*Evolved Analytics LLC*

# Big Insight vs. Big Data

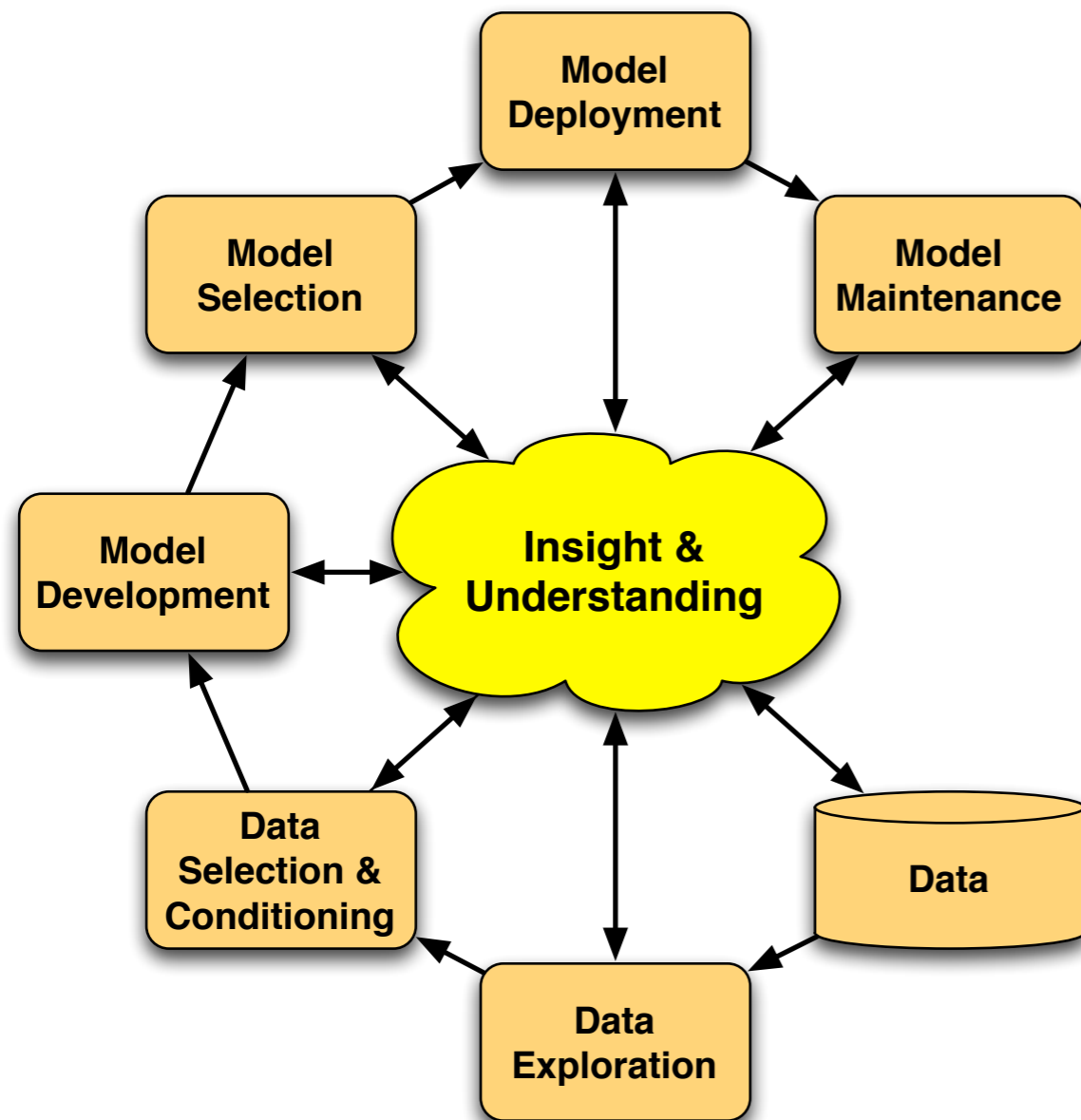Mark Kotanchek

# Modeling Options

Driving variables are known

Model structure is KNOWN to be LINEAR

**LINEAR REGRESSION**

LINEAR SYSTEM & DRIVING VARIABLES ARE KNOWN

---

Driving variables are known

Model structure is KNOWN but NONLINEAR

**NON-LINEAR REGRESSION PARAMETER ESTIMATION**

KNOWN NONLINEAR SYSTEM & DRIVING VARIABLES ARE KNOWN

---

Driving variables are known

Model structure is NOT known

**NEURAL NETWORKS, SVM, RANDOM FORESTS, SYMBOLIC REGRESSION**

UNKNOWN MODEL STRUCTURE BUT DRIVING VARIABLES ARE KNOWN

---

Driving variables are NOT known

Model structure is NOT known

**SYMBOLIC REGRESSION**

MODEL STRUCTURE IS NOT KNOWN & DRIVING VARIABLES ARE NOT KNOWN

# Modeling Workflow is Everything!



- ❖ **Awareness** & **execution** across all aspects
- ❖ Analysis flow is **iterative**
- ❖ Utilize **visualization** to guide analysis
- ❖ An **audit trail** is fundamental

3

# Data in the Real World

Missing Elements

Wide

Lots of Records

Too Little Data

Correlated Variables

Wrong Data

Noisy

Unreliable Sensors

# Key Point
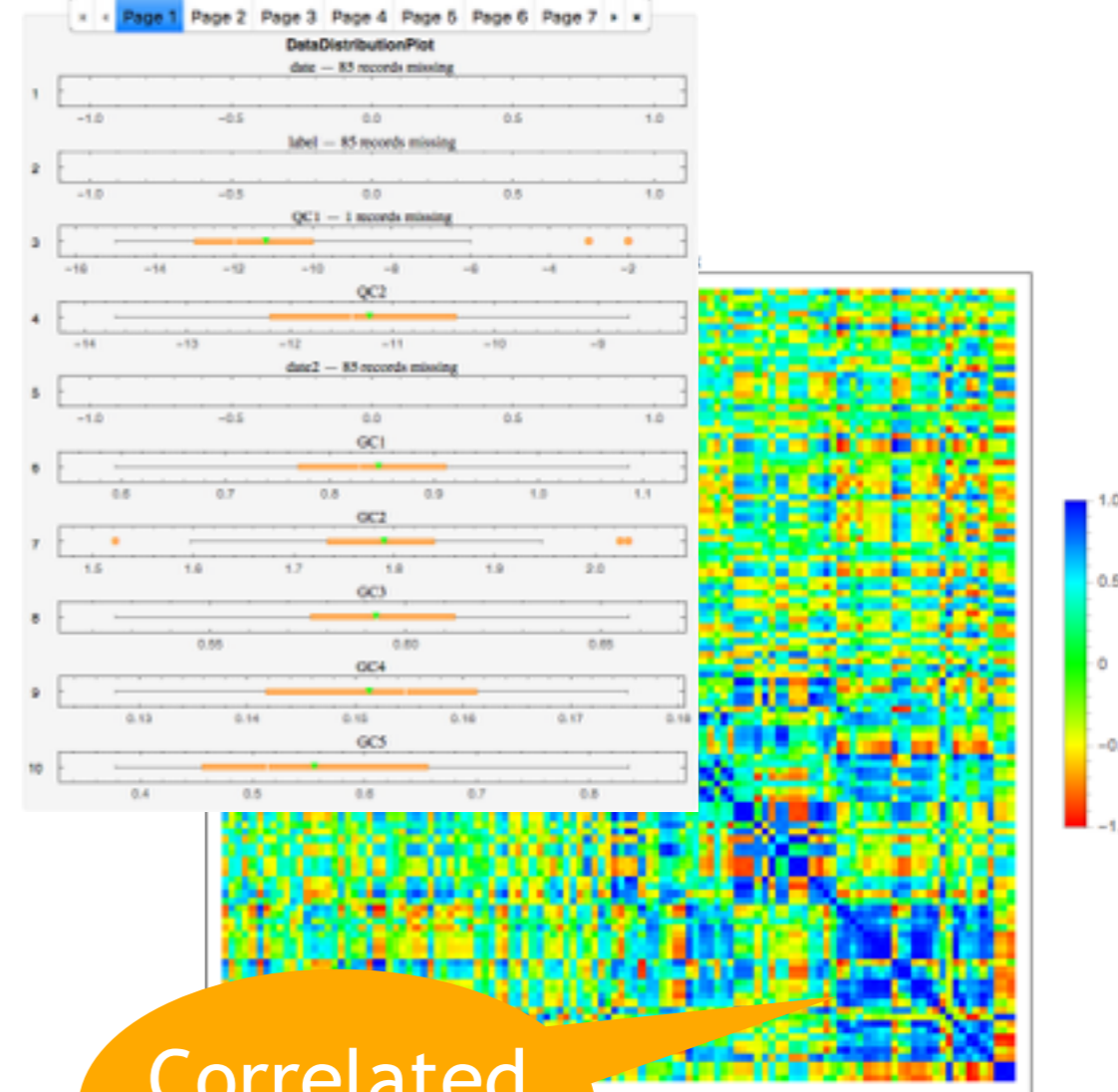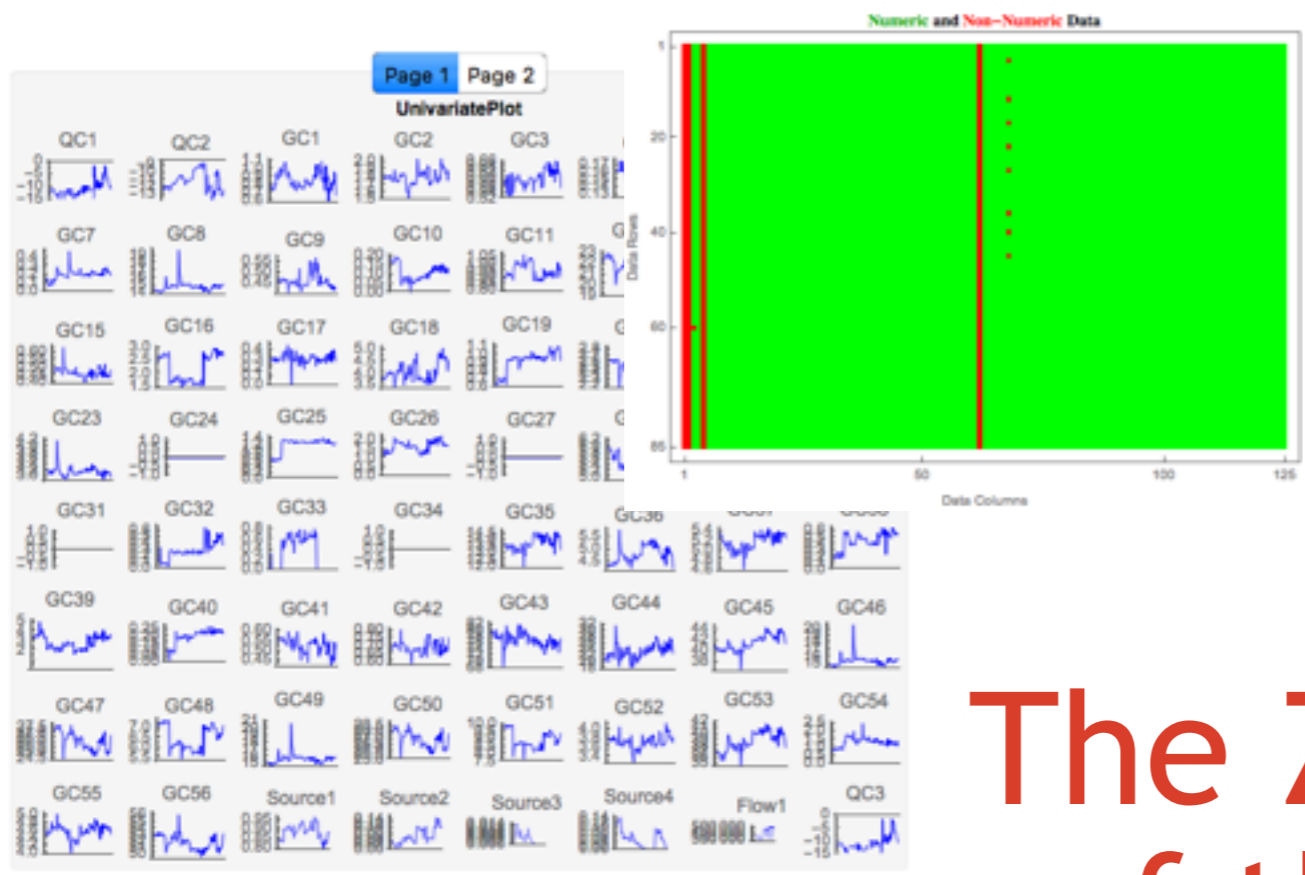
## SYMBOLIC REGRESSION $\Rightarrow$ HYPOTHESIS GENERATOR

The only constraint is the supplied building blocks.

Human limits of imagination & possibility are not imposed!

We can exploit this creativity to produce trustable data models.

# The Illustration Data

* Data from an industrial chemical reactor

* Having problems with product quality (QC)

* 125 variables sampled over three months

* Chemical composition from gas chromatography (GC)

* Process information from plant (flows, temps & source)

* Two users:

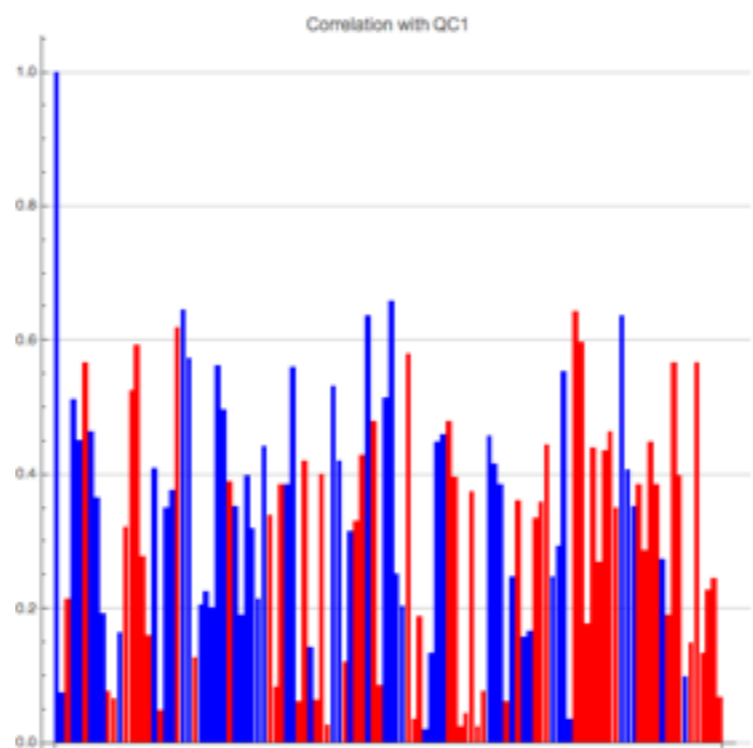  * Analytical Chemists (Why??)

  * Production Engineers (What to do??)

# The Zen of the data

Correlated Inputs

Generate Models

Developed Models

volved-analytics.com

# Supporting the Chemists (Round 2)

**MetaVariableDistributionTable**

| | Rank | # models | MetaVariable | # Evolutions | % Evolutions | Max Count | Max % | Mean % |
|---|---|---|---|---|---|---|---|---|
| ++ | 1 | 1611 | $\frac{1}{gC13}$ | 32 | 100.0 | 133 | 77.3 | 35.3 |
| ++ | 2 | 666 | $gC40^3$ | 19 | 59.4 | 118 | 85.5 | 14.6 |
| ++ | 3 | 121 | $gC18^3$ | 2 | | | | |
| ++ | 4 | 437 | $\sqrt{gC13}$ | 25 | | | | |
| ++ | 5 | 170 | $gC40^2$ | 12 | | | | |
| ++ | 6 | 173 | $\sqrt[3]{gC4}$ | 14 | | | | |
| ++ | 7 | 302 | $gC13^3$ | 23 | | | | |
| ++ | 8 | 398 | $\sqrt[3]{gC13}$ | 28 | | | | |
| ++ | 9 | 281 | $\frac{1}{gC52}$ | 21 | | | | |
| -- | 10 | 125 | $\frac{gC52}{gC13}$ | 13 | | | | |
| ++ | 11 | 123 | $gC40^{4.0}$ | 9 | | | | |
| ++ | 12 | 117 | $gC52^3$ | 16 | | | | |
| ++ | 13 | 55 | $gC18^3\ gC40^3\ gC52^3$ | 1 | | | | |
| ++ | 14 | 235 | $\frac{1}{gC2}$ | 24 | | | | |
| ++ | 15 | 236 | $gC2^3$ | | | | | |
| ++ | 16 | 167 | $gC13^{-3.0}$ | 22 | | | | |
| ++ | 17 | 64 | $gC39^3$ | 4 | | | | |
| ++ | 18 | 47 | $gC39^3\ gC40^3$ | 1 | | | | |
| ++ | 19 | 45 | $\sqrt[3]{gC18}$ | | | | | |
| ++ | 20 | 221 | $\frac{1}{gC4}$ | 16 | | | | |

Variable set saved as GC Driver Variables (Rnd 2 – 20%)

Driver Variables — Off
Variables to Plot — 0.2
Bar Origin — Left
Image Size — 500
Aspect Ratio — 1
Plot Label — Enter text here.
Reset — <<<

**Variable Presence Chart**

gC13 ⟹ 75.7 % ( 495 of the 654 models )

**Variable Presence Chart**

gC40 ⟹ 91.7 % ( 600 of the 654 models )

**Variable subsets can be isolated for focused modeling**

**MetaVariables can be explored and used in subsequent modeling**

**Distribution of MetaVariables from IndependentEvolutions**

Supplied MetaVariables were exploited

gC44–gC50 is in 985 models total

**Model Dimensionality Table**

| # Vars | # Models ⟹ % | VariableCombinationMap | ParetoFrontPlot |
|---|---|---|---|
| 5 | 186 ⟹ 16.3 % | | |
| | 191 ⟹ 16.8 % | | |
| 7 | 455 ⟹ | | |
| 8 | 222 ⟹ | | |

# Chemists (focused)

It looks like 5 variables are required

Page 1  Page 2

**Variable Combination Table**

| num ⟹ % | Variables Used | ParetoFrontPlot |
|---|---|---|
| 1  76 ⟹ 35.3 % | gC4 gC13 gC44 gC50 gC52 | |
| 2  70 ⟹ 32.6 % | gC2 gC13 gC21 gC40 gC52 | |
| | gC4 gC13 gC21 | |

□ −   ☑ +   ☑

This is an interesting combination

**Distribution of VariablePresence from Indepe**

gC2
gC4
gC13
gC18
gC21
gC30
gC40
gC44
gC50
gC52

gC50 ⟹ 86.4 % ( 985 of the 1140 models )

Each search is stochastic

12

www.evolved-analytics.com

Pareto Front Context Plot — 1379 of 1379 selected

Modeling only using PI variables

$$16575.37 - \frac{1849591.70}{temp_{20}} - 37.55\, temp_{20} + 0.48\, temp_{25}$$

**Model Dimensionality Table**

| # Vars | # Models ⟹ % | VariableCombinationMap | ParetoFrontPlot |
|--------|--------------|------------------------|-----------------|
| 2 | 118 ⟹ 8.6 % | temp25 is used by 11.0 % ( 13 of 118 models) | |
| 3 | 211 ⟹ 15.3 % | | |
| 4 | 336 ⟹ 24.4 % | | |
| 5 | 363 ⟹ 26.3 % | | |

3 or 4 vars needed for good-enough models

# Supporting Production

Define inputs for focused model development (vars in at least 20% of models)

Driver Variables — Off

Variables to Plot
Bar Origin
Image Size
Aspect Ratio
Plot Label — Enter text here.
Reset

✓ Automatic
All
0.01
0.03
0.05
0.07
0.09
0.1
0.2
0.3
0.4
0.5
0.6
0.7
0.8
0.9
1
5
10
15
20

# Focused Production Models

# Deployable Models



**Pareto Front Plot — 1921 models**

**Model Properties for qC1**

| | |
|---|---|
| **Model Expression** | $14190.43 + \dfrac{33.68}{flow_{32}} - \dfrac{8.18}{flow_{32} - source_1} - \dfrac{1584564.10}{temp_{20}} - 31.81\, temp_{20}$ |
| **R Squared** | 0.841055 |
| **Adjusted R–Squared** | 0.833007 |

| | | DF | SS | MS | F-Statistic | P-Value |
|---|---|---|---|---|---|---|
| **ANOVA Table** | $\dfrac{1}{flow32}$ | 1 | 7.95431 | 7.95431 | 6.87623 | 0.0104768 |
| | $\dfrac{1}{flow32 - source1}$ | 1 | 77.7843 | 77.7843 | 67.2419 | $3.55065 \times 10^{-12}$ |
| | $\dfrac{1}{temp20}$ | 1 | 237.426 | 237.426 | 205.246 | $1.13695 \times 10^{-23}$ |
| | temp20 | 1 | 160.402 | 160.402 | 138.662 | $4.57606 \times 10^{-19}$ |
| | Error | 79 | 91.3859 | 1.15678 | | |
| | Total | 83 | 574.952 | | | |

- ❖ 80 minutes of model development (32 independent searches of 10 minutes on a quad-core laptop)

- ❖ Models were rewarded for simplicity and accuracy

- ❖ The individual models are good; however, we want trustable models based upon ensembles
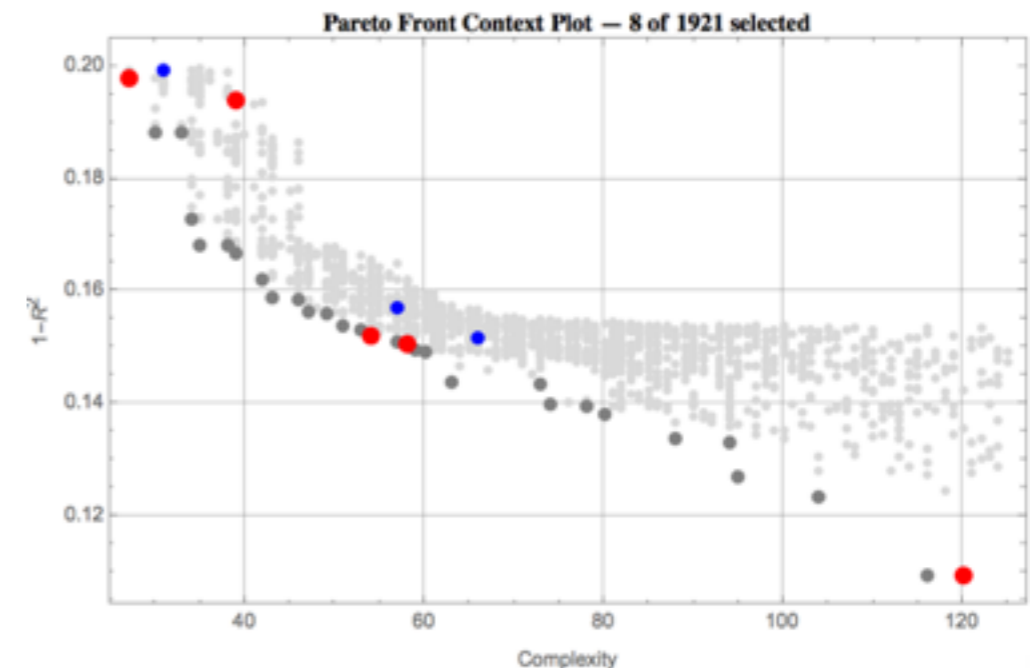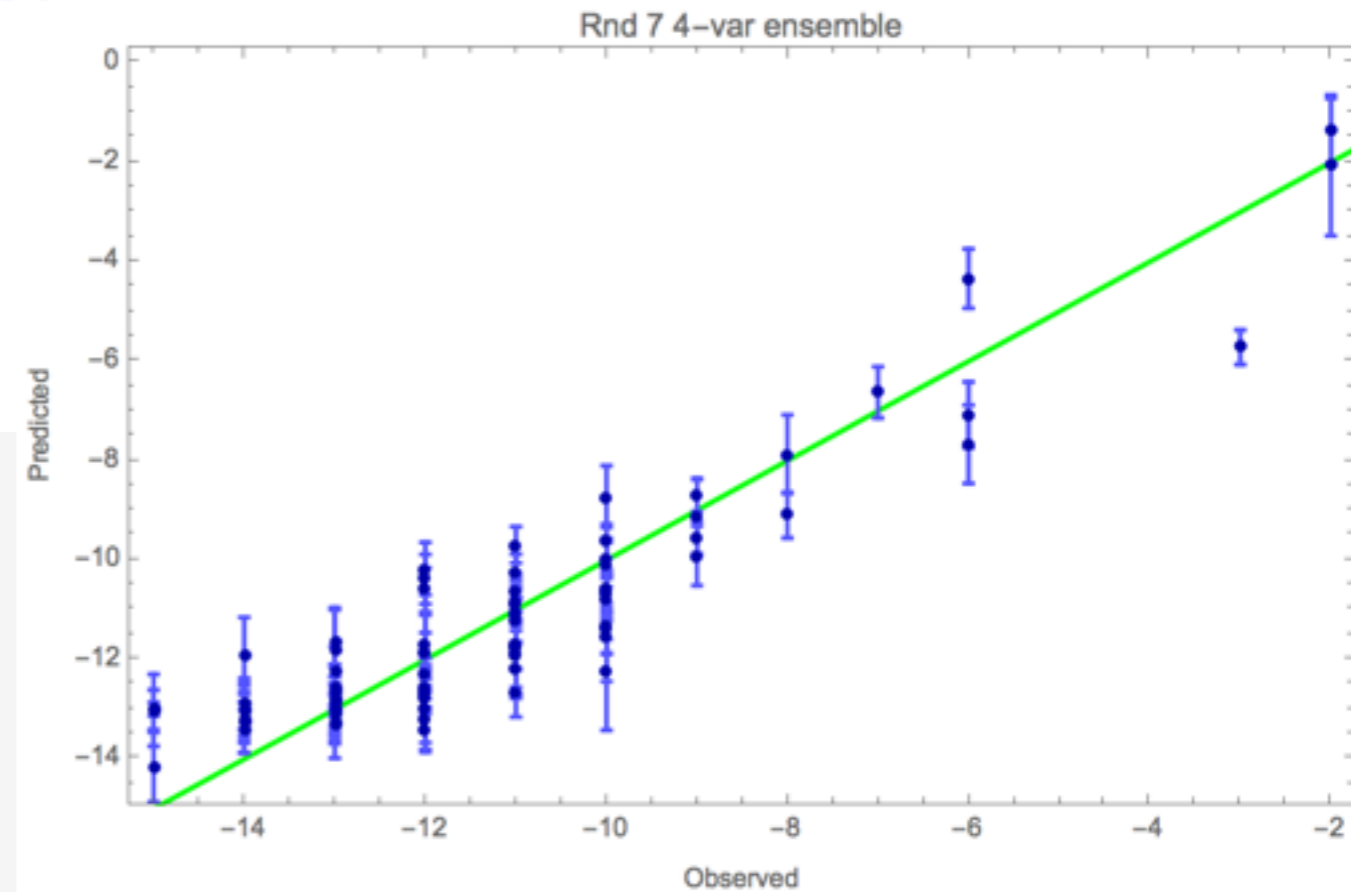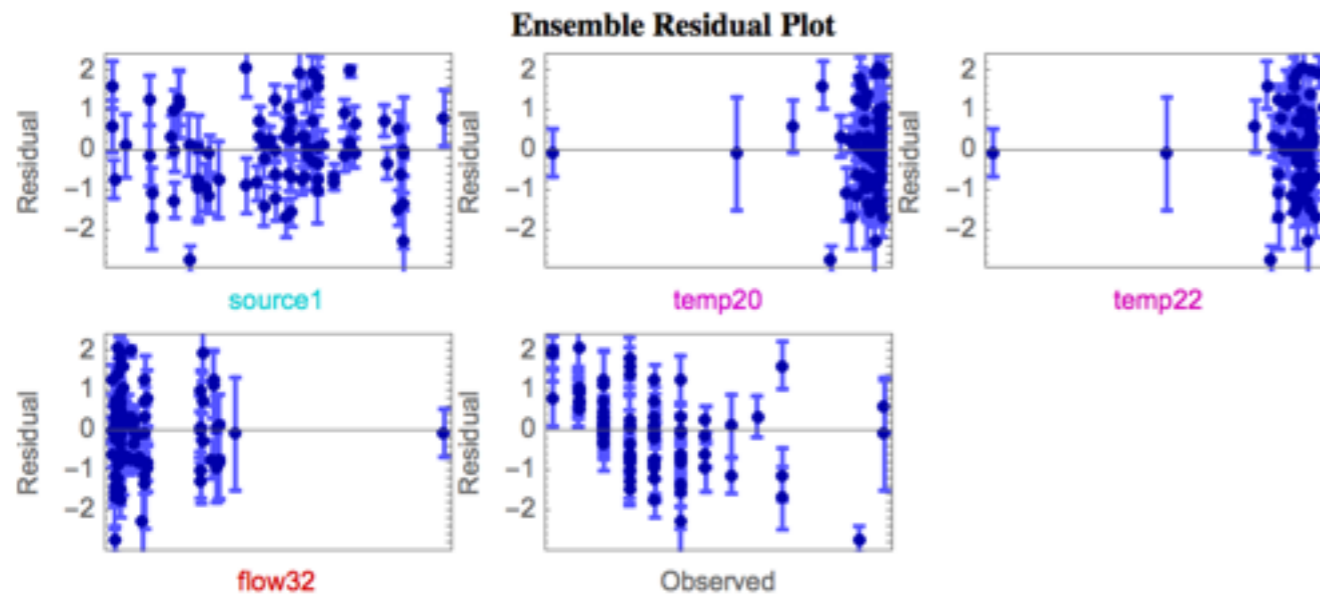
# An Ensemble

- From candidate (accurate and simple) models, models chosen for their diversity

- Ensemble has better prediction accuracy than individual models

- Divergence of models provides a trust metric!

**Model Selection Table**

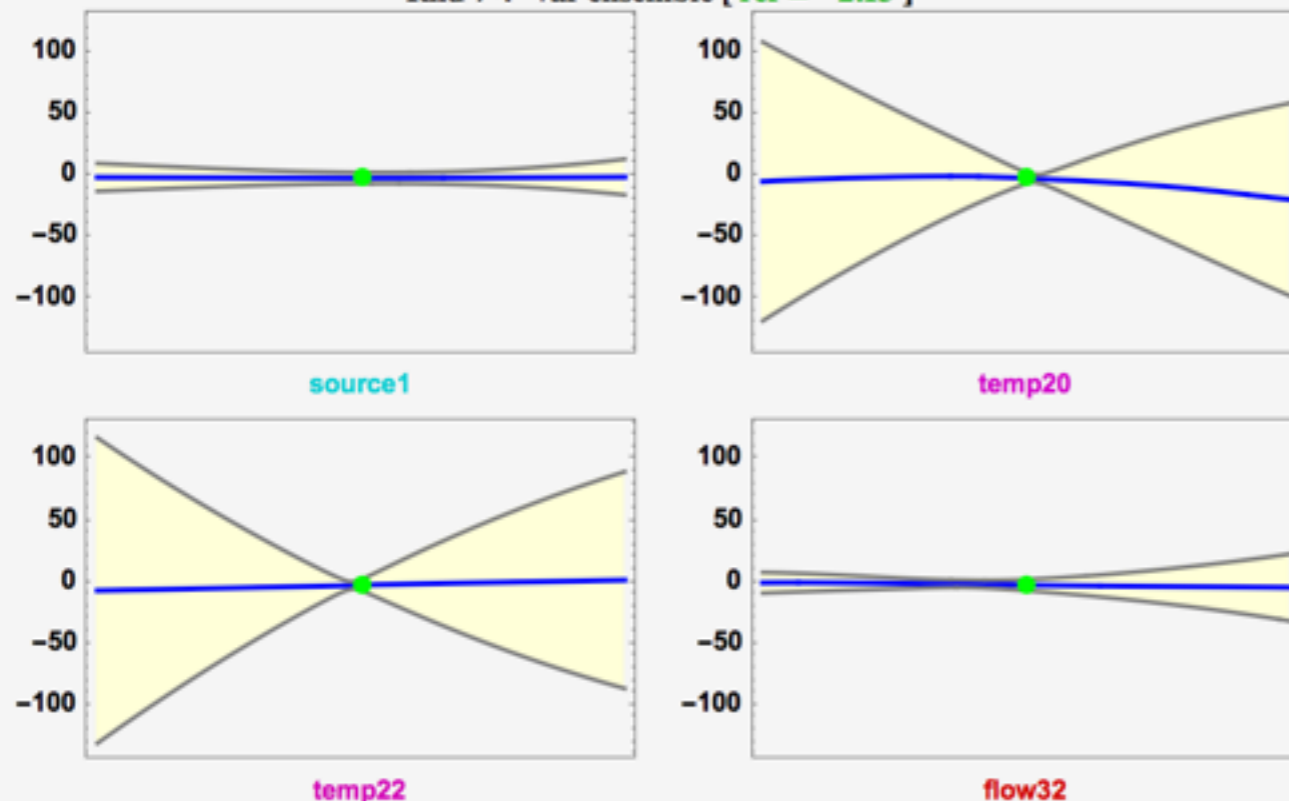| | Complexity | $1-R^2$ | Vars | Function |
|---|---|---|---|---|
| 1 | 27 | 0.198 | temp20, flow32 | $16237.31 - 1.04\,\text{flow}_{32} - \frac{1814657.90}{\text{temp}_{20}} - 36.31\,\text{temp}_{20}$ |
| 2 | 31 | 0.199 | source1, temp20, flow32 | $434.45 - 0.19\,\text{flow}_{32}^3 - 23.20\,\text{source}_1 - 1.82\,\text{temp}_{20}$ |
| 3 | 39 | 0.194 | temp20, temp22 | $-211858.98 + \frac{15344821.00}{\text{temp}_{20}} + 973.50\,\text{temp}_{20} - 1.49\,\text{temp}_{20}^2 + 0.53\,\text{temp}_{22}$ |
| 4 | 54 | 0.152 | source1, temp20, temp22, flow32 | $14648.69 - 21.78\,\text{flow}_{32} + \frac{44.01}{\text{source}_1} + 23.73\,\text{flow}_{32}\,\text{source}_1 - \frac{1643486.50}{\text{temp}_{20}} - 33.23\,\text{temp}_{20} + 0.40\,\text{temp}_{22}$ |
| 5 | 57 | 0.157 | source1, temp20, temp22, flow32 | $-3664.72 - 34.66\,\text{flow}_{32}^{1/3} + 25.20\,\text{temp}_{20} - (1.71\times10^{-4})\text{temp}_{20}^3 - \frac{5822.82\,\text{source}_1}{\text{flow}_{32}\,\text{temp}_{22}}$ |
| 6 | 58 | 0.151 | source1, temp20, temp22, flow32 | $13518.39 + \frac{21.80}{\text{flow}_{32}} - \frac{31.04\,\text{source}_1^3}{\text{flow}_{32}^3} - \frac{1504966.60}{\text{temp}_{20}} - 30.64\,\text{temp}_{20} + 0.28\,\text{temp}_{22}$ |
| 7 | 66 | 0.152 | source1, temp20, temp22, flow32 | $14059.10 - \frac{4.09}{\text{flow}_{32}} - 11.64\,\text{flow}_{32} - 54.13\,\text{source}_1 + 12.50\,\text{flow}_{32}\,\text{source}_1^2 - \frac{1562528.70}{\text{temp}_{20}} - 31.73\,\text{temp}_{20} + 0.38\,\text{temp}_{22}$ |
| 8 | 120 | 0.109 | source1, temp20, temp22, flow32 | $10870.11 + 471.98\,\text{flow}_{32} - \frac{2345.05}{\text{source}_1} - 5241.64\,\text{source}_1 + 891.67\,\text{source}_1^3 - 1161.84\sqrt{\text{temp}_{20}} + 42.50\,\text{temp}_{20} - 4.89\sqrt{\text{flow}_{32}^2}\,\text{source}_1\,\text{temp}_{20} + 37.30\,\text{temp}_{22} + 3.16\,\text{flow}_{32}\,\text{source}_1^2\,\text{temp}_{22} - (9.50\times10^{-2})\text{temp}_{22}^2$ |

**Model Ensemble Properties for Rnd 7 4-var ensemble**

| Model Expression | MedianAverage[...] |
|---|---|
| R Squared | 0.85464 |
| Adjusted R-Squared | 0.852867 |
| ANOVA Table | Anova is not available for model ensembles. |

Pareto Front Context Plot — 8 of 1921 selected

# Ensemble Performance



Ensemble Residual Plot

Rnd 7 4-var ensemble

Rnd 7 4-var ensemble [ ref = −2.13 ]

- ❖ Ensemble predictions have a trust metric based upon divergence.

- ❖ Temperatures are coupled so they cannot be varied independently $\implies$ prediction spread greatly increases if we try to do that!

# Conclusions

❖ One more thing …

   ❖ Modeling and data results can be archived to analysis report at the click of a button

   ❖ A function package is available for use in a notebook front end as well as to facilitate automated analysis flows

❖ For more information or trial licenses for DataModeler, contact

   ❖ info@evolved-analytics.com

❖ We also do consulting and custom analysis system development

❖ We have offices in the US and Europe