

Vehicle Fuel Economy

CSCI 4146 Data Science Project

July, 2015

Derek Neil B00163969 dneil@cs.dal.ca

Table of Contents

Abstract.....	3
Background.....	3
Data Collection.....	3
Sources.....	3
Real World Data.....	4
Institutional Data	6
Complimentary Data.....	7
Analytical Questions	8
Questions.....	8
Additional Questions	8
Data Cleaning.....	9
Cleaning	9
Subsets.....	9
Methods.....	9
Manual Random Sampling.....	9
Visualizations	9
Predictive models	9
Results & Evaluation	9
Visualizations	9
Real World vs. Institutional Data	12
Prediction Models.....	16
Further Questions	17
Conclusion	18
Future Work.....	18
Appendices	19
Appendix 1 – Raw data, processed data, and source code.....	19

Abstract

Understanding and comparing fuel economy for a large number of vehicles can often be difficult due to the large number of factors that influence the performance. High level visualizations aren't common, even among enthusiast sites such as [fueley.com](#) which only visualize individual vehicle models.

In an attempt to create several different visualizations comparing different vehicle models, real world fuel economy data was collected from user based sites such as [fueley.com](#) and combined with other online sources ([Edmunds.com](#)) to fill in attributes such as weight and horsepower. Official fuel economy ratings, both the old 2-cycle, and the newer 5-cycle data sets, were obtained from the [open.canada.ca/data](#) website. Once concatenated and the institutional 5-cycle data was explored for composition, and vehicle configuration trends such as the adoption of 6 speed transmissions over the previously popular 5 speed transmissions.

The correlation between fuel economy and engine size, as well as vehicle weight were visualized along with other factors to create high level fuel economy distribution scatter plots.

Specific models that had a large number of real world fuel economy samples were then selected and filtered for outliers and incomplete features. Once cleaned, real world medians for vehicle model configurations were compared to the closest model available in the 2-cycle and 5-cycle institutional data sets. The real world data examined for two vehicle models more closely approximated the newer 5-cycle test that was recently introduced to better approximate real world driving conditions.

Finally prediction models were generated, and discussion of modest increases in 'fleet' fuel economy were discussed in relation to specific vehicle models. Varying configurations, even within the same manufacturer and model of car made it difficult to consistently observe whether fuel economy is increasing relative to previous generations.

Background

Fuel economy has been a social political issue for the last two decades, with rising fuel prices, resistance from manufacturers regarding mandates for raising fleet fuel economy, and changing government testing standards. Owners of vehicles have been able to post public their own fuel economy on purpose made websites providing a collection of data to compare against new and old government fuel economy testing standards.

Previously real world motorcycle fuel economy was examined for generating high level visualizations comparing makes and models. Focusing on four wheeled vehicles for this project provided a larger dataset for examination.

Data Collection

Sources

Data was collected from several sources. Real world data was scrapped from [fueley.com/cars](#), and government certified fuel economy test results were obtained from [open.canada.ca/data](#). Additional data from [edmunds.com](#) was used to compliment both the real world and institutional data sets where attributes such as weight, horsepower and towing capacity were noted.

Real World Data

Real world fuel economy data was collected from the [fuely.com](#). The ‘cars’ section of the site contains some 177,841 vehicles, spread across 5028 distinct model years. The distribution of owners per model is shown below in figure 1 with a long tail distribution.

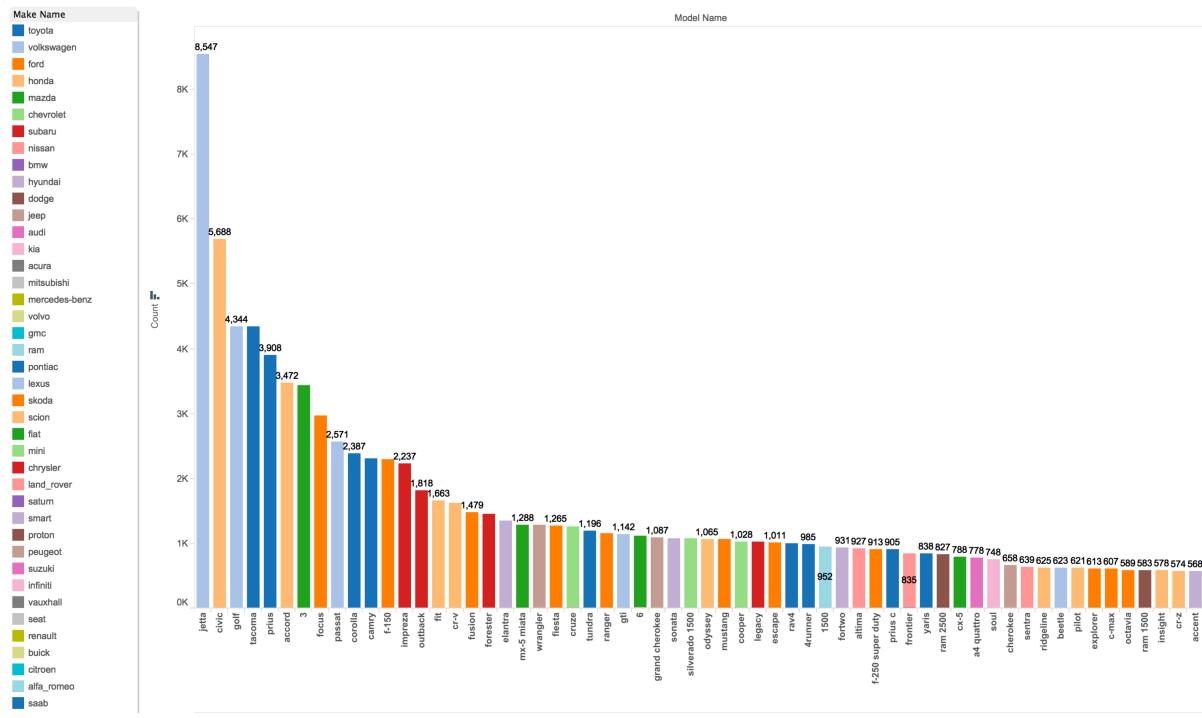


Figure 1 Long tail

distribution of number of owners per model on fuelly.com

Due to the nature of the fuelly.com site (designed to track fuel economy), this may not represent a random sample of the driving population, and could be skewed to drivers that are interested in driving as fuel efficiently as possible.

With this possible bias, we can still make use of the most popular models such as the Jetta, Civic, Golf and Prius which have a high number of users and compare those to the government standards for fuel economy.

Again taking into account the bias the site might be introducing in attracting users, we can also look at the distribution of vehicle manufacturers. Figure 2 shows a more condensed long tail with Toyota and Volkswagen, Ford and Honda being the most popular (in that order).

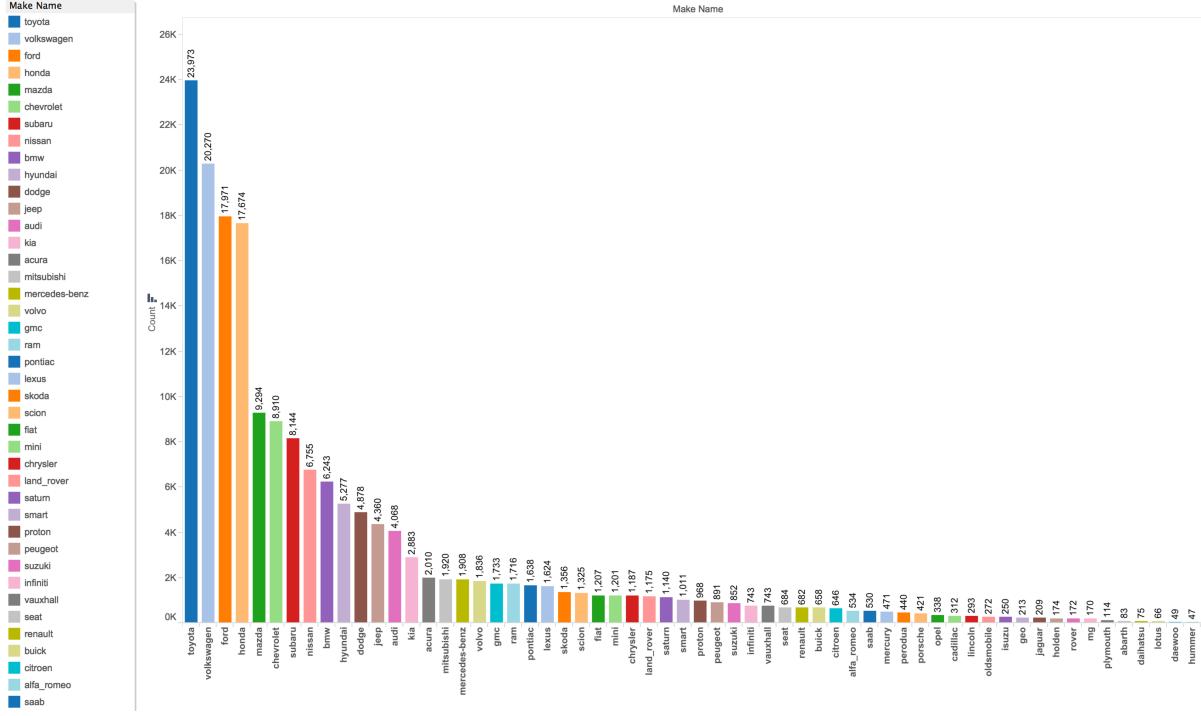


Figure 2 Long tail distribution of number of owners per make on fuely.com

Obtaining the data was done with a python script, one to scan make model years at a high level, taking the site averages, but lacking any fuel type, engine size, and location data; another slower script recording each users average fuel economy for their respective vehicle, allowing for fine grain details such as engine size, fuel type, and location (when provided).

Using PyQuery was very similar in terms of syntax to jQuery, and allowed the scripts to read a bit more clearly then other lower level URL parsing modules in Python. Figure 3 shows several examples of using PyQuery to select web elements by class, type, order and attribute, abstracting away the use of regular expressions for some of the more properly displayed content.

```

for link in pages:
    modelYearListing = pq(link)

    for ownerCar in modelYearListing('.browse-by-vehicle-display').items():
        ownerLink = ownerCar('.browse-by-vehicle-display').attr('data-clickable')
        country = ownerCar('.browse-details p:last a img').attr('src')
        country = countryRE.findall(country)[0]
        country = safeInt(country)
        country = nationalities.get(country, country)
    
```

Figure 3 Example usage of PyQuery to select elements on webpage by class, type, order and attribute

High level data contains the fuel economy average, make, model and year. Secondary attributes for further analysis and cleaning are the number of fuelups for the vehicle, the number of miles, number of owners, and link to the page the information was extracted from. Additional fields such as horsepower, type and weight were added from Edmunds.com as detailed in a following section.

avgMpg	fuelups	horsepower	make	miles	model	owners	type	weight	year	link
36.6	1343	115	volkswagen	584877	jetta	115	SEDAN	2859	2015	http://www.fuelly.com/car/volkswagen/jetta/2015
36.9	14383	115	volkswagen	6325111	jetta	532	SEDAN	2804	2014	http://www.fuelly.com/car/volkswagen/jetta/2014
37.5	36605	115	volkswagen	16367179	jetta	926	SEDAN	2842	2013	http://www.fuelly.com/car/volkswagen/jetta/2013
37.6	45853	115	volkswagen	20997540	jetta	1049	SEDAN	2012		http://www.fuelly.com/car/volkswagen/jetta/2012
37.2	39577	115	volkswagen	18278781	jetta	807	SEDAN	2804	2011	http://www.fuelly.com/car/volkswagen/jetta/2011
38.2	52799	170	volkswagen	24981883	jetta	953	SEDAN	3230	2010	http://www.fuelly.com/car/volkswagen/jetta/2010
37.2	45428	170	volkswagen	20593409	jetta	733	SEDAN	3230	2009	http://www.fuelly.com/car/volkswagen/jetta/2009
25.6	4025		volkswagen	1244890	jetta	95	SEDAN	3230	2008	http://www.fuelly.com/car/volkswagen/jetta/2008
25.4	2606	150	volkswagen	816911	jetta	94	SEDAN	3230	2007	http://www.fuelly.com/car/volkswagen/jetta/2007
40.3	29721	150	volkswagen	14747611	jetta	608	SEDAN	3230	2006	http://www.fuelly.com/car/volkswagen/jetta/2006
38.8	10248	115	volkswagen	5022986	jetta	221	SEDAN	2895	2005	http://www.fuelly.com/car/volkswagen/jetta/2005
38.5	17735	115	volkswagen	8469475	jetta	346	SEDAN	2892	2004	http://www.fuelly.com/car/volkswagen/jetta/2004
41	28420	115	volkswagen	14775171	jetta	591	SEDAN	2892	2003	http://www.fuelly.com/car/volkswagen/jetta/2003
38.1	16601	115	volkswagen	7630784	jetta	392	SEDAN	2892	2002	http://www.fuelly.com/car/volkswagen/jetta/2002

Figure 4 Sample output data from python site scrapper

Institutional Data

The government of Canada, through it's open data initiative, have made both the old and new datasets of standardized fuel economy publicly available directly from their [open.canada.ca/data](http://open.canada.ca/data/en/dataset/98f1a129-f628-4ce4-b24d-6f16bf24dd64) website. The older 2-cycle tests were often regarded by the public as unrealistically high figures not seen by most drivers in day to day driving. Addressing this concern, "The new test methods... integrate three additional test cycles that account for air conditioner use, cold temperature operation and driving at higher speeds with more rapid acceleration and braking. In most cases... the new test ratings are 10 to 20 percent higher than the old ratings because they ... better approximate everyday driving." (source: <http://open.canada.ca/data/en/dataset/98f1a129-f628-4ce4-b24d-6f16bf24dd64>, retrieved June 25th 2015). In this case, higher refers to the amount of fuel used for Liters / 100km testing, this would translate to a reducing using the Miles per Gallon measurement.

MODEL ^Y	MAKE	MODEL(# = high output engine)	VEHICLE	ENGINE	CYLINDERDEF	ION	TRANSMISS	FUEL TYP	FUEL CONSUMPTION	HWY (L/100 km)	COMB (L/100 km)	COMB (mpg)	CO2 EMISSIONS (g/km)
2014 ACURA	ILX	COMPACT	2	4 AS5	Z				8.6	5.6	7.2	39	166
2014 ACURA	ILX	COMPACT	2.4	4 M6	Z				9.8	6.5	8.3	34	191
2014 ACURA	ILX HYBRID	COMPACT	1.5	4 AV7	Z				5	4.8	4.9	58	113
2014 ACURA	MDX 4WD	SUV - SMALL	3.5	6 AS6	Z				11.2	7.7	9.6	29	221
2014 ACURA	RDX AWD	SUV - SMALL	3.5	6 AS6	Z				10.7	7.3	9.2	31	212
2014 ACURA	RLX	MID-SIZE	3.5	6 AS6	Z				10.5	6.4	8.6	33	198
2014 ACURA	TL	MID-SIZE	3.5	6 AS6	Z				10.4	6.8	8.8	32	202
2014 ACURA	TL AWD	MID-SIZE	3.7	6 AS6	Z				11.4	7.6	9.7	29	223
2014 ACURA	TL AWD	MID-SIZE	3.7	6 M6	Z				11.9	8	10.2	28	235
2014 ACURA	TSX	COMPACT	2.4	4 AS5	Z				9.3	6.2	7.9	36	182
2014 ACURA	TSX	COMPACT	2.4	4 M6	Z				9.9	6.8	8.5	33	196
2014 ACURA	TSX	COMPACT	3.5	6 AS5	Z				10.7	7	9	31	207

Figure 5 Sample from standardized fuel economy data

Examining the composition of the data set we can see a couple of interesting characteristics and trends come through already. Shown in figure 6, we can see the prevalence of vehicles with automatic four speed and manual five speed, and their recent decline starting in 2009 and 2006 respectively. Replacement systems such as Automatic (optionally with 'S' Select shift) and manual six speed transmissions show complimentary growth since the mid 2000's.

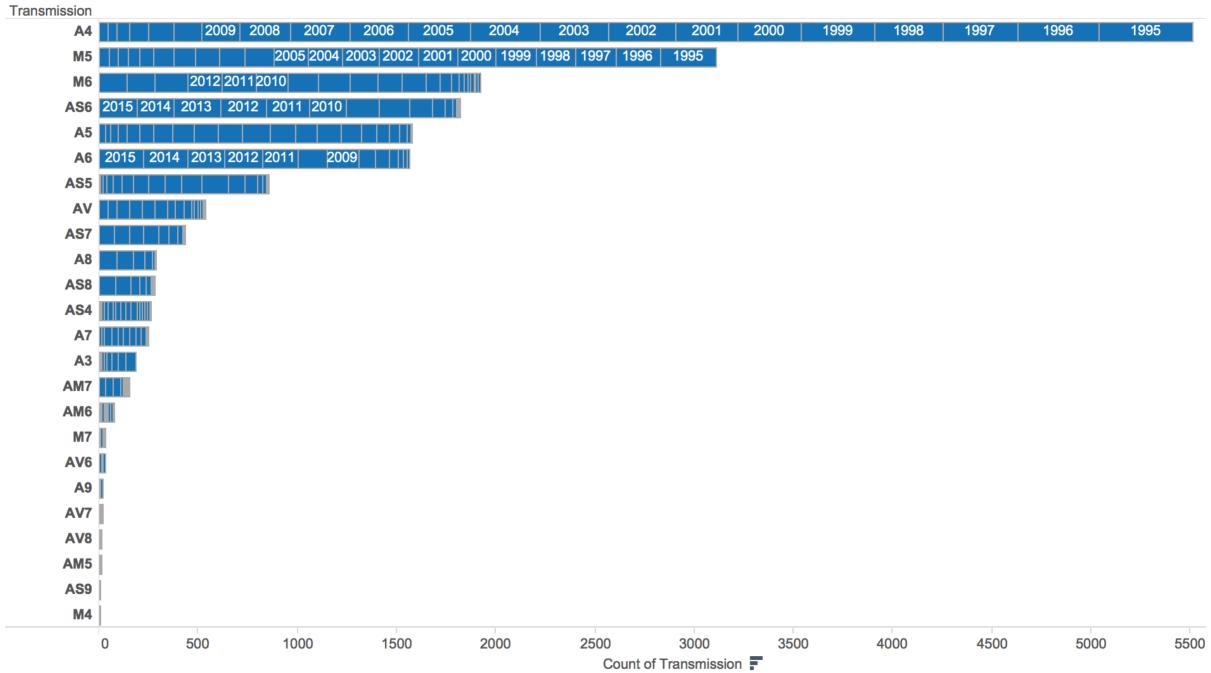


Figure 6 Distribution of vehicle transmission with yearly divisions

Turning our attention to figure 7 showing the distribution of different class vehicles for each year, we can also see an interesting increase in SUV's from 2002 to 2008, while most other models show more modest increases, or fairly standard yearly variation.

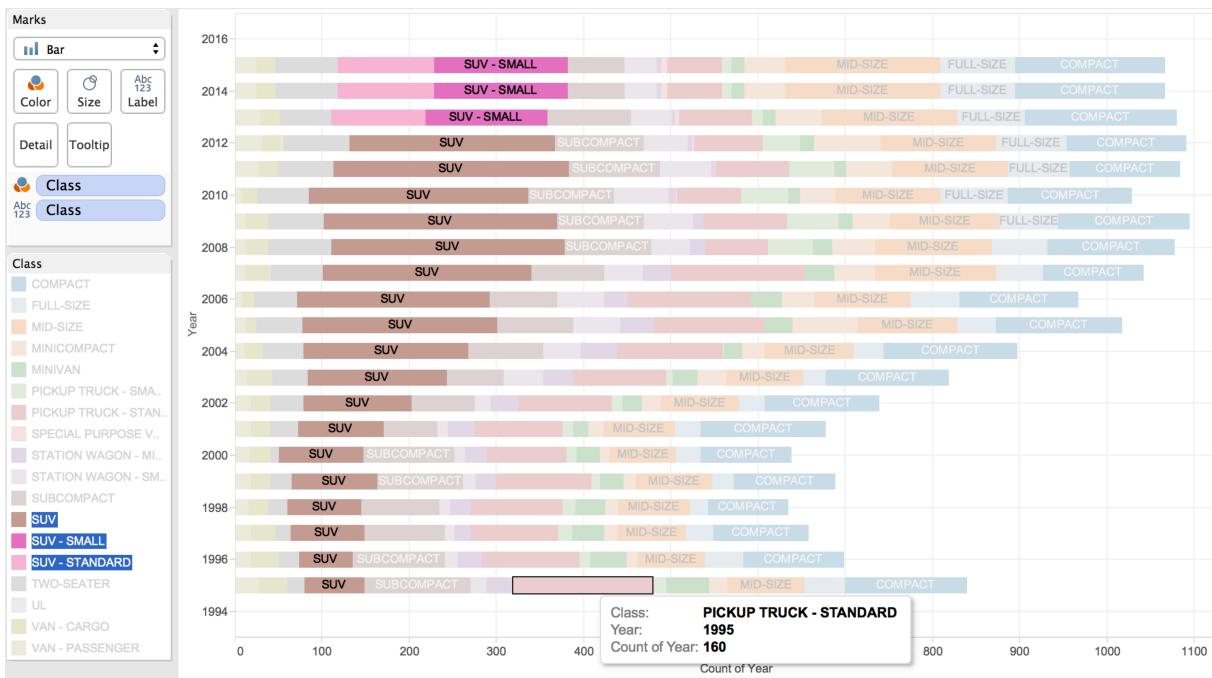


Figure 7 Distribution of vehicle type by year

Complimentary Data

Both the government, and real world data sets lacked important measures of vehicles such as weight, and horsepower. Edmunds.com provides an API for querying vehicle information, including details about each make, model, year and configuration. Unfortunately, due to apparent gaps in their vehicle weight information, this attribute was not offered in the API, instead it was optionally shown on

<http://www.edmunds.com/{MAKE}/{MODEL}/{YEAR}/features-specs/> which could then be scrapped using a method that was called once per make model year. This method provides base model information regarding weight and horsepower that are sufficient for our analysis.

```
curbWeightRE = re.compile('(\d+) lbs', flags=re.IGNORECASE)
horsepowerRE = re.compile('(\d+).*hp.*', flags=re.IGNORECASE)
def api(make, model, year):
    model = model.strip().replace(" ", "-")

    type = ''
    curbWeight = ''
    horsepower = ''
    url="http://www.edmunds.com/"+make+"/"+model+"/"+year+"/features-specs/"

    try:
        if DEBUG: print '\t\tedmunds api call to',url
        modelFeatures = pq(url=url)

        for highlight in modelFeatures("#highlights-pod .data-table li").items():
            span = highlight("span").text()

            if "CAR TYPE" in span:
                type = highlight("em").text().strip().upper()

        for feature in modelFeatures(".feature-spec .items td").items():
            label = feature("label").text()

            if "CURB WEIGHT" in label:
                curbWeight = feature("span").text()
                curbWeight = curbWeightRE.findall(curbWeight)[0]
            elif "HORSEPOWER" in label:
                horsepower = feature("span").text()
                horsepower = horsepowerRE.findall(horsepower)[0]

        if DEBUG: print '\t\ttype:',type, 'curbWeight:',curbWeight, 'horsepower',horsepower
    except:
        if DEBUG: '\t\tedmunds api error'

    return type, curbWeight, horsepower
```

Figure 8 Pseudo api call to edmunds.com to obtain vehicle type, curb weight, and horsepower for a make, model and year.

In retrospective, the vehicle type could have been looked up from the government data, and unknown vehicles discarded. Future work might add towing capacity as it can be a buying decision, and attribute which might affect the design, and thus fuel economy of the vehicle.

Analytical Questions

After initial sampling of real world and institutional data, several analytical questions were formed to guide the development of the project and explore some of the data science methods.

Questions

1. Does fuel economy continue to increase year over year for all makes and models?
2. Does fuel economy increased until some common year and show signs of a plateau since?
3. Does fuel economy differ significantly by location?
4. Is real world combined fuel economy better represented by 2-cycle test or 5-cycle test institutional data?

Additional Questions

Based on feedback from an initial presentation, vehicle weight (gross weight) was added as an additional data attribute on the dataset. In addition to engine size, transmission type, and fuel type, the vehicle weight provides an additional continuous attribute that influences the vehicle fuel economy that was missing from both the real world and institutional data sets. Our analysis will also consider the correlation of weight to fuel economy.

Data Cleaning

Cleaning

The real world data set required filtering outliers and entries with missing values. Filters were primarily done by specifying minimum and maximum ranges for data, as well as using medians for any small clusters of data to mitigate the influence of single low or high values in an otherwise normally distributed cluster of data.

Subsets

Focusing on subsets of the data (specific vehicle models) was used to compare the two cycle, five cycle, and real world data. Comparing identical make, model, year, transmission, and fuel types allowed for easy direct comparison of fuel economies with only engine size and ratio of city to highway driving being uncontrolled variables.

Methods

At each involved step, rapid, and high level methods or tools were used to identify, summarize, visualize, validate, and vary the significant attributes that influence vehicle fuel economy.

Manual Random Sampling

During the development of web scrapers for the real world data, and scripts to concatenate the open data sources, sampling was performed on the range of values using the Dataframe.describe() function. The statistics provided by the describe function included average, median, mode, variance, min, and max which allowed for identifying skewed datasets, and spotting extreme outliers.

Handling parser errors in developing the web scrapper also allowed for eventual online handling of entries with null values in a certain attributes, or identifying fields that would require filtering in visualization and modeling software.

Visualizations

Exploring the correlations between various attributes and fuel economy was assisted by the use of Tableau (<http://www.tableau.com/>), allowing for rapid generation of visualizations in various configurations of real world, and institutional data to identify similarities, differentiated instances, and possible trends.

Predictive models

Finally, Wizard (<http://www.wizardmac.com/>) statistical software was used to quickly investigate significant covariates for the fuel economy outcome, and generate a predictive model where the effects of year, weight, horsepower, engine size, transmission, fuel type and cylinders could be observed.

Results & Evaluation

Starting from an overly aggregated view of the data and moving down to higher dimensional representations we can start to investigate trends in the data, comparing datasets, and finally arriving at predictive models using the significant attributes.

Visualizations

Grouping all vehicles regardless of their class and configuration is obviously not a statistically significant measure of the change in fuel economy but rather looks at the change in 'fleet' fuel economy of all the vehicles offered from year to year that are captured by the institutional data set. Recall that the institutional data set also contains a significant

growth of SUV's during the mid 2000's which, coincidentally, aligns with the decade long plateau in 'fleet' fuel economy.

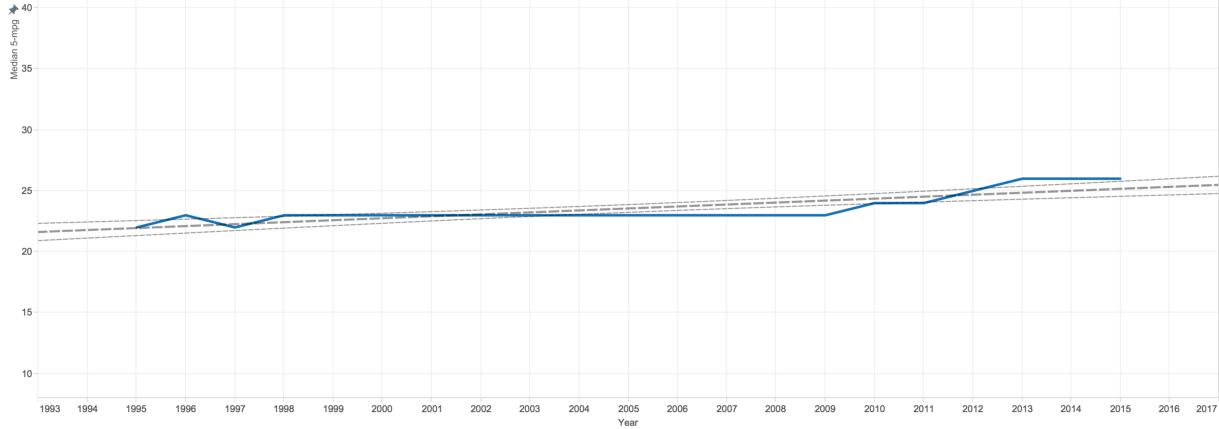


Figure 9 Median fuel 5-cycle economy over time for all vehicles captured in the institutional data set

Moving to a higher dimensional visualization shown in figure 10, with individual points showing unique make-model-year-transmission-fueltype-economy entries, we can start to see some general trends such as higher fuel economy for some diesel (blue) vehicles, and lower fuel economy for ethanol (green) vehicles. Using the tooltips for various outlier points we can also see that hybrid vehicles such as the Toyota Prius are able to achieve higher than average fuel economy.

In terms of fuel type, there is a clear set of vehicles that were tested with both regular (X – brown) and premium (Z – yellow), and we can see in figure 10 that regular fuel achieves a consistently higher average fuel economy across these identical subset of vehicles.

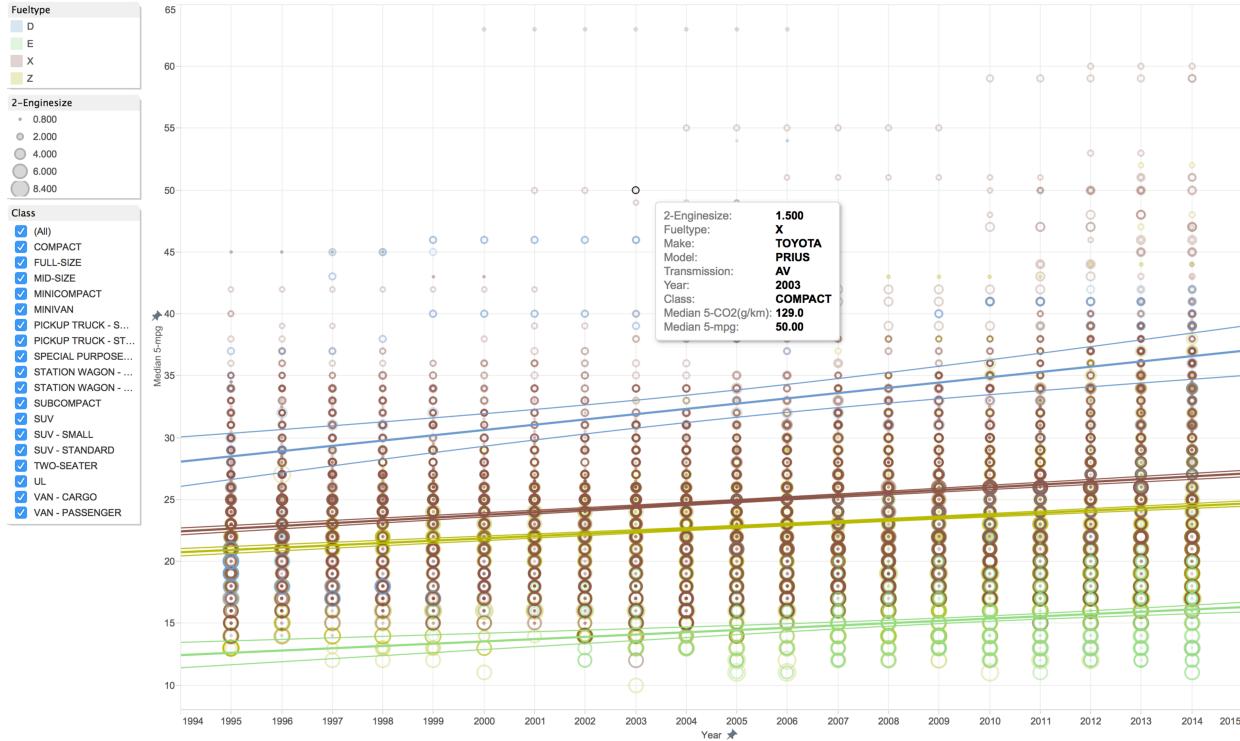


Figure 10 Median 5-cycle fuel economy over time for each individual model configuration

Since the institutional dataset contains a large portion of SUV's we can use that class of vehicle to thin out our previous visualization. Figure 11 shows the trends for SUV's grouped by fuel types. Again, we can see the outliers

such as the Ford Escape Hybrid deviating from general trends. We can also see a general correlation between engine size (circle size) and fuel economy, with larger engines having lower fuel economy than vehicles with smaller engine sizes.

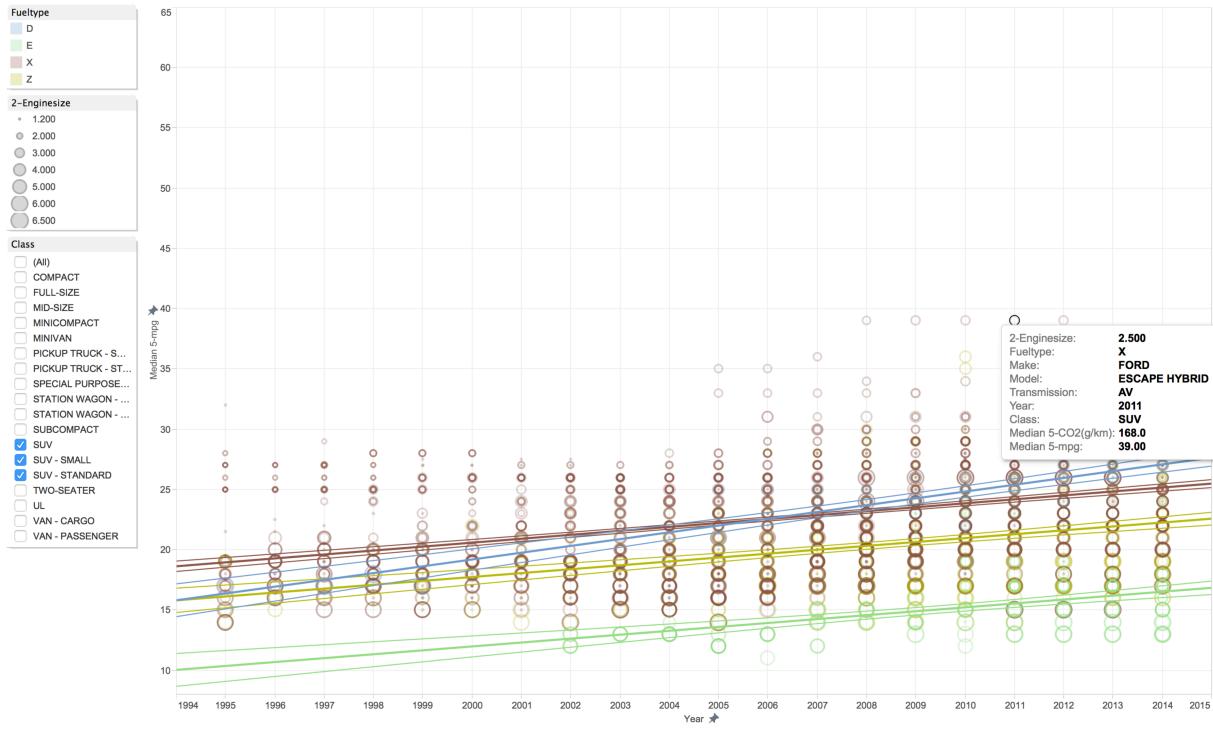


Figure 11 Median 5-cycle fuel economy over time for each individual model configuration

Shifting the visualization from fuel economy over years, to fuel economy vs engine size, we can see in figure 12 there is a moderate correlation of decreasing fuel economy as engine size increases. We also see that the number of cylinders (size of points) has somewhat of a correlation to fuel economy. As with previous figures, we can see small diesel engines and hybrids achieving high fuel economy above the general logarithmic trend lines.

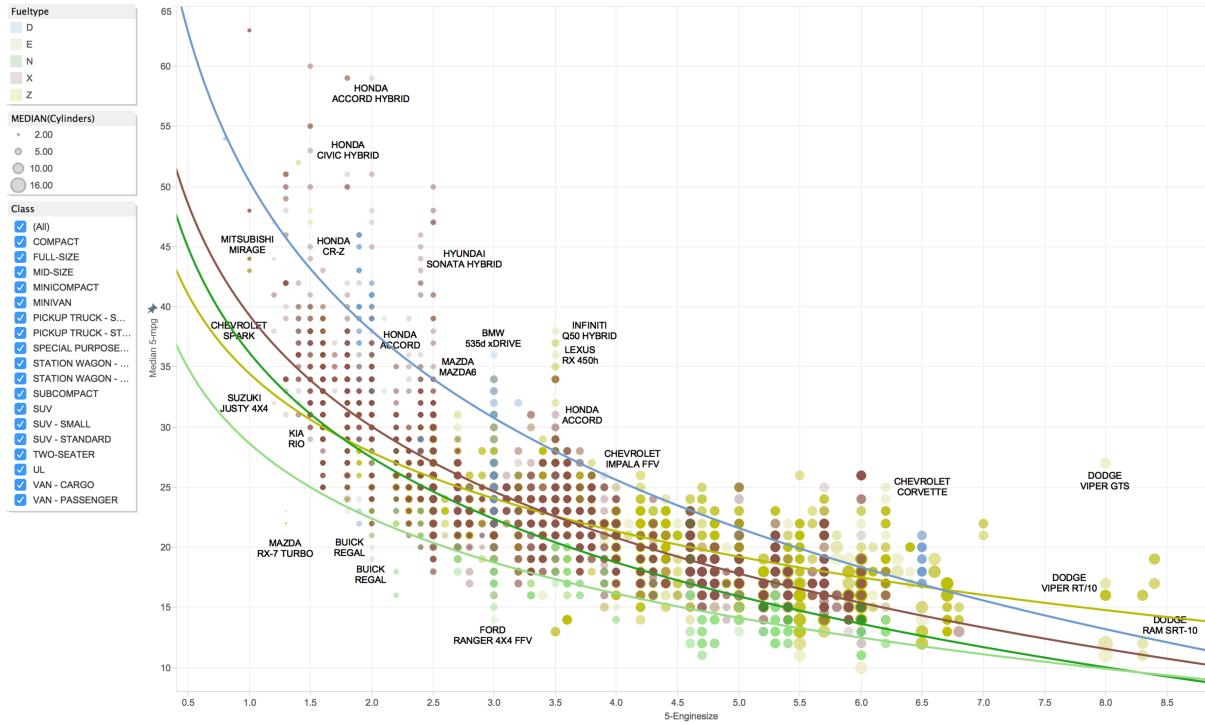


Figure 12 Median 5-cycle fuel economy vs engine size

The real world dataset combining high level data from fuelly and Edmunds contains a wider range of vehicle models, along with their weight (colour) as shown in figure 13. From the coloration we can see a moderate correlation between vehicle weight, and fuel economy. Examining the two ends of the spectrum for example, we can see no ~7000lb vehicles with a fuel economy above 25 mpg, and all but one or two outliers of the lightest ~1700lb vehicles are all above 25mpg. The size of the points represents the ownership on fuelly.com, which shows that users of the site might tend to buy more fuel efficient vehicles, and they may not represent a random sample of all drivers.

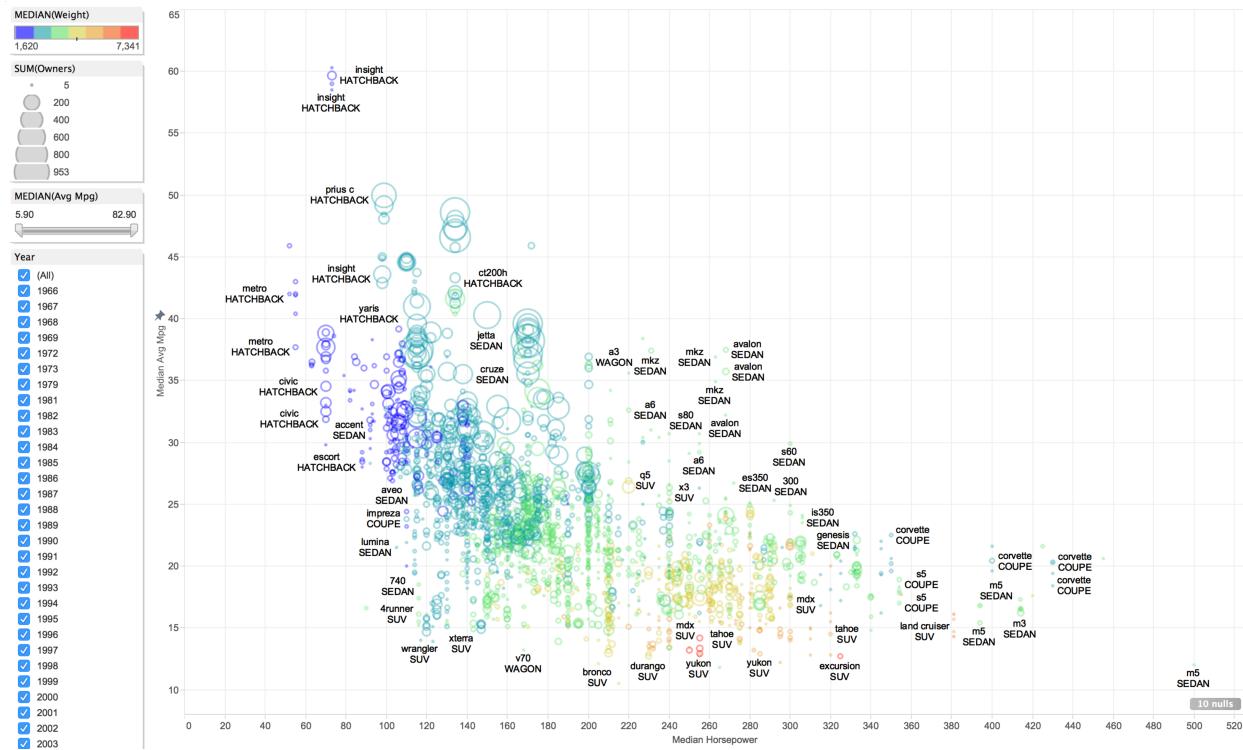


Figure 13 Median real world fuel economy vs horsepower

Real World vs. Institutional Data

After experimenting with visualizations that encompassed the entire data set, focus then shifted to individual vehicle models. Referring back to our distribution of real world models in figure 1, we see that the Volkswagen Jetta is the most popular model on fuelly.com followed closely by several other models that also happen to have been high on our fuel economy visualizations such as the Toyota Prius. Having a large number of data points for a single model carries with it the likelihood that there are good representations for different fuel types, engine sizes, and geographic locations for our analysis.

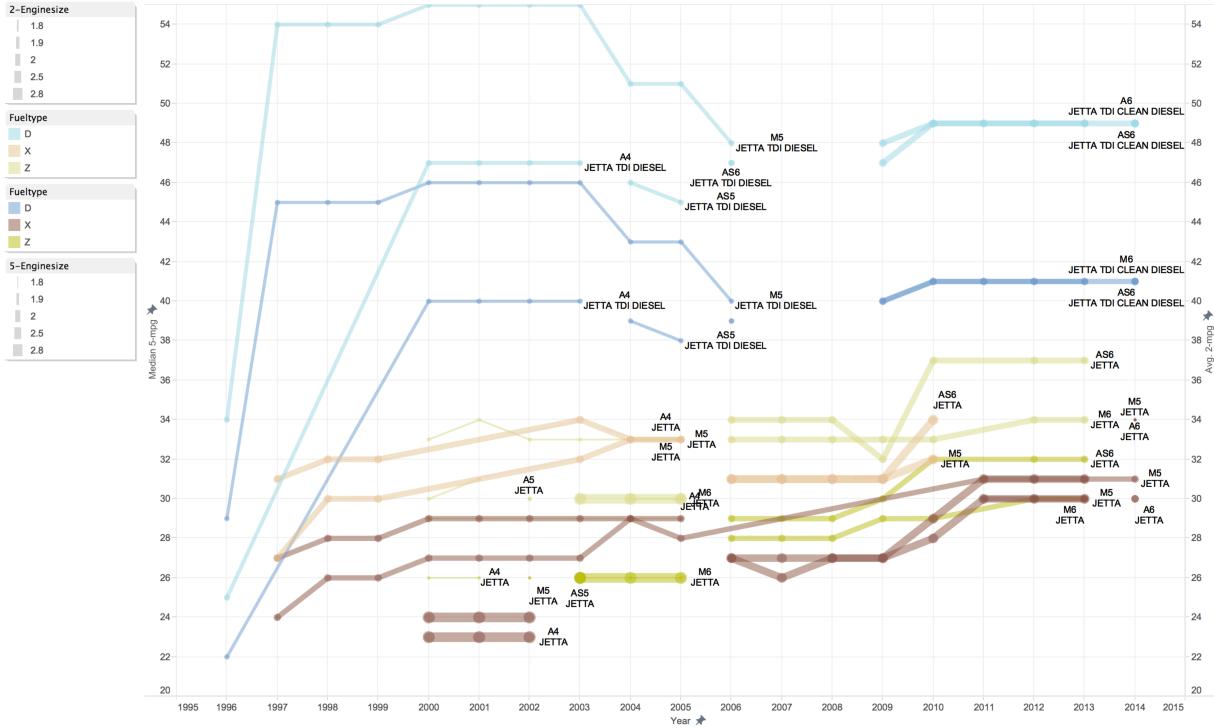


Figure 14 Median fuel economy (2-cycle light shade, 5-cycle dark shade) over time per engine size and transmission

Starting with the Jetta, shown in figure 14, we can see the old 2 cycle test fuel economy (light shade), and the newer 5-cycle test fuel economy (dark shade). Some data points appear to be simple offsets from the old test to the new one, while some configurations show slightly different characteristics such as the AS6 model with automatic select shift. In figure 15 to add the real world data in bright colours and find that it falls closer to, or oscillates around, the 5-cycle test fuel economy data. Although in the case of the M5 diesel (blue & green) configuration, we see the real world (green) data is still slightly below the 5- cycle test, but obviously closer than it would have been to the much higher 2-cycle test fuel economy shown in light blue in figure 14. Other information we can gain from figure 15 are the a mix of configurations that show steady improvement over time (M5 gas in brown), but others shown signs of plateaus from one year to the next in (A4, M5 diesel, M6 diesel). In the case of the manual 6 speed diesel, we can see this plateau even roughly represented in the real world data (green).

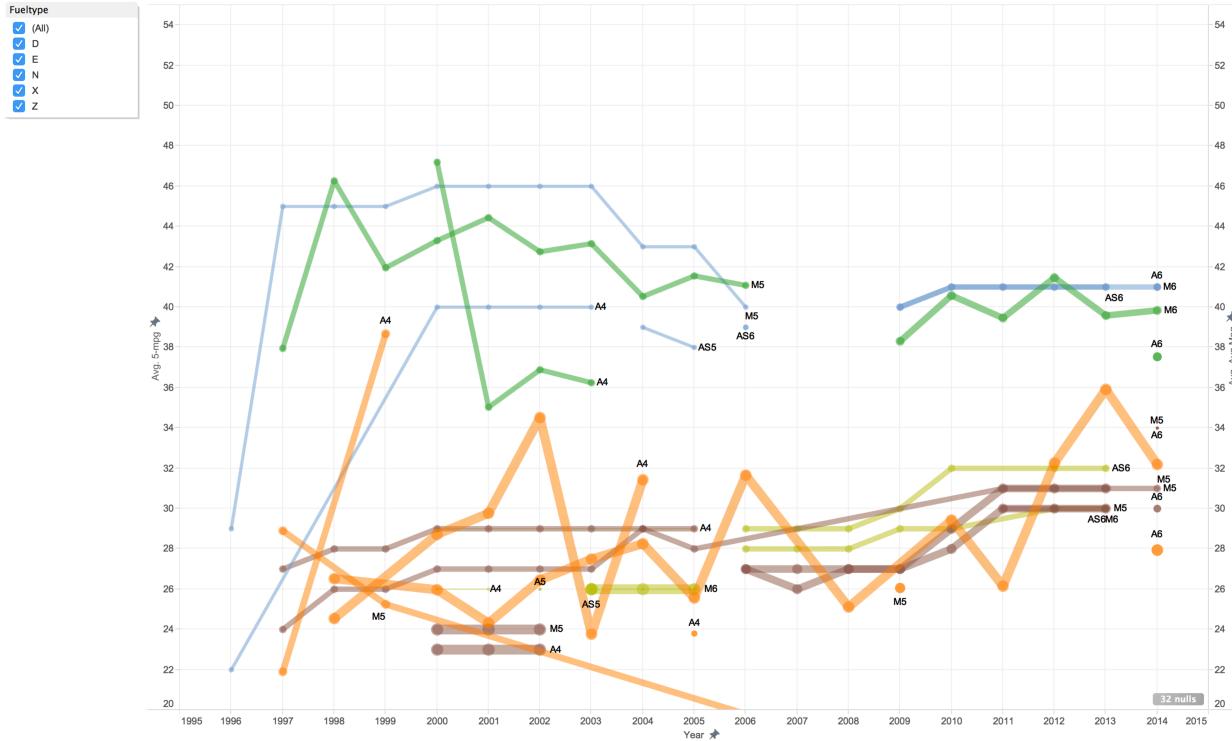


Figure 15 Volkswagen Jetta median fuel economy (5-cycle dark shade, real world bright shade) over time per engine size and transmission

Remaining with the Volkswagen Jetta, we can see in figure 16 there's also some disparity between fuel type usage around the world. Countries such as Brazil, and Costa Rica only have gasoline engines in our real world data set. Conversely, European countries such as Spain, Poland, and Germany only have diesel engines in our real world data set. This has an impact on fuel economy per country for this particular model due to the inherent difference between diesel and gas engine operation, while comparing similar gas types tends to be very similar in terms of fuel economy with the exception of Brazil which is known to use near equal amounts of ethanol (thus lower fuel economy) and gasoline which aren't captured well in our real world dataset.

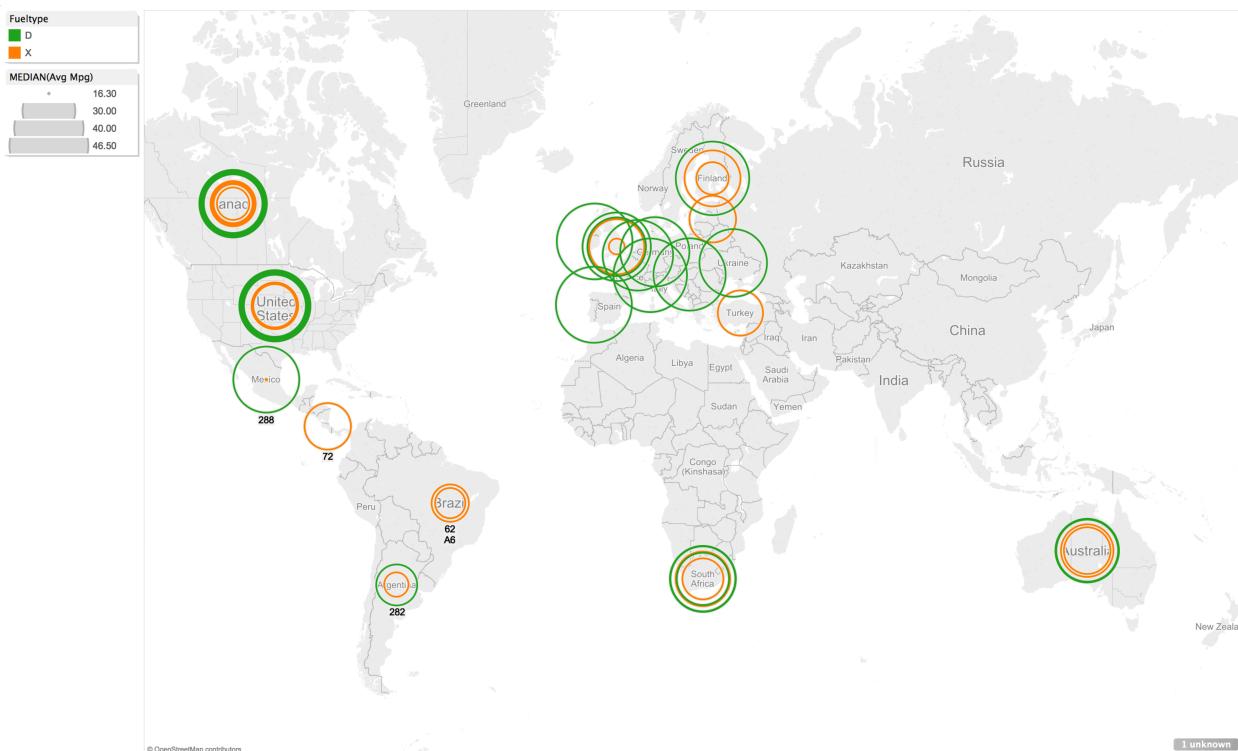


Figure 16 Volkswagen Jetta median fuel economy (real world) per fuel type and country

The Toyota Prius was another exception vehicle in terms of fuel economy, however instead of using diesel fuel to achieve its higher than expected fuel economy, the gas engine is assisted by an electric power unit. With only one fuel type, and two engine sizes represented in our real world data set, the visualization is much clearer with all three (2-cycle light brown, 5-cycle dark brown, and real world orange) fuel economy data sources shown on one graph.

Interestingly, the hybrid Prius was one of the few vehicles to show a higher 5-cycle test score compared to the old 2-cycle test data. The real world results end up somewhere in the middle, suggesting that perhaps additional testing cycles might be required for hybrid vehicles to produce better approximations of real world fuel economy.

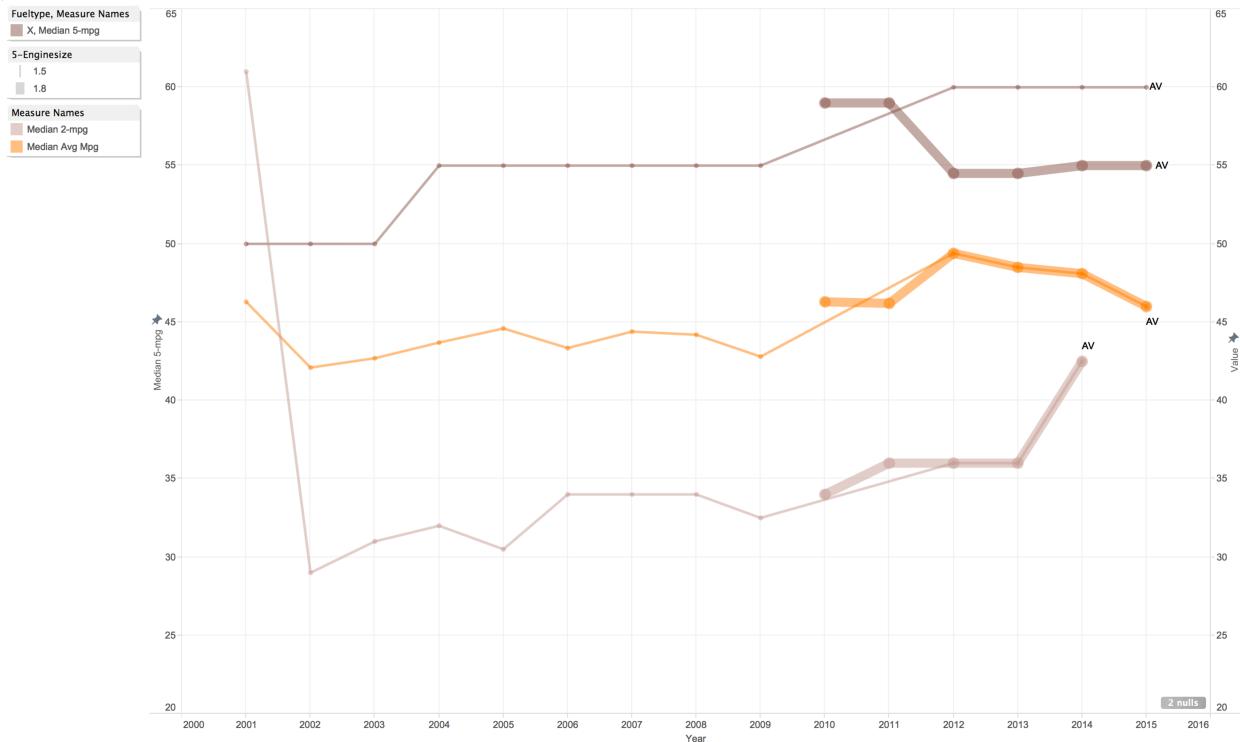


Figure 17 Toyota Prius median fuel economy (2-cycle light shade, real world bright orange, 5-cycle dark shade) over time



Figure 18 Toyota Prius owners by region, fuel type, and transmission

With a single fuel type, and nearly identical fuel economy from one country to the next, the Toyota Prius world map instead shows large contrast in US owners (that happen to be fuelly.com members) relative to everywhere else in the world. Although more price sensitive countries such as UK are second, the user distribution may be skewed to the US since fuelly.com is itself based in North America.

Prediction Models

Using the high level real world data, 5-cycle institutional fuel economy, and a secondary version of the 5-cycle institutional fuel economy with the weight added for available models, we are able to generate three predictive models.

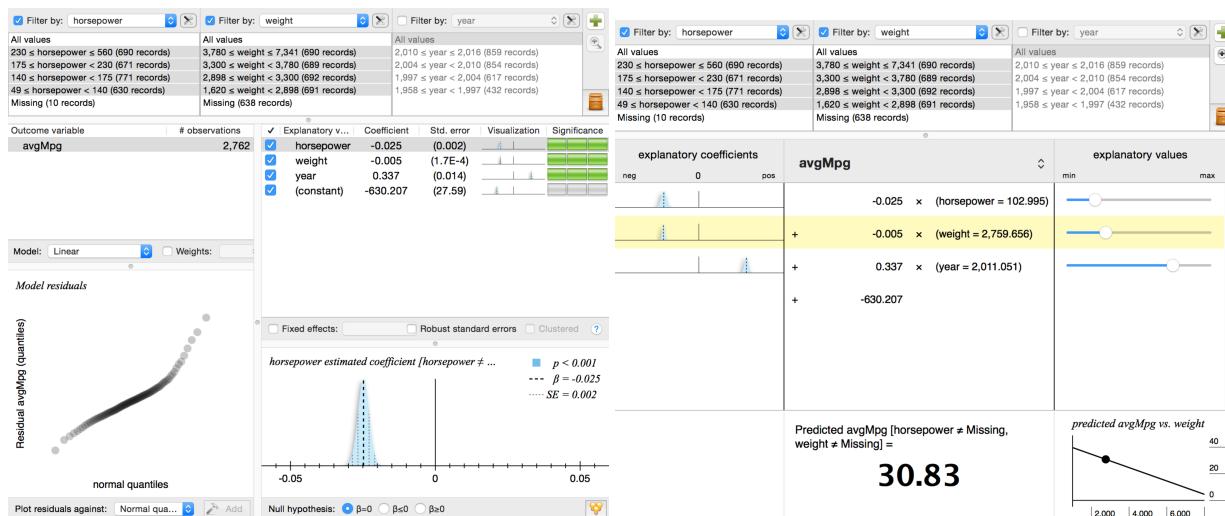


Figure 19 Real world fuel economy filtered prediction model, QQ plot, and coefficients

Although the real world data has some deviation in the QQ plot we can determine coefficients for our model. Notice that although the coefficient for weight is small, the attribute spans a large range of values and does indeed have a significant impact on our prediction model as shown by the fuel economy preview graph in the lower right hand corner.

Moving to 5-cycle institutional data (without weight) we repeat the process, and observe a better QQ plot for our model since the institutional data has more reliable engine size, fuel type, cylinder, and transmission type for each vehicle.

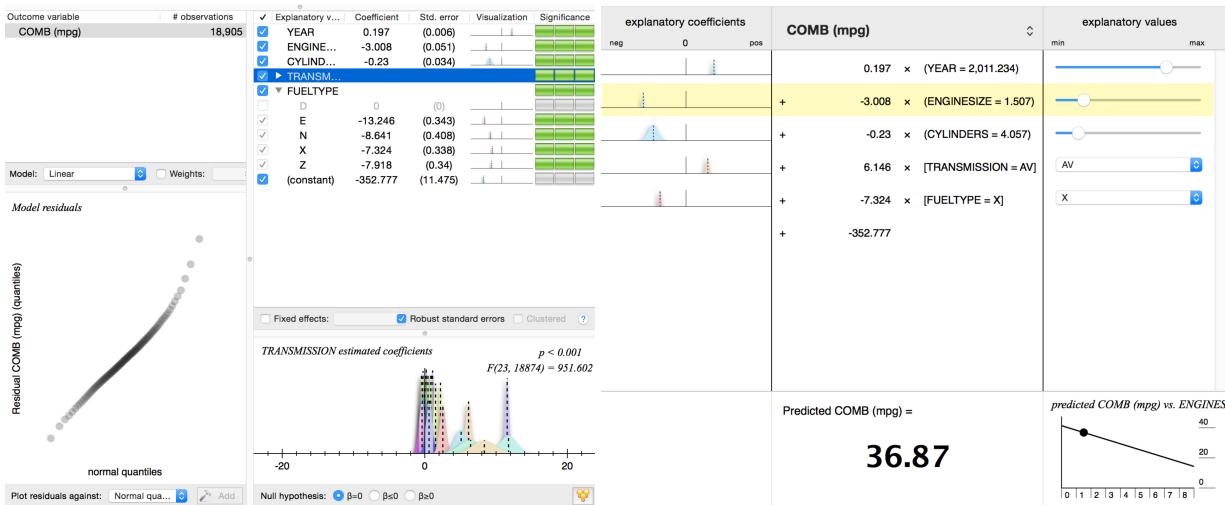


Figure 20 5-cycle institutional fuel economy prediction model, QQ plot, and coefficients

Finally generating a model for the 5-cycle institutional data with weight (where available) provided the most accurate model (near linear QQ plot). Notice that Transmission was kept as a categorical value, each with its own coefficient. We could have optionally separated the number of gears, and type of transmission, but this may have introduced too much interaction in our model, and resulted in combinations that don't exist in the real world such as AV3 which would be automatically variable 3 speed transmission, which would not make sense.

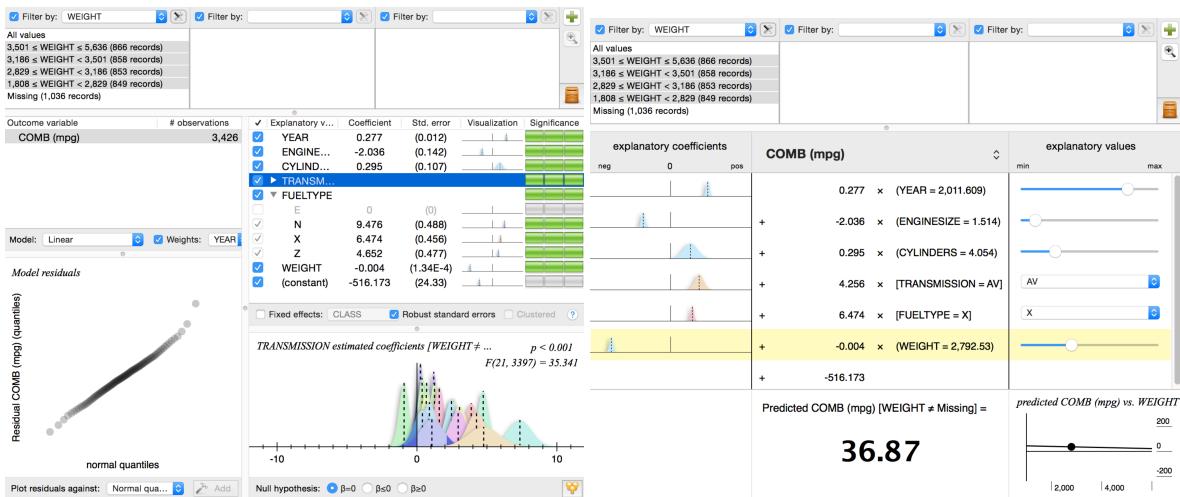


Figure 21 5-cycle institutional fuel economy with vehicle weight added and filtered, prediction model, QQ plot, and coefficients

Further Questions

Having looked closely at the Jetta and Prius models, both of which increase fuel economy substantially compared to our global trend lines, questions such as "why aren't all vehicles hybrids?" and "why aren't there more diesel

hybrids?". After some brief searching, both turned out to be cost related, although as fuel costs rise in various parts of the world, automakers have been showcasing lightweight diesel hybrid concept vehicles that might catch on in the near future.

Conclusion

Each of our analytical questions can now be re-visited based on our results and evaluations.

1. Does fuel economy continue to increase year over year for all makes and models?

For the most part Yes. Some models showed decreases, but where often discontinued, and new models added that were sometimes more fuel efficient. It is a difficult question to answer outside of global 'fleet' averages that are influenced by their model composition from year to year.

2. Does fuel economy increased until some common year and show signs of a plateau since.

For some models, yes. The Volkswagen Jetta showed signs of plateauing for certain models that were identified by their transmission and fuel type combination. The Toyota Prius, another fuel efficient vehicle also had periods of nearly unchanged fuel economy.

3. Does fuel economy differ significantly by location?

According to the sample of data from fuelly.com, no (not for same make and model), but fuel type does. Countries such as Brazil showed almost 100% gasoline usage instead of a mix of gas and diesel like other countries, this had the effect of lower fuel economy due to the type of fuel, but ultimately countries that had the same selection in fuel were generally equivalent in terms of fuel economy from the detailed data examined in this report.

4. Is real world combined fuel economy better represented by 2-cycle test or 5-cycle test institutional data?

Institutional data from the new 5-cycle test better approximates combined "real world" usage than data from the old 2-cycle test. Real world data for some vehicle configurations such as the Jetta M6 were very close to the 5-cycle test fuel economy (further statistical analysis should be conducted for more objective determination).

Future Work

Both real world and institutional data sets had important attributes such as vehicle weight, horsepower and towing capacity. The former two were sourced for this project while the latter wasn't considered until after the data had been collected. Adding towing capacity could increase the utility of the analysis, allowing for an additional filter, and analysis of vehicles that have at least some given amount of towing capacity.

While investigating the individual real world fuel economy entries, some vehicles showed evidence of seasonal fluctuations in fuel economy that could be related to local weather conditions of the vehicle owner. Figure 22 shows approximately a 8 mpg difference between August and December month fuel economy averages. Increasing the depth of the web scrapper to collect individual fuel ups from each user would allow for more detailed analysis, including the exact date of the fuel up. While limited location information was found in the real world data, there is enough geographical distinction at the country level to observe some correlation (for example Canada vs. Brazil).

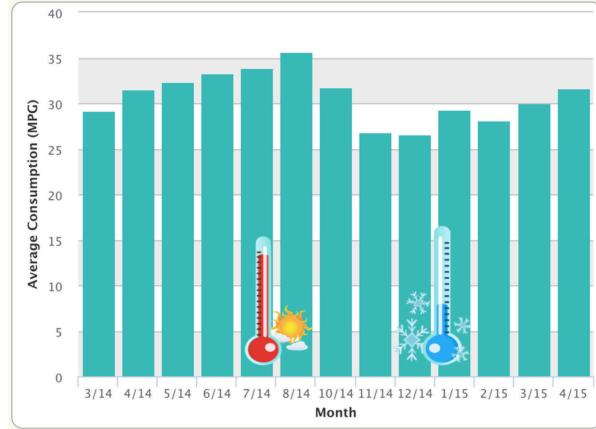


Figure 22 Seasonal variance of monthly average fuel economy for a Canadian owned Toyota Corolla

Appendices

Appendix 1 – Raw data, processed data, and source code

See https://bitbucket.org/derek_neil/fueleconomy for a full listing of source code, raw and processed data