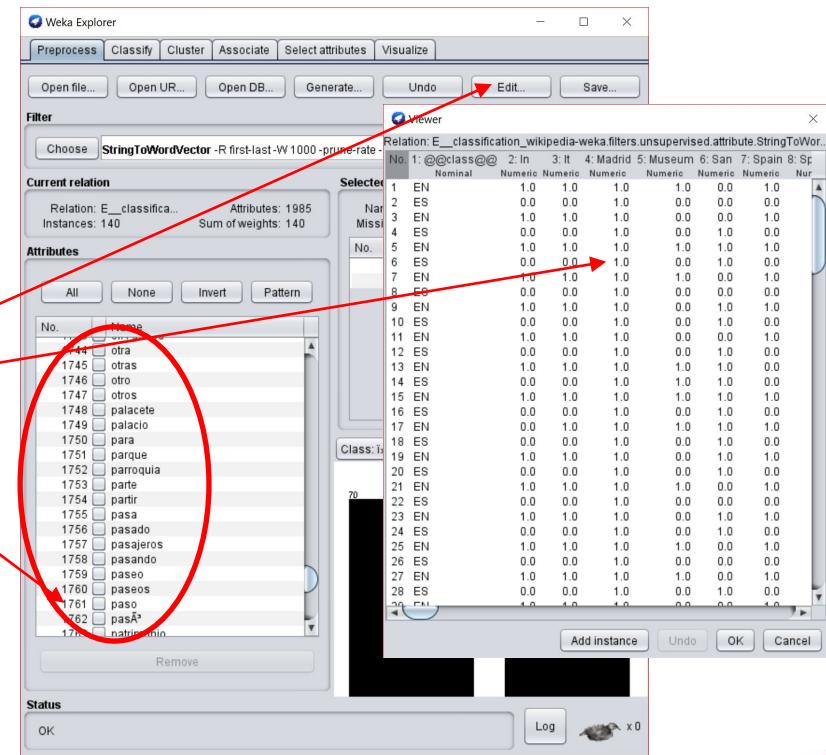
PREPARE DATASET AND EXTRACT FEATURES

For now leave the default options and click on apply to convert the textual data into numbers



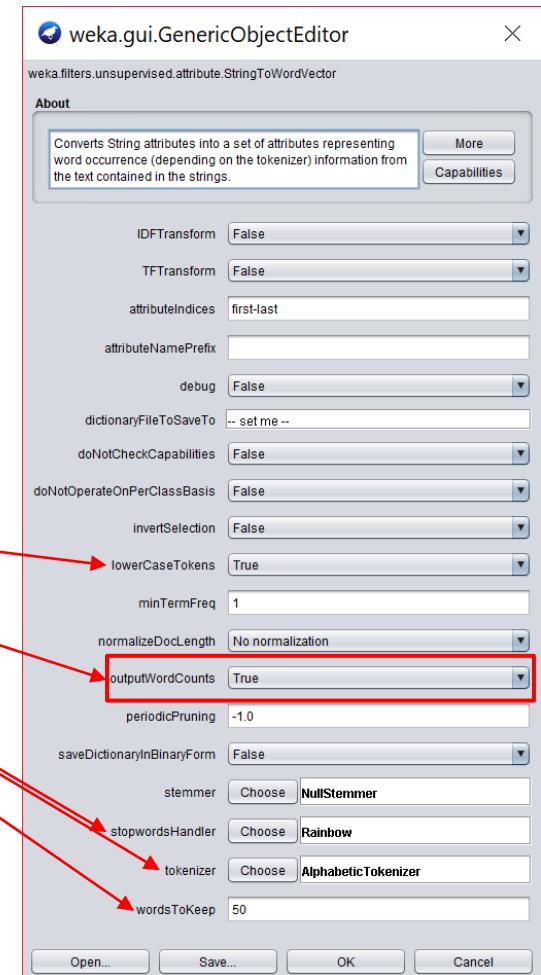
PREPARE DATASET AND EXTRACT FEATURES

- As soon as you click apply you'll notice an ordered list of numbers and terms appearing instead of the “text” attribute.
 - Those are now your numeric features which algorithms can use to train a classifier.
 - Note the number of features extracted



PREPARE DATASET AND EXTRACT FEATURES

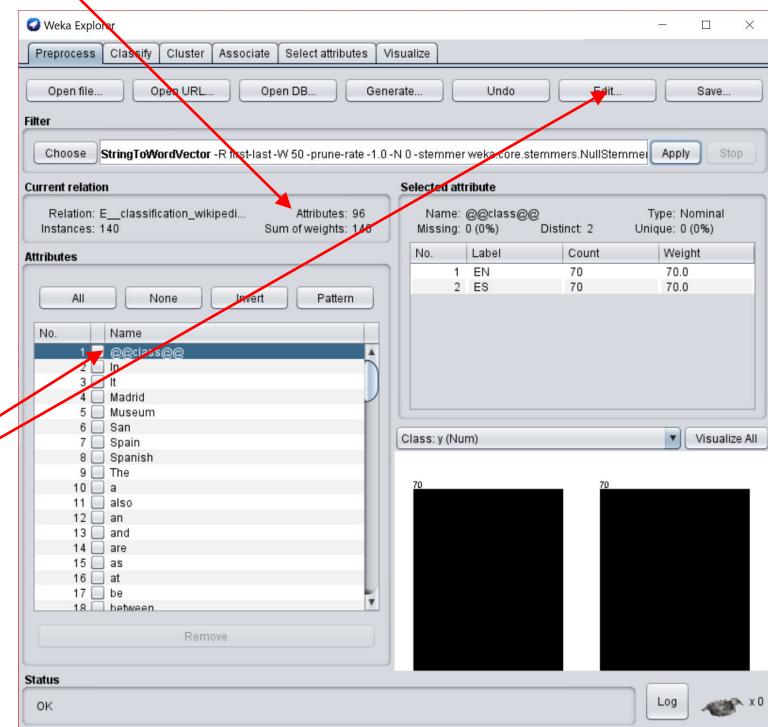
- Let's try reducing number of features.
 - Click “Undo” to cancel the previous operation.
 - Click on StringToWordVector and set the following.
 - Click OK
 - Then Click Apply.



PREPARE DATASET AND EXTRACT FEATURES

- We now have “approx.” 50-100 words (attributes).

- Notice the class
@@class@@ is now
the first in the list.
- We need to push it back
to the end as it's not a
feature.
- Click on “Edit...”



PREPARE DATASET AND EXTRACT FEATURES

- You'll notice:
 1. The @@class@@ is the first attribute.
 2. Word count for each of the ~50 words.
- Right click the @@class@@ column and click “Attribute as Class”, you'll notice it's the last attribute now.

Viewer

Relation: E_classification_wikipedia-weka.filters.unsupervised.attribute.St...

89: su	90: sus	91: también	92: un	93: una	94: y	95: @@class@@
Numeric	Numeric	Numeric	Numeric	Numeric	Numeric	Nominal
0.0	0.0	0.0	0.0	0.0	0.0	EN
5.0	2.0	2.0	5.0	5.0	23.0	ES
0.0	0.0	0.0	0.0	0.0	0.0	EN
7.0	2.0	5.0	12.0	9.0	33.0	ES
0.0	0.0	0.0	0.0	0.0	2.0	EN
9.0	5.0	0.0	15.0	11.0	40.0	ES
0.0	0.0	0.0	0.0	0.0	0.0	EN
11.0	3.0	1.0	6.0	3.0	22.0	ES

Viewer

Relation: E_classification_wikipedia-weka.filters.unsupervised.attribute.St...

No. 1: @@class@@	2: a	3: also	4: an	5: and	6: are	7: as
Nom	Nomeric	Numeric	Numeric	Numeric	Numeric	Numeric
1 EN		14.0	4.0	2.0	27.0	1.0
2 ES		21.0	0.0	0.0	0.0	0.0
3 EN		16.0	5.0	8.0	33.0	1.0
4 ES		14.0	0.0	0.0	0.0	0.0
5 EN		28.0	1.0	4.0	41.0	1.0
6 ES		15.0	0.0	0.0	0.0	0.0
7 EN						
8 ES						
9 EN						
10 ES						
11 EN						
12 ES						
13 EN						
14 ES						
15 EN						
16 ES						
17 EN						
18 ES						
19 EN						
20 ES						
21 EN						
22 ES						
23 EN						
24 ES						
25 EN						
26 ES						
27 EN						
28 ES						
29 EN						

Viewer

Relation: E_classification_wikipedia-weka.filters.unsupervised.attribute.St...

No. 1: @@class@@	2: a	3: also	4: an	5: and	6: are	7: as
Nom	Nomeric	Numeric	Numeric	Numeric	Numeric	Numeric
1 EN		Get mean...				
2 ES		Set all values to...				
3 EN		Set missing values to...				
4 ES		Replace values with...				
5 EN		Rename attribute...				
6 ES		Set attribute weight...				
7 EN		Attribute as class				
8 ES		Delete attribute				
9 EN		Delete attributes...				
10 ES		Sort data (ascending)				
11 EN		Optimal column width (current)				
12 ES		Optimal column width (all)				
13 EN						
14 ES						
15 EN						
16 ES						
17 EN						
18 ES						
19 EN						
20 ES						
21 EN						
22 ES						
23 EN						
24 ES						
25 EN						
26 ES						
27 EN						
28 ES						
29 EN						

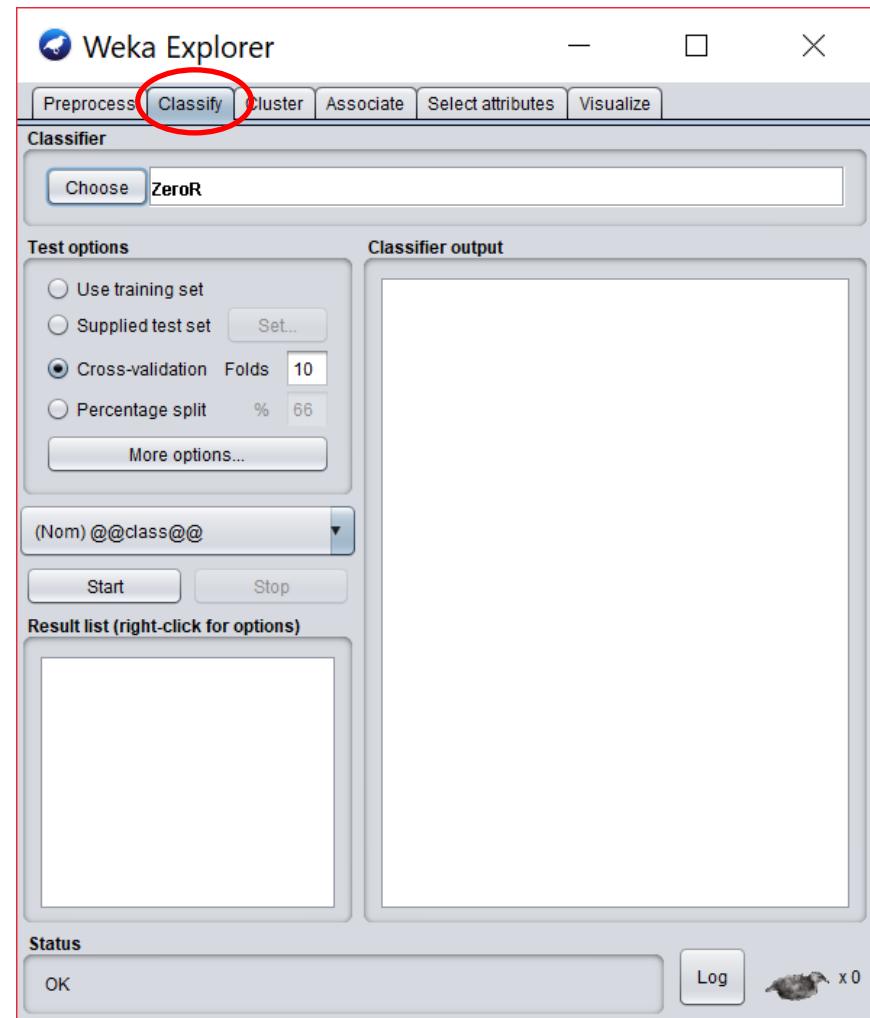
Add instance Undo OK Cancel

STEPS

1. Convert dataset into ARFF file format.
2. Prepare dataset and extract features
3. **Train a classifier using training data**
4. Create a classification model
5. Test the classification model

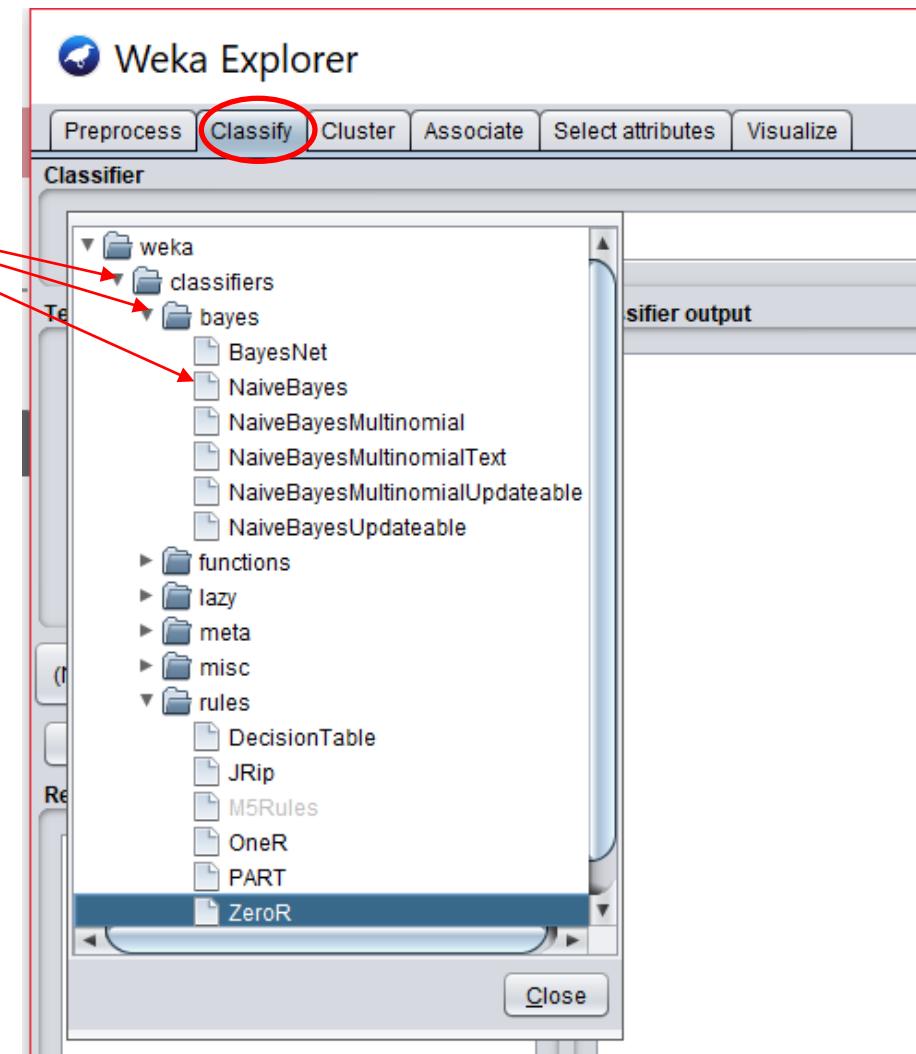
TRAIN A CLASSIFIER USING TRAINING DATA

- Now that we have extracted the features from our dataset we can now start to train a classifier.
- Start by selecting the “Classify” tab in the WEKA Explorer window.



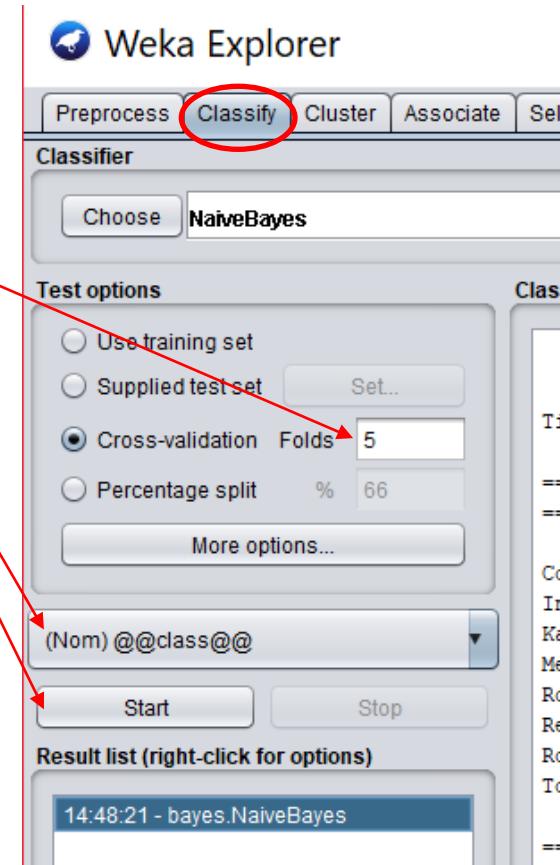
TRAIN A CLASSIFIER USING TRAINING DATA

- Choose a classifier to begin with. Let's start with a Naïve Bayes Classifier
- Click to select the classifier.



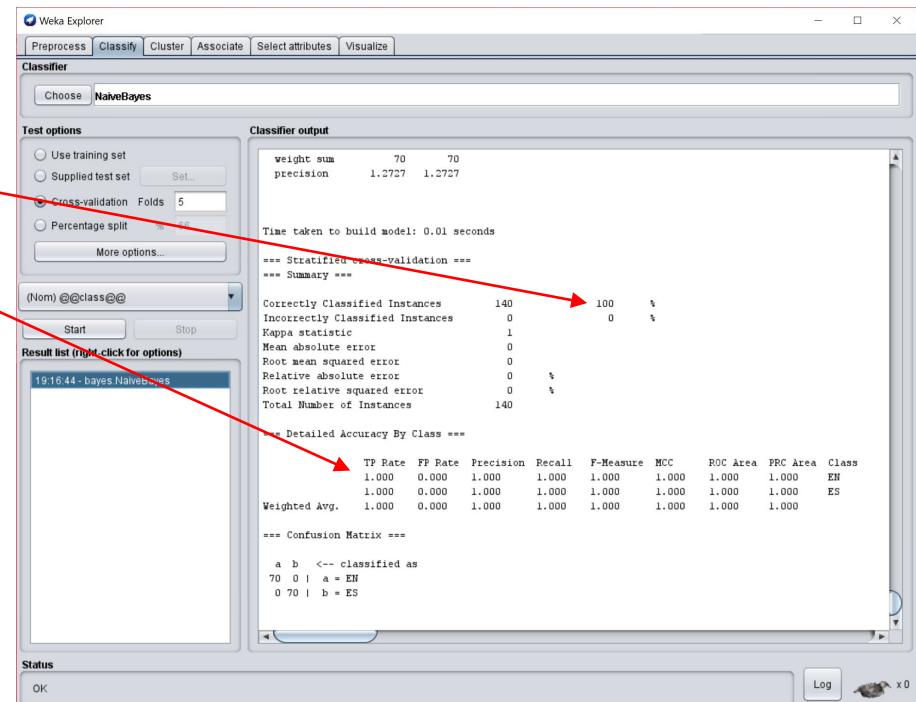
TRAIN A CLASSIFIER USING TRAINING DATA

- Select the number of folds (training k-folds). Select 5 and make sure the (Nom) @@class@@ is selected (ES and EN)
- Then click on Start to start the training process.



TRAIN A CLASSIFIER USING TRAINING DATA

- You'll see a summary of the training process using Naïve Bayes such as:
- Classifier's accuracy
- And a confusion matrix
- This is basically saying that training process expects a 100% accuracy when classifying new unseen samples.





Harder task



Chairman



Governance



Remuneration

Classifying UK financial sections (3 classes)

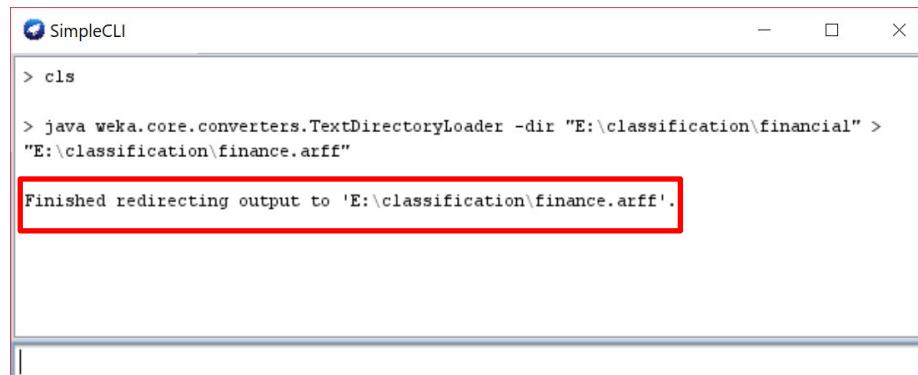
REPEAT STEPS

1. Convert dataset into ARFF file format.
2. Prepare dataset and extract features
3. Train a classifier using training data

REPEAT STEPS

1. Convert dataset into ARFF file format.

```
java weka.core.converters.TextDirectoryLoader -dir  
"E:\classification\financial" > "E:\classification\finance.arff"
```



```
> cls  
  
> java weka.core.converters.TextDirectoryLoader -dir "E:\classification\financial" >  
"E:\classification\finance.arff"  
  
Finished redirecting output to 'E:\classification\finance.arff'.
```



finance.arff

REPEAT STEPS

2. Prepare dataset and extract features

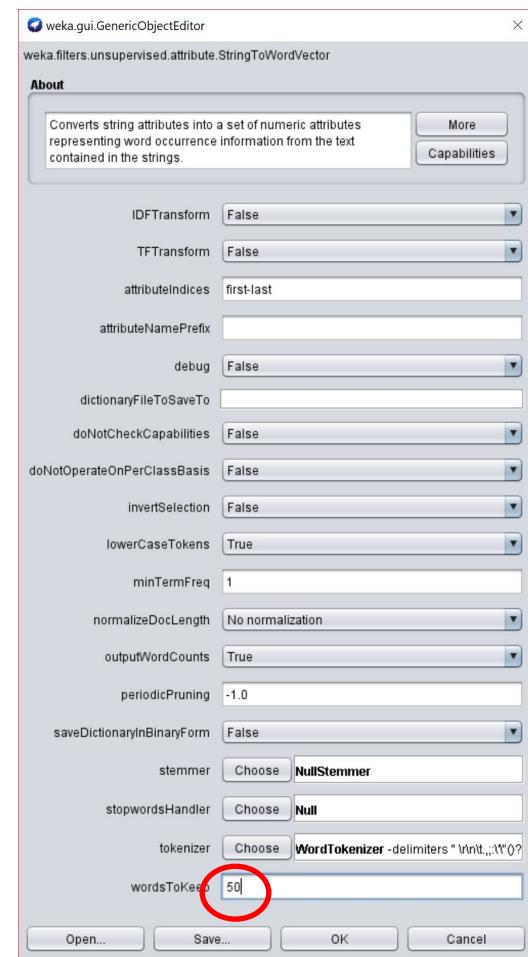
Viewer

Relation: E__classification_financial-weka.filters.unsupervised.attribute.String...

900: until 901: upon 902: vest 903: vested 904: vesting 905: @@class@@

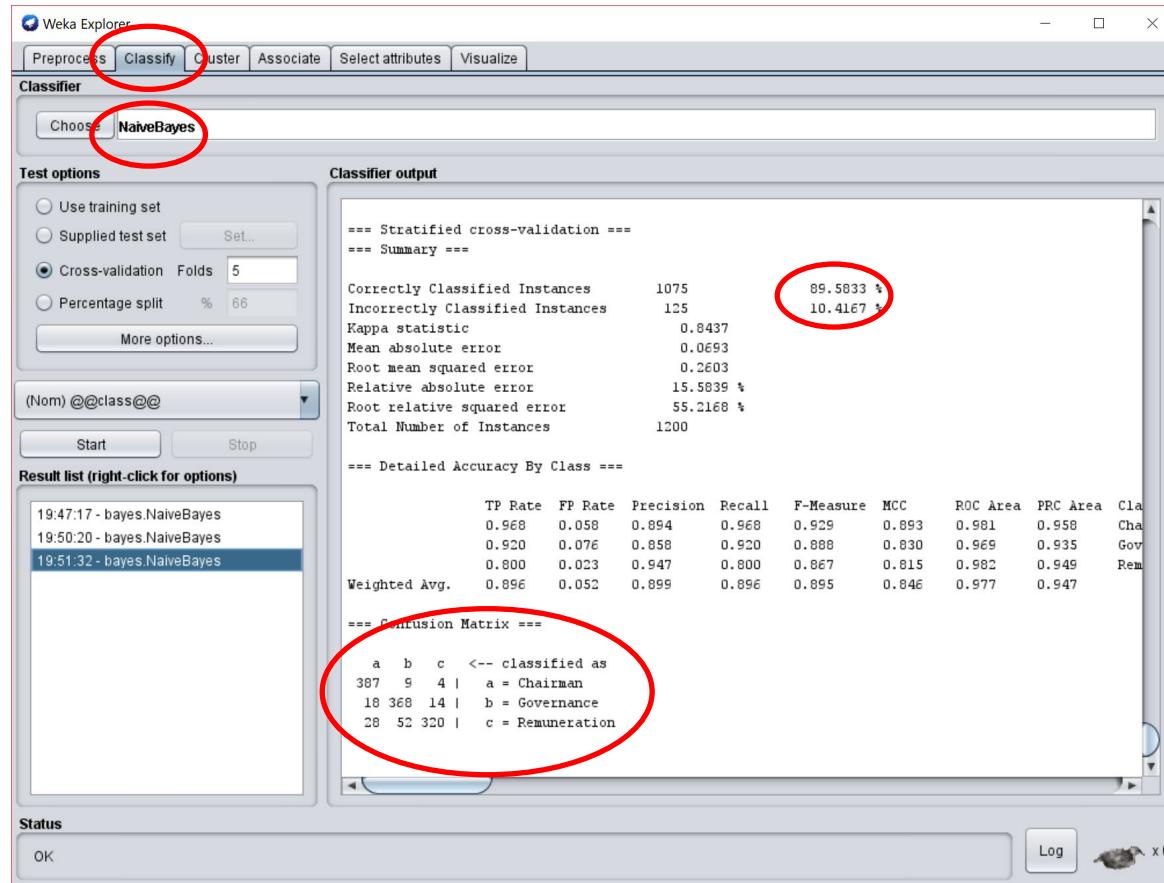
Numeric	Numeric	Numeric	Numeric	Numeric	Nominal
0.0	0.0	0.0	0.0	0.0	Chairman
0.0	0.0	0.0	0.0	0.0	Governance
2.0	0.0	7.0	3.0	14.0	Remuneration
0.0	0.0	0.0	0.0	0.0	Chairman
0.0	0.0	0.0	0.0	0.0	Governance
5.0	0.0	7.0	4.0	15.0	Remuneration
0.0	0.0	0.0	0.0	0.0	Chairman
0.0	0.0	0.0	0.0	0.0	Governance
1.0	0.0	4.0	4.0	16.0	Remuneration
0.0	0.0	0.0	0.0	0.0	Chairman
0.0	0.0	0.0	0.0	0.0	Governance
1.0	0.0	7.0	3.0	53.0	Remuneration
0.0	2.0	0.0	0.0	0.0	Chairman
0.0	0.0	0.0	0.0	0.0	Governance
1.0	4.0	10.0	2.0	4.0	Remuneration
2.0	0.0	0.0	0.0	0.0	Chairman
0.0	0.0	0.0	0.0	0.0	Governance
3.0	7.0	12.0	3.0	7.0	Remuneration
1.0	0.0	0.0	0.0	0.0	Chairman
0.0	0.0	0.0	0.0	0.0	Governance
10.0	3.0	4.0	2.0	1.0	Remuneration
0.0	0.0	0.0	0.0	0.0	Chairman
0.0	1.0	0.0	0.0	0.0	Governance
1.0	0.0	0.0	1.0	0.0	Remuneration
0.0	0.0	0.0	0.0	0.0	Chairman
0.0	0.0	0.0	1.0	0.0	Governance
1.0	0.0	1.0	1.0	1.0	Remuneration
0.0	0.0	0.0	0.0	0.0	Chairman
0.0	0.0	0.0	0.0	0.0	Governance

Add instance Undo OK Cancel



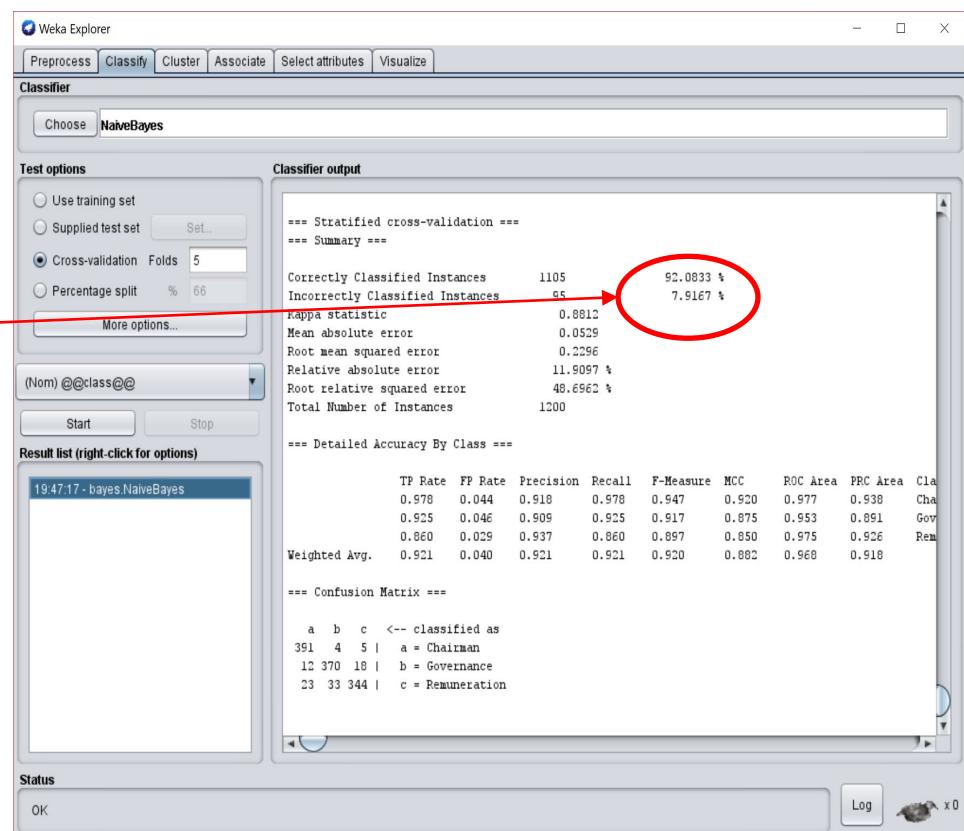
REPEAT STEPS

3. Train a classifier using training data



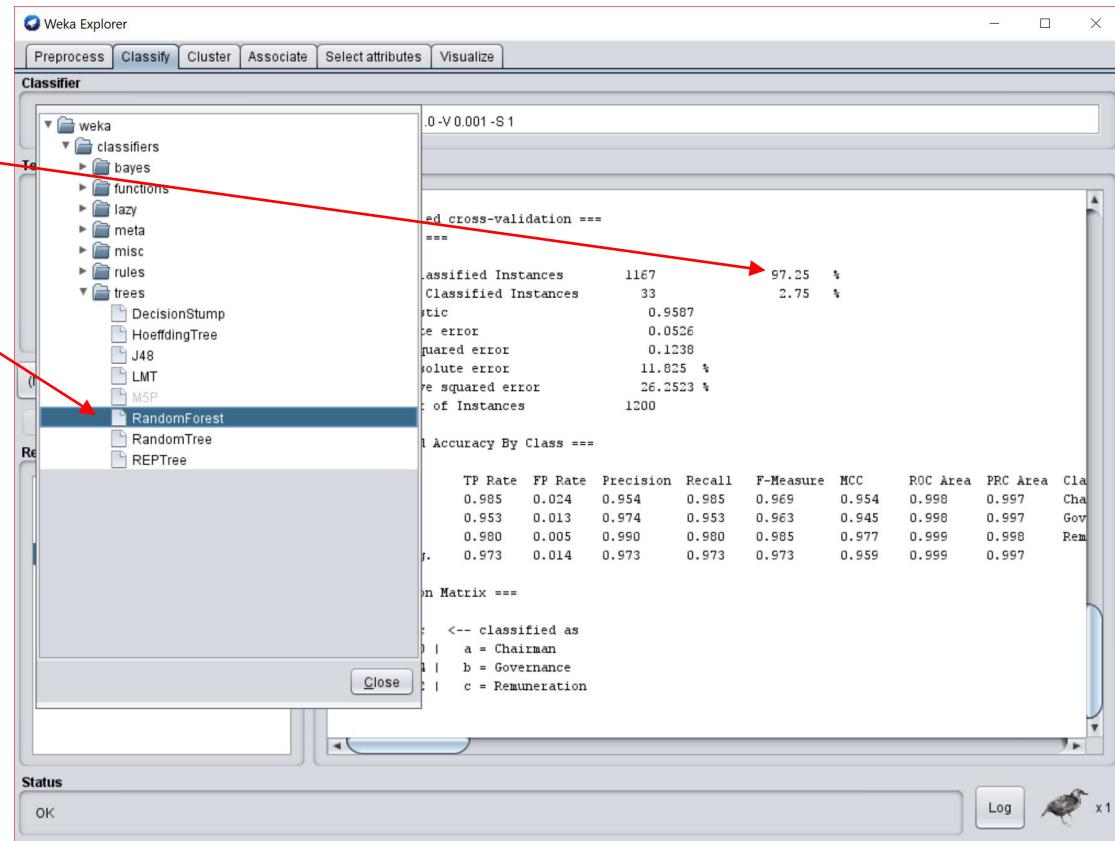
TRAIN A CLASSIFIER USING TRAINING DATA

- Can we increase accuracy?
- Let's go back to the Preprocess tab, undo until we are back to the beginning and then ask StringToWordVector to keep 500 words instead.
- Apply, then go back to Classify and run the classifier again.
- This time I got 92.08% accuracy.



TRAIN A CLASSIFIER USING TRAINING DATA

- Let's try Random Forest
- Basically, you'll have to experiment with few classifiers and different features to arrive at your best model and still avoid overfitting.



STEPS

1. Convert dataset into ARFF file format.
2. Prepare dataset and extract features
3. Train a classifier using training data
4. **Create a classification model**
5. Test the classification model

CREATE A CLASSIFICATION MODEL

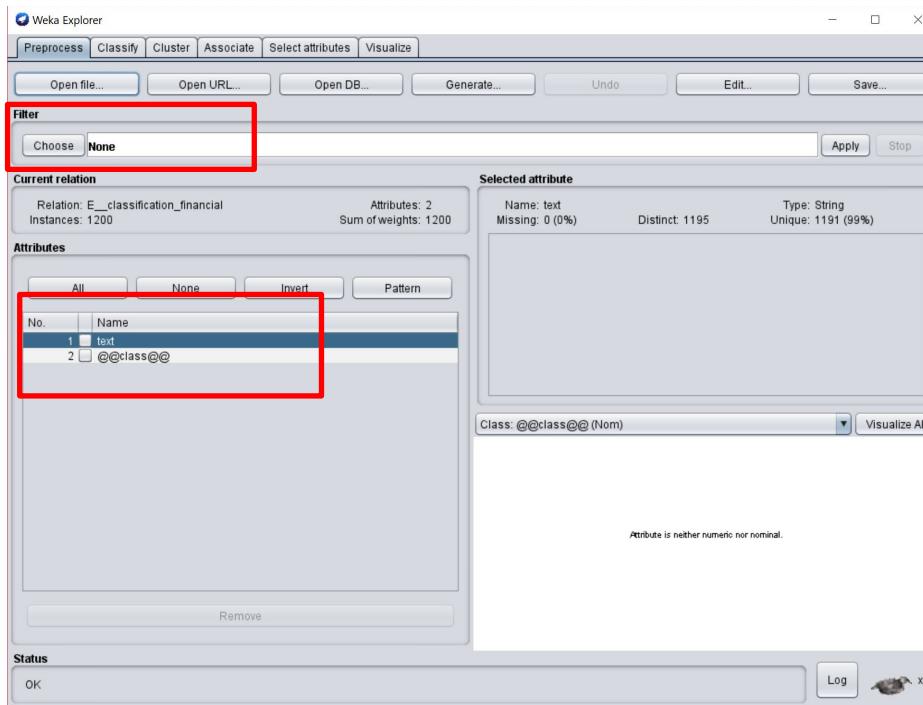
- In order to measure the real accuracy of your classifier you need to test it on unseen examples (those are examples that were not included in the training process and that the algorithm doesn't have a prior knowledge about (?)).
- To be able to do so we need to create a classification model.

CREATE A CLASSIFICATION MODEL

- WEKA requires that you apply the same filter to training and testing datasets when extracting features.
- As mentioned before, you must use **identical features** for both training and testing as otherwise the algorithms may not work.
- To do this we need to create a “**filtered classifier**”. That is a classifier combined with a filter to extract the same features from both the **training** and **testing** datasets.

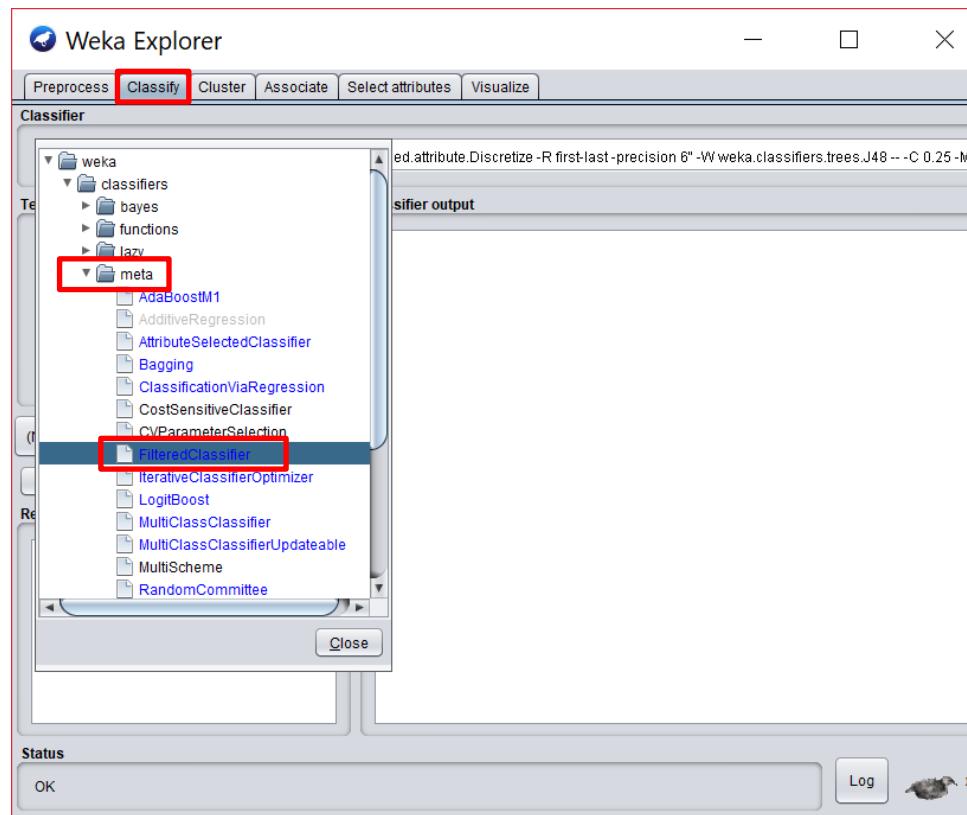
CREATE A CLASSIFICATION MODEL

- Start by loading your training arff file (**finance.arff**).
- This time **do not apply** the StringToWordVector filter, yet!



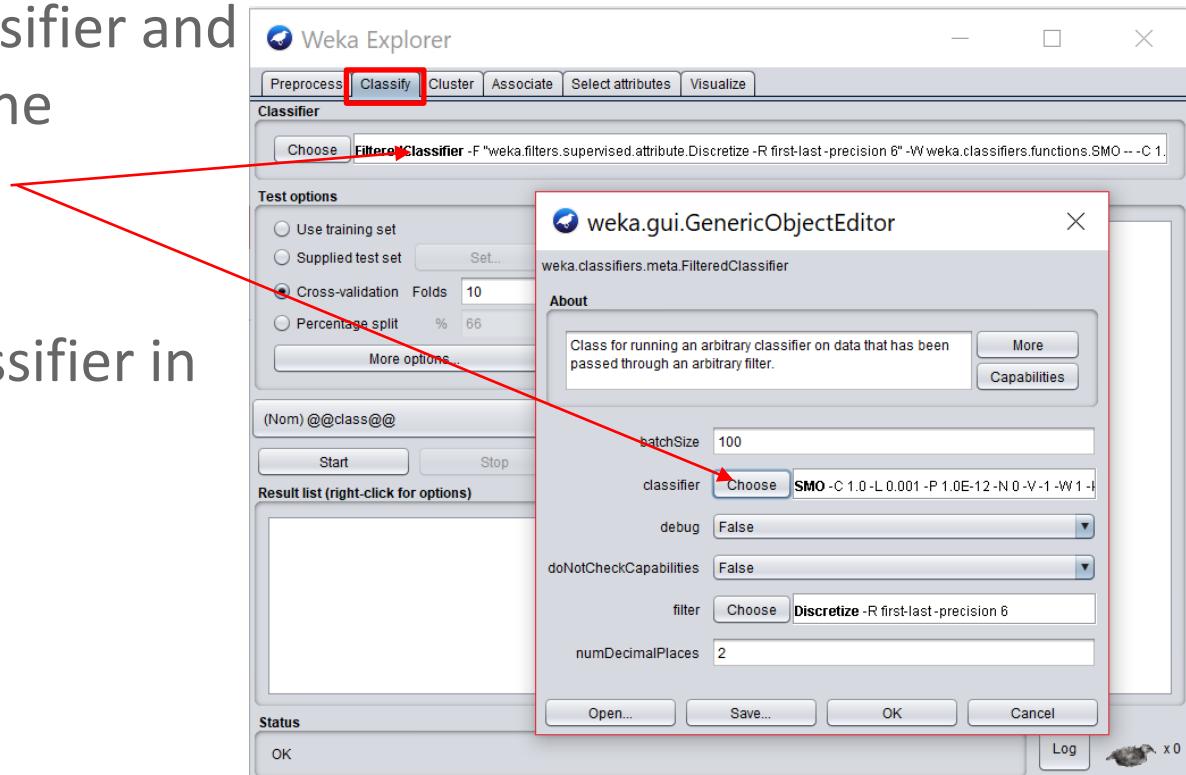
CREATE A CLASSIFICATION MODEL

- Move on to the Classify Tab and choose the “Filtered Classifier” from the Classifier Chooser menu.



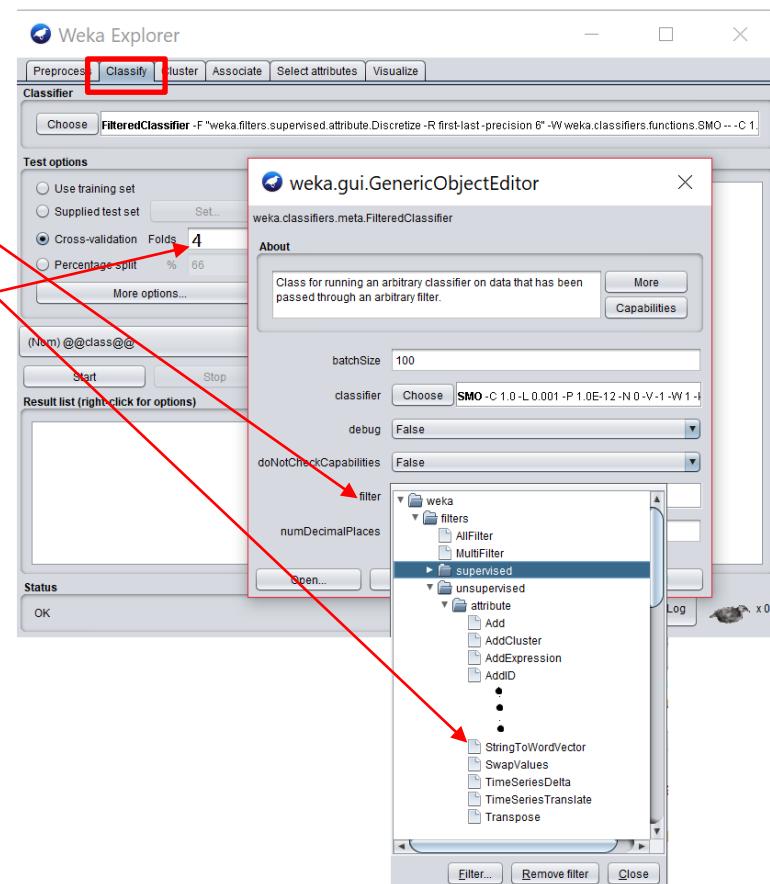
CREATE A CLASSIFICATION MODEL

- Click on FilteredClassifier and Choose SMO from the functions classifiers.
- SMO is the SVM classifier in WEKA.



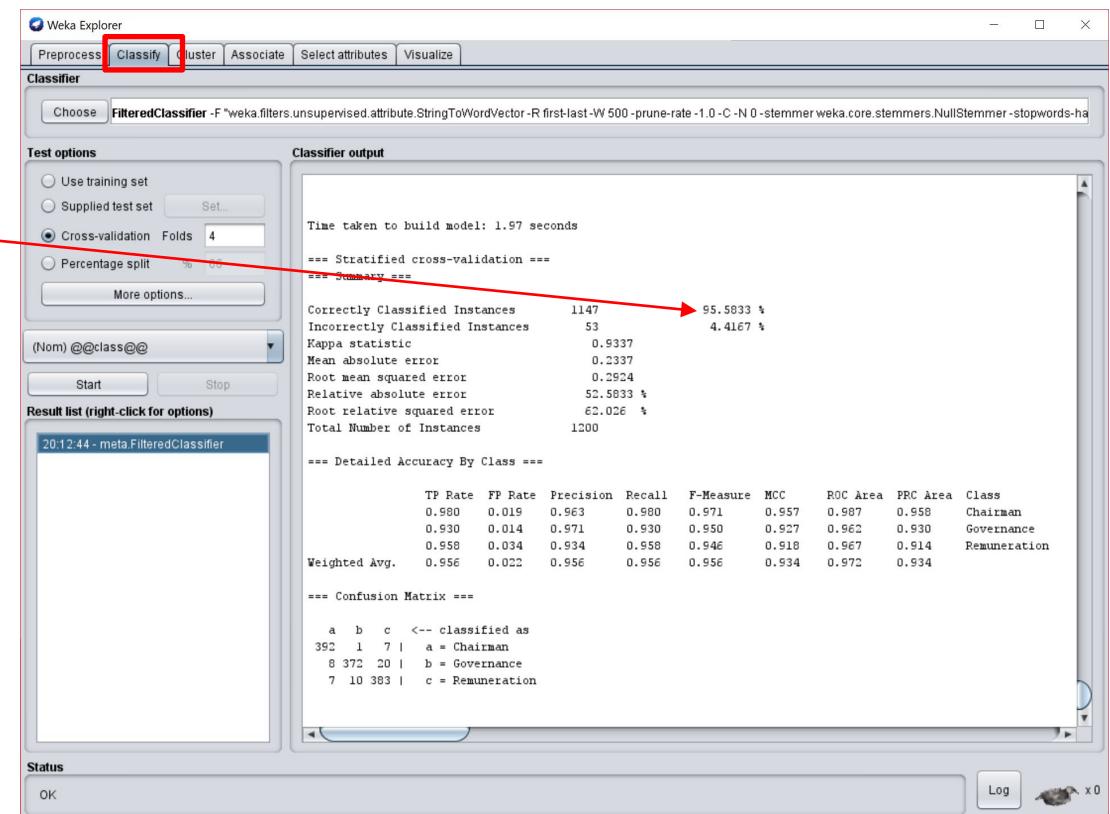
CREATE A CLASSIFICATION MODEL

- Click on filter to choose StringToWordVector (keeping 500 alphabetical words)
- Click Ok
- Choose 4 Cross Validation Folds
- Click Start



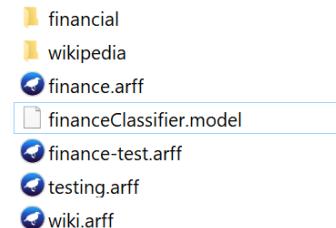
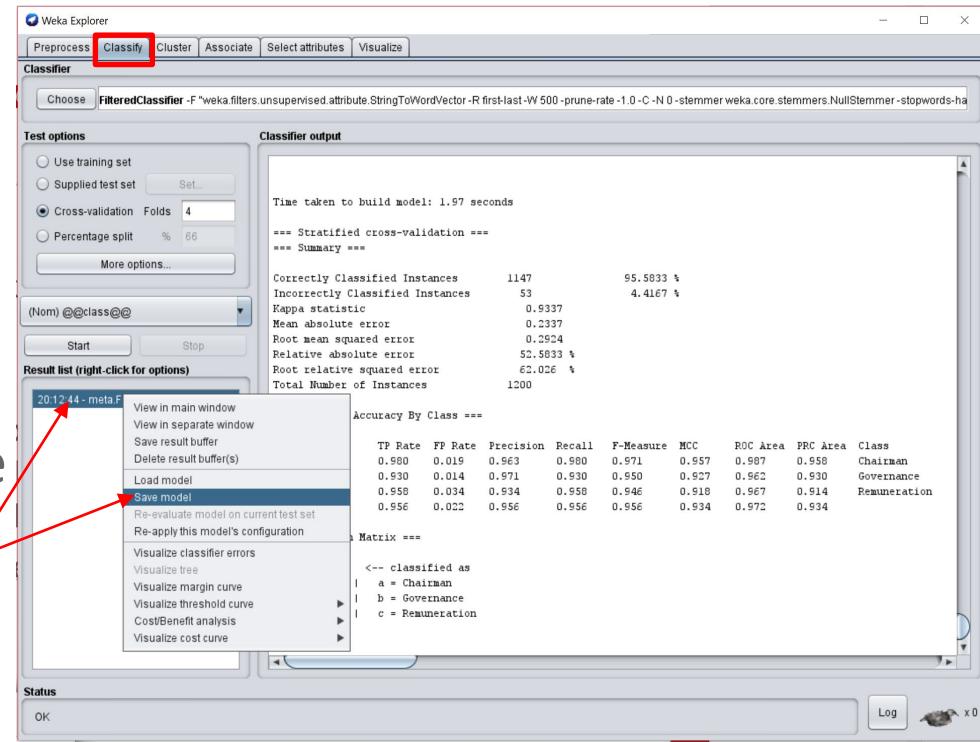
CREATE A CLASSIFICATION MODEL

- Classification Results will show in the Classifier output window



CREATE A CLASSIFICATION MODEL

- Right click the SMO Filtered classifier and click on Save model
- Save the model to your machine next to the arff file so you can find it later on.
- Call it aaerClassifier and click Save.



STEPS

1. Convert dataset into ARFF file format.
2. Prepare dataset and extract features
3. Train a classifier using training data
4. Create a classification model
5. **Test the classification model**



FINANCIAL TESTING DATA

- To save your time I created a testing file for you



finance-test.arff

1-5: chairman statements
6-10: governance statements
11-15: remuneration statements

```
@relation E__classification_financial

@attribute text string
@attribute @@class@@ {Chairman,Governance,Remuneration}

@data

'Euan Worthington Paul LoudonChairman Chief Executive Officer5',?
'outlook with optimism.Otto SchmidChairman 10 October 2002',?
'Financial statements Other information Overview Strategic report',?
```

TEST THE CLASSIFICATION MODEL

- Now that we have a SVM Classifier (Model) we want to test it on unseen data.
- We'll use the **finance-test.arff** file.
- The classifier doesn't know the output and only you know so you can judge the classifier's real accuracy.

FINANCIAL TESTING DATA

- To save your time I created a testing file for you



finance-test.arff

1-5: chairman statements
6-10: governance statements
11-15: remuneration statements

```
@relation E_classification_financial

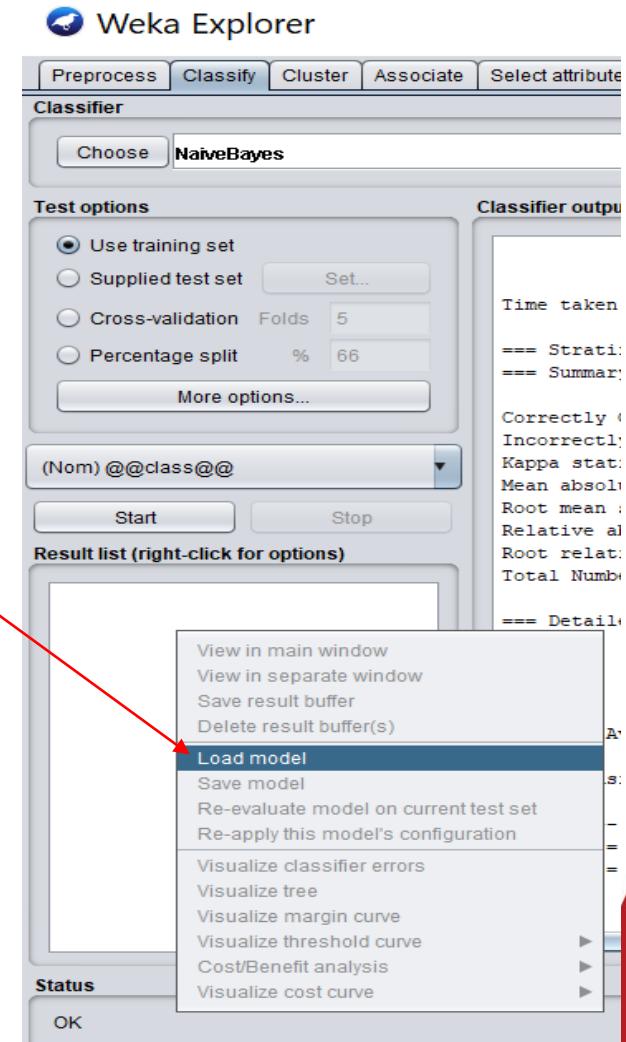
@attribute text string
@attribute @@class@@ {Chairman,Governance,Remuneration}

@data

'Euan Worthington Paul LoudonChairman Chief Executive Officer5',?
'outlook with optimism.Otto SchmidChairman 10 October 2002',?
'Financial statements Other information Overview Strategic report',?
```

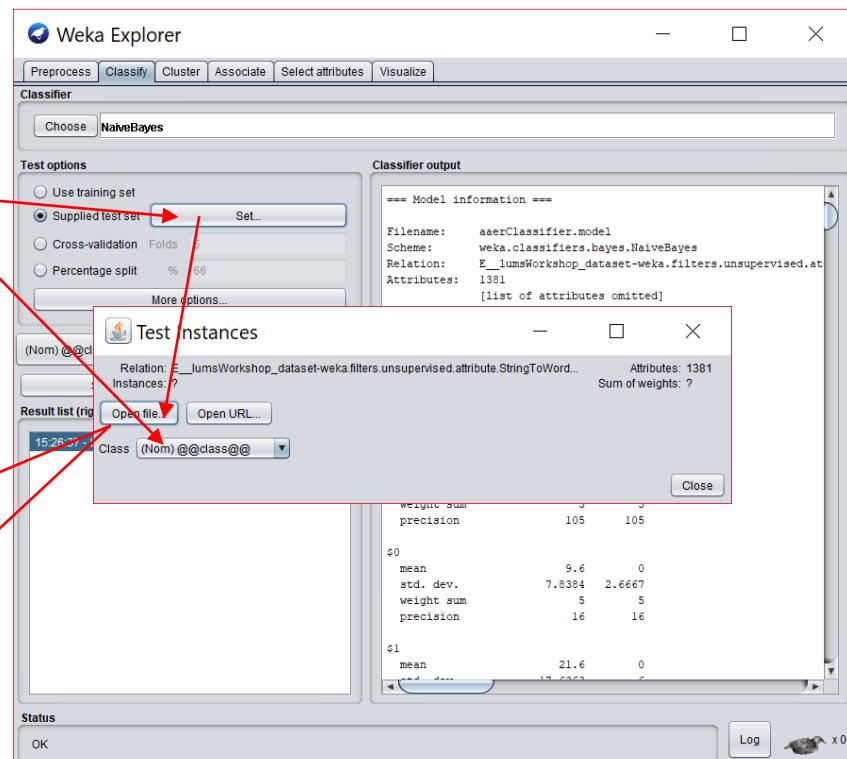
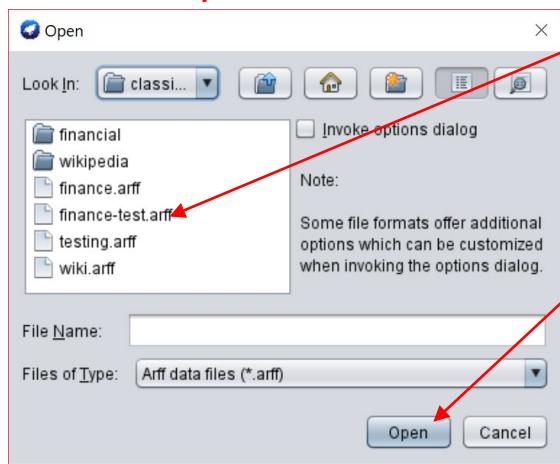
TEST THE CLASSIFICATION MODEL

On Weka Classify right click in the result list to load your financeClassifier.model (only if it's not already there).



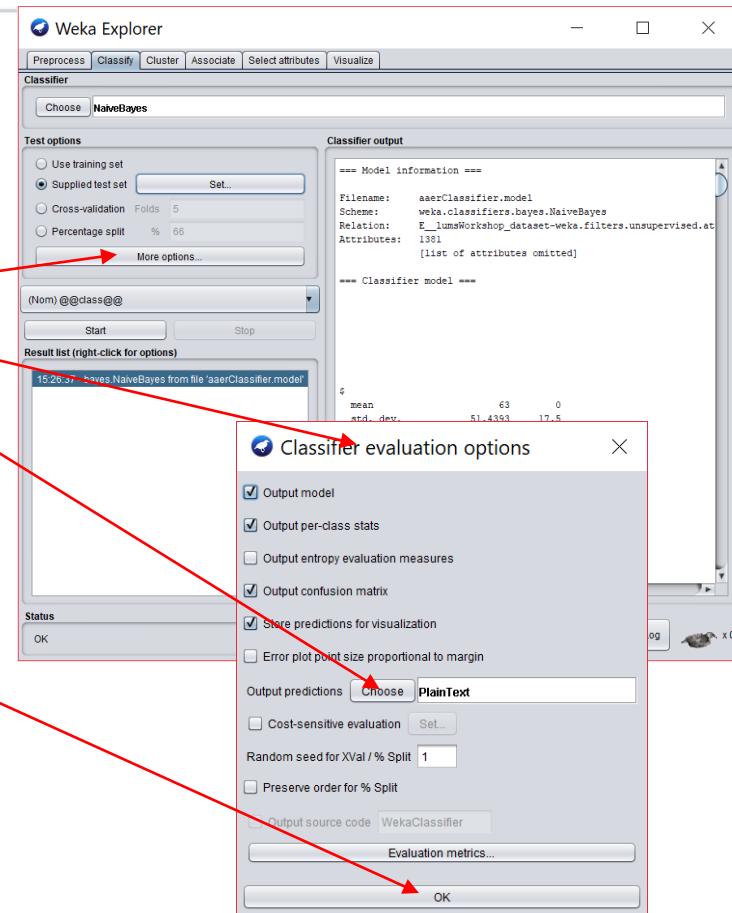
TEST THE CLASSIFICATION MODEL

- Now you need to load your testing sample (finance-test.arff).
- Choose Supplied test set
- Make sure (Nom) @@class@@ is selected
- Open your testing arff file and click open then close the dialogue.
- If you get an error message here then your training and testing arff files are incompatible.



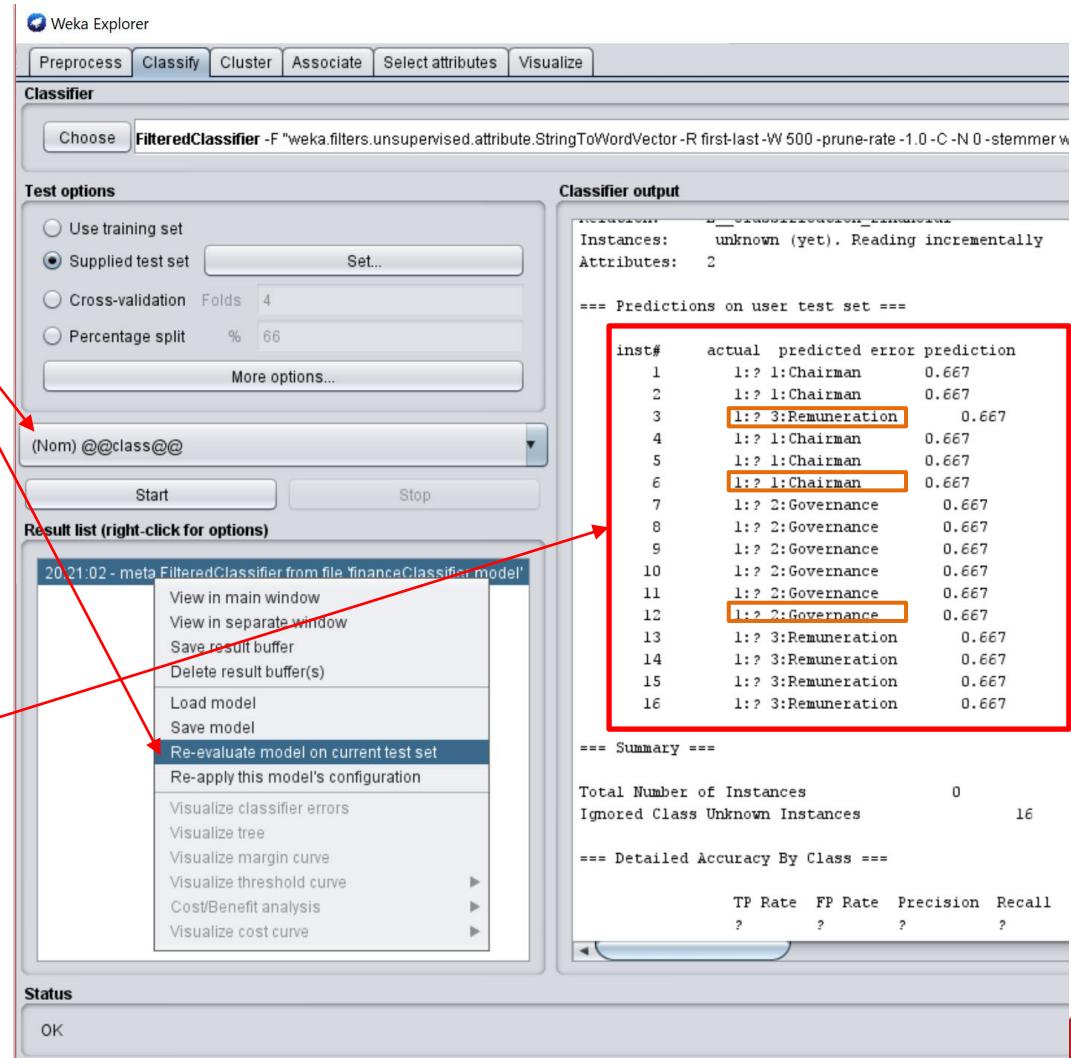
TEST THE CLASSIFICATION MODEL

- Click on More options..
- Make sure Output predictions is PlainText so you can see the output as Chairman, Governance or Remuneration.
- Then click OK



TEST THE CLASSIFICATION MODEL

- Make sure (Nom) @@@class@@@ is selected
- Right Click your model and click Re-evaluate model on current test set.
- You'll notice the predictions in the Classifier output window.





TRY YOURSELF
