

# An artificial intelligence that can build and discuss statistical models

## Abstract

Lots of things are automated, but the interpretable discussion of data is not, until now...

## 1. Introduction

Recent wins for AI (from wikipedia and brain)

- Sometime - Medical diagnosis (see Heckermann 1991)
- Sometime - Robot in surgery
- 1997 - Deep Blue beats Kasparov
- 2000 - NASA's remote agent program
- 2005 - Stanford win DARPA grand challenge
- 2007 - CMU win DARPA urban challenge
- 2009 - Robot scientist
- 2011 - Watson defeats two greatest Jeopardy! champions
- 2013 - Maths

A lot of these are search based.

Remember programming by optimisation.

We have implemented a search based artificial intelligence. Fully automating statistics involves a sequence of models, discussion of fit, looking at residuals and other model checks and revising the model based on these problems. However, model checks are ultimately used for two purposes

- Checking whether or not the conclusions of the model can be trusted
- Inspiring new models

The second of these is just a good search heuristic in the space of models. <sup>(1)</sup>

## 2. Related work

### 2.1. Random list of things

**Structure learning in Bayesian networks** Similar idea of discovering semantics via model search. Semantics are more vague though i.e. a probability table is not an entirely concise summary

**Linear model** These discover highly interpretable semantics but are limited in expressivity

**Nonparametric additive models** Highly flexible but semantics are vague i.e. can only talk about smooth functions

**Equation learning** Very flexible but semantics of equations do not map onto human understanding e.g. saw tooth vs Fourier decomposition of a saw tooth - which is more human understandable? How would you explain a sensor error with Eureqa style equations. <sup>(2)</sup>

**Deep learning** Again very flexible but the semantics are not usually human interpretable. How can we understand the output of complex representation learning algorithms without human intervention (e.g. recognising that your deep net has become a cat classifier).

**Kernel search** Can use the precise semantics of linear models or the vague semantics of nonparametric additive models and other components along this spectrum. Flexible modelling with components that a human might use to describe what is going on.

### 2.2. What to use when?

**Lots of data and goal is interpolation** Any smoothing device e.g. random forest.

**Highly structured and high dimensional input or output** Use dimensionality reduction or any other method of representation learning. The task is then reduced to an easier regression.

**Parametric modelling of the regression function** Linear models, symbolic regression etc.

(1) If we can get some model checks that are useful that would be great - some sort of prior predictive marginal likelihood check i.e. is this data (un)likely?

(2) Try Eureqa on the solar dataset

110	<b>Nonparametric modelling of the regression</b>	165
111	<b>function but more structured than a smooth-</b>	166
112	<b>ing device</b> Various semi-parametric models, GAM	167
113		168
114	<b>Easily interpretable nonparametric modelling</b>	169
115	This work	170
116		171
117	<b>3. Contributions</b>	172
118		173
119	• A very expressive language of statistical models	174
120	with a concise algebraic structure	175
121		176
122	• Automatic construction of appropriate statistical	177
123	models (search heuristic based on the structure of	178
124	the language)	179
125		180
126	• Automatic discussion of the selected model in nat-	181
127	ural language with tables, figures and text i.e. a	182
128	full statistical report	183
129		184
130	<b>3.1. Things we are not doing</b>	185
131	Producing a system that no human will understand.	186
132		187
133	<b>4. Example analyses</b>	188
134		189
135	<b>5. Discussion and conclusions</b>	190
136		191
137	A Jaynes quote again.	192
138		193
139	Refer to philosophy such as Chinese room to emphasise	194
140	that the language means the system is operating with	195
141	semantic representation and could therefore be said to	196
142	understand what it is saying?	197
143		198
144		199
145		200
146		201
147		202
148		203
149		204
150		205
151		206
152		207
153		208
154		209
155		210
156		211
157		212
158		213
159		214
160		215
161		216
162		217
163		218
164		219