

# Data science automation

Using our prepared churn data from week 2:

- Load data
- Use pycaret to find an ML algorithm that performs best on the data
  - Choose a metric you think is best to use for finding the best model. by default, it is accuracy.
- Save the model to disk
- Create a Python script/file/module with a function that takes a pandas dataframe as an input and returns the probability of churn for each row in the dataframe
  - Python file/function should print out the predictions for new data (new\_churn\_data.csv)
  - The true values for the new data are [1, 0, 0, 1, 0]
  - Test your Python module and function with the new data, new\_churn\_data.csv
- Write a short summary of the process and results at the end of this notebook

## Load data

Install Required Libraries. You can do this in your terminal. !conda create -n msds python=3.10.14 -y !conda init !conda activate msds !pip install --upgrade pycaret

Load our prepared data from week 2 where everything has been converted to numbers. Many autoML packages can handle non-numeric data (they usually convert it to numeric with various methods).

```
In [1]: import pandas as pd # Import pandas for data manipulation and analysis
df = pd.read_excel('./prepared_churn_data.xlsx', index_col='customerID')
df
```

Out[1]:

	tenure	PhoneService	Contract	PaymentMethod	MonthlyCharges	TotalCharg
customerID						
7590-VHVEG	1	0	0	3	29.85	29.
5575-GNVDE	34	1	1	2	56.95	1889.
3668-QPYBK	2	1	0	2	53.85	108.
7795-CFOCW	45	0	1	1	42.30	1840.
9237-HQITU	2	1	0	3	70.70	151.
...	...	...	...	...	...	...
6840-RESVB	24	1	1	2	84.80	1990.
2234-XADUH	72	1	1	0	103.20	7362.
4801-JJAZL	11	0	0	3	29.60	346.
8361-LTMKD	4	1	0	2	74.40	306.
3186-AJIEK	66	1	2	1	105.65	6844.

7043 rows × 8 columns



## AutoML with pycaret - to find an ML algorithm that performs best on the data

Use pycaret for autoML. Install the Python package with conda or pip: `conda install -c conda-forge pycaret -y`. Then we can import the functions we need:

```
In [2]: from pycaret.classification import setup, compare_models, predict_model, save_model
        automl = setup(df, target='Churn') # Initialize the PyCaret environment with the Da
```

	Description	Value
0	Session id	3345
1	Target	Churn
2	Target type	Binary
3	Original data shape	(7043, 8)
4	Transformed data shape	(7043, 8)
5	Transformed train set shape	(4930, 8)
6	Transformed test set shape	(2113, 8)
7	Numeric features	7
8	Preprocess	True
9	Imputation type	simple
10	Numeric imputation	mean
11	Categorical imputation	mode
12	Fold Generator	StratifiedKFold
13	Fold Number	10
14	CPU Jobs	-1
15	Use GPU	False
16	Log Experiment	False
17	Experiment Name	clf-default-name
18	USI	9030

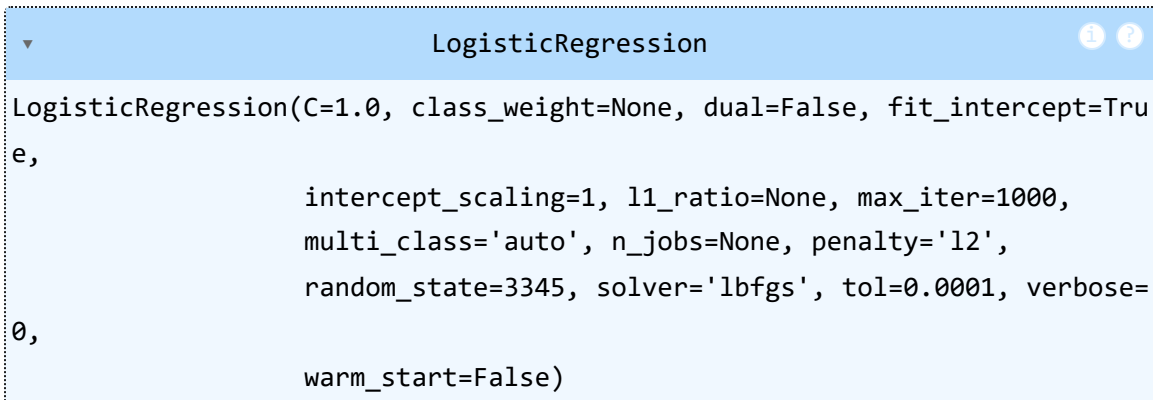
```
In [3]: best_model = compare_models() # Compare different models and select the best one ba
```

	Model	Accuracy	AUC	Recall	Prec.	F1	Kappa	MCC	TT (Sec)
<b>lr</b>	Logistic Regression	0.7923	0.8330	0.5069	0.6371	0.5642	0.4302	0.4353	0.8400
<b>gbc</b>	Gradient Boosting Classifier	0.7915	0.8384	0.4862	0.6413	0.5521	0.4199	0.4272	0.2910
<b>ada</b>	Ada Boost Classifier	0.7892	0.8322	0.4863	0.6350	0.5500	0.4157	0.4225	0.1310
<b>ridge</b>	Ridge Classifier	0.7890	0.8204	0.4404	0.6518	0.5254	0.3963	0.4090	0.0310
<b>lda</b>	Linear Discriminant Analysis	0.7862	0.8204	0.4878	0.6244	0.5475	0.4103	0.4158	0.0270
<b>lightgbm</b>	Light Gradient Boosting Machine	0.7840	0.8276	0.5030	0.6131	0.5520	0.4117	0.4155	0.5430
<b>rf</b>	Random Forest Classifier	0.7763	0.8081	0.4809	0.5967	0.5315	0.3872	0.3916	0.2730
<b>svm</b>	SVM - Linear Kernel	0.7698	0.7477	0.4204	0.6065	0.4781	0.3425	0.3594	0.0350
<b>knn</b>	K Neighbors Classifier	0.7661	0.7520	0.4388	0.5779	0.4981	0.3496	0.3556	0.8150
<b>et</b>	Extra Trees Classifier	0.7643	0.7849	0.4801	0.5654	0.5185	0.3641	0.3666	0.1950
<b>qda</b>	Quadratic Discriminant Analysis	0.7477	0.8154	0.7110	0.5180	0.5992	0.4217	0.4331	0.0460
<b>nb</b>	Naive Bayes	0.7414	0.8125	0.6926	0.5094	0.5868	0.4050	0.4153	0.0350
<b>dummy</b>	Dummy Classifier	0.7347	0.5000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0280
<b>dt</b>	Decision Tree Classifier	0.7296	0.6645	0.5129	0.4916	0.5017	0.3164	0.3168	0.0440

Our `best_model` object now holds the highest-scoring model. We can also set an argument `sort` in `compare_models` to choose another metric as our scoring metric. By default, it uses accuracy (and we can see the table above is sorted by accuracy). We could set this to `sort='Precision'` to use precision ( $TP / (TP + FN)$ ), for example.

```
In [4]: best_model # Display the best-performing model selected by PyCaret, which in this c
```

Out[4]:



```
LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
                  intercept_scaling=1, l1_ratio=None, max_iter=1000,
                  multi_class='auto', n_jobs=None, penalty='l2',
                  random_state=3345, solver='lbfgs', tol=0.0001, verbose=
0,
                  warm_start=False)
```

It looks like our best model is LR, closely followed by some others. This may change when you re-run this - there is some randomness built in that we are not fixing (e.g. for the cross-validation splits possibly), so the top model may be different each time this is run since the accuracy scores are so similar between models.

We can now use the model to make predictions. If our data is not being preprocessed, we can simply use the `best_model` object, which is an sklearn model, to make predictions:

In [5]: `df.iloc[-2:-1].shape`

Out[5]: (1, 8)

We are selecting the last row, but using the indexing `[-2:-1]` to make it a 2D array instead of 1D (which throws an error). Try running `df.iloc[-1].shape` and `df.iloc[-2:-1].shape` to see how they differ.

However, this only works if we set `preprocess=False` in our setup function. Otherwise the order of features may be different

A more robust way (in case we are using preprocessing with autoML) is to use pycaret's `predict_model` function:

In [6]: `predict_model(best_model, df.iloc[-2:-1]) # Predict the output using the best model`

	Model	Accuracy	AUC	Recall	Prec.	F1	Kappa	MCC
0	Logistic Regression	1.0000	0	1.0000	1.0000	1.0000	nan	0.0000

customerID	tenure	PhoneService	Contract	PaymentMethod	MonthlyCharges	TotalCharg
8361-LTMKD	4	1	0	2	74.400002	306.6000

We can see this creates a new column, 'Score', with the probability of class 1. It also creates a 'Label' column with the predicted label, where it rounds up if score is  $\geq 0.5$  (greater than or equal to 0.5).

## Saving and loading our model

Save our trained model so we can use it in a Python file later. pycaret has a handy function for this, which saves the model as a pickle file:

```
In [7]: save_model(best_model, 'best_churn_model') # Save the best model as a .pkl file for
```

Transformation Pipeline and Model Successfully Saved

```
Out[7]: (Pipeline(memory=Memory(location=None),
                 steps=[('numerical_imputer',
                        TransformerWrapper(exclude=None,
                                           include=['tenure', 'PhoneService',
                                                  'Contract', 'PaymentMethod',
                                                  'MonthlyCharges', 'TotalCharges',
                                                  'charge_per_tenure'],
                                           transformer=SimpleImputer(add_indicator=False,
                                                                      copy=True,
                                                                      fill_value=None,
                                                                      keep_empty_features=False,
                                                                      missing_values=nan,
                                                                      strategy='mean'))),
                        ('c...
                        fill_value=None,
                        keep_empty_features=False,
                        missing_values=nan,
                        strategy='most_freq
                        uent'))),
          ('trained_model',
           LogisticRegression(C=1.0, class_weight=None, dual=False,
                              fit_intercept=True, intercept_scaling=1,
                              l1_ratio=None, max_iter=1000,
                              multi_class='auto', n_jobs=None,
                              penalty='l2', random_state=3345,
                              solver='lbfgs', tol=0.0001, verbose=0,
                              warm_start=False)),
          verbose=False),
         'best_churn_model.pkl')
```

`pickle` is a built-in module in the Python standard library which allows for saving and loading of binary data. It's data that's been encoded (usually using hexadecimal encoding) to a file, and we can store any Python object as-is in a pickle file. Then we can load the data from the file and be right back where we left off.

```
In [8]: import pickle # Imports the pickle module for serializing and deserializing Python

with open('best_model.pk', 'wb') as f: # Opens a file named 'best_model.pk' in write mode
    pickle.dump(best_model, f) # Saves the best model to the file using pickle
```

Here, we use the built-in `open` function to open a file with the name `best_model.pk`, then open it for writing with `'w'` and in a binary format using `'b'`. We save that file object in the variable `f`. The `with` statement automatically closes the file after we exit the `with` statement, otherwise, we should call the function `close` from the file object `f`. Then we use `pickle` to save our data to the file. We could reload it like this:

```
In [9]: with open('best_model.pk', 'rb') as f: # Opens the saved file in read-binary mode.
        loaded_model = pickle.load(f) # Loads the model back into memory.
```

```
In [10]: new_data = df.iloc[-2:-1].copy() # Select the second-to-last row for prediction and
new_data.drop('Churn', axis=1, inplace=True) # Remove the 'Churn' column from the data
loaded_model.predict(new_data) # Make a prediction using the loaded model.
```

```
Out[10]: array([1], dtype=int8)
```

Loading it is almost the same, except we use `rb` for "read binary" and use `pickle`'s `load` function.

Under the hood, `pycaret` is doing something similar, but we can use it with the `save_model` function as we saw above.

Once we have our saved `pycaret` model, we can test loading it and making predictions to make sure it works:

```
In [11]: loaded_model = load_model('best_churn_model') # Load the saved model from the .pkl
```

Transformation Pipeline and Model Successfully Loaded

```
In [12]: predict_model(loaded_model, new_data) # Make a prediction using the loaded model
```

```
Out[12]:
```

	tenure	PhoneService	Contract	PaymentMethod	MonthlyCharges	TotalCharg
customerID						
8361-LTMKD	4	1	0	2	74.400002	306.6000

## Making a Python module to make predictions

We can now use this model in a Python file to take in new data and make a prediction. We will first need to compose a Python file. We can do this in many ways:

- Jupyter and Jupyter Lab

- VS Code
- Atom
- Notepad++
- Other text editors or IDEs (integrated development environments)

The benefit of using a code editor or IDE is that it will have lots of bells and whistles, like syntax highlighting, autocomplete, and many other things depending on the code editor or IDE. VS Code is one of the top-most used editors by data scientists and software developers, although you can try any IDE or code editor for Python that you like. You can easily install VS Code through Anaconda Navigator or by visiting the VS Code website. VS Code is developed by Microsoft, and there is also an IDE Visual Studio Code.

The file we've created is show below:

```
In [13]: from IPython.display import Code # Import the Code function to display code files i
Code('predict_churn.py') # Display the contents of the 'predict_churn.py' file in a
```



```

Out[13]: import pandas as pd
from pycaret.classification import predict_model, load_model

def load_data(filepath):
    """
    Loads churn data into a DataFrame from a string filepath.
    """
    df = pd.read_excel('./prepared_churn_data.xlsx', index_col='customerID')
    return df

def make_predictions(df):
    """
    Uses the pycaret best model to make predictions on data in the df dataframe.
    """
    model = load_model('best_churn_model')
    predictions = predict_model(model, data=df)

    # Check the column names
    print(predictions.columns)

    # Rename 'prediction_label' to 'Churn_prediction' if it exists
    if 'prediction_label' in predictions.columns:
        predictions.rename(columns={'prediction_label': 'Churn_prediction'}, inplace=True)

    # Replace values in the new column
    predictions['Churn_prediction']
    ].replace({1: 'Churn', 0: 'No Churn'}, inplace=True)

    return predictions['Churn_prediction']
    else:
        raise KeyError("The 'prediction_label' column was not found in the predictions DataFrame")

if __name__ == "__main__":
    df = load_data('./new_Churn_data.csv')
    predictions = make_predictions(df)
    print('predictions:')
    print(predictions)

```

```

In [15]: #!/dir # Execut the 'dir' command to list the files and directories in the current w

```

In [16]: `%run predict_churn.py` # Execute the 'predict\_churn.py' script within the Jupyter No

Transformation Pipeline and Model Successfully Loaded

	Model	Accuracy	AUC	Recall	Prec.	F1	Kappa	MCC
0	Logistic Regression	0.7933	0.8356	0.5067	0.6394	0.5654	0.4321	0.4372

```
Index(['tenure', 'PhoneService', 'Contract', 'PaymentMethod', 'MonthlyCharges',
      'TotalCharges', 'charge_per_tenure', 'Churn', 'prediction_label',
      'prediction_score'],
      dtype='object')
```

predictions:

customerID

7590-VHVEG Churn

5575-GNVDE No Churn

3668-QPYBK No Churn

7795-CFOCW No Churn

9237-HQITU Churn

...

6840-RESVB No Churn

2234-XADUH No Churn

4801-JZAZL No Churn

8361-LTMKD Churn

3186-AJIEK No Churn

Name: Churn\_prediction, Length: 7043, dtype: object

We can test out running the file with the Jupyter "magic" command `%run`:

We can run the above line over and over after making changes to the file while we are writing it. The true values are 1, 0, 0, 1, 1, so our model is working OK but not perfect. We have 2 false positives in the new data. However, this new data was synthesized based on existing data, so it is a little random.

## Summary

Write a short summary of the process and results

In this project we have utilized PyCaret to identify the best model for predicting customer churn.

- **Data Preparation:** The cleaned and formatted churn dataset for analysis from week 2 is loaded. This included handling missing values, encoding categorical variables, and normalizing numerical features.
- **AutoML with pycaret**
  - **Setting Up PyCaret:** The PyCaret environment using the `setup` function is initialized, specifying the target variable (churn) and selecting a suitable metric for model evaluation. The output from the `setup` function in PyCaret gives a summary for the churn prediction model. The dataset has 7043 rows and 8 columns

with `Churn` as the target variable. It splits the data into 4930 training samples and 2113 testing samples. There are 7 numeric features.

- **Model Comparison:** PyCaret's `compare_models` is used to evaluate different models for predicting customer churn. The `'compare_models()'` output shows that the Gradient Boosting Classifier has the highest accuracy at 79.15% and an AUC of 0.8384 showing it is good at distinguishing churners from non-churners. Logistic Regression also performed well with 79.23% accuracy while models like the Decision Tree and Dummy Classifier scored below 75%. The best model output shows that **Logistic Regression** is best model which can handle up to 1000 iterations. This model is reliable for binary classification and is expected to predict churn effectively. The models output for `predict_model(best_model, df.iloc[-2:-1])` predicts that the customer would churn with a prediction score of 0.5643 indicating a moderate chance of churn. The model shows strong potential for accurately predicting customer churn.
- **Model Saving:** Once the best model is determined it is saved using `save_model` storing it as `best_churn_model.pkl`. Then the `pickle` library is used to save and load the model from a file named `best_model.pk`. After loading the model a new data sample is prepared by selecting the second-to-last row and removing the target column. Finally the model is reloaded with `load_model` and make predictions again to show its effectiveness with new data. The models predicts that the customer would churn with a prediction score of 0.5643 indicating a moderate chance of churn.
- **Creating the Prediction Module:** A Python script is created which includes a function to take a Pandas DataFrame as input and return churn probabilities for each row. This function uses a saved model (`best_churn_model- LR` in this case) to generate predictions. The script imports necessary libraries and defines a function to load data from an Excel file. Within the `make_predictions` function the model is loaded and churn predictions are made based on the input DataFrame. The `prediction_label` column is renamed to `Churn_prediction` clearly indicating whether each customer is likely to churn. The main block of the script runs the predictions on new data from a CSV file and prints the results. The prediction function uses a new dataset (`new_churn_data.csv`) and print the churn probabilities. The true values for this dataset are known to be [1, 0, 0, 1, 0] allowing to evaluate the model's performance against actual outcomes. The output confirms that the model is loaded successfully and shows performance metrics for the Logistic Regression model. Finally it displays the churn predictions for all customers in the new dataset making it easy to identify those likely to churn.
- **Conclusion:** This project shows how effective PyCaret is for model selection and evaluation. The model comparison helped in determining that Logistic Regression is the best model for predicting customer churn. The churn dataset prepared ensuring data quality. PyCaret is used to identify the best-performing model Then a Python module is

created to streamline predictions allowing to load the model and make predictions on new data. The model is saved to be reused later. The output shows that the Logistic Regression model was loaded successfully with an accuracy of 79.23% and an AUC of 0.8330 indicating it performs well in predicting churn. The final predictions highlighted customers likely to churn which could be useful for retention efforts. Overall this project emphasizes the importance of data-driven decisions which can be expanded for future analyses.