**DATA***lysis* 



# When Machine Learning meets Graph Databases



💓 @G\_Ceresa

W: www.datalysis.ch

E: info@datalysis.ch

### **Gianni Ceresa**

Managing Director of DATAlysis GmbH (Switzerland)

Working with Business Analytics, EPM tools and "data" for more than 10 years Oracle ACE Director

▲ DATAlysis

Part-time blogger on **gianniceresa.com** 

Full-time IRC (freenode | #obihackers) resident

Same group on Telegram http://telegram.me/obihackers

ODC (ex OTN) forums addict

Technology geek (or just geek in general)

		Gianni's v@rld	Home About Glanni 💭	year's What's your New
Dverview Activity	BROWSE Intelligence ~ Content People Subspaces	Scripts to use Agents in OBIEE 12c: mass enabling and execution from a page by Gami Ceresa   15 March 2017   Under the hood Few days ago I worke about Agents in OBIEE 12c and how it was possible to enable them by script, automating that part of the process which was still manual. This time i'm going to write about 2 examples of how I used my findings about Agents, providing the code	Bearch Recent Posts Scripts to use Agents in OBITE 12c mass enabling and execution from a page OBITE 12c Agents: enable them by code I OBITE 12c Lotom Style using	Virtual Syddi Heeve Some Koas for 2017 b) c) iupdate to OBIE 12c iupdat
CATEGORIES	ALL CONTENT (2765) 📃 BLOG POSTS (1) 🛅 DOCUMENTS (5) 💬 DISCUSSIONS (2762) 🚮 POLL Filter by action: [None 🔹 🖛 Filter by shared content	OBIEE 12c Agents: enable them by code ! by Gauni Ceresa   13 March 2017   Under the hood	shared folder (analyticsRes) Sequence Numbers for Time Dimensions: new in OBIEE	
ACTIONS	Type to filter by text FILTERY DG Sort by latest activity: newest first •	Agents in OBIEE 12c and 11g, known as IBots back in OBIEE 10g, are a component of the BI platform providing few interesting things like pushing data to users instead of requiring users to connect and get it themselves, have alerts based on data to inform users about	12c Time Series Functions: how, why, what, where	
VIEW THE BLOG	Q Issue with 12c and IE11         ArijitC           Q OTBI Data Lineage.         3467422	OBIEE 12c Custom Style using shared folder (analyticsRes) by Gianni Ceresa   27 February 2017   Under the hood	Categories Community Conferences Hacks	





## 450+ Technical Experts Helping Peers Globally





#### bit.ly/OracleACEProgram

Nominate yourself or someone you know: acenomination.oracle.com



Copyright © 2017, Oracle and/or its affiliates. All rights reserved.

**Graph Database: what's that?** 





#### **Property Graph Database - What's that?**





**W**@G\_Ceresa

#### **Property Graph Database - What's that?**





#### **Oracle Property Graph**

#### Oracle PGX (Parallel Graph AnalyticX) architecture





#### **Property Graph Database - Examples: data lineage**



**DATA** *Iysis* 

**W**@G Ceresa

#### **Property Graph Database - Examples: data lineage**





#### **Property Graph Database - Examples: detecting frauds**



**W**@G\_Ceresa

#### **Graph Database - Examples: shortest path**

Examples of graphs and graphs analytics can be seen when traveling from a location A to a location B :

Finding shortest path between 2 nodes of a graph





🕑 @G\_Ceresa

**Back to the topic: Machine Learning** 















### **Machine Learning**

- R and/or Python
- SQL and NoSQL
- Hadoop / Spark
- Parallel Database
- Functional/Object Orientated Programming

٠

٠

٠

٠

۲

٠

Machine Learning Computer **Statistics** Science Unicorn Research Software III. Development Psychology Sociology Subject Matter **Business Understanding** Expert Hacker Mindset **Problem Solver** Strategic

- Statistics
- Algebra
- Probabilities
- Bayesian



## Simplifying ... a lot

## Lot of linear algebra

- Matrices
- Vectors

## Bunch of algorithms

## Two definitions:

- Arthur Samuel (1959). Machine Learning: Field of study that gives computers the ability to learn without being explicitly programmed.
- Tom Mitchell (1998) Well-posed Learning Problem: 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.





Machine Learning can only be as good as what you feed it with!

For example, estimating the price of a house:



How much?





Machine Learning can only be as good as what you feed it with!

For example, estimating the price of a house:



How much?

Estimated price: N/A

Colour

yellow





Machine Learning can only be as good as what you feed it with!

For example, estimating the price of a house:



Colour	Size	Year	Location
yellow	120 m2	1975	Bern

How much?





Machine Learning can only be as good as what you feed it with!

For example, estimating the price of a house:



## Estimated price: N/A

Machine Learning isn't Machine Guessing





## It's all about the features ... and training

For example: estimating the price of a house:

Colour	Size	Year	Location	•••	Public transport	Groceries	Price
green	100 m2	1970	Bern		150m	200m	800K CHF
blue	180 m2	2010	Zürich		500m	1000m	2.4M CHF
red	75 m2	1950	Basel		50m	50m	500K CHF
(100'000 rows of houses with a price)							

Training

Colour	Size	Year	Location	• • •	Public transport	Groceries	
yellow	120 m2	1975	Bern	•••	100m	250m	

## Estimated price: 1'217'542 CHF





## It's all about the features ... and training

Training? Is training really mandatory in Machine Learning?

A common global understanding of ML is that a model needs to be trained on some data before to be applied to some new data and get the extra insight we want the "machine" to find.

This isn't always the case ...



#### Finally back, for real, to the topic: Machine Learning & Graph Databases





#### Machine Learning with Graph databases

ML can be split into 3 main buckets when related to Graphs

- Graphs algorithms generating new insights
- Export features from a graph into a "classical" ML pipeline
- Bring ML into the graph





### **Graphs algorithms generating new insights**

Example: finding similar customers based on what products they bought (clustering)

#### Personalized PageRank:

The Personalized PageRank allows to select a particular vertex or a set of vertices from the given graph in order to give them a greater importance when computing the ranking score, which will have as result a personalized PageRank score and reveal relevant (or similar) vertices to the ones chosen at the beginning.



Personalized PageRank for 'Customer 1' : Customer 1 : 0.3401

Customer 3 : 0.1082 Customer 2 : 0.0920

Customer 1 is more similar to Customer 3 than Customer 2



### **Graphs algorithms generating new insights**

Oracle PGX comes with a "set" of algorithms available out of the box

• 59 (in PGX 19.1.0 counting all the sub-versions)

It's possible to write new algorithm

- Using Green-Marl and "compiling" them in PGX
- Using Java

<sup>3</sup> GX includes a wide selection of optimized graph algorithms that can be invoked through the Analyst. The following table provides an overview of the available algorithms, grouped by category.				
Category	Algorithms			
Classic graph algorithms	Prim's Algorithm			
Community detection	Conductance Minimization (Soman and Narang Algorithm), Infomap, Label Propagation			
Connected components	Strongly Connected Components, Weakly Connected Components (WCC)			
Link predition	WTF (Whom To Follow) Algorithm			
Matrix factorization	Matrix Factorization			
Other	Graph Traversal Algorithms			
Path finding	Bellman-Ford Algorithms, Bidirectional Dijkstra Algorithms, Dijkstra Algorithms, Fattest Path, Hop Distance Algorithms			
Ranking and walking	Closeness Centrality Algorithms, Degree Centrality Algorithms, Eigenvector Centrality, Hyperlink-Induced Topic Search (HITS), PageRank Algorithms, Random Walk with Restart, Stochastic Approach for Link-Structure Analysis (SALSA) Algorithms, Verter Betweenness Centrality Algorithms			
Structure evaluation	Adamic-Adar index, Conductance, Cycle Detection Algorithms, Degree Distribution Algorithms, Eccentricity Algorithms, K-Core, Local Clustering Coefficient (LCC), Modularity, Partition Conductance, Reachability Algorithms, Topological Ordering Algorithms, Triangle Counting			

#### Ruilt-In Algorithms



#### **Export features from a graph**





#### **Export features from a graph**

Graphs can be queried, results exported as "datasets" (rows and columns)

- Run graphs algorithms to generate "features"
  - PageRank
  - Number of in-out edges
  - etc.
- Query the graph to select all the features you need, including the newly generated properties of the algorithms
- Export (as CSV for example) and feed it into your standard ML pipeline

# The power of classical ML increased by extra features representing the graph (some components of the structure around data, enriching the basic features)





### Bring ML into the graph

To bring ML into the graph, the graph needs to be made compatible with ML techniques





## **Bring ML into the graph**

Oracle PGX has, at the moment (PGX 19.3.1), 2 built-in algorithms for ML in the graph:

- DeepWalk (Vertex embeddings)
- Pg2vec (Graph embeddings)









http://www.perozzi.net/publications/14\_kdd\_deepwalk.pdf



## DeepWalk (Random Walks + Word2vec)

From a node:

- Do random walks
- Build "sentences of words" by putting together the visited nodes
- Each node will have many "sentences of words"

Apply Word2vec:

• Transform the sentences into a vector for each node (applying NLP techniques, distance between words, NN etc.)

Result: each node has a vector representation



(It's way more complicated than that, there are many theoretical and practical papers explaining these algorithms)



### Let's imagine...



#### Let's imagine...





#### Let's imagine...





#### Director













#### Interactions' graph





Director



Teacher B







**W**@G\_Ceresa



#### Students

Students



**W**@G\_Ceresa

#### Interactions' graph





💓 @G\_Ceresa









Some "walking" parameters:

- number of walks per node
- walk length

Result?

A list of "sentences"

- size: ( "number of nodes" \* "number of walks per node") x "walk length"
- sentences of words: the nodes can be identified by numbers or words, doesn't matter



#### Some "walking" parameters:

- number of walks per node
- walk length

Result?

## A list of "sentences"

- size: ( "number of nodes" \* "n
- sentences of words: the node

In our example:

- the graph has 13 nodes
- with "number of walks per node" = 80, "walk length" = 40

The result is a list of 1'040 sentences of 40 words each





(in our example) We have a list of 1'040 sentences of 40 words each, how does that become a vector of features for each node?

Using word embedding, a NLP technique

- NLP: Natural Language Processing
- Word embedding: capturing the context of a word in a document, semantic and syntactic similarity, relation with other words etc.

Word2vec model:

- a word embedding technique generating a vector representation of every single unique word
- The dimensionality of the embedding is a parameter to the model

words to dictionary (for a numeric, integer indices, representation of words)
 feed the "magic" of the model (based on a neural network, a "simple" 3 layers one)
 the result is a kind of lookup table that maps from integer indices (which stand for specific words) to dense vectors (their embeddings)



Word2vec looks at which words are nearby to which other, using a neural network

A Word2vec "key" parameter:

• window size: how many words behind and ahead to consider

There are many other parameters which influence the neural network and the model itself, things like batch size, negative sampling, number of epochs, sampling rate etc.

All the settings of the Oracle PGX implementation of DeepWalk: <u>https://docs.oracle.com/cd/E56133\_01/latest/javadocs/oracle/pgx/api/beta/mllib/DeepWalkModelBuilder.html</u>



sentence: "The quick brown fox jumps over the lazy dog" window size: 2

#### Source Text

#### Training Samples

(the, quick) (the, brown)

The quick brown fox jumps over the lazy dog.  $\Longrightarrow$ 

The quick brown fox jumps over the lazy dog.  $\longrightarrow$  (quick for a set of the lazy dog.  $\longrightarrow$ 

The quick brown fox jumps over the lazy dog.  $\Longrightarrow$ 

The quick brown fox jumps over the lazy dog.  $\longrightarrow$ 

(quick, the) (quick, brown) (quick, fox)

(brown, the) (brown, quick) (brown, fox) (brown, jumps)

(fox, quick) (fox, brown) (fox, jumps) (fox, over)



**9** @G\_Ceresa

Word2vec uses the pairs of words obtained by applying the "windowing"

The objective of the model is to maximise the probability of a "context word" to be predicted as context for a "target word"

#### Result:

. . .

**W**@G\_Ceresa

vertex 1 = 
$$[v_1, v_2, v_3, ..., v_n]$$
  
vertex 2 =  $[v_1, v_2, v_3, ..., v_n]$   
vertex 3 =  $[v_1, v_2, v_3, ..., v_n]$   
vertex 4 =  $[v_1, v_2, v_3, ..., v_n]$   
vertex 5 =  $[v_1, v_2, v_3, ..., v_n]$ 

n = layer size (by default 200 for Deepwalk in PGX)



(the details of the word.2vec implementation are of out of the scope of this presentation and would take too long to cover)



**DeepWalk ... for real!** 









Using DeepWalk every node of the graph has a vector representation

- Putting together all the nodes we get a matrix (1 row = 1 node, 1 column = 1 component of the vector)
  - Size of the vectors is a parameter of DeepWalk, by default it is 200
- This can be used as source to apply "classical" machine learning algorithms
- Or directly in PGX by calling the "computeSimilars" method on the model



DeepWalk on the DBpedia graph as an example (with 8'637'721 vertices and 165'049'964 edges):

```
pgx> var similars = model.computeSimilars("Albert_Einstein", 10)
pgx> similars.print()
```

dstVertex	similarity
<pre>  Albert_Einstein   Physics   Werner_Heisenberg   Richard_Feynman   List_of_physicists   Physicist   Max_Planck   Niels_Bohr   Quantum_mechanics   Special_relativity</pre>	1.0000001192092896 0.8664291501045227 0.8625140190124512 0.8496938943862915 0.8415523767471313 0.8384397625923157 0.8370327353477478 0.8340970873832703 0.8331197500228882 0.8280861973762512
+	+



Using DeepWalk every node of the graph has a vector representation

- Putting together all the nodes we get a matrix (1 row = 1 node, 1 column = 1 component of the vector)
  - Size of the vectors is a parameter of DeepWalk, by default it is 200
- This can be used as source to apply "classical" machine learning algorithms
- Or directly in PGX by calling the "computeSimilars" method on the model

### But ...

What about properties? In a Property Graph every node can have a dynamic set of properties. Edges can have labels and properties too.

The "original" definition of DeepWalk doesn't care at all about that.

## DeepWalk focus on the structure rather than the content...



### Paragraph2Vec

Similar to DeepWalk, but on a sub-graph level instead of node

- Each sub-graph will be considered as a paragraph
- Generate random walks on each sub-graphs

Oracle improved Paragraph2Vec with their own "secret sauce" :

- Consider multiple properties (rather than a single one)
- When generated traces, consider edges (instead of vertices) as words
- Attach global properties of each sub-graph (like the size etc.)



#### Conclusion

- PGX is actively developed by Oracle Labs
  - Comes with great tutorials and documentation
- Has ~60 graphs algorithms out of the box
- ML algorithms are being implemented (DeepWalk, Oracle-custom-pg2vec)
  - Beta since PGX 19.1.0 (still beta at the moment)

#### The best results are achieved by a mix of techniques:

Use the graph to generate features and feed the whole result set to a ML Framework (TensorFlow, scikit-learn etc.) At least for now...

