

Does Anybody Really Know What Time It Is?

Automating the Extraction of Date Scalars

Richard Wesley

Vidya Setlur

Dan Cory

ABSTRACT

With the advent of modern data visualization tools, data preparation has become a bottleneck for analytic workflows. Users who are already engaged in data analysis prefer to stay in the visualization environment for cleaning and preparation tasks whenever possible, both to preserve analytic flow and to take advantage of the visualization environment to spot inconsistent data. This has led many visualization environments to include simple data preparation functions such as scalar parsing, pattern matching and categorical binning in their analytic toolkits.

One of the most common parsing tasks is extracting date and time data from string representations. Several databases include date parsing “mini-languages” to cover the wide range of possible formats. Tableau has recently provided a DATEPARSE function with a single syntax that is translated into the nearest equivalent in the underlying database. Subsequent analysis of customer usage of this function shows that the parsing language syntax was difficult for users to master.

In this paper, we present two algorithms for automatically deriving date format strings from a column of data in our syntax, one based on minimum entropy and one based on natural language modeling. Both have similar accuracies of over 90% on a large corpus of date columns extracted from Tableau Public and are in substantial agreement with each other. The minimal entropy approach can also produce results below the user’s perceptual threshold, making it suitable for interactive work.

Categories and Subject Descriptors

D.3.3 [Data Processing]: Data Cleaning, Natural Language Processing

General Terms

Algorithms, Performance

Keywords

Date parsing, automated cleaning, structure extraction

1. INTRODUCTION

In recent years, there has been a growth of interest in data visualization technologies for human-assisted data analysis using systems such as [1,10,11]. While computers can provide high-speed and high-volume data processing, humans have domain knowledge and the ability to process data in parallel, notably by using our visual systems. Most importantly, humans provide the definition of what is valuable in an analysis. Accordingly, human/computer analytic systems are essential to extracting knowledge of value to humans and their societies from the large amounts of data being generated today.

1.1 Interactivity

Visualization systems are most effective when they are interactive, thereby allowing a user to explore data and connect it to their domain knowledge and sense of what is important without breaking cognitive flow. In recent years, a number of such systems have been developed, both by the academic community and by the commercial sector. Exploration of data consists not only in creating visual displays, but also in creating and modifying domain-specific computations. Consequently, most data visualization systems include facilities for defining such calculations as part of the data model being analyzed. The most effective systems allow users to define these calculations as part of the analytic interaction, which permits the user to stay in the flow of analysis [9].

During the analytic process, a user may discover that parts of the data are not yet suitable for analysis. Solutions to this problem are often provided by data preparation tools external to the visual analysis environment, which requires the user to break their analytic flow, launch another tool and reprocess their data before returning to their analysis. If the user does not own this process (e.g. it is the responsibility of another department), then there can be significant delays (including “never.”) More subtly, the result of updated external processing may not be compatible with the user’s existing work, which can lead to more time lost reconciling the new data model with the existing analysis.

From the user’s perspective, the boundary between preparation and analysis is not nearly so clean cut. Bad data is often discovered using visual analysis techniques (e.g. histograms or scatter plots) and it is most natural for the user

to “clean what she sees” instead of switching to a second tool. This leads to an “adaptive” process whereby users will prefer leveraging existing tools in the analytics environment (no matter how well suited to the task) over switching to another application. Thus a well-designed interactive visual analysis environment will provide tools that enable users to perform such cleaning tasks as interactively as possible.

1.2 The Tableau Ecosystem

In this paper, we shall be looking at an example of this problem in the context of the Tableau system. Tableau is a commercial visual analysis environment derived from the Polaris [1] system developed at Stanford. In addition to an interactive desktop application for creating data visualizations, the Tableau ecosystem also includes servers for publishing and sharing interactive visualizations. These servers can be privately operated by individual customers or hosted in the cloud.

1.2.1 Tableau Public

Of particular relevance in the present work is a free version of the server environment called Tableau Public [12]. Tableau Public (or “Public”) allows authors to publish data visualizations that can be shared with the general public. The published data sets are limited to 100K rows and must be stored using the Tableau Data Engine (or “TDE”) [2, 3], but are otherwise free to use the full analytic capabilities of the Tableau system. In particular, the TDE provides native support for all analytic calculations defined by the authoring and publishing components. Moreover, visualizations published to Public are uploaded as complete Tableau workbooks, including the data model developed as part of the analysis. This makes it a rich source of data on the analytic habits and frustrations of Tableau users.

1.2.2 Functions

While the existing standards for query languages are helpful for defining a core set of row-level functions, there are many useful functions beyond the standards that are only explicitly supported by a subset of RDBMSes, often with different names. And even when a database does not provide an implementation of a particular function, it is usually possible to generate it as an inline combination of existing functionality. Tableau’s calculation language is designed to be a lingua franca that tames this Babel of dialects by providing a common syntax for as many of these functions as possible.

1.3 DATEPARSE

Among the recent additions to the Tableau function library is a function called DATEPARSE, the usability of which is the focus of the work in this paper.

1.3.1 Scalar Dates

The SQL-99 standard defines three temporal scalar types: DATE, TIMESTAMP and TIME. They are typically implemented as fixed-point types containing an offset from some epoch (*e.g.* Julian Days.) This makes them compact to store and allows some temporal operations to be implemented very efficiently using arithmetic. Thus from the RDBMS

perspective, representing dates in scalar form provides numerous benefits for users, both in terms of available operations and query performance.

Tableau models the first two of these types as “Date” and “Date & Time”; the third (pure time) is folded into Date & Time by appending it to a fixed date of 1899-12-30. From the analytic perspective, date types are dimensional (*i.e.* independent variables) and can be used as either categorical (simply ordered) or quantitative (ordered with a distance metric) fields. Categorical dates have a natural hierarchy associated with them generated by calendar binning. The visualization in Figure 11 shows an example of a bar chart employing binned categorical dates in a year/quarter hierarchy.

1.3.2 Parsing

One of the oldest data preparation problems observed in Tableau is the parsing of date scalars so that users can perform these kinds of analyses. Training materials from the early days of the company included examples of how to convert columns of integers in the form yyyyMMdd to date scalars. The documented solution at the time was to convert the integer to a string and perform some locale-dependent string operations before casting the string back to a date. Unfortunately, this approach had a number of problems:

- String operations are notoriously slow compared to scalar operations (typically 10-100x slower)
- Default parsing of date formats is locale-dependent, and may not work when the workbook is shared across an international organization (*e.g.* between the US and European offices)
- The expression code was hard to understand and maintain because it used a verbose, general-purpose string-handling syntax instead of a domain language.

This is but a single format. Our studies of the workbooks on Public suggest that there are hundreds of distinct temporal date formats in user data sets. Some are common, but others can be quite idiosyncratic. Table 1 shows a selection of unusual date formats found on Public.

Our first solution to this problem was to add a new function to the Tableau calculation language called DATEPARSE. This function would take a string column and convert it to a datetime using a special purpose domain language for describing dates. Such functions exist in a number of databases (*e.g.* Oracle, Postgres and MySQL) and adding it to the TDE was straightforward, so it was a natural addition to the function library. We chose the date parsing syntax defined by the International Components for Unicode (or ICU) project [8] for the common domain language because it mirrored what was available in the Tableau code base, and translating this syntax into the formats used by various database vendors (*e.g.* MySQL, Oracle, Postgres) was relatively painless. There are a few patches of non-overlapping functionality, but they tend to be obscure and we provide warnings when functionality is missing.

1.3.3 Usability

After shipping this new function to our user base, we investigated how it was being applied in the field. We examined a few months of workbooks that had been uploaded to Public after the release to see if and how it was being used. We found that although there were a number of new calculations using DATEPARSE, the error rate for the domain language syntax was about 15%. DATEPARSE was capable of solving the problem, and some users were able to discover it, but the syntax did not appear to be easy to use reliably.

1.4 Automated Format Extraction

One solution might have been to design a graphical environment that enabled users to construct valid patterns. This would have involved a substantial development investment with no guarantee that the result would be correct if the user misunderstood the environment. Instead, we have developed two algorithms for automating the derivation of the format string, each of which uses a different machine learning technique. Both algorithms result in over 90% parsing accuracy and 100% syntactic correctness because they are machine generated.

1.5 Related Work

Data preparation has been a known analytic bottleneck since at least the description of the Potter's Wheel system [4]. Since then, several other interactive data preparation systems have been proposed [5,6]. While effective, these systems all make assumptions about possible date formats (e.g. the domain system in [4]), which we suggest are too restrictive for real world data. The Minimum Descriptive Length technique was first proposed in [7]. We have built on the version described in [4]. Related work on parsing languages, outlier detection.

1.6 Overview

The rest of this paper is organized as follows: The next section introduces the parameters of the problem space. The following two sections describe the two different algorithms, one using Minimum Descriptive Length and the other using Natural Language Processing. In section 5, we evaluate the algorithms on a corpus of 30K columns, both by sampling the outputs manually and then by using the algorithms to validate each other. We then discuss future work in section 6 and conclude in section 7.

2. PARAMETERS

Before delving into detailed descriptions of the two algorithms, we describe the common elements of the problem they are intended to solve. These are the input data, the output language and the interpretation of the output language for partial dates.

2.1 Input Data

The training data and test data are taken from Public data sets. We selected a large number of string columns whose names suggested that they might contain temporal data. We included words from multiple languages besides English and extracted the column's locale (actually the collation) along with the data.

2.1.1 NULL Filtering

Each column was then reduced to a maximum sample set of 32. We excluded domain values that appeared to be representations of NULL:

- Values containing the substring NULL
- Values containing no digits and at most one non-whitespace character (e.g. " / ")
- Values on a list of empirically determined common NULL values (e.g. 0000-00-00, NaN)

Columns that had no remaining valid samples were discarded.

2.1.2 Sampling

The remaining samples were hashed and sorted on the hash value, with the top 32 values being retained as the column's sample. We can increase the sample size if needed, but for dates, it appears to be an adequate number. In any case, the median number of non-null rows per domain in our test data is 50, so increasing the sample size would have little benefit.

2.1.3 Numeric Timestamps

After NULL filtering, there remained one common class of date representation that was unsuited for parsing via our date format syntax, namely numeric timestamps. This included Unix epoch timestamps (expressed as second, millisecond or microsecond counts from 1970-01-01) and Microsoft Excel timestamps (expressed fractional days since 1900-01-01). A simple test to check that the column was numeric and inside a specific range of values representing recent dates allowed us to tag these columns to avoid analyzing them further (and incidentally to identify them for generating a simple date extraction calculation for the user.)

2.2 The ICU Date Format Language

The date format syntax we selected is the one provided by the open-source project. We chose it because we were already using ICU in the code base, we had access to the source code and it provides localized date part data for a large number of languages.

The syntax is documented at the ICU web site [8]. While fairly complete, it has a few limitations that we ran into when evaluating the algorithms:

- No support for 4 letter year abbreviations (e.g. Sept.)
- No support for ordinal days (e.g. July 4th)
- No support for quarter postfix notation (e.g. 2Q)
- No support for variant meridian markers (e.g. a.m.)

These limitations did not affect the results significantly, and in the future we hope to submit ICU extensions to handle some of these issues.

One other quirk of the ICU syntax may be a contributing factor to the user confusion around writing correct ICU date formats. The use of lexicographical case in the meta-symbols of the format language can be confusing (e.g. `y` is used for years, but `M` is used for months while `m` is used for minutes.) Another advantage of an automated algorithm is that it hides such problems from most users, significantly improving the usability of the function.

2.3 Partial Dates

Many of the date formats that we encountered were incomplete dates, which necessitated creating rules for what date scalar they represented.

ICU's date parsing APIs allow the specification of default values for parts when a format does not contain them. In our implementation, all time fields are set to 0 (midnight) and the date fields are set based on whether the format contains any date part specifications. When date parts are present, we use 2000-01-01 as the set of default date parts as it is the start of a leap year. When dealing with pure time formats, we use 1899-12-30, (which is the date used internally by Tableau to signal pure time values.)

ICU will also parse Time Zones (which we recognize but do not use in Tableau) and Quarters (which we interpret as the first month of the period.) RDBMSes such as Oracle and Postgres that support time zones will be able to take advantage of the generated time zone field.

3. MINIMAL DESCRIPTIVE LENGTH

The first algorithm is a Minimum Descriptive Length [7] approach derived from the domain system presented in [4]. We describe a number of extensions to their structure extraction system to support more complex redundancy, non-English locales, improved performance and date-specific pruning.

3.1 Domains

[4] presents an algorithm for deriving a common structure for a set of strings by breaking each string down into a sequence of domains. These domains are described by an interface that includes:

- A required inclusion function to test for membership in the domain (`match`)
- An optional function to compute the number of values in the domain with a given length (`cardinality`)
- An optional function to update statistics for the domain based on a given value (`updateStatistics`)
- An optional function to prevent consideration of a domain that is redundant (`isRedundantAfter`).

In our approach, we implement all of these functions, but with significant changes to the last one, which we will describe below.

With this interface, we can now define a set of domains for each date part that we wish to be able to parse. These are mostly straightforward enumerations and numeric ranges,

each tagged with the ICU format code. Since the ICU parser is flexible about parsing single or double digit formats, we use double-digit formats, but accept one or two digits. One important exception to this rule is for years, which are fixed width fields (2 or 4).

We also included some enumerated domains for handling constant strings and some simple regular expression domains for delimiters such as spaces, punctuation or alphabetic characters. We found that the inclusion of arbitrary numeric domains caused the run time to grow exponentially as the number of possible matches could not be pruned intelligently. This restriction extends to domains that can contain arbitrary digit sequences (such as `any`). Because of this restriction, the algorithm cannot extract non-date numeric fields.

3.2 Redundancy Extensions

A difficulty in using this kind of structure extraction is that the algorithm for enumerating structures is exponential in the number of domains. This is especially true in the date format problem because there are identical domains (e.g. months and meridian hours), nearly identical domains (e.g. days and hours) and there are often no field delimiters (e.g. 2012Mar06134427). To handle this, we have extended the existing pruning API with two other sets of domain identifiers:

- A set of *prunable* identifiers, which are not allowed to precede the domain. For example, once we have a month field, no other month fields should be generated. Each month domain therefore lists all the month domains in its prunable set.
- A set of *context* identifiers, one of which must have been previously generated before the domain is considered. For example, a meridian domain can only be generated once an hour field has been found, but there may be other intervening fields.

3.3 Performance

Structure enumeration is computationally expensive, so we added a number of enhancements to the original domain extraction algorithm to keep the run time low enough for interactivity.

3.3.1 Domain Characteristics

Date domains typically have small widths, so we found it advantageous to provide the shortest and longest match sizes for use in structure enumeration and matching.

Date domains are also often uniform in that adding more characters to a mismatch will not help. For example, a 2-digit day domain that does not match a 1-letter substring will not be able to generate a match by adding more characters.

3.3.2 Parallel Evaluation

We have also identified two opportunities for parallel computation during structure enumeration.

The first computationally expensive operation is the enumeration of structures for a given sample. To parallelize

this step, each thread is given a subset of the samples and independently produces a set of candidate domains. When all threads have completed, the duplicates are removed to produce a single list of candidate structures.

The second computationally expensive operation is the evaluation of each generated structure over the entire sample set. This includes computation of the MDL, recording of domain statistics and parameterizing the structure. These tasks are simple to parallelize, because there are no overlaps between the data for each structure.

3.4 Unparameterization

Domain parameterization is an important part of generating compact representations via MDL, but it creates problems for date recognition. For example, if a set of dates contains a constant month string (e.g. all values are in September) it is important to keep track of the month name domain. Consequently, when we parameterize a constant generic $\langle \text{Word} \rangle$ domain, we tag it with the date part domain that it matches (if any). We then need to apply an additional pruning step to remove any structures that also found an equivalent domain (e.g. two-digit month). These rules are equivalent to the context-based redundancy rules above, but have to be applied again after parameterization of generics.

3.5 Global Pruning

The pruning rules used for the structure extraction reduce the search space dramatically, but they are also contextual and can only look backwards. The domains also contain a fair amount of ambiguity that requires the application of domain knowledge. We therefore found it necessary to add some post-generation global pruning rules:

- The set of date parts cannot contain place value gaps (e.g. structures that have year and day without month are removed.)
- Similarly, the set of time parts cannot contain place value gaps and must also be in place value order (times are never written in orders such as mhs.)
- The existence of time parts cannot make dates incomplete (e.g. patterns like year-day-hour are removed.)
- Two digit years require special handling. In particular, they cannot appear adjacent to a two-digit field if the structure contains punctuation. (This can come up in some small early 21st century year domains where a two-digit year can masquerade as almost any numeric field.)

These global pruning rules are simple and intuitive, but are essential to further reducing the search space.

3.6 Locale Sensitivity

Providing an acceptable international user experience requires correctly handling the locale of the column text. Accordingly, at the start of the structure extraction we use the column locale to create a set of domains containing locale-sensitive strings such as month names. We also use the

locale to map these strings to upper- and lower-case in addition to the ICU mixed case strings. This enables us to accurately compute MDL statistics without having to map the candidate strings at runtime (which would be slow).

Knowing the locale of a string is not always helpful. In our data set, we found numerous cases where the locale was specified (e.g. Sweden) but the data was actually in English. Accordingly, we test both locales and rank the combined results. Tableau currently does not support non-Gregorian calendars. Therefore, we have built this system for the Gregorian calendar only. ICU does support non-Gregorian calendars and we expect this system could be applied to them as well.

3.7 Ranking

MDL structure analysis naturally produces a ranked list of format candidates, but we have found that a number of other properties of the formats should be preferred over simple compactness:

- Since we have a set of samples, we can apply the candidate format to the strings to see how well it performs. Formats with fewer parse errors are preferred.
- Date parts can be considered a place-value system, so we prefer more significant components (e.g. month-day-year over hour-minute-second).
- If two formats from different locales give the same results, prefer the original column locale. The sample set may have missed an example where this could be important.
- If the format has an ambiguous date order (e.g. all days are less than 12), then prefer the default date order of the locale. Again, the sample set may have missed a counterexample, so this is the best option.
- Once these semantic preferences have been considered, we then prefer the more compact (MDL) representation. The output of the algorithm is now an ordered list of formats and associated locales. These can then be used to drive a user interface that allows the user to choose between the possibilities or the top-ranking format can simply be used automatically.

4. NATURAL LANGUAGE PROCESSING

4.1 Context-Free Grammar

ICU date-time formats are well defined both structurally and semantically, and can be defined by a context-free grammar (CFG). A CFG is commonly defined as a set of productions or rules of the form $A \rightarrow \alpha$ where A is a variable, α is a sequence of variables and terminal symbols (the tokens that make up the alphabet of the language) plus null (ϵ), and the production symbol (\rightarrow) indicates that the variable A can be expanded into α . A CFG can be formally specified with four components: V , T , P , and S , where V is the set of variables, T the set of terminal symbols, P the set of productions, and S the set of available start symbols (a non-empty subset of V) [cite].

Pattern-recognition problems such as parsing date and time formats initiate from observations generated by some structured stochastic process. In other words, even if the initial higher-level production rule of the grammar is known (i.e. date, time or date-time), there could be several directions that the parser resolve to. For example, in a date string 5/6/2015, the pattern could either be M/d/yyyy or d/M/yyyy.

Probabilistic context-free grammars (PCFGs) have provided a useful method for modeling such uncertainty [cite]. Once we have created a PCFG model of a process, we can apply existing PCFG parsing algorithms to identify a variety of date-time formats. However, the parser's success is often limited in the types of the dominant patterns that it can identify. In addition, the standard parsing techniques generally require specification of a complete observation sequence. In many contexts, we may have only a partial sequence available (e.g. an incomplete entry). Finally, we may be interested in computing the probabilities of date-time patterns that the grammar may not explicitly define. To extend the forms of evidence, inferences, and pattern distributions supported, we need a flexible and expressive representation for the distribution of structures generated by the grammar. We adopt Bayesian networks for this purpose, and define an algorithm to generate a probabilistic distribution of possible parse trees corresponding to a set of date-time patterns as opposed to individual ones.

5. ACKNOWLEDGMENTS

This section is optional; it is a location for you to acknowledge grants, funding, editing assistance and what have you. In the present case, for example, the authors would like to thank Gerald Murray of ACM for his help in codifying this *Author's Guide* and the `.cls` and `.tex` files that it describes.

APPENDIX

A. HEADINGS IN APPENDICES

The rules about hierarchical headings discussed above for the body of the article are different in the appendices. In the `appendix` environment, the command `section` is used to indicate the start of each Appendix, with alphabetic order designation (i.e. the first is A, the second B, etc.) and a title (if you include one). So, if you need hierarchical structure *within* an Appendix, start with `subsection` as the highest level. Here is an outline of the body of this document in Appendix-appropriate form:

A.1 Introduction

A.2 The Body of the Paper

A.2.1 Type Changes and Special Characters

A.2.2 Math Equations

Inline (In-text) Equations

Display Equations

A.2.3 Citations

A.2.4 Tables

A.2.5 Figures

A.2.6 Theorem-like Constructs

A Caveat for the T_EX Expert

A.3 Conclusions

A.4 Acknowledgments

A.5 Additional Authors

This section is inserted by L^AT_EX; you do not insert it. You just add the names and information in the `\additionalauthors` command at the start of the document.

A.6 References

Generated by bibtex from your `.bib` file. Run latex, then bibtex, then latex twice (to resolve references) to create the `.bbl` file. Insert that `.bbl` file into the `.tex` source file and comment out the command `\thebibliography`.

B. MORE HELP FOR THE HARDY

The `acm_proc_article-sp` document class file itself is chock-full of succinct and helpful comments. If you consider yourself a moderately experienced to expert user of L^AT_EX, you may find reading it useful but please remember not to change it.