

Τεχνικές Εξορυξης Δεδομενων

2η Ασκηση

Φοιτητες:

Όνομα: ΓΚΑΡΑΓΚΑΝΗΣ ΕΥΑΓΓΕΛΟΣ

A.M. : 1115201400033

Όνομα: ΚΩΣΤΑΣ ΚΟΤΡΩΝΗΣ

A.M. : 1115201400074

Το καθε ερωτημα και ζητουμενο της ασκησης υλοποιειται σε 4 αρχεια πηγαιου κωδικα σε pythοn 3.5.

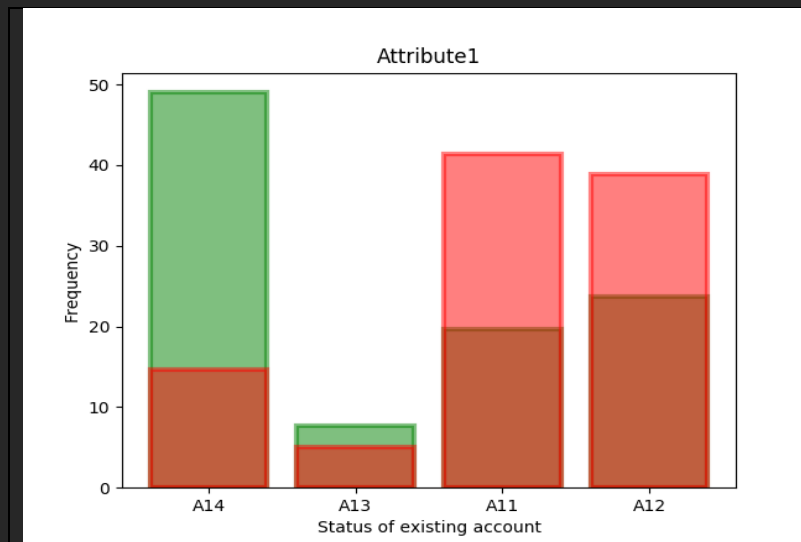
1. Οπτικοποιηση Δεδομενων:
[DataVisualization.py](#)
2. Υλοποιηση Κατηγοριοποιησης:
[Classification.py](#)
3. Επιλογη features (Χαρακτηριστικων):
[FeatureSelection.py](#)
4. TestSet_Predictions:
[Predict.py](#)

Στα αρχεια αυτα υπαρχουν καθοδηγητικα σχολια που εξηγουν την δομη και τη λειτουργια καθε βηματος για την υλοποιηση των ερωτηματων.

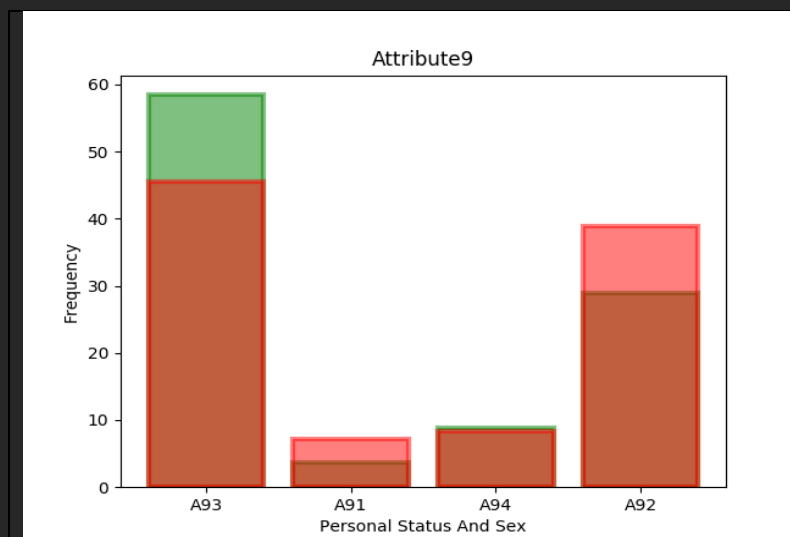
Οπτικοποίηση Δεδομένων (Data Visualization)

Σε αυτό το ερώτημα χρησιμοποιήθηκαν bar plots για την οπτικοποίηση των **categorical** δεδομένων και box plots για **numerical** . Πιο συγκεκριμένα , χρησιμοποιήθηκε η βιβλιοθήκη matplotlib , με την οποία κατασκευάστηκαν τα εξής είδους σχεδιαγράμματα :

Παραδείγματα bar blots



Παραδειγμα histogram , που απεικονίζει το πρώτο feature. Κάθε κατηγορία , για την κατάσταση του τρεχόντος λογαριασμού , κωδικοποιείται όπως μας δίνεται από την άσκηση , σε μια κωδική ονομασία [A14 , A13 , A11 , A12] , το καθένα με την δική του σημασία.

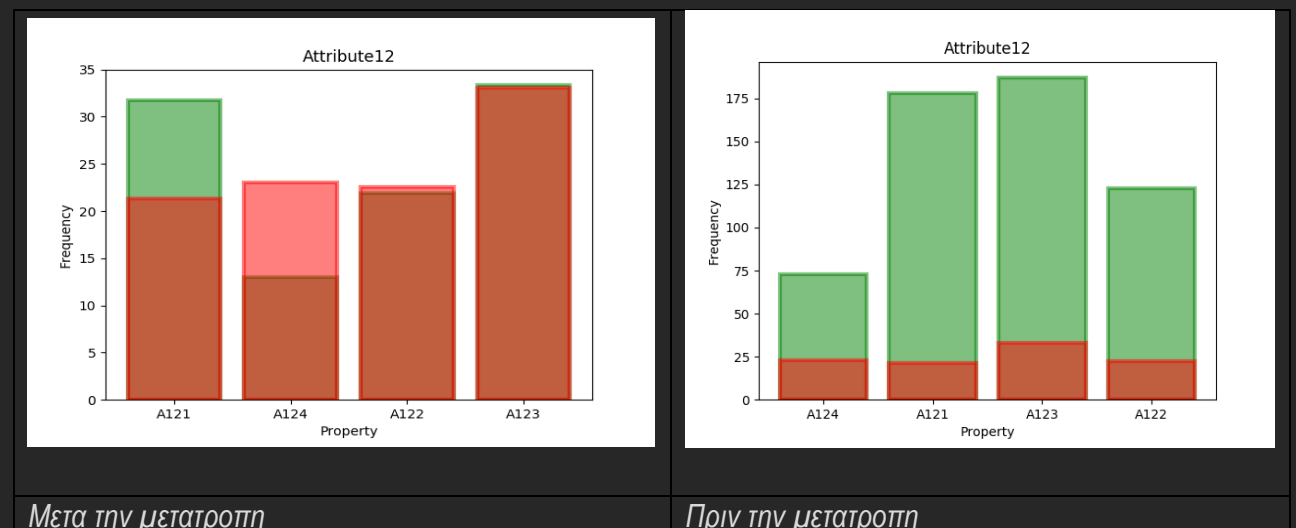


Αντιστοίχα και εδώ. Εξηγούμε:
Attribute 9: (qualitative)
Personal status and sex
A91 : male : divorced/separated
A92 : female : divorced/separated/married
A93 : male : single
A94 : male : married/widowed
A95 : female : single

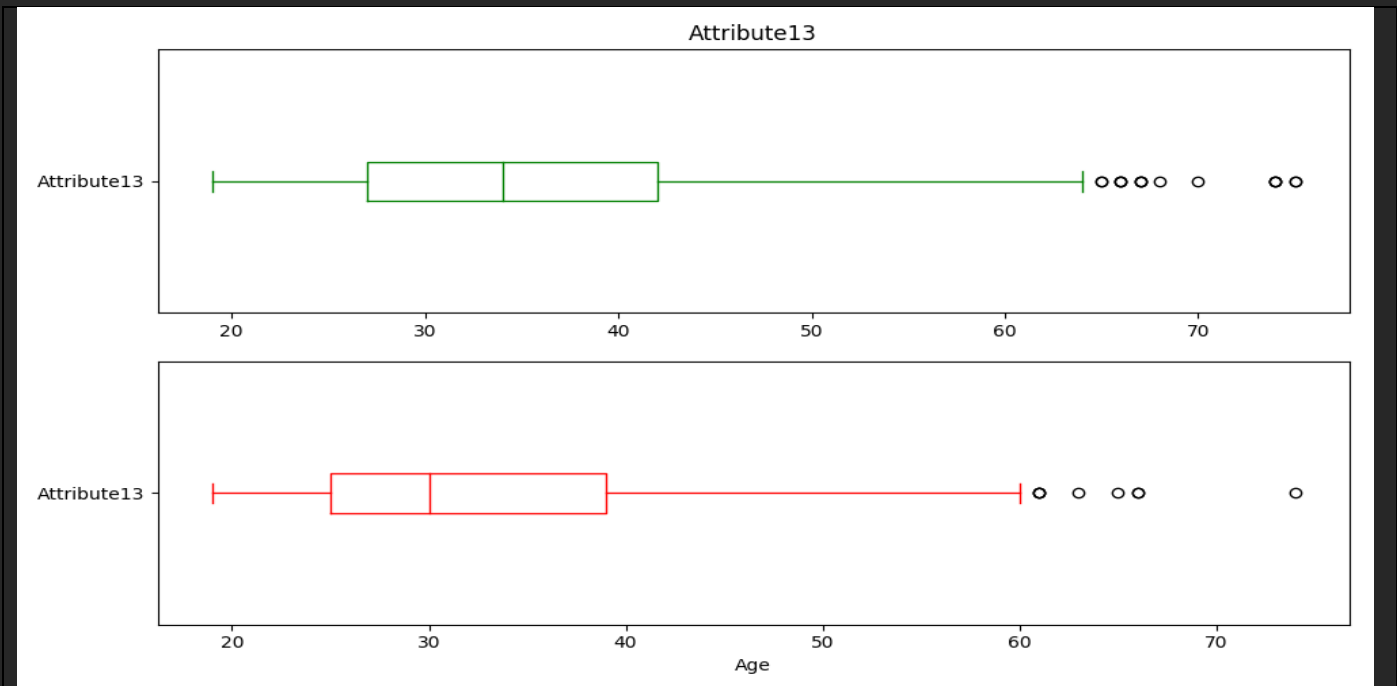
Επεξήγηση σχετικά με την σχεδιαστική επιλογή των bar plot και του τρόπου με τον οποίο αντλούμε πληροφορίες από αυτά :

Στο επάνω μέρος των plot μας , αναγράφεται ο τίτλος του features που οπτικοποιείται . Αριστερά , στον άξονα y , έχουμε την συχνότητα την οποία παρουσιάζει η κάθε κατηγορία για κάθε είδος πελάτη. Στο άξονα x, έχουμε την κατηγορία κάθε κατηγορηματικού feature , ενώ από κάτω επεξηγείται τι ακριβώς συμβολίζει αυτό feature. Ως προς τα plots , με πράσινο χρώμα παριστανονται οι καλοί πελάτες , ενώ με κόκκινο οι κακοί. Το κάθε χρώμα δημιουργείται λόγω του ότι τα bars , πεφτούν το ένα πάνω στο άλλο , το οποίο επελέξαμε για λόγους ευδιακριτής σύγκρισης.

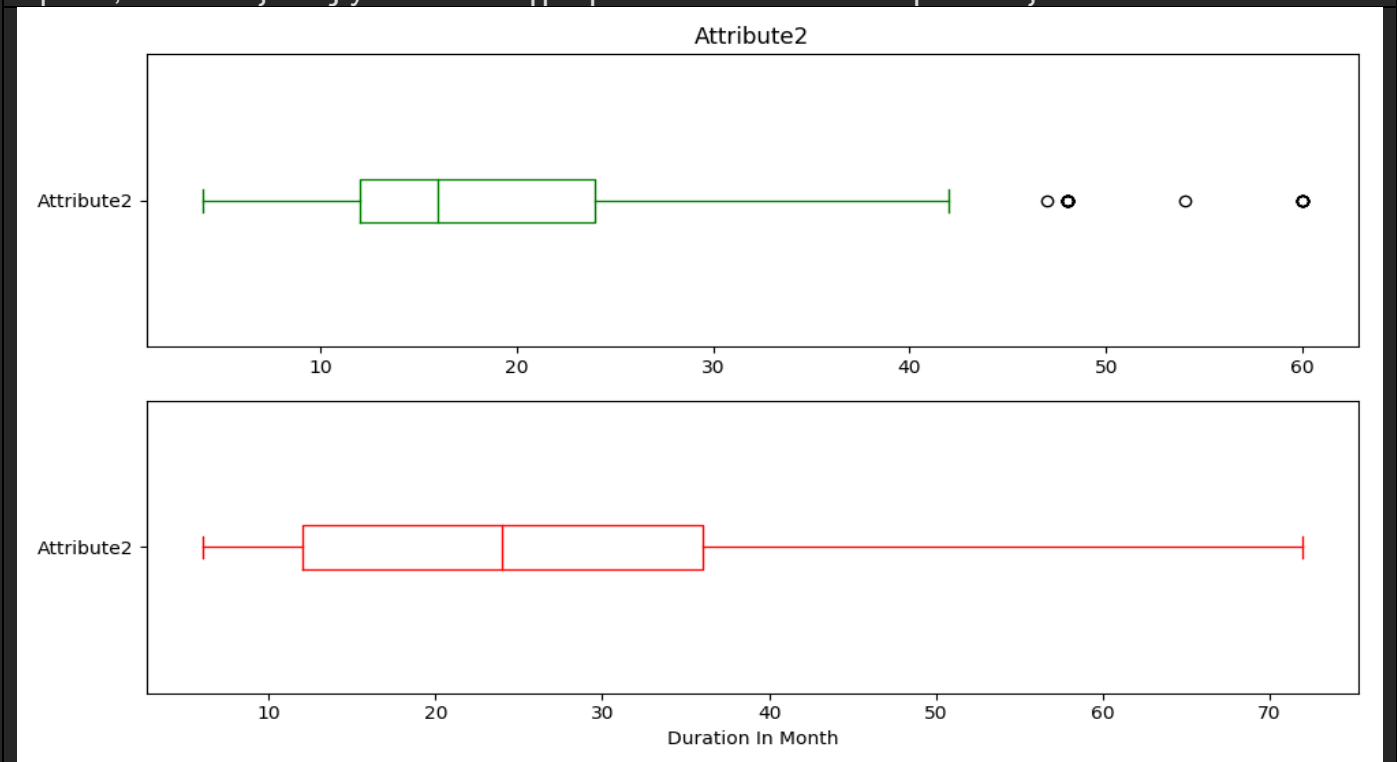
Τα bars δημιουργήκαν με βάση την παρακάτω λογική. Χωρίσαμε το train data set σε καλούς και κακούς πελάτες. (~561 good , 239 bad). Στην συνέχεια για το κάθε είδος συνόλου πελατών , υπολογίσαμε την συχνότητα της κάθε κατηγορίας του feature σε αυτήν. Στην αρχή , επιλέξαμε να προβάλουμε αυτά τα αποτελέσματα αυτούσια στα γραφήματά μας , δηλαδή το ύψος κάθε κατηγορίας να αγγίζει τον ακριβή αριθμό ατόμων που ανήκουν σε αυτή. Αυτό όμως , δεν μας παρέδιδε ξεκαθαρά κάποια εμφανής πληροφορία , αφού οι πλειοψηφία των πελατών μας είναι good . Έτσι επιλέξαμε , να υπολογίζουμε το ποσοστό τις 100 , των πελατών που ανήκουν σε κάθε κατηγορία , για το κάθε σύνολο πελατών , προκειμένου να έχουμε πιο αντιπροσωπευτικά αποτελέσματα. Σύγκριση των bar plot , πριν και μετά την μετατροπή %.



Αντιστοιχη επεξηγηση των box plots



Τα box plots παρουσιάζουν αντιστοιχα , ακολουθοντας τον ιδιο διαχωρισμο των δεδομενων , την αντιστοιχη συγκριση των χαρακτηριστικων μεταξυ των καλων και των κακων πελατων , αυτη την φορα για τα νουμερικα χαρακτηριστικα των πελατων. Εδω , δεν εχουμε % αναπαρασταση των δεδομενων , αφου τα box bots , μας δειχνουν που υπαρχει η μεγαλυτερη συγκεντρωση ατομων . Ο αξονας x περιεχει το ευρος των τιμων , ενω ο αξονας y απλα αναγραφει το feature που παρουσιάζεται.



Παραδειγματα box plots

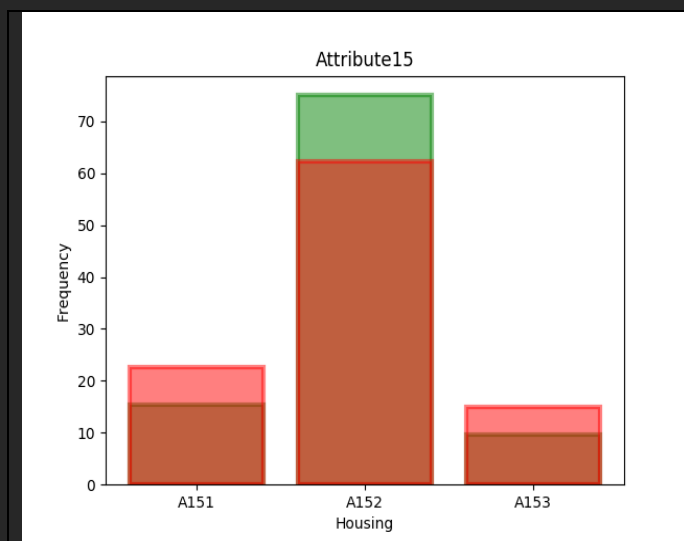
Παρατηρήσεις / Προβλεψεις πανω στα plots

Τα δεδομένα που οπτικοποιήθηκαν επαληθευσαν , για το μεγαλύτερο ποσοστό , τις προβλεψεις μας , καθώς οι κατηγορίες που περιμέναμε να είναι πιο σημαντικές και στις οποίες ήταν αναμενόμενο να δούμε διαφορές ανάμεσα σετους καλούς και κακούς πελάτες , πραγματι είχαν αντιστοιχα γραφηματα. Τα features ανταποκρινονται στην πραγματικοτητα , εξαιρουμενου καποιον εκπληξεων , και παρουσιαζουν μια πολυ καλη εικονα για το ποια χαρακτηριστικα διακρινουν εναν καλο και ενα κακο πελατη.

Πιο συγκεκριμενα , μελετηθηκαν ολα τα features ξεχωριστα , τοσο ως προς τις δικες μας προβλεψεις και εκτιμησεις , οσο και ως προς τις παρατησεις και τα συμπερασματα που απορρεουν απο τα γραφηματα. Ετσι διακρινουμε τα features αναλογα με την σημασια τους για το ποιον ενος πελατη και για το ποσο αναμενομενη ηταν η προβλεψη μας. Τα παρουσιαζουμε παρακατω :

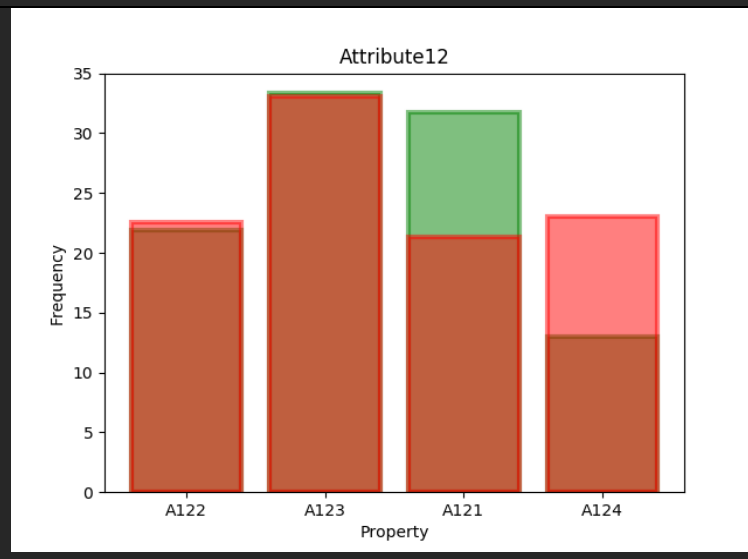
○ **Σημαντικα features :**

Attribute 15 , Housing :



Αναμενομενο . Είναι ακρως λογικο μια τραπεζα να θεωρει την κατοχη σπιτιου καλο στοιχειο για ενα πελατη , λογω της σταθεροτητας ενεπνεει για την ζωη του και για την παροχη εγγυησεων που ισως τις παρεχει , εναντι ενος πελατη που δεν εχει δικο του σπιτι ή μενει καπου δωρεαν. Αυτο το αποδुकνυει και το αντιστοιχο ιστογραμμα για το συγκεκριμενο feature.

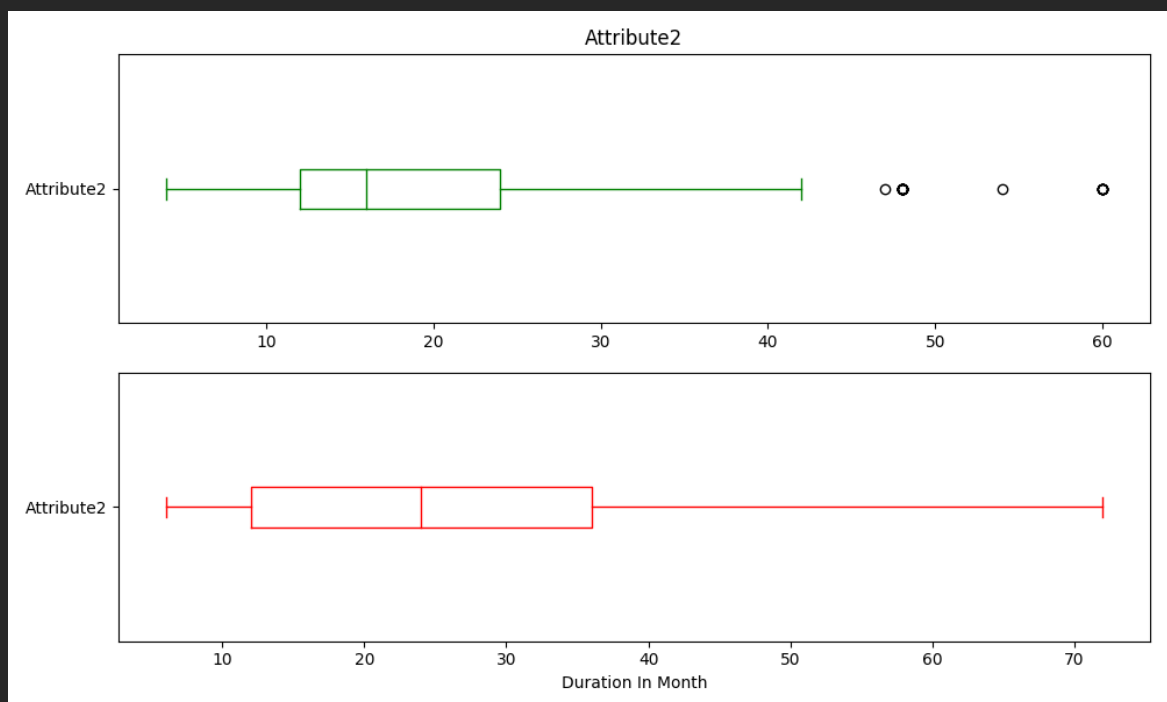
Attribute 12 , Property :

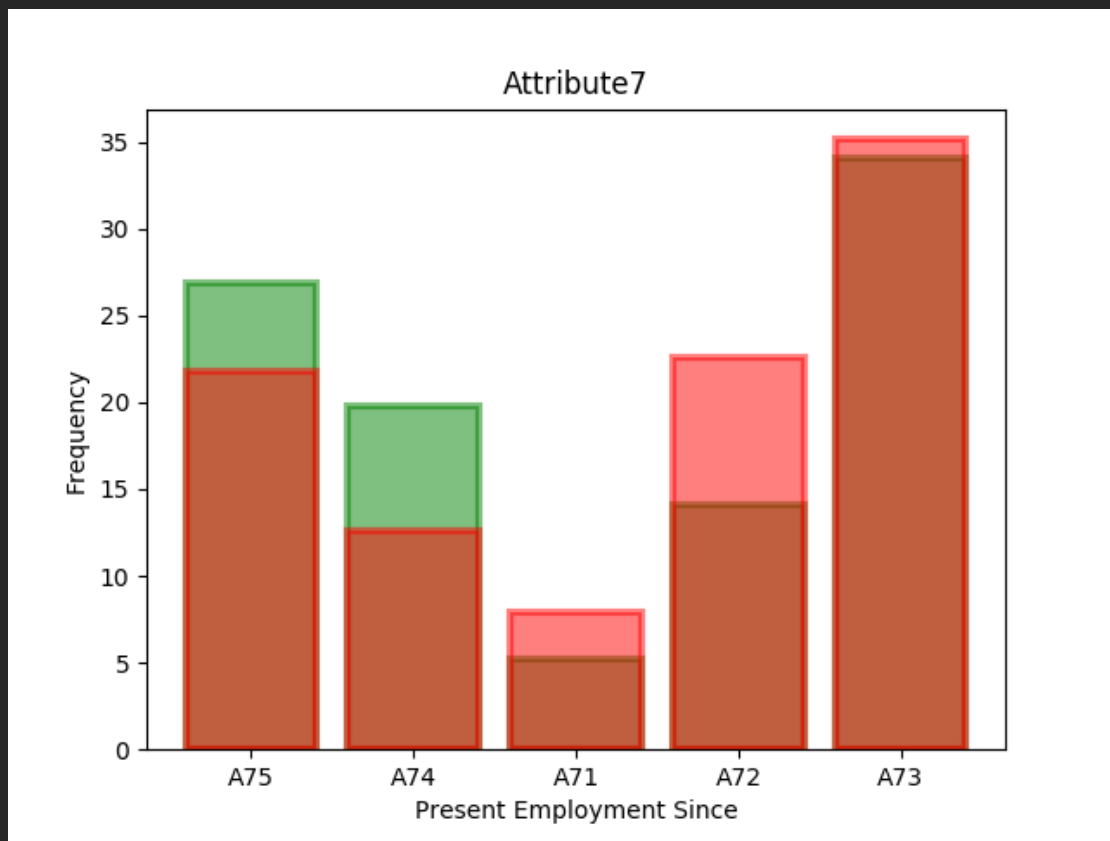


Ανεμενομενο και επίσης λογικο. Πελατες οι οποιοι φροντιζουν για το μελλον τους με ασφαλειες ζωης και οικονομίες για κοινωνικες αναγκες , ειναι πιο ωριμοι ως προς την διαχειριση των χρηματων του , οποτε οι τραπεζα τους θεωρει καλους. Αντιθετως , με τους κακους πελατες , που χρησιμοποιουν τις αναγκες για υλιστικες αγορες

Attribute 2:Duration in month, Attribute 7: Employment Since:

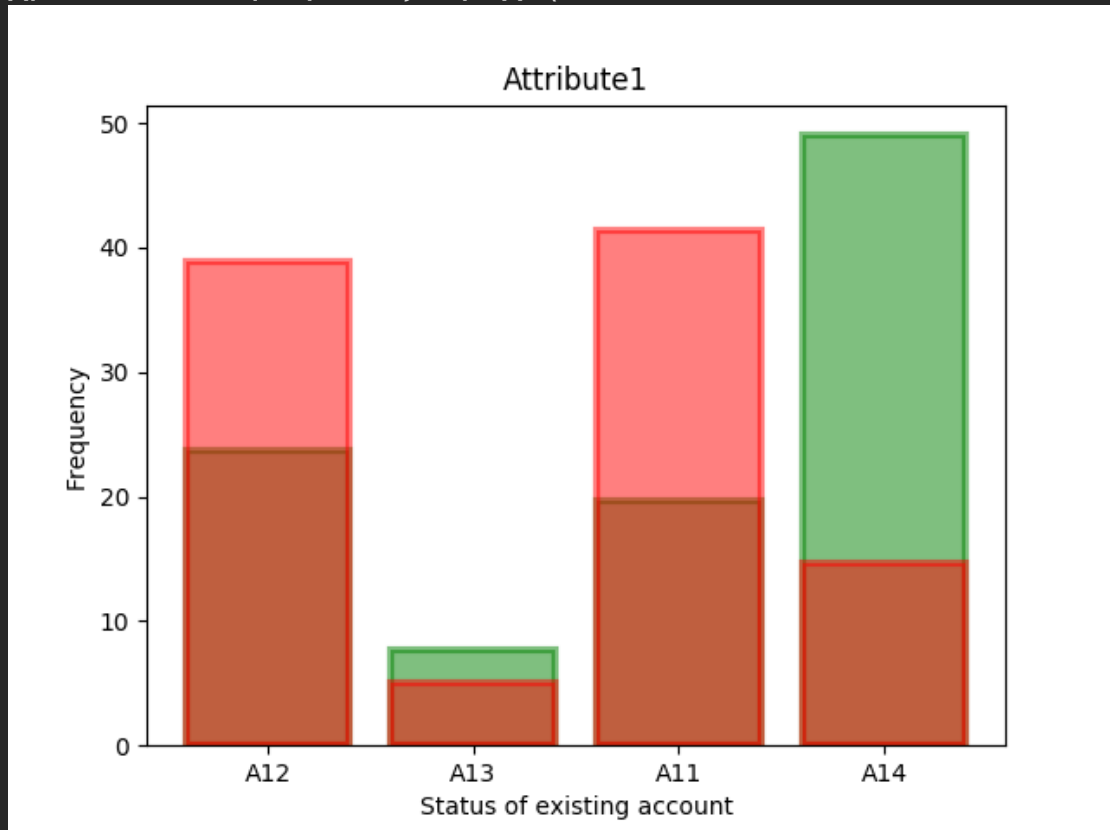
Τα χαρακτηριστικα σταθεροτητα και ωριμοτητα συνεχιζουν να διακρινουν τους καλους πελατες . Οπως βλεπουμε απο τα παρακατω γραφηματα :





Συνεπεις πελατες με μικρα δανεια , σταθεροι ως προς την ζωη τους αφου δουλεουν πανω απο 4 χρονια , ειναι λογικο να ειναι σημαντικοι.

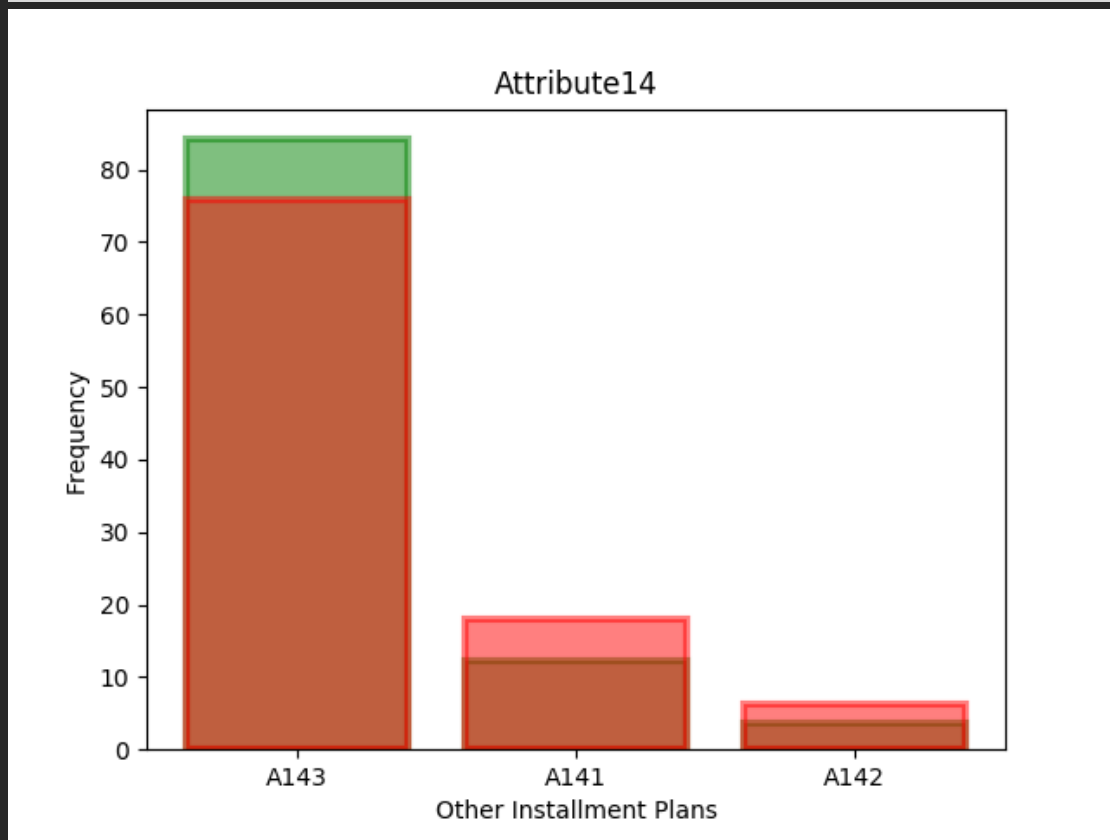
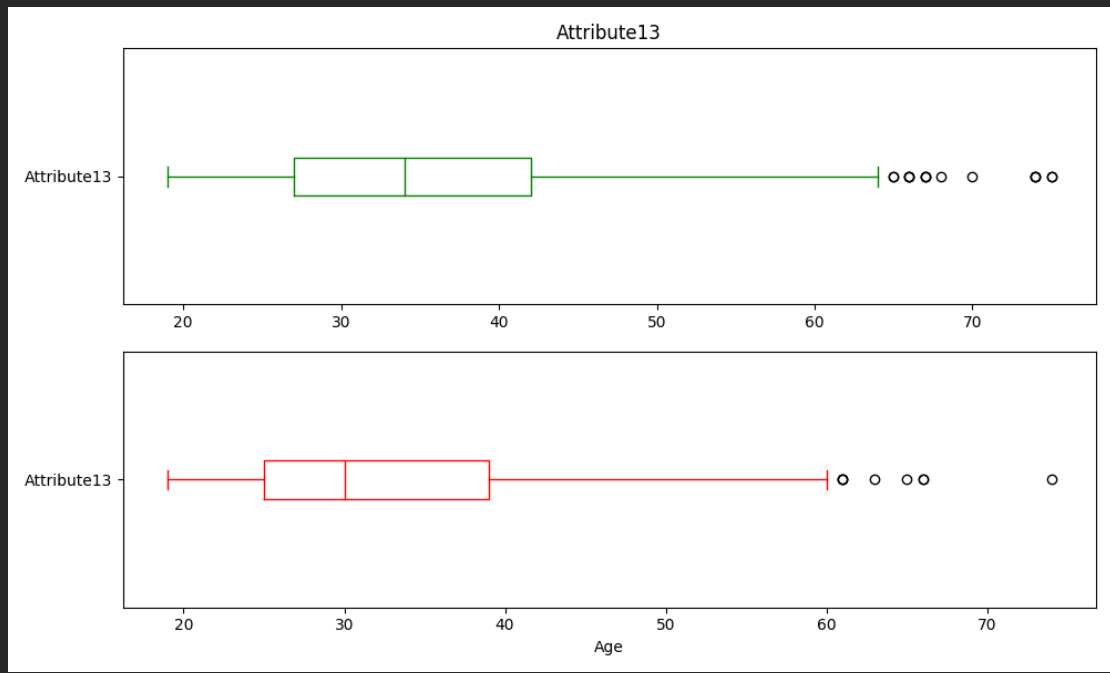
Αντιστοιχη λογικη ακολουθουν και αλλα features οπως ειναι η κατασταση του τρεχοντος λογαριασμου (Att 6), εαν χρωστουν στην τραπεζα ή οχι (**Attribute 1** :



) , των οποιων τα γραφηματα μιλουν απο μονα τους.

ο Μετρια features :

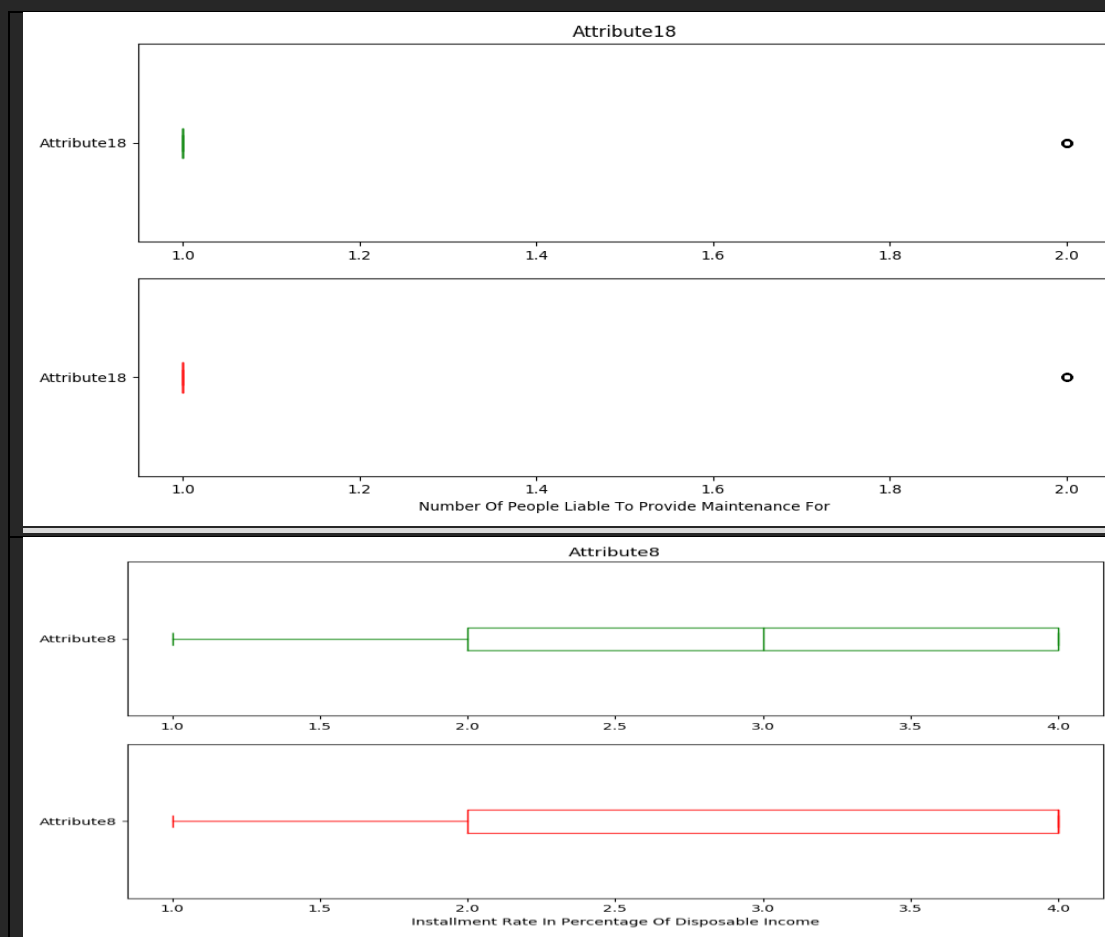
Θα αναφερουμε να χαρακτηριστικα που εχουν μετρια σημασια. Attr 14 , τροπος εξοφλησης , Attr 13 η ηλικια , Attr 15 τα χρηματα στην τραπεζα , Attr 4 , ο σκοπος . Ολα αυτα ειναι πολυ σχετικες ιδιοτητες που με την λογικη και μονο δεν εχουν καποια αμεση σχεση με το αν εισαι καλος πελατης ή οχι. Παρατηρουμε οτι πελατες με μεγαλυτερη ηλικια , με μεσαια εισοδηματα στην τραπεζα , η δουλεια τους (Attr 17) , δεν επηρεαζει αμεσα την κατασταση. Αυτο φαινεται και στα διαγραμματα , με την μετρια διαφορα τους .



○ Ασημαντα features :

Ασημαντα ειναι τα features τα οποια εχουν απειροελαχιστη διαφορα στα plot τους , πραγμα που δειχνει

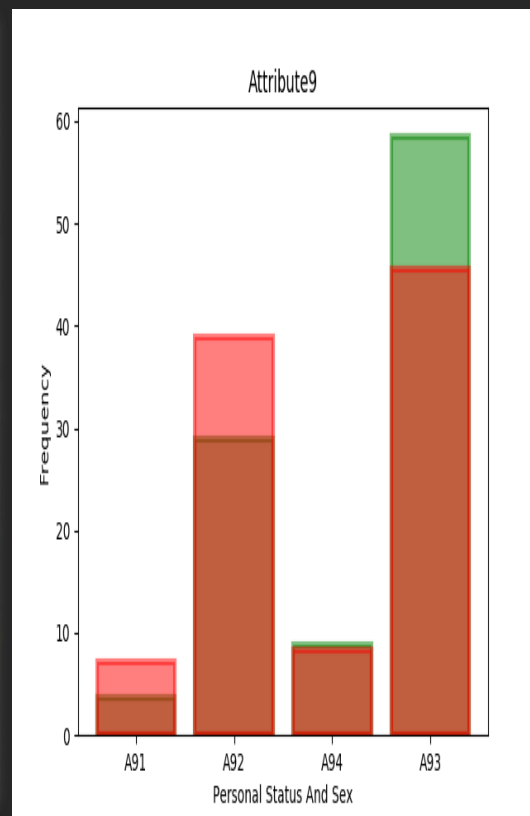
ΟΤΙ ΟΙ ΠΕΛΑΤΕΣ ΤΩΝ ΣΥΓΓΕΚΡΙΜΕΝΩΝ feature δεν παρουσιάζουν κάποια διαφορά που να παίζει ρόλο για το αν είναι καλοί πελάτες ή όχι. Τέτοια features είναι Attr 19 , Attr 18 , Attr 8 , Attr 20 κλπ. Δηλαδή η καταγωγή του πελάτη , το αν έχει τηλέφωνο ή όχι , το ποσο λογριασμούς έχει στην τραπεζα , δηλαδή αυτά τα feature που απλά είναι θέμα επιλογής / τυχαιότητας για τον πελάτη, δεν έχουν τόσο σημασία. Χαρακτηριστικά είναι τα γραφήματα των Attributes 18 και 8 , με τα γραφήματα τους να μην παρέχουν κάποια προφανή πληροφορία ή κάποιο πορίσμα.



Συμπεραινουμε λοιπον , πως τα στοιχεια που ειναι σημαντικα για την διακριση ενος πελατη σε καλο ή κακο , ειναι αυτα που μας λεει και η λογικη. Αυτα ειναι η ωριμοτητα , η

συνεπεια του , το ποσο φροντιζει για το μελλον του , η σταθεροτητα της ζωης του , η επιλογες του κτλπ. Στοιχεια οπως η τρεχον κατασταση του πελατη, οπως ειναι το αν εχει τηλεφωνο ή οχι , τα λεφτα που εχει στην τραπεζα , το ποσο χρονων ειναι , δεν παιζουν τοσο σημαντικο ρολο. Ασημαντα χαρακτηριστικα θεωρουντα αυτα για τα οποια δεν ειναι θεμα επιλογης του πελατη, οπως η καταγογη , και αυτα που δεν εχουν τοσο σημαντικο ρολο.

Fun Fact:



Οπου A93 = Single

Υλοποίηση Κατηγοριοποίησης (Classification)

Για την υλοποίηση της κατηγοριοποίησης χρησιμοποιήθηκε η scikit με ήδη γνωστούς αλγορίθμους , που χρησιμοποιήθηκαν και σχολιαστήκαν εκτενώς στην προηγούμενη άσκηση . Σε αυτό το ερώτημα , χρειάστηκε αναλογία προεπεξεργασία των δεδομένων , προκειμένου να μετατραπουν όλα τα κατηγορηματικά feature σε νουμερικά .

(Η προεπεξεργασία και μετατροπή αυτών των δεδομένων έγινε με την βοήθεια του dummy coding . το οποίο χρησιμοποιείται συνήθως για τη μετατροπή μιας μεταβλητής κατηγορικής εισόδου σε συνεχή μεταβλητή. Το "Dummy", όπως υποδηλώνει το όνομα, είναι μια διπλή μεταβλητή που αντιπροσωπεύει ένα επίπεδο μιας κατηγορικής μεταβλητής. Η παρουσία ενός επιπέδου είναι αντιπροσωπευτική κατά 1 και η απουσία αντιπροσωπεύεται από το 0. Για κάθε υπάρχον επίπεδο, θα δημιουργηθεί μια ψευδομεταβλητή.)

. Η μέθοδοι classification που προτιμήθηκαν είναι οι Random Forest , Naïve Bayes και SVM , πιο συγκεκριμένα ο linear svm λόγο της καλύτερης απόδοσης του. Οι αλγόριθμοι επαληθεύτηκαν με 10 fold cross validation , αμε αντιστοιχες συναρτήσεις που παρέχονται από την sci-kit.

Οι παρακάτω πληροφορίες βασίστηκαν πάνω στις εξής πηγές :

- <https://www.analyticsvidhya.com/blog/2015/11/easy-methods-deal-categorical-variables-predictive-modeling/>
- <https://www.quora.com/Which-algorithm-fits-best-for-categorical-and-continuous-independent-variables-with-categorical-response-in-Machine-Learning>
- <https://users.cs.duke.edu/~rvt/msthesis.pdf>

Ακολουθούν τα αποτελέσματα του evaluation metric μας:

	A	B	C	D	E
1	Statistic Measure	Naive Bayes	Random Forest	SVM	
2	Accuracy	0.61375	0.7287500000000001	0.6637500000000001	
3					

Static Measure	Naïve Bayes	Random Forest	SVM
Accuracy	0.61375	0.72875	0.66375

Κατοπιν πολλων προσπαθειων και δοκιμων αλγοριθμων κατηγοριοποιησης πανω στα δεδομενα , οπως παρατηρουμε και στο αρχειο εξοδου , με εμφανη διαφορα καλυτερος ειναι ο Random Forest Classifier. Η διαφορα αυτη , στην ικανοποιητικη αποδοση του αποδιδεται στα δεντρα αποφασης που δημιουργει ο random forest , πραγμα που καθοδηγει τον αλγοριθμο σε μονοπατια μεγαλυτερης αποκτησης πληροφοριας οπως θα δουμε και παρακατω. Παρα το γεγονος οτι ο Random Forest ειναι ισως ‘κουραστικος’ για την μηχανη στην οποια εκτελειται , αφου αυξανουν αρκετα την υπολογιστη πολυπλοκοτητα προκειμενου να δημιουργησουν τα δεντρα , οπως αναφερεται και σε αντιστοιχο αρθρο του Wikipedia:

“For data including categorical variables with different number of levels, random forests are biased in favor of those attributes with more levels”

Οσον αναφορα τους υπολοιπους αλγοριθμους φαινεται επισης να ειναι σχετικα ικονοποητικοι , και δουλεουν αρκετα καλα με τα categorical data , οσο και με τα numerical.

Τα αποτελέσματα των μετρήσεων μας ανταποκρίνονται στην πραγματικότητα και στο κατά ποσο ευκολο είναι να βγάλουμε συμπεράσματα πάνω στο συγκεκριμένο data set που μας δίνεται . Όπως παρατηρήσαμε στη παραπάνω γενική ιδέα που μας προσέφερε η οπτικοποίηση των δεδομένων , και όπως επίσης θα καταλάβουμε στην συνέχεια από το feature selection και το ποσο πληροφορία εκλαμβάνουμε από τα feature μας , καταλαβαίνουμε ότι η κατηγοριοποίηση δεν είναι αρκετά ευκολή. Αυτό συμβαίνει λόγω των χαμηλών συγκρίσεων που μας παρέχει το data set , αφού οι κακοί και οι καλοί πελάτες δεν έχουν πάντα εμφανής διαφορές και δεν ακολουθούν κάποια πολύ συγκεκριμένα μοτίβα και συμπεριφορές. Αυτό αποδεικνύεται και παρακάτω , με τον υπολογισμό του information gain του κάθε feature.

Τα παραπάνω , πέρα από προσωπικές μας λογικές εκτιμήσεις και υπολογισμούς , φαίνεται να αποδεικνύονται και παρακάτω:

What is Imbalanced Data?

Imbalanced data typically refers to a problem with classification problems where the classes are not represented equally. For example, you may have a 2-class (binary) classification problem with 100 instances (rows). A total of 80 instances are labeled with Class-1 and the remaining 20 instances are labeled with Class-2. This is an imbalanced dataset and the ratio of Class-1 to Class-2 instances is 80:20 or more concisely 4:1.

Το data set μας , με ratio περίπου 1:3 , εκφράζει μια σχετική ασταθία , η οποία εκδηλώνεται τόσο στα plots μας , όσο και στις μεθόδους classification και το information gain που μας παρέχουν

Οι αλγόριθμοι που χρησιμοποιούμε και ειδικά ο random forest , θεωρείται πως έχει υψηλά επιπεδα ευστοχίας για το συγκεκριμένο data set , αφού μονίμα μας παρέχει ένα 70% **ΕΥΣΤΟΧΙΑΣ**.

Πηγή: <http://machinelearningmastery.com/tactics-to-combat-imbalanced-classes-in-your-machine-learning-dataset/>

Επιλογή Feature (Feature Selection)

Για την επιλογή των feature , το ερωτημα αναλυθηκε σε 2 σταδια και συνδυαστηκε στην συνεχεια ,οποτε θα εξηγησουμε το καθενα τοσο ξεχωριστα , οσο και συνδυαστικα στην συνεχεια.

- Υπολογισμος Information Gain για καθε feature:

Επειδη το ερωτημα αυτο υλοποιηθηκε κατοπιν εκτενης ερευνας και κατασκευαστηκε απο εμας , με δικες μας υλοποιησεις , θα ξεκινήσουμε με τον ορισμο του information gain.

This measure of *purity* is called the **information**. It represents the expected amount of information that would be needed to specify a prediction.

Entropy on the other hand is a measure of *impurity* (the opposite).

Η αποκτηση της πληροφοριας ειναι αλληλενδετη με την εντροπια , οπως εξηγηται και παραπανω. Ο μαθηματικος τυπος της εντροπιας ειναι ο παρακατω :

$$H(S) = - \sum_{x \in X} p(x) \log_2 p(x)$$

Where,

- S - The current (data) set for which entropy is being calculated (changes every iteration of the ID3 algorithm)
- X - Set of classes in S
- $p(x)$ - The proportion of the number of elements in class x to the number of elements in set S

Η συνδεση της εντροπιας με την πληροφορια εξηγηται στην παρακατω συλλογιστικη πορεια.

- is a measure of the amount of uncertainty in the (data) set S (i.e. entropy characterizes the (data) set S).
- in other words, it is the average amount of information contained in each message received (message here stands for an event, sample or character drawn from a distribution or data stream)
- it characterizes the uncertainty about our source of information (Entropy is best understood as a measure of uncertainty rather than certainty, as entropy is larger for more random sources)
- a data source is also characterized by the probability distribution of the samples drawn from it (the less likely an event is, the more information it provides when it occurs)
- it makes sense to define information as the negative of the logarithm of the probability distribution (the probability distribution of the events, coupled with the information amount of every event, forms a random variable whose average (expected) value is the average amount of information (entropy) generated by this distribution).

Με βάση αυτή την σχέση, και μετά από αρκετές προσαρμογές, δημιουργήσαμε την δική μας συνάρτηση σε γλώσσα python η οποία υπολογίζει την εντροπία ενός attribute σε ένα συγκεκριμένο dataframe. Ο κωδικός μαζί με την επεξήγηση του υπάρχει στα αρχεία παραδοσης (συγκεκριμένο στο InformationGain.py).

Συνεχίζοντας, το information gain:

$$IG(A, S) = H(S) - \sum_{t \in T} p(t)H(t)$$

Where,

- $H(S)$ - Entropy of set S
- T - The subsets created from splitting set S by attribute A such that $S = \bigcup_{t \in T} t$
- $p(t)$ - The proportion of the number of elements in t to the number of elements in set S
- $H(t)$ - Entropy of subset t

- is the measure of the difference in entropy from before to after the data set S is split on an attribute A
- in other words, how much uncertainty in S was reduced after splitting data set S on attribute A
- it is a synonym for Kullback–Leibler divergence (in the context of decision trees, the term is sometimes used synonymously with mutual information, which is the expectation value of the Kullback–Leibler divergence of a conditional probability distribution. The expected value of the information gain is the mutual information $I(X; A)$ of X and A – i.e. the reduction in the entropy of X achieved by learning the state of the random variable A. In machine learning, this concept is used to define a preferred sequence of attributes to investigate to most rapidly narrow down the state of X. Such a sequence (which depends on the outcome of the investigation of previous attributes at each stage) is called a decision tree. Usually an attribute with high mutual information should be preferred to other attributes).

Παρομοια , δημιουργησαμε την αντιστοιχη συναρτηση
information gain , βασιζομενοι σε αυτες τις πληροφοριες.

(πηγες: <https://stackoverflow.com/questions/1859554/what-is-entropy-and-information-gain>

<http://christianherta.de/lehre/dataScience/machineLearning/decision-trees.php>)

Το information gain για καθε feature ειναι το παρακατω :

```
[('Attribute1', 0.715, 0.09383),
 ('Attribute2', 0.70856, 0.03179),
 ('Attribute3', 0.69375, 0.03789),
 ('Attribute4', 0.68376, 0.0269),
 ('Attribute5', 0.70875, 0.0153),
 ('Attribute6', 0.6875, 0.0222),
 ('Attribute7', 0.6775, 0.01455),
 ('Attribute8', 0.67125, 0.00734),
 ('Attribute9', 0.67375, 0.01275),
 ('Attribute10', 0.67125, 0.00568),
 ('Attribute11', 0.65625, 0.00023),
 ('Attribute12', 0.6425, 0.01491),
 ('Attribute13', 0.62551, 0.01175),
 ('Attribute14', 0.63125, 0.00705),
 ('Attribute15', 0.615, 0.01162),
 ('Attribute16', 0.5875, 0.0024),
 ('Attribute17', 0.59625, 0.00295),
 ('Attribute18', 0.59625, 0.00013),
 ('Attribute19', 0.5775, 0.00121),
 ('Attribute20', 0.59125, 0.00771)]
```

Το Information gain καθε Attribute αντιστοιχει τοσο στις προβλεψεις και εκτιμησεις που καναμε με βαση την λογικη μας , οσο και στα διαγραμματα τα οποια εξηγησαμε.

Χαρακτηρικα βλεπουμε το πως τα features:

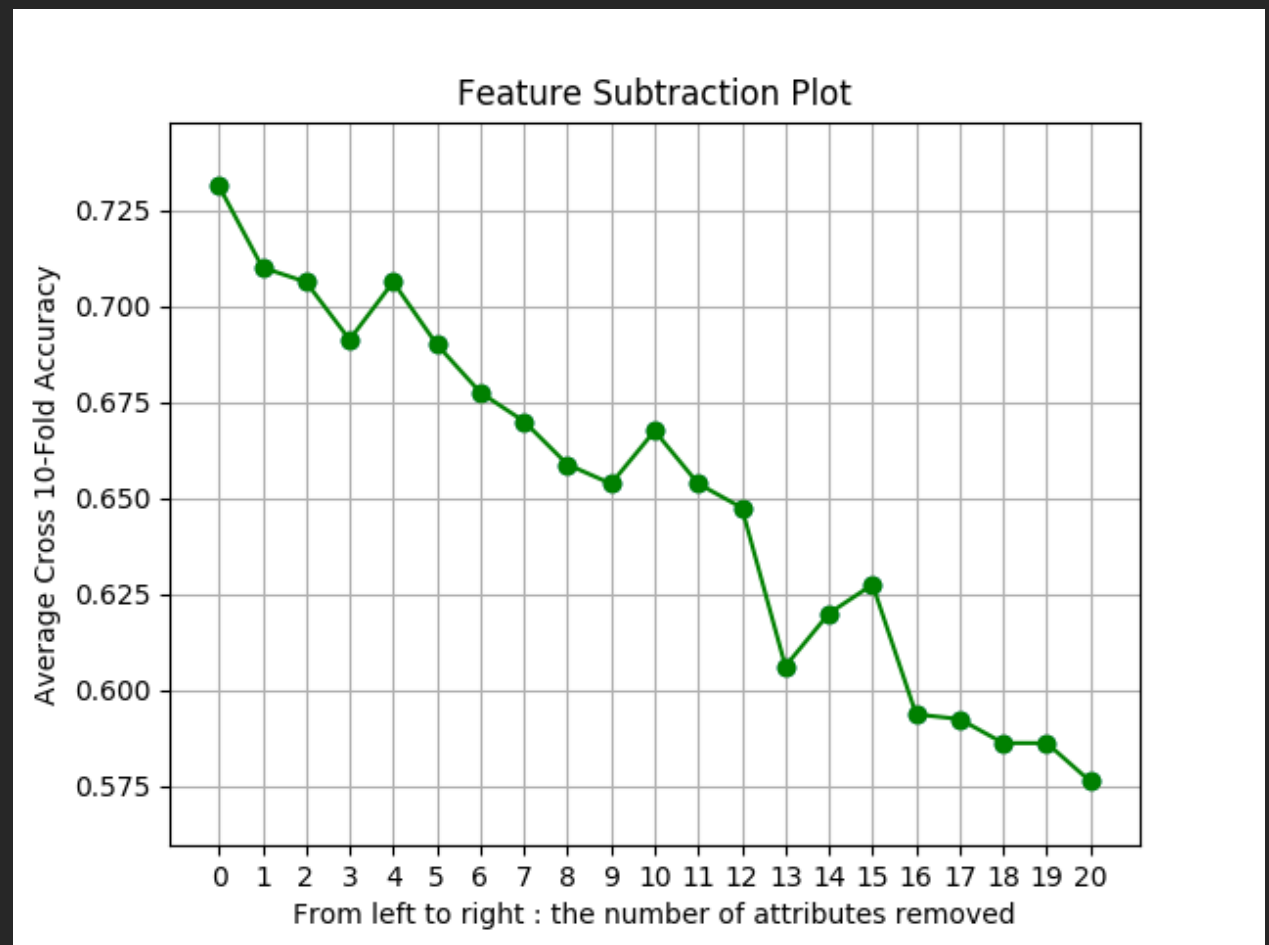
Attribute 1 , 7 , 2 , 12 ,15 μας παρεχουν υψηλα επιπεδα information gain. Αυτο δικαιολει και δικαιολογεται απο τις παρατηρησεις που καναμε . Αντιστοιχα , μεταβλητες που εκτιμησαμε να εχουν μικρη σημασια , οπως ειναι οι Attributes 8 και 18 , εχουν πολυ μικρο IG. Το ιδιο συμβαινει και για αυτες με μεσαια σημασια.

- Σταδιακη αφαιρεση features:

Διαλεξαμε τον random forest , ως καλυτερος classifier απο προηγουμενο ερωτημα , και εκτελεσαμε , ακριβως οπως εξηγησαμε παραπανω , Classification καθε φορα αφαιρωντας ενα feature.

Παρουσίαση Plot / Πινάκα σχέσης μεταξύ IG και Accuracy του Classifier:

Στο παρακατω plot παρουσιαζω το πως αλλαζει το μεσο accuracy για 10 fold cross validation , καθώς αφαιρω feature απο τον classifier



Attribute 1	0.71	0.93
Attribute 2	0.70625	0.03179
Attribute3	0.69125	0.03789
Attribute4	0.70625	0.0269

Attribute5	0.69	0.0153
Attribute6	0.6775	0.0222
Attribute7	0.67	0.01455
Attribute8	0.65875	0.00734
Attribute9	0.65375	0.01275
Attribute10	0.6675	0.00568
Attribute11	0.65375	0.00023
Attribute12	0.65375	0.01491
Attribute13	0.6475	0.01175
Attribute14	0.60625	0.00705
Attribute15	0.62	0.01162
Attribute16	0.62751	0.0024
Attribute17	0.59375	0.00295
Attribute18	0.5925	0.00013
Attribute19	0.58625	0.00121
Attribute20	0.57626	0.00771

Παρουσιάζεται ο πίνακας με το κάθε feature

Που αφαιρούμε και το αντίστοιχο information gain του

Τέλος , παρουσιάζουμε σε ένα plot , το πως το μέσο accuracy με cross fold 10 validation , του random forest classifier , φθίνει κάθε φορά που αφαιρούμε και ένα attribute. Αυτό συμβαίνει , γιατί κάθε φορά χάνεται «υλικό» με το οποίο ο κατηγοριοποιητής μας κάνει τις προβλέψεις του. Παρατηρούμε ότι οι πτώσεις από αφαίρεση σε αφαίρεση , είναι αναλογές του information gain που παρέχουμε.

Τελος , παρουσιαζουμε το Prediction του random forest ,
πανω στο test set:

Client_ID	Predicted_Label
11100	1
11099	2
11098	1
11097	1
11096	1
11095	1
11094	2
11093	1
11092	1
11091	1
11090	1
11089	1
11088	1
11087	1
11086	2
11085	1
11084	1
11083	2
11082	2
11081	1
11080	1
11079	1
11078	1
11077	1
11076	1
11075	1
11074	2
11073	2
.....	.

Συνοψίζοντας

Εχοντας ολοκληρήσει την εργασία , έχουμε ασχοληθεί με την εξερεύνηση των δεδομένων μιας τραπεζας , την αξιολόγηση των χαρακτηριστικών τους και πως αυτά μπορούν να εκμεταλλευτούν για το δικό της κέρδος.

Απο όλη αυτή την εξερεύνηση και μελέτη των δεδομένων(γραφικά , αλγοριθμικά , μαθηματικά) μας κάνει εντύπωση η ανταπόκριση των αποτελεσμάτων στην πραγματικότητα και πήραμε μια εικόνα εις βάθος για τις πτυχές και το βάθος στο οποίο υπάρχουν δεδομένα.

Στην προηγούμενη άσκηση είδαμε το πως οι αλγόριθμοι βοηθούν εμάς , για την εξορυξη των δεδομένων , σε αυτήν , το πως εμείς μπορούμε να τους βοηθήσουμε.