

## Old Cars Data

### A Case Study in Data Cleaning and Provenance

Rahul Varma  
Dept. of Computer Science  
University of Illinois-UC  
Champaign, IL, USA  
[Rvarma4@illinois.edu](mailto:Rvarma4@illinois.edu)

Hemendra Choudhary  
Dept. of Computer Science  
University of Illinois-UC  
Champaign, IL, USA  
[hsc4@illinois.edu](mailto:hsc4@illinois.edu)

### Final Project Report

#### Abstract

**Data Preparation** - Data cleansing is the process of spotting and rectifying inaccurate or corrupt data from a database. The process is mainly used in databases where incorrect, incomplete, inaccurate or irrelevant part of the data are identified and then modified, replaced or deleted. Business enterprises largely rely on data whether it is the integrity of customers' addresses or ensuring accurate invoices are emailed or posted to the recipients. To ensure that the customer data is used in the most productive and meaningful manner that can increase the intrinsic value of the brand, business enterprises must give importance to data quality.

.

#### 1 INTRODUCTION

This report summarizes our experience with an end- to-end data preparation work-flow; in practice of data Cleaning and Provenance establishment. We use tools and techniques introduced in CS598 Theory and Practice of Data Cleaning with a real world dataset and document the whole work-flow along with findings. Tools used include OpenRefine, SQLite and YesWorkflow.

#### 2 DATASET OVERVIEW AND INITIAL ASSESSMENT

##### 2.1 The Dataset

The data was scraped from several websites in Czech Republic and Germany over a period of more than a year. Originally I wanted to build a model for estimating whether a car is a good buy or a bad buy based on the posting. But I was unable to create a model I could be satisfied with and now have no use for this data. I'm a great believer in open data, so here goes.

##### 2.2 Content

The scrapers were tuned slowly over the course of the year and some of the sources were completely unstructured, so as a result the data is dirty, there are missing values and some values are very obviously wrong (e.g. phone numbers scraped as mileage etc.)

There are roughly 3,5 Million rows and the following columns:

1. maker - normalized all lowercase
2. model - normalized all lowercase
3. mileage - in KM
4. manufacture\_year

5. engine\_displacement - in ccm
6. engine\_power - in kW
7. body\_type - almost never present, but I scraped only personal cars, no motorcycles or utility vehicles
8. color\_slug - also almost never present
9. stk\_year - year of the last emission control
10. transmission - automatic or manual
11. door\_count
12. seat\_count
13. fuel\_type - gasoline, diesel, cng, lpg, electric
14. date\_created - when the ad was scraped
15. date\_last\_seen - when the ad was last seen. Our policy was to remove all ads older than 60 days
16. price\_eur - list price converted to EUR

### 2.3 Potential use-cases of cleaned Data

- a) Which factors determine the price of a car?
- b) With what accuracy can the price be predicted?
- c) Can a model trained on all cars be used to accurately predict prices of models with only a few samples?

## 3 DATA CLEANING WITH OPENREFINE

1. TRIM and COLLAPSE WHITE SPACES. It is very common to see unnecessary white spaces in datasets. A lot of times white spaces are hidden at the beginning or the end of a string, and sometimes they are hidden as two consecutive white spaces in a phrase. Here's what you can do to help clean up white spaces.

- Trim all the leading and trailing white spaces in ALL columns that are texts (strings). This includes the maker, model, body\_type, transmission columns.
- Collapse consecutive white spaces in ALL columns that are texts (strings).

2. NUMBER. Incorrect data types is almost always the second thing you inspect in a dataset. Usually numeric data will be seen (or converted to) as text data in a lot of platforms. To correct these, you can do the following:

- Transform all columns that should be in numeric form to number
- Note that whatever you have converted to number will be shown in green

3. Fix uppercase/lowercaseUpper case on columns (eg. maker id.)

4. Convert the date related columns into datetime instead of string.

5. Check for clustering and see if better results can be achieved.

6. DELETE IRRELEVANT COLUMN. we noticed that there are almost no values on the color\_slug column, and we decided that this is an irrelevant column for further analysis.

7. Using power of facets- find crazy mileage ( eg> 200k ) and create a column which points if the mileage is crazy high

8. Similarly create a column for crazy\_new if the mileage is less that a threshold say 100 km.

9. Rectify the manufacturing date of outliers ( eg. made in 1300 ) and use averaging techniques to predict mfg. date.

10. Backfill missing data on engine transmission type (automatic / manual )