

---

# Solution - Big Buck Challenge by IEOR@IITB and McKinsey Knowledge Center India

Presented by  
- Aakash Kerawat  
Final year UG Student  
IIT Roorkee

---

---

# Outline

## 1. Predicting Withdrawal

- a. First thought
- b. Data preprocessing
- c. EDA
- d. Feature Engineering
- e. Model Selection

## 2. Optimizing replenishment

- a. Cost function
- b. Optimization
- c. Plots

---

# Predicting Withdrawal

---

# First Thought

Looking at the data there can be two approaches -

1. Time series forecasting
2. Using machine learning models.

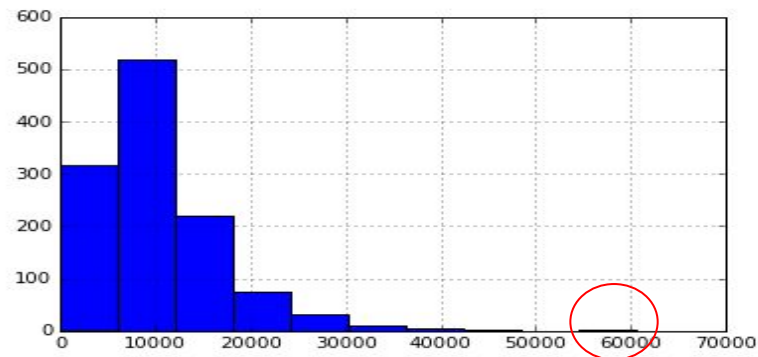
Chose option 2 as -

- A. There was missing data.
- B. ML models can perform at par (or better) than TS models with right feature set.

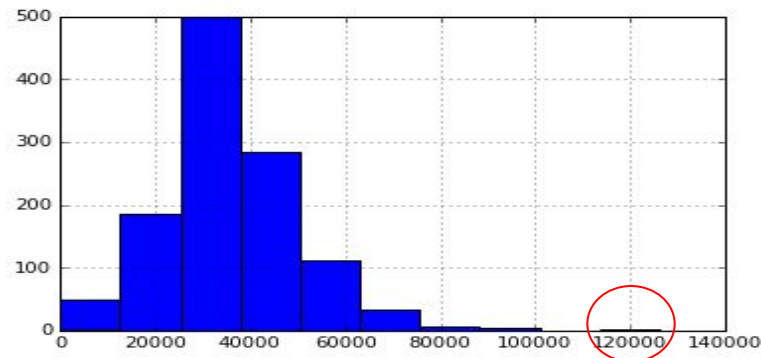
# Data Preprocessing

- Distribution of ATM\_ID specific withdrawal indicated high-value outliers.
- Removed the data points with withdrawal value > 99 percentile of that ATM\_ID..

```
df_train.loc[df_train.ATM_ID=='SRN00279', :].Withd  
plt.show()
```



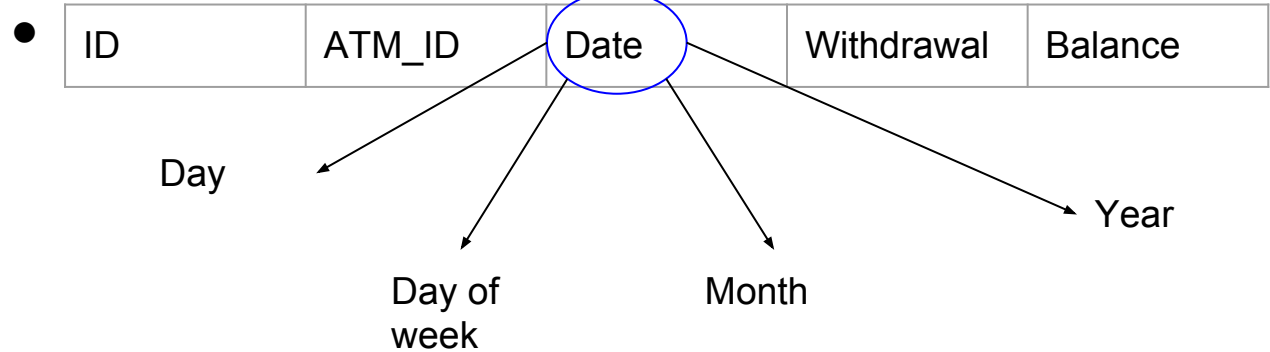
```
] : df_train.loc[df_train.ATM_ID=='SRN004989', :].Withd  
plt.show()
```



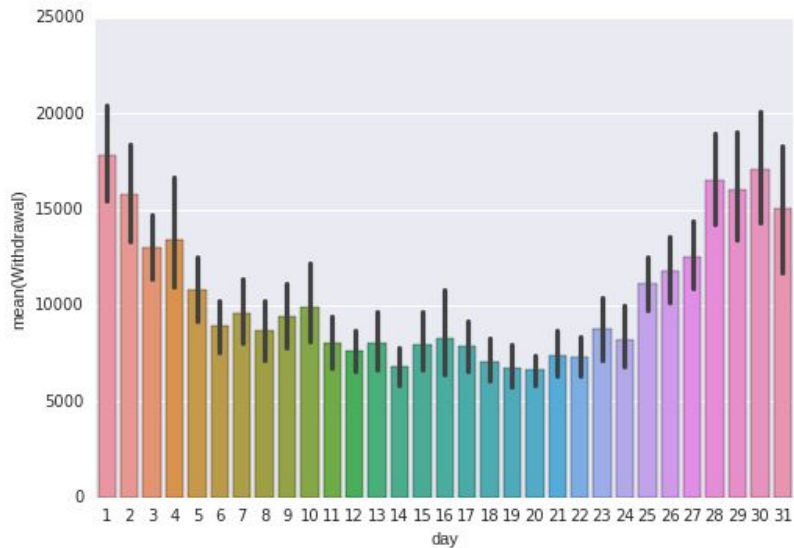
---

# EDA

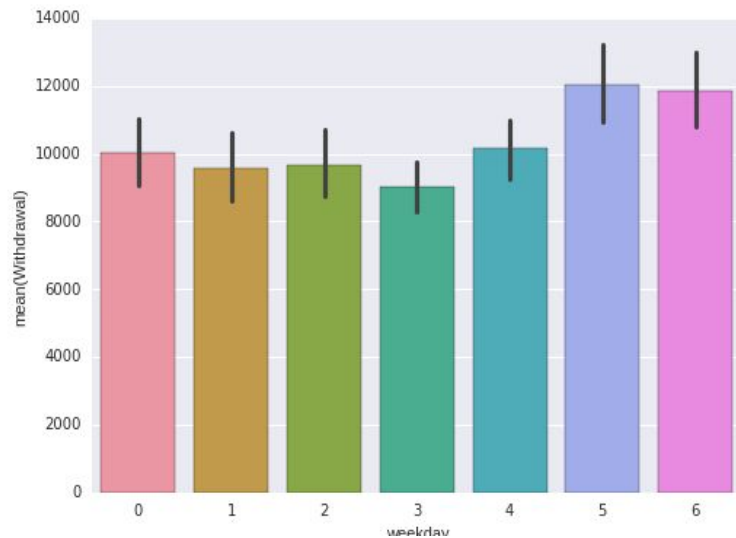
- The given data -



```
sns.barplot(df_train.loc[df_train.ATM_ID=='SRN000279', 'day'],  
plt.show())
```

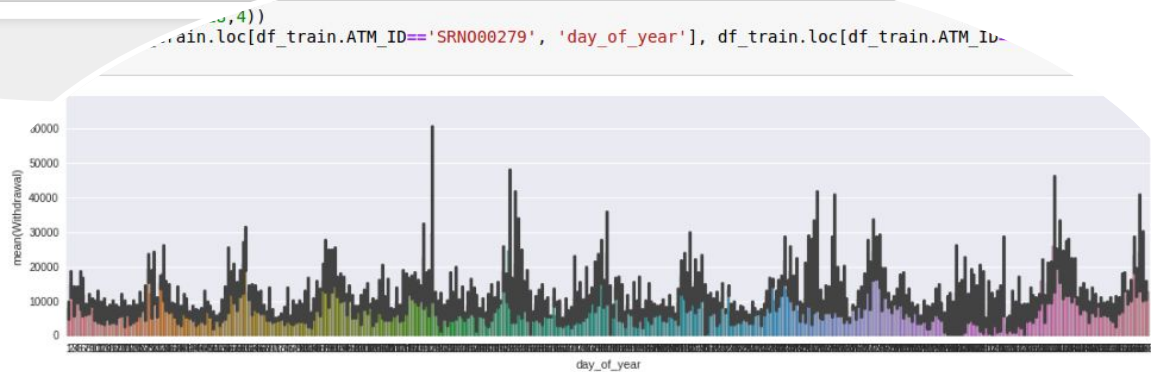
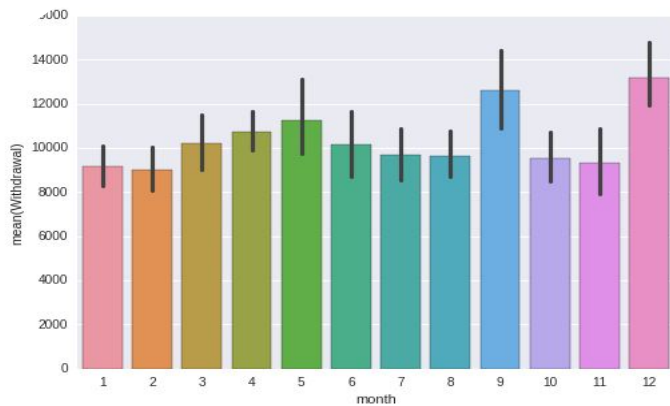


```
sns.barplot(df_train.loc[df_train.ATM_ID=='SRN000279', 'weekday'],  
plt.show())
```



## Trends

# Trends





# Final Date features

- Day
- Day of week
- Month
- Day of year
- ~~Year~~
  - Removed as test data was only of 1 particular year.
  - Also there wasn't a significant trend.

# Additional Features

Merged the ATM\_info data which  
provided -

- Type
- Facility

---

# Engineered Features

Motivation - Initially thought of making different model for each ATM, but this has many limitations.

Following features were added -

- Withdrawal\_mean - Mean value of withdrawal amount for a particular ATM\_ID
- Withdrawal\_std - Std. deviation value of withdrawal amount for a particular ATM\_ID

# Engineered Features

These features gave a sense of distribution of withdrawal for that ATM to the model as a continuous value feature.

- Withdrawal\_uq - Upper quartile of withdrawal amount for a particular ATM\_ID
- Withdrawal\_lq - Lower quartile value of withdrawal amount for a particular ATM\_ID

---

# Model Selection

- Tried a range of algorithms from Linear regression, Decision Tree regressor and Random Forest regressor to Gradient Boosting regressor.
- Chose Gradient boosting as the final model.

# Reasons

- All the date features can be considered as categorical and one-hot encoding would result in more than 300 features.
- Data had both categorical and continuous features.
- Tree based algorithms are good at handling both at the same time.

---

# Reasons

- Tree based algorithms can learn from categorical features without one-hot encoding.
- **Another major reason** - Checked the cross-validation score for all the models and gradient boosting came out to be the best.

---

# Optimizing Replenishment

---



---

---

# Approach

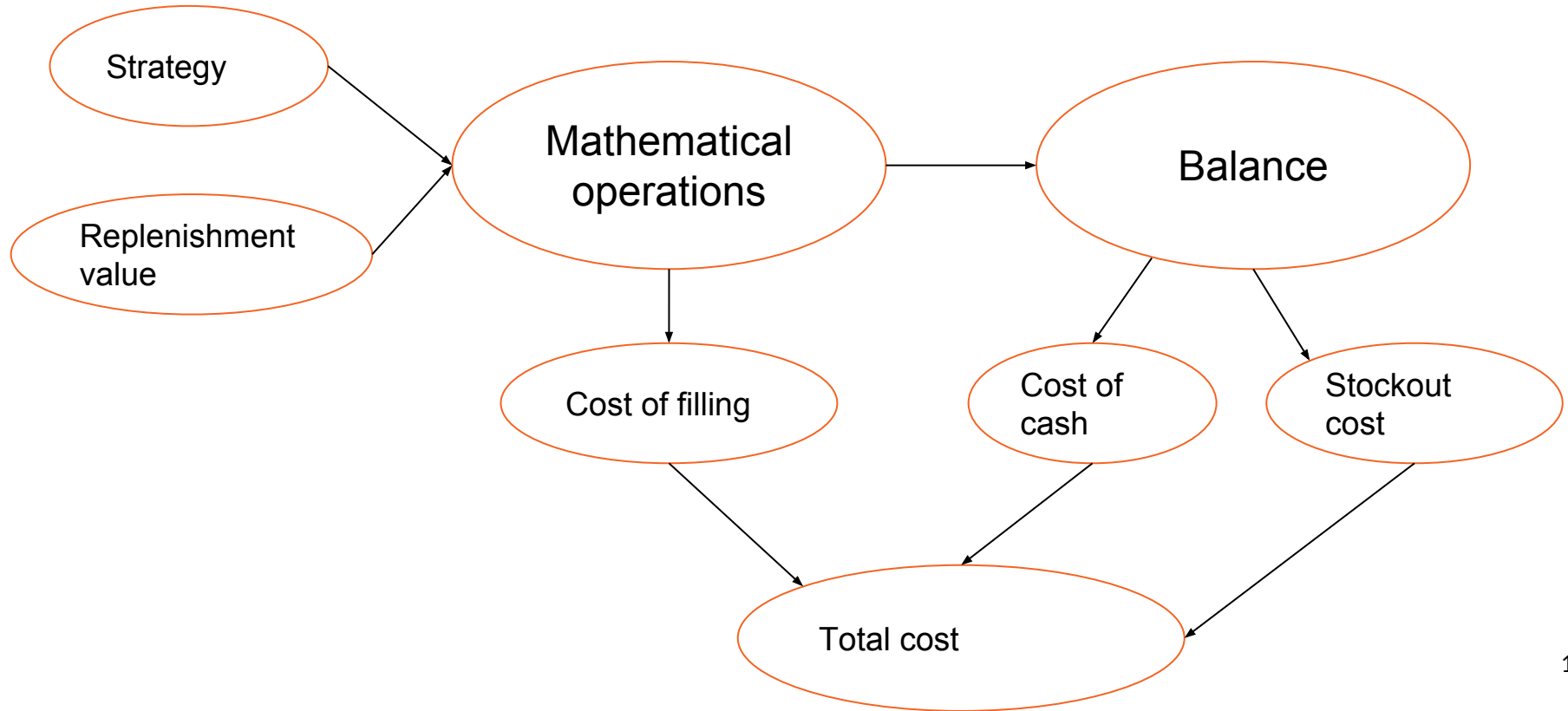
Calculated the total cost for each strategy for a range of replenishment values and finally chose the pair for which the cost was minimum.

# Defined 2 functions

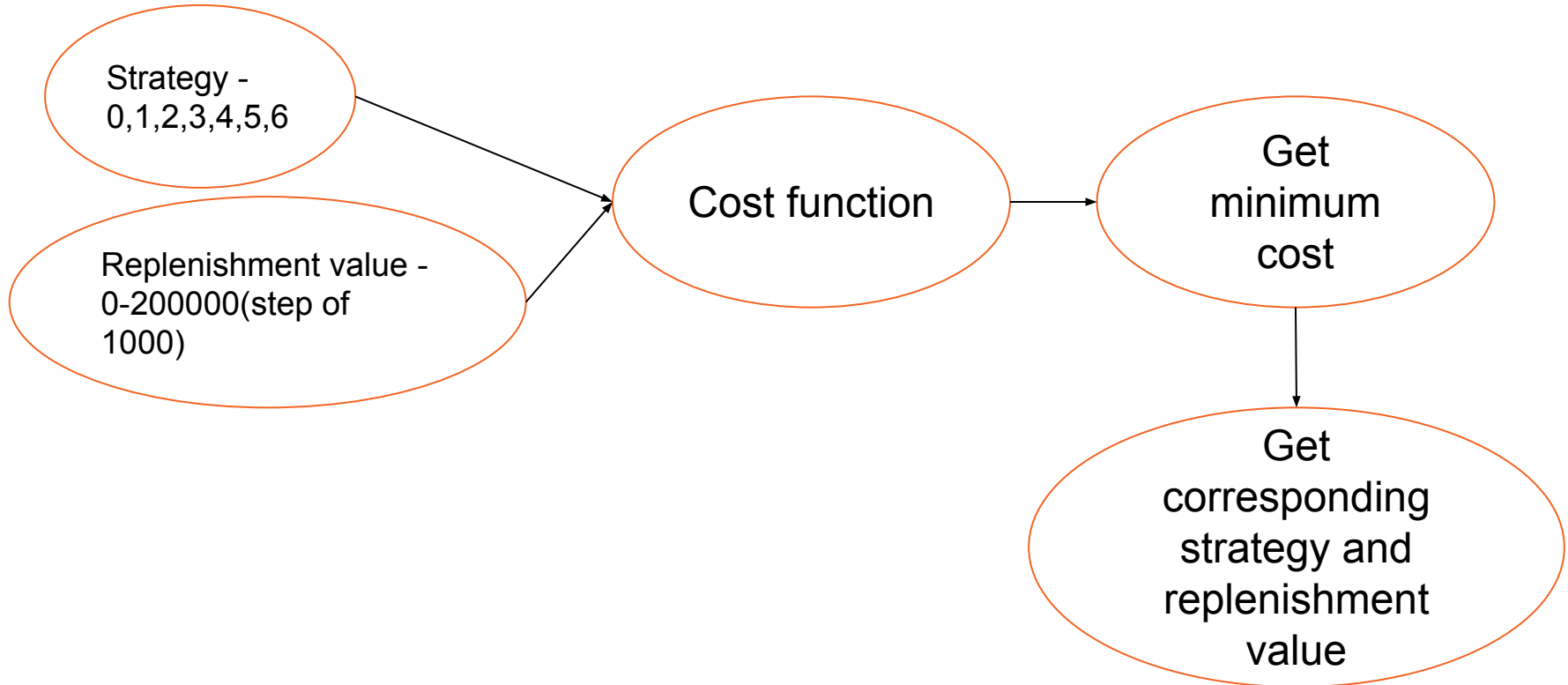
- Cost function
- Optimization function

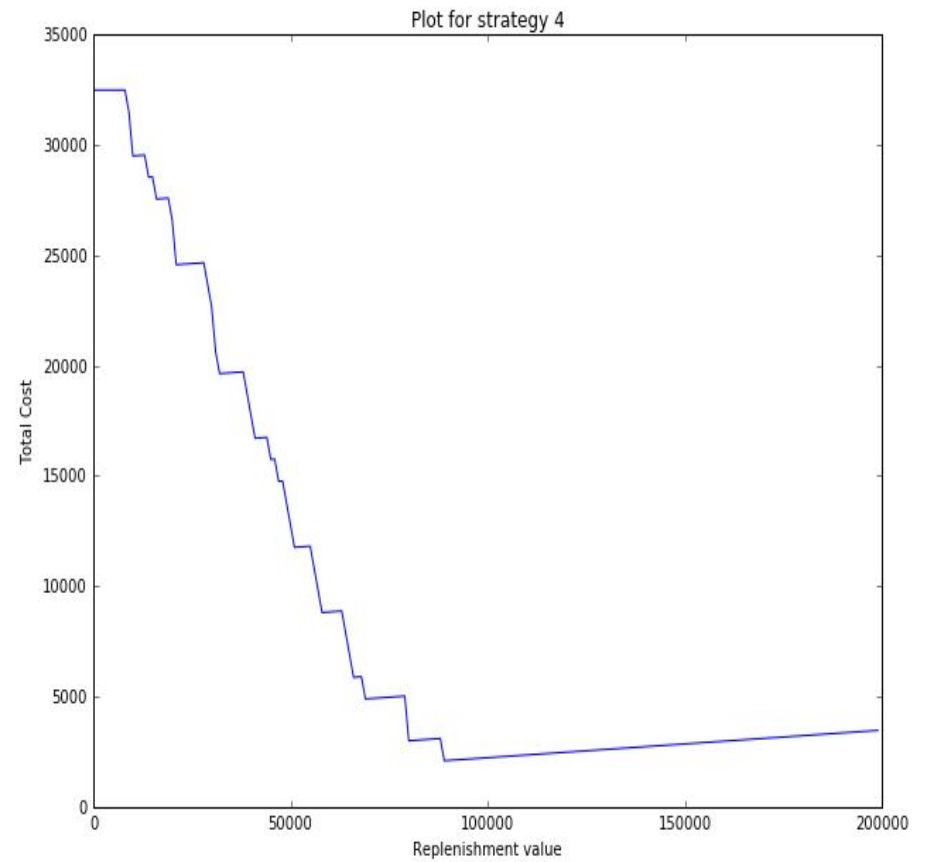
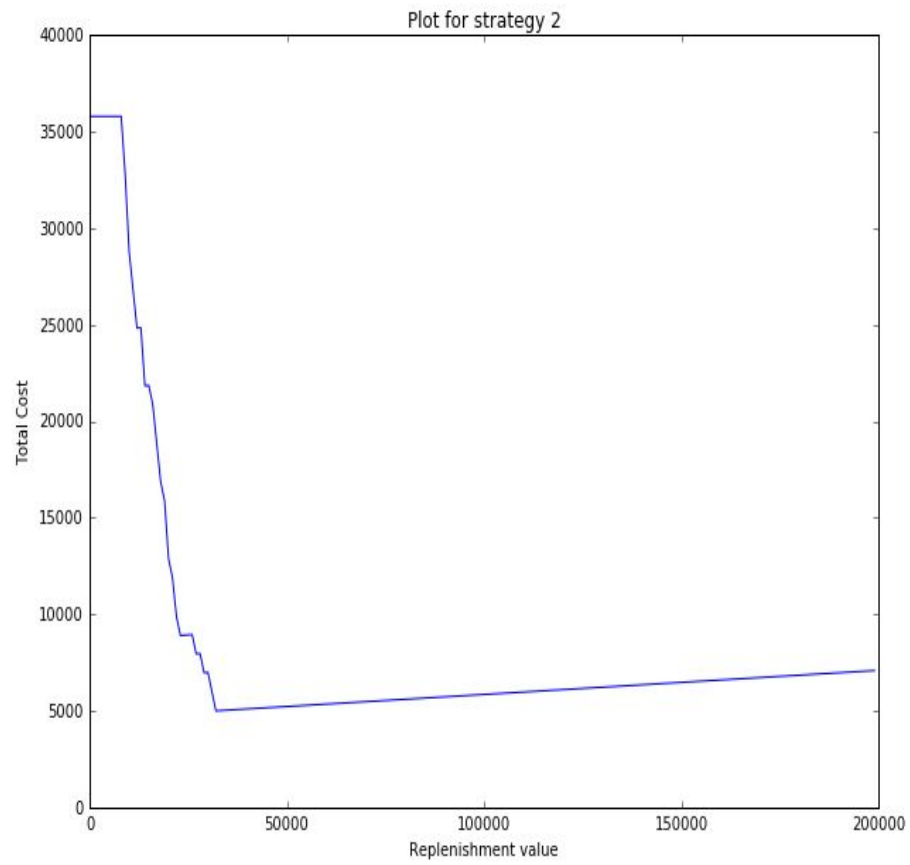
---

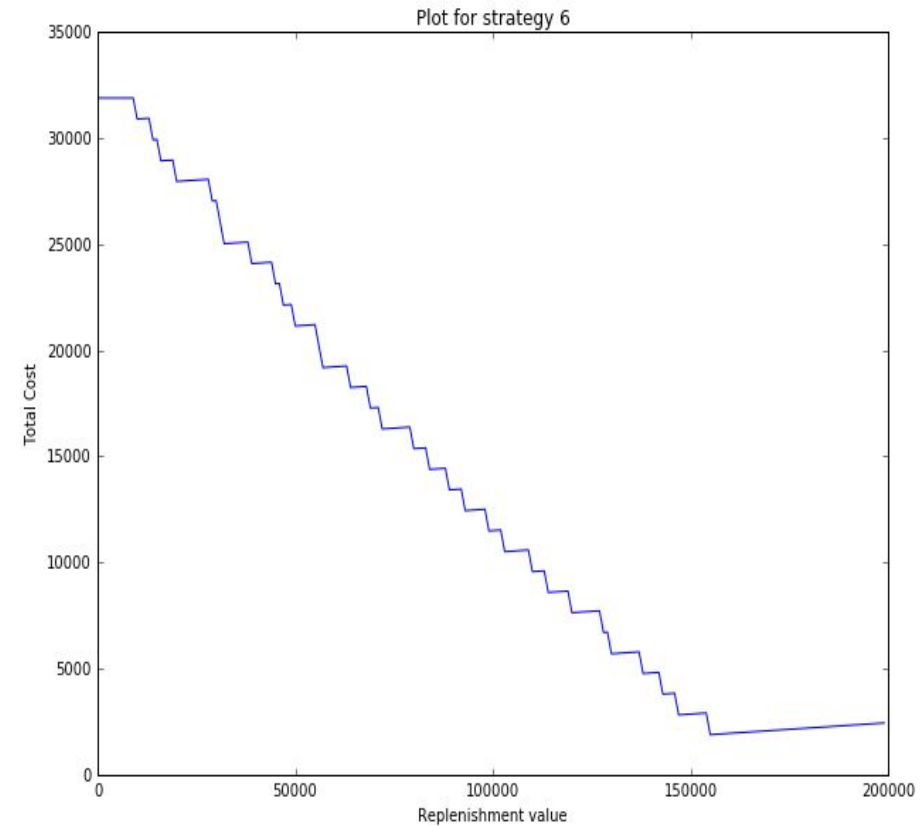
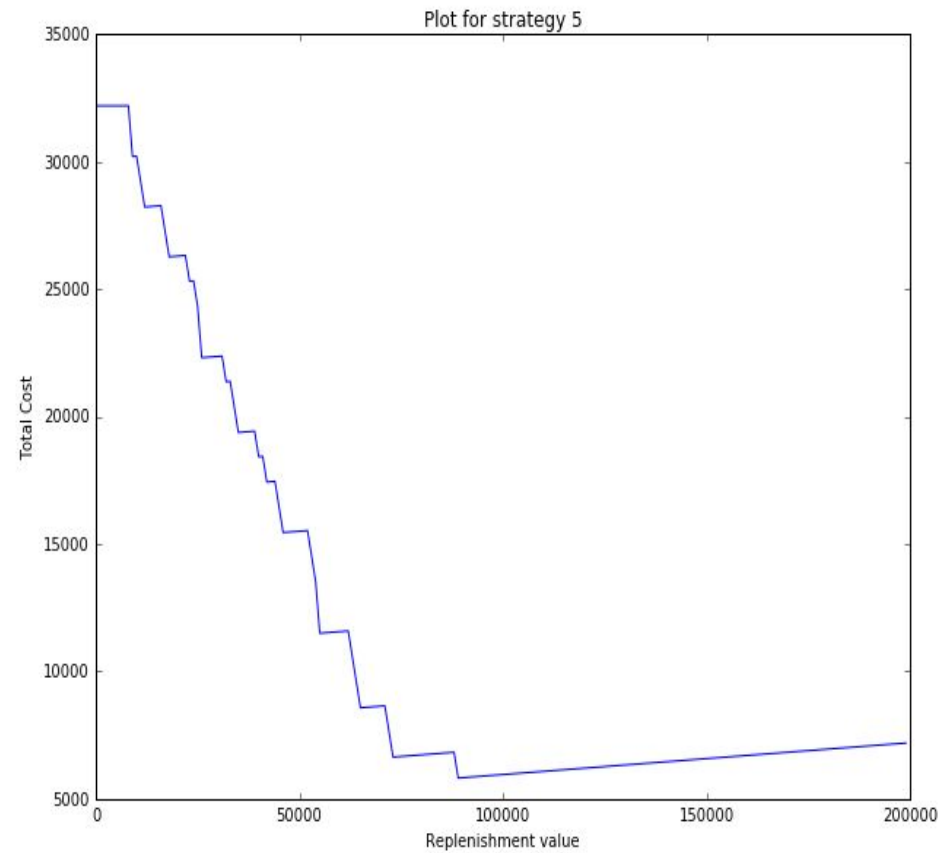
# Cost function



# Optimize







---

# Thank You