

Understanding Maritime Commercial Data through Machine Learning

Jordan L. Taylor, 2018

jordantaylor.io

Summary

The following document describes the application of machine learning to maritime commercial data in the following context,

1. Simplified overview
2. Practical usage in a real-world application

Overview

The ability to apply machine learning to conventional data has become easier in the last decade due in part to the increased speed of central processing units and the ubiquity of off-the-shelf machine learning packages. Many of the algorithms used in these packages were developed in the mid to late 20th century¹ and only until recently have been able to express themselves fully as effective data management tools².

Maritime companies and third parties are now collecting data that can be applied to machine learning in a variety of different ways due in part to the proliferation of voyage management software. Using this data to provide substantive commercial outcomes is not only possible, but where before this type of development might have been cost prohibitive³ in the past is now relatively inexpensive. The process is even within reach of line managers⁴.

Efficiencies that will result from such applications will present a competitive advantage to companies that use these technologies⁵.

Predictive Analytics

¹ For instance support vector machines, common used today in Python's SciKit Learn, was developed by [Vladimir N. Vapnik](#) and [Alexey Ya. Chervonenkis](#) in 1963.

² Mostly due to accessibility of pre-packaged machine learning platforms to data scientists and generalists. A decade ago more labor intensive solutions were required for predictive analysis.

³ And possibly wasteful. Predictive outcomes could very well have no practical application, and this might only be understood after a great deal of cost has been exhausted.

⁴ Microsoft Azure Machine Learning (Azure ML)

⁵ Demonstrated by the proliferation of Amazon, which is fundamentally a data science company.

A simple example of how outcomes could be derived from predictive analytics could be demonstrated by the following question,

Q: How much demurrage⁶ exposure will a vessel have by calling Baytown Refinery?

In a conventional sense one could examine how much demurrage was recovered for vessels that called Baytown, and use the average to obtain an outcome,

M/T RIDGEBURY NICHOLAS A	USD 8.500
M/T AQUASUZARO	USD 9.500
M/T ATLANTIC SUMMER	USD 7.500

...

In this case the average exposure is USD 8.500. However the outcome could be considered scant given amount of factors that can manifest during a voyage. Furthermore there is not much information about what may have caused the exposure in the first place. So a follow-up question might be 'Can we apply the average to the next vessel we are fixing?'. The answer would be probably not.

With predictive analytics tools and a robust data set the texture of the outcomes changes significantly. We can start to use models and machine learning classifiers that might accurately predict demurrage exposure.

For example, given the following⁷ data on past voyages,

Vessel name	Built	Cargo M3	PDPR	Month	Exposure USD
M/T RIDGEBURY NICHOLAS A	2016	500000	40000	January	8.500
M/T AQUASUZARO	2015	300000	35000	June	9.500
M/T ATLANTIC SUMMER	2007	350000	25000	August	7.500

Let's see how an analyst might use this type of data effectively.

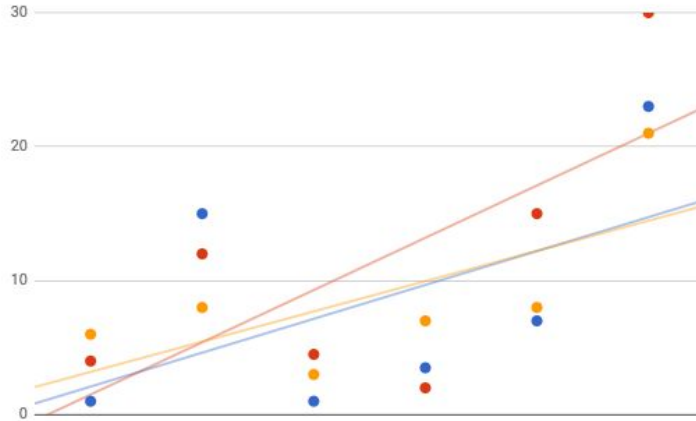
We can see right away that we may be able to draw dependencies around the time of year (i.e. fog season) and perhaps the year the vessel was built (i.e. vessel unable to meet charter party minimum at rail). This is an important observation, and one that may factor into exposure to certain charter party clauses.

⁶ A charge payable to the owner of a chartered ship in respect of failure to load or discharge the ship within the time agreed.

⁷ This data is not exhaustive and is for demonstration purposes only. More robust data points, with at least a few thousand datasets as a minimum would be required to derive actionable outcomes.

We will skip the data cleaning element of pre-processing, which is an element to obtaining outcomes.

For now it might be helpful to understand what the computer will see using the above data⁸,



Each color represents a vessel, and each point indicates a single recorded data point for that vessel. Although the numbers shown are dummy numbers and are not representative of the data⁹ - one can see that patterns emerge¹⁰.

Now we want to know what kind of demurrage a new vessel might be exposed to,

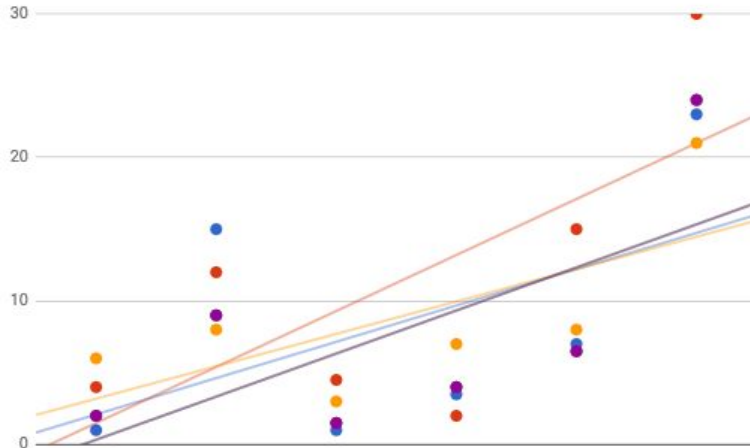
Vessel name	Built	Cargo GSV	PDPR	Month	Exposure USD
M/T PARAMOUNT HYDRA	2010	200000	30000	February	?

And now we have a new graph,

⁸ What we see here is a variation of k-means clustering. This is one method of many.

⁹ Categorical data, like the month, are converted to a binary. This is called one-hot encoding. Interestingly this is also where natural language processing becomes relevant, where categorical data can be presented as - for example - synonyms, thus perhaps improving modelling accuracy. Data that is of a integer or float type are compressed, generally.

¹⁰ Much like how a human might discriminate between a mouse, cat and dog. Since the human has seen all three in the past, an encounter with a mouse will match a pattern that the human has experienced prior.



Notice how the new vessel's magenta data points and trendline somewhat resembles two other vessels. Measuring the distance between each data point will provide a pathway to a prediction target.

Interestingly even if we are missing a prediction target, like quantity to load, a line will still be created and a target outcome produced¹¹,



This type of analysis would be helpful if the platform is missing data, or if we want to add additional data. For instance if the platform tries to guess the port costs of a new port that hasn't been called before.

Conventionally one would call an agent and get an estimated port cost. With a robust enough dataset, the platform could use as much criteria as is available (i.e. geography, location, pilot costs, &c.) to derive an outcome.

Furthermore, if the adjusted error was accurate enough, the system could also flag unusual costs associated with a port call and have at the ready evidence-based rationale for approving these costs.

¹¹ This is a powerful aspect of this type of data analysis.

In conclusion we can see that a trend line can be effective in predicting the demurrage amount, taking in as many factors as can be gathered. Furthermore as more data is applied to the model, the prediction targets will narrow. We also have a second example (port costs) that can be mined by these methods using different predictors. Both could eventually have direct commercial impact on chartering managers.

Training Versus Testing

The above example of predictive demurrage exposure might be interesting, but how could we prove that such an analysis is accurate? If the classifier says that vessel X will have USD 13.000 of demurrage exposure what evidence-based rationale has been used to make that decision?

This is a question that should be asked before setting out to build a machine learning platform¹². If there is not enough data the computer will not be able to guess accurately. Furthermore if the data is not relevant or exhibits leakage¹³ the model will be faulty and the predictive outcomes useless.

We also want the platform to gather new information and in a sense become smarter.

The first step will be to create a model¹⁴ with data that we will use as predictors. In the example above the predictors are the following,

1. Vessel name
2. Year built
3. Cargo quantity
4. Demurrage rate
5. Month vessel calls the port

The prediction target will be the following,

1. Exposure

We will also need data, and a lot of it. Let's say we have five years of data that includes the above predictors and how much demurrage was collected. There are 5,000 records for the five-year period.

Now we train the algorithm on 4,000 random records. This means we create a profile of each dataset like the scatter plots above.

¹² Note, however, that for an expert data scientist the cleaning and proliferation of outcomes should not be a labor intensive task. For a robust data set, swinging around predictors and adjusting quantiles, could represent as little as an hour of labor.

¹³ Data that can only be obtained if the predictive target is obtained

¹⁴ This is an involved process that is generally handled by specialists.

Then we run the remaining 1,000 predictors and allow the algorithm to attempt to predict outcomes based on the data. The system will know if the outcome is correct, as it will have a real-world outcome to compare itself to. A generally accepted metric is mean adjusted error (MAE), which is the error is a number that the prediction target is off by.

Once we run enough random tests we might get an MAE like the following,

USD 2.500

This will mean, on average, the predictions are off by plus or minus USD 2.500. Keep in mind there may have been outcomes that were substantially off, by say USD 10.000, but these outliers were offset by more accurate outcomes. However we can now be somewhat confident that our predicted demurrage of USD 13.000 is accurate + / - 2.500.

The MAE will change based on the following in relation to predictors,

1. Quantity
2. Quality
3. Diversity

For instance in our demurrage example above, the following will likely narrow the MAE,

1. Viscosity of cargo
2. Dataset corpus increased to 6,000
3. Weather delay data obtained from Houston Pilots
4. Currency adjustment

Furthermore if a company collects data¹⁵ in quantity and diversity - and continually applies it to the machine learning classifier - it will by its nature become more accurate.

Implementation

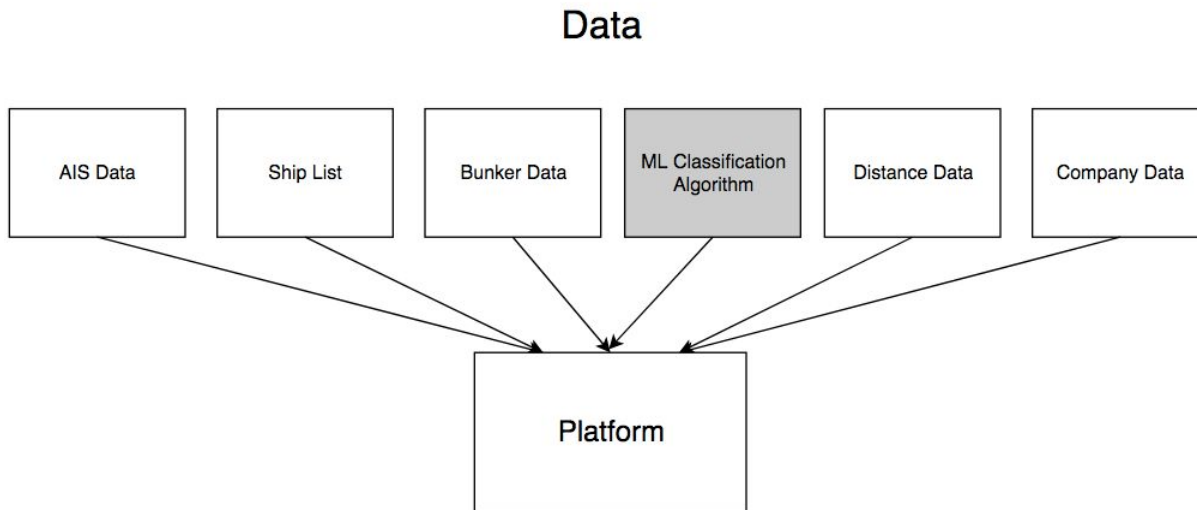
How could the prior example be implemented as a management system?

A data scientist that might want to prove that an outcome has commercial applicability might run the dataset through an integrated development environment and produce outcomes using a terminal. This is not much help to a line manager.

In order to use the information in the real-world a user interface will need to be developed, and some sort of useful dynamic wrapper created around whatever predictive feature that the company wants.

¹⁵ Or expands the number of predictors based on a proven model

Knowing how much demurrage exposure a vessel could experience might be helpful in estimating future voyage costs. Given the following data points¹⁶,



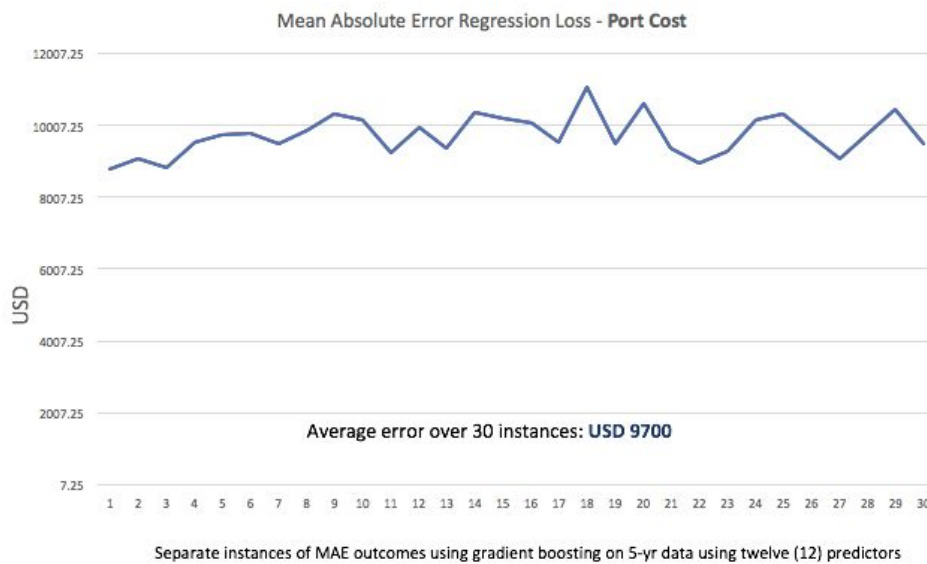
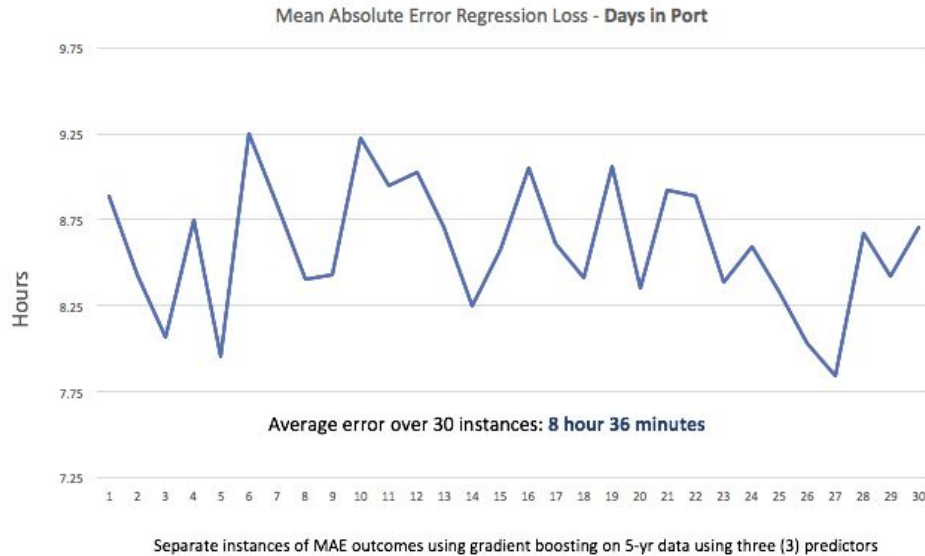
We might be able to estimate voyage costs for not only the company's vessels, but also vessels that may compete for the same cargo. Interestingly, given a large enough dataset, the demurrage predictor may not be the same for all vessels. For instance, if the management company and age of the vessel factors as dependency on the prediction target, the actual time charter equivalent will vary for a vessel competing for the same cargo.

The following are actual¹⁷ MAEs obtained for the following,

1. Estimated time in port
2. Estimated port cost

¹⁶ 'data point' is synonymous with API, or application program interfaces

¹⁷ Using gradient boosting



We can see an immediate application to the time charter equivalent,

$$\frac{\text{Net Freight (i. e. Freight lump sum - variable costs)}}{\text{number of days *}} = \text{TCE/day}$$

**from ship leaving port after discharge previous cargo until discharge current cargo*

The outcome, TCE/day, now contains evidence-based logic for the variables 'variable costs' and 'number of days', thereby narrowing the accuracy of TCE calculations and enabling a line manager to make more informed decisions about the outcome of a voyage.

Conclusion

The above examples provide a way forward for real-world application of machine learning to company proprietary data.

Many companies today are de-facto data management companies. The technologies discussed above are fundamental to the way these technology companies do business. Maritime companies will at some point need to make this transition to become competitive, and in fact many already have.

And even if efforts in this direction are for naught, company experience in implementing forward-thinking data management techniques are important to creating pathways forward for this type of technology.

[END]