## **REAL-WORLD DATA MODELING**

MARK KOTANCHEK

**EVOLVED ANALYTICS LLC** 

WWW.EVOLVED-ANALYTICS.COM

MONEY

Specifically, we want to convert data into money

GECCO'12 COMPANION, JULY 7-11, 2012, PHILADELPHIA, PA, USA. ACM 978-1-4503-1178-6/12/07.

## WHY ARE WE HERE?





#### DEALING WITH THE DATA DELUGE

Conventional technologies struggle ...

- too much data
- too many variables
- ugly data
- too little time
- the simplifying assumptions are not valid
- models are big piles of "trust-me"

Multi-objective symbolic regression ...

- can help to make sense ... and money ... out of the data deluge.
- trustable models are possible
- cheaper-better-faster really IS possible!

### DATA, DATA EVERYWHERE



### THE ESSENCE IS NUMBERS (AND SOME CONTEXT)

symbol	name	BookValu	FloatSha	rForwardP	MarketCa	a PEGRatic P	PERatio	PriceTarg	PriceToBo	PriceToSa	ShortRatin	YearPER:	EarningsF	EBITDA	ForwardE	QuarterFc	YearEami	Dividend'	DividendF	Dividend
A	Agilent Te	7.311	3.4E+08	9.23	5.7E+09	0.85	8.74	26.86	2.23	0.99	1.2	11.11	1.868	1.0E+09	1.77	0.33	1.47		0	2.057
AA	Alcoa Inc	18.755	7.9E+08	16.69	7.8E+09	0.35	4.9	17.27	0.55	0.28	1.9	6.68	2.113	4.0E+09	0.62	-0.15	1.55	0.0657	0.68	0.17
AAPL	Apple Inc	23.674	8.8E+08	13.82	7.9E+10	0.87	16.64	142.48	3.77	2.44	0.4	17.08	5.358	6.7E+09	6.45	1.17	5.22		0	
ABC	Amerisou	17.348	1.5E+08	9.76	5.3E+09	0.94	22.09	39.04	1.96	0.08	1.3	10.85	1.542	9.2E+08	3.49	0.89	3.14	0.0095	0.325	0.1
ABT	Abbott La	12.468	1.5E+09	14.19	8.1E+10	1.31	17.87	61	4.18	2.8	1.2	15.68	2.914	7.6E+09	3.67	0.76	3.32	0.027	1.405	0.36
ACAS	American	24.426	2.0E+08	1.29	6.4E+08	0.19		11.65	0.12	0.53	5.7	1.07	-8.386		2.37	0.62	2.85	1.0131	3.09	1.05
ACS	Affiliated (	24.304	8.8E+07	9.81	4.2E+09	0.84	11.77	51.44	1.71	0.65	1.6	11.37	3.526	1.1E+09	4.23	0.93	3.65		0	
ADBE	Adobe Sy	8.225	5.3E+08	10.74	1.1E+10	0.75	14.36	26.42	2.64	3.22	0.9	11.98	1.51	1.3E+09	2.02	0.45	1.81		0	0.0065
ADI	Analog De	8.312	2.8E+08	15.14	5.3E+09	1.02	7.44	20.82	2.37	2.22	0.9	18.74	2.646	7.8E+08	1.3	0.24	1.05	0.0396	0.78	0.2
ADM	Archer Da	21.097	6.2E+08	9.84	1.8E+10	0.68	7.59	29.36	1.34	0.23		8.24	3.735	4.1E+09	2.88	0.64	3.44	0.0183	0.52	0.13
ADP	Automatic	9.551	5.0E+08	14.27	1.9E+10	1.17	15.99	42.68	3.9	2.11	2.2	15.65	2.33	2.1E+09	2.61	0.81	2.38	0.0322	1.2	0.33
ADSK	Autodesk	6.187	2.2E+08	11.93	4.2E+09	0.86	11.64	18.36	3.1	1.79	2	9.85	1.651	6.0E+08	1.61	0.37	1.95		0	0.015
AEE	Ameren C	33.395	2.1E+08	9.77	6.9E+09	2.81	10.4	32.06	0.98	0.91	1.3	11.17	3.136	2.1E+09	3.34	0.63	2.92	0.0779	2.54	0.635
AEP	American	27.089	4.0E+08	9.26	1.3E+10	1.82	8.34	38.11	1.11	0.84		9.5	3.623	4.2E+09	3.26	0.94	3.18	0.0543	1.64	0.41
AES	The Aes (	6.387	6.6E+08	6.49	5.3E+09	0.47	3.98	14.6	1.18	0.31	1.1	6.91	1.891	4.4E+09	1.16	0.28	1.09		0	
AET	Aetna Inc	20.162	4.6E+08	5.99	1.2E+10	0.44	7.46	40.09	1.21	0.37	0.7	6.21	3.274	3.1E+09	4.08	0.95	3.93	0.0016	0.04	0.04
AFL	Aflac Inco	13.64	4.5E+08	10.14	2.1E+10	0.78	15.66	55.89	3.41	1.33	3.2	11.6	2.971	2.3E+09	4.59	1.11	4.01	0.0206	0.96	0.24
AGN	Allergan li	13.147	3.0E+08	13.51	1.1E+10	1	19.3	48.25	2.81	2.53	1.6	14.47	1.911	1.2E+09	2.73	0.61	2.55	0.0041	0.15	0.05
AIG	American	26.463	2.4E+09	1.35	4.5E+09			5.96	0.07	0.09			-16.523	-4E+10	1.3	0.03	-3.62	0.3543	0.62	0.22
AIV	Apartmen	9.843	8.2E+07	5.11	1.1E+09	1.03	3.7	13.87	1.43	0.68	3	4.59	3.803	9.0E+08	2.76	0.7	3.07	0.5614	7.91	1.8
AIZ	Assurant	31.797	1.2E+08	3.86	3.1E+09	0.39	7.75	47.19	0.79	0.35	0.3	4.43	3.238	6.5E+08	6.51	1.65	5.67	0.0215	0.54	0.14
AKAM	Akamai Te	8.904	1.6E+08	9.28	2.6E+09	0.59	21.1	20.64	1.8	3.57	2.1	9.98	0.761	3.3E+08	1.73	0.41	1.61		0	
AKS	Ak Steel F	12.787	1.1E+08	4.52	1.1E+09	0.27	2.3	23.79	0.87	0.16	1	2.7	4.815	1.1E+09	2.45	0.4	4.1	0.0181	0.2	0.05
ALL	Allstate C	31.601	5.3E+08	5.57	1.6E+10	1.19	78.86	39.73	0.94	0.5	1.1	8.3	0.378		5.35	1.45	3.59	0.055	1.64	0.41
ALTR	Altera Cor	2.902	2.9E+08	18.15	4.8E+09	0.92	15.31	19.15	5.75	3.61	2.7	14.27	1.091	4.3E+08	0.92	0.22	1.17	0.0114	0.19	0.05
AMAT	Applied M	5.673	1.3E+09	14.42	1.3E+10	3.29	15.26	13.09	1.88	1.74	1.1	39.52	0.699	1.7E+09	0.74	0.04	0.27	0.0225	0.24	0.06
AMD	Advanced	2.202	4.7E+08		1.4E+09			3.86	1.04	0.21	4.5		-5.743	8.5E+08	-1.76	-0.58	-1.55		0	
AMGN	Amgen In	18,727	1.1E+09	12.48	6.1E+10	1.28	15.54	71.39	3.12	4.13	1.9	12.84	3.76	6.7E+09	4.68	1.16	4.55		0	
AMP	Ameripris	31.005	2.2E+08	6.61	4.6E+09	0.43	8.56	35.25	0.71	0.58	0.8	6.9	2.563	1.1E+09	3.32	0.73	3.18	0.0292	0.64	0.17
AMT	American	7.245	3.9E+08	41.57	1.1E+10	4.89	47.33	42.07	3.96	7.28	1.6	56.24	0.606	1.0E+09	0.69	0.16	0.51		0	
AMZN	Amazon (	5.89	3.2E+08	35.93	2.2E+10	1.62	36.47	55.45	9.03	1.26	2.5	38.54	1.458	1.1E+09	1.48	0.31	1.38		0	
AN	Autonatio	12.028	1.6E+08	9.89	1.7E+09	0.95		7.21	0.81	0.11	8.1	9.06	-7.021	6.4E+08	0.98	0.25	1.07		0	0.68
ANF	Abercrom	20.686	8.0E+07	10.67	2.0E+09	0.58	4.85	24	1.1	0.52	1.5	6.84	4.684	8.6E+08	2.13	0.22	3.32	0.0308	0.7	0.175
AOC	Aon Corp	22.044	2.5E+08	12.11	1.1E+10	1.36	7.84	51.9	1.92	1.46	0.8	14.87	5,404	1.3E+09	3.5	0.99	2.85	0.0142	0.6	0.15
APA	Apache C	55.77	3.3E+08	10.57	2.3E+10	0.63	5.28	99.35	1.33	1.84	1	5.99	14.023	1.0E+10	7	1.46	12.36	0.0095	0.7	0.15
APC	Anadarko	39 672	4 6E+08	21.55	1 7E+10	1.21	6.74	58.2	0.99	1.24	13	7.52	5 822	8.5E+09	1.82	0.34	5.22	0.0092	0.36	0.09

www.evolved-analytics.com

### WE WANT TO CONVERT NUMBERS INTO ...

IDENTIFY OUTLIERS VARIABLE SELECTION

PREDICTION

GUIDANCE

## MONEY

VARIABLE RELATIONSHIPS

INSIGHT & UNDERSTANDING

EMULATORS & OPTIMIZATION

EARLY WARNING & RISK MANAGEMENT



#### VARIABLE SELECTION

Simply knowing which of the inputs REALLY matter can be important and valuable

www.evolved-analytics.con



### VARIABLE RELATIONSHIPS

Identifying key variable combinations (metavariables) can lead to insights and better models



#### PREDICTION

Inferring from observables the current or future quality, price, temperature, demand, etc. can be valuable.

www.evolved-analytics.com



x1, x2, x3, x5, x6 x4 - reflux flow y - process quality variable x4. y x1 and x2 – pressures x3 and x5 – temperatures x6 – feed flow x8, x9, x10 x11, x12 x8 - x12 and x14 - x19 - temperatures x13 x13 - material balance x14, x15, x16 x17, x18, x19 x7 and x22 - flows x23 - bottom temperatur x20, x21 x7, x22, x23 x20 and x21 - flows

#### **EMULATORS & OPTIMIZATION**

A surrogate for a real system can be useful for coarse optimization or training or exploring what-if scenarios

www.evolved-analytics.com

nugget of information. The first step in that

assessment is to identify them



#### **RISK MANAGEMENT**

Knowing when the the system has changed or the model SHOULDN'T be trusted can be very valuable.

www.evolved-analytics.com



Where should we look for better solutions? How can we drive uncertainty out of the models? What are the key factors?

> DATA IS NEVER TO BE TRUSTED

www.evolved-analytics.com



#### INSIGHT & UNDERSTANDING

'Eureka!' (I found it!) but 'That's funny ...'

### THE HUMAN FACTOR

IF YOU DON'T KNOW WHERE YOU ARE GOING ANY PATH WILL GET YOU THERE

- Modeling Objective
- Data Characteristics
- Context & a priori Knowledge

**CONTEXT-FREE ANALYSIS** LEADS TO CONFIDENTLY WRONG CONCLUSIONS

www.evolved-analytics.com

— Isaac Asimov (1920 - 1992) www.evolved-analytics.com

### WHY MODELING?

### WHY MODELING?





### WHAT IS SYMBOLIC REGRESSION?

- Must discover BOTH the model structure AND the embedded coefficients
- This is a MUCH more difficult problem than conventional regression or neural network modeling — the search space is infinite & an infinite number of models will fit the data
- The Basic Problem:
  - How do we efficiently search for models

www.evolved-analytics.com

• How do we know when we are done?

### Agenda

- Modeling & Motivations
- Conventional GP Overview & Limitations
- The ParetoGP (Multi-Objective) Perspective
- Ensembles & Trustable Models
- Trustability & Active DOE
- Outlier Identification
- Case Studies

### **TUTORIAL GOALS**

- Understanding of the **modeling context** & motivation
- Understanding of GP-based symbolic regression
  - vanilla GP
  - Pareto GP
  - other algorithmic variants
  - common issues, problems, good practices and rules-of-thumb
- Foundations and implications of **ensembles & trustable models**
- Make it tangible
  - case studies using industrial data
- Awareness that GP-based symbolic regression (ParetoGP) is something truly special





### WHAT IS TRUTH?

Each of these models fits the data EXACTLY

We need more information to make an intelligent model selection



### **CONFOUNDING FACTORS**

With many variables and limited data, determining which inputs are truly drivers and which correlations are accidental is difficult



ww.evolved-analytics

### **CONFOUNDING FACTORS**

### **TRUST & TREPIDATION**

Again, all models fit the data EXACTLY

How do we determine which variables are real?





- Using data-driven models is like driving fast with the windshield painted over and using only the rear-view mirrors
- This might work as long as the model is perfectly accurate and there aren't any curves
- Are you feeling lucky?

www.evolved-analytics.com

### KEY POINT

All of the previous response surfaces were generated via symbolic regression.

The only constraint is the supplied building blocks.

### SYMBOLIC REGRESSION

Human limits of imagination & possibility are not imposed!

### HYPOTHESIS GENERATOR

We can exploit this creativity to produce trustable data models.

### WHAT ARE GOOD MODELS?



www.evolved-analytics.com



### SYMBOLIC REGRESSION VIA GP



#### SUMMARY

THERE ARE MANY POSSIBLE VARIANTS OF SYMBOLIC REGRESSION. THESE CAN EITHER MAKE THE MODEL SEARCH MORE EFFICIENT OR THEY CAN MAKE THE SEARCH INEFFECTIVE AND SLOW.

### CONVENTIONAL GENETIC PROGRAMMING PROBLEMS

- RELATIVELY SLOW DISCOVERY
  - Computational demands are intense
- SELECTION OF "QUALITY" SOLUTIONS
  - Trade-off of Complexity vs. Performance
- GOOD-BUT-NOT-GREAT SOLUTIONS
  - Other nonlinear techniques outperform in raw performance
- BLOAT
  - The "best" model explodes in complexity
- STAGNATION
  - The search cannot innovate out of local optima.

www.evolved-analytics.com

### REPRESENTATIONS

- There are many GP variants which use different underlying genome structures
  - Tree-Based (the Koza original)
  - GE Grammatical Evolution
  - Cartesian GP graph based
  - GEP Gene Expression Programming (closer to a GA representation)
  - LinearGP (evolve machine code directly)
  - Stack-based (e.g., PUSH)
- Functionally, they are pretty similar despite the passion of some of their proponents
- HOWEVER, implementation details can produce orders-ofmagnitude differences in performance & model quality

### THE WAY WE USED TO DO DATA MODELING

- 1. Use "stacked analytic networks" to identify the driving variables
- 2. Use support vector regression (SVR) to identify the key data records
- 3. Apply GP (slowly) to the reduced data set
- 4. Painfully search through the developed models for simple and high-quality expressions or to extract insight.
- 5. Think that there HAD to be a better way

# PARETO • Accurate

- Circula
- Simple
- No spurious variables
- Robust
- Limited nonlinearity
- Dimensional consistency
- Smoothness
- etc.

## Accuracy, of course, is dominant

If we focus on model simplicity, then most of the other goals will be achieved as a side-effect

However, we have and do use more than just simplicity as the additional goal

### MODEL COMPLEXITY

(MULTI-OBJECTIVE)

**3+ ORDERS-OF-MAGNITUDE SPEED IMPROVEMENT** 

**OPENS LOTS OF POSSIBILITIES!** 



- What is complexity?
  - # of nodes?
  - Tree depth?
  - Included functions?
  - Number of variables?
  - Combinations?
- Chosen function is sum of sum of node counts of genome
- Provides more resolution at low end of complexity than simply using node count
- Rewards flatter structures with fewer layers
- Maarten Keijzer has renamed this
   "visitation length"
   <u>www.evolved-analytics.com</u>



WHAT DO WE REALLY

WANT?



Models at the "knee" of the Pareto front are generally preferred since they have the best balance; however, we don't know the complexity vs. accuracy trade-off until AFTER we have built the models.

- Identifies trade-off surface between competing objectives
  - e.g., accuracy vs. complexity
- Pareto front solutions are the best "bang-for-the-buck"
- Unwarranted complexity is punished automatically
  - spurious variables
  - introns
  - unnecessary terms

### THE PARETO FRONT



How do we EXPLOIT THE PARETO FRONT?

- Identifies trade-off surface between competing objectives
- e.g., accuracy vs. complexity
- Pareto front solutions are the best "bang-for-the-buck"
- Unwarranted complexity is punished automatically
- spurious variables
- introns
- unnecessary terms

### PARETO PERFORMANCE



- Multiple ways to characterize Pareto Performance
- Computational Issues
  - Brute force is MN<sup>2</sup>
  - Can do M N log<sub>M-1</sub>(N) or M N log<sub>M-2</sub>(N) if clever
    - M = # of objectives
    - N = population size
- Global comparison is expensive & doesn't scale well
  - Pareto tournaments are an alternative
    - www.evolved-analytics.com

www.evolved-analytics.con

### **PARETO TOURNAMENTS**



- Pareto Tournaments
- 1. select a random pool of N models to compete
- 2. models on the multi-objective Pareto front getting breeding rights
- 3. repeat until sufficient parents are selected

- Tournaments are the standard in EC
  - simple
  - robust
  - tunable selectivity
  - scalable NO global awareness
- Pareto tournaments have the same properties

### PARETO TOURNAMENTS

.



- 1. select a random pool of N models to compete
- models on the multi-objective Pareto front getting breeding rights
   repeat until sufficient parents are selected

- Identifying Pareto front of a local tourney is easier & faster than establishing the global pecking order
- The selection focus is automatically the knee of the Pareto front
- Note that this is NOT the commonly used approach of pairwise comparison which is NOT tunable

### CONVENTIONAL GP EVOLUTION

## PARETO TOURNAMENT



- Inferential sensor data set
- Single objective tournaments based upon accuracy
- Upper bound on complexity
- Best model survives across generations
- "Best" models are WAY too complex www.evolved-analytics.com





- Same data set
- Minor changes:
  - Pareto Tourney (30)
  - Pareto front survives across generations
- 50% more model search in same time
- Continual
   improvement
- MUCH better models
   <u>www.evolved-analytics.com</u>

### THE PARETOGP ESSENCE

## BUT THAT'S NOT

ALL...

BEHIND DOOR #2, WE HAVE ...



- Multi-Objective Selection
- Multi-Objective Archive
- Niching
- Continual New
   Genetics Influx

### **GOOD IDEAS**

- ARCHIVES (METHUSELAHS)
  - preserve good models

#### **MULTIPLE OBJECTIVES**

- focus on the real goal
- primary/secondary objectives
- alternating objectives

#### NICHING

model protection during
 development

#### SCALE & TRANSLATION INVARIANCE

- e.g., absolute correlation
- calculating scale & translation is straightforward

#### ELIMINATE REDUNDANT MODELS

• clean out duplicates

#### RESCALE

- mapping variables to a common range can make them more interchangeable
- this will create bigger expressions when returning to the natural scaling

#### LAMARCKIAN EVOLUTION

• simplify & optimize models during evolution

#### INTERVAL ARITHMETIC

• focus on models which can be deployed safely

#### **ORDINALGP & ESSENCE**

- dynamic environment
- data balancing

### MULTIPLE OBJECTIVES



### INTERVAL ARITHMETIC & ROBUSTNESS

### NICHING



- Niching is a foundation of diversity and continuous innovation
- There are many ways to implement niching
  - Multiple Objectives (ParetoGP)
  - HFC Hierarchical Fair Competition
  - ALPS Age-Layered Population Structure
  - Simple Geographies
  - Islands
  - etc.

.

• The key is to let models develop in a protected way

rangeMap = {
 x1 → Interval[{1, 4}],
 x2 → Interval[{2, 10}]
 };
  $\frac{x1}{x2}$  /. rangeMap
Interval[{ $\frac{1}{10}$ , 2}]

```
x1^2 - x1 /. rangeMap
```

Interval[{-3, 15}]

- Rapid test for model stability
  - two-orders of magnitude faster than nonlinear search
  - assumes hypercube of parameter ranges
- Test is conservative rejects well behaved models
  - restricts model structure
  - can be implemented during model evolution for rapid pathology rejection.



### FAT DATA SETS



### LARGE DATA SET STRATEGIES

### CORRELATED DATA

Uniformly distributed data records are not necessarily of equal value



#### BalanceData prioritizes the data records



- More CPU Time
- Balanced Data
   Subset
- OrdinalGP (randomly changing subsets)
- ESSENCE (incrementally adding balanced data)



15 Seconds of Modeling Complexity 1–R<sup>2</sup> Vars Function  $\frac{0.334}{10.633}$  + 0.633 1 11.000 0.564 x<sub>7</sub> X7 0.620 x5 2 19.000 0.167  $\frac{x_2}{x_5}$  -0.482 + x<sub>2</sub>  $\frac{0.106 x_5^2}{x_2^2} + 0.294$  $3 \quad 27.000 \quad 0.089 \quad \begin{array}{c} x_2 \\ x_5 \end{array}$  $X_3$ 1.000 x5 4 29.000 0.000 x<sub>5</sub> x3 x6 x

www.evolved-analytics.com

NOISY DATA





NOISY DATA



### **TRUSTABLE MODELS**







### EXTRAPOLATION



### EXTRAPOLATION



www.evolved-analytics.com

1366



### How do we choose THE model?



- Many models are developed during the evolutionary search
- These models have different structures despite having comparable performance
- Different evolutions will discover different model forms

### How do we choose THE model?

# WE DON'T!

www.evolved-analytics.cor

### How do we choose THE model?

• Instead, we exploit the DIVERSITY of structural forms via model ensembles

#### Recall that the models are constrained by the data and NOT by physics or preconceived assumptions!

### WHAT IS AN ENSEMBLE?

- A collection of models which are ...
  - diverse and
  - "good enough" complexity and accuracy

### OTHER THINGS TO CONSIDER

- model forms for INSIGHT ...
- variable presence for DRIVERS
- variable combinations (metavariables)
- Pareto front shape (modeling potential & difficulty)

### WHY IS AN ENSEMBLE SPECIAL?

www.evolved-analytics.com

www.evolved-analytics.com

- The constituent models ...
  - **AGREE** where constrained by the observed data (otherwise, they would not be good models)
  - **DIVERGE** where they are not constrained by experience (otherwise, they would not be diverse)
- Thus, we have a high-quality prediction **AND** a trust metric to accompany that prediction!

www.evolved-analytics.com

### **TRUSTABLE MODELS!!**

- No more "Trust me"
  - This is a classic issue for use of data-derived models
- This is VERY unique and made possible by the diverse structures coming out of ParetoGP
- Knowing when NOT to trust the prediction can be very valuable
- We can also exploit ensembles for adaptive data collection
- Identified systems also tend to extrapolate well

### ENSEMBLES ARE/ARE NOT ...

ARE NOT ...

- collections of weaklearners
- diverse due to using different training data
- collections of locally-good models

ARE ....

- extrapolate reasonably
- white-box
- robust predictors (use median average for prediction)
- diverse due to uncorrelated prediction residuals

www.evolved-analytics.com

### **DIVERSE MODELS**

### **Diversity Aspects**

- model structure
- constituent variables
- (lack of) error correlation observed data
- prediction disagreement
  - synthetic data

#### A Definition Algorithm

- build covariance matrix
- select most uncorrelated pair
- delete all models not uncorrelated (within threshold) to pair
- repeat until no models left
- Divide-and-conquer for large numbers of models to avoid scaling problems

ALTERNATE: choose a reference and work from there





The definition of "uncorrelated" needs to be relaxed from the conventional statistics criteria.

- A strategy that seems to work is:
  - uncorrelated models from model population
  - uncorrelated models from Pareto front
  - "most typical" model from Pareto front

www.evolved-analytics.com

www.evolved-analytics.con



### EVALUATING THE ENSEMBLE

• Response:

- median, median average or mean
- Consensus:
  - extrema range, inter-quantile range or standard deviation



www.evolved-analytics.com

### SUMMARY: **TRUSTABLE MODELS!**

- Ensembles of diverse models address a fundamental problem with empirical models
- We have outlined a straight-forward approach to assembling those models
- We can now use models with an awareness of the prediction risk
- The windshield has been cleaned!

### FOUNDATIONS SUMMARY

#### ParetoGP

- Appropriate complexity models .
- Automatic variable selection
- Able to handle "fat" data sets
- continual innovation .

Predictions vs. Test & Training Dat

#### **Ensembles**

- More aware models (no need) for data partitioning during model development?)
- Trustable predictions
- Interesting models for expert insight — "hypothesis selection"
- Ability to implement Active DOE

### Summary

- \* A smart ParetoGP implementation is >1,000x faster than vanilla GP and produces better models
- That speed improvement CHANGES what is possible and redefines best practices!

www.evolved-analytics.con

### THE CURSE OF DIMENSIONALITY

**HOW TO GET THERE?** 

## ACTIVE DESIGN-OF-EXPERIMENTS

ACCELERATING LEARNING VIA FOCUSED DATA COLLECTION



- Adding more variables makes the problem difficulty explode
- It is often difficult or impossible to initially collect enough of the right data

www.evolved-analytics.com

### WHAT TO DO???



Target data collection to:

- identify the true driving variables
- reject spurious models
- increase model fidelity (global accuracy)

o: Adaptive Data Collection



- ParetoGP lets us:
  - explore diverse model structures and
  - different variable combinations
- Diverse model ensembles let us:
  - identify regions of lower model fidelity (the models are presumably accurate in the sense that they fit the observed data)

Guiding Principle: We want a diverse collection of accurate but simple models

### MOTIVATIONS & SCENARIOS: OPERATIONS OPTIMIZATION

- Operating plants "walk away" from their optimal set-points
- Many available variables with unknown interactions
- Most historical data is essentially identical
- New data is precious & risky

Data collection is severely constrained

www.evolved-analytics.com

### MOTIVATIONS & SCENARIOS: DESIGN & EMULATION



- The REAL system takes a long time to generate data
- Want an emulator for coarse optimization & human insight
- Want better products and faster time-to-market



- LOTS of options being explored
- Want to quickly identify key factors & relationships
- Want to search for global rather than local optimum

Many variables of unknown significance

**MOTIVATIONS & SCENARIOS:** 

**HIGH-THROUGHPUT RESEARCH** 

www.evolved-analytics.cor

### MOTIVATIONS & SCENARIOS: BIOREACTOR STARTUP

- Nonlinear dynamics
- Lots of control settings
- Full-scale reactors VERY different from lab-scale used for development
- New yeast strain likely behaves differently than previous

www.evolved-analytic

Time is Money

Limited ability to collect data

### MOTIVATIONS & SCENARIOS: MAXIMIZING MODEL VALIDITY



- Knowing that the model is not a good predictor is not enough
- We want the model to be accurate over the entire operating range
- Example: modeling price-volume elasticity

Need to maximize the model fidelity & range

ADAPTIVE DATA COLLECTION SEQUENCE



- 1. build models from the available data
- 2. assemble ensembles from those models
- 3. identify the locations of model uncertainty and potential optima
- 4. collect new data

5. repeat

### MOTIVATIONS & SCENARIOS: COMMON THEMES

- Data is precious
- Many variables of unknown significance
- Data collection is difficult and/or expensive
- BOTH variable selection AND a good model is desired
- Delivery time pressure (Aren't you done YET?)

### WHAT WE HAVEN'T DONE

- assumed which variables are important
- harshly constrained the number of variables
- assumed variables are independent & uncorrelated
- assumed a model form
- collected data blindly

www.evolved-analytics.com 1373



### A 1-D ADAPTIVE DOE

 Using only the divergence in the constituent models for guidance, the adaptive DOE process was able to collect data to produce a high-quality predictor



Note that the ensemble models agree where there is data and diverge where not constrained by data www.evolved-analytics.com

### MAKING THE PROBLEM HARDER

- We have 10 dimensions (variables)
- Only two variables are real
- We start with 12 data points
- After each round of model building, we collect data at the predicted maximum, predicted minimum and the point of maximum ensemble uncertainty

### 20 ROUNDS OF ACTIVE DOE



- Only the data projection into the true variables is shown
- All data inputs are 10-D
- Note that the early stages feature high accuracy but low fidelity
- Response surface only shown for variables which are in >50% of models in the "interesting region"

www.evolved-analytics.con

www.evolved-analytics.com

 $e^{-(x_2-1)^2}$ 

 $(-2.5)^2 + 1.2$ 



### SELECTED ROUNDS

- In the early stages the new data tends to shatter many candidate
- In later stages the incremental is less candidate models
- Spurious variables may appear dominant in the early going but additional data leads to their rejection

### WHAT HAVE WE SEEN?

- Exploiting the DIVERSITY of candidate models generated by ParetoGP,
- we can use a model ensemble prediction divergence to identify regions of UNCERTAINTY
- and target data collection at reducing that uncertainty.
- This results in HIGH-FIDELITY and TRUSTABLE models and the REJECTION of spurious inputs.

## CASE STUDIES

THE PROOF IS IN THE PUDDING

	CountryIndex	AdultPopulation	Airports	AMRadioStations		
	AnnualBirths	AnnualDeaths	ArableLandArea	ArableLandFraction		
	Area	BirthRateFraction	BoundaryLength	CallingCode		
	CellularPhones	ChildPopulation	CoastlineLength	CropsLandArea		
	CropsLandFraction	DeathRateFraction	EconomicAid	ElderlyPopulation		
	ElectricityConsumption	ElectricityExports	ElectricityImports	ElectricityProduction		
	ExchangeRate	ExportValue	ExternalDebt	Female Adult Population		
	FemaleChildPopulation	FemaleElderlyPopulation	FemaleInfantMortalityFraction	FemaleLifeExpectancy		
	FemaleLiteracyFraction	FemaleMedianAge	FemalePopulation	FixedInvestment		
	FMRadioStations	GovernmentConsumption	GovernmentExpenditures	GovernmentReceipts		
	GovernmentSurplus	Grossinvestment	HighestElevation	HouseholdConsumption		
	ImportValue	InfantMortalityFraction	InflationRate	InternetHosts		
	InternetUsers	InventoryChange	IrrigatedLandArea	IrrigatedLandFraction		
	LaborForce	LandArea	LifeExpectancy	LiteracyFraction		
	LowestElevation	MaleAdultPopulation	MaleChildPopulation	MaleElderlyPopulation		
	MaleInfantMortalityFraction	MaleLifeExpectancy	MaleLiteracyFraction	MaleMedianAge		
	MalePopulation	MedianAge	MigrationRateFraction	MilitaryAgeMales		
	MilitaryExpenditureFraction	MilitaryFitMales	NaturalGasConsumption	NaturalGasExports		
	NaturalGasImports	NaturalGasProduction	NaturalGasReserves	OliConsumption		
	OilExports	OilImports	OilProduction	PavedAirports		
	PavedRoadLength	PhoneLines	Population	PopulationGrowth		
	PriceIndex	RadioStations	RoadLength	ShortWaveRadioStations		
	TelevisionStations	TotalConsumption	TotalFertilityRate	UnemploymentFraction		
1	UNNumber	WaterArea	ExpenditureFractions- ExportValue	ExpenditureFractions- FixedInvestment		
G	ExpenditureFractions- GovernmentConsumption	ExpenditureFractions- GrossInvestment	ExpenditureFractions- HouseholdConsumption	ExpenditureFractions- ImportValue		
X	ExpenditureFractions- InventoryChange	ExpenditureFractions-TotalConsumption	PavedAirportLengths-3000To5000Feet	PavedAirportLengths-5000To8000Feet		
	PavedAirportLengths- 8000To10000Feet	PavedAirportLengths- Over10000Feet	PavedAirportLengths-Total	PavedAirportLengths- Under3000Feet		
tics	GDPPerCapita					

The selected data consists of 109 attributes from 132 countries

### **GDP PER** CAPITA

www.evolved-analytics.com

- Source: Mathematica CountryData
- Looking at GDP Per Capita
- Using numeric input and output



### GDP PER CAPITA INFORMATION CONTENT



### GDP PER CAPITA ATTRIBUTE INFORMATION CONTENT





- Scores using SMITS algorithm
- Ranking of attributes for discriminating between countries

### SUMMARY: INFORMATION CONTENT & DATA BALANCING

- We have a means to create a **balanced data** set from unbalanced data
- Balanced data avoids inappropriate weighting of data regions ⇒ better models AND faster modeling
- We can also detect (some) outliers during the data balancing
- The SMITS algorithm does not assume a data model (just an information metric)
- Identifying the dominant records can be insightful

   these are the prototypes

### GDP PER CAPITA MODELING RESULTS





- 1207 different variable
- combinations in developed models



### GDP PER CAPITA VARIABLE SELECTION & NICHING



Despite the fact that we have a **fat array and correlated inputs**, we managed to identify and isolate driving variables

www.evolved-analytics.com

- Diverse model structures were uncovered that fit the data
- The discovered models feature **diversity** in numbers of models (dimensionality) and different variable combinations (subspaces)
- We can select models for use which have limited dimensionality or targeted variable combinations



### GDP PER CAPITA OUTLIER DETECTION

## Outlier Detection HAS to be model-based.

	Complexity	$1-R^2$	Function
1	27	0.030	-150.723 + 0.652 TotaX Consumption 4508 200 -FemalePopulation
2	22	0.023	-135.975 + 0.315 (FixedInvestment+TotalConsumption) FemaleAdultPopulation
3	34	0.017	110.237 + 0.815 (Fixed Investment + GovernmentReceipts + HouseholdConsumption) Population
4	34	0.028	10.236 + 0.469 [Female Adult Population -GrowInvestment - TotalConsumption] Female Population
5	34	0.029	370.370 + 0.325 (GovernmentReceipts-GrowInvextment-TotalConsumption) Female Population
6	34	0.030	120.922 + 0.472 (GovernmentReceipts-GrowInvestment-TotalConsumption) AdultPopulation
7	37	0.030	-123.700 + 1.334 TotalConsumption 10342.930-2 Female Population
8	37	0.030	-93.049 + -0.954 (FixedInvestment+TotalConsumption) 6.271-Population
9	22	0.023	-22.434 + 0.642(2+GovernmentSurplus+GrowsInvestment+TotalConsumption) AdultPopulation
10	41	0.025	54.309 + 0.165 (GovernmentReceipts+2 TatalConsumption) FemaleAdultPopulation
11	45	0.014	36.600 + 0.603 (2.136 Fixed Investment - Government Receipts - Household Consumption) Population
12	45	0.020	-79:205 + 0.307 (5 FixedInvestment+GovernmentReceipts+HouseholdConsumption) AdultPopulation
13	62	0.014	0.682 3FixedInvestment+GovernmentReceipte-GrossInvestment+HouseholdConsumption+Population 33.122 + Population
14	63	0.018	-371.255 + 0.164 (AdukPopulation-Export Value+5 GrowInvestment+3 TotalConsumption) -371.255 + AdukPopulation
15	64	0.017	-416.607 +
16	80	0.012	-00.204 + 0.007/ExportValue-14 FixedInvestment+5 GovernmentReceipts-Population+7 TotalConsumption Population
17	ao	0.014	70.673 + 0.122 (Export Value +5 Government Receipts + 0 Grood avestment + Honsehold Consumption + 3 Total Consump Population
10	87	0.014	-52.551 + 0.552 (7 FixedInvestment+2 GovernmentReceipts+2 GrossInvestment+4 Total Consumption) Proportion



### SYMBOLIC REGRESSION: SUMMARY BENEFITS

- Compact Nonlinear Models
  - Compact empirical models can be suitable for online implementation
  - Model(s) can be used as an emulator for coarse system optimization
- Driving Variable Selection & Identification
  - Appropriate models may be developed from poorly structured data sets (too many variables & not enough measurements)
  - Identified driving variables may be used as inputs into other modeling tools
- Metasensor (Variable Transform) Identification
  - Identifying variable couplings can give insight into underlying physical mechanisms
  - Identified metavariables can enable linearizing transforms to meld symbolic regression and more traditional statistical analysis
- Diverse Model Ensembles

- The independent evolutions will produce independent models. Independent (but comparable) models may be stacked into ensembles whose divergence in prediction may be an indicator of extrapolation & model trustworthiness. This is an issue in high dimensional parameter spaces.
- Human Insight
  - The transparency of the evolved models as well as the explicit identification of the model complexity-accuracy trade-off is very compelling
  - Examining an expression can be viewed as a visualization technique for highdimensional data
- Rapid Modeling
- Exploitation of the Pareto front has resulted in orders-of-magnitude in the symbolic regression performance relative to more traditional GP. This greatly increases the range of possible applications.
- Metavariables can also be used as inputs into other modeling tools

EC IN DOW CHEMICAL

Application Domains	Examples
Material Design	<ul> <li>Color Matching</li> <li>Appearance Engineering</li> <li>Polymer Design</li> <li>Synthetic Leather</li> </ul>
Materials Research	<ul> <li>Diverse Chemical Library Selection</li> <li>Fundamental Model Building</li> <li>Reaction Kinetics Modeling</li> <li>Combi-Chem Catalyst Exploration</li> <li>Combi-Chem Data Analysis</li> </ul>
Production Design	• Acicular Mullite Emulator • EDC/VCM Nonlinear DOE • Bioreactor Optimization
Production Monitoring & Analysis	<ul> <li>Epoxy Holdup Monitoring</li> <li>Isocyanate Level Estimation</li> <li>FTIR Calibration Variable Selection</li> <li>Poly-3 Volatile Emission Monitoring</li> <li>Percet Emulator for Online Optimization</li> <li>Emissions Monitoring</li> </ul>
Business Modeling	Diffusion of Innovation     Hydrocarbon Trading & Energy Systems Optimization     Scheduling Heuristics     Plant Capacity Drivers

www.evolved-analytics.com

## **MORE INFORMATION**

## ACKNOWLEDGEMENT

www.evolved-analytics.com

Guido Smits Katya Vladislavleva Arthur Kordon



