# Machine learning/ Natural language processing

**INFO319 – Research Topics in Big data**

Vimala Nunavath

vimala.nunavath@uia.no

**04.10.2019**
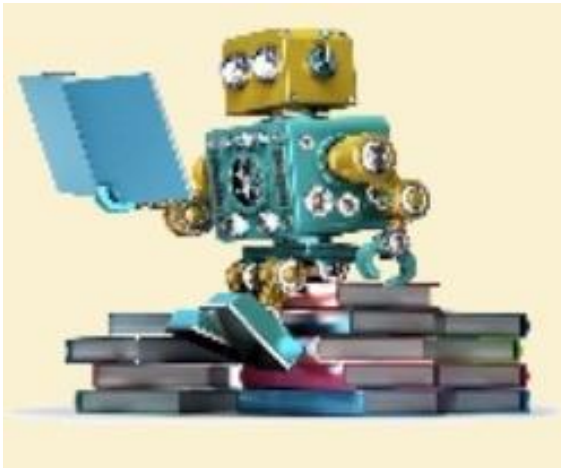
# Agenda

- Machine learning

- Natural language processing

- Presentations by you.

- Practical session

# Programming project

- Deadline is today!!

INFO319, autumn 2019, session 6

# Machine learning

# Machine learning



"Programming computer to optimize the performance using example data and past experience"

# Machine learning



Field of study that gives "computers the ability to learn without being explicitly programmed"

---- Arthur Samuel, 1959

# Machine learning

- Subset of Artificial Intelligence
- Design and development of algorithms
- Computers evolve behavior based on empirical data

# Machine learning

- The **input data** to a machine learning system can be **numerical, textual, audio, visual, or multimedia.**

- The corresponding **output data** of the system can be a **floating-point number.**

- **E.g., -** the velocity of a rocket, an integer representing a category or a class,

  - a pigeon or a sunflower from image recognition.

# Machine learning - Applications

- Recommend friends, dates and products to end-users.

# Machine learning - Applications

- Classify content into predefined groups

# Machine learning - Applications

- Identify key topics in large collections of text

# Machine learning - Applications

- Natural language processing

# Machine learning - Applications

- Detect Anomalies within given data

- Ranking Search Results with user feedback learning

- Classifying DNA sequences

- Sentiment analysis/ opinion mining

- Bioinformatics

- Speech and Handwriting Recognition

# Machine learning types

Machine learning

Supervised
Deals with labeled training data, learn to map inputs to outputs

Unsupervised
Deals with unlabeled data, observe structure in the data and find patterns.

Reinforcement learning
Dynamic environment, perform a certain goal

# Machine learning types

# Machine learning - Classification

# Classification

# Machine learning - Classification

Check Email

Spam? → No

Yes

**SPAM DETECTION**

New email → Spam classifier system → Spam / Non spam

# Machine learning - Regression

# Machine learning - Regression

- Regression models are used to predict a continuous value.
  - E.g., predicting prices of a house given the features of house like size, price etc

- In liner regression, we draw all possible lines going through the points such that it is closest to all.

Linear Regression Example Data

| House Price in $1000s (Y) | Square Feet (X) |
| --- | --- |
| 245 | 1400 |
| 312 | 1600 |
| 279 | 1700 |
| 308 | 1875 |
| 199 | 1100 |
| 219 | 1550 |
| 405 | 2350 |
| 324 | 2450 |
| 319 | 1425 |
| 255 | 1700 |

Statistics for Managers Using Microsoft
Excel, 5e © 2008 Prentice-Hall, Inc.

Chap 13-1

# Machine learning types

# Machine learning - Clustering

- **Clustering** is the task of dividing the population or data points into a number of groups such that data points in the same groups are more similar to other data points in the same group and dissimilar to the data points in other groups.

- It is basically a collection of objects on the basis of similarity and dissimilarity between them.

# Machine Learning Workflow

# **Natural language processing**

# What is Natural language processing?

- Shortened as NLP,

- **Field in ML** that **deals with** the **interaction between computers** and **humans** using the **natural language**.

- The ultimate objective of NLP is to **read**, **decipher**, **understand**, and **make sense of the human languages** in a manner that is valuable.

- Most NLP techniques rely on machine learning to derive meaning from human languages.

# The goal of NLP

- The goal of NLP is to "accomplish human-like language processing".

- A full NLP System would be able to:

    1. Paraphrase an input text

    2. Translate the text into another language

    3. Answer questions about the contents of the text

    4. Draw inferences from the text

# Why NLP?

- Aids communication between two humans

  - Machine translation

  - Speech-to-speech translation

  - Speech-to-text & text-to-speech

  - Editorial aids (spelling & grammar checkers)

  - Aids communication between human and machine

  - Personal assistants

  - Interactive Voice Response systems

  - Aids communication between two machines

# Why NLP for Social Media?

- Social Media generates BIG UNSTRUCTURED NATURAL LANGUAGE DATA

  - Volume: 1.3 Billion monthly active FB users

  - Velocity: 5700 Tweets/sec. 2500 FB-msg/sec

  - Variety: scripts, languages, style, topic, …

- Today's world resides in social media

- It is impossible to process (consume, understand or summarize) this information manually.

# Why NLP?

- Trending Topic Detection

- Information Retrieval & Extraction

- Information Summarization

- **Sentiment Detection**

- Rumor Detection

# Challenges of NLP

- Traditionally, NLP systems are designed to handle input that is

  - Grammatically correct

  - No spelling errors

  - Single language

  - The right script

# Sentiment Analysis

- Sentiment analysis refers to the class of computational and natural language processing-based techniques **used to identify, extract or characterize subjective information**, such as opinions, expressed in a given piece of text.

- The main purpose of sentiment analysis is to classify a writer's attitude towards various topics into positive, negative or neutral categories.

# Sentiment Analysis

- Sentiment analysis has many applications in different domains e.g., business intelligence, politics, sociology, social networking websites, microblogs, wikis and Web applications

- Data such as web-postings, Tweets, videos, etc., all express opinions on various topics and events, offer immense opportunities to study and analyze human opinions and sentiment.

# Latent Dirichlet Allocation (LDA)

- It is a "generative probabilistic model" of a collection of composites made up of parts.

- Here composites are documents and the parts are words and/or phrases.

- Allows a word to simultaneously belong to several clusters with varying degree.

# Latent Dirichlet Allocation (LDA)

- **Latent**: refers to everything that we don't know a priori and are hidden in the data.
  - E.g., the themes or topics that document consists of are unknown, but they are believed to be present as the text is generated based on those topics.

- **Dirichlet:** It is a 'distribution of distributions.
  - E.g.,in the context of topic modeling, the Dirichlet is the distribution of topics in documents and distribution of words in the topic.

- **Allocation**: This means that once we have Dirichlet, we will allocate topics to the documents and words of the document to topics.

# Spark MLib

INFO319, autumn 2019, session 6

# Spark ecosystem

# MLib

- It is Spark's library of machine learning functions designed to run in parallel on clusters.

- Contains variety of learning algorithms

- MLib invokes various algorithms on RDDs

- Soma basic ML algorithms are not included with Spark Mlib as they are not designed for parallel.

# Spark MLib

# Spark MLlib Overview

- Divided into two packages:

    - *spark.mllib* contains the original API built on top of RDDs.

    - *spark.ml* provides higher level API built on top of DataFrames.

- Using spark.ml is recommended because with DataFrames, the API is more versatile and flexible.

# MLlib Algorithms

- The popular algorithms and utilities in Spark MLlib are:
  - Basic Statistics
  - Regression
  - Classification
  - Clustering
  - Recommendation System
  - Dimensionality Reduction
  - Feature Extraction

# MLlib Algorithms - Basic Statistics

- Basic Statistics includes the most basic machine learning techniques.

  - **Summary Statistics**: provides column summary statistics for RDD[Vector]
    - e.g., include mean, variance, count, max, min and numNonZeros.
  - **Correlations**: pairwise correlations, Spearman and Pearson are some ways to find correlation.
  - **Stratified Sampling**: two methods: sampleBykey and sampleByKeyExact.
  - **Hypothesis Testing**: Pearson's chi-squared test is an example of hypothesis testing.
  - **Random Data Generation**: RandomRDDs, Normal and Poisson are used to generate random data.

# MLlib Algorithms - Recommendation System

- A recommendation system provides suggestions to the users through a filtering process that is based on user preferences and browsing history.

- A recommendation system is a subclass of **information filtering system** that seeks to predict the "rating" or "preference" that a user would give to an item.

- Recommender systems have become increasingly popular in recent years and are utilized in a variety of areas.
  - e.g., movies, music, news, books, research articles, search queries, social tags, and products in general.

# MLlib Algorithms - Recommendation System



COLLABORATIVE FILTERING

Read by both users

Similar users

Read by her,
recommended to him!

CONTENT-BASED FILTERING

Read by user

Similar articles

Recommended
to user

# MLlib Algorithms - Dimensionality Reduction

- Dimensionality Reduction is the process of reducing the number of random variables under consideration, via obtaining a set of principal variables. It can be divided into **feature selection and feature extraction.**

  - **Feature Selection:** Feature selection finds a subset of the original variables (also called features or attributes).
  - **Feature Extraction:** transforms the data in the high-dimensional space to a space of fewer dimensions.

# MLlib Algorithms - Feature Extraction



(Image reference O'Reilly Learning Spark)

# MLlib Algorithms - Feature Extraction

- Feature Extraction starts from **an initial set of measured data** and builds derived values (features) intended to be informative and non-redundant, facilitating the **subsequent learning** and generalization steps, and in some cases leading to better human interpretations.

- Algorithms for working with features, roughly divided into these groups:
  1. Extraction: Extracting features from "raw" data
  2. Transformation: Scaling, converting, or modifying features
  3. Selection: Selecting a subset from a larger set of features
  4. Locality Sensitive Hashing (LSH): This class of algorithms combines aspects of feature transformation with other algorithms.

# Feature Extraction and transformation

- It contains 5 techniques:

    - Term frequency (TF) and inverse document frequency (IDF)

    - **Word2Vec (**skip-gram and continuous bag of words**)**

    - Standard scaler

    - Normalizer

    - Chi-square Selector

# Word2Vec- display the top 40 synonyms of the specified word.

```python
from pyspark import SparkContext
from pyspark.mllib.feature import Word2Vec

sc = SparkContext(appName='Word2Vec')
inp = sc.textFile("text8_lines").map(lambda row: row.split(" "))

word2vec = Word2Vec()
model = word2vec.fit(inp)

synonyms = model.findSynonyms('china', 40)

for word, cosine_distance in synonyms:
    print "{}: {}".format(word, cosine_distance)
```

# Sentiment Analysis in Disaster Relief

- Sentiment analysis of disaster related posts in social media is one of the techniques that could gear up **detecting posts for situational awareness.**

- Useful to better **understand** the dynamics of the network including **users' feelings, panics and concerns** as it is **used to** identify polarity **of sentiments** expressed by users during disaster events to improve decision making.

- It **helps** authorities to **find answers** to their questions and make better decisions regarding the event assistance without paying the cost as the traditional public surveys.

# Sentiment Analysis in Disaster Relief

- Sentiment information could also be **used to project** the information regarding **the devastation and recovery situation** and **donation requests to the crowd** in better ways.

- Using the results obtained from sentiment analysis, **authorities can figure out** where they should look for particular information regarding the disaster such as the most affected areas, types of emergency needs.

# Use case – Sentiment analysis from Twitter data

# Sentiment analysis

- Sentiment refers to the emotion behind a social media mention online.

- Sentiment analysis is categorizing the tweets related to particular topic and performing data mining using sentiment automation analytics tools.

- We will be performing Twitter Sentiment analysis as our use case for Spark streaming.

- Sentiment analysis helps in disaster management.

# Problem statement

- To design a Twitter Sentiment analysis system where we populate real time sentiments for disaster management.

- Sentiment analysis is used to:

    - Identify, extract or characterize subjective information, such as opinions, expressed in a given piece of text.

    - Classify a writer's attitude towards various topics into positive, negative or neutral categories.

    - To study and analyze human opinions and sentiment.

# Performance metrics for classification

- Confusion Matrix

- Accuracy

- Precision

- Recall or Sensitivity

- Specificity

# Confusion matrix



Confusion Matrix

**False Positives (FP):** False positives are the cases when the actual class of the data point was 0(False) and the predicted is 1(True). False is because the model has predicted incorrectly and positive because the class predicted was a positive one. Ex: A person NOT having cancer and the model classifying his case as cancer comes under False Positives.

**True Positives (TP):** True positives are the cases when the actual class of the data point was 1(True) and the predicted is also 1(True)
*Ex: The case where a person is actually having cancer(1) and the model classifying his case as cancer(1) comes under True positive*

**True Negatives (TN):** True negatives are the cases when the actual class of the data point was 0(False) and the predicted is also 0(False).
*Ex: The case where a person NOT having cancer and the model classifying his case as Not cancer comes under True Negatives.*

**False Negatives (FN):** False negatives are the cases when the actual class of the data point was 1(True) and the predicted is 0(False). False is because the model has predicted incorrectly and negative because the class predicted was a negative one. (0)
Ex: A person having cancer and the model classifying his case as No-cancer comes under False Negatives.

# Accuracy

- **When to use Accuracy:**

Accuracy is a good measure when the target variable classes in the data are nearly balanced. Ex:60% classes in our fruits' images data are apple and 40% are oranges.
A model which predicts whether a new image is Apple or an Orange, 97% of times correctly is a very good measure in this example.

- **When NOT to use Accuracy:**

Accuracy should NEVER be used as a measure when the target variable classes in the data are a majority of one class.



$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN}$$

Accuracy

# Precision

Precision is a measure that tells us what proportion of patients that we diagnosed as having cancer, actually had cancer. The predicted positives (People predicted as cancerous are TP and FP) and the people actually having a cancer are TP.



$$\text{Precision} = \frac{TP}{TP + FP}$$

Precision

# Recall or Sensitivity

Recall is a measure that tells us what proportion of patients that actually had cancer was diagnosed by the algorithm as having cancer. The actual positives (People having cancer are TP and FN) and the people diagnosed by the model having a cancer are TP. (Note: FN is included because the Person actually had a cancer even though the model predicted otherwise).



Recall or Sensitivity

# Specificity

Specificity is a measure that tells us what proportion of patients that did NOT have cancer, were predicted by the model as non-cancerous. The actual negatives (People actually NOT having cancer are FP and TN) and the people diagnosed by us not having cancer are TN. (Note: FP is included because the Person did NOT actually have cancer even though the model predicted otherwise).



$$\text{Specificity} = \frac{TN}{TN + FP}$$

# Use Case – Importing packages

```
//Import the necessary packages into the Spark Program
import org.apache.spark.streaming.{Seconds, StreamingContext}
import org.apache.spark.SparkContext._
import org.apache.spark.streaming.twitter._
import org.apache.spark.SparkConf
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._
import org.apache.spark._
import org.apache.spark.rdd._
import org.apache.spark.rdd.RDD
import org.apache.spark.SparkContext._
import org.apache.spark.sql
import org.apache.spark.storage.StorageLevel
import scala.io.Source
import scala.collection.mutable.HashMap
import java.io.File
```

# Use case- Twitter token authorization

```scala
object mapr {

  def main(args: Array[String]) {
    if (args.length < 4) {
      System.err.println("Usage: TwitterPopularTags <consumer key>
<consumer secret> " +
      "<access token> <access token secret> [<filters>]")
      System.exit(1)
    }

    StreamingExamples.setStreamingLogLevels()
    //Passing our Twitter keys and tokens as arguments for authorization
    val Array(consumerKey, consumerSecret, accessToken,
accessTokenSecret) = args.take(4)
    val filters = args.takeRight(args.length - 4)
```

# Use case- Dstream transformation

```scala
// Set the system properties so that Twitter4j library used by twitter stream
// Use them to generate OAuth credentials
System.setProperty("twitter4j.oauth.consumerKey", consumerKey)
System.setProperty("twitter4j.oauth.consumerSecret", consumerSecret)
System.setProperty("twitter4j.oauth.accessToken", accessToken)
System.setProperty("twitter4j.oauth.accessTokenSecret",
accessTokenSecret)

val sparkConf = new
SparkConf().setAppName("Sentiments").setMaster("local[2]")
val ssc = new StreamingContext(sparkConf, Seconds(5))
val stream = TwitterUtils.createStream(ssc, None, filters)

//Input DStream transformation using flatMap
val tags = stream.flatMap { status =>
status.getHashtagEntities.map(_.getText)}
```

# Use case – Generating tweet data

```scala
//RDD transformation using sortBy and then map function
tags.countByValue()
 .foreachRDD { rdd =>
 val now = org.joda.time.DateTime.now()
 rdd
 .sortBy(_._2)
 .map(x => (x, now))
 //Saving our output at ~/twitter/ directory
 .saveAsTextFile(s"~/twitter/$now")
 }

//DStream transformation using filter and map functions
val tweets = stream.filter {t =>
 val tags = t.getText.split("
").filter(_.startsWith("#")).map(_.toLowerCase)
 tags.exists { x => true }
}
```

# Use case – Extracting Sentiments

```scala
val data = tweets.map { status =>
val sentiment = SentimentAnalysisUtils.detectSentiment(status.getText)
val tagss = status.getHashtagEntities.map(_.getText.toLowerCase)
(status.getText, sentiment.toString, tagss.toString())
}

data.print()
//Saving our output at ~/ with filenames starting like twitterss
data.saveAsTextFiles("~/twitterss","20000")

ssc.start()
ssc.awaitTermination()
 }
}
```

# Use case - Results



Figure: Output file containing tweet and its sentiment

# Conclusions

- Machine learning

- Natural language processing

- Spark MLib

- Usecase