

ISOM 3360 Assignment 2

Joshua Chang, Fung Ho Yin Matthew, Chan Shun Hang

{jchangad, hymfung, shchanax}@connect.ust.hk

1. Exploring and preprocessing the dataset

Before we build any models, we first need to explore and reformat the dataset. The dataframe consisted of 7032 Telco customers, and had 19 columns to each customer, with the last column (churn) being our target variable. This means the dataset had shape (7032, 19).

We checked and confirmed that the dataset did not contain any na values, and used `np.unique()` to find out the possible values to each column. After exploring the dataset, we are able to group the columns of the dataframe into three main categories:

- General information

These columns contained basic information that every customer had, for example their gender, whether they have a partner, their monthly charges and their contract length.

- Phone service related information

These columns had information on whether the customers were subscribed to their phone service, and whether they had multiple lines with Telco.

- Internet service related information

These columns had information on whether the customers were subscribed to their internet service, had online security, online backup and so on.

Since most columns contained only a few possible values, we converted the dataframe into a NumPy array for easier modelling later on. The detailed conversion rules can be found in the python notebook.

2. Building decision trees (Task 1, Question 1)

We built **FILL THIS IN LATER** decision tree models in total. Each model was trained and validated using 10-fold cross validation on 70% of the whole dataset. To ensure our trees are reproducible for fair comparison between different trees, we fixed “random_state” as “2211”, Evaluation of the trees were done using **accuracy**, **precision**, **recall**, and **AUC** as metrics, where each metric was taken as the simple average over all 10 folds. The summary of our models can be found in table 1.

Tree	Accuracy	Precision	Recall	AUC
Basic (§2.1)	0.738	0.505	0.517	0.668
Adj. params (§2.2)	0.0	0.0	0.0	0.0
Ensemble (§2.3)	0.0	0.0	0.0	0.0

Table 1. Comparison between trees

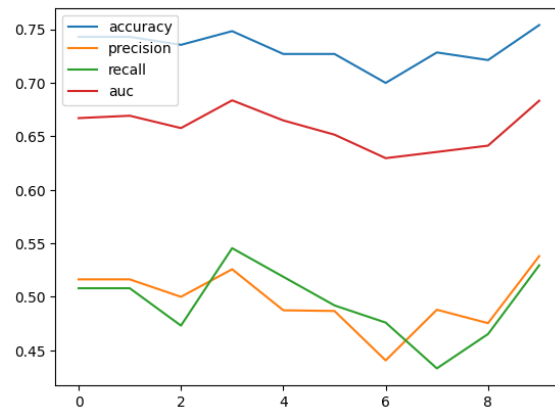


Figure 1. different metrics across the folds

2.1. The basic tree

The first and most basic model we built was a basic decision tree with all the features used for training, and all the default parameters `sklearn` provides. This means using “entropy” to decide which feature to split, splitting with the best feature, and allowing the tree to grow until all leaves were pure or contained less than 2 samples. During each split, the model considers *all* attributes.

Figure 1 shows the metrics across different folds during cross validation, and figure 2 shows the AUC curve. The metrics seems normal, however the AUC curve is too good to be true. Indeed, upon inspection we found that

```
model.predict_proba(X)
>>> array([[0., 1.], [0., 1.],
..., [1., 0.]])
```

which indicates that the model predicts the probabilities for each sample with 100% certainty. In fact, when comparing

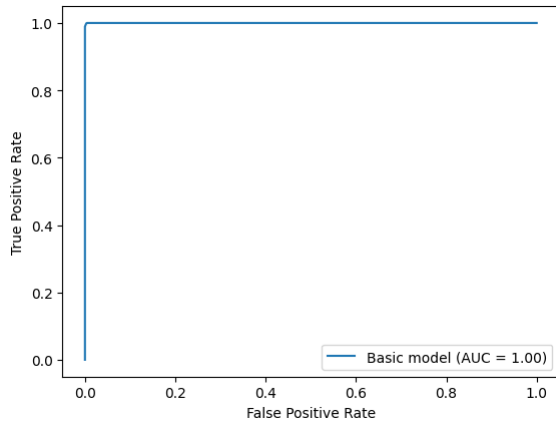


Figure 2. AUC curve for the basic tree

it with the ground truth, out of 4,992 samples it correctly classifies 4,889 of them. This means adjusting the threshold won't affect the prediction results, which explains the absurd AUC curve.

The model also had a depth of 34 with 1,467 leaves. The reason the model overfits the dataset is because by default, `sklearn` will continue splitting until each leaf is pure or contains 2 or less samples, so these are the first parameters we are going to change in our next iteration.

2.2. A tree with adjusted parameters

The next model we will build aims to solve the issue above. Specifically, we will adjust the following parameters, and find the best model by `GridSearchCV` using 10 folds.

- `max_depth`

Search range is `np.arange(5, 35, 5)`.

We chose this range because the overfitted model in section 2.1 had a depth of 34, so setting `max_depth` to 30 should be sufficient.

- `min_samples_split`

Search range is `np.arange(10, 110, 10)`.

We chose this range because based on empirical testing. Even when we set lower / higher values, the best estimator found by `GridSearchCV` still mostly took on values in this range.

- `min_samples_leaf`

Search range is `np.arange(10, 110, 10)`.

We chose this range based on the same empirical testing reason.

After deciding how to build the tree, the next step is to decide how

Criteria	Accuracy	Precision	Recall	AUC
Accuracy	0.738	0.505	0.517	0.668
Precision	0.0	0.0	0.0	0.0
Recall	0.0	0.0	0.0	0.0
F1	0.0	0.0	0.0	0.0

Table 2. The best trees based on different criterion