

Die Labels aus den Annotations habe ich mit dem LabelEncoder von sklearn encoded.

Das nicht-numerische Feature „dns\_query\_name“ enthält als Werte Listen, die nur 1 String beinhalten. Daher müssen diese Strings erstmal extrahiert werden.

Für nicht-numerische Features, die als Werte Strings enthalten, habe ich einen eigenen Feature-Encoder „My\_Feature\_Encoder“ programmiert. Im Grunde genommen ist es ein LabelEncoder, der aber nach dem Fitting neuhinzugekommene Werte hinzufitten kann und auch mit fehlenden Werten bzw. nans umgehen kann. Dieser Encoder encodiert Strings auf numerische Werte.

Ich habe mich dafür entschieden, fehlende Werte bzw. nans mit der Zahl -1 zu ersetzen. Denn sowohl in Trainings- als auch die Test-Daten ist die kleinste vorkommende Zahl eine 0 gewesen. Diese Idee wird im ganzen Programm umgesetzt.

Danach werden sowohl die Trainings- als auch die Test-Daten preprocesset. Hier gilt es fehlende Werte mit -1 zu ergänzen und nicht-numerische Features in numerische Features zu verwandeln.

In den Trainings- und Test-Dataframes werden diese Features entfernt:

- 'da', 'sa': Können nur 1 Wert annehmen. Daher sind sie irrelevant
- "rev\_intervals\_ccnt", "pld\_ccnt", "rev\_hdr\_ccnt", "rev\_pld\_ccnt", "rev\_ack\_psh\_rst\_syn\_fin\_cnt", "ack\_psh\_rst\_syn\_fin\_cnt", "intervals\_ccnt", "hdr\_ccnt": Diese Histogramme haben nach dem RandomForest-Fitting bei Ausgabe der Importances so gut wie gar keine Wichtigkeit aufgewiesen. Außerdem ist es sehr aufwendig diese Histogramme zu encodieren. Denn jedes dieser Histogramme müsste in 5 bis 16 Features aufgeteilt werden. Das verbraucht bei 8 Histogrammen sehr viel Speicher und Leistung, obwohl sie so oder so unwichtig scheinen.
- 'dns\_query\_name\_len', 'dns\_answer\_ttl', 'dns\_query\_class', 'dns\_answer\_cnt', 'dns\_query\_type', 'dns\_query\_cnt', 'http\_uri': Aus logischer Sicht für mich unwichtig für die Klassifizierung.

Alle TLS-Arrays sind nicht-numerische Features. Die Arrays bestehen aus numerischen Werten und werden zu Mittelwerten transformiert. Dabei werden alle numerischen Werte addiert und durch die zugehörigen Counting-Features dividiert. Denn diese TLS-Arrays haben keine gleiche Länge. Um für jedes Feature ein Array auf die maximale Länge zu füllen, müsste man manche Feature in 80 bis 192 Features aufteilen. Das würde sehr viel Speicher und Leistung verbrauchen.

Die nicht-numerischen Features 'dns\_query\_name', 'http\_content\_type', 'http\_host' haben als Werte Strings. Diese werden mit den jeweiligen My\_Feature\_Encodern encodiert.

Anschließend haben wir noch das nicht-numerische Feature „dns\_answer\_ip“. Es enthält Listen ohne fixe Länge. Diese Listen enthalten IP-Adressen. Wie bei den TLS-Arrays würde es sich nicht lohnen das Feature auf die maximale Länge x einer Liste in x Features aufzuteilen. Stattdessen berechnet das Programm die gemeinsame Netzwerkadresse als Dezimalzahl und speichert sie in „dns\_answer\_common\_network\_ip“.

Sowohl das Trainings- als auch das Test-Dataframe müssen für das Fitting und das Predicten die gleiche Feature-Reihenfolge haben. Daher erhält das Test-Dataframe die Feature-Reihenfolge des Trainings-Dataframes.

Ich habe mich für einen RandomForest als Classifier entschieden. Als Anzahl der Estimators habe ich 100 gewählt, weil es sich bei meinen Tests als besten Wert herausgestellt hat.