

## Project Name: Telco Churn Classification Project

**Project Description:** develop and optimize machine learning model that will predict churn using telco dataset

### Project Objectives:

- **Data Pipeline:** Acquire data from SQL database and convert into a panda dataframe, prepare the data for exploration, explore the data and document key takeaways, visualize feature attributes and incorporate into machine learning model.
- Document data preparation/exploration progress, document key takeaways and present results in Jupyter Notebook
- Create modules (acquire.py, prepare.py) that make process repeatable for 3rd party
- Present thought process and modeling results to cohort utilizing Jupyter Notebook as presentation material
- Be prepared to handle questions about code, process, findings, key takeaways, and model.

### Project Goals:

- Find drivers for customer churn in the telco data that can be utilized to construct a predictive model
- Construct and optimize a machine learning classification model that accurately predicts churn while maximizing for recall
- Document the process so that 3rd party can read like a report and easily follow along/replicate

### Project Goals:

Stage	Tools	Brief Description of Process	Takeaways/Lessons Learned
Acquire	<ul style="list-style-type: none"><li>• python</li><li>• pandas</li><li>• SQL</li><li>• user-defined functions</li></ul>	<ul style="list-style-type: none"><li>• Defined two functions that established connection, ran query and joined tables in SQL ACE and imported into panda df</li></ul>	<ul style="list-style-type: none"><li>• Consider importing multiple queries into separate dfs as opposed to bringing the data in on one joined query</li></ul>

<b>Prepare</b>	<ul style="list-style-type: none"> <li>python</li> <li>pandas</li> <li>user-defined functions</li> </ul>	<ul style="list-style-type: none"> <li>Cleaned raw data (dealt with na, null, blank, values), dropped redundant columns, imputed and encoded data</li> <li>Defined function that automated the process</li> </ul>	<ul style="list-style-type: none"> <li>Keeping the raw data fully intact and untouched in a separate df is a good practice</li> </ul>
<b>Explore</b>	<ul style="list-style-type: none"> <li>python</li> <li>pandas</li> <li>matplotlib</li> <li>seaborn</li> <li>SciPy</li> <li>sklearn</li> </ul>	<ul style="list-style-type: none"> <li>Plotted and visualized features independently and visualized feature relationships vs churn, tenure &amp; monthly charges</li> <li>Tested selected features for statistical relevance vs churn for model selection</li> </ul>	<ul style="list-style-type: none"> <li>In my opinion, the most important part of the pipeline</li> <li>Some of the best insights into the data came from experimentation and unplanned exploration.</li> <li>A methodically predetermined start didn't yield results but sparked fruitful thoughts</li> <li>The power of the visuals can not be overstated and ultimately framed the statistical and model results</li> </ul>
<b>Model</b>	<ul style="list-style-type: none"> <li>sklearn</li> </ul>	<ul style="list-style-type: none"> <li>Defined 3 models and evaluated them against each other</li> <li>I used Decision Tree, Random Forest &amp; Logistic Regression machine learning models</li> </ul>	<ul style="list-style-type: none"> <li>Think of modeling as an inverted pyramid. The Test set is where everything is done. The Validate set is touched with every model and the Test is only touched with gloves</li> </ul>
<b>Evaluate</b>	<ul style="list-style-type: none"> <li>python</li> <li>pandas</li> <li>sklearn</li> </ul>	<ul style="list-style-type: none"> <li>Established baseline</li> <li>Fit the models to the train data and then evaluated initially on accuracy</li> <li>Used sklearn confusion matrix and classification</li> </ul>	<ul style="list-style-type: none"> <li>Understanding what common evaluation metrics like (Precision, Recall, Accuracy) measured was key to optimizing</li> </ul>

		reports to compare model performance	models <ul style="list-style-type: none"> <li>Understanding where and how the predictive value addressed the root problem was key to selecting which metrics to measure which model outperformed</li> </ul>
<b>Model Explanation</b>	<b>How does your algorithm work?</b>	<p>All 3 of my models appeared to perform equally well on my accuracy metric. Initially, my Logistic Regression appeared to significantly outperform the Decision Tree and Random Forest models on Recall. However, given that I had set out maximizing Recall as a goal, I continued to tweek hyperparameters and re-evaluated my models. Accuracy metrics remained fairly consistent but by adjusting the max_depth on my Decision Tree I was able to keep my accuracy performance in line with previous test but maximize Recall score.</p> <p>For this reason, I chose my Decision Tree model.</p>	<p>Had I not prioritized my Recall score as a goal beforehand, I probably would have chosen my Logistic Regression after my initial tests. The Recall scores on the other models in that initial round badly underperformed so much I would have felt comfortable choosing my logit model at that point. In fact, I did not expect my next iterations to be fruitful for the DT &amp; RF models and was mainly trying to improve Recall on my Logit model.</p> <p>Recommended Good Practices:</p> <ul style="list-style-type: none"> <li>Test all models at each iteration</li> <li>Understand the issues you're trying to address by developing a predictive model and establish evaluation metric goals before evaluating models.</li> </ul>