

Code how to use document loaders

- load environment variables, like Open AI API key
- Για PDF's θα πρέπει να ανεβάσουμε τον PDF loader
- Να φορτώσουμε τα data χαλώντας την μέθοδο φόρτωσης

We load a bunch of
PDF's in the folder in
workspace what the
ween by that. 2.21win
video 1

Document splitting

Split documents in to smaller chunks.

Δεν είναι εύκολη διαδικασία, είναι κάποιες λεπτομέρειες που έχουν μεγάλη επιπτώση.

Το splitting γίνεται μετά την φόρτωση των data, στο document format και ΠΡΙΝ πάει στο vector store.

! Το πως θα σπάσεις τα κείμενα σε chunks είναι ιδιαίτερα σημαντικό. Πολλές φορές αν δεν μπορέσεις να κάνεις καλό χωρισμό στα chunks, μπορεί να καταλήξεις με ένα μέρος μιας πρότασης σε ένα chunk και με το άλλο μέρος της πρότασης σε ένα άλλο chunk. Και μετά όταν προσπαθήσεις να απαντήσεις σε ερωτήσεις το chunk δεν θα περιέχει την σωστή πληροφορία και δεν θα μπορεί να απαντήσει την ερώτηση σωστά.

Example Splitter

Μια σωστή μέθοδος για να σπας τα chunks είναι να είναι ενός συχνευρισμένου αριθμού, όπου ένα κομμάτι του τέλους του ενός chunk να είναι η αρχή του άλλου chunk (overlap)

Αυτό βοηθάει στο να δημιουργηθεί μια έννοια συνέχειας καθώς υπάρχει το ίδιο κομμάτι περιεχομένου στο τέλος ενός κομματιού και στην αρχή του άλλου.

Στο langChain υπάρχουν πολλοί τύποι splitters που ο κάθε ένας έχει διαφορετική λειτουργικότητα ή χαρακτηριστικά. Επίσης μπορεί να υπάρχει ποικιλία στον τρόπο που μετρούν το μήκος των κομματιών.

Επίσης υπάρχουν κάποια που χρησιμοποιούν μικρότερα μοντέλα για να προσδιορίσουν πότε μπορεί να τελειώσει μια πρόταση.

Ακόμα ένας άλλος παράγοντας στο πως να χωρίσεις τα

chunks είναι τα metadata. Η διατήρηση των ίδιων metadata σε όλα τα chunks, Αλλά και η προσθήκη νέων δομημένων metadata όταν είναι σχετικό είναι κάτι που υιοθετεί κάποιοι document splitting.

Η διαίρεση των κομματιών μπορεί συχνά να είναι ειδική ανάλογα με το τύπο εγγράφων με τον οποίο εργαζόμαστε. Τέλος έχουμε splitters που διαχωρίζουν με διαφορετικούς separators για διαφορετικές γλώσσες προγραμματισμού όπως Python, Ruby, C και αυτοί λαμβάνουν υπόψη για την κάθε διαφορετική γλώσσα, τους διαφορετικούς separators που έχουν.

Metadata

Είναι σαν μια επιμεθορία ή μια κατηγορία που περιγράφει & παρέχει πληροφορίες για το εύρος δεδομένων.

Σε ένα κείμενο που έχει χωριστεί σε chunk τα metadata αυτά λειτουργούν ως κοινά χαρακτηριστικά που διατηρούν συνέπεια σε όλα τα κομμάτια.

Μπορούμε να συνδυάσουμε
τεχνικές splitter;