

Last updated 2/27/2020

Created by Patrick de Guzman

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

<http://patrickdeguzman.me/>

	Stage	Steps	ADDITIONAL INFO	Useful Functions/Methods
<input type="checkbox"/>	Prereq: Business & Data Understanding	<u>Understand the data, your questions, and your goals</u> <ul style="list-style-type: none"> Are you simply exploring the data? Are you preparing it for machine learning? Is it in a tabular format? How many features should I expect? 	<ul style="list-style-type: none"> Get a Data Dictionary or schema if possible Understand what rows represent in your data Studying the dataset for 1-2 hours will save you a ton of headache, especially if the dataset has >50 features 	
<input type="checkbox"/>	I. Import Data & Libraries	Download the data and make it available in your coding environment	<ul style="list-style-type: none"> Import important libraries (pandas, numpy, matplotlib, seaborn, datetime), then import others as needed Multiple datasets? Combine if you are concatenating (union). Otherwise, join when you understand them and are ready 	<ul style="list-style-type: none"> pd.concat pd.merge
<input type="checkbox"/>		Check for duplicates	<ul style="list-style-type: none"> We don't need to keep any rows that are pure duplicates of each other 	<ul style="list-style-type: none"> df.drop_duplicates()
<input type="checkbox"/>		Separate Data Types (Take an inventory of what data types you have)	<ul style="list-style-type: none"> Numerical <ul style="list-style-type: none"> Discrete Continuous Categorical <ul style="list-style-type: none"> Ordinal Nominal Binary Date/Time (time-stamps) Text data (tweets/reviews) Image Sound 	<ul style="list-style-type: none"> df.select_dtypes(['object', 'bool']) df.select_dtypes(['float', 'int']) dtale.show() df.info()
<input type="checkbox"/>		Initial Data Cleaning <ul style="list-style-type: none"> Clean anything that would prevent you from exploring the data 	Examples of things to consider... <ul style="list-style-type: none"> Are there categorical columns that should be numerical? Is the data in the first few rows consistent with the name of the feature? Are there lists or dictionaries packed into one feature? Are dates in the date data type? 	<ul style="list-style-type: none"> pd.Series.str.replace() pd.Series.astype() pd.Series.map() pd.Series.apply() lambda functions pd.cut() sklearn.preprocessing.MultiLabelBinarizer pd.to_datetime()
<input type="checkbox"/>		II. Exploratory Data Analysis <ul style="list-style-type: none"> Visualize & Understand <ul style="list-style-type: none"> Understand how your data is distributed (numerical & categorical) How are the columns related? (Find correlations or other relationships) Are there any outliers? Note them (but don't remove them yet!) This can also be a good time to do any statistical tests (T-tests maybe?) if you're interested 	Some ideas <ul style="list-style-type: none"> Numerical: Histograms & Scatter Plots Categorical: Bar plots Both: Box plots, violin plots, colored histograms Date/Time: Line plots What data can tell you <ul style="list-style-type: none"> Change Over Time Hierarchy Drill Down Zoom in and out of granularity Contrasting Values Intersections Different Factors contributing to a larger phenomena Outliers Correlation 	<ul style="list-style-type: none"> df.value_counts() seaborn.distplot() seaborn.countplot() matplotlib.pyplot.bar() seaborn.FacetGrid() df.groupby() scipy.stats.ttest_ind()

<input type="checkbox"/>		<p>Assess Missing Values (Don't fill/impute yet!)</p> <ul style="list-style-type: none"> The goal here is to figure out your strategy for dealing with missing values since most ML algorithms cannot handle them. You have 2 options: <u>impute/fill</u> them or <u>remove</u> them <ul style="list-style-type: none"> For <u>Imputing</u>: skip below under IV for some imputation strategies For <u>Removing</u>: try your best to critically think if removing is the best option for you <ul style="list-style-type: none"> Are there many missing values in one column? Are there many missing values in one row? Is a row missing the column you want to predict? 	<p>Things to consider when working with missing data...</p> <ul style="list-style-type: none"> How many per row? How many per column? Are they encoded as something else? 	<ul style="list-style-type: none"> df.isna().any() df.drop() np.isinf()
<input type="checkbox"/>	III. Train/Test Split	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"> sklearn.model_selection.train_test_split sklearn.model_selection.StratifiedShuffleSplit
<input type="checkbox"/>	IV. Prepare for ML	<p>Dealing with Missing Data (Many options)</p> <ul style="list-style-type: none"> Mean/Median/Mode Find similar columns and fill Fill with a unique value (like zero) Predict Missing Values with ML <ul style="list-style-type: none"> KNN (categorical) Linear Regression (numerical) Multiple Imputation or MICE for advanced methods Maximum Likelihood Estimation 	<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"> Are there rows or columns you're okay with dropping? Can you infer the value from other columns? Categorical: most frequent may be a good option Numerical: mean or median may be good options See IterativeImputer for one method of using ML to fill multiple NA values <ul style="list-style-type: none"> Key tradeoff between ML imputation and simple imputation... <ul style="list-style-type: none"> ML imputation gives you greater variability and precision in your features Simple imputation is much easier and less costly in production 	<ul style="list-style-type: none"> sklearn.impute.SimpleImputer sklearn.impute.IterativeImputer df.fillna() fancyimpute.IterativeImputer
		<p>Feature Engineering</p> <ul style="list-style-type: none"> What columns/features can you make to add value & information to your data? 	<p>Some ideas</p> <ul style="list-style-type: none"> Aggregations (across groups or dates) Ratios (divide) Interactions (multiply) Frequency (counts) Pull parts from dates (months/days/hours) 	<ul style="list-style-type: none"> sum mean / (divide) df.groupby
<input type="checkbox"/>		<p>Transform Data</p> <ul style="list-style-type: none"> Numerical <ul style="list-style-type: none"> Normalize or Standardize Log-transform Remove outliers Categorical <ul style="list-style-type: none"> One-hot encode (nominal) Label encoder (ordinal) Binarize (binary) Text <ul style="list-style-type: none"> Tokenize Stem/Lemma TF-IDF (and much more NLP techniques) 	<p>Considerations:</p> <ul style="list-style-type: none"> Numerical <ul style="list-style-type: none"> Some ML models perform better when features are all on the same scale log-transforming can make numerical features seem more normal removing outliers may increase your models' performance Categorical <ul style="list-style-type: none"> Try to avoid using pd.get_dummies if you want to replicate the transformation you fit during training onto your testing set Use OneHotEncoder or other sklearn transformers instead 	<ul style="list-style-type: none"> sklearn.preprocessing.StandardScaler sklearn.preprocessing.MinMaxScaler sklearn.preprocessing.normalize sklearn.preprocessing.LabelBinarizer sklearn.preprocessing.MultiLabelBinarizer sklearn.preprocessing.OneHotEncoder pd.get_dummies nltk.tokenize.word_tokenize nltk.corpus.stopwords nltk.stem.porter.PorterStemmer nltk.stem.wordnet.WordNetLemmatizer text.lower() text.split() sklearn.feature_extraction.text.CountVectorizer sklearn.feature_extraction.text.TfidfVectorizer
<input type="checkbox"/>		<p>Feature Selection</p> <ul style="list-style-type: none"> Numerical: Correlation (Pearson or Spearman) or ANOVA Categorical: Chi-Square test Domain Knowledge Recursive Feature Elimination (Like Forward Selection) Low importance features (calculated via permutation_importance or feature_importance) 	<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"> Features you should consider removing... <ul style="list-style-type: none"> Low variance (low variance = low information) One of two highly correlated features (maybe corr > 0.95)? <ul style="list-style-type: none"> Pearson, Spearman, or ANOVA F-value If categorical, high Chi-Squared statistic 	<ul style="list-style-type: none"> df.corr().abs() sklearn.feature_selection.VarianceThreshold sklearn.feature_selection.SelectKBest sklearn.feature_selection.chi2 sklearn.feature_selection.f_classif sklearn.feature_selection.RFECV

<input type="checkbox"/>	V. Pick your Models	<ul style="list-style-type: none"> Some Regression Examples <ul style="list-style-type: none"> - Linear Regression - Support Vector Regressor - Random Forest - Boosted Trees - Neural Networks Some Classification Examples <ul style="list-style-type: none"> - Support Vector Classifier - Random Forest - Logistic Regression - Boosted Trees - Neural Networks 	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"> <code>sklearn.model_selection.train_test_split</code> <code>sklearn.model_selection.KFold</code> <code>sklearn.model_selection.StratifiedKFold</code> <code>yellowbrick.classifier.roc_auc</code> <code>yellowbrick.classifier.ClassificationReport</code> <code>yellowbrick.regressor.ResidualsPlot</code>
<input type="checkbox"/>	VII. Model Tuning	Tune model hyperparameters <ul style="list-style-type: none"> Ideally use Cross-Validation again to choose your hyperparameters 	Some examples you can use <ul style="list-style-type: none"> Grid Search Random Search (Faster Grid Search) Bayesian Optimization (Smarter Randomized Search) Also identify a good decision boundary (AKA discrimination threshold) if using classification <ul style="list-style-type: none"> Can be done with Yellowbrick's quick <code>DiscriminationThreshold</code> viz 	<ul style="list-style-type: none"> <code>sklearn.model_selection.GridSearchCV</code> <code>sklearn.model_selection.RandomizedSearchCV</code> <code>hyperopt</code> library (Bayesian Optimization) <code>yellowbrick.classifier.DiscriminationThreshold</code>
<input type="checkbox"/>	VIII. Pick the best model	Pick the model that performed the best, and you're done!	Woohoo!	