

# Τεχνητή Νοημοσύνη

## 2η Εργασία

Ομάδα:

Κουτσομπίνας Γιώργος 3150251

Βελαώρας Απόστολος 3180249

Καλδής Αργύριος 3160045

### Λογιστική παλινδρόμηση

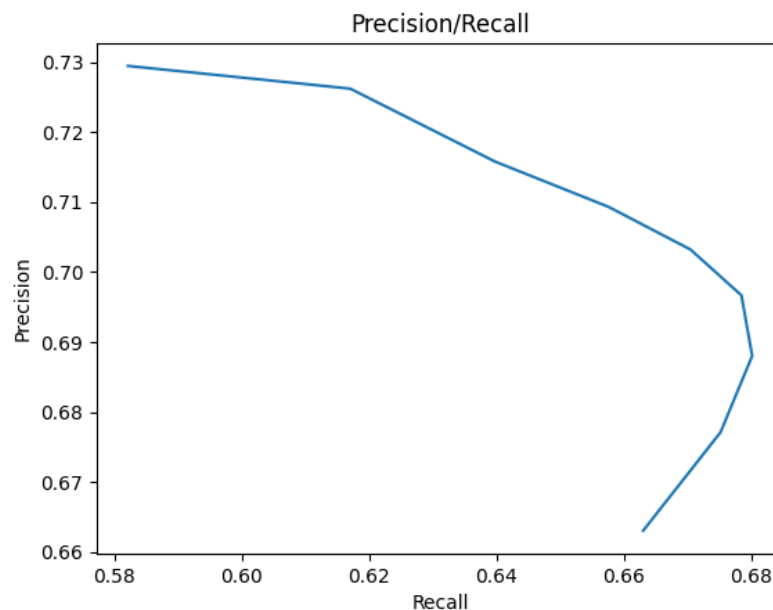
- data.py: χρησιμοποιείται για την δημιουργία λεξικού και ανάγνωση των .txt αρχείων. Για κάθε λέξη που διαβάζεται, επιστρέφει ένα διάνυσμα με τιμές 0 ή 1 ανάλογα με το αν υπάρχει η λέξη στο λεξικό. Στο τέλος κάθε διανύσματος μπαίνει η τιμή 0 αν το σχόλιο ανήκει στην αρνητική κατηγορία, ή η τιμή 1 αν το σχόλιο ανήκει στην θετική κατηγορία
- logistic\_regression.py: περιλαμβάνει τις μεθόδους για την υλοποίηση της λογιστικής παλινδρόμησης:
  - logistic\_function(x, w): σιγμοειδής συνάρτηση, επιστρέφει έναν αριθμό από 0 έως 1.
  - train(train\_total): Διαβάζει τα σχόλια από τον φάκελο train, τα χωρίζει σε train και dev και εκτελεί την στοχαστική ανάβαση κλίσης.
  - test(): Διαβάζει τα σχόλια από τον φάκελο test και δοκιμάζει τα βάρη σε αυτά μέσω της σιγμοειδούς συνάρτησης για να τα αξιολογήσει.
  - test\_external(path): Αξιολογεί ένα μεμονωμένο σχόλιο.
  - cost(w, l\_val): Επιστρέφει το κανονικοποιημένο συνολικό κόστος.
  - single\_cost(x, l\_val): Επιστρέφει το κανονικοποιημένο κόστος ενός μεμονωμένου σχολίου.

- `evaluate(pos)`: Επιστρέφει την κατηγορία στην οποία ανήκει ένα σχόλιο train/test.
  - `evaluate_dev(pos)`: Επιστρέφει την κατηγορία στην οποία ανήκει ένα σχόλιο dev.
  - `select_model()`: Διαλέγει τα βάρη με το μικρότερο error rate στα dev δεδομένα, σε περίπτωση που δοθούν πάνω από μια λ τιμές.
  - `converges(old, new)`: Ελέγχει αν συγκλίνει το κόστος.
  - `get_stats()`: Επιστρέφει διαγράμματα με καμπύλες μάθησης.
- Η επιλογή του όρου κανονικοποίησης  $\lambda = 4$  έγινε μετά από παρακολούθηση των καμπυλών μάθησης και από πολλαπλές δοκιμές σε ποσοστά των δεδομένων εκπαίδευσης και
  - Η επιλογή του learning rate  $\alpha = 0.01$  στην στοχαστική κατάβαση κλίσης έγινε μετά από παρακολούθηση των τιμών του συνολικού κόστους σε κάθε ενημέρωση βαρών.

Καμπύλη μάθησης ορθότητας:



Καμπύλη μάθησης ακρίβειας/ανάκλησης:



### Αφελής ταξινομητής Bayes

- **NayveBayes.py**: Περιέχει την υλοποίηση του αλγορίθμου Αφελής ταξινομητής Bayes.
  - **#TRAIN**
  - **countWords()** : Μετράει τις λέξεις που περιέχονται στο dictionary.txt
  - **countPosNeg()**: Επιστρέφει τον αριθμό το θετικών και αρνητικών reviews στη βάση (train)
  - **calcPossibility()**: υπολογίζει την πιθανότητα  $P(C=1)$  κ  $P(C=0)$  (Οπού 1 θετικό , 0 αρνητικό review)
  - **countWordsPerVector()**: Επιστρέφει το άθροισμα των λέξεων του dictionary για κάθε μια από τις δυο κατηγορίες (Θετικές αρνητικές κριτικές)
  - **calcPossibilityOfWords()**: Επιστρέφει τις πιθανότητες κάθε λέξεις για κάθε μια απ τις κατηγορίες ( $P(X|C=1)$  ,  $P(X|C=0)$  οπού X η κάθε λέξη)
  - **summer()**: Επιστρέφει το άθροισμα όλων των λέξεων σε κάθε κατηγορία.

## TEST

- **getFilesVector():** Επιστρέφει το vector για ένα δοσμένο review αλλά σε διάνυσμα που μετράει το πλήθος των λέξεων (όχι μόνο τιμές 0 και 1)
- **getFilesVector2():** Επιστρέφει το vector για ένα δοσμένο review σε διάνυσμα 0 κ 1 (δεν χρησιμοποιείτε , ελάχιστα καλύτερα αποτελέσματα η προηγούμενη )
- **finalCalc():** Υπολογίζει την τελική πιθανότητα  $(P(X_1|C)*P(X_2|C).....P(X_i|C)*P(C))$  για  $I = \text{πλήθος λέξεων}$  και  $C = 0 \text{ ή } 1$   
Με βάση το vector του review υπολογίζει τη πιθανότητα και ανάλογα πόσες φορές έχει εμφανιστεί ή δεδομένη λέξη υπολογίζει την δύναμη της με την Power() και επιστρέφει 0 αν το review είναι αρνητικό ή 1 αν είναι θετικό.
- **Power(Element,N):** Επιστρέφει το Element στη δύναμη N
- **TestResult():** Επιστρέφει ένα άθροισμα με βάση ποσά από τα review ήταν accurate.
- **Accuracy():** Επιστρέφει το ποσοστό ακριβείας των τεστ.

Καμπύλη μάθησης ορθότητας:

