#### ΕΘΝΙΚΟ ΚΑΙ ΚΑΠΟΔΙΣΤΡΙΑΚΟ ΠΑΝΕΠΙΣΤΗΜΙΟ ΑΘΗΝΩΝ



# Διαχείριση Μεγάλων Δεδομένων

Project 1

Όνομα: Χρήστος - Θρασύβουλος Κούλης

**Αριθμός Μητρώου:** en2190007

**Μεταπτυχιακό:** ΠΜΣ στη Μηχανική Υπολογιστών, Τηλεπικοινωνιών και Δικτύων

(MSc in Computer, Telecommunications and Network Engineering) -Ειδίκευση Τηλεπικοινωνιών και Επεξεργασίας Σήματος

Ημερομηνία Παράδοσης : 12/6/2020

Εγκατέστησα το **Spark** σε λειτουργικό Linux Mint, σύμφωνα με τις οδηγίες. Με *pip install* εγκατέστησα το **pyspark**. Με εντολές από το command window εγκατέστησα το **scala**.

Το πρόγραμμα τρέχει κανονικά αν γράψω *spark-shell* στο command window, χωρίς προβλήματα.

# Ερώτημα 2

Έφτιαξα λογαριασμό στο Kaggle και κατέβασα το  $fake\_job\_postings\ dataset$ . Πρόκειται για ένα .csv αρχείο.

Για αυτό το ερώτημα ξεκίνησα ένα spark shell, φόρτωσα το fake\_job\_postings.csv σε μία απλή μεταβλητή και με κάποιες απλές εντολές άντλησα πληροφορίες από αυτό.

To dataset έγινε load με χρήση της εντολής:

val df =

spark.read.format("csv").option("header", "true").load("/home/akis/Desktop/Big Data Projects/Project 1/fake job postings.csv")

όπου μέσα στο load είναι το directory που είχα αποθηκεύσει το αρχείο .csv.

**a.)** Για να βρω πόσες γραμμές έχει το αρχείο χρησιμοποίησα την εντολή:

df.count()

Που απλά μετράει πόσες γραμμές έχει το dataset. Βρήκα ότι το dataset έχει **17880** γραμμές.

**b.)** Η στήλη που καθορίζει αν μία δουλειά είναι αληθινή ή ψεύτικη είναι η ['fraudulent']. Η στήλη αυτή έπαιρνε μόνο δύο τιμές, 0 και 1. Το 0 αντιστοιχούσε σε αληθινές δουλειές και το 1 σε ψεύτικες. Για να βρω πόσα job postings ήταν fake χρησιμοποίησα την εντολή:

```
val fake = df.filter($"fraudulent" === 1)
```

Που ουσιαστικά κρατάει από το αρχικό dataset μόνο όσες γραμμές ήταν fake job postings.

Βρήκα ότι οι ψεύτικες δουλειές είναι 886 γραμμές (άρα **886 fake job postings**).

**c.)** Ομοίως με πριν, για να βρω πόσα job postings ήταν πραγματικά χρησιμοποίησα την εντολή:

```
val fake = df.filter($"fraudulent" === 0)
```

Που ουσιαστικά κρατάει από το αρχικό dataset μόνο όσες γραμμές ήταν real job postings.

Βρήκα ότι οι αληθινές δουλειές είναι 16080 γραμμές (άρα **16080 real job postings**).

**d.)** Για να βρω τις κορυφαίες απαιτήσεις ασχολήθηκα με την στήλη ['required\_education']. Για να βρω τις 10 κορυφαίες έγραψα αυτήν την εντολή:

```
val fakereq =
fake.groupBy("required education").count().orderBy(desc("count")).limit(10).show()
```

Αρχικά ομαδοποιούμε μαζί όλες τις γραμμές με παρόμοιες τιμές (groupBy()), και ύστερα μετράμε (count()) πόσες φορές εμφανίζεται η κάθε μία. Με orderBy(desc("count")) τα ταξινομούμε με φθίνουσα σειρά, ενώ με limit(10) κρατάμε μόνο τις 10 κορυφαίες. Με το show() τα εμφανίζουμε στην οθόνη. Με το όρισμα show(false) εμφανίζουμε ολόκληρο το όνομα τους. Παρακάτω είναι τα αποτελέσματα:

++				
l — — — — — — — — — — — — — — — — — — —				
required_education count				
null  438				
High School or eq  164				
Bachelor's Degree  96				
Unspecified  60				
Master's Degree  30				
Some High School  20				
Certification  18				
We provide:* Comp  5				
Associate Degree  5				
GoldLeaf provides  3				
BENEFITSGoldLeaf  3				
Some College Cour  3				
Professional  3				
and/or Swedish" 2				
We provide:* Comp  2				
tackling operati 2				
We provide:* Our  1				
BBA Accounting; C  1				
and having the a  1				
Passion for digit  1				
rassion for digit				
only showing top 20 rows				
only showing top 20 rows				

Παρατηρούμε ότι στα fake job postings πολλές γραμμές δεν έχουν σωστό data για αυτήν την στήλη, ενώ υπάρχουν πολλές fake και null τιμές, πράγμα που ενισχύει τις υποψίες μας ότι μία δουλειά μπορεί να είναι ψεύτικη. Το όρισμα *limit(10)* δεν έχει χρησιμοποιηθεί για τον παραπάνω πίνακα για να έχουμε πιο ολοκληρωμένη εικόνα.

**e.)** Ομοίως με πριν, ασχολήθηκα με την στήλη ['required\_education']. Για να βρω τις 10 κορυφαίες έγραψα αυτήν την εντολή:

```
val realreq =
real.groupBy("required_education").count().orderBy(desc("count")).limit(10).show(false)
```

Ουσιαστικά κάνουμε ακριβώς το ίδιο με πριν αλλά μόνο στο dataframe με τις αληθινές δουλειές. Τα αποτελέσματα ήταν τα παρακάτω:

```
required education|count|
    null| 7223|
Bachelor's Degree| 4758|
High School or eq...| 1799|
Unspecified| 1237|
      Master's Degree
                            356 İ
     Associate Degree
                            253
Certification
Some College Cour...
                            150 j
                             931
          Professional
                              64
                             47
            Vocational|
             Doctorate|
                              22
Vocational - HS D...
                              9 j
Some High School ...
                               6 j
 Vocational - Degree
                               4
 strong negotiati...
                               3 |
We provide:* Comp...
                               2
Competitive baseA...|
Medical,Dental,Vi...|
                               2
Be part of the Ta...
                               11
 We also party a ...
only showing top 20 rows
```

Παρατηρούμε ότι στα real job postings έχουμε πολύ πιο ξεκάθαρη εικόνα για το τι χρειαζόμαστε, παρόλα αυτά υπάρχουν ακόμα πολλές fake και null τιμές. Το όρισμα limit(10) δεν έχει χρησιμοποιηθεί για τον παραπάνω πίνακα για να έχουμε πιο ολοκληρωμένη εικόνα.

Όλες οι εντολές που χρησιμοποιήθηκαν για αυτό το ερώτημα βρίσκονται στο αρχείο Erotima1\_3.txt. Απλά κάντε copy-paste τα περιεχόμενα του μέσα σε ένα spark-shell για να εμφανιστούν όλα τα αποτελέσματα.

To pyspark τρέχει κανονικά γράφοντας pyspark στο command window.

```
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties

Setting default log level to "WARN".

To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).

Welcome to

/ _ / _ / _ / _ / _ / _ / _ / _ /

_ \ \ / _ / _ / _ / _ / _ / _ / _ /

Using Python version 3.6.9 (default, Nov 7 2019 10:44:02)

SparkSession available as 'spark'.
```

Μέσα στο pyspark shell γράφουμε κώδικα python. Ο κώδικας βρίσκεται στο αρχείο **Erotima1\_4.py**. Ο κώδικας έχει αρκετά σχόλια που επεξηγούν για κάθε ερώτημα αναλυτικά τι έκανα. Παρακάτω παραθέτω τα αποτελέσματα μου:

f.) Για τα fake job postings βρήκα τα ακόλουθα:

$$average = 55904.10$$

 $standard\ deviation = 30735.32$ 

g.) Για τα real job postings βρήκα τα ακόλουθα:

$$median = 40$$

Παρ' όλα αυτά, το νούμερο αυτό δεν είναι σωστό καθώς η στήλη των minimum salary είναι string value οπότε δεν βρίσκεται σωστά ο median. Αφού έκανα την στήλη typecast σε double, βρήκα ότι ο σωστό median είναι ο:

$$median = 40000$$

Για να βρω τα δύο αυτά ερωτήματα χρειάστηκε να χωρίσω την στήλη ['salary\_range'] σε δύο στήλες, ['min\_salary'] και ['max\_salary']. Όλες οι γραμμές με *null* τιμές αγνοήθηκαν (*df.na.drop()*).

h.) Τα outliers τα βρήκα λίγο διαισθητικά. Αφού είδα όλες τις τιμές για τα fake και τα real job postings (εύκολο, καθώς non-null τιμές δεν ήταν παραπάνω από λίγες εκατοντάδες), βρήκα 3 ειδών outliers.

Πρώτα, βρήκα ότι για τα maximum salaries στα fake job postings υπάρχει μία τιμή 180000, που είναι πάνω από 50% της αμέσως μεγαλύτερης τιμής, οπότε λογικά ίσως είναι outlier. Η τιμή αυτή αφαιρέθηκε από το dataframe.

Για τα minimum salaries των real job postings υπάρχουν πιο προφανή outlier. Είδα ότι υπάρχουν αρκετές τιμές σχεδόν μηδενικές, κάτι που δεν βγάζει νόημα, οπότε αφαίρεσα όλα τα job postings με minimum μισθό λιγότερο των 1000 ευρώ. Ακόμα, υπάρχουν 2 τιμές μισθών άνω των 100 εκατομμυρίων που προφανώς δεν βγάζουν νόημα. Οι τιμές λοιπόν με minimum μισθό άνω των 600000 ευρώ αφαιρέθηκαν.

Για το f.), τα νέα αποτελέσματα είναι:

$$average = 54180.55$$

#### $standard\ deviation = 27167.15$

Ενώ για το g.):

$$median = 40000$$

Βλέπουμε ότι δεν άλλαξε, κάτι που είναι λογικό καθώς ο median δεν εξαρτάται από outliers. Ο τρόπος που υπολογίστηκαν τα παραπάνω βρίσκονται στον κώδικα.

i.) Για να βρω τα bigram και τα trigram συγκέντρωσα όλα τα description σε list (και ύστερα string) μεταβλητές και χρησιμοποίησα τις βιβλιοθήκες *re, ngrams* και *collections* της python για να τα απαριθμίσω. Τα αποτελέσματα ήταν τα παρακάτω:

#### Για τις ψεύτικες δουλειές:

Bigram	Occurences	Trigram	Occurences
of , the	462	are , looking , for	184
we , are	412	we , are , looking	145
in , the	397	oil , and , gas	128
looking, for	344	looking , for , a	124
to , the	288	and , gas , industry	98
for , the	285	we , are , seeking	91
to , work	276	be , responsible ,	85
		for	
in , a	237	be , able , to	74
is, a	235	as , well , as	70
are , looking	218	we , are , a	64

# Για τις αληθινές δουλειές:

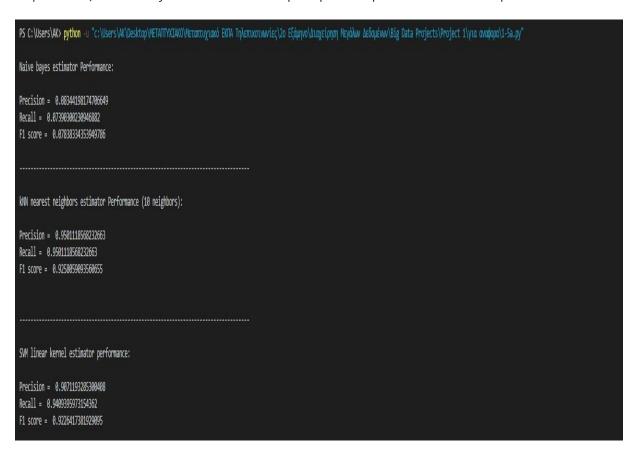
Bigram	Occurences	Trigram	Occurences
of , the	9722	you , will , be	2923
in , the	8611	are , looking , for	2920
will , be	8352	we , are , looking	2678
you , will	7724	looking , for , a	2576
we, are	7140	as , well , as	2160
looking , for	6404	be , responsible ,	1636
		for	
in , a	5184	to , join , our	1584
is, a	5010	be , able , to	1530
to , the	4811	will , be ,	1376
		responsible	
for , the	4794	this , is , a	1365

Ο κώδικας βρίσκεται στο αρχείο Erotima1\_4.py. Πρόκειται για ένα απλό python αρχείο.

Για το ερώτημα αυτό χρησιμοποιήθηκε η python, όχι όμως σε περιβάλλον pyspark.

j.) Αρχικά θα χρησιμοποιήσω μόνο την στήλη ['telecommuting'] ως feature σε έναν naive bayes classifier. Τα αποτελέσματα εδώ όπως θα δούμε και στην παρακάτω εικόνα δεν είναι καθόλου ικανοποιητικά, με precision, recall και  $f1\_score$  κάτω του 10%. Αντίθετα, όταν ενσωματώσουμε στα feature list και τις στήλες ['has\_company\_logo'] και ['has\_questions'], το ποσοστό επιτυχίας αυξάνεται ραγδαία (πάνω του 90%), όταν χρησιμοποίησουμε k-Nearest Neighbors classifier με N=10 γείτονες ή SVM με sigmoid kernel. Για το SVM το precision είναι μικρότερο, καθώς όπως είναι γνωστό τα sigmoid kernel δεν βγάζουν τόσο καλά αποτελέσματα.

Τα precision, recall και f1-score που υπολόγισα για τα 3 μοντέλα είναι τα παρακάτω:



Ο κώδικας βρίσκεται στο αρχείο 1-5a.py Πρόκειται για ένα απλό python αρχείο.