

## Kaggle Model Description

**Team: Kitty - Members: Minjie Yang (Leader), Yu Ren, Peter Xie**

**One attempt is to encode metal into its corresponding density** in order to give it both a numeric value and an order. Doing this could reach a score around 300 on Kaggle.

**Another attempt is to split up the dataset according to the cost column** for the following motivation. Doing this could reach a score around 200 on Kaggle.

The toy's **weight** can be evaluated in the following two ways:

**1.weight = cost/metal\_cost**

**2.weight = density\*volume (realized by the helpful function calculate\_weight)**

After analyzing the dataset and running some experiments, we find out that the first way has a high accuracy in predicting **weight** (RMSE on the holdout set is around 5), while the second way has a low accuracy (around 400, and this result is about the same among all three shapes). Unfortunately, there is a lot of missing value in the **cost** column. Simply encoding the missing data of the **cost** column into its median will lead to a bad performance, because it will destroy the nice feature of it. So in order to good full use of **cost**, we separate the dataset into two parts:

**If the cost value is complete: com\_train and com\_test**

**If the cost value is missing: mis\_train and mis\_test**

On the **com\_train** and **com\_test**, we add an approximation column predicting the **weight** by the first way to enhance the performance. On the **mis\_train** and **mis\_test**, we add an approximation column predicting the **weight** by the second way.

Eventually, we choose the two predictions on **com\_test** and **mis\_test** from the two models with the best holdout score respectively. Then we concatenate them together into one final prediction **df\_submit**.

So the code of the model will be organized by:

**0.com\_train, com\_test, mis\_train and mis\_test preparing**

**1.com\_train and com\_test feature encoding > com\_train training > com\_test predicting : df\_com**

**2.mis\_train and mis\_test feature encoding > mis\_train training > mis\_test predicting : df\_mis**

**3.concatenate df\_com and df\_mis together into df\_submit**

More information and remarks are written on 'Kaggle\_Kitty.ipynb'.