

Εργασία 2 Γραμμικά μοντέλα στο scikit-learn

Τεχνητή Νοημοσύνη

Νικόλαος Γουρνάκης, it22023

13 louvíou 2024

Ερωτημα 1)

1)

Για να απαντήσω στο ποσες εγγραφες εχουμε στο συνολο δεδομενων χρησιμοποιω το shape attribute του dataframe

```
print(data.shape)
```

```
(12330, 18)
```

Οπου το πρωτο νουμερο ειναι το συνολο εγγραφων και και το δευτερο ειναι το συνολο στηλων/features

Αρα

```
print(data.shape[0])
```

εχουμε 12330 εγγραφες.

2)

ποσοστό από αυτές οι χρήστες αγόρασαν τελικά = συνολο ατομων που αγορασαν / συνολο ατομων * 100. Τα ατομα που αγορασαν ειναι αυτα που στο 'Revenue' ειναι True

```
print(data[data['Revenue'] = True].shape[0] / data.shape[0] * 100)
15.474452554744525
```

αρα το ποσοστό από αυτές οι χρήστες αγόρασαν τελικά ειναι περιπου 15.47%

3)

Για αυτο το σεναριο ξερουμε οτι το μοντελο παντα προβλεπει False αρα για να υπολογίσουμε την ευστοχια πρεπει να κανουμε = συνολο ατομων που **ΔΕΝ** αγορασαν / συνολο ατομων * 100. Μια εναλλακτικη θα ηταν να κανουμε αφαιρέσουμε το αποτελεσμα του προηγουμενου ερωτηματος απο το 100 και θα εχουμε την ιδια απαντηση.

```
print(data[data['Revenue'] = False].shape[0] / data.shape[0] * 100)
84.52554744525548
```

Το accuracy ενος μοντελου που παντα επιστρεφει False ειναι περιπου 84.52%

Ερωτημα 5)

Ευστοχια μοντελου στο συνολο εκπαιδευσης:

```
print(model.score(X_train, y_train))
```

0.8767234387672344

Ευστοχια μοντελου στο συνολο δοκιμης:

```
print(model.score(X_test, y_test))
```

0.8745606920789403

Πινακας συγχυσης:

```
print(confusion_matrix(y_test, model.predict(X_test)))
```

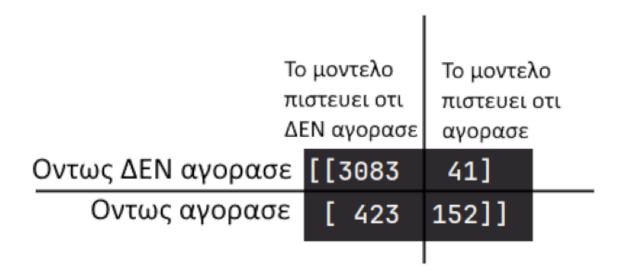
```
[[3083 41]
[ 423 152]]
```

Ερμηνεια πινακα συγχυσης:

Αμα δουμε τα docs της συναρτησης confusion_matrix

Returns: Confusion matrix whose i-th row and j-th column entry indicates the number of samples with true label being i-th class and predicted label being j-th class.

Βλεπουμε οτι η γραμμες ειναι true classes και οι στηλες predicted classes αρα αυτο που μας λεει ο πινακας συγχυσης ειναι το ακολουθο



Αρα 3083 ατομα στο τεστ το μοντελο καταφερε να μαντέψει σωστα οτι δεν θα αγορασουν και 423 λαθος οτι δεν θα αγορασουν ενω στην πραγματικοτητα θα αγορασουν. Και 152 ατομα στο τεστ το μοντελο καταφερε να μαντέψει σωστα οτι θα αγορασουν και 41 λαθος οτι θα αγορασουν ενω πραγματικα δεν θα αγορασουν. Το accuracy του μοντελου βγαινει απο την διαγωνιο / το συνολο δηλαδη = (3083 + 152) / (3083 + 152 + 423 + 41) = 3235 / $3699 \approx 0.87$. Βλεπω οτι το μοντελο δινει παραπανω βαρος στο να μαντέψω κατι σαν "μη αγορα" καθως εχουμε 423 False Negatives, αυτο λογικα γινεται γιατι η κατανομη τον κλασεων ειναι παρα πολυ ανησωροπη.

Τροποποιησεις:

Το κυριο προβλημα ειναι αυτη η ανισορροπια μεταξυ τον δυο κλασεων, ξερουμε οτι στο συνολο δεδομενων τα ατομα που θα αγορασουν ειναι 15.47% αρα η True class ειναι 15.47% και η False class ειναι το υπολοιπο 84.53% που ειναι τεραστια διαφορα.

Αυτο μπορει να προκαλεσει προβληματα:

- Το μοντελο δινει πιο πολυ βαρος στο False class καθως ειναι πιο πιθανο να ειναι σωστο
- Με το μικρο ποσοστο του True class εχουμε υψηλη πιθανοτητα να γινει καταστροφικος διαχωρισμος των δεδομενων, δηλαδη ολα τα False class να πεσουν στο training set και ολα τα True class να πεσουν στο Testing set, που αυτο εχει ως αποτελεσμα το μοντελο να μαντεύει παντα false.

Αρα τροποι βελτιωσεις:

- Καταρχας θα ηταν καλο να κανουμε KFold validation για να σιγουρευτούμε ότι δεν κανουμε overfit.
- Θα μπορουσαμε να βαλουμε βαρος στην καθε κλαση ετσι ωστε το μοντελο να κρινει πιο σημαντικο να πετυχει σωστα την κλαση μειονοτητα.
- Θα μπορουσαμε να κανουμε stratification στο train_test_split ετσι ωστε να υπαρχουν αρκετα δεδομενα και απο τις δυο κλασεις και στα δυο σετ.
- Θα μπορουσαμε να κανουμε undersampling ετσι ωστε να φερουμε το νουμερο το κλασεων στο 50/50
- Θα μπορουσαμε να κανουμε <u>SMOTE</u> για να αυξήσουμε συνθετικα το νούμερο της κλασης μειονοτητα