

Last updated 1/18/2020

Created by Patrick de Guzman

To download: Toolbar > File > Download > (Desired File Type)

<http://patrickdeguzman.me/>

<input checked="" type="checkbox"/>	Stage	Steps	ADDITIONAL INFO	Useful Functions/Methods
<input type="checkbox"/>	<b>Prereq: Business &amp; Data Understanding</b>	Understand the data, your questions, and your goals <ul style="list-style-type: none"> <li>• Are you simply exploring the data?</li> <li>• Are you preparing it for machine learning?</li> <li>• Is it in a tabular format?</li> <li>• How many features should I expect?</li> </ul>	<ul style="list-style-type: none"> <li>• Get a Data Dictionary or schema if possible</li> <li>• Understand what rows represent in your data</li> <li>• Studying the dataset for 1-2 hours will save you a ton of headache, especially if the dataset has &gt;50 features</li> </ul>	
<input type="checkbox"/>	<b>I. Import Data &amp; Libraries</b>	Download the data and make it available in your coding environment	<ul style="list-style-type: none"> <li>• Import important libraries (pandas, numpy, matplotlib, seaborn, datetime), then import others as needed</li> <li>• Multiple datasets? Combine if you are concatenating (union). Otherwise, join when you understand them and are ready</li> </ul>	<ul style="list-style-type: none"> <li>• pd.concat</li> <li>• pd.merge</li> </ul>
<input type="checkbox"/>	<b>II. Exploratory Data Analysis</b>	Check for duplicates	<ul style="list-style-type: none"> <li>• We don't need to keep any rows that are pure duplicates of each other</li> </ul>	<ul style="list-style-type: none"> <li>• df.drop_duplicates()</li> </ul>
<input type="checkbox"/>		Separate Data Types (Take an inventory of what data types you have)	<ul style="list-style-type: none"> <li>• Numerical <ul style="list-style-type: none"> <li>- Discrete</li> <li>- Continuous</li> </ul> </li> <li>• Categorical <ul style="list-style-type: none"> <li>- Ordinal</li> <li>- Nominal</li> <li>- Binary</li> </ul> </li> <li>• Date/Time (time-stamps)</li> <li>• Text data (tweets/reviews)</li> <li>• Image</li> <li>• Sound</li> </ul>	<ul style="list-style-type: none"> <li>• df.select_dtypes(['object', 'bool'])</li> <li>• df.select_dtypes(['float', 'int'])</li> </ul>
<input type="checkbox"/>		Initial Data Cleaning <ul style="list-style-type: none"> <li>• Clean anything that would prevent you from exploring the data</li> </ul>	<p>Examples of things to consider...</p> <ul style="list-style-type: none"> <li>• Are there categorical columns that should be numerical?</li> <li>• Is the data in the first few rows consistent with the name of the feature?</li> <li>• Are there lists or dictionaries packed into one feature?</li> <li>• Are dates in the date data type?</li> </ul>	<ul style="list-style-type: none"> <li>• pd.Series.str.replace()</li> <li>• pd.Series.astype()</li> <li>• pd.Series.map()</li> <li>• pd.Series.apply()</li> <li>• lambda functions</li> <li>• pd.cut()</li> <li>• sklearn.preprocessing.MultiLabelBinarizer</li> <li>• pd.to_datetime()</li> </ul>
<input type="checkbox"/>		Visualize & Understand <ul style="list-style-type: none"> <li>• Understand how your data is distributed (numerical &amp; categorical)</li> <li>• How are the columns related? (Find correlations or other relationships)</li> <li>• Are there any outliers? Note them (but don't remove them yet!)</li> <li>• This can also be a good time to do any statistical tests (T-tests maybe?) if you're interested</li> </ul>	<p>Some ideas</p> <ul style="list-style-type: none"> <li>• Numerical: Histograms &amp; Scatter Plots</li> <li>• Categorical: Bar plots</li> <li>• Both: Box plots, violin plots, colored histograms</li> <li>• Date/Time: Line plots</li> </ul>	<ul style="list-style-type: none"> <li>• df.value_counts()</li> <li>• seaborn.distplot()</li> <li>• seaborn.countplot()</li> <li>• matplotlib.pyplot.bar()</li> <li>• seaborn.FacetGrid()</li> <li>• df.groupby()</li> <li>• scipy.stats.ttest_ind()</li> </ul>
<input type="checkbox"/>		Assess Missing Values ( <b>Don't fill/impute yet!</b> ) <ul style="list-style-type: none"> <li>• The goal here is to figure out your strategy for dealing with missing values since most ML algorithms cannot handle them.</li> <li>• You have 2 options: <b>impute/fill</b> them or <b>remove</b> them <ul style="list-style-type: none"> <li>- For <b>Imputing</b>: skip below under <b>IV</b> for some imputation strategies</li> <li>- For <b>Removing</b>: try your best to critically think if removing is the best option for you <ul style="list-style-type: none"> <li>◦ Are there many missing values in one column?</li> <li>◦ Are there many missing values in one row?</li> <li>◦ Is a row missing the column you want to predict?</li> </ul> </li> </ul> </li> </ul>	<p>Things to consider when working with missing data...</p> <ul style="list-style-type: none"> <li>• How many per row?</li> <li>• How many per column?</li> <li>• Are they encoded as something else?</li> </ul>	<ul style="list-style-type: none"> <li>• df.isna().any()</li> <li>• df.drop()</li> <li>• np.isinf()</li> </ul>
<input type="checkbox"/>	<b>III. Train/Test Split</b>	Set aside some data for testing.	Depending on size of your data, this can be anywhere between 80-90% train.	<ul style="list-style-type: none"> <li>• sklearn.model_selection.train_test_split</li> <li>• sklearn.model_selection.StratifiedShuffleSplit</li> </ul>
<input type="checkbox"/>		Dealing with Missing Data (Many options) <ul style="list-style-type: none"> <li>• Mean/Median/Mode</li> <li>• Find similar columns and fill</li> <li>• Fill with a unique value (like zero)</li> <li>• Predict Missing Values with ML <ul style="list-style-type: none"> <li>- KNN (categorical)</li> <li>- Linear Regression (numerical)</li> </ul> </li> </ul>	<p>The reason we want to deal with missing data after we've split our data is because we want to simulate real world conditions when we test as much as we can.</p> <p>Some ideas:</p> <ul style="list-style-type: none"> <li>• Are there rows or columns you're okay with dropping?</li> <li>• Can you infer the value from other columns?</li> <li>• Categorical: most frequent may be a good option</li> <li>• Numerical: mean or median may be good options</li> <li>• See IterativeImputer for one method of using ML to fill NA</li> </ul>	<ul style="list-style-type: none"> <li>• sklearn.impute.SimpleImputer</li> <li>• sklearn.impute.IterativeImputer</li> <li>• df.fillna()</li> </ul>

		<p>Feature Engineering</p> <ul style="list-style-type: none"> <li>What columns/features can you make to add value &amp; information to your data?</li> </ul>	<p>Some ideas</p> <ul style="list-style-type: none"> <li>Aggregations (across groups or dates)</li> <li>Ratios (divide)</li> <li>Interactions (multiply)</li> <li>Frequency (counts)</li> <li>Pull parts from dates (months/days/hours)</li> </ul>	<ul style="list-style-type: none"> <li>sum</li> <li>mean</li> <li>/ (divide)</li> <li>df.groupby</li> </ul>
<input type="checkbox"/>	IV. Prepare for ML	<p>Transform Data</p> <ul style="list-style-type: none"> <li>Numerical <ul style="list-style-type: none"> <li>Normalize or Standardize</li> <li>Log-transform</li> <li>Remove outliers</li> </ul> </li> <li>Categorical <ul style="list-style-type: none"> <li>One-hot encode (nominal)</li> <li>Label encoder (ordinal)</li> <li>Binarize (binary)</li> </ul> </li> <li>Text <ul style="list-style-type: none"> <li>Tokenize</li> <li>Stem/Lemma</li> <li>TF-IDF</li> <li>(and much more NLP techniques)</li> </ul> </li> </ul>	<p>Considerations:</p> <ul style="list-style-type: none"> <li>Numerical <ul style="list-style-type: none"> <li>Some ML models perform better when features are all on the same scale</li> <li>log-transforming can make numerical features seem more normal</li> <li>removing outliers may increase your models' performance</li> </ul> </li> <li>Categorical <ul style="list-style-type: none"> <li>Try to avoid using pd.get_dummies if you want to replicate the transformation you fit during training onto your testing set</li> <li>Use OneHotEncoder or other sklearn transformers instead</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>sklearn.preprocessing.StandardScaler</li> <li>sklearn.preprocessing.MinMaxScaler</li> <li>sklearn.preprocessing.normalize</li> <li>sklearn.preprocessing.LabelBinarizer</li> <li>sklearn.preprocessing.MultiLabelBinarizer</li> <li>sklearn.preprocessing.OneHotEncoder</li> <li>pd.get_dummies</li> <li>nlTK.tokenize.word_tokenize</li> <li>nlTK.corpus.stopwords</li> <li>nlTK.stem.porter.PorterStemmer</li> <li>nlTK.stem.wordnet.WordNetLemmatizer</li> <li>text.lower()</li> <li>text.split()</li> <li>sklearn.feature_extraction.text.CountVecorizer</li> <li>sklearn.feature_extraction.text.TfidfVecorizer</li> </ul>
<input type="checkbox"/>		<p>Feature Selection</p> <ul style="list-style-type: none"> <li>Numerical: Correlation (Pearson or Spearman) or ANOVA</li> <li>Categorical: Chi-Square test</li> <li>Domain Knowledge</li> <li>Recursive Feature Elimination (Like Forward Selection)</li> <li>Low importance features (calculated via permutation_importance or feature_importance)</li> </ul>	<p>Reducing dimensionality of your data can not only improve runtime, but also the quality of your predictions. Highly correlated or low variance features might work against you.</p> <ul style="list-style-type: none"> <li>Features you should consider removing... <ul style="list-style-type: none"> <li>Low variance (low variance = low information)</li> <li>One of two highly correlated features (maybe corr &gt; 0.95)? <ul style="list-style-type: none"> <li>Pearson, Spearman, or ANOVA F-value</li> </ul> </li> <li>If categorical, high Chi-Squared statistic</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>df.corr().abs()</li> <li>sklearn.feature_selection.VarianceThreshold</li> <li>sklearn.feature_selection.SelectKBest</li> <li>sklearn.feature_selection.chi2</li> <li>sklearn.feature_selection.f_classif</li> <li>sklearn.feature_selection.RFECV</li> </ul>
<input type="checkbox"/>	V. Pick your Models	<ul style="list-style-type: none"> <li>Some Regression Examples <ul style="list-style-type: none"> <li>Linear Regression</li> <li>Support Vector Regressor</li> <li>Random Forest</li> <li>Boosted Trees</li> <li>Neural Networks</li> </ul> </li> <li>Some Classification Examples <ul style="list-style-type: none"> <li>Support Vector Classifier</li> <li>Random Forest</li> <li>Logistic Regression</li> <li>Boosted Trees</li> <li>Neural Networks</li> </ul> </li> </ul>	Go wild.	
<input type="checkbox"/>	VI. Model Selection	Pick one algorithm via some form of Cross-Validation	Cross validation is a great way to estimate how your models will perform out in the wild.	<ul style="list-style-type: none"> <li>sklearn.model_selection.train_test_split</li> <li>sklearn.model_selection.KFold</li> <li>sklearn.model_selection.StratifiedKFold</li> </ul>
<input type="checkbox"/>	VII. Model Tuning	<p>Tune model hyperparameters</p> <ul style="list-style-type: none"> <li>Ideally use Cross-Validation again to choose your hyperparameters</li> </ul>	<p>Some examples you can use</p> <ul style="list-style-type: none"> <li>Grid Search</li> <li>Random Search (Faster Grid Search)</li> <li>Bayesian Optimization (Smarter Randomized Search)</li> </ul>	<ul style="list-style-type: none"> <li>sklearn.model_selection.GridSearchCV</li> <li>sklearn.model_selection.RandomizedSearchCV</li> <li>hyperopt library (Bayesian Optimization)</li> </ul>
<input type="checkbox"/>	VIII. Pick the best model	Pick the model that performed the best, and you're done!	Woohoo!	