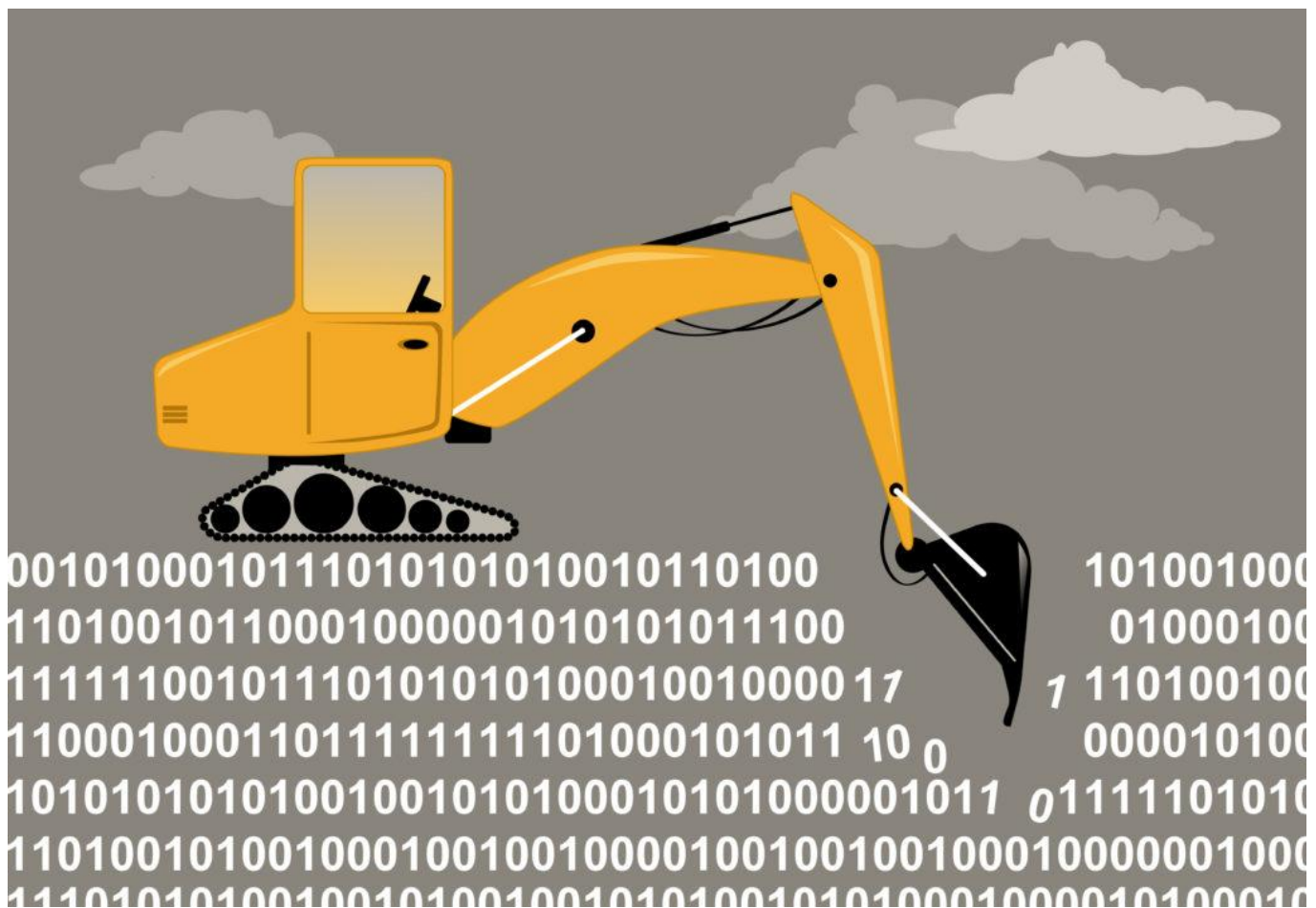


Μάθημα: «Αποθήκες Δεδομένων και Εξόρυξη Γνώσης (6ο εξ.)»

Ομάδα εργασίας: (Π15128, Σκάρος Γεώργιος)

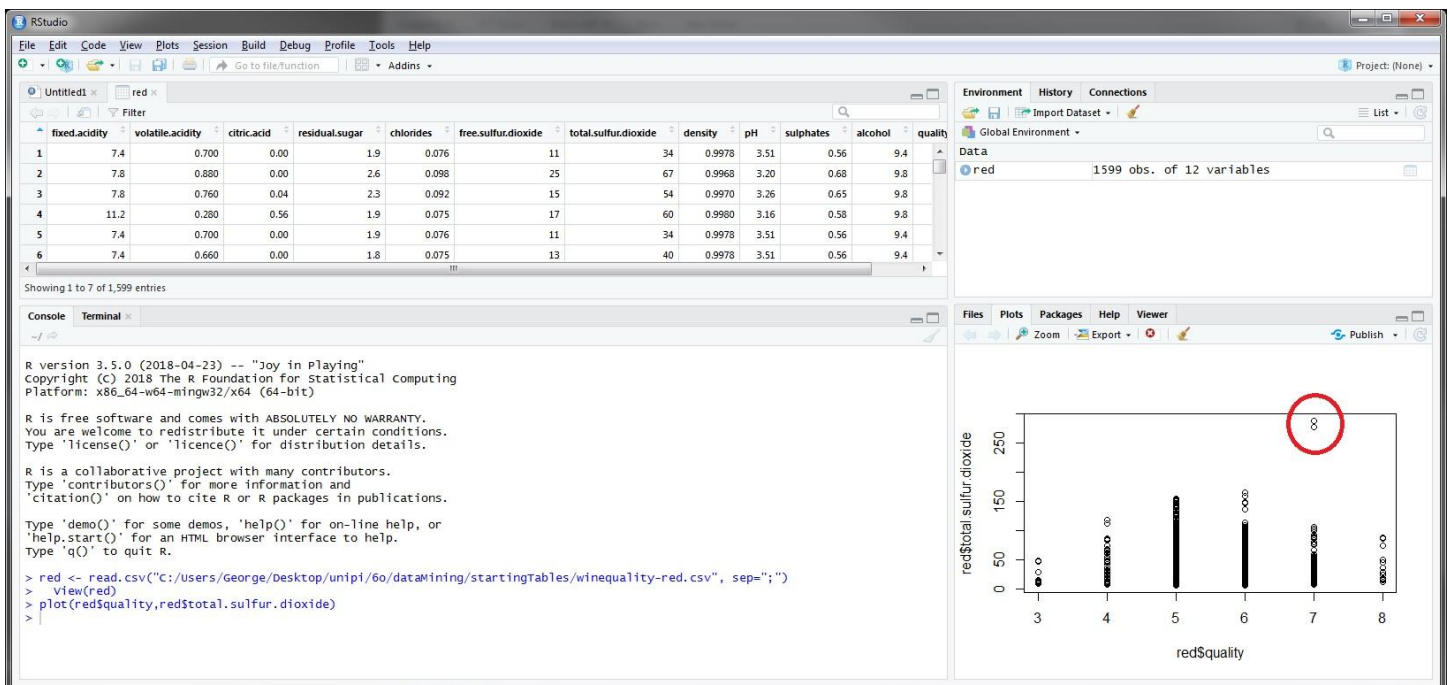


1^ο βήμα

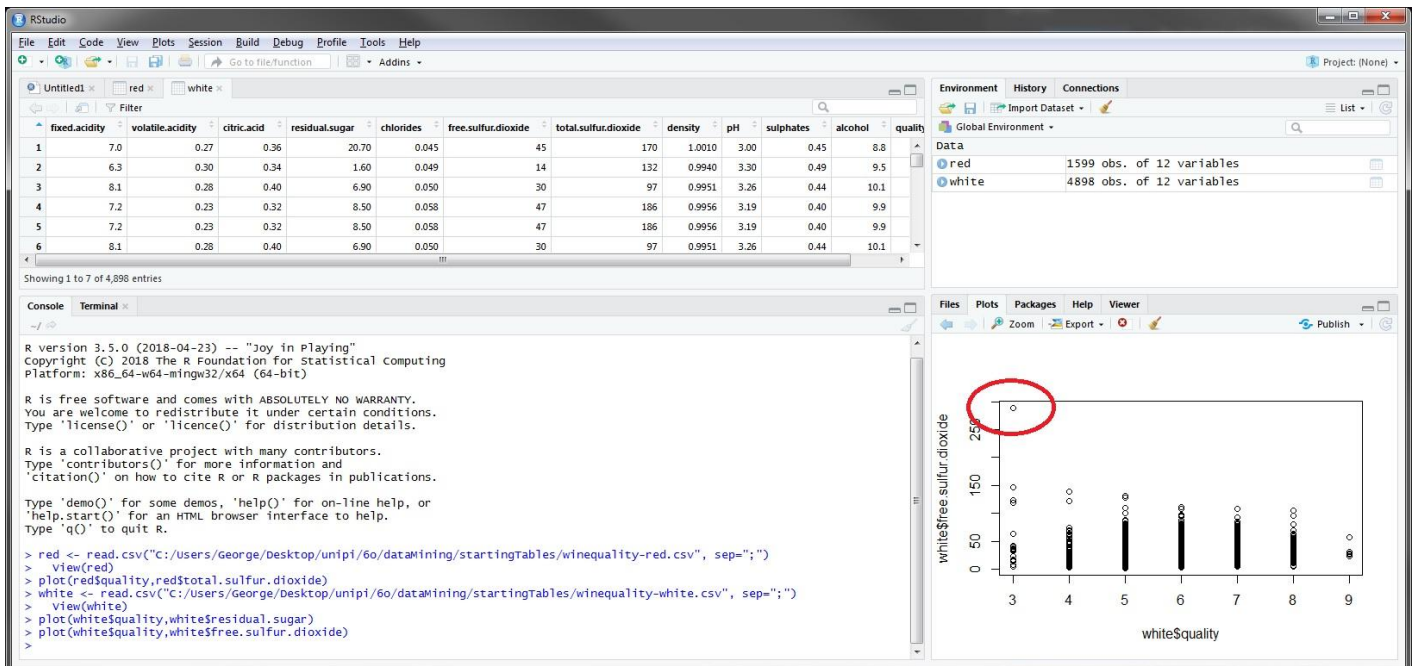
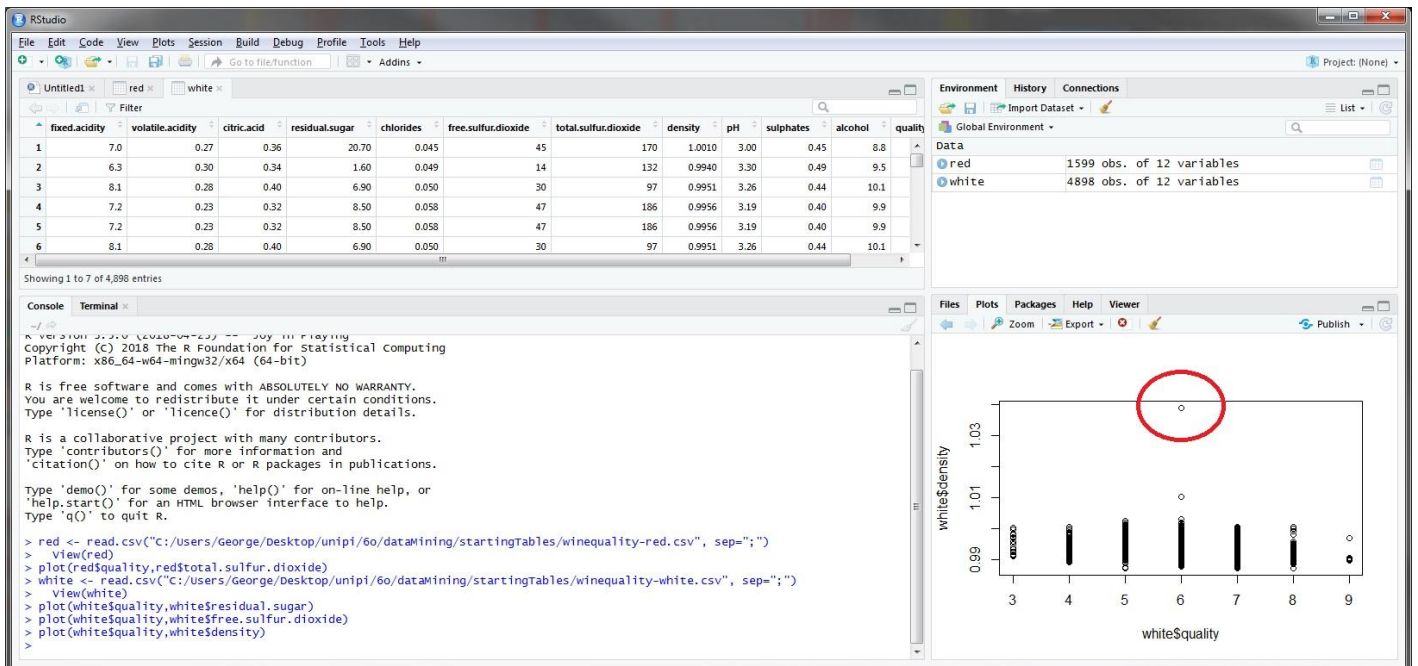
Δημιούργησα δύο πίνακες σε μια βάση δεδομένων υλοποιημένη σε PostgreSQL. Τα queries για τα create table υπάρχουν στον φάκελο createTables που θα βρείτε μέσα στο zip.

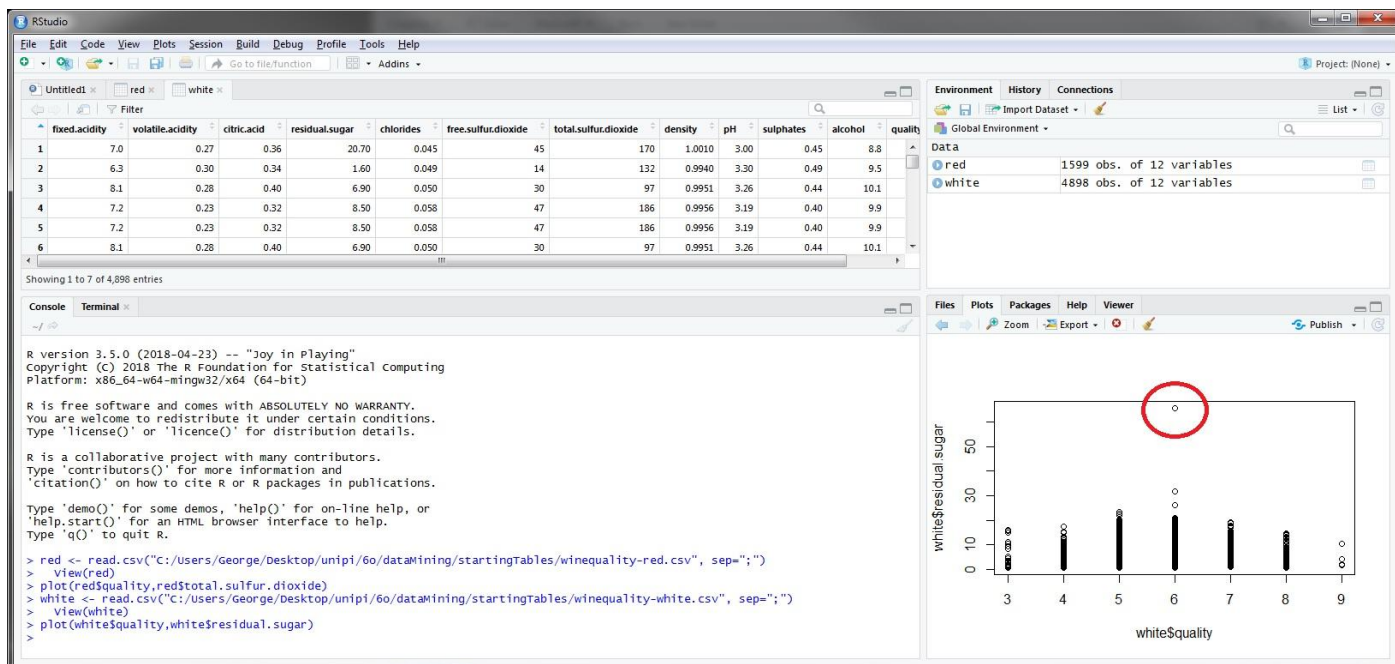
2^ο βήμα

- α) Αρχικά, φόρτωσα τα data στο R Studio. Έπειτα, με την εντολή plot βρήκα στον κάθε πίνακα κάποιες εντολές οι οποίες θεώρησα πως είχαν λανθασμένες τιμές κρίνοντας από το πόσο μεγάλες διαφορές είχαν από τις υπόλοιπες τιμές στις αντίστοιχες στήλες. Πιο συγκεκριμένα, στον πίνακα red :

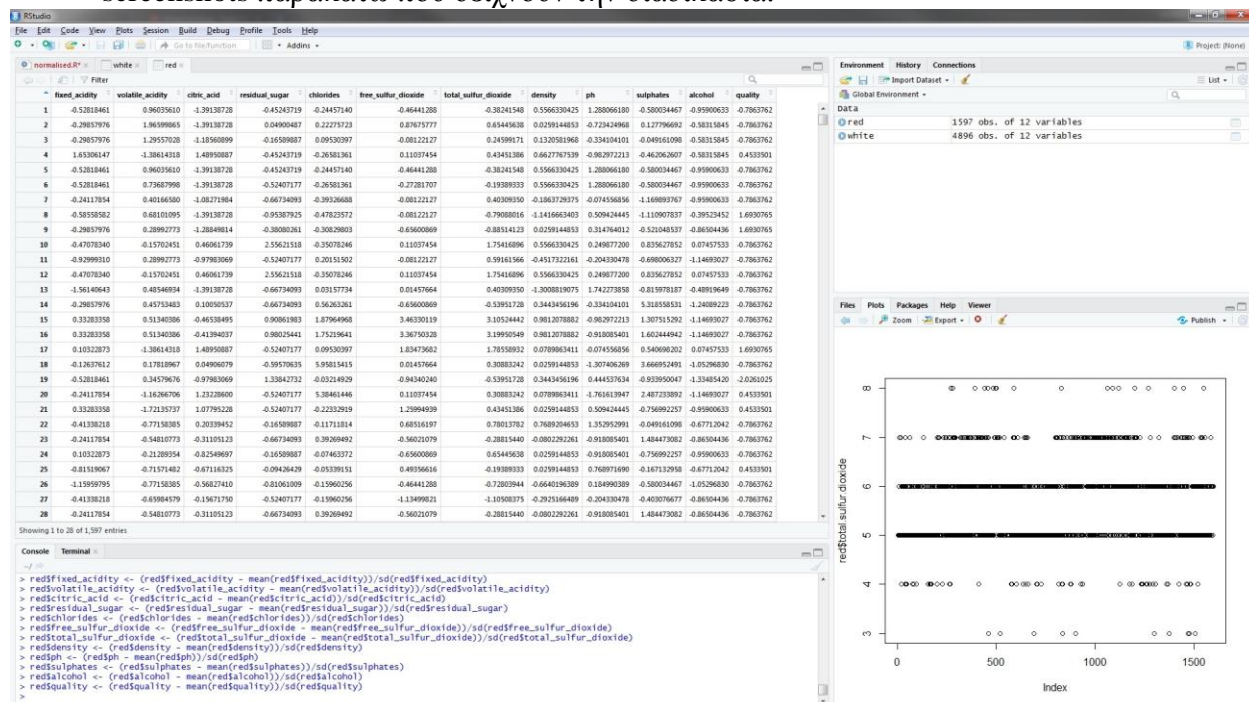


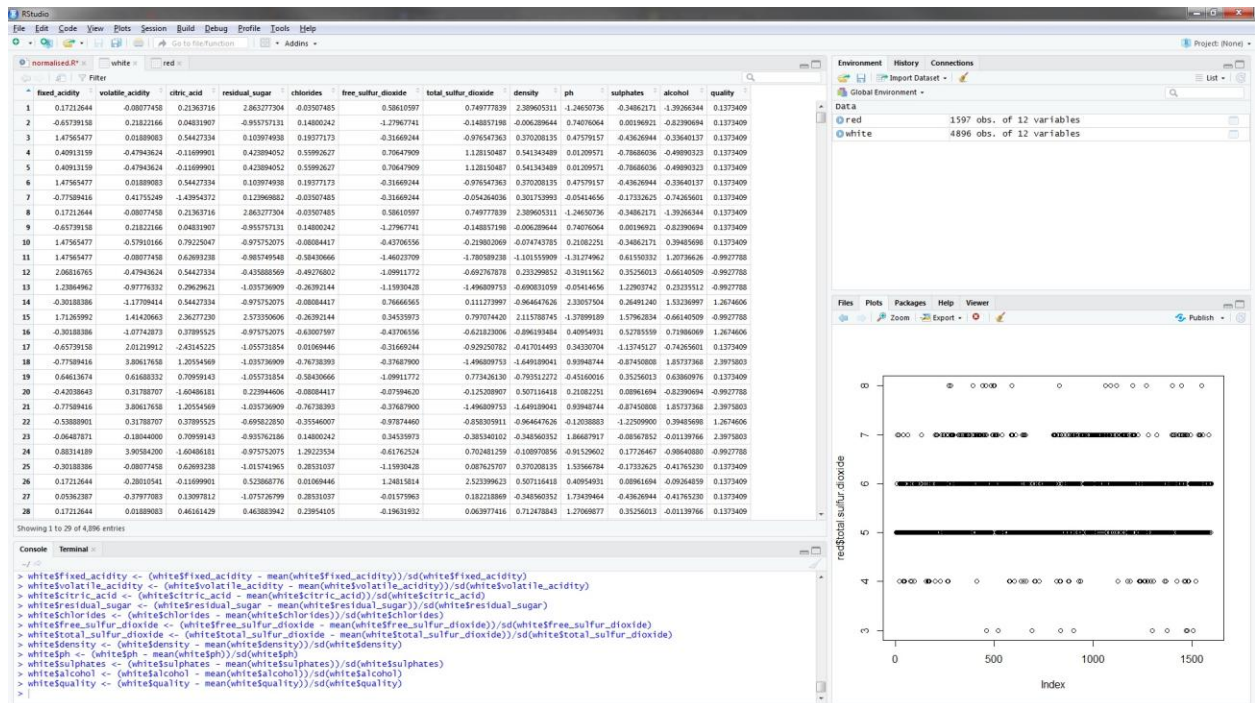
Και για τον white :





b) Στην συνέχεια, κανονικοποίησα τα δεδομένα. Αυτό το έκανα ουσιαστικά βρίσκοντας την μέση τιμή και την τυπική απόκλιση βάζοντας την μέση τιμή στο 0 και το έκανα με τον κώδικα που θα βρείτε στον φάκελο NormaliseFunc. Έχω επισυνάψει και μερικά screenshots παρακάτω που δείχνουν την διαδικασία:





c) Σε αυτό το βήμα χρησιμοποίησα την εντολή cor ανάμεσα στην στήλη quality, που είναι άλλωστε και η στήλη στόχος, και σε κάθε άλλη στήλη στον πίνακα και πήρα τα παρακάτω αποτελέσματα. Γνωρίζοντας ότι όσο πιο κοντά είναι το correlation στο 1 ή στο -1 τόσο πιο γραμμικώς εξαρτημένα είναι τα δύο μεγέθη. Αν είναι θετικό σημαίνει ότι είναι ανάλογα και αν είναι αρνητικό αντιστρόφως ανάλογα. Συνεπώς μπορούμε να ποιες στήλες είναι πιο κοντά στο 0 δεν συνεισφέρουν πολύ για την πρόβλεψη του quality οπότε θα μπορούσα να τα παραλείψω σε επόμενα βήματα.

