

# Problem Identification

Greenfield Refinery Feasibility Validation Presentation

# Greenfield Refinery Problem Statement

**Develop a data-driven methodology to validate economic feasibility study of building and operating a small (under 50kBpd) crude oil refinery in the U.S., using datasets from the U.S. Census and EIA.gov websites. Identify 2-3 key parameters that have high impact on the economic performance of a refinery, located in the pre-selected PAD-2-2 district, optimizing for anticipated gross margin (\$/Bbl crude oil refined), based on at least 8 years of historical data**

## Context

Green Light Co. (GLC), a California startup, wishes to validate its preliminary economical model for building a greenfield crude oil refineries in the U.S. It recognizes that no new refinery with a significant downstream capacity has been built in the U.S. in 40 years; the existing refinery plants are relying on outdated technology to produce petroleum products.

GLC preliminary has selected several potential locations near major sources of domestic crude oil supply and theorizes that these locations will bring the highest profit margins. I developed a methodology that could be used to forecast the margins, based on the geographical location of a proposed refinery site, quality of crude oil (feedstock), plant equipment configuration, and proximity to destination markets.

## Criteria for success

The criterion for success is to validate the anticipated gross profit margins in the range of \$25-\$30 per 1 U.S. barrel of crude oil processed, given the pre-selected sites for new refinery construction, type of crude available, and proximity to destination markets.

## Scope of solution space

Historical data on energy production and consumptions - we will use a 5-year datasets:

## Scope of solution space

Crude Oil production + Net Imports and Crude Oil price indices - WTI for the U.S., and Brent for the global markets, to determine correlation with the prices of petroleum products produced by refineries and net imports

Historical refinery operations parameters: average plant utilization rates across the U.S., product yields by input crude and refinery configuration types, and volumes produced - to determine indicative profit margins for the existing refineries and build comparative analysis

## Constraints within solution space

Some of the constraints on data availability and level of accuracy I had to accept were as follows:

- Much of the actual refinery sale prices for ready products is not publicly available data; there are agencies that sell historical petroleum pricing data, but their costs were prohibitive for an academic exercise; we will just have to rely on the EIA.gov data that is the next best thing
- There's uncertainty on the gasoline and possibly, diesel demand, not due to politics and government, but the free market, i.e., an anticipated proliferation of Electrical Vehicles within the next 10-15 years – the minimum runway for refinery project investment returns
- Finally, the COVID-19 related lockdowns have demonstrated that a significant office-based workforce might successfully transition to remote work attendance practices, leading to lower commute and possibly, lower consumption of petroleum-based fuels

# Greenfield Refinery Problem Statement

**Develop a data-driven methodology to validate economic feasibility study of building and operating a small (under 50kBpd) crude oil refinery in the U.S., using datasets from the U.S. Census and EIA.gov websites. Identify 2-3 key parameters that have high impact on the economic performance of a refinery, located in the pre-selected PAD-2-2 district, optimizing for anticipated gross margin (\$/Bbl crude oil refined), based on at least 8 years of historical data**

## Constraints within solution space

- Great level of uncertainty in the future of gasoline markets due to politically driven decisions by the U.S. government, which may have a negative impact on the U.S. energy sector, e.g., the 'Carbon Free by 2035' plan

I used datasets obtained from the US Census and EIA.gov websites to identify key metrics:

## Key data sources

- major metro areas and specific region population
- historical energy production and consumption rates
- historical energy source composition dynamics (oil & gas, solar, hydro, nuclear, etc.)
- refinery operational throughput dynamics (to determine trends over the last 10 years)
- refinery product yields and utilization rates - to assess their historical profitability
- crude oil actual (refinery) crack-spreads by type of refining configurations,
- finally, predictive price modeling for ready products, in relation to crude oil price indices

I abandoned some items, as they were either not very useful for narrowing down the solution scope, or were becoming too large of a task, i.e., outside of the stated scope:

- (abandoned) Make-up of energy produced and consumed by source - to determine where the sustainable production of crude oil based refineries are the strongest and to analyze the patterns of consumption by regions
- (abandoned) Amount of energy generated, as well as its net imports into the U.S., measured in British Thermal Units
- (abandoned) Amount of energy consumed in the U.S., also measured in British Thermal Units

# Recommendation and Key Findings

Greenfield Refinery Feasibility Validation Presentation

# Key Findings and Recommendations

## Key Findings

The key findings, which were determined upon analyzing the data and running models, are as follows:

- Brent Differential has a stronger impact on the economics than the domestic WTI crude price index. There is also an indication that WTI is dependent on the Brent.
- This was expected, as the heavier crude doesn't automatically mean less valuable products may be extracted from it than the lighter crude – it is mainly a function of refinery process type and the equipment. However, there is more lighter product yields in the lighter crude oil. It is easier to extract than from the heavy crude. This effect caps out at the lower and higher ranges, in this case, the 39-42 API range.
- The model did not rank any of the one-hot encoded 'PAD' district features very well in the 'feature importance' assessment. The reality though supports the theory that the PAD location is a very important factor in economic feasibility of a crude oil refinery – an indication that the model definitely needs further tweaking

## Recommendations

Continue working to produce a useful application for the industry. There is plenty of space for improvement – work on both the data and the models.

# Modeling Results and Analysis

Greenfield Refinery Feasibility Validation Presentation

# Modeling Results and Analysis

## Exploratory Data Analysis

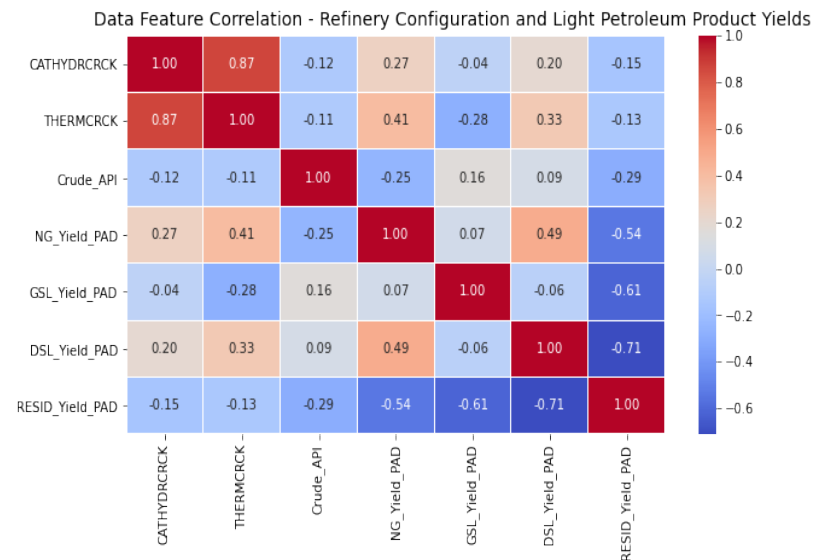
During the exploratory data analysis, it becomes clear that there is relationship between plant configuration, Crude API and the Product Yields.

The analysis showed strong correlation of light product yields with the following data features:

- *Catalytic Hydrocracking*
- *Thermal Cracking*
- *Crude API*

The heatmap and a series of scatter plots further confirmed this pattern.

## Heatmap of features



# Modeling Results and Analysis

## Linear Regression Model

The Linear Model further supported the patterns, initially found during the exploration.

The linear regression model found that there is a strong correlation between the Crack Spreads and the following features:

- *Brent Price*
- *Light Product Yields (Gasoline, Diesel)*
- *Crude API*
- *Crude Price Differentials*

## Linear Regression Model Results

### Linear regression model performance

```
lr_mae_mean = np.mean(-1 * lr_neg_mae['test_score'])
```

```
lr_mae_std = np.std(-1 * lr_neg_mae['test_score'])
```

```
(1.0908461747848313, 0.04277474406145815)
```

```
mean_absolute_error(y_test, lr_grid_cv.best_estimator_.predict(X_test))
```

```
1.1543230831643632
```



# Modeling Results and Analysis

## Random Forest Model

The Random Forest Model also returned top features that were common with the Linear Model:

- *Brent Price Differential*
- *Product Yield*
- *Crude API*

## The Random Forest Model performance

```
rf_mae_mean = np.mean(-1 * rf_neg_mae['test_score'])
```

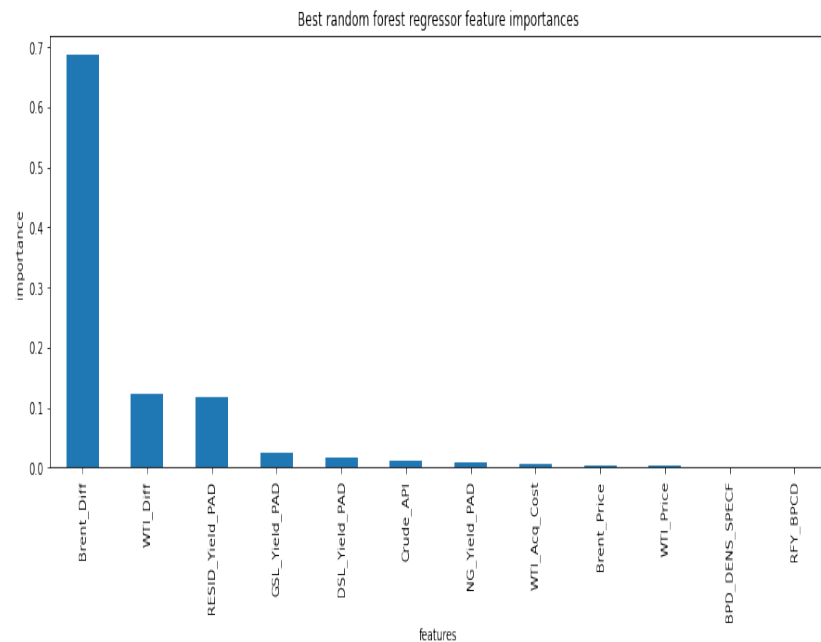
```
rf_mae_std = np.std(-1 * rf_neg_mae['test_score'])
```

```
(0.1334513516888089, 0.05712131291841946)
```

```
mean_absolute_error(y_test, rf_grid_cv.best_estimator_.predict(X_test))
```

```
0.0806464774716703
```

## Random Forest Model Results



# Modeling Results and Analysis

## Analysis

After comparing the cross-validation MAE values for the two models, we concluded that the random forest model was a better choice, because it showed lower CV MAE and less variability. The performance on the test set was also consistent with the cross-validation results.

## Data quantity assessment

Finally, we assessed whether the model would benefit from additional data collection by running the sklearn's `learning_curve` function.

The CV score and Training set size plot showed that no further data collection would be required. As the plot shows, the improvement in the model scores starts flattening out by around a sample size of 250



# Summary and Conclusion

Greenfield Refinery Feasibility Validation Presentation

# Summary and Conclusion

## Summary

The ultimate goal, as outlined in the Problem Statement was to validate the feasibility studies for building a new refinery in PAD-2-2 district by identifying 2-3 key parameters that have the highest impact on the Specific Crack Spread. I had already confirmed the key parameters – i.e., the top features by importance.

### Scenario 1:

Increasing WTI Discount from the Baseline \$(6.05)/Bbl to \$(7.5)/Bbl, while keeping the Brent Discount at zero, added a modest increase in the Crack Spread – from \$20.58/Bbl to \$22.99/Bbl, but still lower than the Baseline of \$25.97/Bbl

Increasing Brent Discount from the Baseline \$(2.08) to \$(5.5), while keeping the WTI Discount at zero had a significant impact, increasing the Crack Spread – to \$34.66/Bbl

**Scenario 2:** Changing the Baseline value of the Crude API from the heavier 26.11 to the lighter 39 resulted in a slightly better outcome – the Crack Spread went from \$25.97/Bbl to \$26.03/Bbl. Further increase of Crude API to 41 and 42 did not add anything – the Crack Spread stayed at the \$26.03/Bbl level

## Summary

**Scenario 3:** Moving the refinery location from PAD-2-2 to PAD-4 resulted in a very negligible change in the predicted Crack Spread – it went from \$25.97/Bbl to \$25.99/Bbl. Similarly, changing the location from PAD-2-2 to PAD-3-3 did not have any further affect

## Conclusion

The conclusion here is that Brent Differential has a stronger impact on the economics than the domestic WTI crude price index. There is also an indication that WTI is dependent on the Brent.

The model did not rank any of the one-hot encoded 'PAD' district features very well in the 'feature importance' assessment. The reality though supports the theory that the PAD location is a very important factor in economic feasibility of a crude oil refinery – an indication that the model definitely needs further tweaking.