



# Data Challenge

## Team

Gkountoumas Filippas

Legkas Sotiris

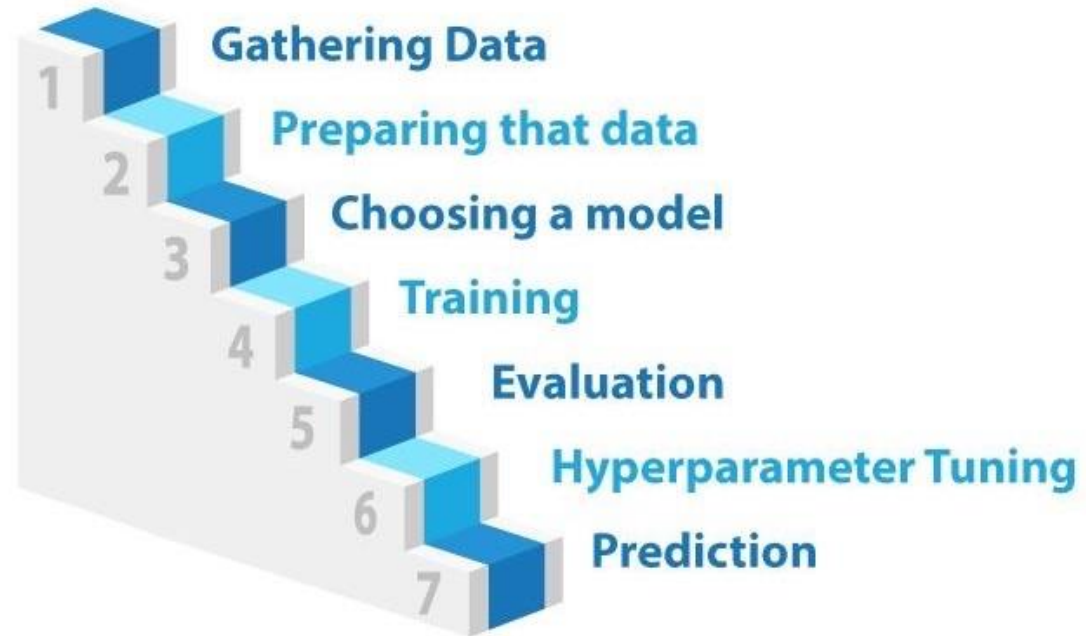
Papadopoulos Nikolaos



## Domain Name Classification Challenge

- Study and Apply ML/ Data mining Techniques
- Assign Domain Names into a Predefined Category
- Textual data & Web Graph
- Goal : Predict the class of unlabeled domains

# Project Pipeline





# Data Exploration



- Unbalanced Data
- Major categories "0"&"3"
- Categories "7","8","9" less than 25 obs



# Data Preprocessing (Data Cleaning)

- **Removal of**
  - HTML tags
  - Punctuation and symbols
  - Numbers
  - Greek stop words
  - Words under 4 letters
  - Multiple spaces
- **Further preprocessing**
  - Lowercasing the text
  - Accent stripping
  - Text length shortening ( 2500 words per text)
  - Stemming ( using Greek Stemmer)

# Data Preprocessing (Data Cleaning)

'http://kollintzas.gr/%CE%94%CE%99%CE%91%CE%93%CE%A9%CE%9D%CE%99%CE%A3%CE%9C%CE%9F%CE%A3%20%CE%95%CE%9A%CE%A0%CE%91%CE%99%CE%94%CE%95%CE%A5%CE%A4%CE%99%CE%9A%CE%A9%CE%9D.htm \*#\* ευρώ\η Καταβολή διδασκτρων Στη η διδασκτική συνάντηση\ηΓια δηλώσεις συμμετοχής στα δωρεάν δοκιμαστικά μαθήματα απαιτείται η δήλωση συμμετοχής τουλάχιστον τρεις ημέρες πριν την ημερομηνία έναρξης\ηΣτους υποψήφιους θα δοθούν εκτός των επιμελημένων δακτυλογραφημένων σημειώσεων που συνοδεύουν την κάθε διδασκτική συνάντηση συνοπτικά ερωτήσεις πολλαπλών επιλογών \ηΑΝΑΚΟΙΝΩΣΕΙΣ ΣΧΕΤΙΚΑ ΜΕ ΤΗΝ ΠΟΡΕΙΑ ΤΟΥ ΠΡΟΣΕΧΟΥΣ ΔΙΑΓΩΝΙΣΜΟΥ \ηΔιαγωνισμός ΑΣΕΠ Εκπαιδευτικών Ν έες Δηλώσεις κ Υπουργού \ηΕφόσον οι πιστωτές της χώρας συναινέσουν στο σχέδιο του υπουργού Παιδείας τότε θα ανοίξει ο δρόμος για τη διενέργεια διαγωνιστού ΑΣΕΠ για την πρόσληψη εκπαιδευτικών στα σχολεία Όπως αποκαλύπτει ο κ Λοβερδός στην Η στόχος είναι ο διαγωνισμός να διενεργηθεί τον χειμώνα και συγκεκριμένα τον Δεκέμβριο του ή τον Ιανουάριο του και να αφορά πάντοτε μεγάλες κατηγορίες εκπαιδευτικών φιλολόγους αθηατικούς φυσικούς δασκάλους και νηπιαγωγούς \ηΣύμφωνα με υπολογισμούς του υπουργείου Παιδείας για την αλήθεια των σχολείων το απαιτούμενο περί τις προσλήψεις τακτικού προσωπικού Σημειώνεται ότι ο τελευταίος διαγωνισμός ο οποίος ΑΣΕΠ είχε διενεργηθεί το και από αυτόν δεν έχουν ακόμη διοριστεί επιτυχόντες εκπαιδευτικοί \ηΟ κ Λοβερδός εκφράζει αισιοδοξία για το αποτέλεσμα των διαπραγματεύσεων που θα έχει με την τρόικα σε τεχνικό επίπεδο βασικό ενός όπως λέει στο γεγονός ότι η κατάσταση στον χώρο της εκπαίδευσης έχει ξεπεράσει τα όρια Το σύστημα κάλυψης των κενών στα σχολεία αναπληρωτές πεθαίνουν Από τη στιγμή στερούνται οι πηγές χρηματοδότησης των αναπληρώσεων έσω του Προγράμματος Δημοσίων Επενδύσεων και του ΕΣΠΑ κι από τη ν άλλη εξαντλούνται τα έσοδα πολιτικής επίσημης είναι ο κ Λοβερδός αποκλείοντας το ενδεχόμενο να προχωρήσει σε νέο κύκλο συγχώνευσης καταργήσεων σχολικών ονδών \ηΑλλάστε όπως αναφέρει ο ίδιος από το σχέδιο Αθηνά για την αναδιάρθρωση των πανεπιστημίων και ΤΕΙ της χώρας δεν εξοικονομήθηκε ούτε ένα ευρώ παρότι η κυβέρνηση στόχευε να κάνει οικονομία της τάξεως των εκατό ευρώ Αυτό αποκάλυψε ο αναπληρωτής υπουργός Οικονομικών Χρήστος Σταϊκούρας σε πρόσφατη σύσκεψη που έγινε στο Μέγαρο Μαξίμου για τα οικονομικά των ΑΕΙ παρουσία του πρωθυπουργού και του αντιπροέδρου της κυβέρνησης \ηΤο αίτημα στην τρόικα\ηΣε προσπάθεια να πείσει τη

Original Text



\*ευρω καταβολή διδασκτρων διδασκτική συνάντηση δηλως συμμετοχή δωρεάν δοκιμαστικά μαθήματα απαιτείται δηλως συμμετοχή τουλάχιστον 3 ημερών πριν την ημερομηνία έναρξης υποψηφίου δοθ εκτός επιμελημένων δακτυλογραφικών σημειώσεων συνοδευόμενες διδασκτική συνάντηση συνοπτικές ερωτήσεις πολλαπλής επιλογής αναφορικά σχετικά μπορεί προσέχ διαγωνισμός διαγωνισμός ασεπ εκπαιδευτικών δηλως υπουργός εφόσον πιστωτές χωρών συναινέσουν σχέδιο υπουργού παιδείας τότε ανοίγει διενεργείται διαγωνισμός ασεπ πρόσληψη εκπαιδευτικών σχολείων όπως αποκαλύπτει λοβερδός στην στοχο είναι διαγωνισμός διενεργηθεί συγκεκριμένα δεκέμβριο ή ιανουάριο αφού πάντοτε μεγάλες κατηγορίες εκπαιδευτικών φιλόλογους αθηατικούς φυσικούς δασκάλους νηπιαγωγών φων υπολογισμούς υπουργείου παιδείας αναφέρει ότι ο τελευταίος διαγωνισμός ο οποίος ασεπ είχε διενεργηθεί το και από αυτόν δεν έχουν διοριστεί επιτυχόντες εκπαιδευτικοί λοβερδός εκφράζει αισιοδοξία για το αποτέλεσμα των διαπραγματεύσεων που θα έχει με την τρόικα σε τεχνικό επίπεδο βασικό ενός όπως λέει στο γεγονός ότι η κατάσταση στον χώρο της εκπαίδευσης έχει ξεπεράσει τα όρια Το σύστημα κάλυψης των κενών στα σχολεία αναπληρωτές πεθαίνουν Από τη στιγμή στερούνται οι πηγές χρηματοδότησης των αναπληρώσεων έσω του Προγράμματος Δημοσίων Επενδύσεων και του ΕΣΠΑ κι από τη ν άλλη εξαντλούνται τα έσοδα πολιτικής επίσημης είναι ο κ Λοβερδός αποκλείοντας το ενδεχόμενο να προχωρήσει σε νέο κύκλο συγχώνευσης καταργήσεων σχολικών ονδών \ηΑλλάστε όπως αναφέρει ο ίδιος από το σχέδιο Αθηνά για την αναδιάρθρωση των πανεπιστημίων και ΤΕΙ της χώρας δεν εξοικονομήθηκε ούτε ένα ευρώ παρότι η κυβέρνηση στόχευε να κάνει οικονομία της τάξεως των εκατό ευρώ Αυτό αποκάλυψε ο αναπληρωτής υπουργός Οικονομικών Χρήστος Σταϊκούρας σε πρόσφατη σύσκεψη που έγινε στο Μέγαρο Μαξίμου για τα οικονομικά των ΑΕΙ παρουσία του πρωθυπουργού και του αντιπροέδρου της κυβέρνησης \ηΤο αίτημα στην τρόικα\ηΣε προσπάθεια να πείσει τη

Processed Text





# Feature Engineering

- **TF-IDF**

Only kept words occurring more than 50 and less than 300 times  
After the data transformation only used the 20.000 most frequent words

- **Greek Word Embeddings**

Used fasttext embeddings from AUEB

- **Dimensionality Reduction**

Used Truncated SVD to reduce feature dimensionality to 128

- **Graph creation**

Created a Graph with all node connections from the given file

- **Node2Vec**

Implemented by combining StellarGraph's random walk generator with the word2vec algorithm from Gensim



# Model Creation

- **GCN**

Approach for semi-supervised learning on graph-structured data.

Spatial Convolution works on a local neighbourhood of nodes and understands

## Model Structure

- 2 Convolution Graph Layers (32 units each, activation : ReLU)
- 2 Dropout Layers ( Dropout Rate 0.5 to avoid overfitting)
- 1 Dense Layer ( 10 units, activation : SoftMax)

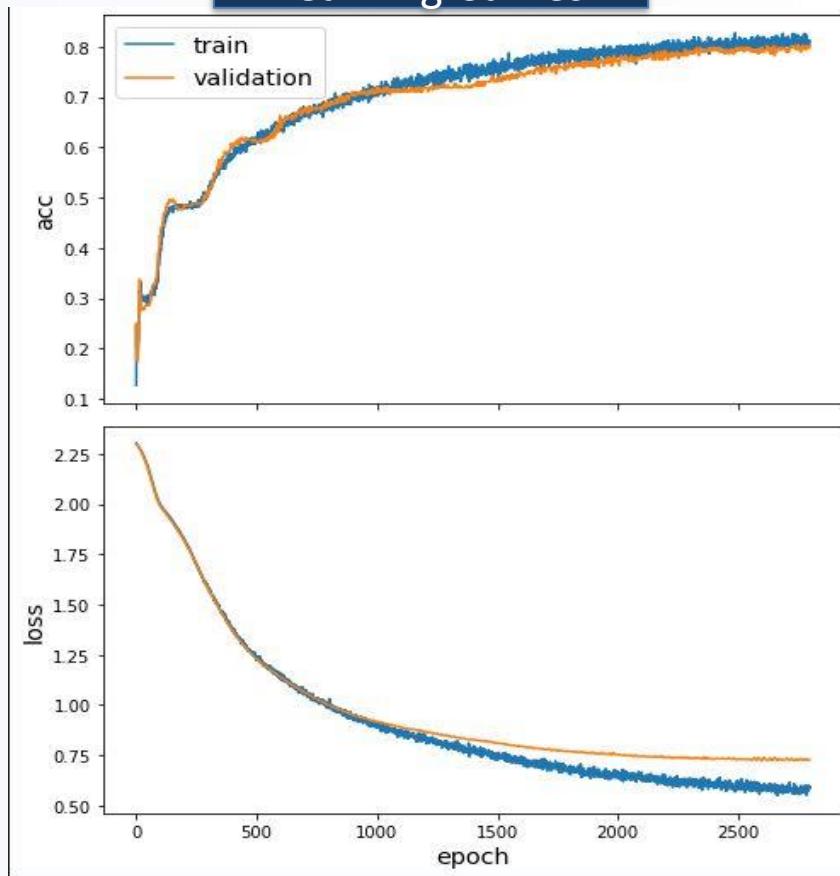
## Model fit

- Loss : Categorical Cross Entropy
- Optimizer : Adam
- Learning Rate :  $10e-3$
- Split : Stratified to maintain class balance



# Model Evaluation

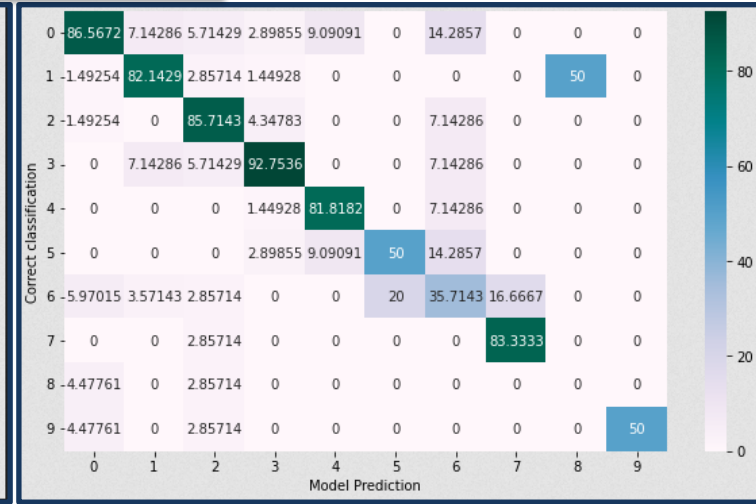
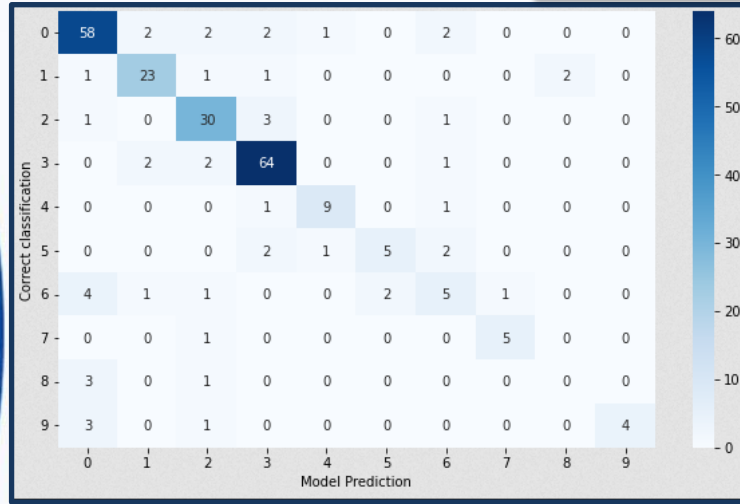
## Learning Curves



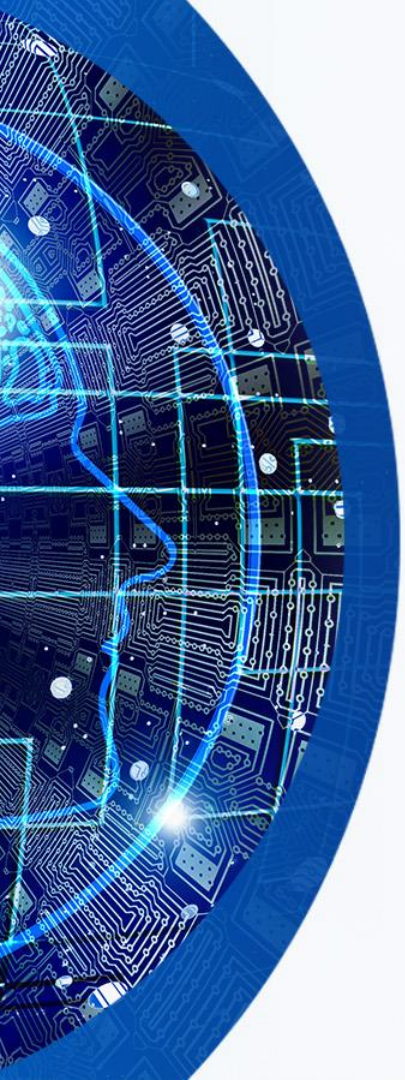
- No overfitting occurred
- Excessive fluctuation might be a sign of unrepresentative dataset

# Error Analysis

## Confusion Matrices



- Domains belonging to the most common categories, get classified with higher accuracy
- Categories with less observation are way more difficult to classify  
Especially categories that even human annotators would have difficulties classifying



Thank you