

## **Αναφορά 2<sup>ης</sup> Άσκησης Τεχνικές Εξόρυξης Δεδομένων**

### **1.Οπτικοποίηση Δεδομένων**

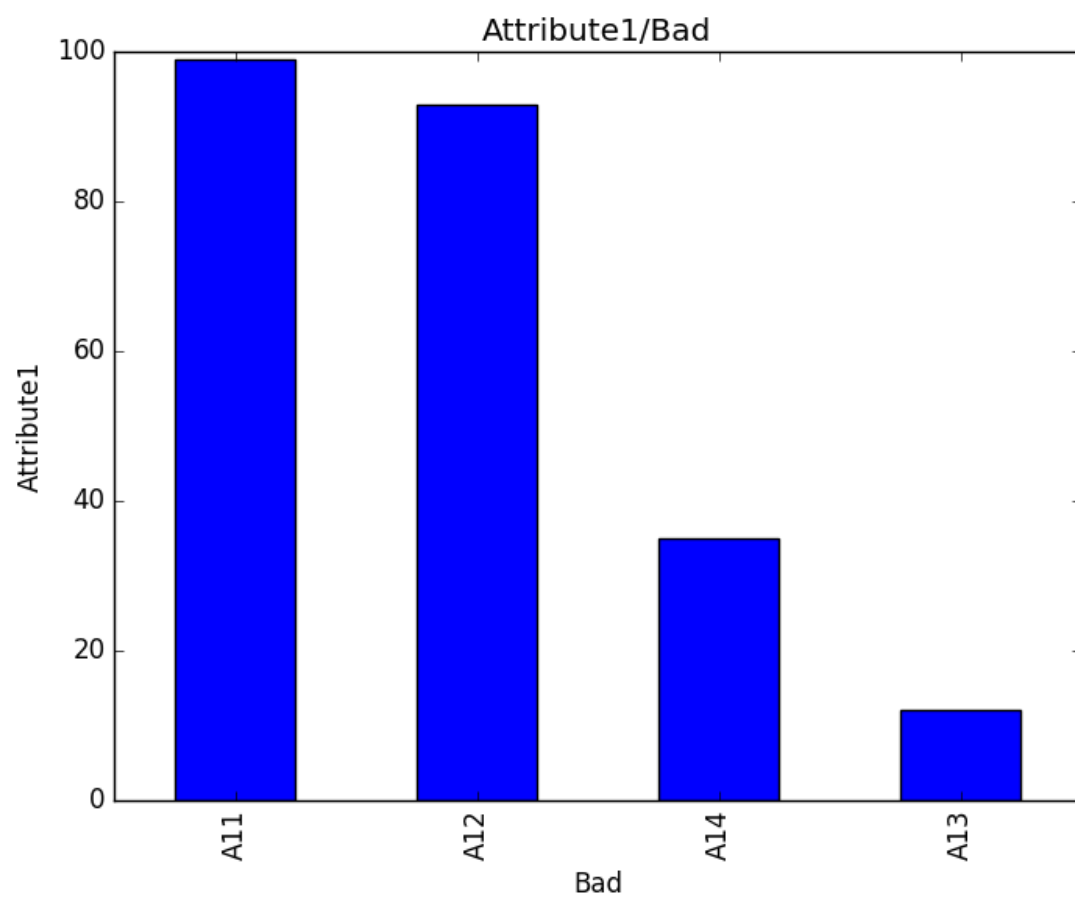
Σε αυτό το ερώτημα υλοποιήσαμε πλήρως τα histograms και box plots για τα categorical και numerical data αντίστοιχα.

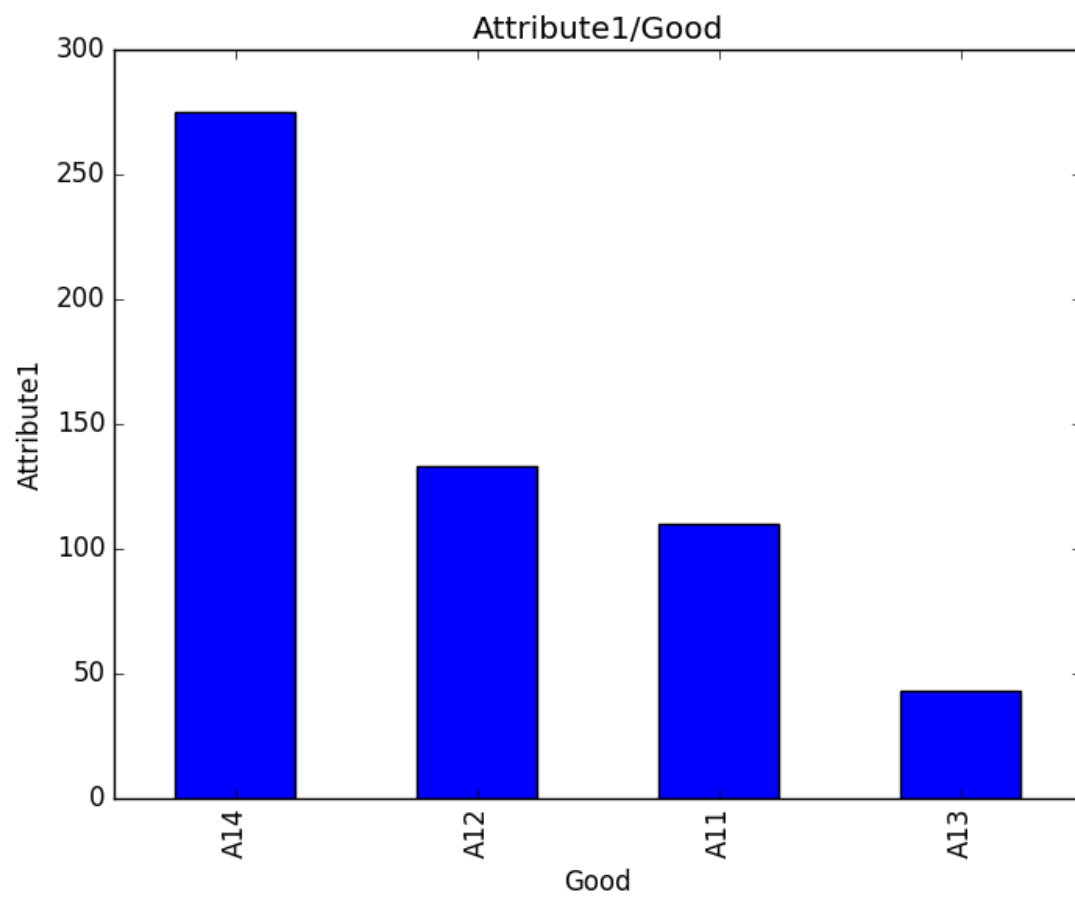
#### **Συμπέρασμα:**

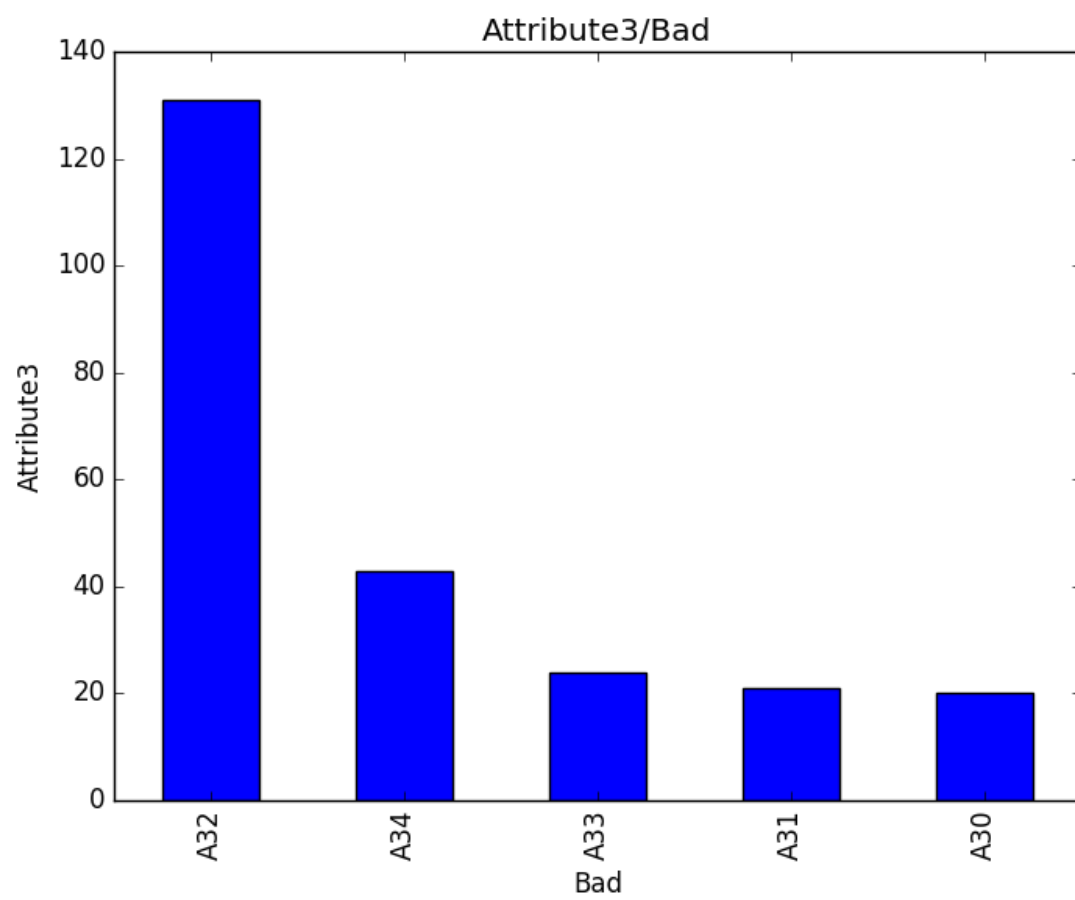
Με βάση την εντροπία, περιμένουμε ως πιο χρήσιμα τα στοιχεία των οποίων τα δεδομένα είναι πιο «μοιρασμένα». Πιο συγκεκριμένα αν οι μπάρες, για παράδειγμα, ενός Attribute είναι ίδιες για όλα τα στοιχεία (entropy = 1), τότε το Attribute αυτό αποτελεί πολύ καλό στοιχείο για το split του δένδρου αποφάσεων.

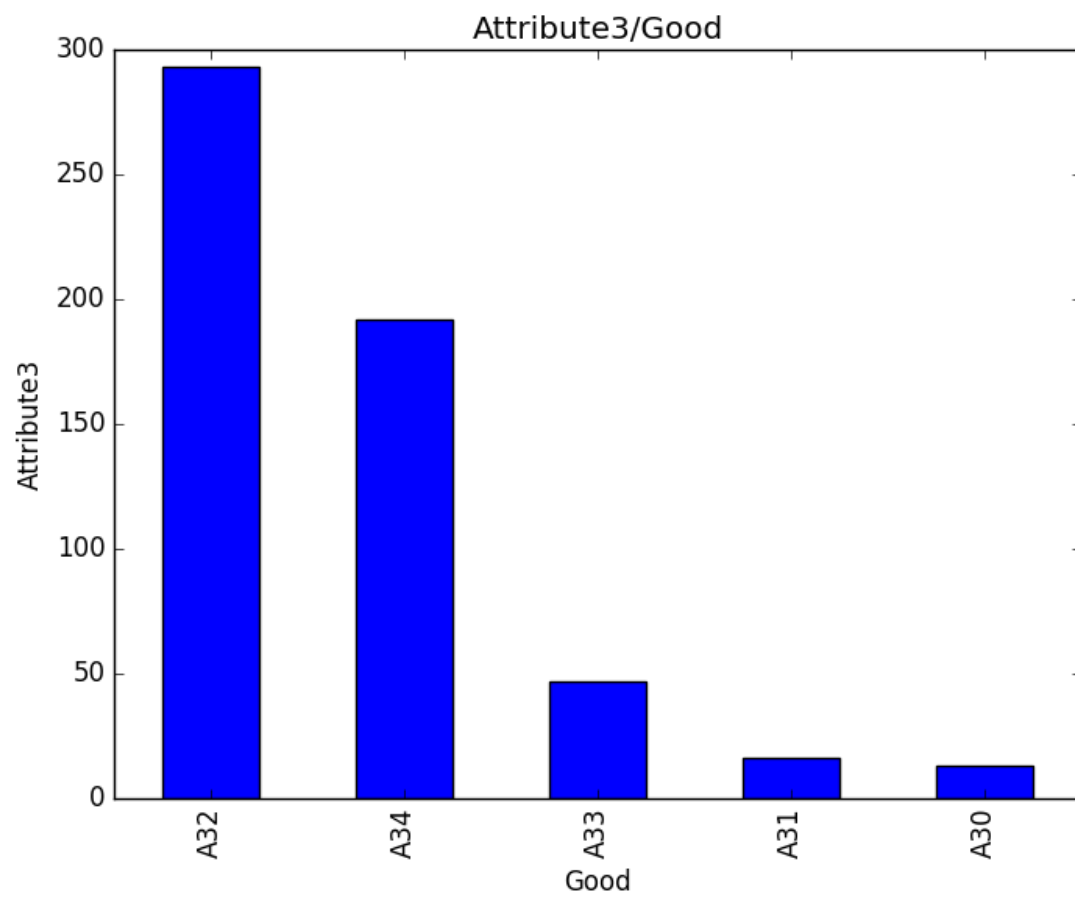
Παραθέτονται όλα τα διαγράμματα (histograms & box plots):

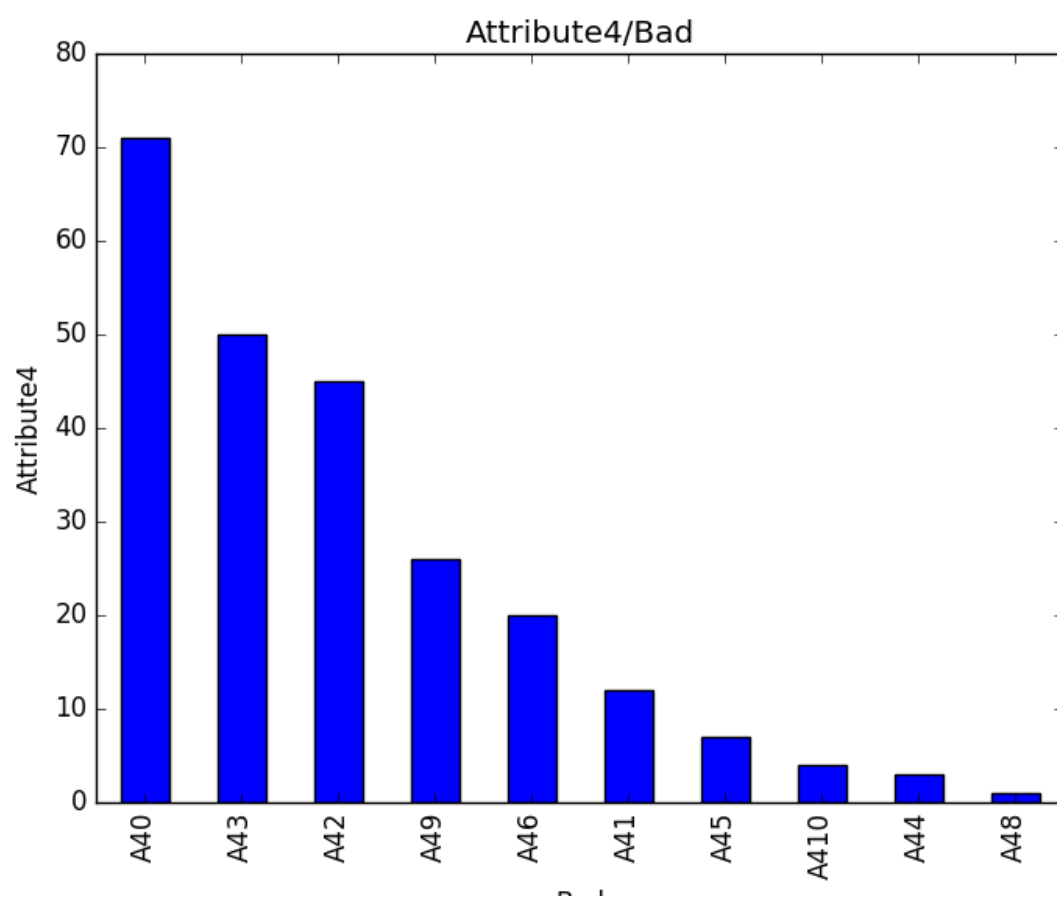
\*Επειδή ζητήσατε όλα τα διαγράμματα είναι προφανές ότι το report ξεπερνάει τις 30 σελίδες.

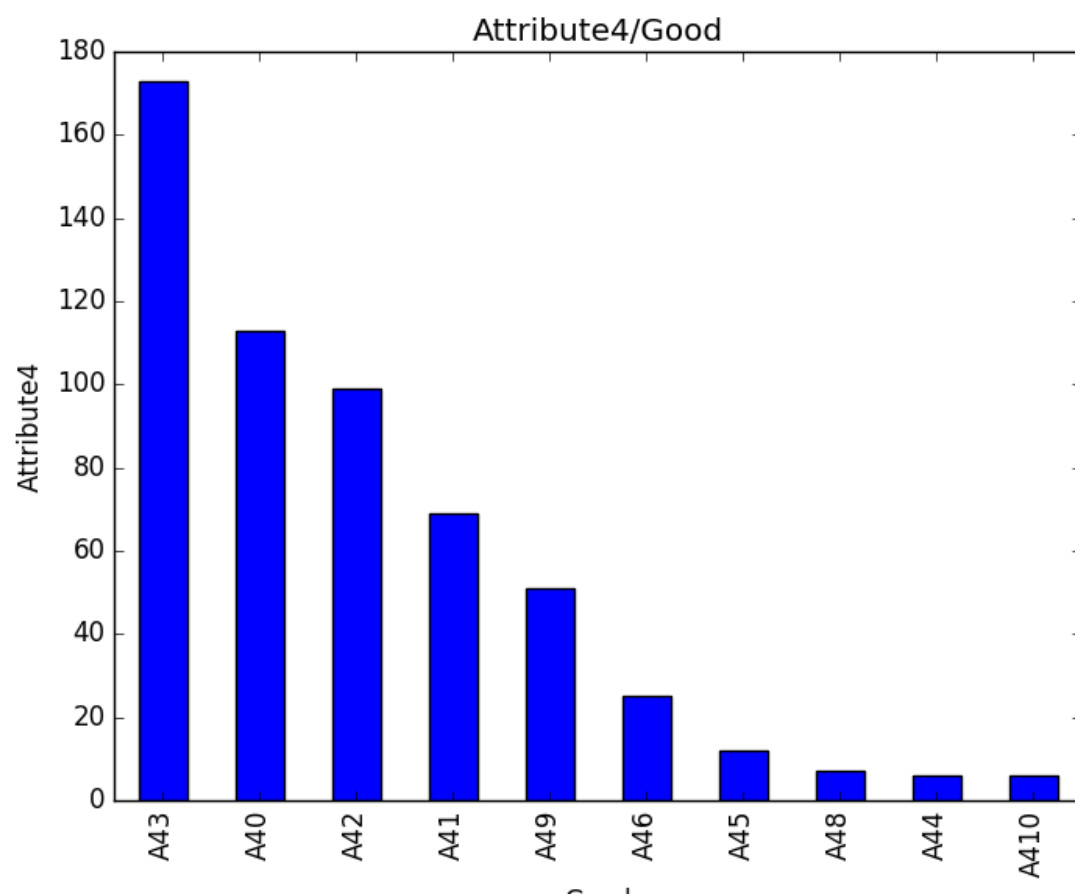


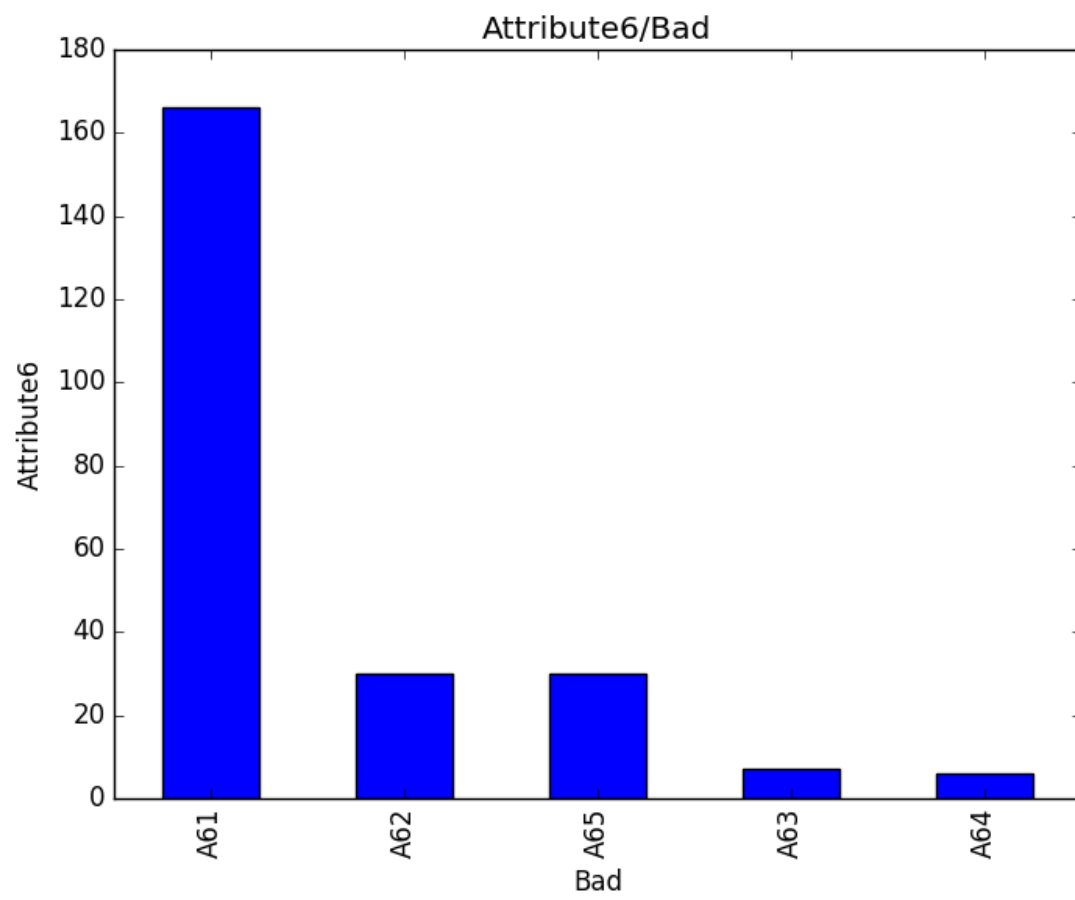




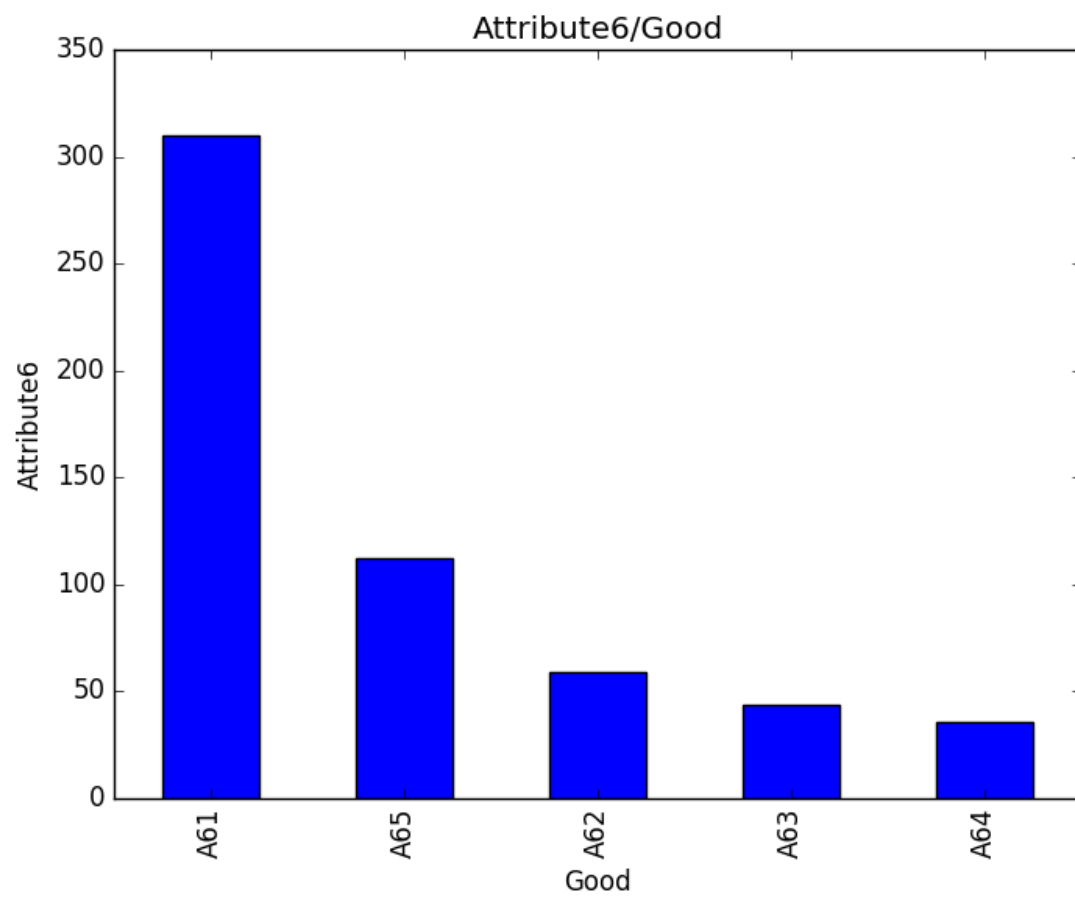


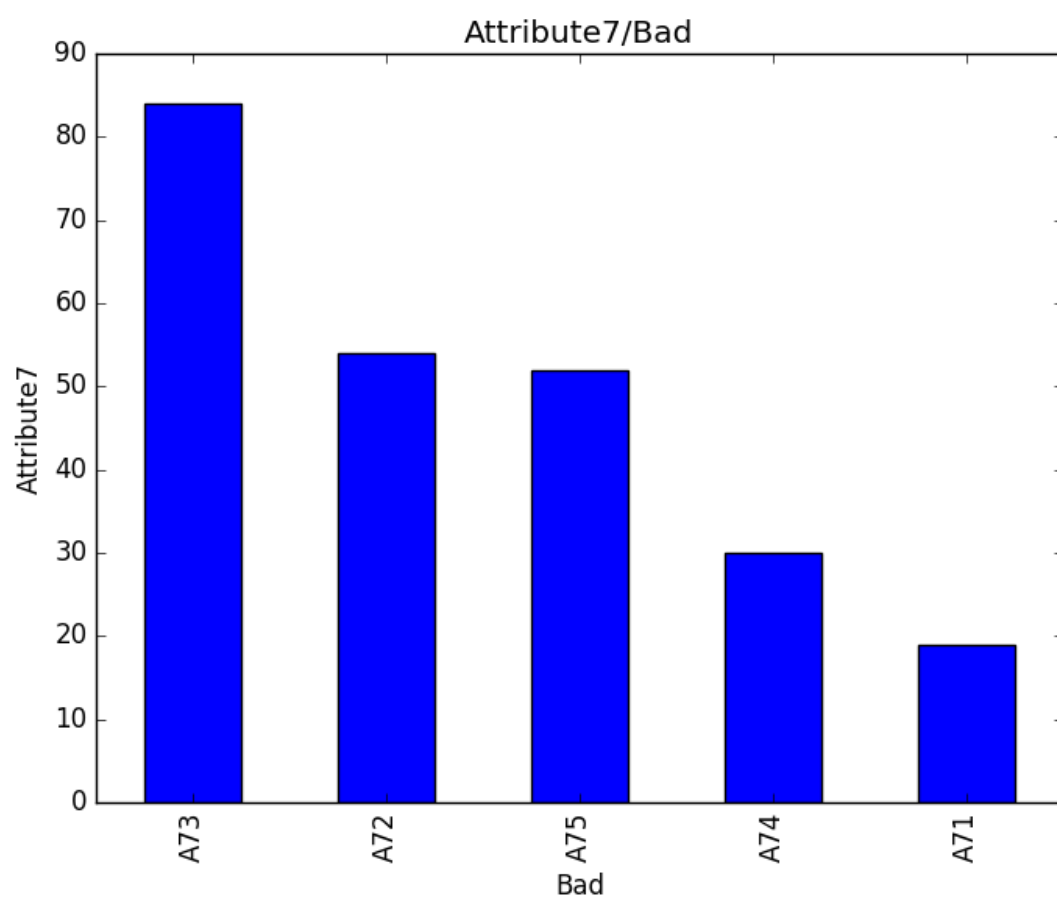


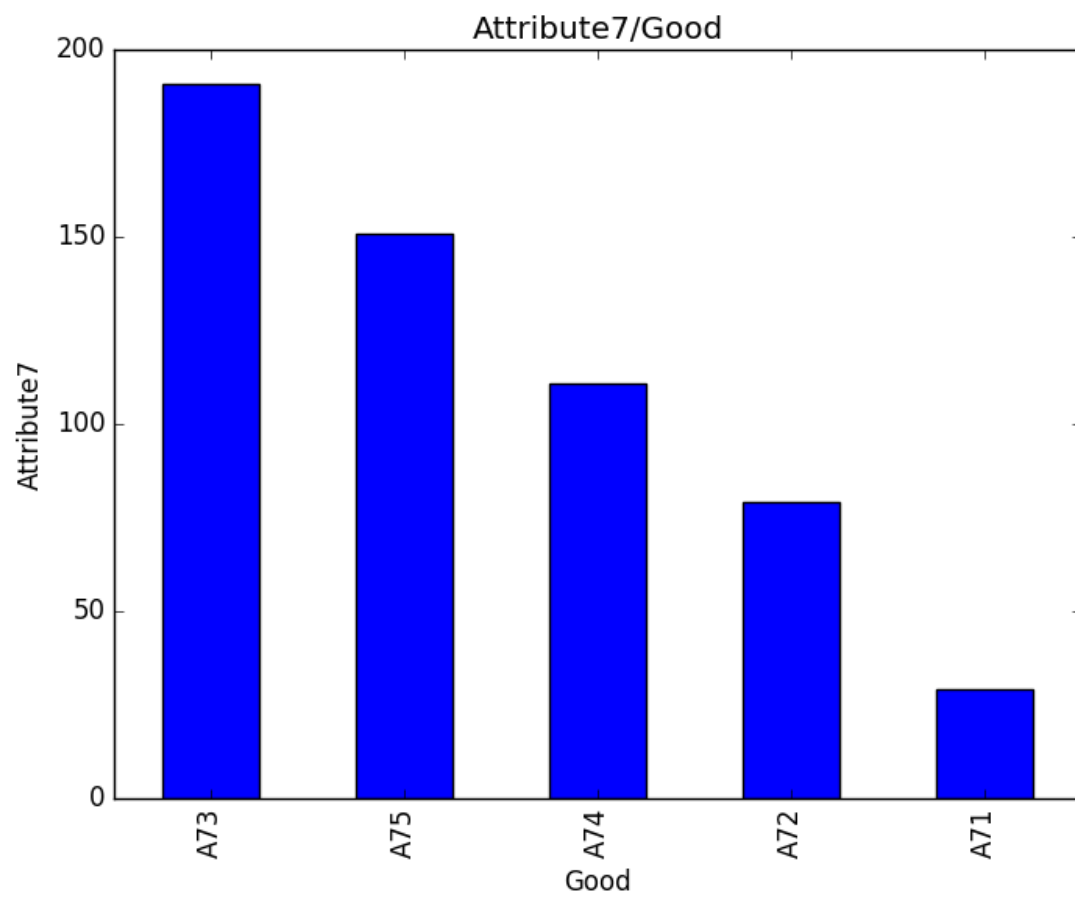


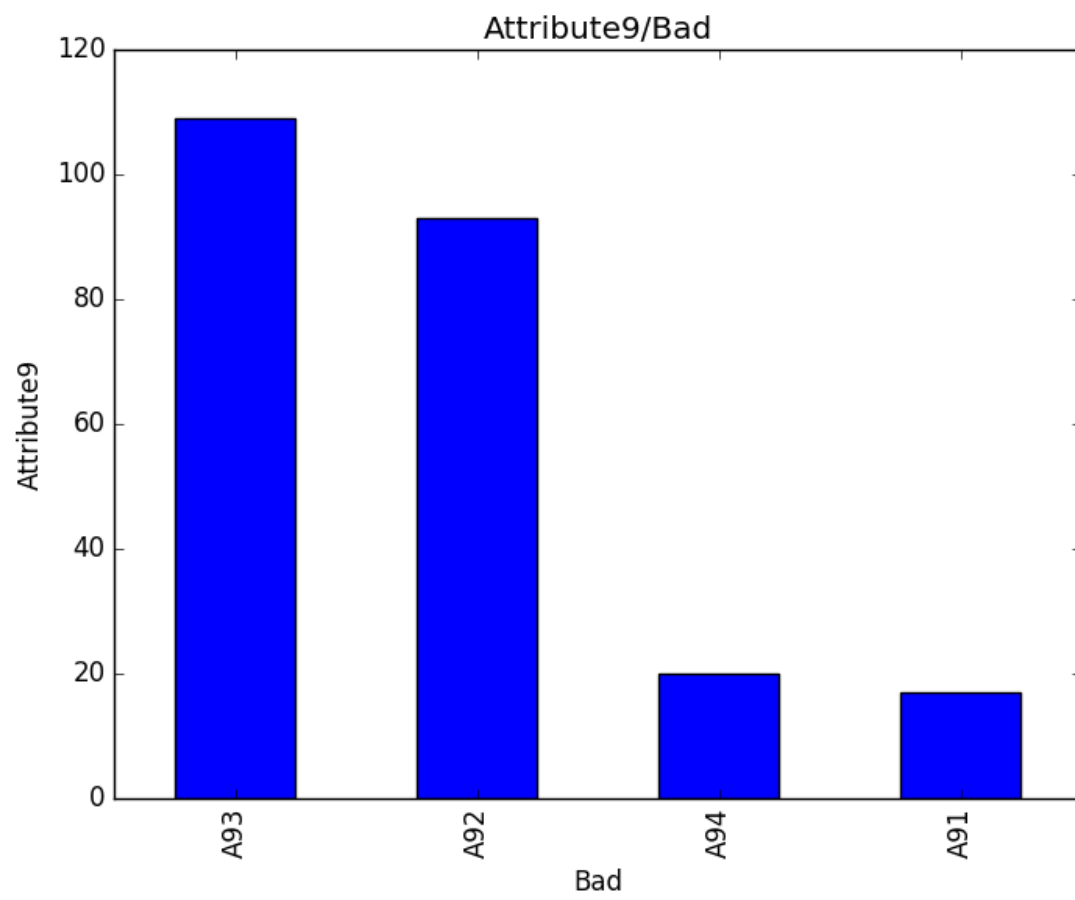


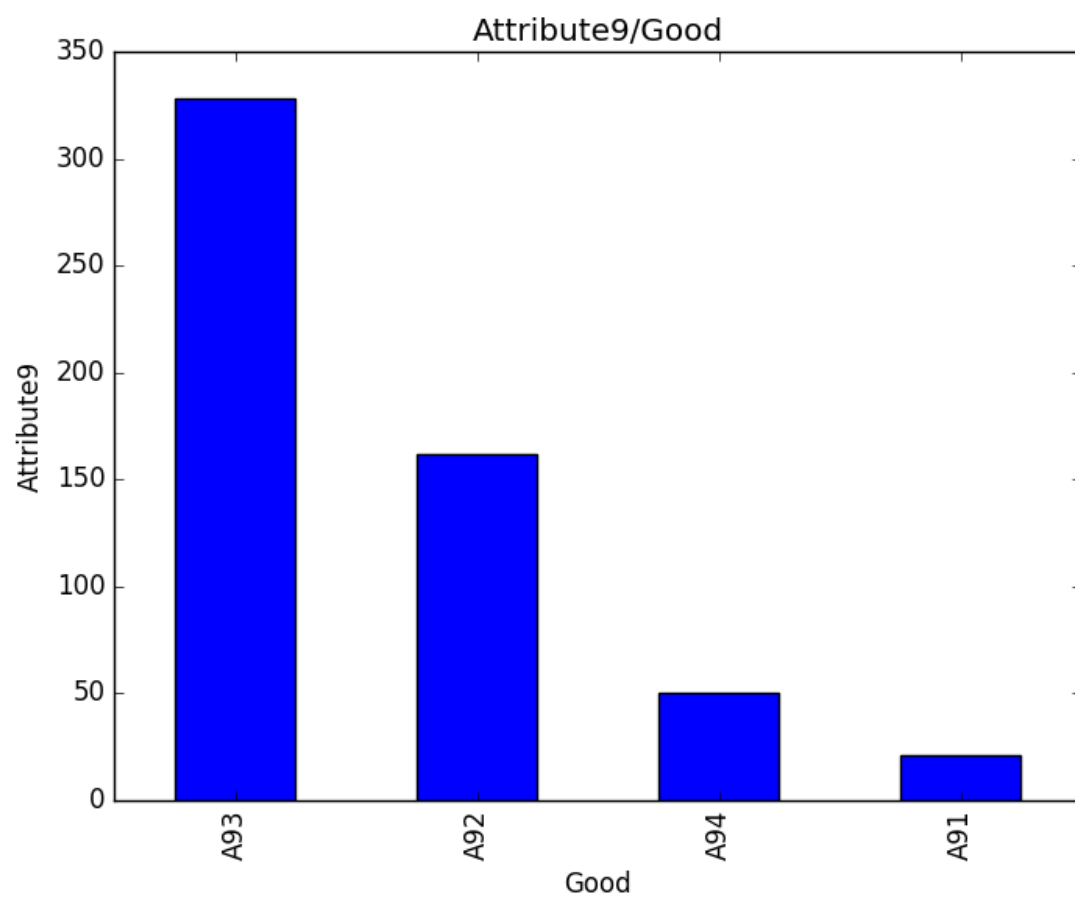


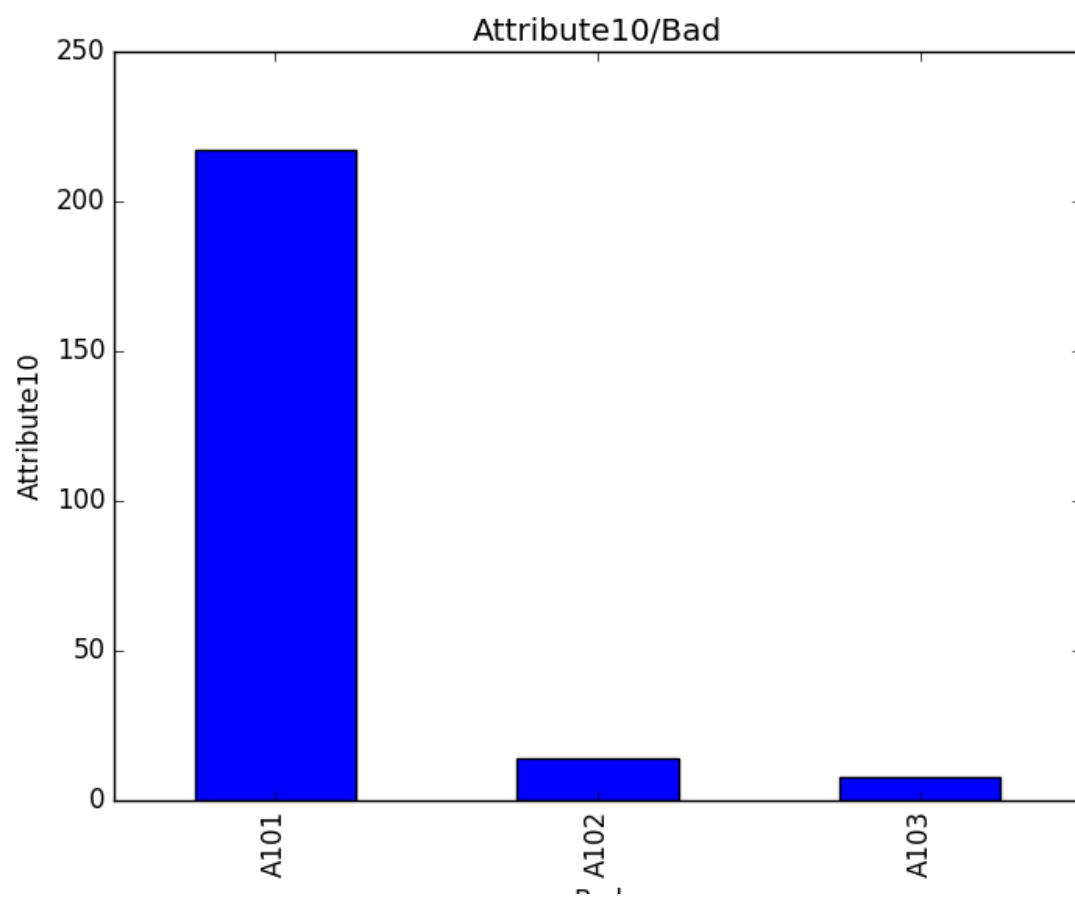


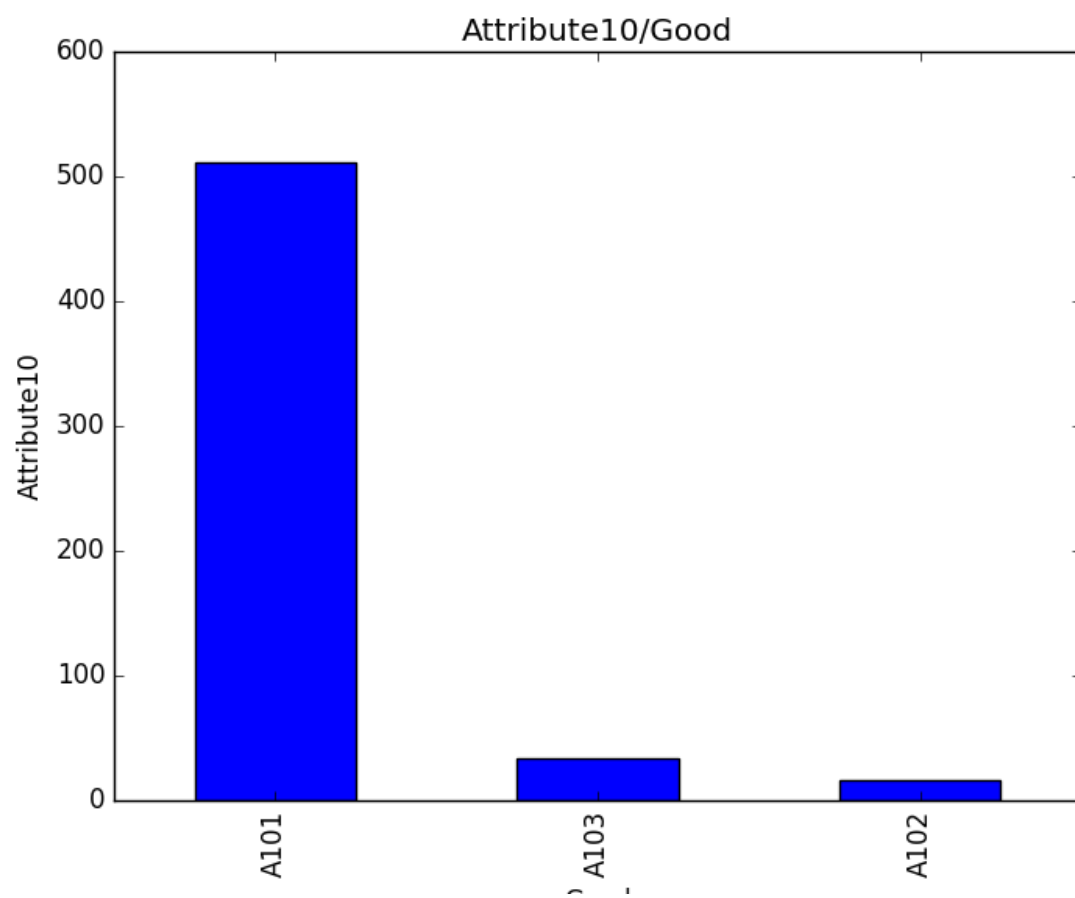


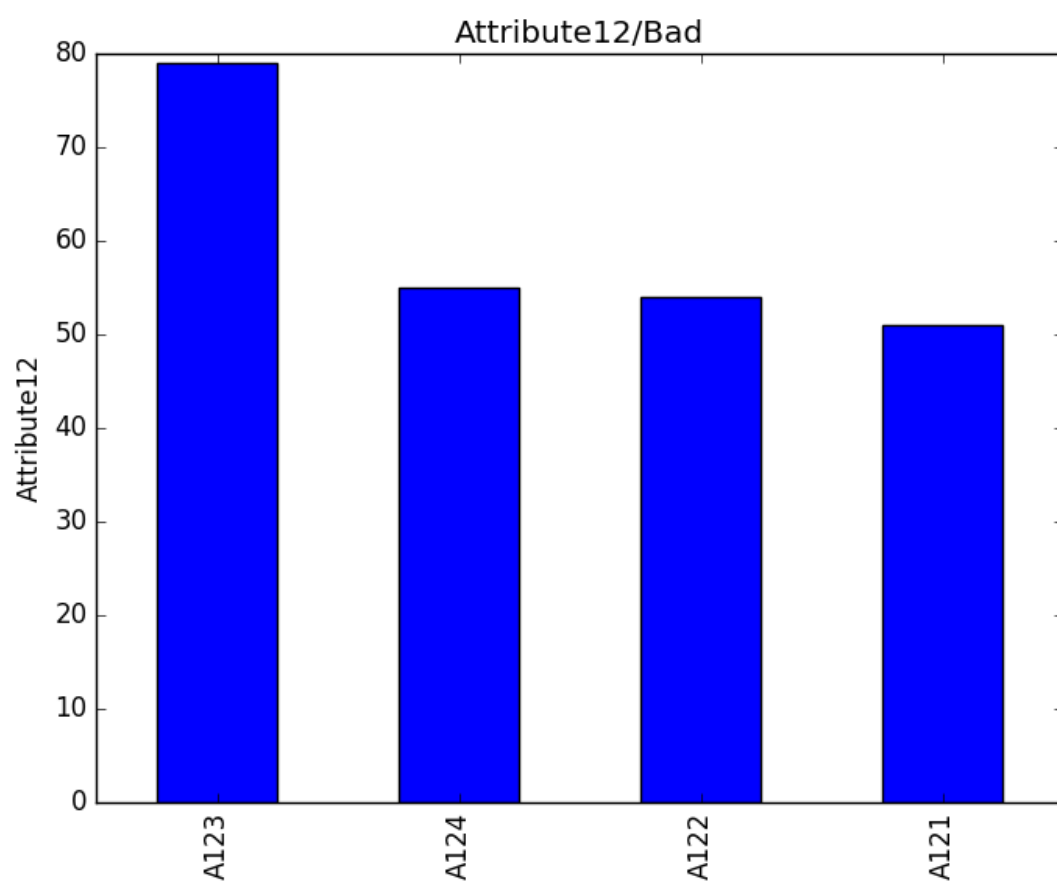




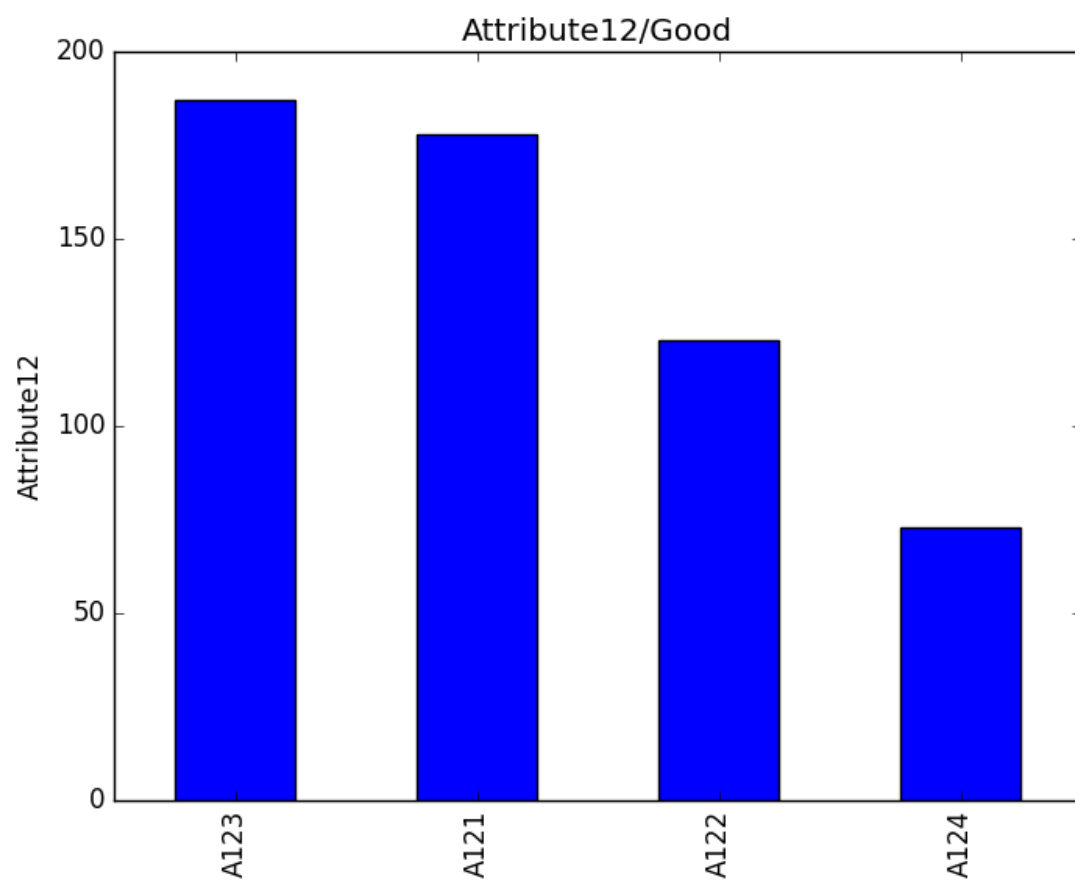


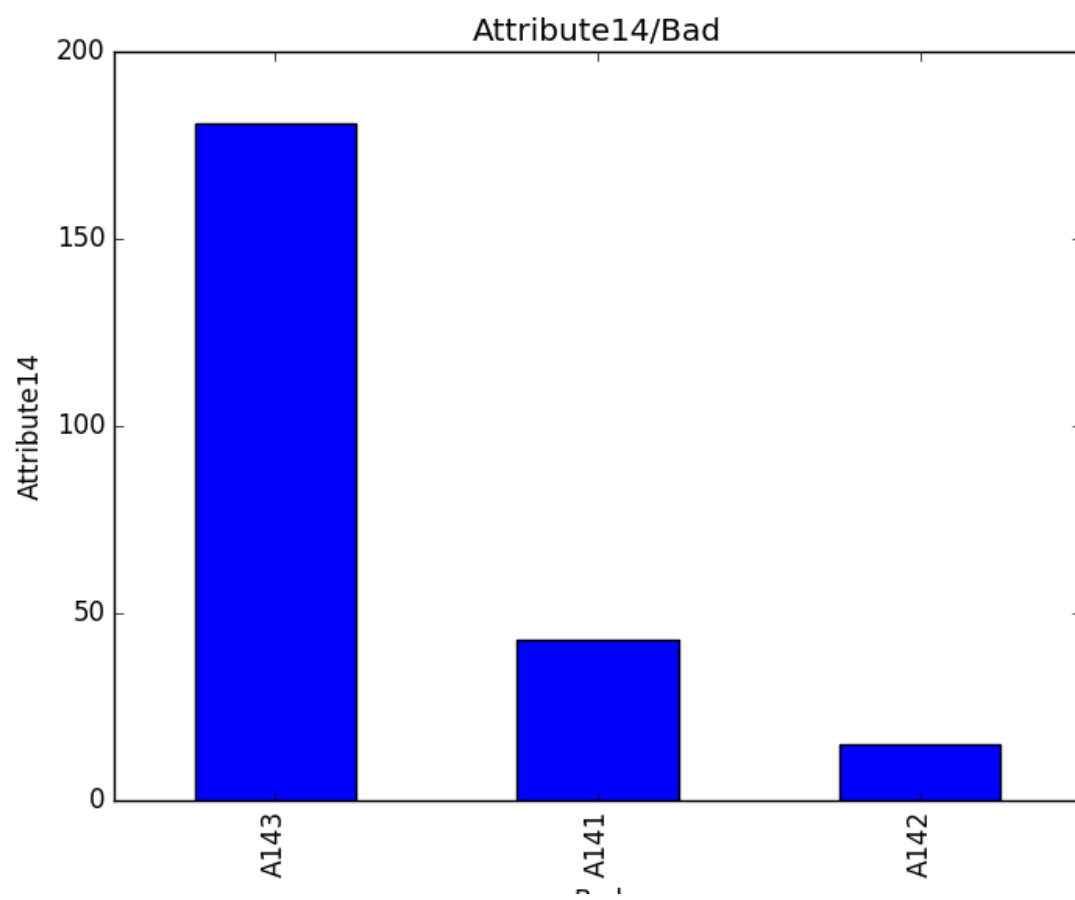


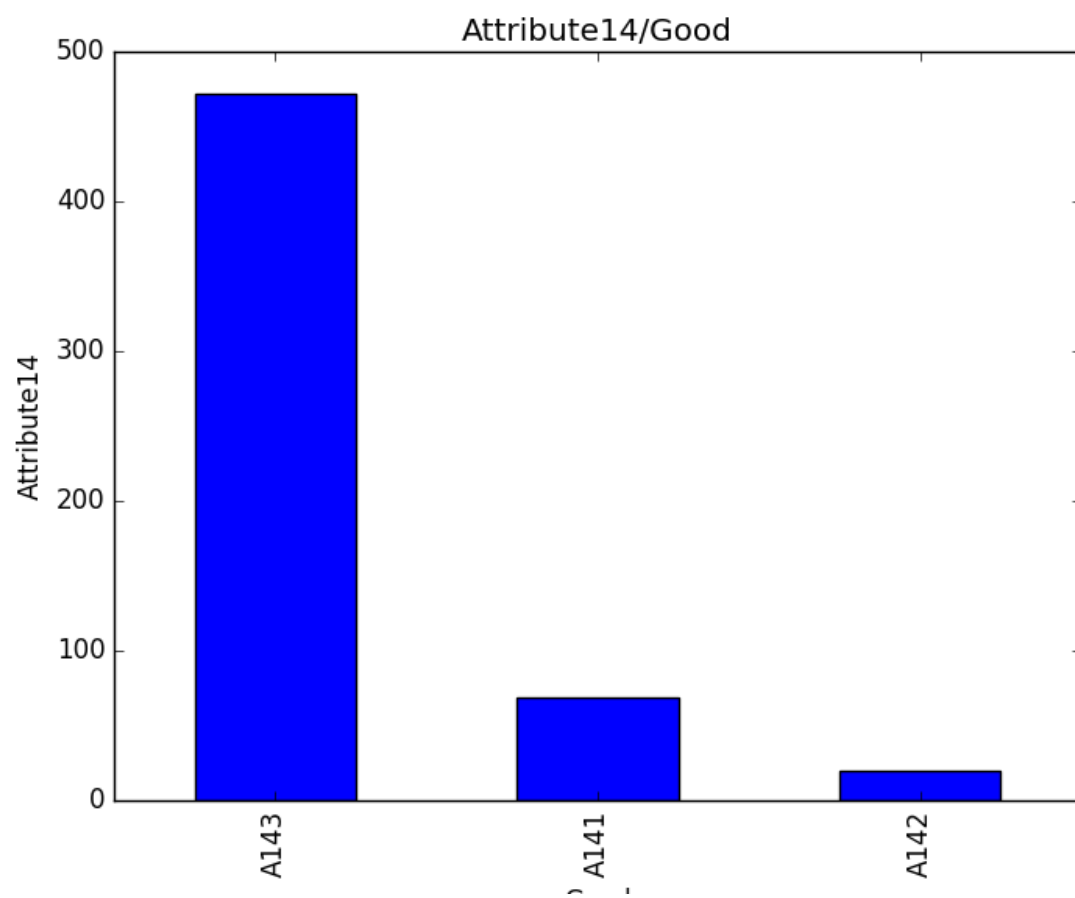


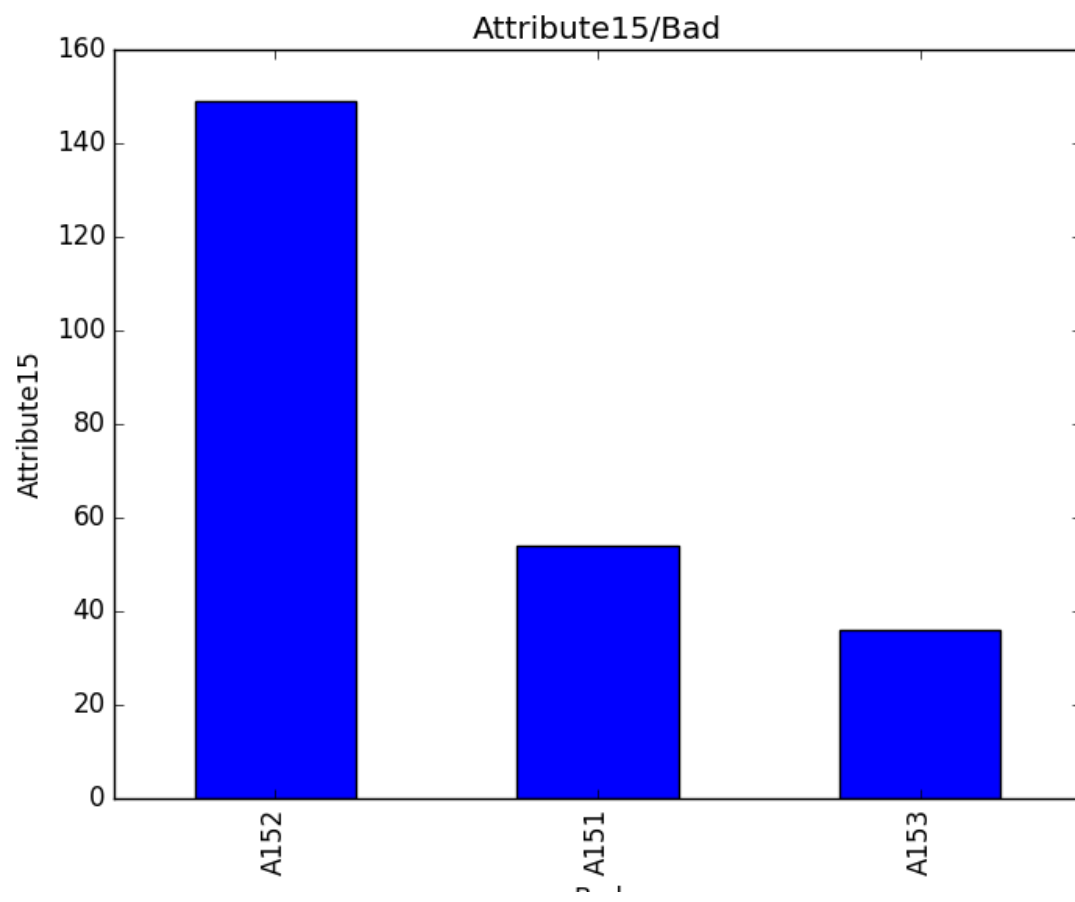


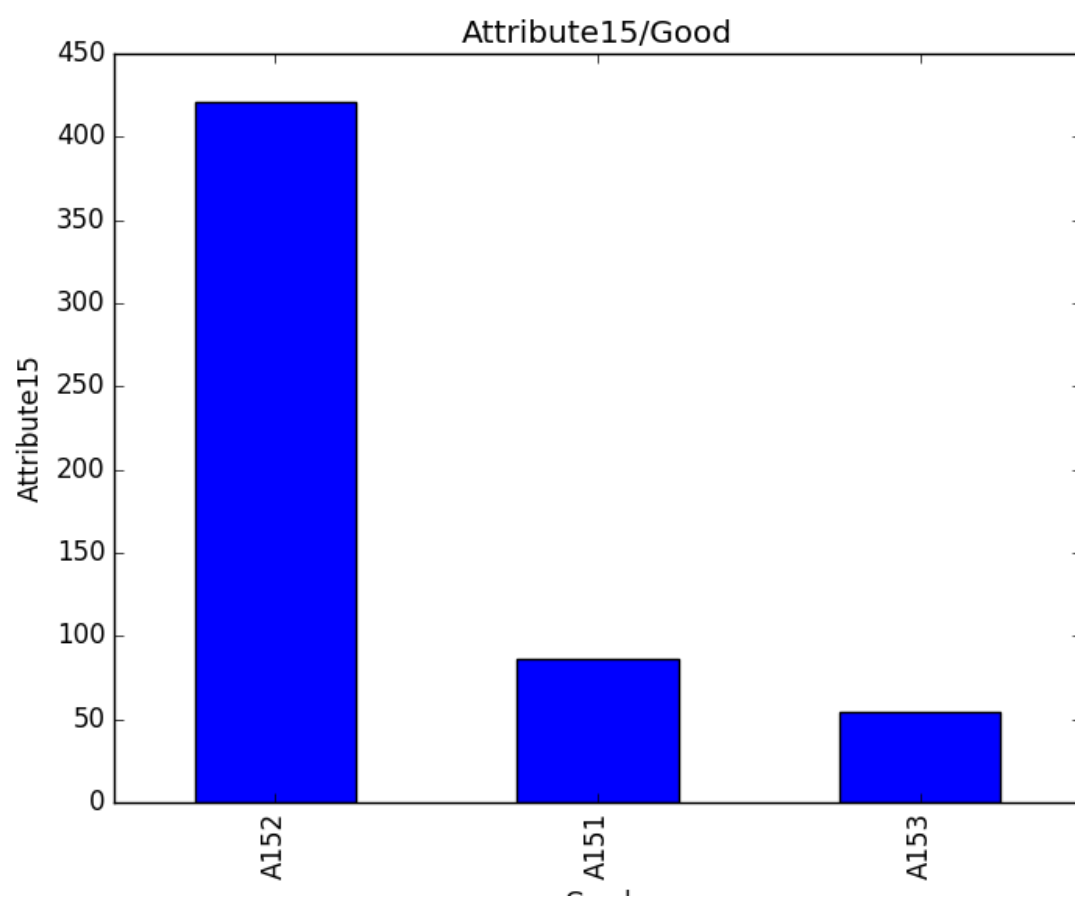


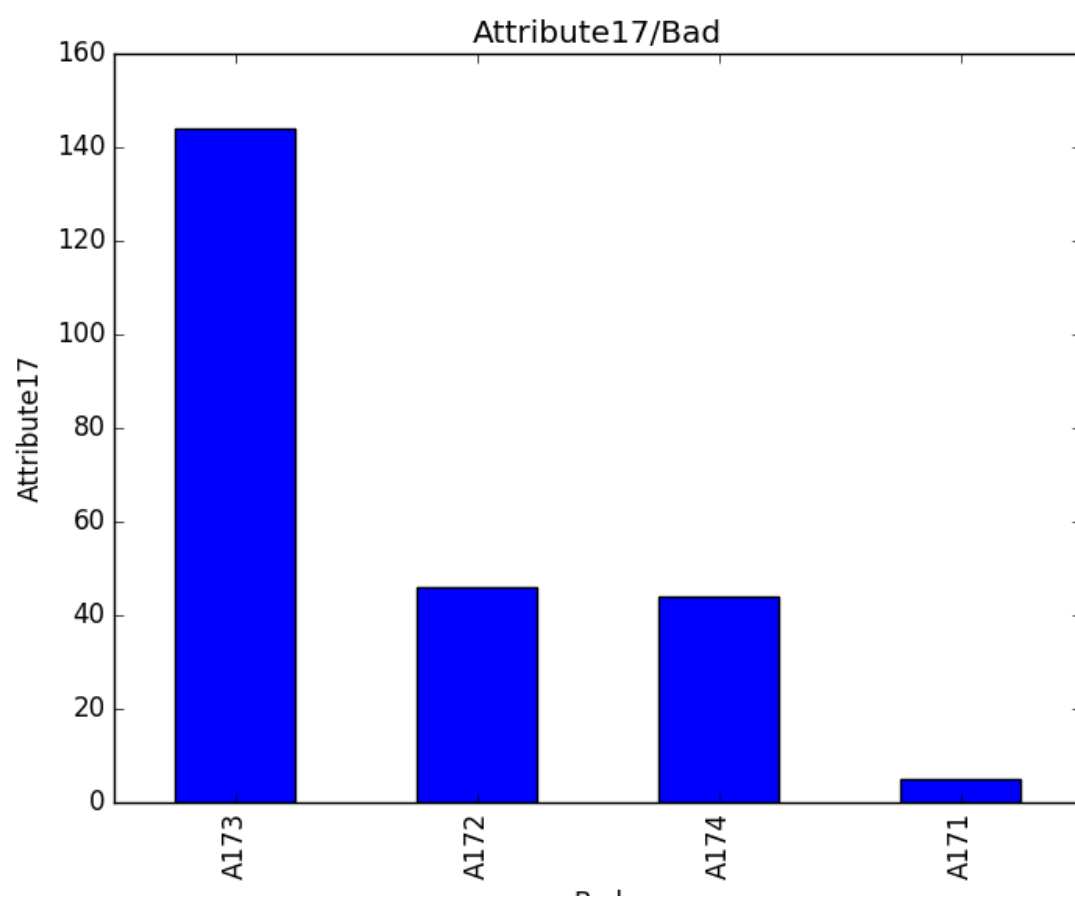


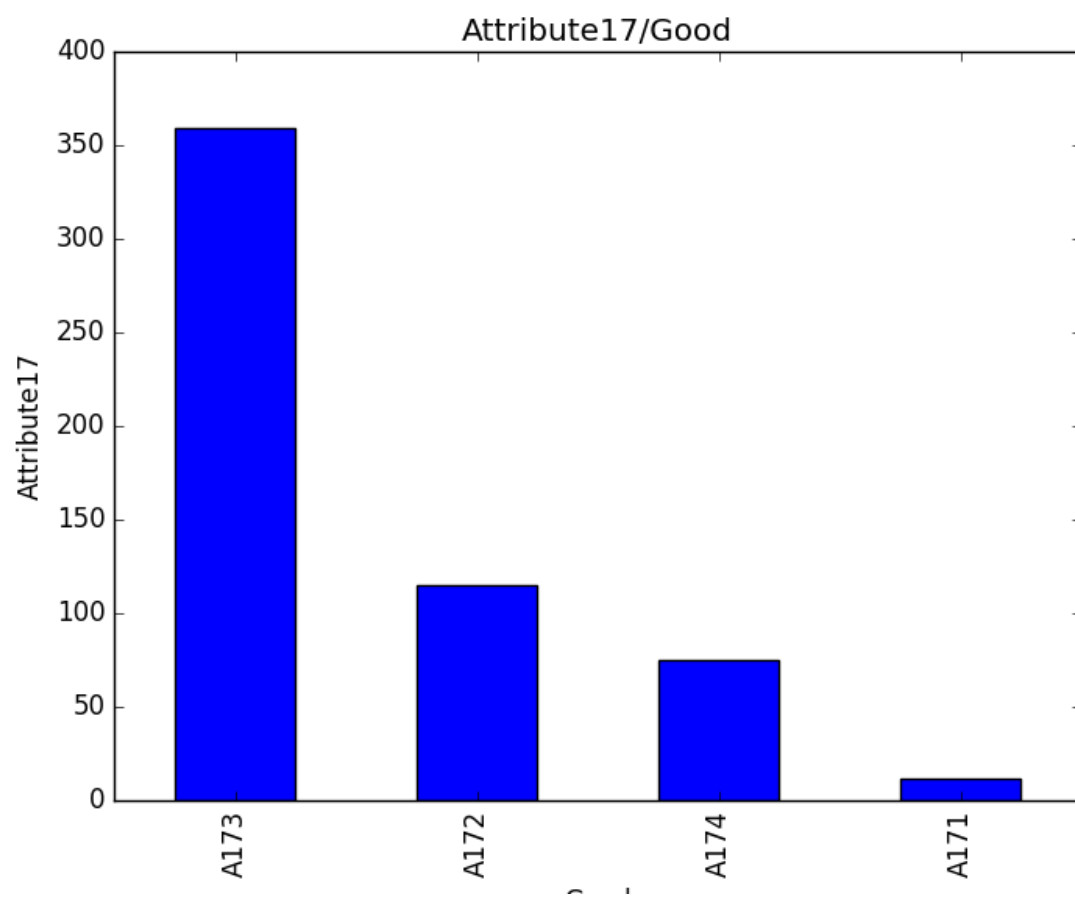


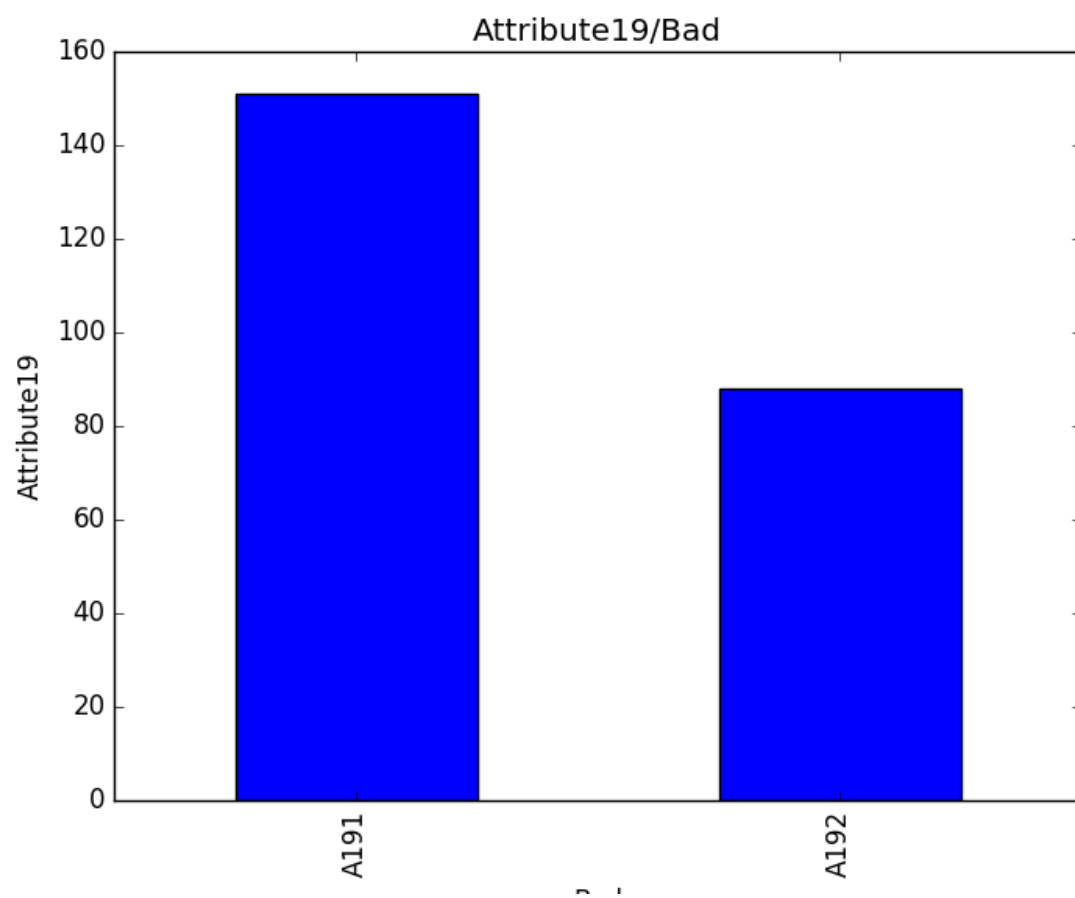




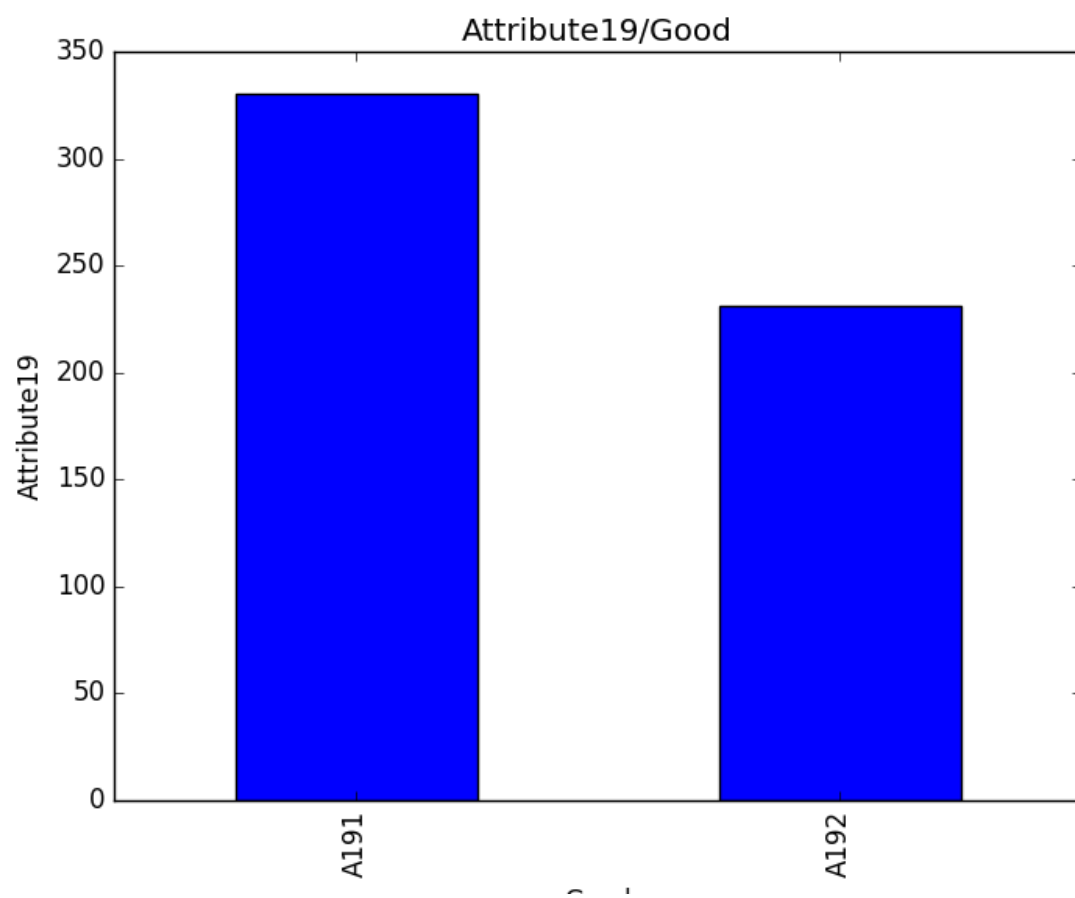


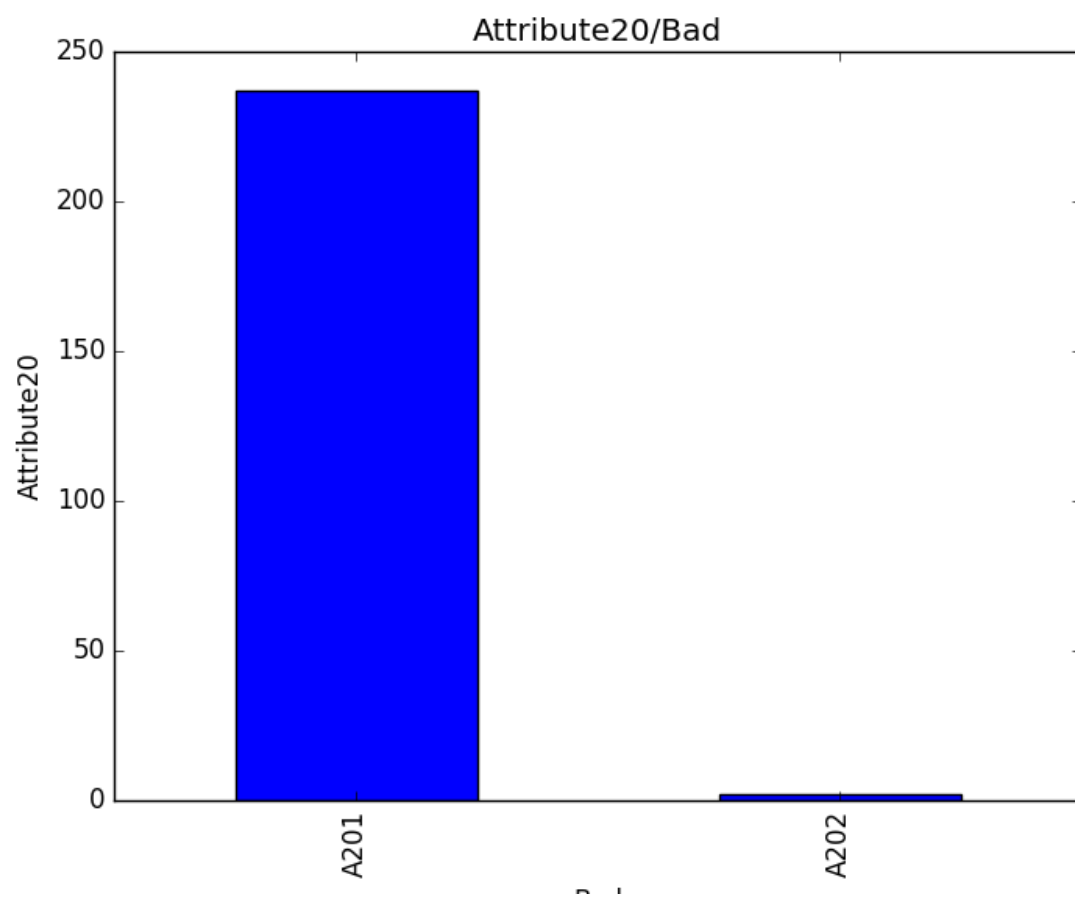


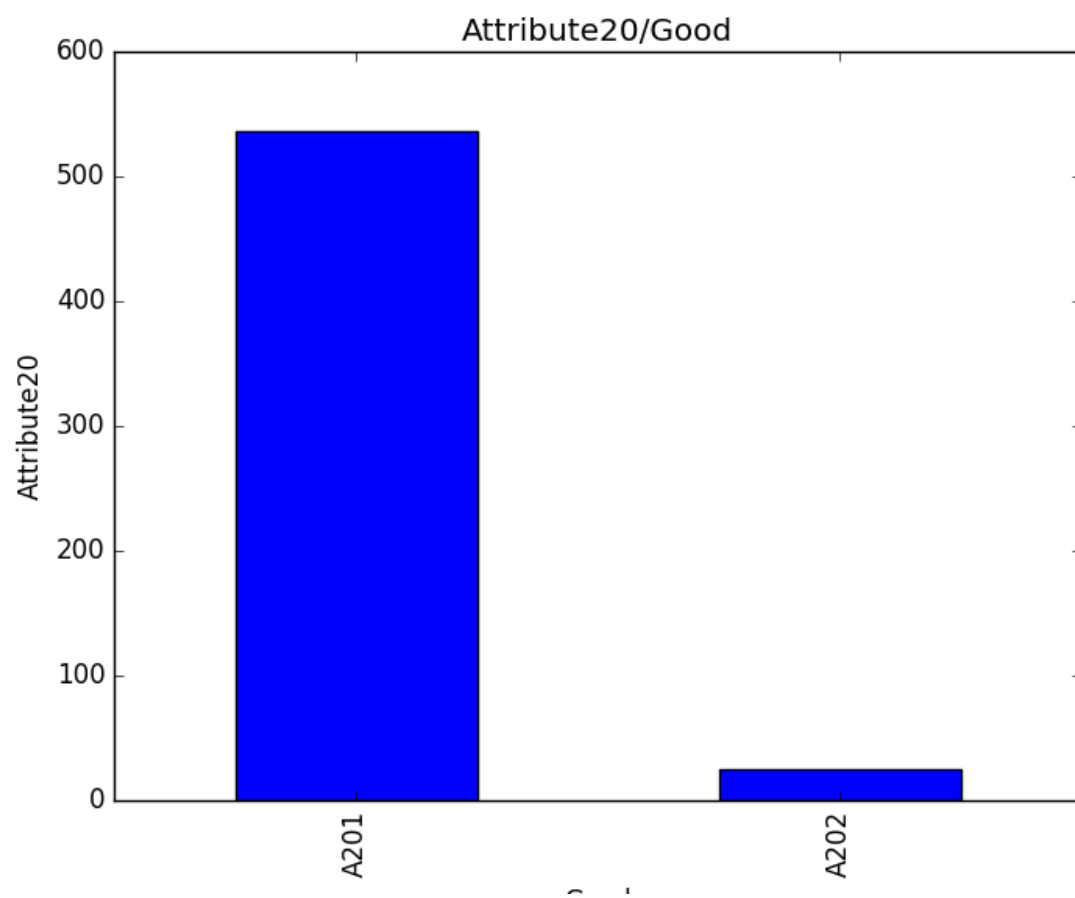


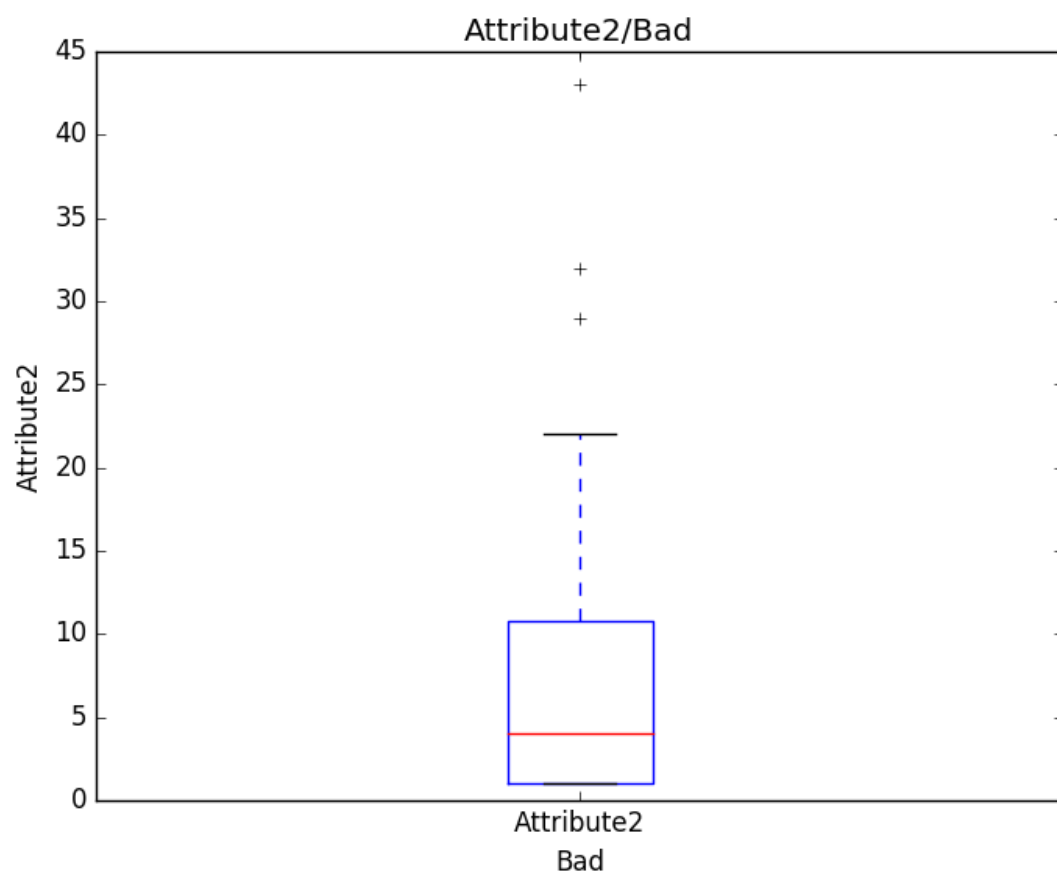


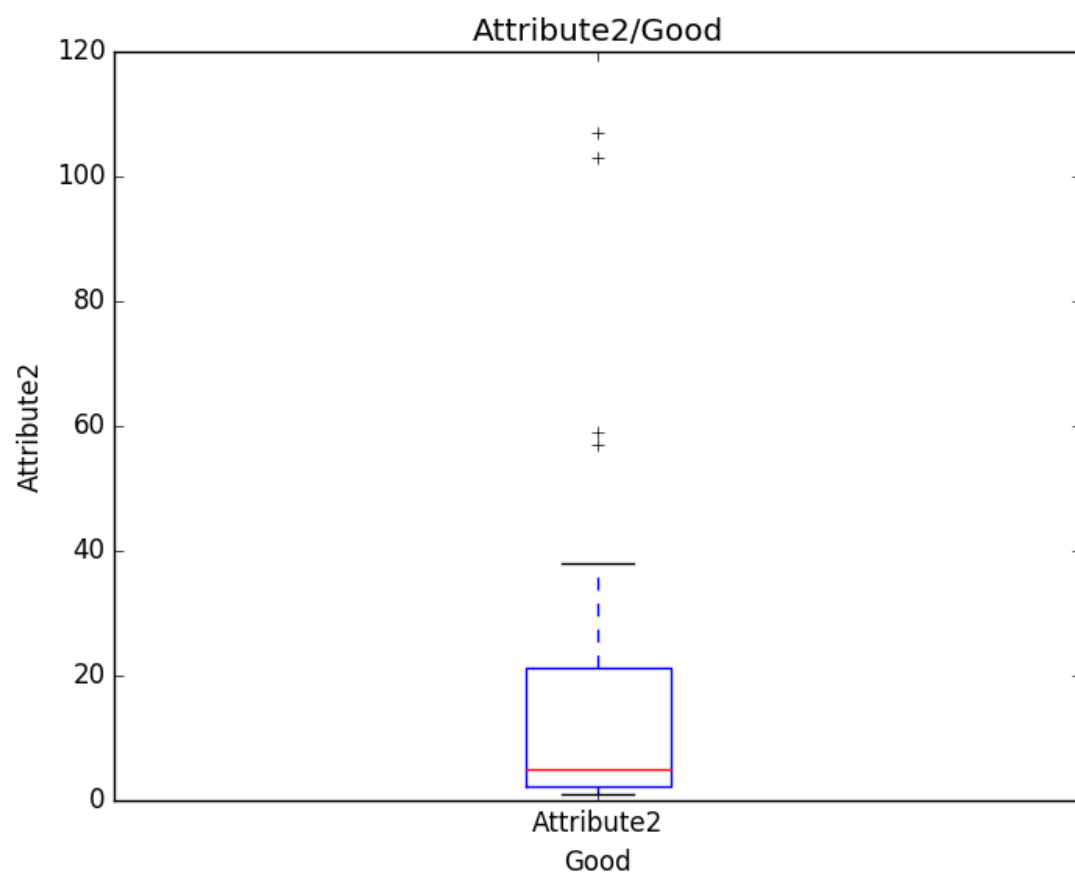


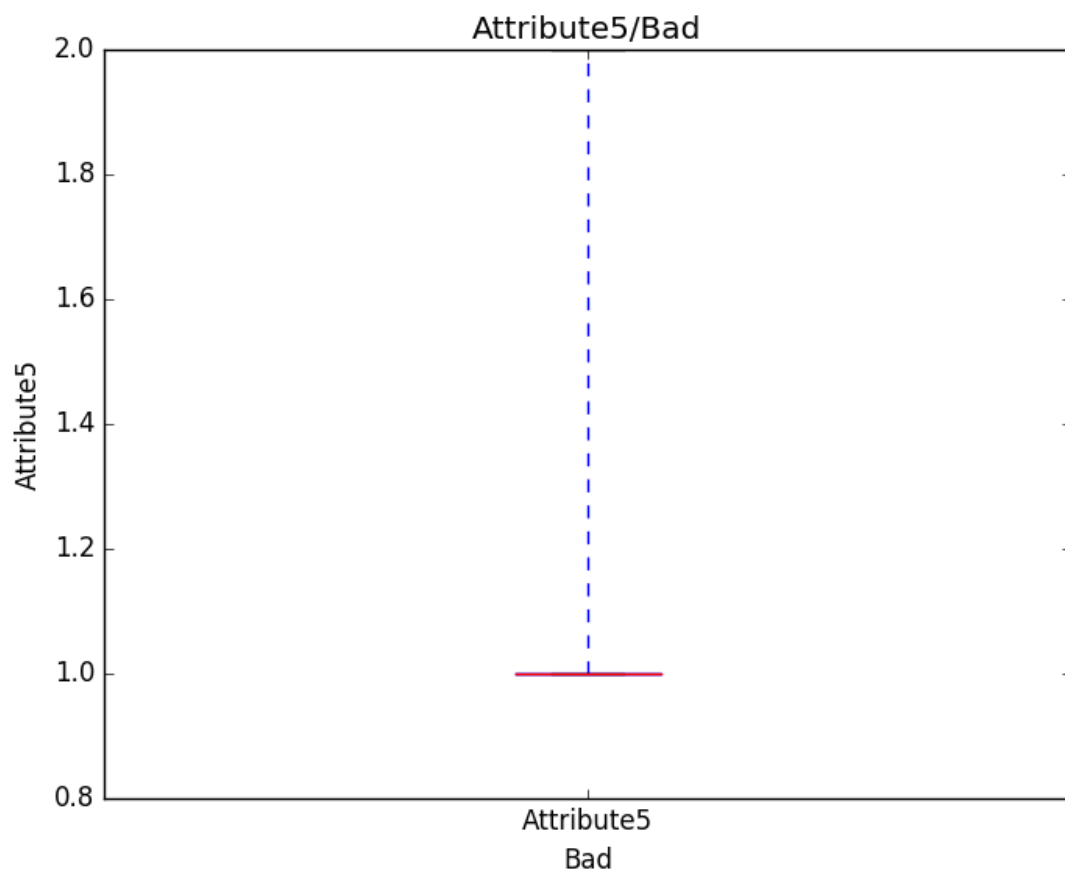


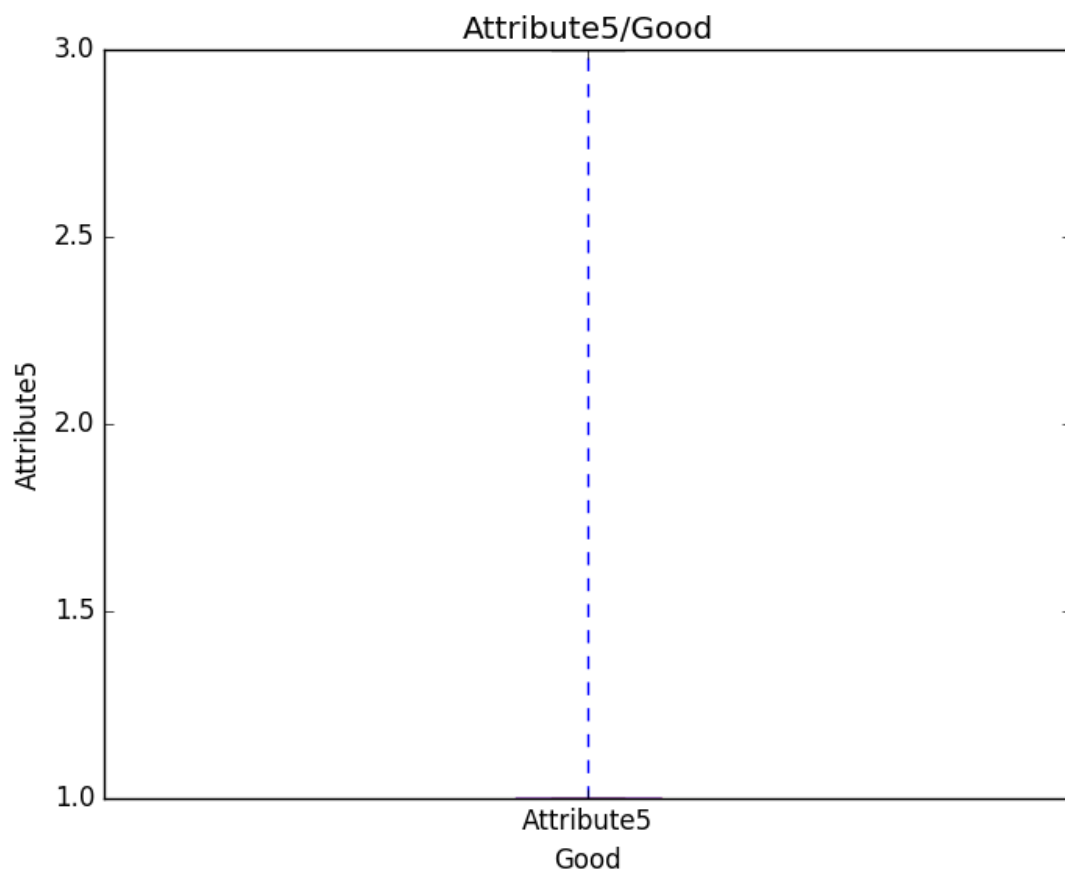


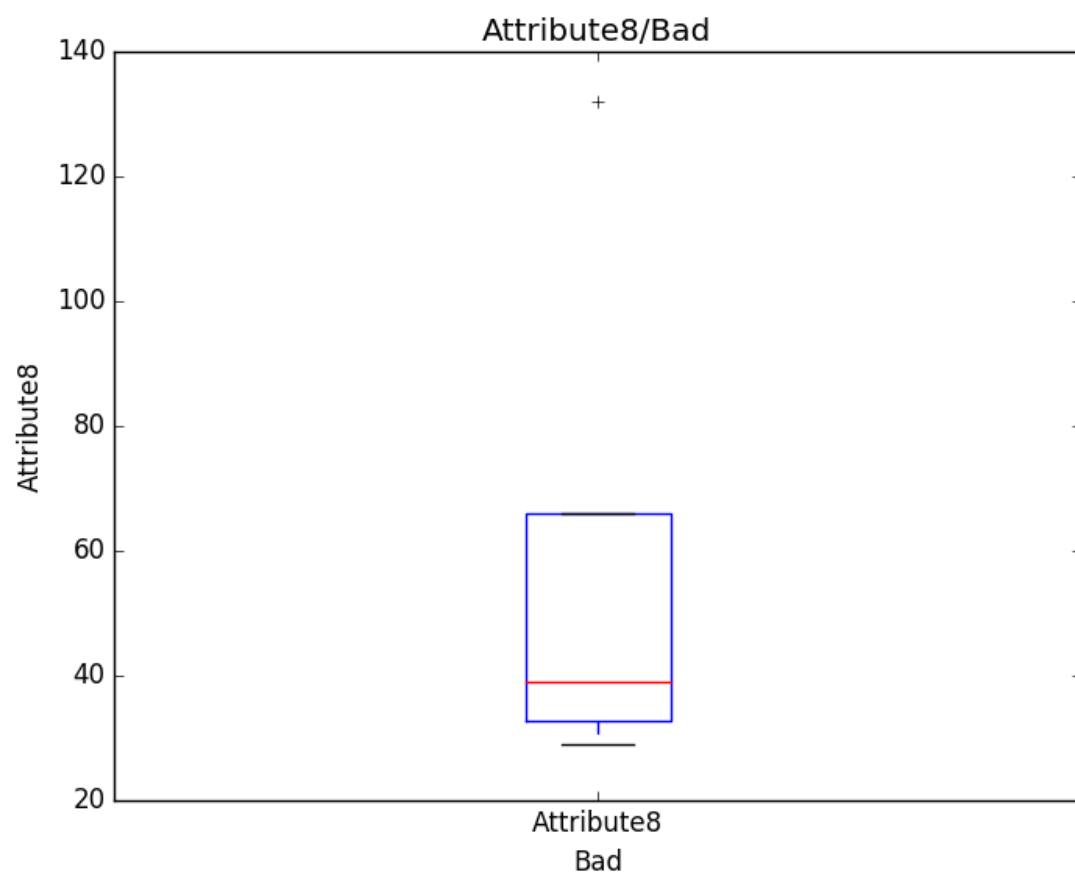




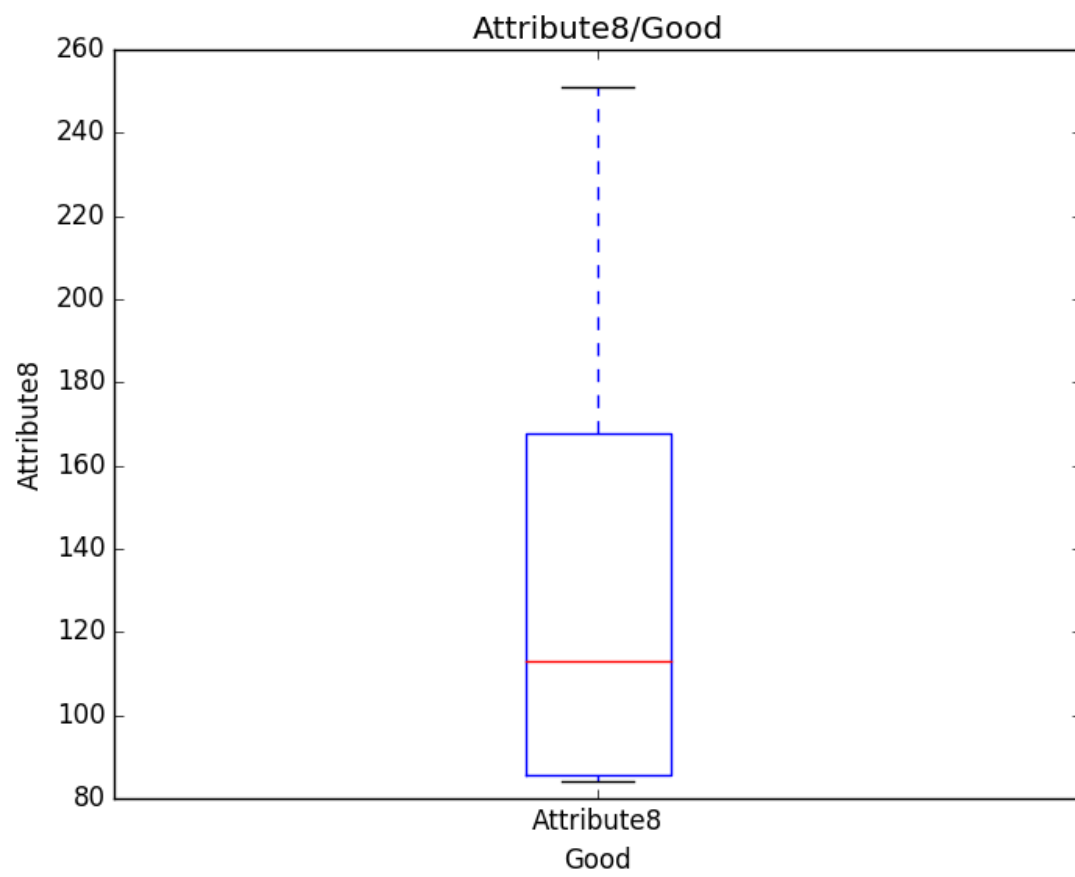


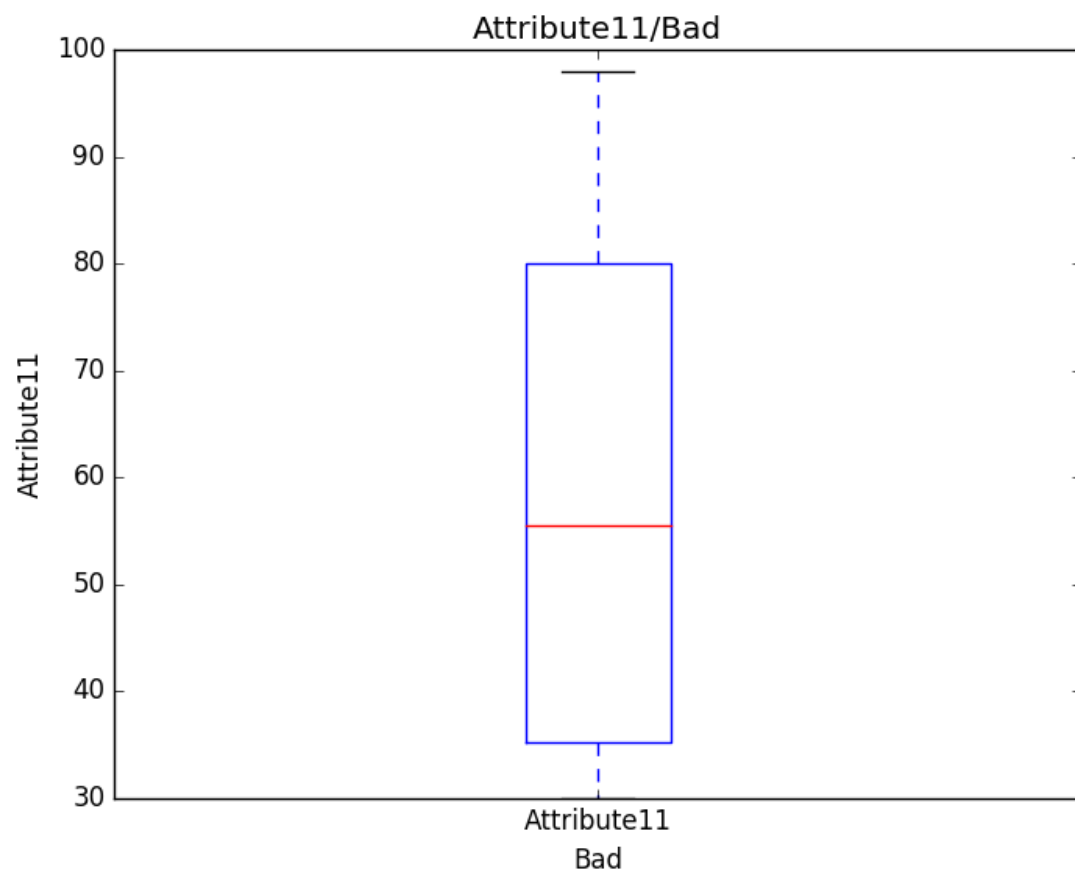


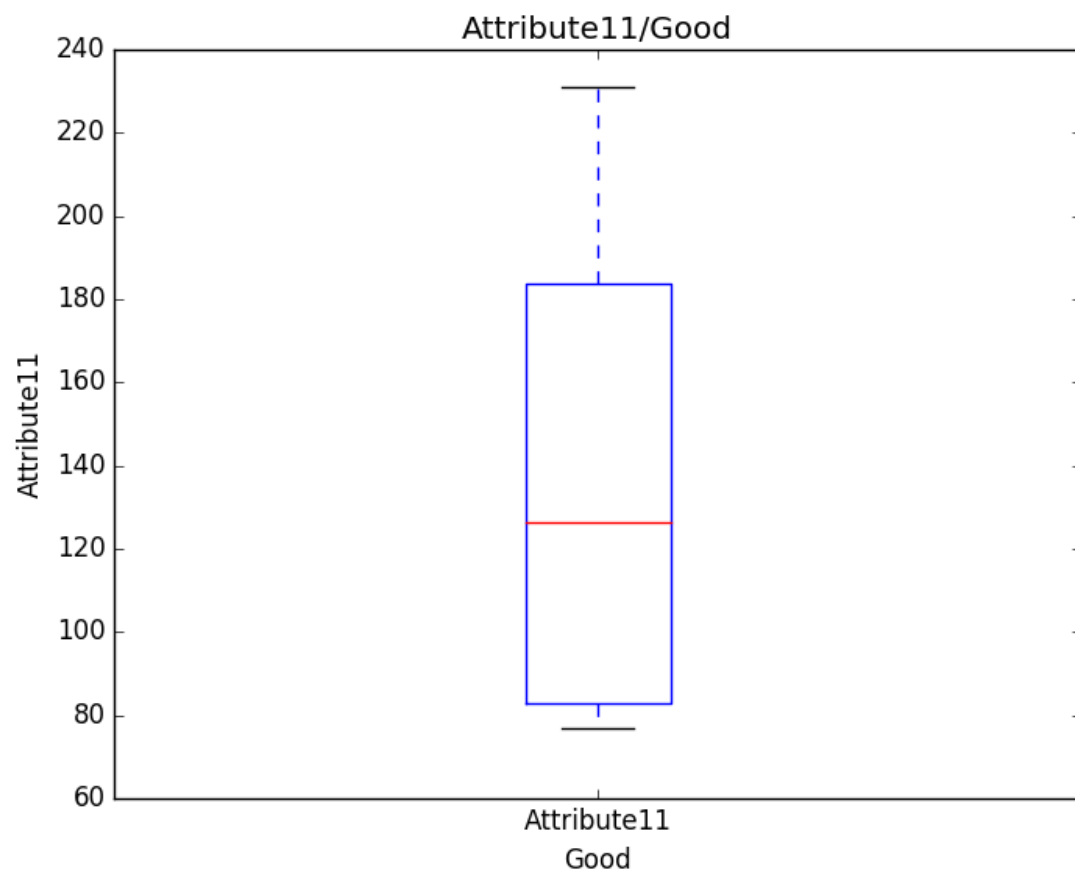


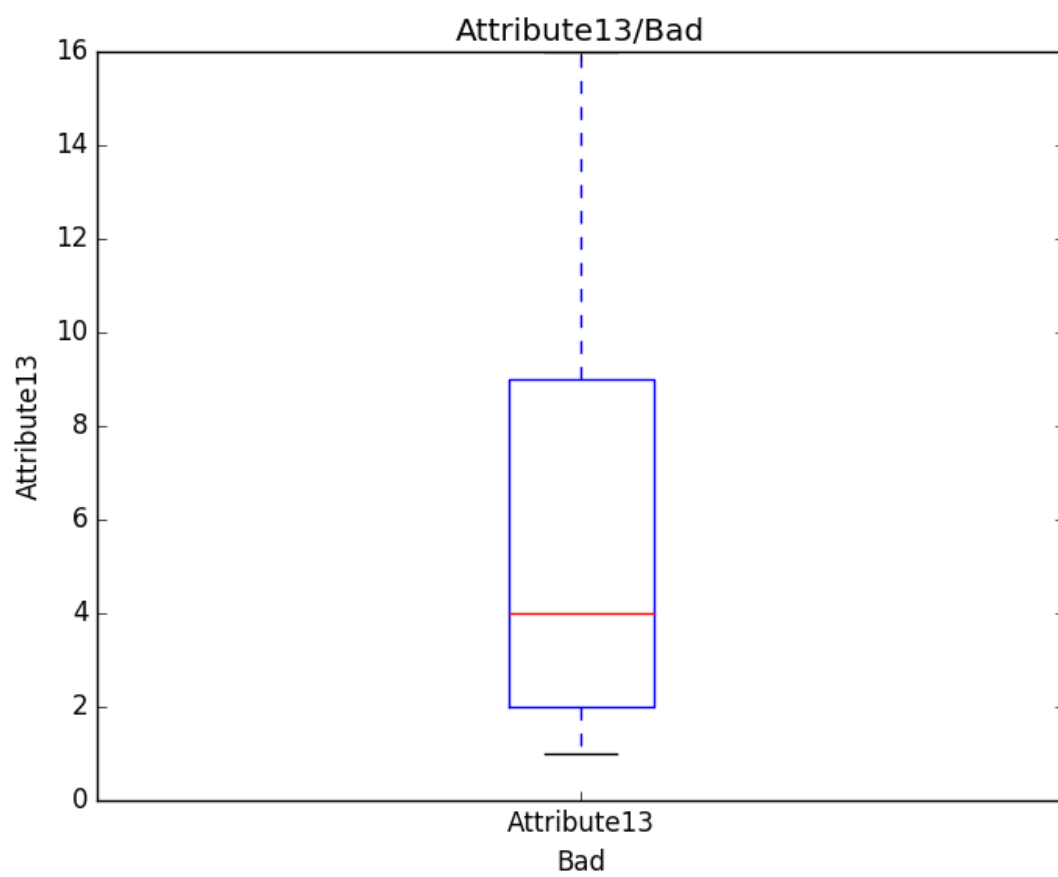


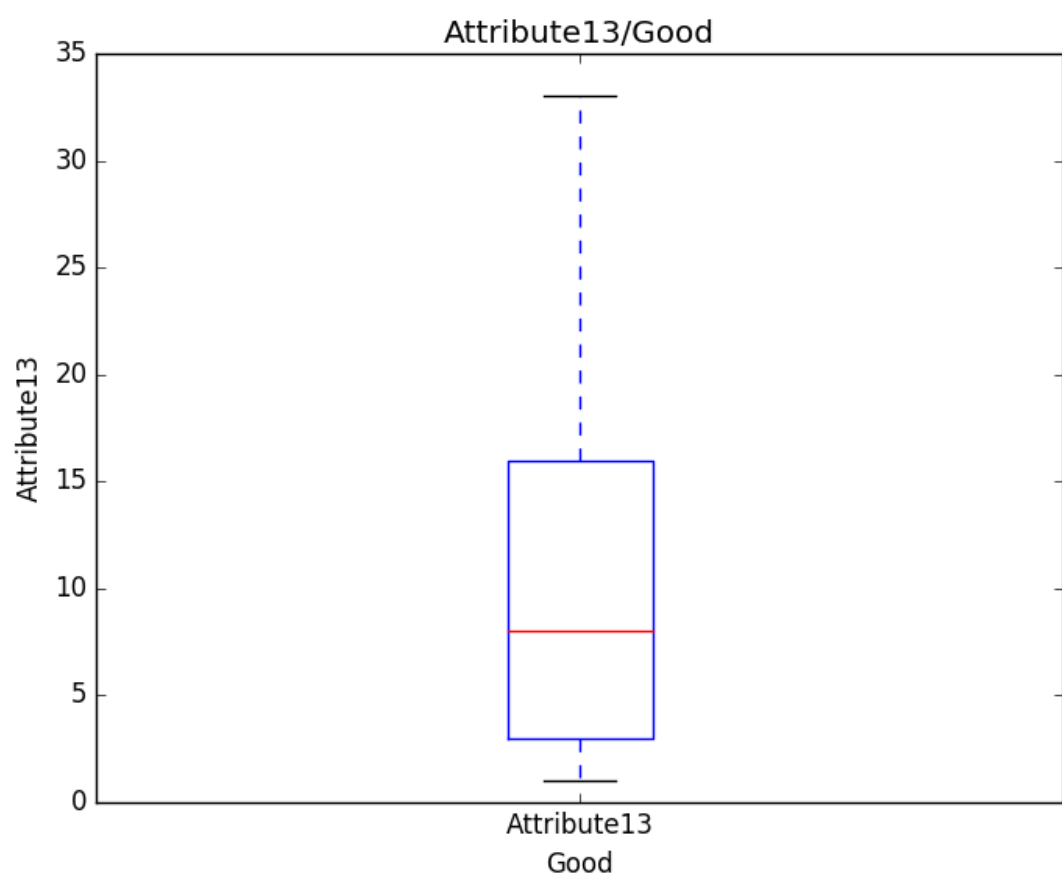


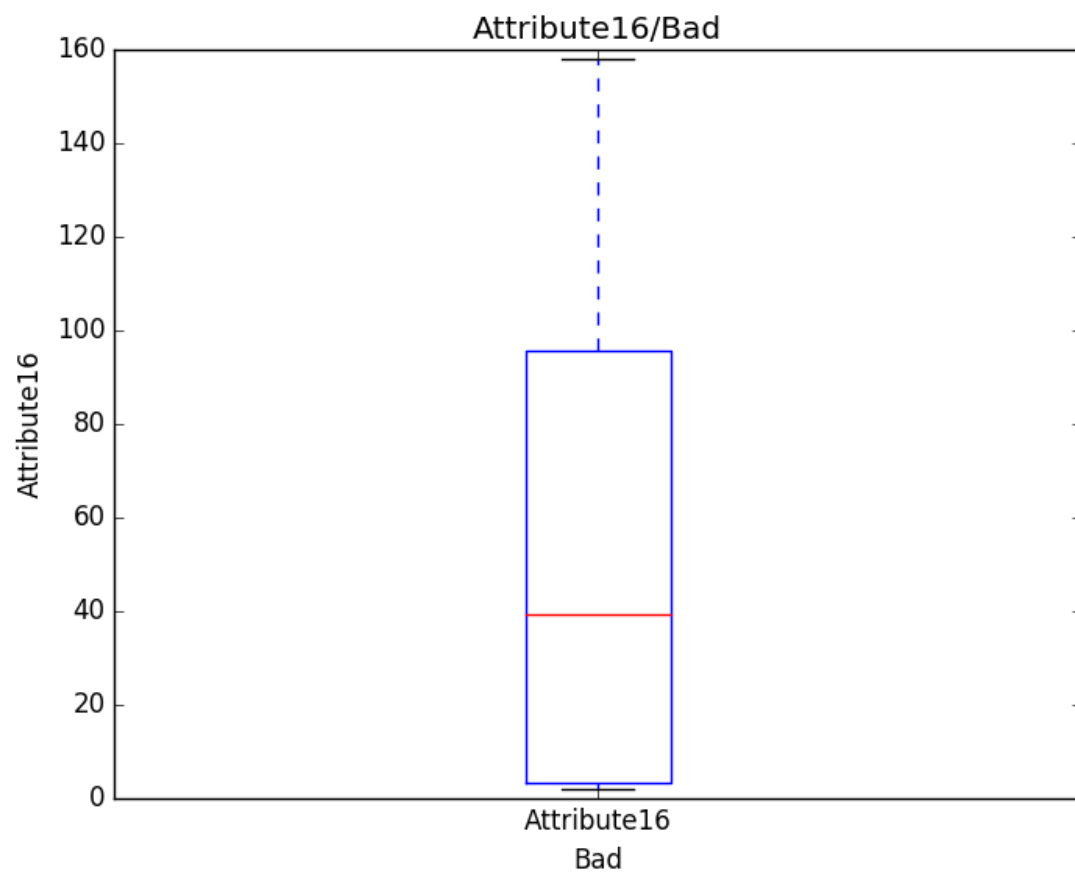


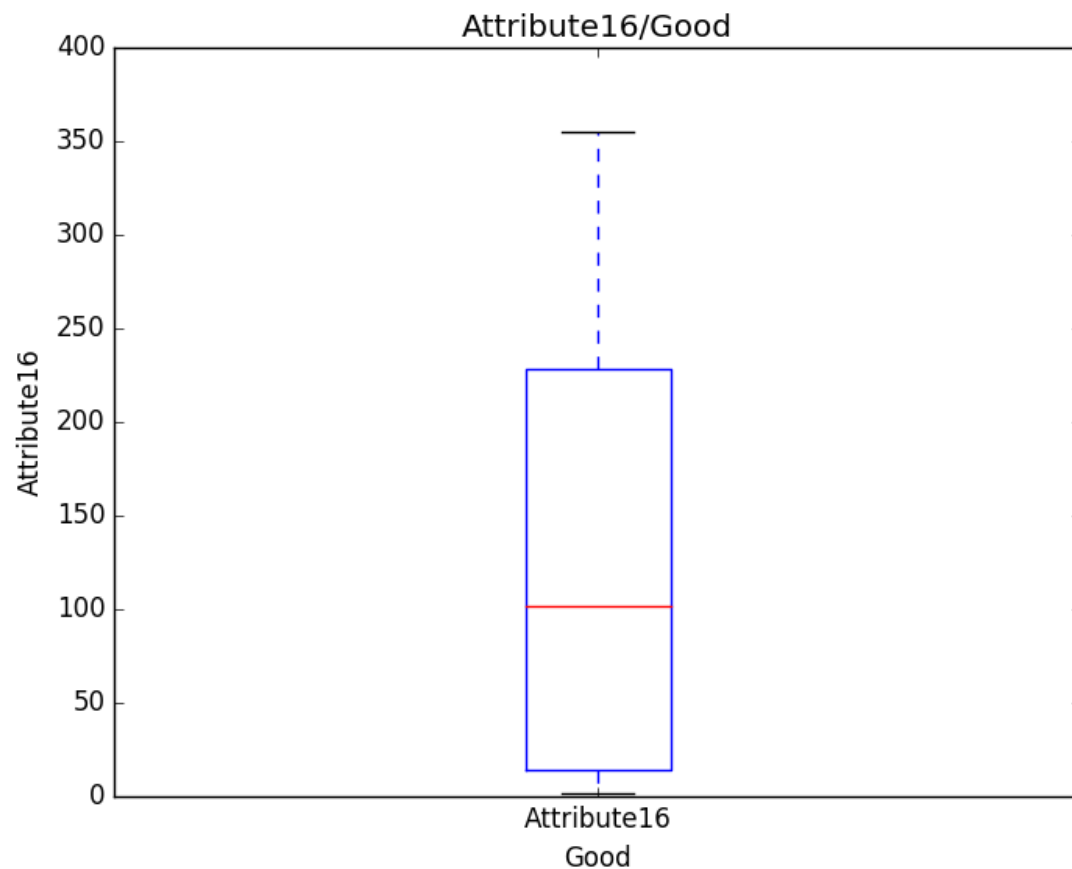












## 2.Υλοποίηση Κατηγοριοποίησης (Classification)

Σε αυτό το ερώτημα υλοποιήσαμε επίσης πλήρως τους αλγόριθμους κατηγοριοποίησης παρατηρώντας accuracy μετά από αρκετές εκτελέσεις περίπου 70-75% και για τους τρεις, με μεγαλύτερη αστάθεια στον random forest περίπου 5%.

Παραθέτεται screenshot του EvaluationMetric10fold(accuracy) και screenshot του test\_setPredictions(δεν αλλαξαμε το 1 σε good και το 2 σε bad):

[illegible]



	A	B	C	D	E	F	G	H	I	J	K	L
1		Client_ID	Predicted_Label									
2	0	10101	1									
3	1	10102	2									
4	2	10103	1									
5	3	10104	2									
6	4	10105	2									
7	5	10106	2									
8	6	10107	1									
9	7	10108	2									
10	8	10109	1									
11	9	10110	2									
12	10	10111	1									
13	11	10112	2									
14	12	10113	1									
15	13	10114	1									
16	14	10115	1									
17	15	10116	1									
18	16	10117	1									
19	17	10118	2									
20	18	10119	2									
21	19	10120	2									
22	20	10121	1									
23	21	10122	1									
24	22	10123	1									
25	23	10124	1									
26	24	10125	1									
27	25	10126	1									
28	26	10127	1									
29	27	10128	1									
30	28	10129	1									
31	29	10130	2									
32	30	10131	1									
33	31	10132	2									
34	32	10133	2									
35	33	10134	1									
36	34	10135	1									
37	35	10136	2									
38	36	10137	2									
39	37	10138	1									
40	38	10139	1									

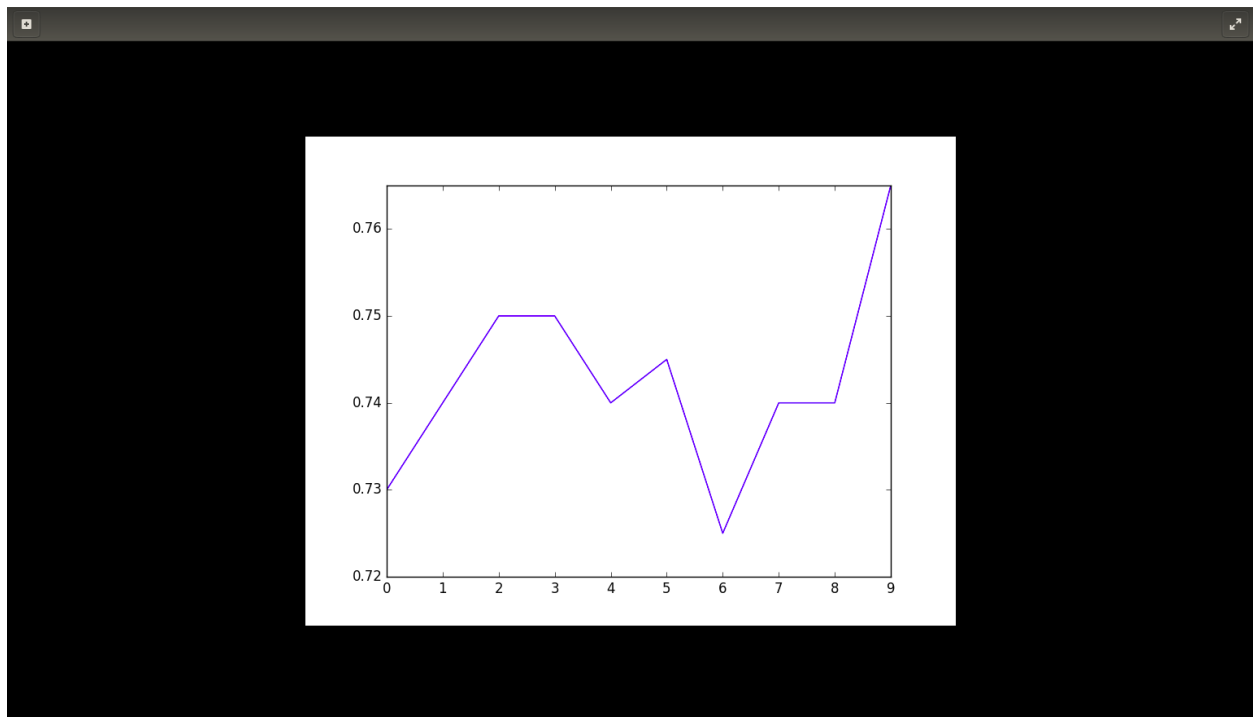
### 3.Επιλογή Features

Σε αυτό το ερώτημα είχαμε το μεγαλύτερο πρόβλημα καθώς πολλά ζητούμενα ήταν ασαφή. Δημιουργήσαμε δικές μας συναρτήσεις για υπολογισμό του entropy ολόκληρου του dataset, υπολογισμό του entropy ενός μόνο Attribute και υπολογισμό του InformationGain ενός Attribute ως προς το dataset. Στη συνέχεια

αφαιρούσαμε το Attribute με το μικρότερο InfoGain και τρέχαμε ξανά τον RandomForest Classifier. Εδώ παρουσιάστηκε το πρόβλημα ότι μετά το 10<sup>ο</sup> Attribute τερμάτιζε, οπότε έχουμε αποτελέσματα για τα 10 χειρότερα ως προς το InfoGain Attributes. Τα αποτελέσματά μας ωστόσο επιβεβαιώνουν τις εκτιμήσεις μας καθώς το accuracy εν τέλει ανεβαίνει. Κανονικά, αν βγάzaμε και τα υπόλοιπα 10 το accuracy θα έπρεπε να πέφτει καθώς μειώνονται τα Attributes και ουσιαστικά ο classifier γίνεται τυχαίος. Σε κάποια σημεία κάνει καμπή, την οποία περιμέναμε, λόγω της αστάθειας του accuracy του αλγορίθμου.

Παραθέτονται το accuracy plot και τα Attributes/InfoGain που αφαιρέσαμε:

(το csv έβαζε το καινούργιο Attribute στην αρχή του αρχείου, οπότε τα στοιχεία είναι ανάποδα. Έχουμε και μια βοηθητική λίστα στην οποία τα στοιχεία έχουν μπει με την σωστή σειρά)



Libération Sans 10 B I U T A

A1 f\_x Σ =

	A	B	C	D	E	F	G	H	I	J
1		Attribute	InformationGain							
2	11	Attribute3	-51.3493661524							
3	12	Attribute9	-51.5402480275							
4	13	Attribute17	-51.6455058487							
5	14	Attribute15	-51.9157988299							
6	15	Attribute16	-51.9398267556							
7	16	Attribute19	-52.0886744223							
8	17	Attribute14	-52.225172854							
9	18	Attribute18	-52.4553405852							
10	19	Attribute10	-52.5342311106							
11	20	Attribute20	-52.8460265277							
12										
13										
14										
15										
16										
17										
18										
19										
20										
21										
22										
23										
24										
25										
26										
27										
28										
29										
30										
31										
32										
33										
34										
35										
36										
37										
38										
39										
40										