

Easier AI

Simpler, smarter, faster, more flexible and understandable

Tim Menzies and the EZR gang North Carolina State University

Package: <https://pypi.org/project/eZR/0.1.0/>

Source: <http://github.com/timm/eZR>

Latex: <http://github.com/timm/eZR-tex>

© 2024 by Tim Menzies and the EZ gang is licensed under [Creative Commons Attribution-ShareAlike 4.0 International](#) 

We show how to simplify the process of AI analytics, which involves extracting high-quality insights from large quantities of data. We advocate for more efficient and accessible analytical methods that require fewer data samples and less complexity. This allows for easier verification and understanding of results. We highlights the benefits of using incremental methods in building models that can provide valuable insights with minimal data.

This work can be viewed as a (polite) protest against the pre-vailing preference for complex solutions in the industry, suggesting that simplicity could offer more practical and appreciable benefits but is often overlooked due to commercial interests. The call is for a shift towards simplicity in analytics, making it faster, smarter, and more flexible, to better serve practical needs and enhance comprehensibility.

Introduction

Suppose we want to use data to make policies— about what to do, what to avoid, what to do better, etc etc. How to do that?

This process is called *analytics*, i.e. the reduction of large amounts of low-quality data into tiny high-quality statements. Think of it like “finding the diamonds in the dust”.

Many people have been doing data-driven analytics for decades. So it seems

the right time to ask how can we make analytics simpler, smarter, faster, more flexible and more understandable?

For example, according to Tom Zimmermann (from Microsoft Research), there are many things we want to do with analytics:

	Past	Present	Future
Exploration Find important conditions.	Trends Quantifies how an artifact is changing. Useful for understanding the direction of a project. <ul style="list-style-type: none">Regression analysis.	Alerts Reports unusual changes in artifacts when they happen. Helps users respond quickly to events. <ul style="list-style-type: none">Anomaly detection.	Forecasting Predicts events based on current trends. Helps users make pro-active decisions. <ul style="list-style-type: none">Extrapolation.
Analysis Explain conditions.	Summarization Succinctly characterizes key aspects of artifacts or groups of artifacts. Quickly maps artifacts to development activities or other project dimensions. <ul style="list-style-type: none">Topic analysis.	Overlays Compares artifacts or development histories interactively. Helps establish guidelines. <ul style="list-style-type: none">Correlation.	Goals Discovers how artifacts are changing with respect to goals. Provides assistance for planning. <ul style="list-style-type: none">Root-cause analysis.
Experimentation Compare alternative conditions.	Modeling Characterizes normal development behavior. Facilitates learning from previous work. <ul style="list-style-type: none">Machine learning.	Benchmarking Compares artifacts to established best practices. Helps with evaluation. <ul style="list-style-type: none">Significance testing.	Simulation Tests decisions before making them. Helps when choosing between decision alternatives. <ul style="list-style-type: none">What-if? analysis.

Note all the different algorithms in all the boxes. Before we knew better:

- We used to study those algorithms as separate things. Now we see them as very similar things, all of which call the same underling structure. This means that once we code one of them, we can quickly code up the rest.
- We ’d explore 100,000s of possibilities to find patterns in the data (e.g. during a "what-if" query). But these days, I can do the same analysis with

30 samples, or less¹ This means if someone wants to check my conclusions, they only need to review a few dozen samples. Such a review was impossible using prior methods since the reasoning was so complicated.

Why can I do things so easily? Well, based on three decades of work on analytics [1] (which includes the work of 20 Ph.D. students, hundreds of research papers and millions of dollars in research funding) I say:

- When building models, there are incremental methods that can find models after very few samples.
- This is because the main message of most models is contained in just a few variables [1].

I'm not the first to say these things². So it is a little strange that someone else has not offer something like this simpler synthesis. But maybe our culture prefers complex solutions:

Simplicity is a great virtue but it requires hard work to achieve it and education to appreciate it. And to make matters worse: complexity sells better.

– Edsger W. Dijkstra

By making things harder than they need to be, companies can motivate the sale of intricate tools to clients who wished there was a simpler way. Well, maybe there is.

I'm not the first to say these things³. So it is a little strange that someone else has not offer something like this simpler synthesis. But maybe our culture prefers complex solutions:

Simplicity is a great virtue but it requires hard work to achieve it and education to appreciate it. And to make matters worse: complexity sells better.

– Edsger W. Dijkstra

By making things harder than they need to be, companies can motivate the sale of intricate tools to clients who wished there was a simpler way. Well, maybe there is.

```
1 import re,ast
2 from typing import Any,Iterable,Callable
3 from fileinput import FileInput as file_or_stdin
4 #-----
5 def coerce(s:str) -> Any:
6     "s is a int,float,bool, or a string"
7     try: return ast.literal_eval(s) #
8     except Exception: return s
9
10 def csv(file=None) -> Iterable[Row]: a
11     "read from file or standard input"
12     with file_or_stdin(file) as src:
13         for line in src:
14             line = re.sub(r'([\n\t\ ]|\.*)', '', line) # no comments,white space
15             if line: yield [coerce(s.strip()) for s in line.split(",")]
16 #-----
17 class COLS(OBJ):
18     """Turns a list of names into NUMs and SYMs columns. All columns are held
19     in i.all. For convenience sake, some are also help in i.x,i.y
20     (for indepedent, dependent cols) as well as i.klass (for the klass goal,
21     if it exists)."""
22     def __init__(i, names: List[str]):
23         i.x, i.y, i.all, i.names, i.klass = [], [], [], names, None
24         for at,txt in enumerate(names):
25             a,z = txt[0], txt[-1] % first and last letter
26             col = (NUM if a.isupper() else SYM)(at=at,txt=txt)
27             i.all.append(col)
28             if z != "X": # if not ignoring, maybe make then klass,x, or y
29                 (i.y if z in "!+-" else i.x).append(col)
30                 if z == "!": i.klass= col
31
32     def add(i,row: Row) -> Row:
33         "summarize a row into the NUMs and SYMs"
34         [col.add(row[col.at]) for col in i.all if row[col.at] != "?"]
35         return row
```

Listing 1: Python example

References

- [1] TJ Menzies, M Dean, JL Black, and JF Fleming. Combining heuristics and simulation models: an expert system for the optimal management of pigs. In *AI'88: 2nd Australian Joint Artificial Intelligence Conference Adelaide, Australia, November 15-18, 1988 Proceedings 2*, pages 48–61. Springer, 1988.

¹Using semi-supervised multi-objective optimization via sequential model optimization (which is all described, later in this document).

²From Wikipedia: The manifold hypothesis posits that many high-dimensional data sets that occur in the real world actually lie along low-dimensional latent manifolds inside that high-dimensional space. As a consequence of the manifold hypothesis, many data sets that appear to initially require many variables to describe, can actually be described by a comparatively small number of variables, likened to the local coordinate system of the underlying manifold.

³From Wikipedia: The manifold hypothesis posits that many high-dimensional data sets that occur in the real world actually lie along low-dimensional latent manifolds inside that high-dimensional space. As a consequence of the manifold hypothesis, many data sets that appear to initially require many variables to describe, can actually be described by a comparatively small number of variables, likened to the local coordinate system of the underlying manifold.