# "Big Data" Analytics in Operations Research

Cara Cuiule and Grady Noll, PRICE Systems L.L.C.

## Abstract

As the world becomes more connected and data-driven, accumulating masses of data is becoming a more common corporate strategy. This paper discusses the current PRICE® Systems data collection methodology within the company that utilizes web scraping/crawling. The team will present both lessons learned and potential shortcomings in the current method of approach, along with plans for future endeavors.

## Introduction

Since new estimates are informed by the qualities of older projects, data is critical for the job of a cost estimator. Data that is plentiful and in an ideal format for analysis is hard to come by, especially because aerospace and defense industry information is often classified.

The internet provides a potential source of data due to its large amounts of accessible information. Online data can be collected manually, but this can quickly turn into a labor-intensive task that is prone to human error. As a result, automated data collection can be incorporated through web crawling and scraping, processes which collect data from internet webpages using special software [1]. Both were utilized by the cost research team at PRICE to collect data from several different websites.

The goal was to build repeatable, automated processes that would streamline the collection of cost-related data. After the data was validated it could then be studied to establish Cost Estimating Relationships (CERs) to inform future estimates. This is beneficial because of the variable nature of the online commercial market. The procedure could be rerun periodically to update the resulting database.

This paper is not intended as an expert's analysis but instead as a review of the team's projects, challenges, and lessons learned. It intends to inform a general procedure for these projects. In addition, it will expand upon the presentation given at ICEAA 2017 called "Automated Data Collection Using Open Source Web Crawling Technology" [2].

## More About Web Scraping/Crawling

Web crawlers and scrapers can be software programs, or scripts written in a programming language. The team used commercial software due to a lack of script writing expertise. Web crawlers and scrapers are very similar, but are not the same thing. While both take the data from a webpage's Hypertext Markup Language (HTML) code, web crawling is focused on

determining the webpage's structure and other web page links through keywords and URLs. Web scraping takes specified pieces of data from web pages. It can be done for a variety of purposes [3]. The team used web crawlers to grab product page URLs and then scrapers to grab data. Both may be done using the same program.

## Methodology

When scraping a website, the team generally followed these steps:

1. Find websites with relevant data
2. Determine if crawling/scraping from the website is legally permissible
3. Choose the appropriate tool
4. Obtain and normalize the data
   a. Obtain a list of product webpages
   b. Extract the data from product webpages
5. Application
6. Data validation

These steps are elaborated below. Throughout the explanation of this procedure, we will reference the creation of a handgun database.

### 1. Find websites with relevant data

Finding cost related data is obviously critical for establishing CERs. This can be problematic when online commercial stores only have prices. Prices have additional fees, special taxes, and profit margins in addition to the cost of manufacturing.

To mitigate the effects of these factors on cost, the team used publicly available financial data from the manufacturers to estimate profit margins on prices. Due to the obscure nature of pricing in commercial markets, this tactic is not flawless. For example, there may be some products that are sold below cost to attract customers, but there is no way to know which particular items this would apply to [4]. In addition, Manufacturer Suggested Retail Price (MSRP) or market pricing also accounts for the wholesaler and retailer profit [5]. Wholesaler and retailer profits are difficult to measure due to commercial markets fluctuating to supply and demand. However, data would still be acceptable if the prices on the vendor website in consideration were similar to the online list prices directly from a publicly traded manufacturing company.

Other critical data gathered included the quantity of the product being sold (if relevant) and unit weight. The team also gathered other specifications. This included characteristics that were identified as potential drivers of cost (e.g., material composition), or helpful for organizing the data (e.g., brand).

While there are many websites that have this type of data for handguns, other factors forced the team to limit which websites we could choose from. These included the ease of scraping and the legality of collecting information from the particular website.

**2. Determine if crawling/scraping from the website is legally permissible**

While it is possible to scrape massive amounts of publicly available data, that does not mean it is guaranteed to be permitted by law. One of the most critical parts of this process was that the team had to determine if the website gave permission to applications for crawling and scraping. Krotov and Silva [6] point out that so far there has been no legislation completely for or against this method of data collection, or even a set of rules that software must abide by. They propose a list that expands upon the legal and ethical concerns, most of which are applicable to those looking to gather cost-related data from commercial websites. For example, not only is using or accessing data a potential legal issue, but users must also be careful that data isn't copyrighted. In addition, users must be careful not to overburden the servers with requests or gather data that can damage the website's owner through the revealing of trade secrets.

Based on these suggestions, the team considered each of these concerns before starting the data collection phase of the project. The volume of scraping being done wasn't overly ambitious. The intention of the data collection was purely academic in nature, and not intended to be harmful. However, team members still needed to ensure that webpages were being accessed legally.

The team consulted each website's version of "The Robots Exclusion Protocol" to check that the data was contained in parts of the website that could be scraped by a robot. This is also called the robots.txt file, which is the part of the site dedicated to noting which directories of the server should not be accessed by a robot. Despite these warnings, some robots and programs may be built to ignore these directions [7]. This is a dubious practice that the team avoided.

For instance, the text below means that no robots are allowed to access the "/welcome.html" directory on this website:

```
User-agent: *
Disallow: /welcome.html
```

*Figure 1: Robots.txt Snapshot*

It is possible that the Robots Exclusion Protocol may change over time to include or exclude certain addresses. When these processes are repeated to update the datasets, the user must recheck these webpages to ensure they are still adhering to these policies.

For creating the handgun database, the only website we found that met all the conditions in the first two steps was Hyatt Guns, a brick-and-mortar gun shop with an extensive online catalog [8].

### 3. Choose the appropriate tool

Once a website was chosen, the team had to select a tool for web crawling/scraping. Two applications have been used so far for the project. The first was RapidMiner, which was introduced in the 2017 ICEAA presentation [9]. RapidMiner breaks its work down into operators that are easily viewed, customized, and moved in its Graphical User Interface (GUI). It has many capabilities such as data transformation, validation, and, modeling [2]. RapidMiner is not primarily a web scraping or crawling tool, but it has a downloadable extension called "Web Mining" that can be used for the task. Unfortunately, it has a steep learning curve. In fact, the team spent a few hours a day learning the tool for several weeks before knowledge of the application was sufficient to extract data from websites like Hyatt Guns.

The second tool used was Octoparse [10]. This tool uses machine learning algorithms to see patterns in the HTML code. All the user needs to do is "point-and-click" on the desired webpages and/or data. Then, the program creates a Microsoft Excel or a .csv file. Due to the ease of use, it's possible to learn how to use this program within a couple of hours. However, because of the limitations of this software, the spreadsheet must be normalized with another program.

These tools are only an example of web scrapers that are available online. They were chosen because of previous use within the company and because there are free versions of the software. For example, RapidMiner and Octoparse will let any file in their respective free versions have up to 10,000 data points. A larger (but still not exhaustive) list of scrapers can be found here [11].

### 4. Obtain and normalize the data

The process described below assumes the use of a more involved program such as RapidMiner versus a simpler program such as Octoparse. This step is typically the most time consuming, so it is described in two sections: grabbing product webpages and then taking the desired data.

### a. Obtain a list of product webpages

This part of the procedure was related to web crawling, not scraping. On a store's webpage the links that contain product prices, names, and specifications are needed for any analysis to be possible. For example, Figure 2 shows a small subset of the list of products you would see on the Hyatt Gun Store [8]:
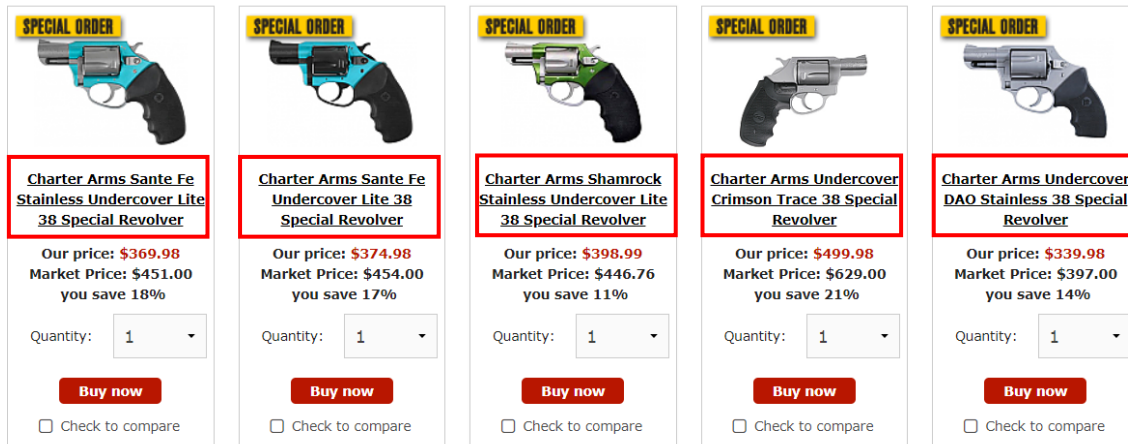
*Figure 2: Hyatt Guns Products*

Outlined in red is the section a user would click on to get to the product page. It is also the part of the HTML code where the URL lies. The ICEAA presentation from 2017 described multiple web crawling methodologies to obtain a list of URLs that lead to each individual product page [2]. After this is completed, the data is parsed from these webpages.

### b. Extract the data from product webpages

For data extraction, only a basic understanding of HTML is necessary. The technique the team took advantage of patterns identified in the code. Once this was set up properly, the program can do the rest.

Web scrapers are able to take HTML code from each page in the list and convert it into a text file which is used to parse the data. Then, the team primarily retrieved information using "String-matching". This method obtains all the text between two different strings as specified by the user. Before using, it is important to check multiple webpages to ensure that the structure of the code is consistent among the products, and that the information desired has unique strings surrounding it. This is because string-matching only grabs the information from the first instance of the exact string indicated in the text.

Figures 3 and 4 demonstrate screenshots from a revolver being sold on the Hyatt Guns website [12]:

## Charter Arms Undercover Crimson Trace 38 Special Revolver

Charter made its name with the classic .38 Undercover. Today, Charter's family of .38 Specials has grown to meet the tastes and demands of a variety of shooters.

At 16 oz., this five-shot .38 Special revolver is compact and lightweight. Its 2" barrel and superior safety features makes it ideal for concealed carry situations.

It's a perfect compromise between size, weight and stopping power!

Qty 1

**5-7 DAYS**

*what's this?*

✔ In Stock ❓

~~$629.00~~ *$499.98*

*Figure 4: Revolver Description*

| Brand | 🟤 |
|---|---|
| SKU | 73824 |
| Quantity in stock | 1 item(s) available |
| Shipping: | All Firearms MUST ship to an FFL Dealer or can be picked up in our Retail Store |
| Condition | NEW |
| Caliber: | .38 Special |
| Capacity: | 5 |
| Barrel Length: | 2 inches |
| Action: | Double Action / Single Action |
| Grips: | Lasergrips |
| Weight (unloaded): | 16 oz |

*Figure 3: Revolver Technical Specifications*

The HTML code in a browser can be located by right-clicking on a web page and selecting "View Page Source". From there, the section of the code where the tables are can be isolated and selected. Note that in the code pictured, there is a regular price and a sale price. To keep consistency for the validation phase, the market price of $629.00 was extracted, not the discounted price offered by the website. Therefore, the string-matching method must be directed to obtain the text between the boxed text below, *class="market-price"><span class="currency">* and *</span>*:

```
class="market-price"><span
class="currency">$629.00</span></span><span
itemprop="offers" itemscope itemtype="http://schema.org/Offer"><meta
itemprop="priceCurrency" content="$" />
<span
class="product-price-value" itemprop="price" content="499.98"><span
class="currency">$<span
id="product_price">499.98</span></span></span><link
```

*Figure 5: HTML String Matching*

This method was repeated to grab the desired parts of the data on every product page such as brand, SKU, condition, caliber, etc. There are situations where string-matching may not work, such as when HTML code surrounding data is non-unique or it varies slightly throughout the webpages. Methods such as X-Path or Regular Expressions are harder to learn than string-matching but may be successful in these special cases [13] [14].

Once this phase of the project is complete, there still may be extra spaces, unwanted characters, or differing metrics in the dataset. These should be cleaned up and items missing critical data should be removed. Data can be easily normalized by using a program that creates an automated process for doing so (such as RapidMiner or a VBA script in Excel). In the case of small datasets that are mostly complete, normalizing by hand may be reasonable.

Regardless of methodology, establishing consistency in data is critical for determining CERs. At this point the dataset is not yet complete, but may look something like Figure 6:

| Name | Price | Item SKU | Brand | Caliber | Capacity | Condition |
|---|---|---|---|---|---|---|
| Charter Arms Undercover Crimson Trace 38 Special Revolver | 629 | 73824 | Charter Arms | .38 Special | 5 | NEW |
| Charter Arms Undercover DAO 38 Special Revolver | 383 | 13811 | Charter Arms | .38 Special | 5 | NEW |
| Charter Arms Undercover DAO Stainless 38 Special Revolver | 397 | 73811 | Charter Arms | .38 Special | 5 | NEW |
| Charter Arms Undercover Green and Black 38 Special Revolver | 419.99 | 23820 | Charter Arms | .38 Special | 5 | NEW |
| Charter Arms Undercover Grey Stainless 38 Special Revolver | 424 | 43820 | Charter Arms | .38 Special | 5 | NEW |
| Charter Arms Undercover Lite 38 Special Revolver | 448 | 53820 | Charter Arms | .38 Special | 5 | NEW |
| Charter Arms Undercover Lite Black and Polished Silver 38 Special Revolver | 423.99 | 53871 | Charter Arms | .38 Special | 5 | NEW |
| Charter Arms Undercover Lite Bronze and Black 38 Special Revolver | 449.99 | 53883 | Charter Arms | .38 Special | 5 | NEW |
| Charter Arms Undercover Lite Red and Black 38 Special Revolver | 451 | 53824 | Charter Arms | .38 Special | 5 | NEW |

*Figure 6: Revolver Dataset*

## 5. Application

As stated earlier, cost is needed to establish CERs, but information on consumer websites typically yields prices. However, cost can be approximated when special taxes and fees such as overhead and profit are considered.

It is assumed that any taxes would be the final addition to price, so it should be the first to be deducted. Firearms have a special excise tax that adds 10-11% based on the category of firearm [15]. Therefore, the team created a new column of data where the price of the gun before the excise tax was calculated.

The next step was determining the percentage of profit and General & Administrative (G&A) fees that are attached to the product's market price. The team reverse engineered these cost factors from publicly available financial records published by publicly traded companies in the relevant industry. Unfortunately, there was only public financial data for one firearm manufacturer: Sturm, Ruger & Company, Inc. Over the course of three years, G&A expenses averaged about 20% of expenditures and total profit averaged about 12% [16]. These values are similar to internal guidance for general government contracts. PRICE's® TruePlanning® software calculates cost automatically based on the indicated G&A and Fee/Profit, values but this can also be done by hand or with other programs.

Now that the team had criteria for approximating the cost of each item, the data was applied to the Hardware cost model in TruePlanning. Weight was the only quality immediately applicable to the model. Other aspects like manufacturing process, percent of new structure, and production numbers were once again interpolated from online sources or subject matter expert (SME) knowledge. These values were industry averages that were used consistently across every item. The only exception to this were the additional production numbers, which were averages based on production surveys that varied among firearm type [17].
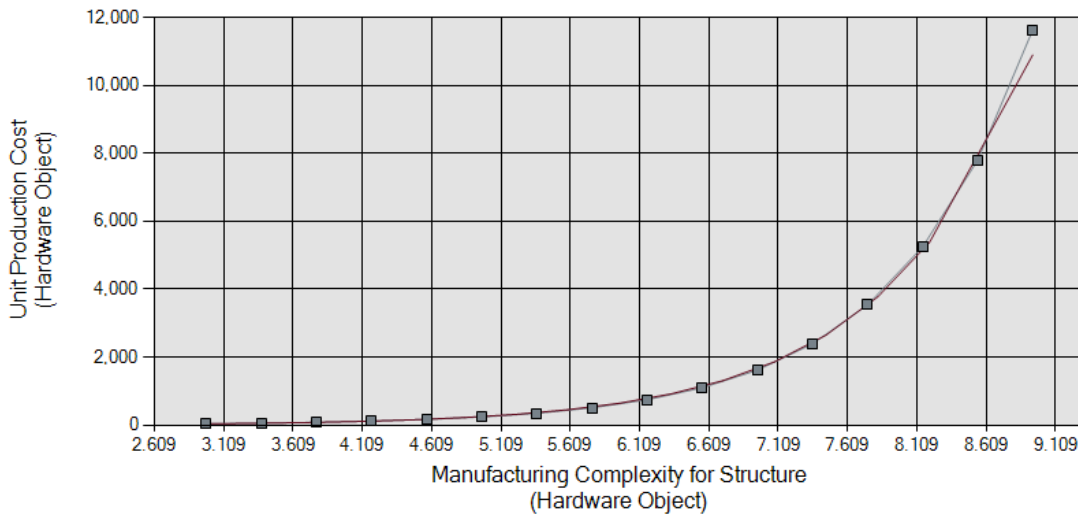
*Figure 7:  Unit Production Cost vs MCPLXS*

After that, each item was calibrated to a Manufacturing Complexity for Structure based on the interpolated cost. Manufacturing Complexity for Structure is a major driver of the Hardware model that serves as a modified measure of cost per weight unit. For the two different types of handguns, pistols and revolvers, the histograms revealed that the distributions of the calibration results approached a normal distribution:
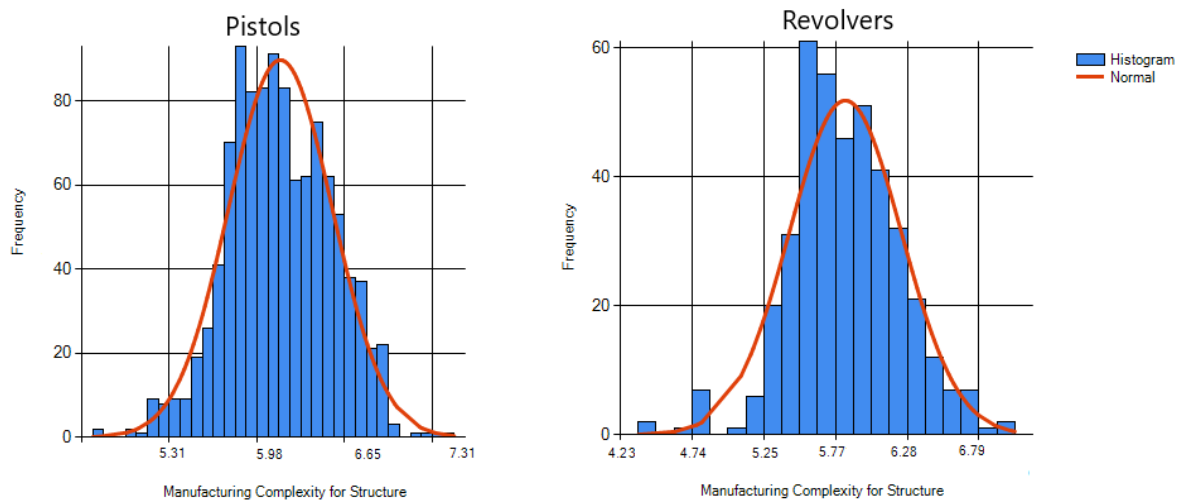


*Figure 8: MCPLXS Histograms vs Normal Distributions*

When the dataset was completed, it was also applied to the TruePlanning framework. Data was sorted by brand. Each brand was represented in the Product Breakdown Structure (PBS) as a single Hardware object with average weights and Manufacturing Complexity for Structure values. Figure 9 shows the PBS of what was now called the component database. On the right is the input sheet that is representative for all of the Smith & Wesson revolvers:
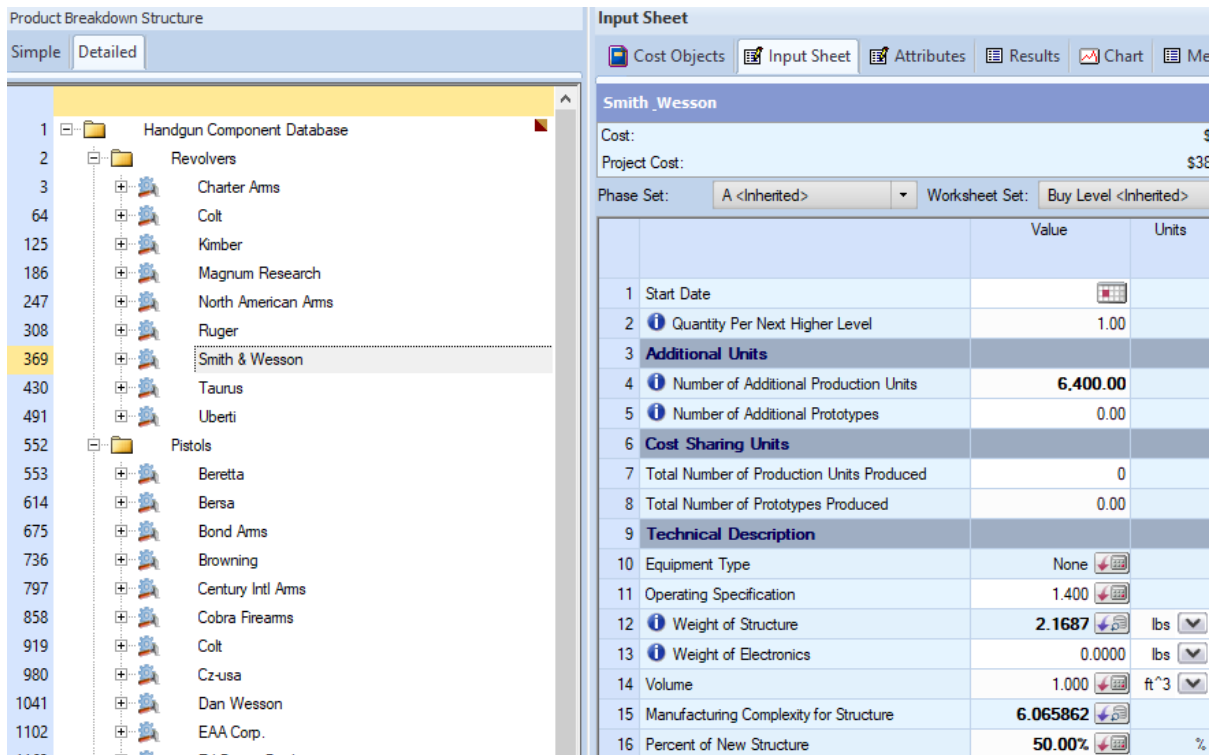
*Figure 9: TruePlanning® Component Database*

## 6. Data validation

Since the project focused on obtaining prices from the internet and interpolating cost from that data, validation is important to ensure that these estimations were reasonable. Finding relevant cost data of other public online sources to cross-check the database was the key to validation. The team called these estimates test cases.

A set of test cases are valid if the average percent error was within one standard deviation of a normal distribution (+/- 33%). The first type of test cases in this particular database involved the firearms competing for the XM17 modular handgun system competition, which was a contract listed at $5.4 million for about 7,100 pistols [18]. Using this information, the team found that the exact firearms in contention for the contract were listed in the database, allowing the team to estimate costs for that contract based on each specific gun's weight and manufacturing complexity. The estimate closest to the contract award amount, the Sig Sauer P320, is the gun that won the contract. The seven estimates were, on average, 10% below actuals:

| Brand | Firearm | TP Estimate | %Error |
|---|---|---|---|
| Beretta | M9A3 | $7,471,667 | 38.4% |
| Beretta | APX | $3,879,026 | -28.2% |
| CZ-USA | P-09 | $4,066,033 | -24.7% |
| FN USA | 509 | $4,438,781 | -17.8% |

| Glock | 19 | $4,675,500 | -13.4% |
|---|---|---|---|
| Sig Sauer | P320 | $4,914,740 | -9.0% |
| Smith & Wesson | M&P 9 | $4,564,371 | -15.5% |

*Figure 10: XM17 Contract Results*

Another test case was created to demonstrate capability of performing early program estimates. Contracts for 65,000 were awarded in 2005 to both Sig Sauer and Heckler and Koch. They ranged from $23.7-$26.2m. The calibers of the guns for both orders were 9mm, 40 S&W, and 0.357 Sig [19]. When the data is filtered to these specific calibers and given the production number that matches the contracts, the estimate was $25.6m, which fell within the range of expected cost.

Both types of test cases demonstrate the ability to use the databases for cross-checking the following: competitor data, early program estimates, or price-to-win scenarios. These databases have already been used within PRICE as either a cross check or as an initial exploration of data that is not easily accessible elsewhere.

**Another Type of Database**

The team also completed a long gun component database that was very similar to the one created for handguns. Both were an example of a PBS of products organized by one or more qualities (object type and brand). Another type of database was a PBS meant to model a system whose components are taken from the data collected. An example of this was the Humvee component database, which is a PBS representative of the military surface vehicle. The information for each cost object was gathered from Kascar, an authorized parts and service provider for the Humvee [20]. Once again, weights and costs for each datapoint were used to calibrate each item to its respective cost, which was then calibrated to a Manufacturing Complexity for Structure. Averages for datapoints of the appropriate product type were taken and applied to the PBS:
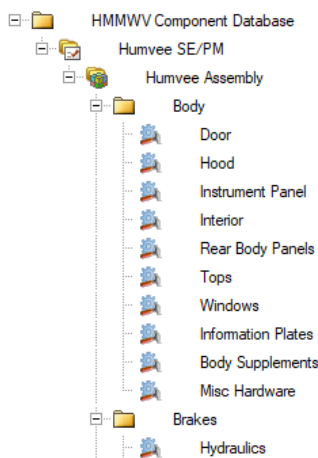


*Figure 81: Humvee Component Database PBS*

Just as with the firearms databases, this data was also verified. The major drivers for the engine were compared against PRICE®'s internal historically-based recommendations. The Kascar engine data was only two pounds more than the actual mass, and when calibrated it yielded a manufacturing Complexity for Structure well within 5% (4.15 compared to the recommended value of 4.30). Additionally, the cost of a non-armored Humvee was quoted at $70k in 2011 [21]. Our results indicated a cost of just over $74k for a system that was about a hundred pounds less than the actual curb weight of 5200 pounds [22]. Results were well within the acceptable range of percent error.

## Limitations

There were many reasons as to why test cases did not match the actuals. Generalizations were made when data was applied to our model. For example, cost is approximated based on price. While the team did develop a methodology for calculating profit for the manufacturer, there is no way to account for how much retailer profit is in MSRP or other types of prices listed online. In addition, all weights were treated as structure and not electronics. While this was not an issue for the firearms data, it was certainly a shortcoming of the Humvee component database. However, the validation of data has demonstrated that even with these approximations and extra fees baked in to the calibrated Manufacturing Complexity for Structure, these estimated costs from the websites can still provide for acceptable cross-checks or price-to-win scenarios.

Another limitation is the team's use of certain tools. A more sophisticated way to obtain large amounts of information is through using an Application Programming Interface (API). These are useful because they are tied directly to the databases in websites [23]. API requests must be written a certain way. Then, the output, usually a large mass of text, would need to be organized into a usable database. In addition, API's are intended to output a stable structure, unlike retailer websites which can change without warning. While one must still be careful to ensure it is legal to store or use the information contained in the API, many APIs are meant to be available to the public [24]. As an example, Best Buy has an API for its product specifications and sale prices. Users must request a key before accessing the information [25]. As of this writing the team is still learning how to use APIs.

## Future Directions

The team would like to expand the scope of the project and create more databases. Prospective projects include cloud computing pricing, microcircuits, and space components. In addition to web scraping, APIs will be used to extract the information. If possible, the current databases will be expanded with more data points and categories. Not only would this be helpful for future estimates, but it would be beneficial to find more sources to cross-check the data.

**Conclusion**

This work has outlined a methodology for obtaining large amounts of data from the internet by utilizing web crawling and scraping. These steps have been successfully applied to various types of equipment of interest to many in the cost community. While imperfect, early results are promising for using the databases as a cross check for early estimates and price-to-win scenarios.  The team intends to continue applying and improving this approach as new sources of publicly available relevant data become available.

**Works Cited**

[1]  "Web Scraping," Techopedia, [Online]. Available: https://www.techopedia.com/definition/5212/web-scraping. [Accessed 13 February 2019].

[2]  A. Foote, "Automated Data Collection Using Open Source Web," 2017. [Online]. Available: http://www.iceaaonline.com/ready/wp-content/uploads/2017/07/MM03-PPT-Foote-Automated-Data-Collection-Web-Crawling.pdf. [Accessed 13 February 2019].

[3]  V. Fedak, "Big Data Scraping vs Web Data Crawling," 15 February 2018. [Online]. Available: https://techburst.io/big-data-scraping-vs-web-data-crawling-4ef5a71d7888. [Accessed 13 February 2019].

[4]  "https://stats.oecd.org/glossary/detail.asp?ID=3309," Organisation for Economic Co-operation and Development, 13 March 2003. [Online]. Available: https://stats.oecd.org/glossary/detail.asp?ID=3309. [Accessed 13 February 2019].

[5]  W. Kenton, "Manufacturer's Suggested Retail Price - MSRP," Investopedia, 21 April 2018. [Online]. Available: https://www.investopedia.com/terms/m/manufacturers-suggested-retail-price-msrp.asp. [Accessed 13 February 2019].

[6]  V. Krotov and L. Silva, "Legality and Ethics of Web Scraping," September 2018. [Online]. Available: https://www.researchgate.net/publication/324907302_Legality_and_Ethics_of_Web_Scraping. [Accessed 13 February 2019].

[7]  "The Web Robots Pages," Web Robot Pages, [Online]. Available: http://www.robotstxt.org/. [Accessed 13 February 2019].

[8]  Hyatt Gun Store, [Online]. Available: https://www.hyattgunstore.com/. [Accessed 2019 February 13].

[9]  RapidMiner, Inc, [Online]. Available: https://rapidminer.com/. [Accessed 2019 February 13].

[10] Octopus Data Inc., [Online]. Available: https://www.octoparse.com/. [Accessed 13 February 2019].

[11] J. Koshy, "7 Best Software tools to Acquire Data Without Coding," [Online]. Available: https://www.promptcloud.com/blog/best-software-tools-acquire-data. [Accessed 13 February 2019].

[12] Hyatt Gun Store, "Charter Arms Undercover Crimson Trace 38 Special Revolver," [Online]. Available: https://www.hyattgunstore.com/charter-arms-undercover-crimson-trace-38-special-revolver.html. [Accessed 13 February 2019].

[13] Mozilla, "XPath," [Online]. Available: https://developer.mozilla.org/en-US/docs/Web/XPath. [Accessed 13 February 2019].

[14] J. Goyvaerts. [Online]. Available: https://www.regular-expressions.info/tutorial.html. [Accessed 13 February 2019].

[15] Bureau of Alcohol, Tobacco, Firearms and Explosives, "Firearms - Guides - Importation & Verification of Firearms, Ammunition and Implements of War - Firearms and Ammunition Excise Tax (FAET)," 22 September 2016. [Online]. Available: https://www.atf.gov/firearms/firearms-guides-importation-verification-firearms-ammunition-and-implements-war-firearms. [Accessed 13 February 2019].

[16] Yahoo! Finance, "Sturm, Ruger & Company, Inc. (RGR)," [Online]. Available: https://finance.yahoo.com/quote/RGR/financials?p=RGR. [Accessed 13 February 2019].

[17] J. Brauer, "The US Firearms Industry: Production and Supply," February 2013. [Online]. Available: http://www.smallarmssurvey.org/fileadmin/docs/F-Working-papers/SAS-WP14-US-Firearms-Industry.pdf. [Accessed 19 February 2019].

[18] J. M. Buol, "The 8 Pistols That Battled to Win the Army's XM17 MHS Competition," Tactical Life, 11 September 2017. [Online]. Available: https://www.tactical-life.com/firearms/xm17-mhs-army-pistol/. [Accessed 13 February 2019].

[19] D. Crane, "SIGARMS and Heckler & Koch/HK Defense Win Major Pistol Contracts with DHS," Defense Review, 16 August 2005. [Online]. Available: http://www.defensereview.com/sigarms-and-heckler-kochhk-defense-win-major-pistol-contracts-with-dhs/. [Accessed 13 February 2019].

[20] Kascar, LLC, [Online]. Available: https://real4wd.com/. [Accessed 13 February 2019].

[21] C. Keyes, "Steep cost of military vehicles outlined in Army report," CNN, 27 January 2011. [Online]. Available: http://www.cnn.com/2011/US/01/27/army.vehicle.costs/index.html. [Accessed 13 February 2019].

[22] Wikipedia, "Humvee," [Online]. Available: https://en.wikipedia.org/wiki/Humvee. [Accessed 25 February 2019].

[23] O. Lam, "Using APIs to collect website data," Medium, 27 June 2018. [Online]. Available: https://medium.com/pew-research-center-decoded/using-apis-to-collect-website-data-b7fc340d59e3. [Accessed 13 February 2019].

[24] C. Wodehouse, "Public APIs vs. Private APIs: What's the Difference?," UpWork, [Online]. Available: https://www.upwork.com/hiring/development/public-apis-vs-private-apis-whats-the-difference/. [Accessed 13 February 2019].

[25] Best Buy, "Overview," [Online]. Available: https://bestbuyapis.github.io/api-documentation/#overview. [Accessed 19 February 2019].