

ΙΟΝΙΟ ΠΑΝΕΠΙΣΤΗΜΙΟ

ΤΜΗΜΑ ΠΛΗΡΟΦΟΡΙΚΗΣ



IONIAN
UNIVERSITY

Ανάλυση Συναισθημάτων comments του Stack Overflow

χρησιμοποιώντας το μοντέλο

SamLowe/roberta-base-go_emotions

Στέργιος Μουτζίκος

ΑΜ: inf2021149

Μάθημα: Σημασιολογικός και Κοινωνικός Ιστός

Διδάσκοντες: κ. Κορφιάτης, κα. Αίγλη

Ιούνιος 2025

Περιεχόμενα

1. Περίληψη.....	3
2. Εισαγωγή.....	3-7
2.1 Εισαγωγή στον Ερευνητικό Χώρο της Εργασίας.....	4
2.2 Βασικές Προσεγγίσεις Επίλυσης.....	4-5
2.3 Συνεισφορά της Εργασίας.....	6-7
3. Ερευνητικός χώρος (με βιβλιογραφική επισκόπηση).....	7-10
3.1 Βιβλιογραφική Επισκόπηση.....	7-10
3.1.1 Αρχικές Προσεγγίσεις.....	7-8
3.1.2 Προσεγγίσεις Μηχανικής Μάθησης.....	8
3.1.3 Deep learning και Transformers.....	8-9
3.1.4 Ανάλυση Συναισθημάτων στον Κοινωνικό Ιστό.....	9
3.1.5 Συσχέτιση Συναισθημάτων με Χαρακτηριστικά Χρήστη:.....	9-10
4. Το Μοντέλο SamLowe/roberta-base-go_emotions.....	10-11
5. Μεθοδολογική Διαδικασία.....	11-15
5.1 Περιγραφή Δεδομένων.....	11-12
5.2 Περιγραφή Προεπεξεργασίας.....	12-14
5.2.1 Καθαρισμός Κειμένου (Text Cleaning).....	12-13
5.2.2 Προετοιμασία Δεδομένων για Ανάλυση NLP.....	13-14
5.2.3 Κατηγοριοποίηση Αριθμητικών Δεδομένων (Binning).....	14
5.3 Πειράματα Μάθησης.....	14-15
5.3.1 Ανάλυση Συναισθημάτων με GoEmotions Model.....	14-15
5.3.2 Εξαγωγή και Δομή Δεδομένων Συναισθημάτων.....	15
5.3.3 Συγχώνευση Δεδομένων για Ανάλυση.....	15
6. Ανάλυση και Αξιολόγηση.....	16-25
6.1 Εξερεύνηση και Χαρακτηριστικά Δεδομένων.....	16-25
6.1.1 Επισκόπηση Συνόλων Δεδομένων.....	16-17
6.1.2 Χαρακτηριστικά Κειμένου των Σχολίων.....	17-19
6.1.3 Γενική Κατανομή Συναισθημάτων.....	19-20
6.1.4 Σχέση Συναισθημάτων με τη Δημοτικότητα του Περιεχομένου.....	20-22
6.1.5 Γεωγραφικές Διαφοροποιήσεις ανά συναίσθημα.....	23-24
6.1.6 Συναισθήματα και Αντίδραση της Κοινότητας.....	24-25

7. Συμπεράσματα.....	25
8. Προτάσεις για Μελλοντικές Βελτιώσεις.....	26
8.1 Εμπλουτισμός Δεδομένων.....	26
8.2 Ποιοτική Ανάλυση και Επαλήθευση.....	26
8.3 Προηγμένες Τεχνικές Ανάλυσης.....	26
9. Βιβλιογραφία.....	27-28
10. Code-appendix (with Jupyter notebook Python).....	29-43

.1 Περίληψη

Ο στόχος της παρούσας μελέτης είναι η ανάλυση των συναισθημάτων των χρηστών σε δεδομένα Κοινωνικού και Σημασιολογικού Ιστού με βάση το μοντέλο Transformer. Χρησιμοποιήθηκε το προ-εκπαιδευμένο μοντέλο **SamLowe/roberta-base-go_emotions** για να αναγνωριστούν 27 συναισθήματα σε σχόλια χρηστών. Στη μεθοδολογία περιλαμβάνεται μια αρχική φάση "εξερεύνησης", ακολουθούμενη από έναν εκτεταμένο προεπεξεργαστικό έλεγχο (αφαίρεση URLs/emojis, λημματοποίηση, ανάλυση συχνότητας λέξεων και bigrams) και την κατηγοριοποίηση metadata (φήμη χρήστη, αξιολογήσεις σχολίων κτλ).

Η ανάλυση έδειξε τη κατανομή των συναισθημάτων, συσχετίζοντάς τα με τη δημοτικότητα του περιεχομένου και με τα γεωγραφικά χαρακτηριστικά των χρηστών. Επιπρόσθετα, δίνεται έμφαση στα σχόλια με τη μεγαλύτερη συναισθηματική ακρίβεια και στη σχέση των συναισθημάτων με τα upvotes/downvotes. Τα ευρήματα αποτελούν πολύτιμη πηγή ενδείξεων για την πολυπλοκότητα των συναισθημάτων στο διαδίκτυο. Προτάσεις για μελλοντικές βελτιώσεις περιλαμβάνουν την προσθήκη δεδομένων από διαφορετικές χρονικές και περιβαλλοντικές διαστάσεις, καθώς και τη χρήση πιο εξειδικευμένων τεχνικών για την ανίχνευση περισσότερων συναισθημάτων.

.2 Εισαγωγή

Στη σύγχρονη ψηφιακή εποχή, όπου ο όγκος των παραγόμενων δεδομένων αυξάνεται εκθετικά, η αποκωδικοποίηση των ανθρώπινων συναισθημάτων και της ανθρώπινης έκφρασης έχει κεντρική σημασία. Ο Κοινωνικός και Σημασιολογικός Ιστός (Social and Semantic Web) διαθέτει τεράστιο όγκο πληροφοριών σχετικά με αυτά, με τα άρθρα στα social media, τις αναρτήσεις σε forum και blogs και τις κριτικές προϊόντων να αποτελούν απλώς ένα μέρος αυτών των πληροφοριών. Η Ανάλυση Συναισθημάτων (Sentiment Analysis), επίσης γνωστή ως Opinion Mining, αποτελεί έναν γρήγορα αναπτυσσόμενο κλάδο της Επεξεργασίας Φυσικής Γλώσσας (Natural Language Processing - NLP) και της Τεχνητής Νοημοσύνης (Artificial Intelligence - AI) και έχει ως στόχο την αυτόματη ανίχνευση, ανάκτηση και κατηγοριοποίηση των συναισθημάτων που περιέχονται σε διάφορα κείμενα, όχι μόνο του Κοινωνικού Ιστού, αλλά και της Σημασιολογικής Ανάλυσης και γενικότερα του Διαδικτύου.

.2.1 Εισαγωγή στον Ερευνητικό Χώρο της Εργασίας

Η έρευνα που εξετάζει αυτή η εργασία ασχολείται με την ανάλυση πολύπλευρων συναισθημάτων στο Κοινωνικό Διαδίκτυο και τον Σημασιολογικό Ιστό. Συγκεκριμένα, η εργασία ασχολείται με τον τρόπο έκφρασης των συναισθημάτων σε σχόλια χρηστών και τον τρόπο με τον οποίο αυτά συνδέονται με πολλαπλά δεδομένα, όπως η δημοτικότητα μιας ανάρτησης, η φήμη του χρήστη που κάνει το σχόλιο, η γεωγραφική τοποθεσία και η απήχηση του σχολίου στην κοινότητα (π.χ. μέσω αναφορών, likes ή dislikes).

Η πολυπλοκότητα της ανθρώπινης γλώσσας, η ύπαρξη ειρωνείας, σαρκασμού, αλλά και η συνεχής εξέλιξη της ορολογίας και των εκφράσεων κάνουν την ανάλυση των συναισθημάτων δύσκολη. Η προσέγγιση της έρευνας στον Κοινωνικό Διαδίκτυο και τον Σημασιολογικό Ιστό είναι πιο λεπτομερής. Με τη χρήση ενός μοντέλου που μπορεί να αναγνωρίζει μια εκτεταμένη γκάμα συναισθημάτων (27 διαφορετικές κατηγορίες συναισθημάτων), η έρευνα καταφέρνει να κατανοήσει βαθύτερα τα συναισθήματα που βιώνουν οι χρήστες, να ξεπεράσει τους περιορισμούς των απλών προσεγγίσεων και να παρέχει πλουσιότερες πληροφορίες για την ανθρώπινη επικοινωνία στο διαδίκτυο. Η έρευνα αυτή δεν αφορά μόνο τα συναισθήματα, αλλά φθάνει μέχρι την ανάλυση των μοτίβων και των συνδέσεων που προκύπτουν ανάμεσα στους χρήστες και το περιεχόμενο.

.2.2 Βασικές Προσεγγίσεις Επίλυσης

Η ανάλυση συναισθημάτων έχει αναπτυχθεί σημαντικά μέσα στον χρόνο χρησιμοποιώντας διάφορες τεχνικές και μεθοδολογίες. Οι βασικές προσεγγίσεις διαχωρίζονται ως εξής:

1. **Προεπεξεργασία του Κειμένου (Text Preprocessing):** Πριν από την ανάλυση οποιουδήποτε είδους, τα un-edited κειμενικά δεδομένα απαιτούν κάποιες επεξεργασίες. Περιλαμβάνουν τα εξής:
 - **Μετατροπή σε πεζά γράμματα (Lowercasing):** Όλες οι λέξεις μετατρέπονται σε πεζά για να αντιμετωπιστούν ως ίδιες (π.χ., "Θυμός" και "θυμός").
 - **Αφαίρεση URLs και Emojis:** Οι ηλεκτρονικοί σύνδεσμοι και τα emojis, αν και μπορούν να μεταδώσουν συναισθήματα, συχνά προσθέτουν "θόρυβο" στην ανάλυση κειμένου ειδικά όταν χρησιμοποιούνται μοντέλα που βασίζονται στις λέξεις. Για την εστίαση στο λεκτικό περιεχόμενο, αυτά αφαιρούνται ή καθαρίζονται.

- **Αφαίρεση Stop Words:** Κοινές λέξεις που δεν προσθέτουν ουσιαστικό νόημα (π.χ., "ο", "η", "και", "είναι") αφαιρούνται για να μειωθεί ο όγκος των δεδομένων και να βελτιωθεί η εστίαση στις πιο σημαντικές λέξεις .

2. **Προσεγγίσεις Βασισμένες σε Μηχανική Μάθηση (Machine Learning-based Approaches):** Αυτές οι προσεγγίσεις χρησιμοποιούν μοντέλα μηχανικής μάθησης (π.χ. Support Vector Machines, Naive Bayes, Logistic Regression), τα οποία έχουν εκπαιδευτεί σε μεγάλα σύνολα δεδομένων με επισημασμένα τα συναισθήματα. Έτσι, το κείμενο μετατρέπεται σε αριθμητικά χαρακτηριστικά (π.χ. TF-IDF, Word Embeddings), τα οποία λαμβάνουν εκπαίδευση μέσω του μοντέλου. Η επίδοσή τους εξαρτάται κυρίως από τον τρόπο που διαμορφώνεται και το μέγεθος του όγκου της εκπαίδευσης. Η βελτίωση σε αυτόν τον τομέα οδήγησε στη δημιουργία βαθύων νευρωνικών δικτύων.
3. **Προσεγγίσεις Βασισμένες σε Βαθιά Μάθηση (Deep Learning-based Approaches):** Η εμφάνιση των σύγχρονων βαθιών νευρωνικών δικτύων, όπως τα Επαναλαμβανόμενα Νευρωνικά Δίκτυα (Recurrent Neural Networks - RNNs), τα Συνελικτικά Νευρωνικά Δίκτυα (Convolutional Neural Networks - CNNs) και κυρίως τα Transformers, ανοίγουν νέους ορίζοντες στην ανάλυση των συναισθημάτων. Αυτά τα μοντέλα μπορούν να δημιουργήσουν πολύπλοκες αναπαραστάσεις του κειμένου και να ανιχνεύουν τις μακροπρόθεσμες δημιουργούμενες, έτσι επιτυγχάνουν πολύ καλύτερες αποτελεσματικότητες. Τα μοντέλα Transformers, όπως το BERT και το RoBERTa, έχουν αποδειχθεί ιδιαίτερα αποτελεσματικά, διότι μπορούν να αντιληφθούν τον γενικό πλαίσιο των λέξεων σε μια πρόταση και να δημιουργήσουν υψηλής ποιότητας αναπαραστάσεις κειμένου (embeddings).

Η παρούσα εργασία εμπνέεται από μια μέθοδο βαθύς μάθησης, που πρόκειται για ένα μοντέλο *Transformer* (*SamLowe/roberta-base-go emotions*) που είναι προεκπαιδευμένο στην αναγνώριση πολλαπλών συναισθημάτων. Η συγκεκριμένη επιλογή δίνει τη δυνατότητα για την καλύτερη αντιμετώπιση των δυσκολιών που έχουν αναφερθεί, καθώς παρέχει μια λεπτομερή και ακριβή ανάλυση. Ο χρήστης μέσω ενός *pipeline* από τη βιβλιοθήκη *Hugging Face* καθιστά εύκολη την εκτέλεση αυτών των σύνθετων μοντέλων, δίνοντας την ευκαιρία να επικεντρωθεί στην ανάλυση των αποτελεσμάτων και όχι στην κατασκευή του μοντέλου.

.2.3 Συνεισφορά της Εργασίας

Η παρούσα εργασία συμβάλλει με λεπτομερειακό τρόπο στον τομέα της ανάλυσης συναισθημάτων, αναλυτικότερα:

1. **Πλήρης και Λεπτομερής Ανάλυση Συναισθημάτων:** Η εργασία δεν περιορίζεται σε λίγες βασικές κατηγορίες, αλλά εφαρμόζει ένα προηγμένο μοντέλο για την αναγνώριση 27 διακριτών συναισθημάτων. Με αυτό το μοντέλο είναι εφικτό να γίνει πιο πλούσια και λεπτομερής κατανόηση των συναισθηματικών αποχρώσεων στα σχόλια των χρηστών που είναι το βασικό συστατικό για να βγουν ουσιαστικά συμπεράσματα.
2. **Συσχέτιση Συναισθημάτων με metadata:** Η εργασία δεν αναζητά μόνο τα συναισθήματα αλλά και συνδιάζει τα αποτελέσματα της ανάλυσης των συναισθημάτων με τα διάφορα μεταδεδομένα (όπως φήμη χρήστη, βαθμολογία ανάρτησης, βαθμολογία σχολίου, γεωγραφική τοποθεσία, upvotes/downvotes). Αυτή η συγχώνευση των δεδομένων επιτρέπει τη διερεύνηση σύνθετων σχέσεων και μοτίβων, παρέχοντας απαντήσεις σε ερωτήματα όπως:
 - Πώς επηρεάζει η φήμη ενός χρήστη τα συναισθήματα που εκφράζει ή λαμβάνει;
 - Υπάρχουν συγκεκριμένα συναισθήματα που σχετίζονται με υψηλότερες ή χαμηλότερες βαθμολογίες αναρτήσεων/σχολίων;
 - Διαφέρουν τα εκφραζόμενα συναισθήματα ανάλογα με την γεωγραφική τοποθεσία των χρηστών;
 - Πώς συσχετίζονται τα upvotes και downvotes με τα συναισθήματα που εκφράζονται στα σχόλια;
 -
3. **Οπτικοποίηση και Ερμηνεία των Αποτελεσμάτων:** Αναλυτικά γραφήματα (όπως ιστογράμματα, ραβδόγραμματα, heatmaps, box plots) αποτελούν σημαντικό τμήμα του έργου καθώς με τη βοήθειά τους γίνεται ευκολότερη η κατανόηση των πολύπλοκων δεδομένων και παρουσιάζονται έτσι τα ευρήματα με έναν πιο κατανοητό τρόπο τόσο στους ειδικούς όσο και στους μη ειδικούς.
4. **Χρήση Προηγμένων Εργαλείων NLP:** Η αξιοποίηση σύγχρονων εργαλείων NLP και βιβλιοθηκών (όπως οι transformers και το NLTK) και οι τεχνικές ληματοποίησης, εξαγωγής bigrams και το TF-IDF μοντέλο αποδεικνύουν ότι είμαστε σε θέση να εφαρμόσουμε προηγμένες μεθόδους στην επίλυση πραγματικών προβλημάτων.

5. Εφαρμογές και Επιπτώσεις: Τα ευρήματα αυτής της εργασίας μπορούν να έχουν σημαντικές πρακτικές εφαρμογές σε διάφορους τομείς, όπως:

- Υπηρεσίες Πελατών: Κατανοώντας τα συναισθήματα των πελατών για τα προϊόντα ή τις υπηρεσίες που λαμβάνουν.
- Μάρκετινγκ: Δίνοντας πληροφορίες σχετικά με την αντίληψη του κοινού για μια εταιρεία ή ένα προϊόν.
- Κοινωνιολογία: Αναλύοντας κοινωνικές τάσεις και αντιδράσεις σε θέματα.

.3 Ερευνητικός χώρος (με βιβλιογραφική επισκόπηση)

Η έρευνα που παρουσιάζεται σε αυτή τη διατριβή εστιάζει στον τομέα όπου συναντιούνται η Ανάλυση Συναισθημάτων (Sentiment Analysis), η Επεξεργασία Φυσικής Γλώσσας (Natural Language Processing - NLP) και το Κοινωνικό και Σημασιολογικό Διαδίκτυο (Social and Semantic Web). Πρόκειται για έναν τομέα που έχει αποκτήσει τεράστια σημασία τα τελευταία χρόνια, καθώς ο όγκος του γραπτού υλικού που παράγεται από τους χρήστες στο διαδίκτυο έχει φτάσει σε εξωπραγματικά επίπεδα. Η δυνατότητα αυτόματης ανάλυσης και κατανόησης των συναισθημάτων που εκφράζονται σε αυτά τα δεδομένα είναι ζωτικής σημασίας για πολλούς τομείς, από την επιχειρηματική νοημοσύνη και το μάρκετινγκ έως την κοινωνιολογική έρευνα και την πρόβλεψη των τάσεων.

.3.1 Βιβλιογραφική Επισκόπηση

Η ανάλυση γνώμης, γνωστή και ως Opinion Mining, έχει βαθιές ρίζες στην επιστημονική κοινότητα. Αρχικά, η έρευνα επικεντρωνόταν στην κατηγοριοποίηση των κειμένων στις διακριτές κλίμακες της πολικότητας (θετικό, αρνητικό, ουδέτερο) [16]. Ωστόσο, η πολυπλοκότητα της ανθρώπινης γλώσσας και η ανάγκη για πιο λεπτομερή κατανόηση των συναισθημάτων οδήγησαν στην ανάπτυξη πιο προηγμένων τεχνικών, οι οποίες δίνουν τη δυνατότητα να αναγνωρίζουν μια ευρύτερη πληθώρα συναισθημάτων (όπως χαρά, λύπη, θυμός, έκπληξη, φόβος, εμπιστοσύνη) [2].

.3.1.1 Αρχικές Προσεγγίσεις

Η πρώτη φάση της ανάλυσης συναισθημάτων βασιζόταν κατά κύριο λόγο σε λεξικά συναισθημάτων, όπου οι λέξεις και οι φράσεις αντιστοιχίζονταν σε συγκεκριμένες συναισθηματικές

κλίμακες[11]. Παραδείγματα τέτοιων λεξικών περιλαμβάνουν το SentiWordNet [5], το οποίο αντιστοιχίζει σε σύνολα synsets από το WordNet συναισθηματικές τιμές και το AFINN [14], ένα λεξικό με βάση τις βαθμολογίες. Παρά το γεγονός ότι οι μέθοδοι αυτές είναι απλές και αποτελεσματικές για βασικού επιπέδου εφαρμογές, συχνά αποτυγχάνουν να συλλάβουν το πλήρες πλαίσιο της λέξης, τον σαρκασμό ή την ειρωνεία, καθώς και τις διαπολιτισμικές ιδιαιτερότητες της γλώσσας. [13].

.3.1.2 Προσεγγίσεις Μηχανικής Μάθησης

Με την αύξηση της ποσότητας των διαθέσιμων επισημασμένων δεδομένων, οι προσεγγίσεις βασισμένες σε μηχανική μάθηση (Machine Learning - ML) ξεπέρασαν τις υπόλοιπες. Παραδοσιακοί αλγόριθμοι όπως οι Support Vector Machines (SVMs)[9], Naive Bayes [15] και η Logistic Regression χρησιμοποιήθηκαν με επιτυχία για την κατηγοριοποίηση συναισθημάτων. Αυτές οι μέθοδοι απαιτούν τη μετατροπή του κειμένου σε αριθμητικές αναπαραστάσεις, όπως το Bag-of-Words ή το TF-IDF (Term Frequency-Inverse Document Frequency) [20]. Παρά τις προόδους, τα παραδοσιακά μοντέλα ML είχαν συχνά δυσκολίες να λύσουν το πρόβλημα της σημασιολογικής πολυπλοκότητας.

.3.1.3 Deep learning και Transformers

Η σημαντικότερη εξέλιξη που έχει σημειωθεί στον τομέα της επεξεργασίας φυσικής γλώσσας (NLP) και, κατ' επέκταση, στην αναγνώριση συναισθημάτων, είναι η είσοδος της Βαθιάς Μάθησης (Deep Learning - DL), και, συγκεκριμένα, των Transformer architectures [22]. Μοντέλα όπως το BERT (Bidirectional Encoder Representations from Transformers) [4], το RoBERTa (Robustly Optimized BERT Pretraining Approach)[12] και το GPT (Generative Pre-trained Transformer) [19] έχουν αλλάξει τον τρόπο με τον οποίο κατανοούμε το φυσικό λόγο. Αυτά τα μοντέλα, που έχουν υποστεί προεκπαίδευση σε τεράστια όγκου δεδομένων, μπορούν να κατανοήσουν τον συνολικό σημασιολογικό πυρήνα των λέξεων έτσι, ώστε να πετυχαίνουν πρωτοφανείς επιδόσεις σε πολλές εργασίες NLP, συμπεριλαμβανομένης της αναγνώρισης συναισθημάτων[23].

Το μοντέλο RoBERTa-base-go_emotions, το οποίο χρησιμοποιείται στην παρούσα έρευνα, ανήκει σε αυτή την κατηγορία. Πρόκειται για μια βελτιστοποιημένη έκδοση του BERT, η οποία έχει προσαρμοστεί για την ταξινόμηση συναισθημάτων σε 27 διακριτές κατηγορίες, όπως αυτές προκύπτουν από το σύνολο δεδομένων GoEmotions[3]. Αυτό το σύνολο δεδομένων περιλαμβάνει σχόλια από το Reddit και έχει σχεδιαστεί ειδικά για την ανίχνευση λεπτομερών

συναισθημάτων, προσφέροντας έτσι μια πιο λεπτομερή ανάλυση σε συναισθηματικό επίπεδο στα δεδομένα του[3].

.3.1.4 Ανάλυση Συναισθημάτων στον Κοινωνικό Ιστό

Τον Κοινωνικό Ιστό, δηλαδή τις πλατφόρμες όπως το Facebook, το X (πρώην Twitter), το Reddit και τα φόρα των συζητήσεων είναι ένας πλούσιος, αλλά και δύσκολος, χώρος για να αναλύσει κανείς τα συναισθήματα [1]. Τα δεδομένα που προέρχονται από αυτές τις πηγές έχουν τα εξής χαρακτηριστικά:

- **Ανεπίσημη γλώσσα:** Συχνή χρήση συντομογραφιών, αργκό, emojis και ανορθόγραφων λέξεων.
- **Πολυπλοκότητα κειμένου:** Το νόημα ενός σχολίου μπορεί να εξαρτάται από προηγούμενα σχόλια ή από τη γενική διάθεση της συζήτησης.
- **Θορυβώδη δεδομένα:** Παρουσία spam, διαφημίσεων ή άσχετου περιεχομένου.

Αρκετές έρευνες έχουν πραγματοποιηθεί για τη χρήση της ανάλυσης συναισθημάτων σε συγκεκριμένες πλατφόρμες κοινωνικών μέσων. Για παράδειγμα, η ανάλυση των tweets έχει εφαρμοστεί για τον προσδιορισμό των εκλογικών αποτελεσμάτων [21] και για τη μέτρηση της δημόσιας γνώμης για προϊόντα [7]. Η ανάλυση των σχολίων σε φόρουμ ή πλατφόρμες όπως το Reddit παρέχει τη δυνατότητα να διερευνηθούν μακροσκελή άρθρα και πιο οργανωμένες συζητήσεις, αποκαλύπτοντας τις συναισθηματικές αντιδράσεις των χρηστών σε συγκεκριμένα θέματα ή αναρτήσεις [6].

.3.1.5 Συσχέτιση Συναισθημάτων με Χαρακτηριστικά Χρήστη

Εκτός από την απλή ανάλυση των συναισθημάτων, η σύγχρονη έρευνα επικεντρώνεται πλέον στην εξέταση του πώς τα συναισθήματα συνδέονται με άλλες συμπεριφορές και χαρακτηριστικά των χρηστών και του περιεχομένου. Έρευνες έχουν δείξει ότι η φήμη ενός χρήστη (π.χ., ο αριθμός των followers, το ιστορικό δημοσιεύσεών του) μπορεί να επηρεάσει τον τρόπο έκφρασης των συναισθημάτων του[8]. Αντίστοιχα, η δημοτικότητα μιας δημοσίευσης (όπως αυτή που μετρά likes, shares, upvotes) σχετίζεται συχνά με το συναισθηματικό πλαίσιο των σχολίων που λαμβάνει[10]. Η γεωγραφική θέση των χρηστών είναι άλλο ένα σημαντικό στοιχείο, καθώς οι διαφορετικές πολιτισμικές ομάδες μπορεί να έχουν διαφορετικούς τρόπους έκφρασης των συναισθημάτων τους [18]. Η ανάλυση των upvotes και downvotes στα σχόλια παρέχει

άμεσα στοιχεία για το αν το συναισθηματικό περιεχόμενο έχει γενική αποδοχή ή απορρίπτεται από τους χρήστες, και μας βοηθά να καταλάβουμε πώς κάποια συναισθήματα επηρεάζουν την κοινότητα [17].

.4 Το Μοντέλο SamLowe/roberta-base-go_emotions

Για την ανάλυση των συναισθημάτων, όπως εξηγήθηκε προηγουμένως, χρησιμοποιήθηκε το μοντέλο SamLowe/roberta-base-go_emotions. Πρόκειται για ένα προηγμένο εργαλείο βαθιάς μάθησης, βασισμένο στην αρχιτεκτονική Transformer, το οποίο αποτελεί έναν από τους ορόσημους στον τομέα της Επεξεργασίας Φυσικής Γλώσσας (NLP).

Το εν λόγω μοντέλο αποτελεί μια βελτιωμένη έκδοση του RoBERTa, το οποίο, από την άλλη, είναι μια περαιτέρω αναβάθμιση του BERT. Οι βελτιώσεις που έχει επιφέρει το RoBERTa περιλαμβάνουν το γεγονός ότι έχει εκπαιδευτεί σε πολύ πιο μεγάλα datasets και ότι χρησιμοποιεί μια πολύ πιο δυναμική στρατηγική "masking" λέξεων, γεγονός που το καθιστά πιο αποτελεσματικό και ισχυρό στην κατανόηση κάθε γλωσσικού πλαισίου. Το καθοριστικής σημασίας στοιχείο για την εργασία είναι ότι το εν λόγω μοντέλο έχει γίνει υπόλοιπο/διάπλαση (fine-tuned) με χρήση υπο-datasets που προέρχονται από το Stack Overflow Dataset. Το συγκεκριμένο σύνολο δεδομένων περιλαμβάνει σχόλια από το Stack Overflow τα οποία έχουν ταξινομηθεί σε 27 διαφορετικές κατηγορίες συναισθημάτων (π.χ., χαρά, θλίψη, θυμός, έκπληξη, ουδέτερο). Η διαδικασία αυτή επιτρέπει στο μοντέλο να αναγνωρίζει με πολύ μεγαλύτερη ακρίβεια μια ευρύτερη γκάμα συναισθημάτων, από το απλό θετικό/αρνητικό/ουδέτερο.

Η λειτουργία του:

Όταν εισάγεται ένα σχόλιο στο μοντέλο μέσω του pipeline της Hugging Face transformers:

- Το κείμενο διασπάται σε tokens.
- Αυτά τα tokens μετατρέπονται σε αριθμητικές αναπαραστάσεις (embeddings) που "συλλαμβάνουν" τη σημασία και το πλαίσιο.
- Οι Transformer στρώσεις επεξεργάζονται αυτές τις αναπαραστάσεις.

- Ένα τελικό επίπεδο ταξινόμησης παράγει πιθανότητες για κάθε ένα από τα 27 συναισθήματα.

Γιατί επιλέχθηκε:

Η επιλογή του SamLowe/roberta-base-go_emotions προσφέρει:

- Υψηλή ακρίβεια στην ανίχνευση συναισθημάτων. Λεπτομερής κατανόηση των συναισθηματικών αποχρώσεων, κάτι που είναι απαραίτητο για την πολύπλοκη ανάλυσή που πραγματοποιείται.
- Την ικανότητα να χειριστεί το πλαίσιο του κειμένου.
- Ευκολία χρήσης χάρη στην ενσωμάτωσή του στο Hugging Face App.
- Με αυτόν τον αλγόριθμο, διενεργείται μια βαθιά και αξιόπιστη ανάλυση των συναισθημάτων στα σχόλια του Κοινωνικού και Σημασιολογικού Ιστού.

.5 Μεθοδολογική Διαδικασία

Η παρούσα μελέτη ακολουθεί μια συστηματική μεθοδολογία για την ανάλυση συναισθημάτων σε δεδομένα Κοινωνικού και Σημασιολογικού Ιστού. Η διαδικασία περιλαμβάνει τρία βασικά στάδια: την περιγραφή των δεδομένων, την προεπεξεργασία αυτών και τη διεξαγωγή πειραμάτων μάθησης για την εξαγωγή συναισθημάτων και την περαιτέρω ανάλυσή τους.

.5.1 Περιγραφή Δεδομένων

Για τους σκοπούς αυτής της ανάλυσης, χρησιμοποιήθηκαν τρία υπό-datasets (comments.csv, users.csv, posts_answers.csv και το παραγόμενο comments_with_emotions.csv) του Stack Overflow Data Dataset, τα οποία φορτώθηκαν με τη χρήση της βιβλιοθήκης pandas.

- **comments.csv:** Αυτό το αρχείο περιέχει τα σχόλια των χρηστών. Αποτελεί το πρωτεύον σύνολο δεδομένων για την ανάλυση συναισθημάτων. Κατά την αρχική εξερεύνηση, ελέγχθηκε για missing-values, duplicates και η μοναδικότητα των user_id και post_id για να διασφαλιστεί η ακεραιότητα των δεδομένων. (Περιέχει στήλες όπως id, text, post_id, user_id και score).
- **posts_answers.csv:** Αυτό το αρχείο περιλαμβάνει δεδομένα σχετικά με τις αναρτήσεις στις οποίες ανήκουν τα σχόλια, καθώς και τις απαντήσεις τους. Ένα κρίσιμο

χαρακτηριστικό εδώ είναι η βαθμολογία (score) της ανάρτησης, η οποία υποδηλώνει τη δημοτικότητα ή τη χρησιμότητά της. (Περιέχει στήλες όπως id, comment_count, favorite_count και score).

- **users.csv:** Αυτό το αρχείο περιέχει πληροφορίες για τους χρήστες που αλληλεπιδρούν στην πλατφόρμα. Βασικά χαρακτηριστικά περιλαμβάνουν τη φήμη (reputation) του χρήστη, η οποία συχνά αντανακλά την επιρροή ή τη δραστηριότητά του στην κοινότητα, καθώς και την τοποθεσία (location) του. Επίσης, αναλύθηκε η κατανομή της φήμης και εντοπίστηκαν οι δέκα κορυφαίες τοποθεσίες χρηστών, παρέχοντας ένα γεωγραφικό πλαίσιο στην ανάλυση. (Περιέχει στήλες όπως id, display_name, location, reputation, up_votes και down_votes).

.5.2 Περιγραφή Προεπεξεργασίας

Η αποτελεσματικότητα οποιασδήποτε ανάλυσης κειμένου εξαρτάται σε μεγάλο βαθμό από την ποιότητα της προεπεξεργασίας των δεδομένων. Στην παρούσα εργασία, εφαρμόστηκε μια σειρά από βήματα "καθαρισμού" και προετοιμασίας των κειμενικών δεδομένων, καθώς και μετασχηματισμών των αριθμητικών δεδομένων, ώστε να είναι κατάλληλα για την ανάλυση συναισθημάτων.

.5.2.1 Καθαρισμός Κειμένου (Text Cleaning)

Τα σχόλια των χρηστών, τα οποία συχνά περιέχουν "θόρυβο" από τον ανεπίσημο χαρακτήρα του διαδικτυακού λόγου, υποβλήθηκαν στους ακόλουθους μετασχηματισμούς:

- **Αφαίρεση URL και Emojis:** Εφαρμόστηκαν τακτικές (**URL_PATTERN**, **EMOJI_PATTERN**) για την αφαίρεση URLs και emojis, καθώς αυτά τα στοιχεία, αν και ενίοτε μεταφέρουν συναισθήματα, μπορούν να περιπλέξουν την επεξεργασία από τα μοντέλα NLP, ειδικά όταν η ανάλυση εστιάζει στο λεκτικό περιεχόμενο.
- **Μετατροπή σε Πεζά (Lowercasing):** Όλα τα γράμματα μετατράπηκαν σε πεζά για να διασφαλιστεί ότι οι ίδιες λέξεις αναγνωρίζονται ως τέτοιες ανεξαρτήτως κεφαλαιοποίησης (π.χ., "Excellent" και "excellent").
- **Φιλτράρισμα "meaningful" Κειμένου:** Χρησιμοποιήθηκε η συνάρτηση **has_meaningful_text** (η οποία βασίζεται στο **MEANINGFUL_TEXT_PATTERN**)

για να διασφαλιστεί ότι κάθε σχόλιο περιέχει τουλάχιστον μία λέξη με τρεις ή περισσότερους αλφαριθμητικούς χαρακτήρες, απορρίπτοντας σχόλια που αποτελούνται μόνο από σύμβολα ή πολύ σύντομες, ανούσιες εκφράσεις.

- **Αφαίρεση Διπλοτύπων και Μικρού Μήκους Σχολίων:** Τα διπλότυπα σχόλια αφαιρέθηκαν για να αποφευχθεί η υπερεκπροσώπηση συγκεκριμένων εκφράσεων. Επιπλέον, σχόλια με λιγότερες από τρεις λέξεις (`min_words=3`) απορρίφθηκαν, καθώς θεωρήθηκαν ότι παρέχουν ανεπαρκές περιεχόμενο για ουσιαστική συναισθηματική ανάλυση.

Αυτά τα βήματα υλοποιήθηκαν μέσω των συναρτήσεων `clean_text` και `clean_comments`, οι οποίες εφαρμόστηκαν στο `DataFrame` των σχολίων.

.5.2.2 Προετοιμασία Δεδομένων για Ανάλυση NLP

Για την περαιτέρω ανάλυση και εξαγωγή χαρακτηριστικών, εφαρμόστηκαν οι ακόλουθες τεχνικές NLP:

- **Tokenization:** Τα σχόλια διαχωρίστηκαν σε μεμονωμένες λέξεις (tokens) χρησιμοποιώντας τη συνάρτηση `word_tokenize` της βιβλιοθήκης `nlTK`.
- **Αφαίρεση Stop Words:** Λέξεις που είναι κοινές στη γλώσσα και δεν προσθέτουν σημαντικό σημασιολογικό νόημα (π.χ., “the”, “a”, “is”, “and”) αφαιρέθηκαν χρησιμοποιώντας τη λίστα `stopwords` της `nlTK`. Αυτό βοηθά στην εστίαση στις πιο ουσιαστικές λέξεις.
- **Λημματοποίηση (Lemmatization):** Κάθε λέξη αναχώρησε στην αρχική της λεξικογραφική μορφή (π.χ., “running”, “ran” γίνονται “run”) με τη χρήση του `WordNetLemmatizer`. Αυτό μειώνει την πολυπλοκότητα του λεξιλογίου και ομαδοποιεί τις σημασιολογικά παρόμοιες λέξεις. Τα λημματοποιημένα tokens αποθηκεύτηκαν στη στήλη `lemmatized_tokens`.
- **Εξαγωγή Bigrams:** Πέραν των μεμονωμένων λέξεων, αναλύθηκαν και τα *bigrams* (ζεύγη γειτονικών λέξεων), καθώς συχνά μεταφέρουν περισσότερο νόημα (π.χ., “not good”). Οι 20 πιο συχνές λέξεις και bigrams οπτικοποιήθηκαν.
- **TF-IDF Vectorization:** Για την ποσοτικοποίηση της σημασίας των λέξεων στα σχόλια, εφαρμόστηκε η διανυσματοποίηση TF-IDF (Term Frequency-Inverse Document Frequency) με χρήση του `TfidfVectorizer`. Ορίστηκαν `max_features=28000`,

`stop_words='english'` και `ngram_range=(1,1)` ώστε η ανάλυση να επικεντρωθεί σε unigrams.

.5.2.3 Κατηγοριοποίηση Αριθμητικών Δεδομένων (Binning)

Για να διευκολυνθεί η ανάλυση των συσχετίσεων μεταξύ των συναισθημάτων και των αριθμητικών χαρακτηριστικών, οι συνεχείς αριθμητικές στήλες κατηγοριοποιήθηκαν σε διακριτές ομάδες (*binning*) χρησιμοποιώντας τη συνάρτηση `pd.cut` μέσω της βοηθητικής συνάρτησης `bin_column`.

- **Ομαδοποίηση Φήμης Χρήστη (`rep_group`):** Η φήμη των χρηστών (`reputation`) χωρίστηκε σε πέντε κατηγορίες: “Very Low”, “Low”, “Medium”, “High”, “Very High”.
- **Ομαδοποίηση Βαθμολογίας Ανάρτησης (`post_score_group`):** Η βαθμολογία της ανάρτησης (`score_post`) κατηγοριοποιήθηκε σε “Negative”, “Low”, “Moderate”, “High”, “Very High”.
- **Ομαδοποίηση Βαθμολογίας Σχολίου (`comment_score_group`):** Η βαθμολογία του σχολίου (`score`) κατηγοριοποιήθηκε ομοίως σε “Negative”, “Low”, “Medium”, “High”, “Very High”.

Αυτές οι κατηγοριοποιήσεις επέτρεψαν την οπτικοποίηση των κατανομών των συναισθημάτων σε σχέση με διαφορετικά επίπεδα φήμης και δημοτικότητας.

.5.3 Πειράματα Μάθησης

Τα πειράματα μάθησης επικεντρώθηκαν στην εφαρμογή του μοντέλου ανάλυσης συναισθημάτων και στην ενοποίηση των αποτελεσμάτων με τα `metadata` για την εξαγωγή ουσιαστικών συμπερασμάτων.

.5.3.1 Ανάλυση Συναισθημάτων με GoEmotions Model

Η βασική ανάλυση συναισθημάτων, όπως αναφερθηκε και προηγουμένος, πραγματοποιήθηκε χρησιμοποιώντας το προεκπαιδευμένο μοντέλο `SamLowe/roberta-base-go_emotions` από τη βιβλιοθήκη `transformers` της Hugging Face.

- **Εφαρμογή Μοντέλου:** Το μοντέλο φορτώθηκε ως `text-classification pipeline` με `top_k=None` για την εξαγωγή πιθανοτήτων για όλες τις 27 κατηγορίες συναισθημάτων.
- **Batch Inference:** Για την επιτάχυνση της πρόβλεψης σε μεγάλο όγκο σχολίων, χρησιμοποιήθηκε η συνάρτηση `batch_infer` με `batch_size=16`. Κάθε σχόλιο περικόπηκε στα 512 tokens, το μέγιστο επιτρεπτό από το μοντέλο.
- **Αποθήκευση Αποτελεσμάτων:** Τα προβλεπόμενα συναισθήματα αποθηκεύτηκαν στη στήλη `emotions` του `DataFrame comments` και εξήχθησαν σε αρχείο `comments_with_emotions.csv` για μελλοντική χρήση.

.5.3.2 Εξαγωγή και Δομή Δεδομένων Συναισθημάτων

Αφού ανιχνεύθηκαν τα συναισθήματα, ακολούθησαν τα εξής βήματα:

- **Φιλτράρισμα με Thresholding:** Η συνάρτηση `extract_labels` εφάρμοσε κατώφλι `threshold=0.3` για να διατηρηθούν μόνο οι ετικέτες με υψηλή βεβαιότητα στη στήλη `emotion_labels`.
- **“Explode” Πολλαπλών Συναισθημάτων:** Χρησιμοποιήθηκε η μέθοδος `explode("emotion_labels")` ώστε κάθε συναίσθημα να αποτελεί ξεχωριστή γραμμή, διευκολύνοντας την ανάλυση συχνοτήτων και συσχετίσεων.

.5.3.3 Συγχώνευση Δεδομένων για Ανάλυση

Για τη μελέτη των σχέσεων μεταξύ συναισθημάτων και μεταδεδομένων, τα επεξεργασμένα σχόλια συγχωνεύτηκαν με τα δεδομένα χρηστών και αναρτήσεων.

- **Ενοποίηση DataFrames:** Η συνάρτηση `merge.all` εφάρμοσε *left merge* του `comments_exp` με τα `users` και `posts`, με βάση τα `user_id` και `post_id`. Πριν τη συγχώνευση, τα `user_id` και `post_id` μετατράπηκαν σε αριθμητικό τύπο για αποφυγή σφαλμάτων.

Αυτή η συγχωνευμένη δομή αποτέλεσε τη βάση για όλες τις περαιτέρω αναλύσεις και οπτικοποιήσεις, επιτρέποντας τη διερεύνηση σύνθετων σχέσεων μεταξύ εκφρασμένων συναισθημάτων και χαρακτηριστικών χρηστών και περιεχομένου.

.6 Ανάλυση και Αξιολόγηση

Η φάση της ανάλυσης και αξιολόγησης επικεντρώνεται στην εξαγωγή και ερμηνεία των συναισθημάτων που προκύπτουν από την ανάλυση κειμένων, καθώς και στη διερεύνηση των συνδέσεων μεταξύ αυτών των συναισθημάτων και των χαρακτηριστικών των χρηστών και των δημοσιεύσεών τους. Μιας και χρησιμοποιήθηκε ένα προ-εκπαιδευμένο μοντέλο για την ανίχνευση των συναισθημάτων, η αξιολόγηση επικεντρώνεται στην ποιότητα των πληροφοριών που ανακτήθηκαν από τα δεδομένα και όχι στην απόδοση του μοντέλου. Η ανάλυση γίνεται πιο κατανοητή μέσω διαφόρων γραφημάτων και διαγραμμάτων.

.6.1 Εξερεύνηση και Χαρακτηριστικά Δεδομένων

.6.1.1 Επισκόπηση Συνόλων Δεδομένων

- **Σχόλια (comments):** Έγινε έλεγχος των τίτλων, των κενών τιμών (missing values) και των διπλο-αναρτήσεων. Έτσι, προέκυψε ο συνολικός αριθμός των σχολίων αλλά και ο αριθμός των μοναδικών χρηστών και των αναρτήσεων.
- **Αναρτήσεις (posts):** Σε αντίστοιχη διαδικασία, εξετάστηκαν οι τίτλοι, τα missing values και τα διπλο-αναρτήσεις. Επιπλέον, δόθηκε έμφαση στη διανομή των βαθμολογιών.

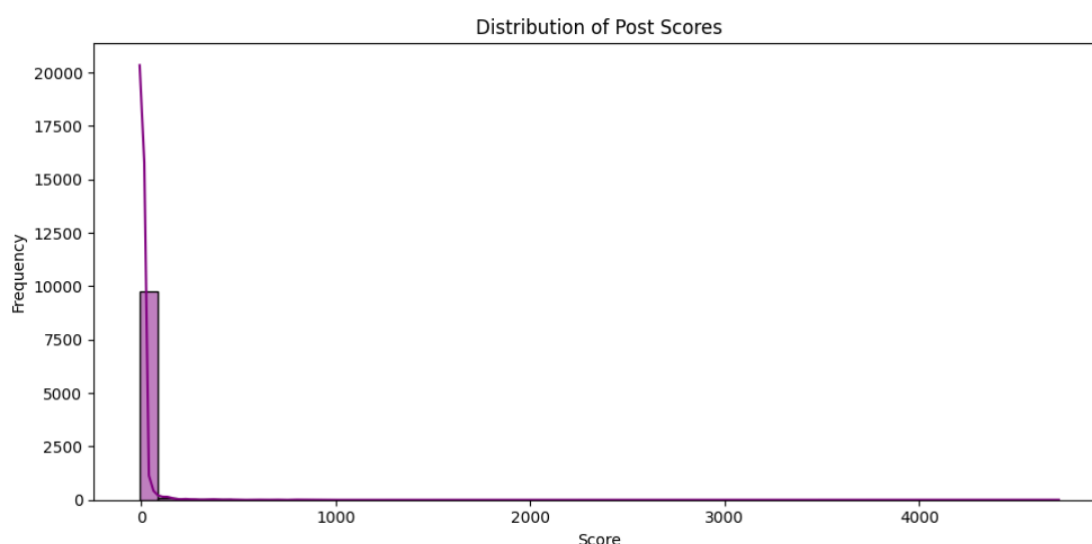


Figure 1: *Distribution of post scores*

- **Χρήστες (users):** Ελέγχθηκαν οι επικεφαλίδες των στηλών, τα κενά και τα διπλό-τυπα των δεδομένων. Αναλύθηκε η κατανομή της φήμης των χρηστών και ορίστηκαν οι δέκα κορυφαίες τοποθεσίες.

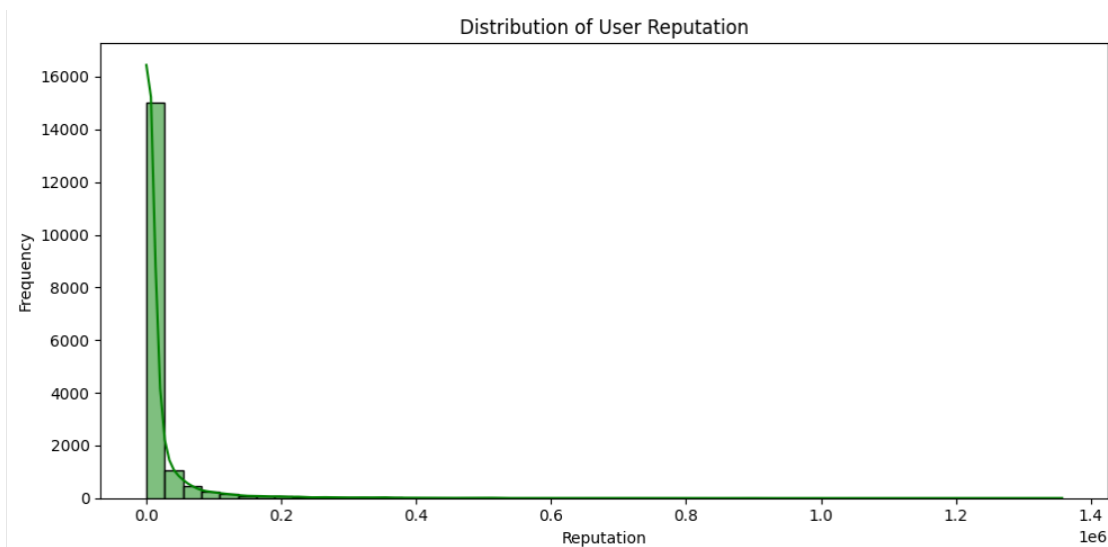


Figure 2: *Distribution of user reputation*

```
Top 10 User Locations:
location
India                248
Germany              230
United States        174
United Kingdom       172
London, United Kingdom 159
Berlin, Germany      106
Netherlands          106
France                99
Bangalore, Karnataka, India 89
Paris, France         77
Name: count, dtype: int64
```

Figure 3: *Top 10 user locations*

.6.1.2 Χαρακτηριστικά Κειμένου των Σχολίων

- **Κατανομή Μήκους Σχολίων:** Ο υπολογισμός του αριθμού των λέξεων σε κάθε σχόλιο και η οπτικοποίησή του.

Number of comments: 20958
 Unique users: 17216
 Missing values in text: 0

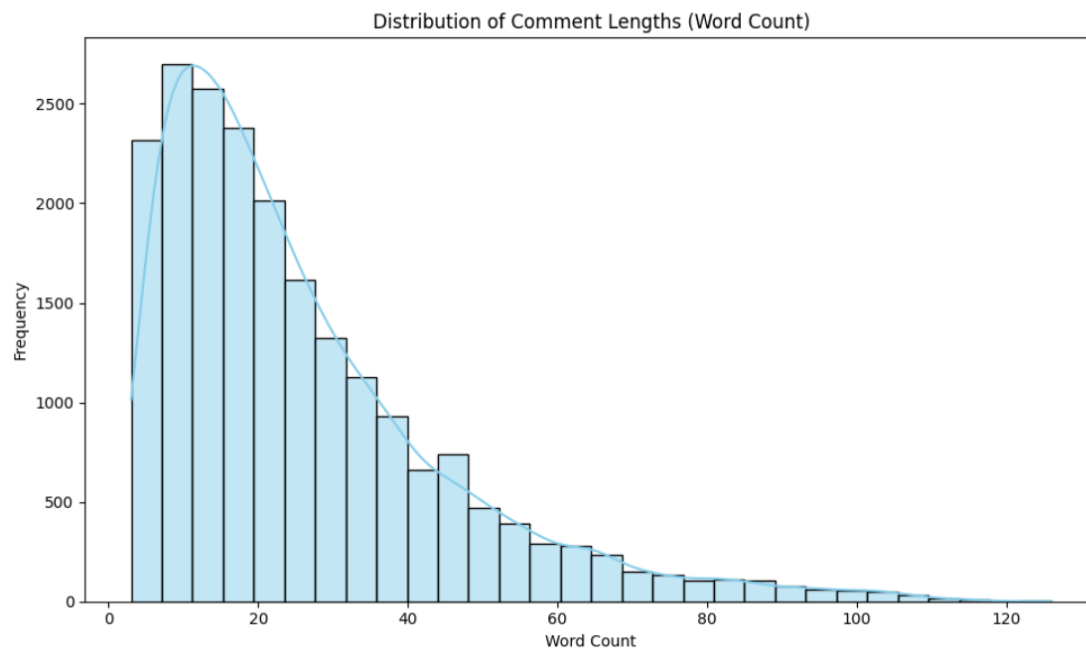


Figure 4: *Word count*

- **Ανάλυση Συχνότητας Λέξεων και Bigrams:** Προσδιορίστηκαν οι πιο συχνές λέξεις και φράσεις.

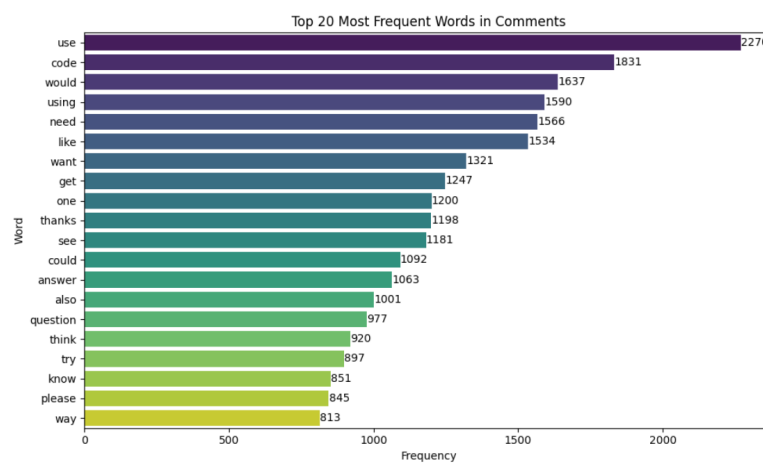


Figure 5: *Συχνότητα λέξεων*

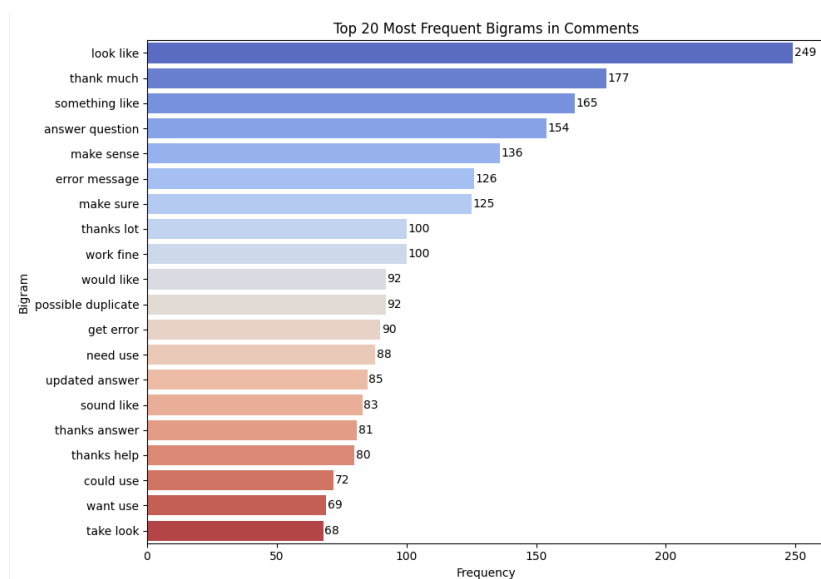


Figure 6: *Most frequent bigrams*

.6.1.3 Γενική Κατανομή Συναισθημάτων

Αρχικά, εξετάστηκε η συνολική κατανομή των 27 διακριτών συναισθημάτων που ανιχνεύθηκαν στα σχόλια.

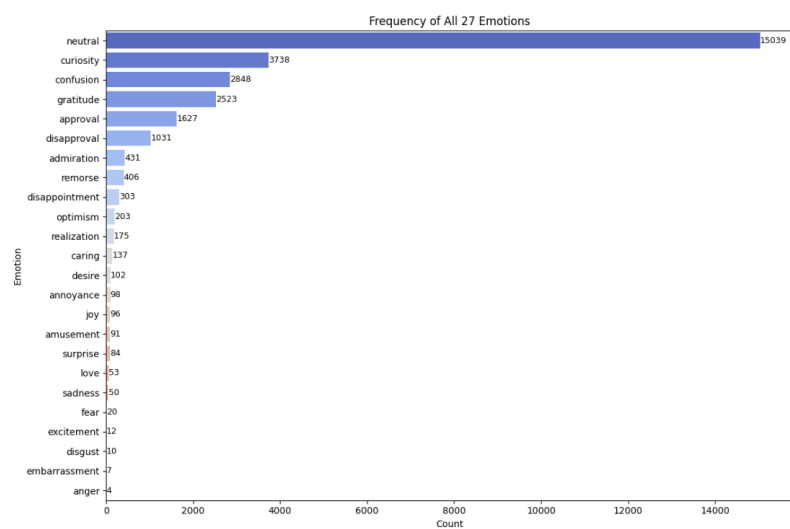


Figure 7: *Frequency of All 27 Emotions (Bar Plot)*

Το διάγραμμα παρουσιάζει ποια συναισθήματα εκφράζονται περισσότερο σε όλα τα σχόλια.

Δίνει μια ιδέα για το ποια συναισθήματα είναι πιο δημοφιλή στις Stack Overflow συζητήσεις. Υπολογίζεται ότι το άκρως αντικειμενικό συναίσθημα θα είναι ανάμεσα στα πιο συχνά, καθώς πολλά σχόλια μπορεί να μην έχουν έντονο συναισθηματικό φορτίο ή να είναι απλά πληροφοριακά.

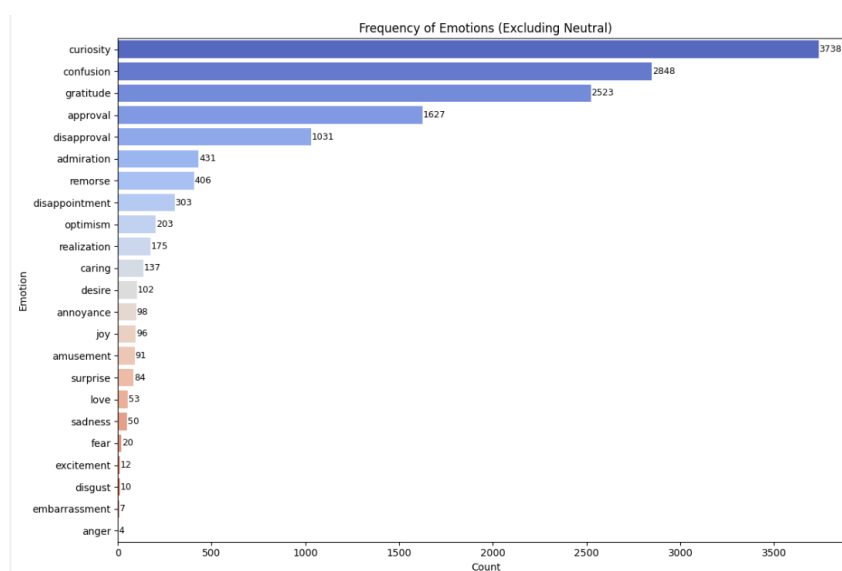


Figure 8: *Frequency of Emotions (Excluding Neutral)*

Αυτό το κείμενο αποκλείει την κατηγορία ουδέτερο, καθώς φανερώνει με κυρίαρχο τρόπο θετικά, αρνητικά ή ακόμα και ανάμεικτα συναισθήματα. Για παράδειγμα, η συχνή αναφορά στην ευγνωμοσύνη ή την έγκριση μπορεί να υποδηλώνει μια γενικά θετική ατμόσφαιρα, ενώ η αναφορά στη δυσaráεσκεια μπορεί να αποτελεί ένδειξη έντασης.

.6.1.4 Σχέση Συναισθημάτων με τη Δημοτικότητα του Περιεχομένου

Μελετήθηκε η σχέση των συναισθημάτων που εκφράζονται στα σχόλια με τη δημοτικότητά τους και τη δημοτικότητα των δημοσιεύσεων, στις οποίες ανήκουν. Τα σκορ (scores) ταξινομήθηκαν σε κατηγορίες (Negative, Low, Medium, High, Very High).

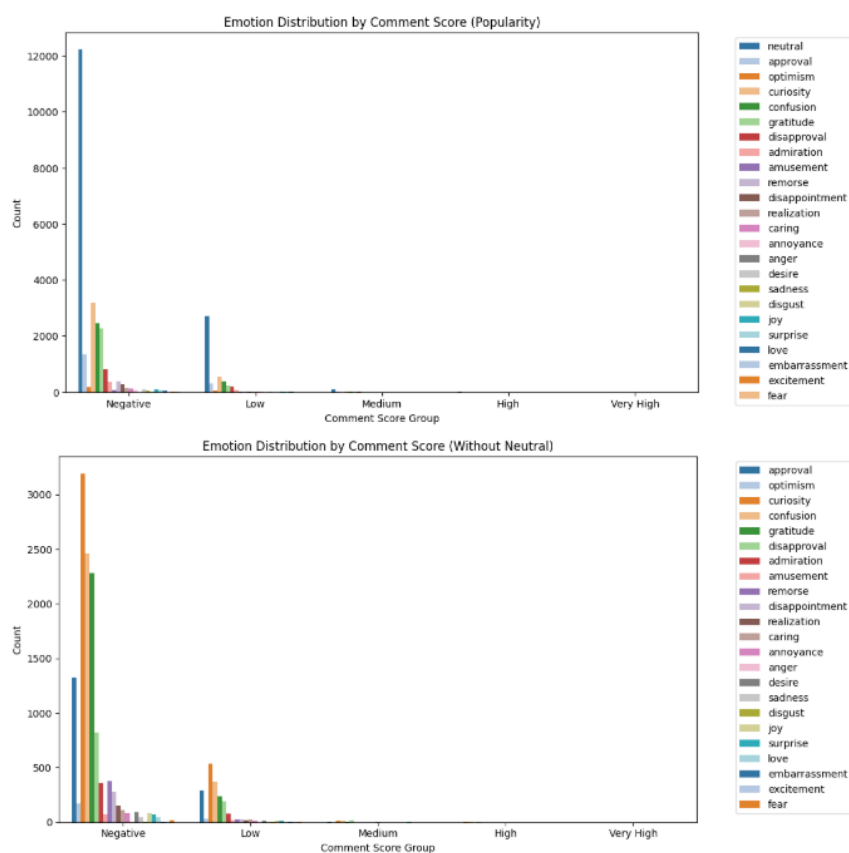


Figure 9: *Emotion Distribution by Comment Score*

Το παραπάνω γράφημα δείχνει την κατανομή συναισθημάτων σε σχέση με τη βαθμολογία των σχολίων.

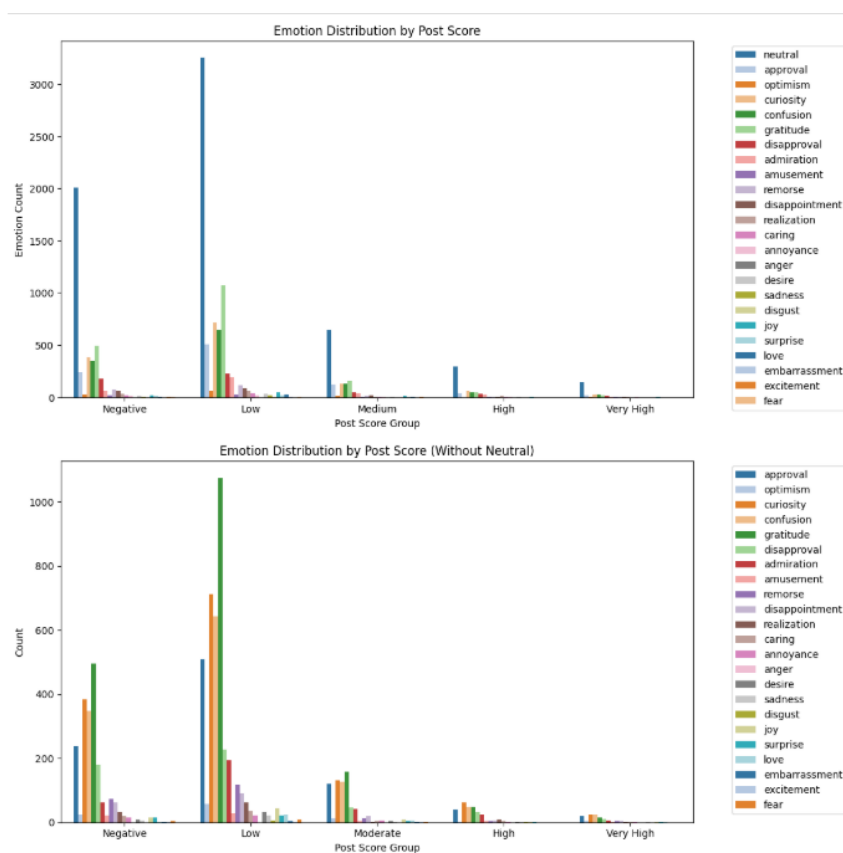


Figure 10: *Emotion Distribution by Post Score*

Αντίστοιχα, αυτό το γράφημα δείχνει την κατανομή συναισθημάτων σε σχέση με τη βαθμολογία των αναρτήσεων στις οποίες ανήκουν τα σχόλια.

.6.1.5 Γεωγραφικές Διαφοροποιήσεις ανά συναίσθημα

Αναζητήθηκαν διαφορές στην έκφραση συναισθημάτων ανάλογα με τη γεωγραφική τοποθεσία των χρηστών, με εστίαση στις 10 κορυφαίες περιοχές.

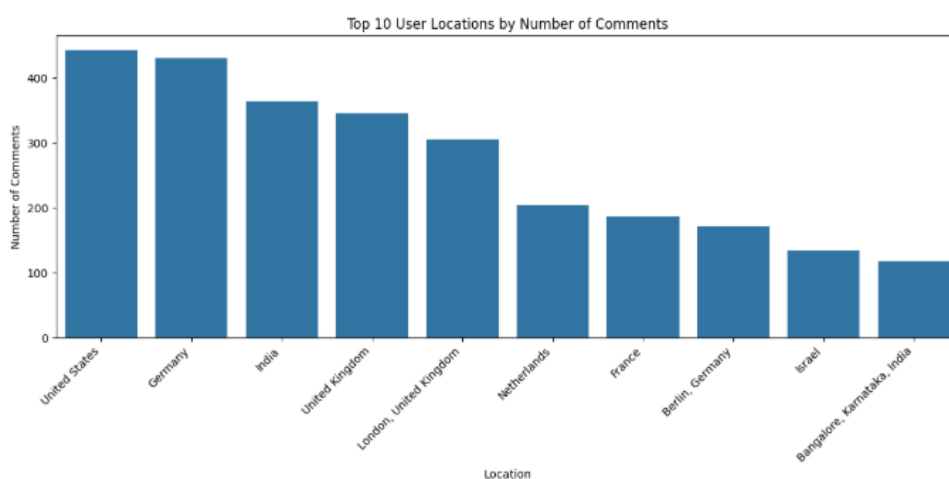


Figure 11: Emotions by User Location (Top 10 Locations)

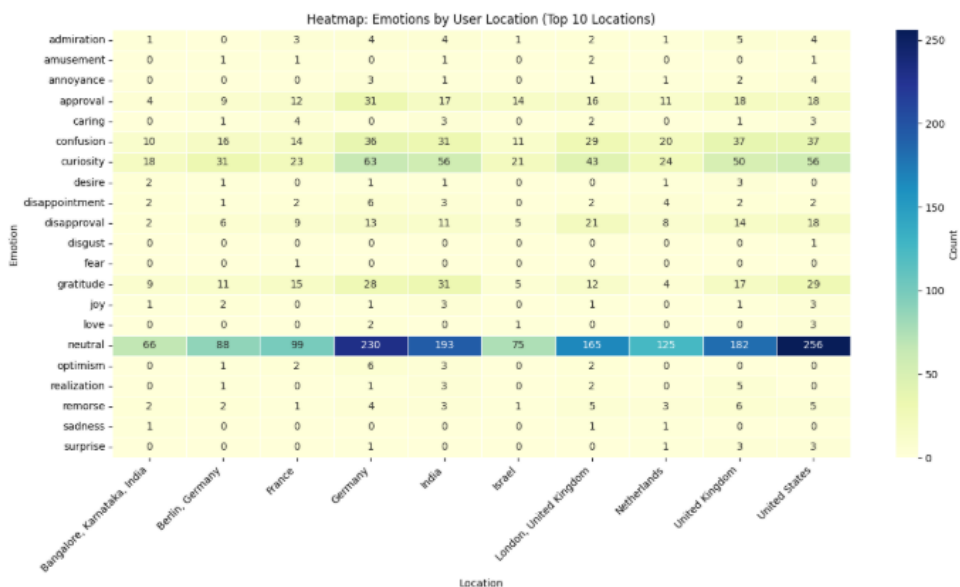


Figure 12: Emotions Heatmap by User Location

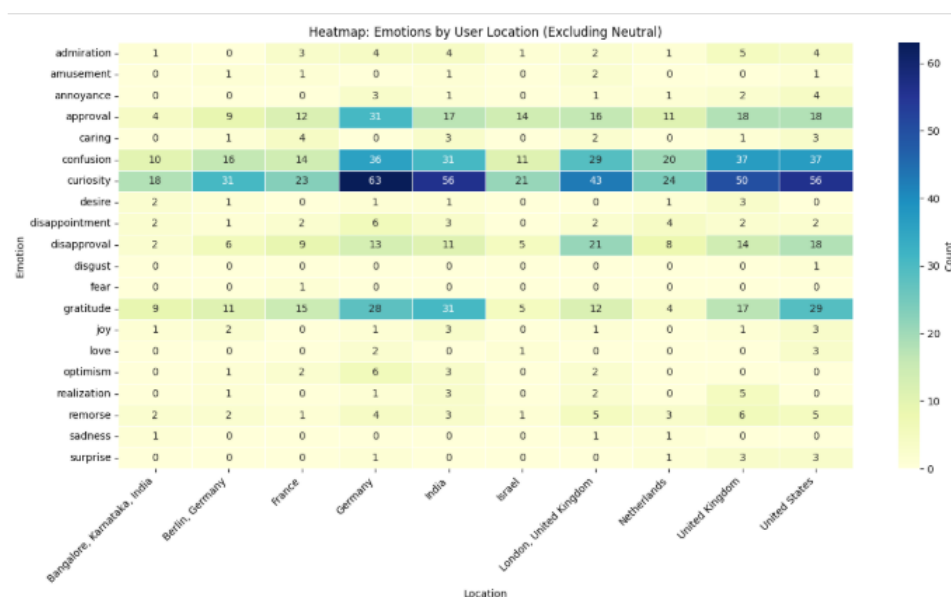


Figure 13: Emotions Heatmap by User Location (Excluding Neutral)

.6.1.6 Συναισθήματα και Αντίδραση της Κοινότητας

Η ανάλυση εξετάζει τη σχέση μεταξύ συναισθημάτων και της αντίδρασης της κοινότητας (upvotes/downvotes):

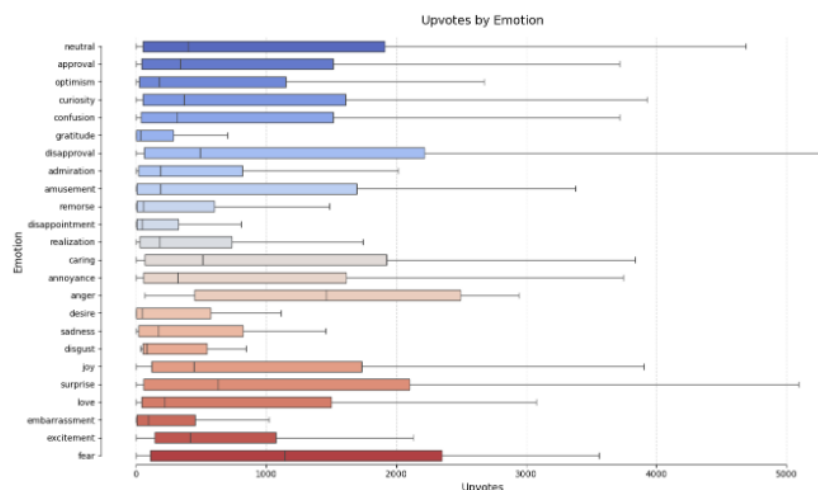


Figure 14: Upvotes by Emotion (Box Plot)

Ποιες συναισθηματικές εκφράσεις τείνουν να επιβραβεύονται με περισσότερα upvotes;

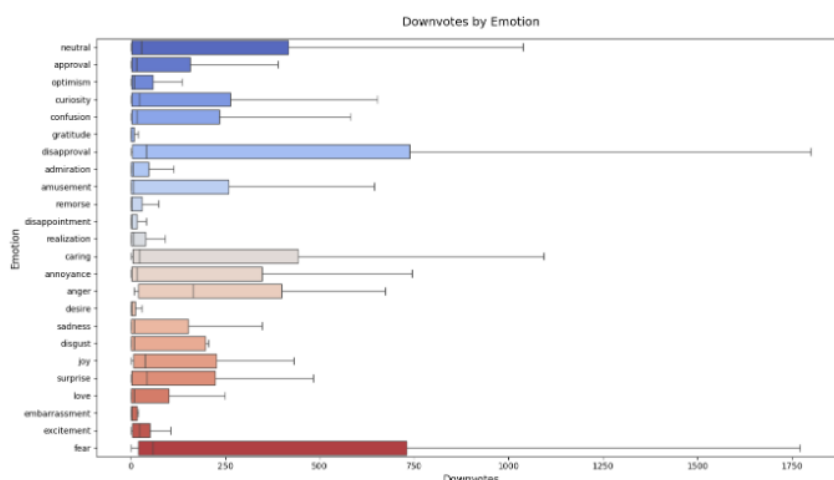


Figure 15: *Downvotes by Emotion (Box Plot)*

Αντίστοιχα, ποιες συναισθηματικές εκφράσεις τείνουν να προκαλούν αρνητική ανατροφοδότηση;

Διαβάζοντας προσεκτικά αυτήν την ανάλυση δεδομένων και με βάση τις οπτικοποιήσεις, μπορούμε να καταλήξουμε σε σημαντικά συμπεράσματα σχετικά με τον τρόπο έκφρασης των συναισθημάτων στον διαδικτυο (Κοινωνικός και Σημασιολογικός Ιστός) και τον τρόπο που σχετίζεται με τα χαρακτηριστικά των χρηστών και το περιεχόμενο.

.7 Συμπεράσματα

Η έρευνα αυτή επέδειξε την ισχύ των μοντέλων Transformer, όπως το *SamLowe/roberta-base-go_emotions*, σε λεπτομερή αναλύσεις συναισθημάτων που βασίζονται σε δεδομένα κοινωνικού δικτύου. Χάρη στην ανάλυση, καταφέρθηκε να απεικονιστεί ο τρόπος διάχυσης των 27 διαφορετικών συναισθημάτων στα σχόλια των χρηστών, αποκαλύπτοντας έτσι τον πλούτο των συναισθημάτων. Διαπιστώθηκε η σχέση των συναισθημάτων με τη δημοτικότητα των σχολίων και των δημοσιεύσεων, καθώς και τα γεωγραφικά χαρακτηριστικά των χρηστών. Η έρευνα σχετικά με τα κορυφαία σχόλια και τις προβλέψεις που πραγματοποιήθηκαν με υψηλή επιτυχία πρόσφερε μια ολοκληρωμένη εικόνα της απόδοσης του μοντέλου. Επίσης, η ανάλυση των upvotes και downvotes ανά συναισθηματικό ήταν πολύτιμη για την κατανόηση της κοινοτικής αντίδρασης σε συγκεκριμένες συναισθηματικές εκφράσεις. Όλα αυτά τα ευρήματα έχουν μεγάλη σημασία για την κατανόηση της δυναμικής της επικοινωνίας στις διαδικτυακές πλατφόρμες.

.8 Προτάσεις για Μελλοντικές Βελτιώσεις

Βασιζόμενοι στα ευρήματα και την εμπειρία από την παρούσα μελέτη, προτείνονται οι ακόλουθες βελτιώσεις και επεκτάσεις για μελλοντική έρευνα:

.8.1 Εμπλουτισμός Δεδομένων

- **Ενσωμάτωση Μεταδεδομένων Χρόνου:** Η ανάλυση της χρονικής διάστασης (π.χ. η ώρα της ημέρας, η ημέρα της εβδομάδας, εποχικές τάσεις) θα μπορούσε να αποκαλύψει ενδιαφέροντα μοτίβα στην έκφραση συναισθημάτων.
- **Πληροφορίες Θέματος/Κατηγορίας Ανάρτησης:** Η κατηγοριοποίηση των αναρτήσεων ανά θέμα θα επέτρεπε τη διερεύνηση συναισθηματικών διαφορών μεταξύ διαφορετικών γνωστικών πεδίων ή ενδιαφερόντων.
- **Ανάλυση Συνομιλίας (Conversational Context):** Αντί της ανάλυσης μεμονωμένων σχολίων, η εξέταση ολόκληρων νημάτων συζητήσεων θα μπορούσε να αποκαλύψει τη συναισθηματική εξέλιξη των αλληλεπιδράσεων.

.8.2 Ποιοτική Ανάλυση & Επαλήθευση

- **Συνέντευξη Χρηστών:** Ποιοτική έρευνα μέσω συνεντεύξεων με χρήστες θα μπορούσε να παρέχει βαθύτερες γνώσεις για το γιατί εκφράζουν συγκεκριμένα συναισθήματα και πώς αντιλαμβάνονται τις αντιδράσεις της κοινότητας.

.8.3 Προηγμένες Τεχνικές Ανάλυσης

- **Μοντέλα Πολλαπλών Γλωσσών:** Για πλατφόρμες με πολλαπλές γλώσσες, η χρήση πολυγλωσσικών μοντέλων Transformer θα επέτρεπε την ανάλυση συναισθημάτων χωρίς περιορισμό στη γλώσσα.
- **Ανίχνευση Σαρκασμού/Ειρωνείας:** Η ενσωμάτωση ειδικών μοντέλων ή τεχνικών για την ανίχνευση σαρκασμού και ειρωνείας θα βελτίωνε σημαντικά την ακρίβεια της συναισθηματικής ανάλυσης.

.9 Βιβλιογραφία

- [1] Charu C. Aggarwal. *Social Network Data Analytics*. Springer, 2011.
- [2] Erik Cambria et al. “New avenues in opinion mining and sentiment analysis”. In: *IEEE Intelligent Systems* 33.2 (2018), pp. 15–23.
- [3] Dorottya Demszky et al. “GoEmotions: A Dataset of Fine-Grained Emotions”. In: *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 2020, pp. 4048–4063.
- [4] Jacob Devlin et al. “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding”. In: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Vol. 1. 2019, pp. 4171–4186.
- [5] Andrea Esuli and Fabrizio Sebastiani. “SentiWordNet: A Lexical Resource for Sentiment Analysis”. In: *Proceedings of LREC*. Vol. 6. 2006, pp. 417–422.
- [6] Jing Fan et al. “Topic-Sentiment Analysis for Online Reviews Based on LDA and SVM”. In: *Journal of Physics: Conference Series* 1683.1 (2020), p. 012078.
- [7] Minqing Hu and Bing Liu. “Mining and Summarizing Customer Reviews”. In: *Proceedings of the 10th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 2004, pp. 168–177.
- [8] C.J. Hutto and Eric E. Gilbert. “VADER: A Parsimonious Rule-Based Model for Sentiment Analysis of Social Media Text”. In: *Proceedings of the International AAAI Conference on Web and Social Media*. Vol. 8. 1. 2014, pp. 216–225.
- [9] Thorsten Joachims. “Text Categorization with Support Vector Machines: Learning with Many Relevant Features”. In: *European Conference on Machine Learning*. 1998.
- [10] Soo-Min Kim and Eduard H. Hovy. “Extracting Opinions with Conditional Random Fields”. In: *Proceedings of the HLT-NAACL*. 2006.
- [11] Bing Liu. *Sentiment Analysis and Opinion Mining*. Morgan & Claypool Publishers, 2012.
- [12] Yinhan Liu et al. “RoBERTa: A Robustly Optimized BERT Pretraining Approach”. In: *arXiv preprint arXiv:1907.11692* (2019).
- [13] Saif M. Mohammad and Peter D. Turney. “Crowdsourcing a Word-Emotion Association Lexicon”. In: *Computational Intelligence* 29.3 (2013), pp. 436–465.

- [14] Finn Årup Nielsen. “A new ANEW: Evaluation of a word list for sentiment analysis in microblogs”. In: *Proceedings of the 1st Workshop on 'Affect in Text and Speech'*. Vol. 11. 1. 2011, pp. 9–14.
- [15] Bo Pang and Lillian Lee. “Opinion Mining and Sentiment Analysis”. In: *Foundations and Trends in Information Retrieval* 2.1–2 (2008), pp. 1–135.
- [16] Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. “Thumbs up? Sentiment Classification using Machine Learning Techniques”. In: *Proceedings of the ACL-02 Conference on Empirical Methods in Natural Language Processing - Volume 10*. 2002, pp. 79–86.
- [17] Hyeonmin Park and Joonas Lee. “Impact of User Feedback on Sentiment Analysis: An Empirical Study of Movie Reviews”. In: *Journal of Big Data* 5.1 (2018), pp. 1–19.
- [18] Matthew Purver. “Location-Based Sentiment Analysis on Twitter: A Comparison of Methods”. In: *Proceedings of the IEEE International Conference on Social Computing (SocialCom)*. 2011, pp. 192–199.
- [19] Alec Radford et al. *Improving Language Understanding by Generative Pre-Training*. OpenAI blog. 2018.
- [20] Gerard Salton and Christopher Buckley. “Term-weighting Approaches in Automatic Text Retrieval”. In: *Information Processing & Management* 24.5 (1988), pp. 513–523.
- [21] Andranik Tumasjan et al. “Predicting Elections with Twitter: What 140 Characters Reveal About Political Sentiment”. In: *Proceedings of the International AAAI Conference on Web and Social Media*. Vol. 4. 1. 2010.
- [22] Ashish Vaswani et al. “Attention Is All You Need”. In: *Advances in Neural Information Processing Systems*. Vol. 30. 2017.
- [23] Tom Young et al. “Recent Trends in Deep Learning for Natural Language Processing”. In: *IEEE Computational Intelligence Magazine* 13.3 (2018), pp. 55–73.

.10 Code-appendix (with Jupyter notebook Python)

Listing 1: *REQUIRED INSTALLATIONS*

```
1 pip install -r requirements.txt
```

Listing 2: *IMPORT LIBRARIES*

```
1 import pandas as pd
2 import numpy as np
3 import ast
4 import re
5 import os
6 import emoji
7 import warnings
8 from collections import Counter
9 from difflib import SequenceMatcher
10
11 # Visualization
12 import seaborn as sns
13 import matplotlib.pyplot as plt
14
15 # NLP & Text
16 import nltk
17 from nltk.stem import WordNetLemmatizer
18 from nltk.tokenize import word_tokenize
19 from nltk.corpus import stopwords
20
21 # Transformers
22 from transformers import pipeline
23
24 # Progress Bar
25 from tqdm.notebook import tqdm
26
27 # TF-IDF and Word2Vec
28 from sklearn.feature_extraction.text import TfidfVectorizer
29 from gensim.models import Word2Vec
```

Listing 3: DOWNLOAD NLTK RESOURCES

```
1 nltk.download('punkt')
2 nltk.download('wordnet')
3 nltk.download('omw-1.4')
4 nltk.download('stopwords')
5
6 # Activate tqdm for pandas
7 tqdm.pandas()
```

Listing 4: LOAD DATA

```
1 comments = pd.read_csv("comments.csv")
2 posts = pd.read_csv("posts_answers.csv")
3 users = pd.read_csv("users.csv")
4
5 warnings.filterwarnings('ignore')
```

Listing 5: INITIAL EXPLORATION (1/2)

```
1 print("--- Initial Exploration: Comments Data ---")
2 print("\nHead of Comments:")
3 display(comments.head())
4
5 print("\nMissing Values:")
6 display(comments.isnull().sum())
7
8 print("\nDuplicate Rows:", comments.duplicated().sum())
9 print("Unique Users:", comments['user_id'].nunique())
10 print("Unique Posts:", comments['post_id'].nunique())
11
12 print("\n\n--- Initial Exploration: Posts_answers Data ---")
13 display(posts.head())
14 display(posts.isnull().sum())
15 print("Duplicate Rows:", posts.duplicated().sum())
16
17 plt.figure(figsize=(10, 5))
18 sns.histplot(posts['score'].dropna(), bins=50, kde=True, color='purple')
19 plt.title('Distribution of Post Scores')
20 plt.xlabel('Score')
21 plt.ylabel('Frequency')
22 plt.tight_layout()
23 plt.show()
```

Listing 6: *INITIAL EXPLORATION (2/2)*

```

1 print("\n\n--- Initial Exploration: Users Data ---")
2 display(users.head())
3 display(users.isnull().sum())
4 print("Duplicate Rows:", users.duplicated().sum())
5
6 plt.figure(figsize=(10, 5))
7 sns.histplot(users['reputation'].dropna(), bins=50, kde=True, color='↵
   green')
8 plt.title('Distribution of User Reputation')
9 plt.xlabel('Reputation')
10 plt.ylabel('Frequency')
11 plt.tight_layout()
12 plt.show()
13
14 print("\nTop 10 User Locations:")
15 display(users['location'].value_counts().head(10))

```

Listing 7: *TEXT CLEANING (1/3)*

```

1 # Regex Patterns
2 URL_PATTERN = re.compile(r"http\S+|www.\S+")
3 MEANINGFUL_TEXT_PATTERN = re.compile(r"\w{3,}")
4 EMOJI_PATTERN = re.compile("[\U0001F600-\U0001F64F"
5                             "\U0001F300-\U0001F5FF"
6                             "\U0001F680-\U0001F6FF"
7                             "\U0001F1E0-\U0001F1FF]+", flags=re.UNICODE)

```

Listing 8: *TEXT CLEANING (2/3)*

```

1 def has_meaningful_text(text):
2     text = URL_PATTERN.sub("", text)
3     return bool(MEANINGFUL_TEXT_PATTERN.search(text))
4
5 def remove_emojis(text):
6     return EMOJI_PATTERN.sub("", text)
7
8 def clean_text(text, lower=True, strip=True, remove_urls=True, ↵
   remove_emojis_flag=True):

```



```

9     if pd.isna(text):
10         return ""
11     if lower:
12         text = text.lower()
13     if remove_urls:
14         text = URL_PATTERN.sub("", text)
15     if remove_emojis_flag:
16         text = remove_emojis(text)
17     if strip:
18         text = text.strip()
19     return text

```

Listing 9: *TEXT CLEANING (3/3)*

```

1 def clean_comments(df, text_col="text", min_words=3, drop_duplicates=True↵
    , remove_emojis_flag=True, lower=True):
2     df = df.copy()
3     df = df[df[text_col].notna()]
4     df = df[df[text_col].str.strip() != ""]
5     df[text_col] = df[text_col].apply(lambda x: clean_text(x, lower=lower↵
        , remove_emojis_flag=remove_emojis_flag))
6     if drop_duplicates:
7         df = df.drop_duplicates(subset=[text_col])
8     df = df[df[text_col].str.split().str.len() >= min_words]
9     df = df[df[text_col].apply(has_meaningful_text)]
10    return df.reset_index(drop=True)
11
12 # Apply cleaning
13 comments = clean_comments(comments, min_words=3, remove_emojis_flag=True)

```

Listing 10: *BASIC STATS VISUALS*

```

1 print("Number of comments:", len(comments))
2 print("Unique users:", comments["user_id"].nunique())
3 print("Missing values in text:", comments["text"].isna().sum())
4
5 # Word Count Distribution
6 comments["word_count"] = comments["text"].apply(lambda x: len(str(x).↵
    split()))
7
8 plt.figure(figsize=(10, 6))
9 sns.histplot(comments["word_count"], bins=30, kde=True, color="skyblue")

```

```

10 plt.title("Distribution of Comment Lengths (Word Count)")
11 plt.xlabel("Word Count")
12 plt.ylabel("Frequency")
13 plt.tight_layout()
14 plt.show()

```

Listing 11: WORD FREQUENCY ANALYSIS

```

1 stop_words = set(stopwords.words("english"))
2
3 def tokenize(text):
4     return [word.lower() for word in str(text).split() if word.lower() ↵
5         not in stop_words and word.isalpha()]
6
7 all_words = sum(comments["text"].apply(tokenize).tolist(), [])
8 top_words = Counter(all_words).most_common(20)
9
10 # Plot Top Words
11 words, counts = zip(*top_words)
12 plt.figure(figsize=(10, 6))
13 ax = sns.barplot(x=list(counts), y=list(words), palette="viridis")
14 plt.title("Top 20 Most Frequent Words in Comments")
15 plt.xlabel("Frequency")
16 plt.ylabel("Word")
17
18 for i, count in enumerate(counts):
19     ax.text(count + 0.5, i, str(count), va='center', fontsize=10)
20 plt.tight_layout()
21 plt.show()

```

Listing 12: FEATURE EXTRACTION

```

1 # Ensure text column is string type
2 comments['text'] = comments['text'].fillna('').astype(str)
3
4 # TF-IDF Vectorization
5 tfidf_vectorizer = TfidfVectorizer(max_features=28000, stop_words='↵
6     english', ngram_range=(1, 1))
7 tfidf_matrix = tfidf_vectorizer.fit_transform(comments['text'])
8
9 print(f"TF-IDF Matrix Shape: {tfidf_matrix.shape}")
10 print(f"Unique Features Extracted: {len(tfidf_vectorizer.↵

```

```
get_feature_names_out())})")
```

Listing 13: *TOKENIZATION*

```
1 lemmatizer = WordNetLemmatizer()
2
3 def tokenize_and_lemmatize(text):
4     tokens = word_tokenize(text.lower())
5     return [lemmatizer.lemmatize(word) for word in tokens if word.isalpha()
6             () and word not in stop_words]
7
8 comments['lemmatized_tokens'] = comments['text'].progress_apply(
9     tokenize_and_lemmatize)
10
11 # Generate bigrams
12 bigrams_list = []
13 for tokens in comments['lemmatized_tokens']:
14     bigrams_list.extend(list(nltk.bigrams(tokens)))
15
16 bigram_counts = Counter(bigrams_list)
17 top_bigrams = bigram_counts.most_common(20)
```

Listing 14: *BIGRAM PLOT*

```
1 # Plot Top Bigrams
2 plt.figure(figsize=(10, 7))
3 bigram_words = [" ".join(bg) for bg, count in top_bigrams]
4 bigram_counts_val = [count for bg, count in top_bigrams]
5
6 ax = sns.barplot(x=bigram_counts_val, y=bigram_words, palette="coolwarm")
7 plt.title("Top 20 Most Frequent Bigrams in Comments")
8 plt.xlabel("Frequency")
9 plt.ylabel("Bigram")
10
11 for i, count in enumerate(bigram_counts_val):
12     ax.text(count + 0.5, i, str(count), va='center', fontsize=10)
13 plt.tight_layout()
14 plt.show()
```

Listing 15: *GoEmotions FROM HUGGING FACE*

```

1 import time
2 # Define batch inference function
3 def batch_infer(df, batch_size=16, model_name="SamLowe/roberta-base-↵
  go_emotions"):
4     classifier = pipeline("text-classification", model=model_name, top_k=↵
      None)
5     results = []
6     texts = df["text"].tolist()
7
8     for i in tqdm(range(0, len(texts), batch_size), desc="Batch Inference↵
      "):
9         batch = [str(t)[:512] for t in texts[i:i + batch_size]]
10        batch_result = classifier(batch)
11        results.extend(batch_result)
12
13    return results
14
15 # Run the function and time it
16 print("Running emotion detection with GoEmotions model...")
17
18 start_time = time.time()
19 comments["emotions"] = batch_infer(comments)
20 comments.to_csv("comments_with_emotions.csv", index=False)
21 print("----> Saved emotions per comment to comments_with_emotions.csv")
22
23 end_time = time.time()
24
25 print(f"----> Emotion inference completed in {end_time - start_time:.2f} ↵
      seconds.")

```

Listing 16: *EMOTION DATA PROCESSING (1/3)*

```

1 # Load Comments with Emotions
2 comments = pd.read_csv("comments_with_emotions.csv")
3 comments["emotions"] = comments["emotions"].apply(ast.literal_eval)
4
5 print(f" ----> Loaded {len(comments)} comments with parsed emotions.")

```

Listing 17: *EMOTION DATA PROCESSING (2/3)*

```

1 # Extract Emotion Labels Above a Score Threshold
2 def extract_labels(e_list, threshold=0.3):
3     """
4     Filters emotion predictions above a given score threshold.
5     """
6     return [e["label"] for e in e_list if e["score"] > threshold]
7
8 comments["emotion_labels"] = comments["emotions"].apply(extract_labels)
9 print(" ---> Extracted emotion labels for each comment.")

```

Listing 18: *EMOTION DATA PROCESSING (3/3)*

```

1 # Explode Multi-Emotion Comments into Rows
2 comments_exp = comments.explode("emotion_labels")
3 comments_exp = comments_exp[comments_exp["emotion_labels"].notna()]
4 print(f" ---> Exploded to {len(comments_exp)} rows (one row per emotion ↔
   per comment).")
5
6 # Ensure ID Columns Are Numeric for Merging
7 comments_exp["user_id"] = pd.to_numeric(comments_exp["user_id"], errors="↔
   coerce")
8 comments_exp["post_id"] = pd.to_numeric(comments_exp["post_id"], errors="↔
   coerce")
9 print(" ---> Converted user_id and post_id to numeric for merging.")

```

Listing 19: *MERGE COMMENTS WITH USERS AND POSTS*

```

1 # Merge Comments with Users and Posts
2 def merge_all(comments_df, users_df, posts_df):
3     """
4     Merges comment data with user and post metadata.
5     """
6     merged = comments_df.merge(users_df, left_on="user_id", right_on="id"↔
   , how="left", suffixes=('', '_user'))
7     merged = merged.merge(posts_df, left_on="post_id", right_on="id", how↔
   ="left", suffixes=('', '_post'))
8     return merged
9
10 df = merge_all(comments_exp, users, posts)
11 print(f" ---> Merged with users and posts. Final shape: {df.shape}")

```

Listing 20: *BIN NUMERIC COLUMNS INTO CATEGORIES*

```

1 def bin_column(df, col, bins, labels, new_col):
2     """
3     Adds a categorical column based on binning of a numeric column.
4     """
5     df[new_col] = pd.cut(df[col], bins=bins, labels=labels)
6     return df

```

Listing 21: *REPUTATION GROUPING*

```

1 df = bin_column(
2     df, "reputation",
3     bins=[-1, 100, 1000, 5000, 20000, float("inf")],
4     labels=["Very Low", "Low", "Medium", "High", "Very High"],
5     new_col="rep_group"
6 )
7
8 # Post Score Grouping
9 df = bin_column(
10     df, "score_post",
11     bins=[-10, 0, 5, 20, 100, float("inf")],
12     labels=["Negative", "Low", "Moderate", "High", "Very High"],
13     new_col="post_score_group"
14 )
15
16 # Comment Score Grouping
17 df = bin_column(
18     df, "score",
19     bins=[-5, 0, 5, 20, 100, float("inf")],
20     labels=["Negative", "Low", "Moderate", "High", "Very High"],
21     new_col="comment_score_group"
22 )
23
24 print(" ---> Created reputation and post score groups.")

```

Listing 22: *DISPLAY SAMPLE OUTPUT*

```

1 # Display Sample Output
2 display(df[[
3     "text", "emotion_labels",
4     "reputation", "rep_group",
5     "score_post", "post_score_group"

```

```
6 ]] .sample(5))
```

Listing 23: *TOP-SCORING COMMENTS BY EMOTION*

```
1 top_comments_per_emotion = (
2     df[df["emotion_labels"].notna()]
3     .groupby("emotion_labels")
4     .apply(lambda g: g.sort_values("score", ascending=False).head(1))
5 )[[ "text", "score", "display_name", "emotion_labels" ]]
6
7 print("----> Top-Scoring Comments by Emotion:")
8 display(top_comments_per_emotion.reset_index(drop=True))
```

Listing 24: *MOST CONFIDENT EMOTION PREDICTIONS*

```
1 # Most Confident Emotion Predictions
2 df["emotion_confidence"] = df["emotions"].apply(lambda x: max([e["score"] for e in x]) if x else 0)
3 df_top_confident = df.sort_values("emotion_confidence", ascending=False).head(10)
4
5 print("----> Most Confident Emotion Predictions:")
6 display(df_top_confident[[ "text", "emotion_labels", "emotion_confidence", "score", "display_name" ]])
```

Listing 25: *BAR PLOT - FREQUENCY OF ALL 27 EMOTIONS*

```
1 from collections import Counter
2 all_emotions = sum(comments["emotion_labels"].tolist(), [])
3 emotion_counts = pd.Series(Counter(all_emotions)).sort_values(ascending=False)
4
5 plt.figure(figsize=(12, 8))
6 ax = sns.barplot(x=emotion_counts.values, y=emotion_counts.index, palette="coolwarm")
7 plt.title("Frequency of All 27 Emotions")
8 plt.xlabel("Count")
9 plt.ylabel("Emotion")
10 for i, (value, label) in enumerate(zip(emotion_counts.values, emotion_counts.index)):
11     ax.text(value + 2, i, str(value), va='center', fontsize=9)
12 plt.tight_layout()
```

```
13 plt.show()
```

Listing 26: *BAR PLOT - FREQUENCY OF ALL 27 EMOTIONS EXCLUDING NEUTRAL*

```
1 # Bar Plot: Emotions (Excluding Neutral)
2 non_neutral_emotions = [e for e in all_emotions if e != "neutral"]
3 non_neutral_counts = pd.Series(Counter(non_neutral_emotions).sort_values(↵
    (ascending=False)
4
5 plt.figure(figsize=(12, 8))
6 ax = sns.barplot(x=non_neutral_counts.values, y=non_neutral_counts.index,↵
    palette="coolwarm")
7 plt.title("Frequency of Emotions (Excluding Neutral)")
8 plt.xlabel("Count")
9 plt.ylabel("Emotion")
10 for i, (value, label) in enumerate(zip(non_neutral_counts.values, ↵
    non_neutral_counts.index)):
11     ax.text(value + 2, i, str(value), va='center', fontsize=9)
12 plt.tight_layout()
13 plt.show()
```

Listing 27: *EMOTION DISTRIBUTION BY COMMENT SCORE*

```
1 # Plot: Emotion Distribution by Comment Score
2 plt.figure(figsize=(12, 6))
3 sns.countplot(data=df, x="comment_score_group", hue="emotion_labels", ↵
    palette="tab20")
4 plt.title("Emotion Distribution by Comment Score (Popularity)")
5 plt.xlabel("Comment Score Group")
6 plt.ylabel("Count")
7 plt.legend(loc="upper right", bbox_to_anchor=(1.25, 1))
8 plt.tight_layout()
9 plt.show()
```


Listing 28: *EMOTION DISTRIBUTION BY COMMENT SCORE EXCLUDING NEUTRAL*

```

1 df_no_neutral = df[df["emotion_labels"] != "neutral"]
2
3 plt.figure(figsize=(12, 6))
4 sns.countplot(data=df_no_neutral, x="comment_score_group", hue="↔
    emotion_labels", palette="tab20")
5 plt.title("Emotion Distribution by Comment Score (Without Neutral)")
6 plt.xlabel("Comment Score Group")
7 plt.ylabel("Count")
8 plt.legend(loc="upper right", bbox_to_anchor=(1.25, 1))
9 plt.tight_layout()
10 plt.show()

```

Listing 29: *RE-BIN POST SCORES FOR PLOTTING*

```

1 # Re-bin Post Scores for Plotting
2 df["post_score_group"] = pd.cut(
3     df["score_post"],
4     bins=[-10, 0, 5, 20, 100, float("inf")],
5     labels=["Negative", "Low", "Medium", "High", "Very High"]
6 )

```

Listing 30: *EMOTION DISTRIBUTION PLOT BY POST SCORE*

```

1 plt.figure(figsize=(12, 6))
2 sns.countplot(data=df, x="post_score_group", hue="emotion_labels", ↔
    palette="tab20")
3 plt.title("Emotion Distribution by Post Score")
4 plt.xlabel("Post Score Group")
5 plt.ylabel("Emotion Count")
6 plt.legend(loc="upper right", bbox_to_anchor=(1.25, 1))
7 plt.tight_layout()
8 plt.show()

```

Listing 31: *EMOTION DISTRIBUTION PLOT BY POST SCORE EXCLUDING NEUTRAL*

```

1 # Exclude Neutral
2
3 plt.figure(figsize=(12, 6))
4 sns.countplot(data=df_no_neutral, x="post_score_group", hue="↵
    emotion_labels", palette="tab20")
5 plt.title("Emotion Distribution by Post Score (Without Neutral)")
6 plt.xlabel("Post Score Group")
7 plt.ylabel("Count")
8 plt.legend(loc="upper right", bbox_to_anchor=(1.25, 1))
9 plt.tight_layout()
10 plt.show()

```

Listing 32: *GEOGRAPHICAL EMOTION ANALYSIS*

```

1 df_clean = df.copy()
2 df_clean["location"] = df_clean["location"].replace(["United States", "↵
    USA"], "United States")
3
4 top_locations_clean = df_clean["location"].value_counts().nlargest(10).↵
    index
5 df_top_loc_clean = df_clean[df_clean["location"].isin(top_locations_clean↵
    )]
6
7 plt.figure(figsize=(12, 6))
8 sns.countplot(data=df_top_loc_clean, x="location", order=↵
    top_locations_clean)
9 plt.title("Top 10 User Locations by Number of Comments")
10 plt.xlabel("Location")
11 plt.ylabel("Number of Comments")
12 plt.xticks(rotation=45, ha="right")
13 plt.tight_layout()
14 plt.show()

```

Listing 33: *EMOTIONS HEATMAP BY LOCATION*

```

1 df["location"] = df["location"].replace({"USA": "United States", "U.S.": ↵
    "United States"})
2 top_locations = df["location"].value_counts().nlargest(10).index
3 df_top_loc = df[df["location"].isin(top_locations)]

```

```

4
5 heatmap_data = df_top_loc.groupby(["emotion_labels", "location"]).size().↵
    unstack(fill_value=0)
6
7 plt.figure(figsize=(14, 8))
8 sns.heatmap(
9     heatmap_data,
10    annot=True,
11    fmt="d",
12    cmap="YlGnBu",
13    linewidths=0.5,
14    cbar_kws={"label": "Count"}
15 )
16 plt.title("Heatmap: Emotions by User Location (Top 10 Locations)")
17 plt.xlabel("Location")
18 plt.ylabel("Emotion")
19 plt.xticks(rotation=45, ha="right")
20 plt.tight_layout()
21 plt.show()

```

Listing 34: *EMOTIONS HEATMAP BY LOCATION EXCLUDING NEUTRAL*

```

1 # Heatmap: Emotions by Location (Excluding Neutral)
2 heatmap_no_neutral = df_top_loc_clean[df_top_loc_clean["emotion_labels"] ↵
    != "neutral"] \
3     .groupby(["emotion_labels", "location"]).size().unstack(fill_value=0)
4
5 plt.figure(figsize=(14, 8))
6 sns.heatmap(
7     heatmap_no_neutral,
8     annot=True,
9     fmt="d",
10    cmap="YlGnBu",
11    linewidths=0.5,
12    cbar_kws={"label": "Count"}
13 )
14 plt.title("Heatmap: Emotions by User Location (Excluding Neutral)")
15 plt.xlabel("Location")
16 plt.ylabel("Emotion")
17 plt.xticks(rotation=45, ha="right")
18 plt.tight_layout()
19 plt.show()

```

Listing 35: *UPVOTES BOX-PLOT BY EMOTIONS*

```
1 # Box Plot: Upvotes by Emotion
2 plt.figure(figsize=(14, 8))
3 sns.boxplot(
4     data=df,
5     y="emotion_labels",
6     x="up_votes",
7     palette="coolwarm",
8     showfliers=False,
9     width=0.6
10 )
11 plt.title("Upvotes by Emotion", fontsize=14, pad=15)
12 plt.xlabel("Upvotes", fontsize=12)
13 plt.ylabel("Emotion", fontsize=12)
14 plt.ticklabel_format(style="plain", axis="x")
15 sns.despine(trim=True)
16 plt.grid(axis="x", linestyle="--", alpha=0.4)
17 plt.tight_layout()
18 plt.show()
```

Listing 36: *DOWNVOTES BOX-PLOT BY EMOTIONS*

```
1 # Box Plot: Downvotes by Emotion
2 plt.figure(figsize=(14, 8))
3 sns.boxplot(
4     data=df,
5     y="emotion_labels",
6     x="down_votes",
7     palette="coolwarm",
8     showfliers=False
9 )
10 plt.title("Downvotes by Emotion", fontsize=14, pad=15)
11 plt.xlabel("Downvotes", fontsize=12)
12 plt.ylabel("Emotion", fontsize=12)
13 plt.tight_layout()
14 plt.show()
```