

Kickstarter is a funding platform where their main goal is to transform ideas into projects. Each project is created and crafted by the person behind it where once they're ready, they can launch their project to the community. Every project creator sets their funding goal and deadline, and if the project gains enough interest, people can pledge money to make it happen. This project is divided into 3 main parts; regression, classification and clustering which will be covered in detail.

## Data Preprocessing

As a starting point, I imported the clean dataset of 18568 observations and 45 columns and kept only the observations where the project state is either successful or failed, reducing the amount of observations to 15685. The following tables shows the predictors that were dropped:

Predictor	Reason Behind Removal
<b>Project_ID, name</b>	Used for identification, have no effect on the model
<b>Goal, static_usd_rate</b>	A new variable was created that multiplies both for a more scalable comparison
<b>Pledged</b>	After the project was submitted
<b>State</b>	After the project was submitted (For regression)
<b>Disable_communication</b>	After the project was submitted
<b>Currency</b>	Highly correlated to country
<b>Deadline, State_changed_at, created_at, launched_at</b>	Their respective dummies were used
<b>Backers_count</b>	After the project was submitted
<b>Spotlight</b>	After the project was submitted
<b>Staff_pick</b>	After the project was submitted
<b>Name_len, blurb_len, Name_len_clean, blurb_len_clean</b>	Due to feature selection they were dropped as they decreased model's performance (they were kept for classification)
<b>State_changed_at_weekday, month, day, year, hour</b>	After the project was submitted
<b>Launch_to_state_change_days</b>	After the project was submitted
<b>Usd_pledged</b>	After the project was submitted (for classification task)

After dropping the previous predictors, dummies were created for the following categorical variables: country, category, deadline\_weekday, created\_at\_weekday, created\_at\_yr, deadline\_yr, deadline\_month, created\_at\_month. The next part of the process involved the removal of anomalies using Isolation Forest. In the regression part, it was realized that when removing outliers for the whole dataset, the target variable once plotted its boxplot still had extreme outliers. Hence, two different Isolation Forest models were built for the target variable and for the predictors respectively, reducing the total amount of observations to 13793. As for the clustering part, anomaly detection was performed on my 3 predictors. To further reduce the number of input variables, feature selection was performed for both regression and classification tasks, using Lasso and Logistic feature selection techniques.

## Model Selection

For both regression and classification tasks, the optimal model was selected based on the least MSE and highest accuracy using cross-validation. Different models were tested such as Lasso, Ridge, Multiple Linear Regression, KNN, ANN, Random Forest and Gradient Boosting. However, in both models Gradient Boosting was selected. The optimal hyperparameters were selected by looping through each hyper parameter. The tables below represent the results:

Regression:

Technique	Ridge	Lasso	Multiple Linear Regression	KNN	ANN	Random Forest	Gradient Boosting
<b>MSE (in billion)</b>	1.322	1.324	1.321	1.354	1.31	1.217	1.206

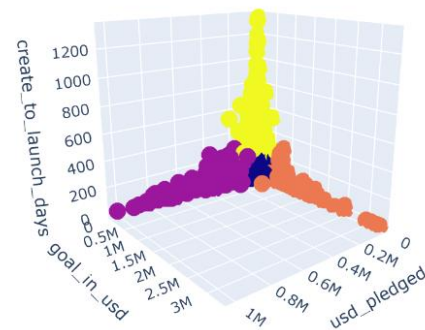
Classification:

Technique	KNN	ANN	Random Forest	Gradient Boosting
<b>CV Score</b>	0.6994	0.6719	0.7473	0.758

	Regression	Classification
<b>Max features</b>	Auto	Auto
<b>Max depth</b>	2	6
<b>Minimum samples split</b>	42	9
<b>Min Samples leaf</b>	13	3
<b>N_estimators</b>	90	250

From a business perspective, the prediction model requires more variables to provide a higher and more reliable model through for example an industry analysis of each category. As for the classification model it can be recommended due to its ability to classify the state of a project at a 75.6% accuracy.

For the clustering model, it was performed on USD\_pledged, Goal in USD and Create to lunch days. After dropping NAs and removing anomalies using Isolation Forest, K-Means Clustering was selected, where using elbow method the optimal number of clusters was 4. As we can see from the 3D scatter using Plotly, it is easy to visualize our 4 clusters. Yellow Cluster represents creators who took too long to launch with low goals and ended with low USD\_pledged. Purple cluster represents creators who had low to intermediate goals and time to launch but exceeded their own expectations and ended up with a high USD\_pledged. Orange Cluster represents creators with intermediate-high goals, low-intermediate time to launch and ended with low USD\_pledged. Finally, the blue cluster represents creators with low goals, launched quickly and ended with very low USD\_pledged.



The model was then tested using silhouette score (0.775), Calinski-Harabasz score (8551.06) and p-value ( $1.11e-16$ ), meaning that our data was well clustered.