## Life (Data Science) Doesn't Have to be Hard

Mark Kotanchek Evolved Analytics LLC www.evolved-analytics.com



### I never guess. It is a capital mistake to theorize before one has data. Insensibly one begins to twist facts to suit theories, instead of theories to suit facts.

– Sir Arthur Conan Doyle (1859 - 1930), The Sign of Four, A Scandal in Bohemia

### AI ⇒ Augmented Intelligence

- Recognized ~20 years ago that we needed better tools to extract actionable insight from real-world multivariate data
- The human needs to be in the loop with the machine serving not the converse
- If an analysis component can be automated, it should be automated!

## **Real-World Data**

- Multivariate
- Correlated variables
- Unbalanced (closed-loop)
- Nonlinear
- Missing Data
- Wrong Data

### **Questions & Needs**

- Exploratory Data Analysis
- Which variables matter?
- Are there important metavariables?
- What variable combinations are useful?
- How to handle extrapolation and system changes
- How to deploy models?

# GUI and/or Notebook

	DataModeler Tutorials – Wolfram Mathematica 11.2		
< > ▼ _ ♠	DataModeler/tutorial/00Overview	۹	S
DataModeler			
		LIPL -	

### - DataModeler Tutorials

- Preliminaries
  - Preface Release Notes Introduction
- Getting Started
  - An Overview of the Modeling Process FAQ: Frequently Asked Questions (and some illustrations) Symbolic Regression is Not Enough — Context Matters A Meandering Journey through DataModeler Overview of the Function Taxonomy
- Data Exploration & Selection
  - Data Exploration Function Taxonomy Data Visualization Functions Data Statistics Functions Data Subset Selection Functions
- Symbolic Regression Modeling
  - Symbolic Regression Function Taxonomy GPModel Data Structure and Accessor Methods GPModel Creation Evaluating Model Quality Selecting Models Introducing Diversity into the Model Set Evolving Models
- Exploring Models & Model Sets
   Reviewing Model Sets
   Exploring Models
  - Genetic Engineering (Model Optimization)
  - Building Trustable Models (Model Ensembles) Model Response Exploration
  - Model Prediction Performances

### DataModeler Design

 User can choose between a GUI or the classic Mathematica notebook interface

 The functions used by the GUI are exposed for power users and custom workflows

### ~450 functions are exposed and (well)

	Detailed and Malfred Methods to 2	Evolved Analytics
<>• A	DataModeler/guide/DataModeler	Abort Search
lfram Language	Func	Contact Support

### - DataModeler

DataModeler is an incredibly powerful tool. The best way to get started is to explore the tutorials (which include a number of case studies illustrating key features using industrial data). In the short-term, the function overview, quick start and FAQ are good places to get started as well as see some of the capabilities.

Finally, www.evolved-analytics.com is a good source for publications and additional information.

#### Data Exploration

#### Data Adjustments

MakeDataNumeric = MakeDataNumericMapping = MergeInputResponseData = SplitInputResponseData = RescaleData = AugmentData

#### Data Visualization

SmallPlot • CorrelationChart • UnivariatePlot • BivariatePlot • CorrelationMatrixPlot • DataSummaryTable • DataDistributionPlot • DataCompletenessPlot • DataCompletenessMap

#### Statistics

RobustCorrelationMatrix • ConfidenceEllipsoid • AbsCorrelation • MedianAverage • NoisePower • ScaleInvariantNoisePower • SummaryStatistics • WeightedMean • WeightedStandardDeviation

#### Data Subset Selection

NearestDataRecord = ExtractDataSubset = UncorrelatedVariables = SubSample = ConfidenceEllipsoidSelection = ConfidenceEllipsoidSelectionIndices = NumericDataRecords = NumericDataRecordIndices = NonNumericDataRecords = NonNumericDataRecordIndices

#### Outlier Detection

DataStrangeness • DataOutlierAnalysis • DataOutliers • DataOutlierIndices • DataOutlierTable

Model Development (Symbolic Regression)

#### Model Creation

- BuildFunctionPatterns RandomModels CreateModelFromExpression RandomGenomes CreateModelFromGenome • ExtractGenomeSubtrees • MetaVariableModels
- Ensemble Creation
- CreateModelEnsemble
- Evaluating Model and Ensemble Quality

EvaluateModelQuality = UpdateModelQuality = RearrangeModelQuality = EvaluateModelQualityVsMultipleDataSets =

DataModeler.nb

DataModeler

Version Information

Functions

Tutorials

Launch GUI

New Notebook

# GUI used for Demo



### Project co-located with data



## Choose a Target

2 Continent

6 AMRadioStations

14 CellularPhones

18 CropsLandArea

22 ElderlyPopula

Elec

126

130 Male Tobacco Ferce
134 Boys Cigarette Use
138 Gun Homicides
142 Guns per 100 reside

146 Subregion 150 Total 154 Wine (%)

158 Prisoner Rate
162 Iodized Salt Consumption
166 Household Health Fract

170 Health expenditure, priv
174 Hospital Stay Duration
178 Contraception Met Rate

 182
 Female Adult Obesity

 186
 Female Smoking Rate

 190
 Tuberculosis Death Ra

 194
 Male Unemployment R

 198
 Urban Population

 202
 Income share held by I

 206
 Income share held by I

 210
 Government Effectiven

 214
 Voice and Accountabili

 215
 Ilitercy Percent Female

222 Secondary School Exp

226 Education Fraction of G

230 Secondary per Pupil Pe

10 ArableLandFraction

Generate Models Analyze Models Test & Validate What If? Reports

All Variables Selected Data

🜻 DataModeler 9.3 - [Death Age Gap

3 AdultPopulation

15 ChildPopulation

19 CropsLandFraction

7 AnnualBirths

11 Area

Selected Variables

Predictive Analytics – response target is selected for modeling

Select Variables Explore Data

Variable Sets

Data Record Label

232 Death Age Gap

👩 Sele.

÷

Variables

1 1 Country

2 5 Airports

3 9 ArableLandArea

6 21 EconomicAid

7 25 ElectricityImports

4 13 Boundaryl enoth

5 17 ConstructionValueAdde

Launch Project

Advanced Search & Selection

Response

Target Response

Multiple response targets can be selected within a given project

	8	29	ExternalDebt
	9	33	FemaleInfantMortalityFraction
	10	37	FemalePopulation
1	11	41	GDPAtParity
	5	45	GovernmentExpenditures
		49	HighestElevation
		53	IndustrialValueAdded
		57	InternetUsers
		61	LifeExpectancy
		65	MaleChildPopulation
		69	MaleLiteracyFraction
		73	MedianAge
		77	MilitaryExpenditureFraction
		81	MilitaryFitPopulation
		85	NaturalGasExports
	23	89	OilConsumption
	24	93	OilReserves
	25	97	Population
	26	101	RoadLength
	27	105	TotalFertilityRate
	28	109	UNNumber
	29	117	LaborFraction+Agriculture
	35	.17	ExpenditureFraction+GovernmentConsumption
	1	121	ExpenditureFraction•ImportValue
ľ	32	125	Male Smoking Deaths
	33	129	Male Cigarette Smoker Percent
	34	133	Total Smoking Prevaence
	35	137	Firearm Death Rate
	36	141	Gun Deaths TBD
	37	145	Region
	38	149	NonReligiosity
	39	153	Beer (%)
	40	157	2015 projection
	41	161	Non-communicable Death Fraction
	42	165	Breastfeeding Percent
	43	169	Health expenditure per capita
	44	173	Hospital Bed Rate
	45	177	Female Labor Fraction
	46	181	Male Adult Obesity
	47	185	Female-Male Literacy Ratio
	48	189	Teen Pregnancy Rate
	49	193	Female Unemployment Rate
	50	197	Urban Population Percent
	51	201	GINI index
	52	205	Income share held by lowest 10%
	53	209	Control of Corruption
	54	213	Rule of Law
	55	217	Adult Illiterate Population
	56	221	Primary School Expense Fraction

57 225 Education Fraction of GDP

58 229 Primary per Pupil Percent of GDP

The GUI implicitly supports a workflow moving from data insights to variable selection to model development and selection to definition of trustable models and creation of deployable models

4 AgriculturalValueAdded

8 AnnualDeaths

12 BirthRateFraction

16 CoastlineLength

20 DeathRateFraction

	Tot Tomaio organosto Forcons	
	135 Girls Cigarette Use	136 Youth Cigarette Use
	139 Gun Suicides	140 Gun Deaths Accidental
(2014)	143 Rate	144 Count
	147 Year listed	148 Religiosity
	151 Recorded consumption	152 Unrecorded consumption
	155 Spirits (%)	156 Other (%)
	159 Communicable Disease Death Rate	160 Injury Death Percent
n	163 Contraceptive Prevalence	164 Diabetes Prevalence
ion	167 Female-Headed Households	168 GNI per capita
ate (GDP)	171 Health expenditure, private (total)	172 Health expenditure, total (GDP)
	175 Improved Sanitation Access	176 Improved Water Access
1	179 Physician Availability	180 Male Childhood Obesity
	183 Female Childhood Obesity	184 Childhood Obesity
	187 Male Smoking Rate	188 Tuberculosis Detection Rate
e	191 Tuberculosis Prevalence Rate	192 Turberculosis Treament Success Rate
ite	195 Total Unemployment Rate	196 Unmet Contraception Need
	199 Vitamin A Coverage	200 Wanted Fertility Rate
ourth 20%	203 Income share held by highest 10%	204 Income share held by highest 20%
owest 20%	207 Income share held by second 20%	208 Income share held by third 20%
ISS	211 Political Stability	212 Regulatory Quality
у	215 Female Illiterate Population	216 Male Illiterate Population
	219 Education Expenditure Percent	220 Compulsory Education Duration
ense Fraction	223 Tertiary School Expense Fraction	224 Preschool Expense Fraction
DP (I)	227 Public Education Expenditure Fraction	228 Total per Pupil Percent of GDP
rcent of GDP	231 Tertiany per Pupil Percent of GDP	232 Death Age Gap

### Analysis now supports loworder categoricals

								🧶 Data	Model	er 9.3 -	[Fisher	lris]			
Launch	Project	Select Variables	Explore Data	Gene	erate Models	A	nalyze M	odels	Te	st & Validate	What If?	Rep	oorts		
			0	▶ (	7				Datas	SummaryTable					X
Target Respo	nse		Species of iris 5	Col	Label	Туре	Uniformity	Class	Unique	Distribution Plot	Zero-Cross	Min	Mean	Median	Max
Variable Set		All	0	1	Sepal length in cm. Sepal width in cm.	123	100%	~ ~	29 20		. <b>₽</b>	4.4	5.8 3.0	5.7	7.7 4.4
Data File	Summary			3	Petal length in cm.	123	100%	had	37		. →	1.0	3.7	4.4	6.9
Exploring	g Methods				Potal width in cm	123	100%	A	21			0.1	1.2	13	25
Max Explor	ers	1 🗘		-		120	100%	~	21		. "		1.2	viselelee	- declates
🗸 DataSu	mmaryTable			5	Species of iris	ABC	100%	111	3		···	setosa	virginica	virginica	virginica
View D	ata to Explore	9								virg	ginica				
Import	ed Data Repor	rt													
Univari	atePlot														
DataDi	stributionPlot														
DataCo	mpletenessM	lap													
DataCo	mpletenessPl	lot													
Bivaria	tePlot		Nou	i ki	addad	ic	tha								
Correla	tionChart		INEV	/IY	auueu	15	uie								
Correla	tionMatrixPlo	t	ability	to	autom	at	ically								
					rt low_	orc	lor								
		als to r	nur	nerio	CS										

## Getting the Zen of the Data



#### •••

DataModeler 9.3 - [Death Age Gap]

Launch Project Select Variables Explore Data Generate Models Analyze Models Test & Validate What If? Reports







The default rule for model development is that a model must be supported by having at least 75% of the data records be completely numeric in the variables used by the model. Incomplete records are not considered in the quality assessment.

▶ 🏹											X
					Page	1 Page 2					
					DataSun	nmaryTable					
Col	Label	Туре	Uniformity	Class	Unique	Distribution Plot	Zero-Cross	Min	Mean	Median	Max
20	DeathRateFraction	123	94%	$\sim$	202	Jan	↔	0.0	0.0	0.0	0.0
21	EconomicAid	123	98%	$\sim$	213		. ↔	-2.4×10 <sup>10</sup>	3.6×10 <sup>7</sup>	8.7×10 <sup>7</sup>	2.2×10 <sup>10</sup>
33	FemaleInfantMortalityFraction	123	92%	$\sim$	214		↔	0.0	0.0	0.0	0.2
35	FemaleLiteracyFraction	123	86%	$\sim$	142		↔	0.1	0.8	0.9	1.0
42	GDPPerCapita	123	96%	$\sim$	231		↔	137.6	15952.0	5514.7	211500.0
43	GDPRealGrowth	123	91%	$\sim$	218	<b>h</b>	↔	-0.1	0.0	0.0	0.2
61	LifeExpectancy	123	96%	$\sim$	228		↔	45.6	71.5	74.0	84.4
62	LiteracyFraction	123	92%	h	144			0.2	0.9	0.9	1.0

Missing data is handled seamlessly for both modeling and data exploration and model development

				DataS	ummary	Table				
Col	Label	Туре	Uniformity	Class	Unique	Distribution Plot	Zero-Cross	Min	Mean	Median
1	cap-shape	ABC	100%	Ш	6	00	II	bell	convex	convex
2	cap-surface	ABC	100%	111	4		11	fibrous	scaly	
3	cap-color	ABC	100%	л	10	00000	II	brown	brow	
4	bruises?	ABC	100%	=	2		П	False	Fals	
5	odor	ABC	100%	л	9		11	almond	non	
6	gill-attachment	ABC	100%		2		П	attached	free	A
7	gill-spacing	ABC	100%	$\equiv$	2		П	close	clos	is t
8	gill-size	ABC	100%	$\equiv$	2		П	broad	broa	ро
9	gill-color	ABC	100%	л	12	.0000_00_0_	11	black	buff	Go
10	stalk-shape	ABC	100%	$\equiv$	2		П	enlarging	taperi	att
11	stalk-root	ABC	100%	Ш	5		II		bulbo	pre
12	stalk-surface-above-ring	ABC	100%	111	4		II	fibrous	smoo	
13	stalk-surface-below-ring	ABC	100%	111	4		II	fibrous	smoc	We
14	stalk-color-above-ring	ABC	100%	л	9	fibrous	II	brown	whit	ca
15	stalk-color-below-ring	ABC	100%	л	9	O_O_	II	brown	whit	tho
16	veil-type	ABC	100%	_	1		•	partial	parti	0+:
17	veil-color	ABC	100%		4		II	brown	whit	rec
18	ring-number	ABC	100%	_	3	_ 🛛 _	II	none	one	nre
19	ring-type	ABC	100%	Ш	5		II	evanescent	penda	
20	spore-print-color	ABC	100%	Ш	9	00_00_	II	black	white	
21	population	ABC	100%	л	6		II	abundant	several	several
22	habitat	ABC	100%	л	7	D	II	grasses	woods	woods
23	edibility of mushroom (either edible or poisonous)	ABC	100%	$\equiv$	2		П	edible	edible	edible

► Ÿ

### **Categorical Analysis**

Max

sunken

- A purely textual data set is the mushroom edible or poisonous?
- Goal is to identify key attributes & develop predictive models
- We cap the number of categories converted with those not converted ignored
- Still doing symbolic regression to produce a predictive model

solitary

woods

poisonous

### Model Search



- The assumption is that not all variables are equal and that we can consolidate to a relative handful of inputs
  - Models are evolved rewarding those which are simple and accurate with breeding rights
- Developed models are simple algebraic
   expressions human interpretable!
- We can view the search as an automated hypothesis generation & refinement
- Lots of diverse accurate-but-simple models are developed from which we can extract insight.





www.evolved-analytics.com

• • •										DataModeler	9.3 - [Mushroo	m Edibility]					
Launch Project	Select Var	iables	Explore Data	Ger	nerate Models	Analyz	e Models	Test	& Valid	ate What If? R	Reports						
Target Response	edibility of mushr	oom (eith	her edible or poisonous) 23	Þ	<b>⊽</b>	ec	libility•of•mus	shroom•eith	ner•edible	or•poisonous	X -1		Н	ere	we look at the		
Models Analyz	ze				capsurface -	a •	<b>_</b> ,	•				varia	able	pre	esence in each of th	e	
Select Models 307/3	2813		Reset Filters		bruises -		•				-	indep r	ena ours	ent ues	their own path to	acn	/
	Candidates Focu	s Ense	mbles		gillspacing -	•	-				1				success		
Quality Box			All		gillsize -					•	-						
			1-R <sup>2</sup> 0.2		stalkshape -		-				-				We		
Selection Fraction			50% 🗘	st	talksurfacebelowring			_			]		car	n als	so look for model		
Variables Required Variables		b N	Max		ringnumber -			•				sur explo	ostru oitec	JOTU 1. Th	ires which have been here metavariables	en can	
Allowed Variables		► A			ringtype -							provid	e in	sigh	it as well as be exp	loited	
Excluded Variables		► N	one		sporeprintcolor -					•	-		in s	ubs	equent rounds of		, ,
Power Limit Robust Models			4 V False			0	20	40	60	80 100					modeling		
ANOVA Trim			Apply														
Selected Model(s)			Save Model Set		► V	edibi	lity•of•mush	room•eith	er•edible	or•poisonous		► V		edibil	ity•of•mushroom	onous	
Model Ensemble			Create		$num \Rightarrow \%$	v	ariables Used		0.20	ParetoFrontPlot			Rank	Count	MetaVa riable	6 of models	% of MetaVariables
Explore Mode	lere we	loo	k at the				ille le e		0.15	il in the second		++	1	53	gillsize + √sporeprintcolor	17.3	35.1
Max Explor MO	ost preva	alen	t variable		1 64 ⇒ 20.8 %	y ri s	ngnumber poreprintcolor		0.10			++	3	28	$\sqrt{\text{gillsize}} + \sqrt{\text{sporeprintcolor}}$	9.1	18.5
Model Di. ♥ Variable Pres.	mbinatio perfo	ons rma	and their ince			<b>-</b>	☑ +		0.05			++	4	22	gillsize <sup>2</sup> stalksurfacebelowring	7.2	14.6
VariablePresence	Table								0.00	40 50 60 70	80 90 100	++	5	22	$\left( \text{gillsize} + \text{ringnumber} + \sqrt{\text{sporeprintcolor}} \right)^{1/3}$	7.2	14.6
VariablePresence	DistributionChart								0.20	12.			6	21	gillsize ringtype	6.8	13.9
VariablePresence	Map					9	illsize talkshape		0.15			++	7	14	1 -1.03928+sporeprintcolor	4.6	9.3
▼ Meta Variables (1)					<b>2</b> 28 ⇒ 9.1 %	ri	ngnumber		0.10			++	8	14	gillsize <sup>2</sup> + √sporeprintcolor	4.6	9.3
ModelBasisSetTa	ible*						⊠ +		0.05				9	13	$(\text{gillsize}^2 + \sqrt{\text{sporeprintcolor}})^{1/3}$	4.2	8.6
ModelBasisSetDis	stributionTable*								0.00						(ginaizo + V aporoprintocolor )	7.2	0.0
ModelBasisSetDis	stributionChart								0.20	40 50 60 70	80 90 100	++	10	12	$\sqrt{\text{sporeprintcolor}^3}$	3.9	7.9
🗸 MetaVariableTabl	le*					g	illsize		0.15								
MetaVariableDist	ributionTable*					r	ngnumber										
MetaVariableDist	ributionChart				3 21 ⇒ 6.8 %	s	poreprintcolor		0.10						_		
variable Combination	s (1) tionTable*					<b>-</b>	☑ +		0.05	2.53 × 10 <sup>-3</sup> – 1.35 ringnu	mber + 0.37 gillsize	e ringtype + 0.95	√ spore	eprintco	lor		
VariableCombinat	tionMap								0.00	40 50 60 70	80 90 100						
											50 50 100						
www.evolved-an	ialytics.com																







0.5

0.0

0.2

Each graphic has controls which can be exposed and used to tweak the display

- Here we are looking at one of the candidate models
- Unfortunately, THE model does not typically exist in data-derived models
- We can use ensembles of models to create a trustable model to detect extrapolation or underlying system changes

edibility+of+mushroom+either+edible+or+poisonous Residual Plot



DataModeler 9.3 - [Mushroom Edibility]

Tooltips are used to show the variable mapping being used

-1.0-0.50.0 0.5 1.0 1.5

ringtype

ringtype

evanescent → -1 flaring → 0

> large  $\rightarrow$  1 pendant  $\rightarrow$  2

## Test & Validate



www.evolved-analytics.com





www.evolved-analytics.com



Instead of trying to find THE model, we can exploit the abundance of explored model forms and select a diverse set from the good-andsimple category to form a TRUSTABLE model



Expression

Ensemble

Quality

Model Quality

Model	Complexity	1-R <sup>2</sup>
1	26	0.0229559
2	29	0.0212007
3	29	0.029529
4	33	0.0137302
5	34	0.0182884
6	34	0.0235753
7	38	0.0166514
8	38	0.02229
9	38	0.025075
10	41	0.0109854
11	41	0.0220358

### **Ensemble Performance**

The (algorithmically generated) ensemble gives pretty good composite results when asked
 DataModeler 9.3 - to extrapolate







Round 3 Ensemble (3 Vars) [DistTower\_training] [ ref = 0.66 ]







Since the variables in this case are coupled, changing one without a corresponding change of the others would venture into unknown regions of parameter space. This is detected by the ensemble since the constituent models (light gray lines) diverge if we try to change, for example, the tray temperature independently from the observed reference point (green dot)



Round 3 Ensemble (3 Vars) [DistTower\_test] [ ref = -0.11 ]

- Here we are looking at the model prediction in the lowest observed value in the test data set.
- The constituent models diverge when away from known data records which is the desired behavior — even though the actual prediction degrades reasonably gracefully.
- Graceful degradation for extrapolation is a unique benefit of symbolic regression since other datadriven modeling techniques tend to fail spectacularly when asked to extrapolate
- Since these data points were explicitly chosen to test the ensemble ability to detect extrapolation, we can conclude that the ensemble is, in fact, a trustable model



## Spiffy Feature





### AI ⇒ Augmented Intelligence

- DataModeler easily handles exploration and analysis of multivariate data with correlated inputs, missing values and (now) low-order categoricals
- Searching for simple & accurate models allows identifying driving variables as well as variable combinations useful for developing quality models for deployment
- Algebraic models allows for analyst inspection and insight development
- Ensembles of simple & accurate models can be used as trustable models and obviate the need for data set partitioning — which is important when data is sparse
- Ease-of-use is enhanced with graphical tools supporting an effective analysis workflow.

# For More Information

- www.evolved-analytics.com
- mark@evolved-analytics.com

- GUI Development Acknowledgement:
  - Ariel Sepúlveda
  - ariel.sepulveda@prontoanalytics.com
  - <u>www.prontoanalytics.com</u>

