

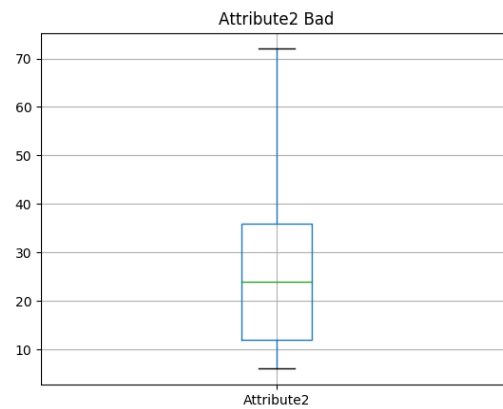
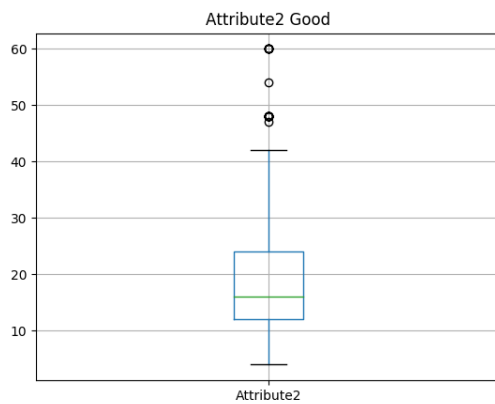
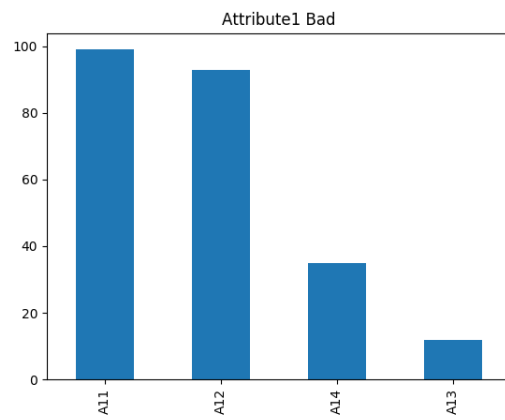
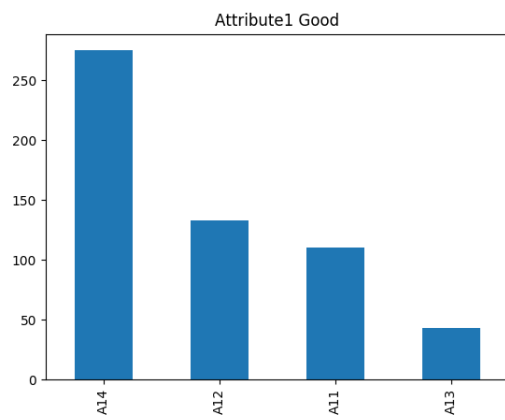
# ΤΕΧΝΙΚΕΣ ΕΞΟΡΥΞΗΣ ΔΕΔΟΜΕΝΩΝ ΕΡΓΑΣΙΑ 2η

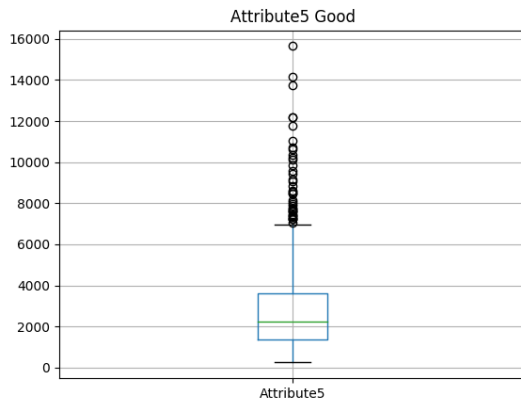
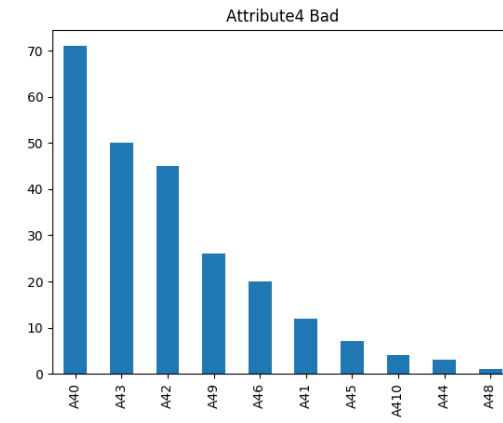
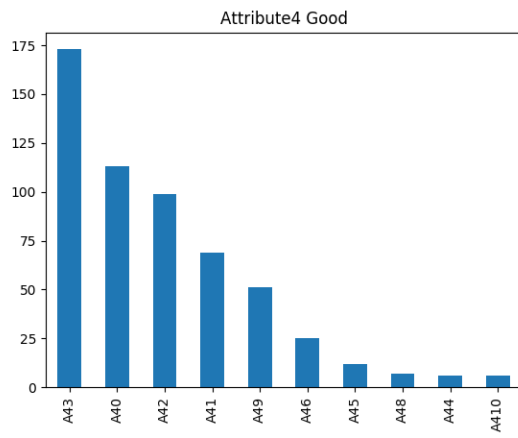
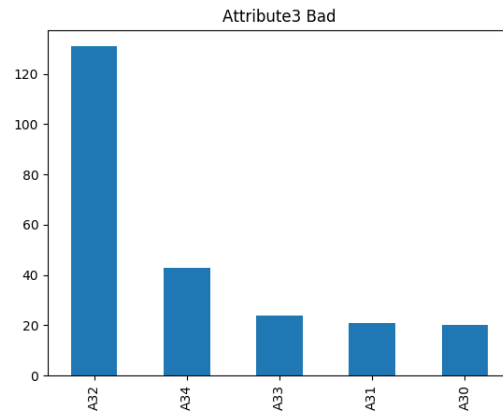
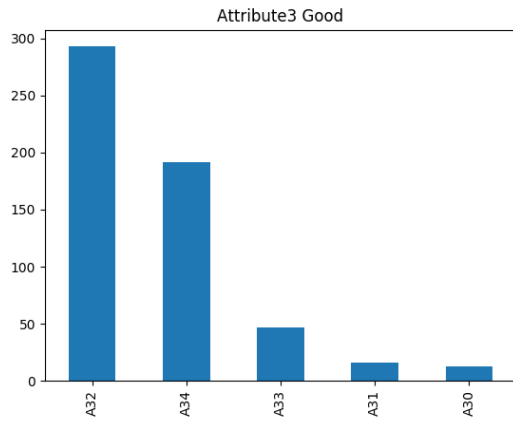
Συμμετέχοντες:

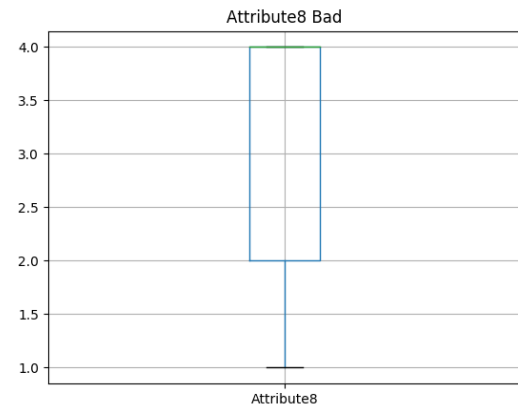
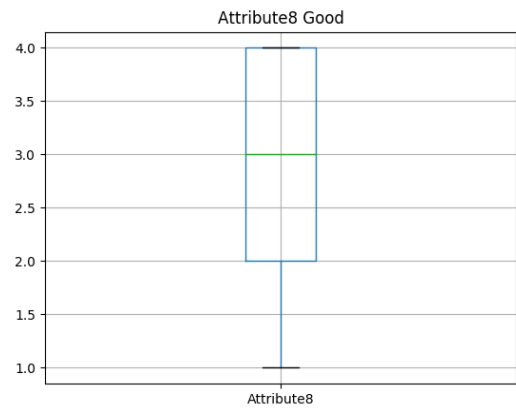
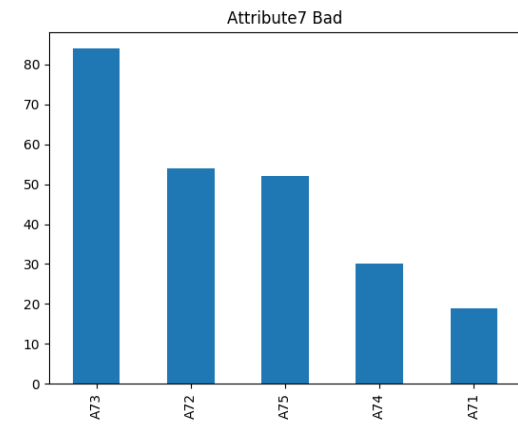
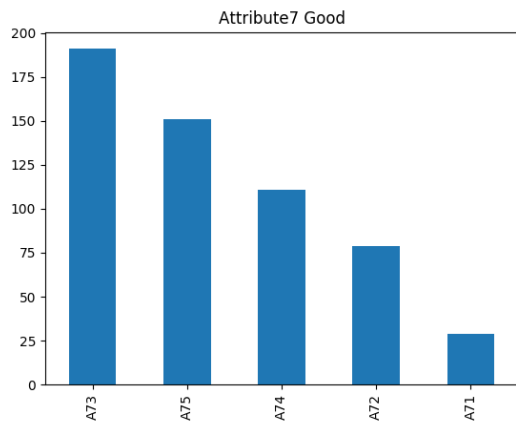
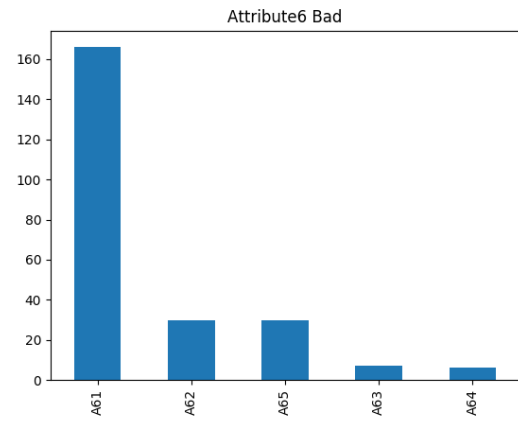
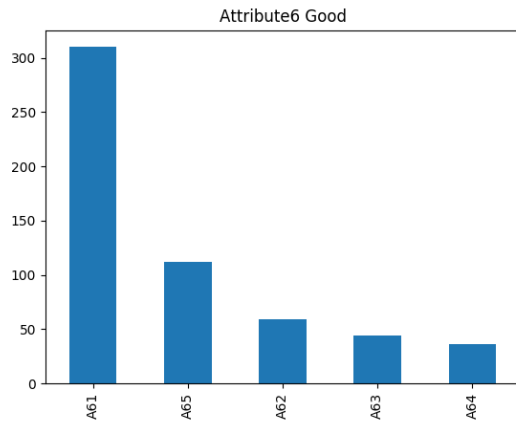
Βαγγέλης Τσιατούρας (1115201200185)

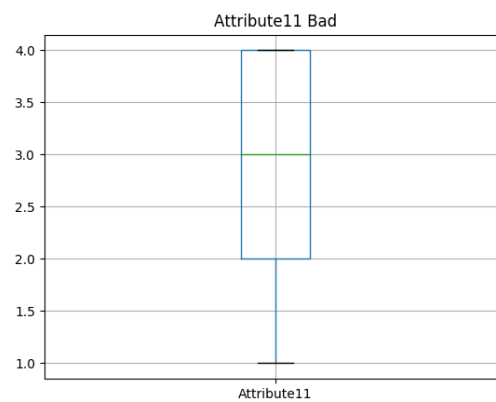
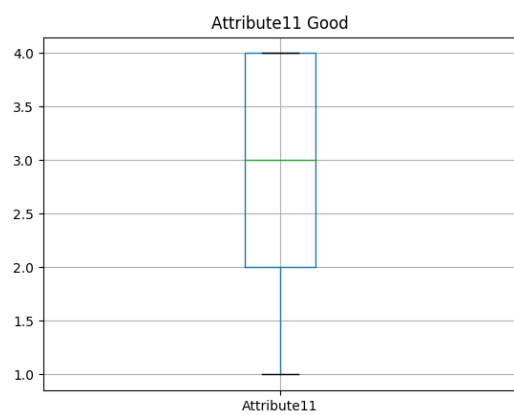
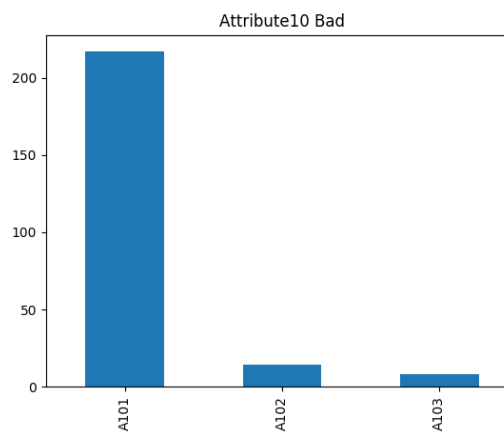
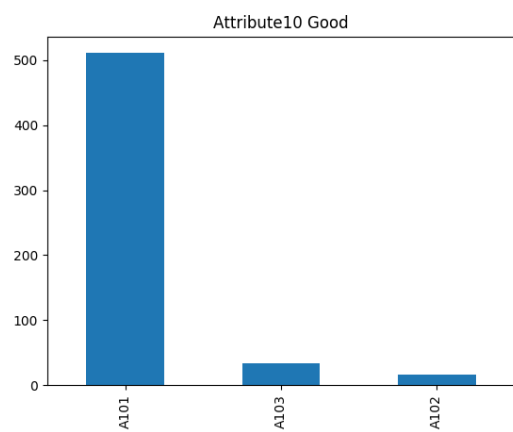
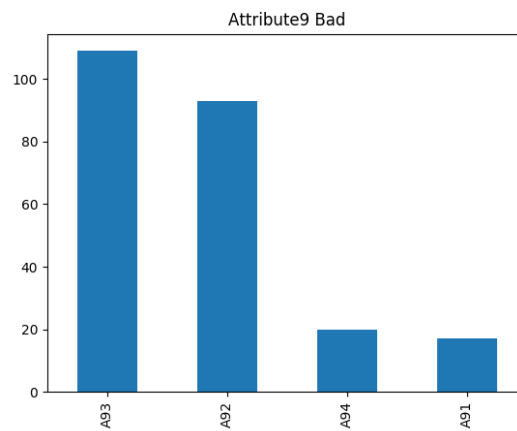
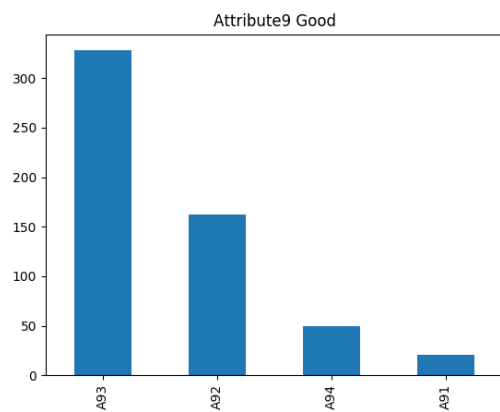
Νίκος Σοφράς (1115201200168)

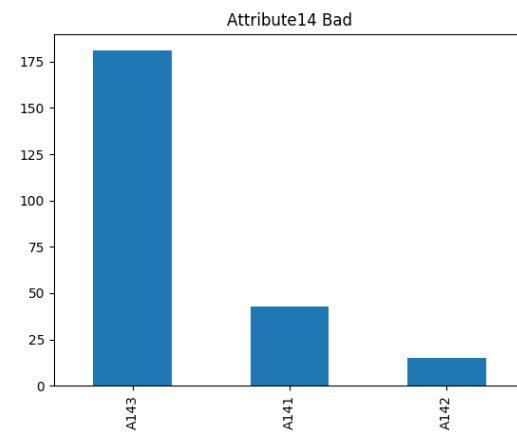
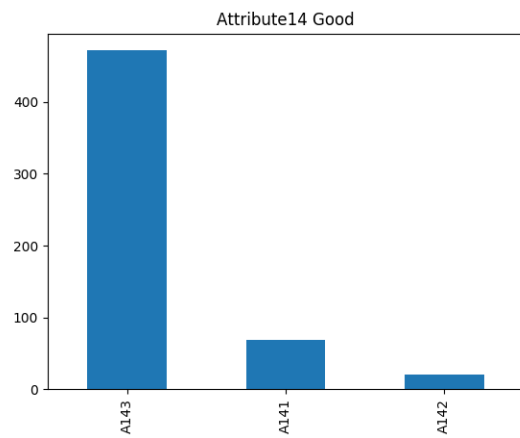
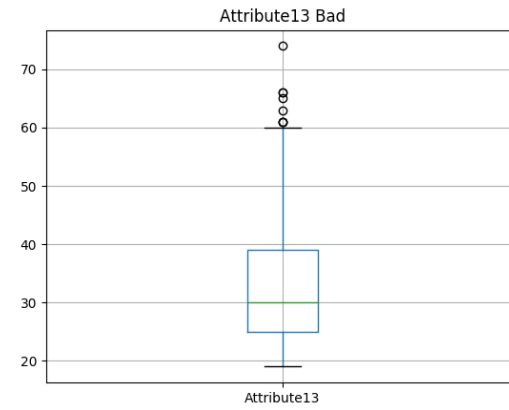
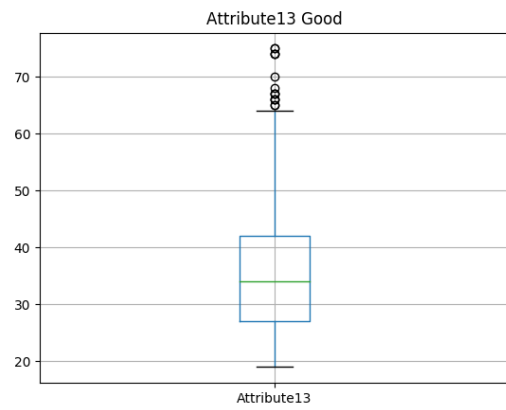
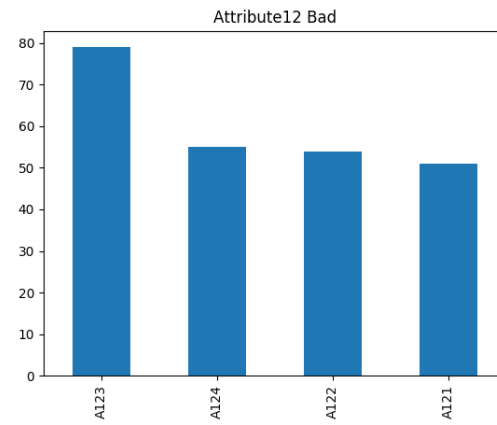
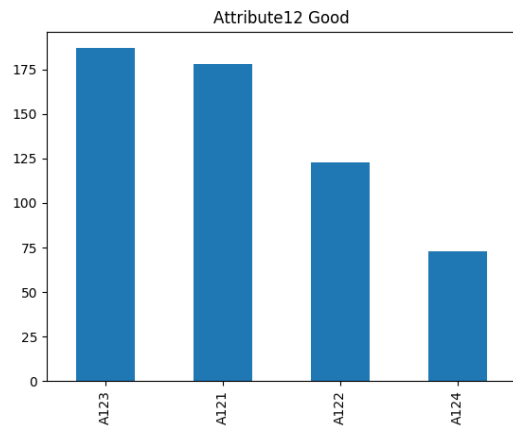
# Οπτικοποίηση Δεδομένων

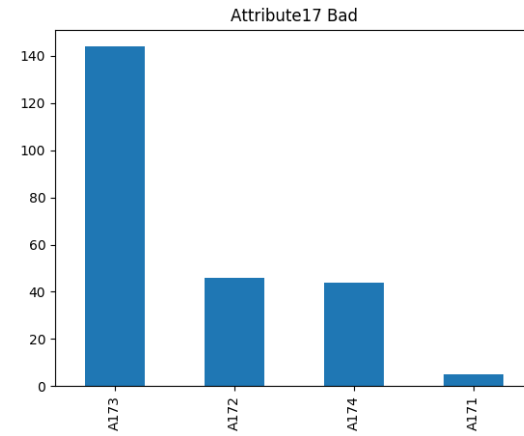
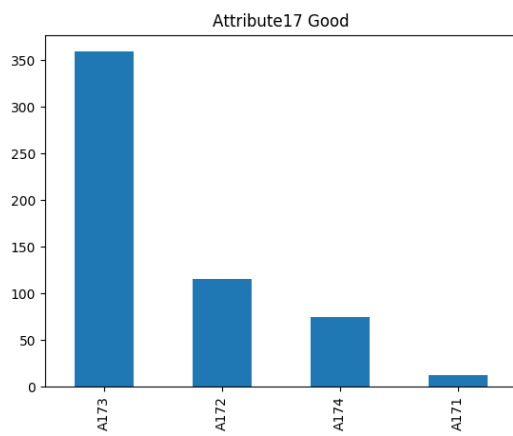
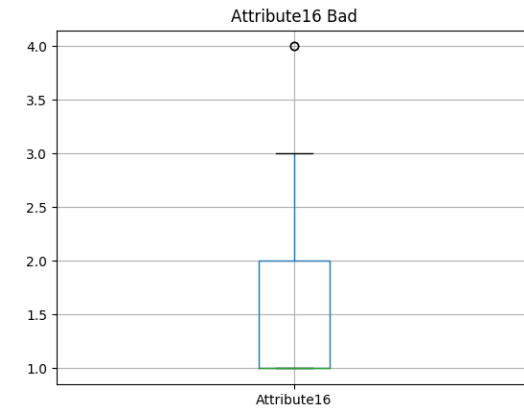
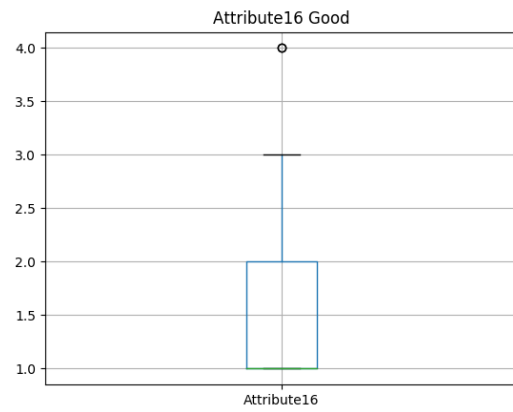
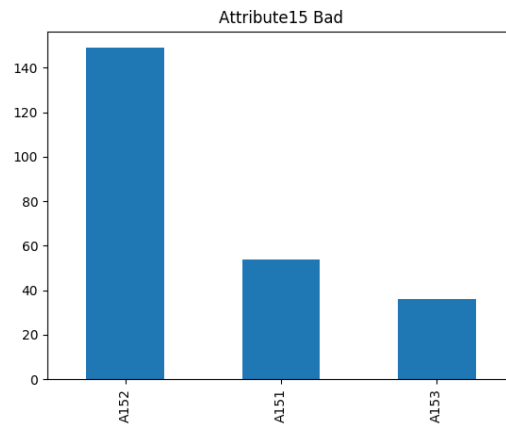
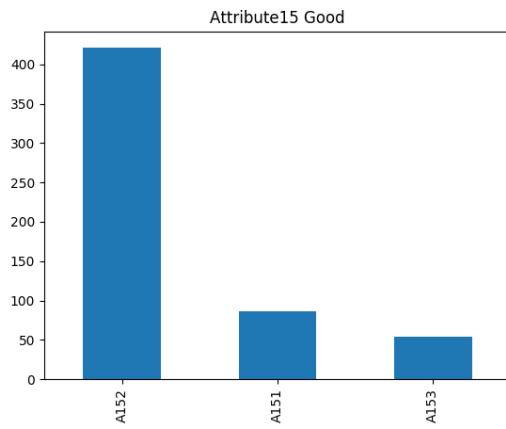


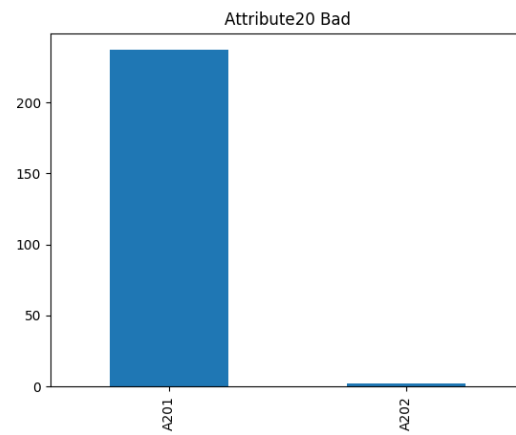
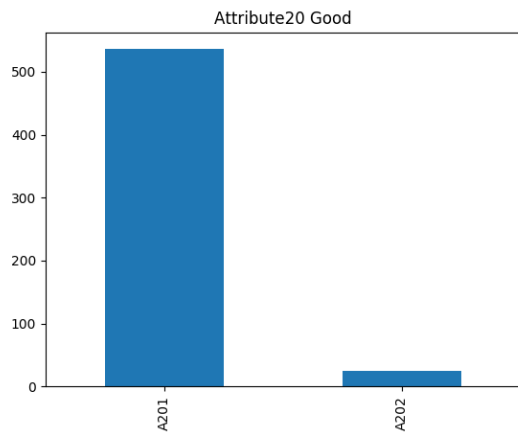
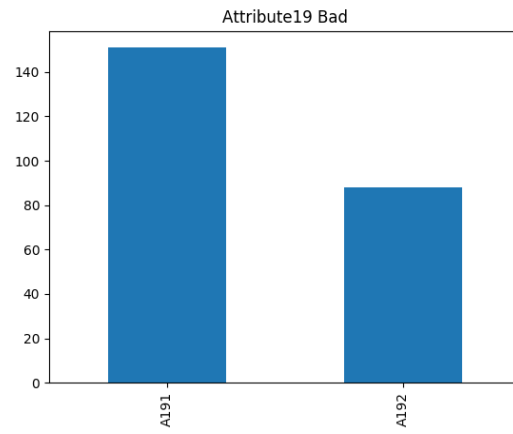
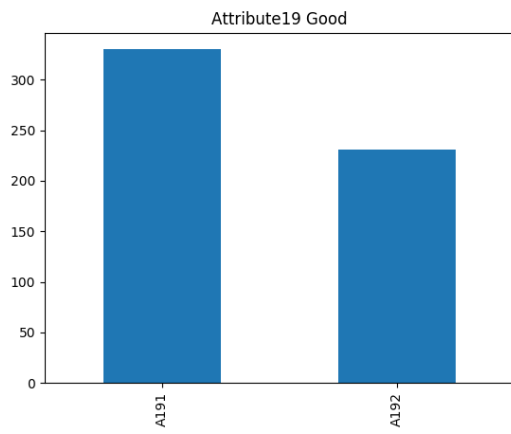
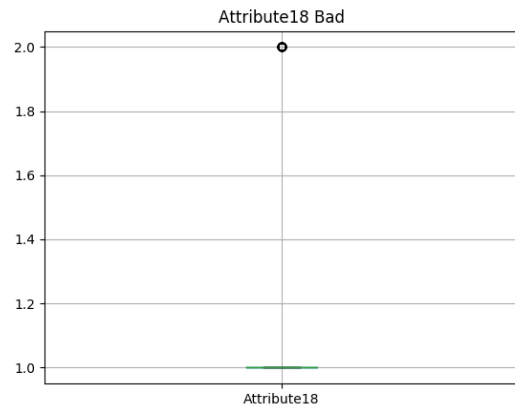
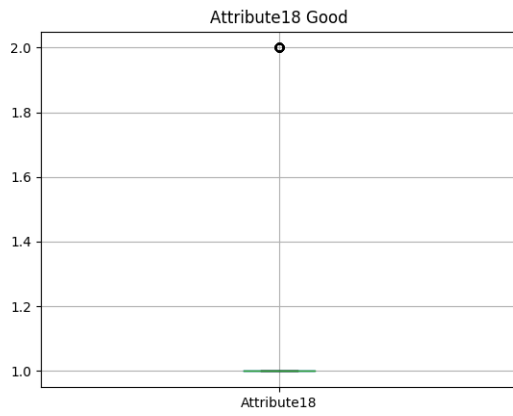














Παρατηρούμε στα παραπάνω διαγράμματα ότι σε κάποια attributes τα γραφήματα έχουν ελάχιστες διαφορές μεταξύ good και bad ενώ σε κάποια άλλα έχουν σημαντικές διαφορές. Πιστεύουμε ότι όσο πιο μεγάλες οι διαφορές ανάμεσα στα διαγράμματα, τόσο πιο χρήσιμη η πληροφορία που προσφέρει το συγκεκριμένο attribute στην κατηγοριοποίηση των δανειοληπτών.

Πιο συγκεκριμένα έντονες διαφορές παρατηρούνται στα attributes: 1, 2, 3, 5, 9, 12, 13 και 15. Αντιθέτως τα attributes: 8, 10, 11, 16, 18 και 19 έχουν αμελητέες διαφορές. Αναμένουμε τα attributes με τις μεγαλύτερες διαφορές να έχουν και μεγαλύτερο information gain.

## Classification

Εκτελέσαμε τους κατηγοριοποιητές Naive Bayes, Random Forest και SVM στο train.tsv με 10 fold cross validation και τα αποτελέσματα τα οποία προέκυψαν είναι τα εξής:

Statistic Measure	Naive Bayes	Random Forest	SVM
Accuracy	0.6375	0.74625	0.70125

## Επιλογή Features

Με βάση τα προηγούμενα αποτελέσματα συμπεραίναμε ότι ο Random Forest είναι η καλύτερη επιλογή για την κατηγοριοποίηση των πελατών. Επομένως το αρχείο testSet\_Predictions.csv συμπληρώθηκε με τη χρήση του παραπάνω κατηγοριοποιητή.

Στον παρακάτω πίνακα παρουσιάζεται το information gain σε αύξουσα ταξινόμηση για κάθε attribute.

Attribute	Information Gain
Attribute20	0.21
Attribute10	0.52
Attribute18	0.6
Attribute14	0.83
Attribute19	0.97
Attribute16	1.12
Attribute15	1.14
Attribute17	1.41
Attribute9	1.52
Attribute3	1.71
Attribute6	1.72
Attribute8	1.81
Attribute1	1.82
Attribute11	1.85
Attribute12	1.95
Attribute7	2.15
Attribute4	2.69
Attribute2	3.73
Attribute13	5.26
Attribute5	9.52

Όντως τα γραφήματα με μεγάλες διαφορές ανάμεσα σε good και bad στο προηγούμενο ερώτημα βρίσκονται χαμηλά στον πίνακα, δηλαδή έχουν μεγάλο Information Gain. Αντίθετα τα γραφήματα με μικρές διαφορές έχουν μικρό Information Gain.

Τέλος παρουσιάζεται η διακύμανση του μέσου accuracy για 10 fold cross validation στον κατηγοριοποιητή Random Forest καθώς αφαιρούνται feature. Αξίζει να σημειωθεί ότι η σειρά με την οποία αφαιρούνται τα features ακολουθούν την σειρά με την οποία είναι στον προηγούμενο πίνακα, δηλαδή αφαιρούνται πρώτα τα attributes με το μικρότερο information gain.

