# SYSC 4700 - Lecture 10 Carleton University

# Part 1: Introduction to Big Data Analytics Part 2: Introduction to Machine Learning

Presenter: André Brandão abran012@uottawa.ca February 08, 2018

email: abran012@uottawa.ca

### **SYSC 4700 - Lecture 10 - Part 1**

# **Introduction to Big Data Analytics**

#### We will discuss:

- Basic concepts
- Examples of analytics
- Tapping into the database
- Techniques for data aggregation



#### **Basic Concepts**

- Data Analytics refers to the extraction of knowledge from large data sets. A large data set contains information about many things.
- The analyst job is to acquire knowledge by linking information that results in a theoretical or practical understanding about something.



#### **Knowledge, more than Information**

An exercise in analytics: **W**hat makes a fisherman successful? You collected lots of data during a fishing season: boat type, nets, bait type, water temperature, wind, cloud cover, water turbidity, time of the day, number of people per boat, jacket colours, fishing rod length, etc. Here you are with a nice table full of information.



Is that possible fish can see people on the boat when they wear yellow?

# Your table gets bigger

An exercise in analytics: **W**ell, you decided to move on and collect data about what is happening on the fishermen's village: the village size, number of churches, stores, theatres, schools, cars, fish processing facilities, transportation, restaurants, sport's arena, community centre, gender distribution, income distribution, family sizes, language and culture in general, etc.

Wind	Jacket Colour	# of Fish	People/boat	Villages	Stores	Churches
5 knots	5 knots blue		3	6400	8	3
N/A	blue	8	2	9000	12	3
7	red	3	4	568	2	1
2	yellow	6	6	N/A	-	-
8 brown		7	2	125	1	-

You research the stores in small villages selling fishermen's clothing. Affordable clothing are in blue, all the same brand, made in China.

# **Bigger table, bigger picture**

<u>An exercise in analytics</u>: **F**inally, the extra information (i.e. facts) tells you: villagers in that region live within medium to low income communities and buy at local stores. Those who caught more fish wear the same jacket brand from local stores.

v	Vind	Jacket Colour	Fish Caught	People per boat	Village size	Stores	Churches	Schools	Income Distribution
51	nots	blue 🚽	- +2 -		12000	32	4	5	
06/	06/17	blue	8	8	9000 🥆	26	6	° —	
	7	red	3	2	250000	>100	25		
-	2	yellow		1	300000	>100			
	WA.	blue	7	5	N/A	N/A	A second evaluation shows		
	8	red	3	3	56000	>			
	2	red	2	3	98000		those blue jacket fishermen		
		yellow	5	1	120000		are locals, from small villages.		
red		red		3			a. e . o caio)		aBesi

You have not concluded what makes a good fisherman yet. You started by noting some blue-jacket guys wearing the same brand, from the same store, catching more fish.

# **Bigger Picture poses a Dilemma**

<u>An exercise in analytics</u>: **A**fter evaluating the local environment as well as all the technical aspects of fishing you encounter a phenomenal Aristotelian dilemma to your analysis:

For any hypothesis there can be an infinite number of thesis and no scientific method exists that can find deterministically a thesis that proves a certain hypothesis.

This means: as your database increases, you gain more "insight", but at the same time the number of variables increase taking with it the number of possibilities for a certain cause-effect event.

### **Analytics and Intuition**

An exercise in analytics: At the end, you use your Intuition.



And finally you may conclude:

Local fishermen know best how to catch fish in their own backyards.

# **Modelling and Prediction**

We have shown the fishermen story and how the analyst acquired some *awareness* about the village and that local environment.

Apart of awareness, the analyst may also **predict** behaviour if the amount or "quality" of information in the database is relevant.

The history of *modelling*, *prediction* and *pattern recognition* within *Information Theory* is not recent. It gained force with giants like Widrow, Kalman, Abramson and others during the 1960's and matured in the 1980's with the advent of affordable digital signal processors and computers.

**T**oday, the *prediction of customer behaviour is big business in the Internet*. Banks also use prediction techniques for investors in the stock market.

# **Modelling and Prediction**

**P**rediction requires a model. A model in telecommunications is the mathematical transfer function of a system that outputs the desired response when excited with a proper signal.



When  $error \rightarrow 0$  the model reconstructs s(n) by using the information contained in the past series  $s(n-\tau)$ .

# **Constructing a Model**

There are many ways to construct a prediction model for big data analytics. One particular technique, based on the LMS (least-mean-square), is known as the "steepest descent" adaptive algorithm.

The steepest descent is a learning algorithm that adapts its model dynamically as it is fed by a numeric sequence.



For example, consider the following numeric sequence:

André L. Brandão, M.Sc., Ph.D. - Research Engineer - SYSC4700 Carleton University - Feb.2018

# **Constructing a Model**

**G**iven 800 samples (source is your database), what would be the value of sample 801? or 802 and so on.

As the predictor acquires data its engine works out to minimize the prediction error. When the error is minimized the resulting model is ready.



# **Constructing a Model**

The model gives you a mathematical description for the system that generates the original sequence s(n).

In this example our learning algorithm converges into a polynomial which represents the prediction model (transfer function H(z)) below:



#### **Adaptive Linear Prediction**

**B**elow is the code that I used for the learning algorithm example. It generates a numeric sequence (suppose this is from a database) and applies the algorithm to find the prediction model H(z). You can play with the parameters and see how the model changes.





#### **Adaptive Linear Prediction**



#### **Database in the Cloud**

Access to big data is normally done in the internet cloud. Many commercial companies offer data repositories over the cloud with servers that respond to SQL (Structured Query Language) enquires from client terminals.

For example, suppose you have data stored in the cloud (i.e. Amazon, Google, etc.) and want to retrieve information about fishermen wearing blue jackets with records of catching lots of fish. SQL is powerful and intuitive. The client would send a message to the SQL database like this:

#### "Select from fisherman\_table where jacket\_colour = 'blue' and fish\_caught >6"

SQL queries allow you to access big tables and retrieve from the cloud only what is necessary. This saves you local memory space and searching time. An excellent site that allows you to practice with SQL statements is this: http://www.w3schools.com/sql/sql\_select.asp

# **Working with Big Data**

Although SQL selects only those necessary fields you want from the database, when dealing with large data sets even a few fields may contain more

information than the memory of your local computer can handle.

Minimize reading. (bring data)

Minimize writing. (push data)



#### **The Power of Parallelism**

Working with virtual machines (VM), physically located close to your data repository, is the way to go for cluster or parallel computing in the cloud today.

Spark is one tool that offers parallelism for big data analytics. Spark is an open source software framework that evolved from Apache Hadoop.

see:

http://spark.apache.org/



#### **How Spark Works**



#### **How Spark Works**



#### **How Spark Works**



# **Concluding Remarks**

**B**ig data analytics is part of the digital economy. It is present in all aspects of everyday life. Examples of how diverse this is include:

- Government agencies (statistics, revenue, service delivery, etc.)
- Retailers (understanding customer behaviour)
- Business processes (optimize warehouse stock, deliveries)
- Health care (clustering diseases, etc.)
- Sports (game stats, crowd attendance, etc.)
- Science and research (molecular data analysis, genetics, etc.)
- Safety and Security (clustering for safety and security, etc.)
- Banks (stock markets, etc.)

Apart of Spark there are other frameworks and development platforms devoted to Big Data Analytics (Cloudera, Hortonworks, Databricks, direct SQL, etc.).