# Industrial Evolutionary Computing

Arthur Kordon[#], Guido Smits[#], and Mark Kotanchek[+]

**The Dow Chemical Company [#]**
**Evolved Analytics [+]**

**GECCO 2007**

---

## Overview

*In theory, there is no difference between theory and practice. In practice, there is.*
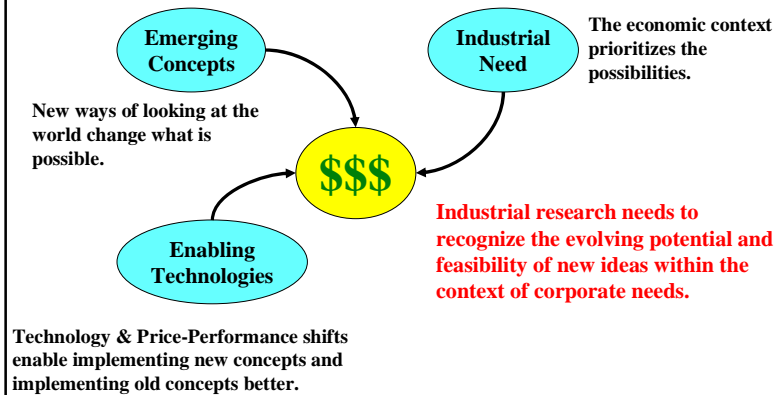- Jan L.A. van de Snepscheut

- Evolutionary Computing and the business model
- Key Technologies
  - Analytic Neural Networks + Support Vector Machines + Genetic Programming + Particle Swarms + …
- Implementation Guidelines
- Integrate & Conquer
- Key Application Areas
- Open Issues & Research areas

Kordon, Smits & Kotanchek

GECCO 2007                                                                2

---

## Data Modeling
### At the Intersection of Opportunity & Need



**Emerging Concepts**

New ways of looking at the world change what is possible.

**Industrial Need**

The economic context prioritizes the possibilities.

**$$$**

**Enabling Technologies**

Industrial research needs to recognize the evolving potential and feasibility of new ideas within the context of corporate needs.

Technology & Price-Performance shifts enable implementing new concepts and implementing old concepts better.

Kordon, Smits & Kotanchek
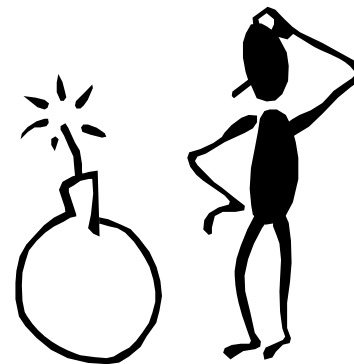
GECCO 2007                                                                3

---

## Motivation



- Industry is great at collecting data … and then performing records retention
- Extracting insight from multivariate data is hard
- Time and money is being wasted

*"We are drowning in information and starving for knowledge" –*
*R.D. Roger*

Kordon, Smits & Kotanchek

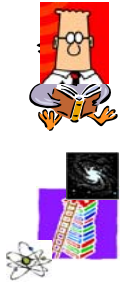GECCO 2007                                                                4

## Academic vs. industrial data analysis

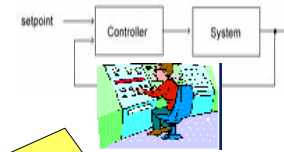Transfer data into knowledge        Transfer data into value

## Special Features of Industrial Data Analysis
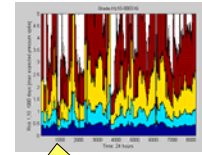
**Operators intervention**          **Curse of closed loops**

Operators manually modify the process

setpoint — Controller — System

The majority of process variables are in closed loops and depend on controller adjustments

**Multiple time scales**            **Real-time pressure**

Time scales vary from milliseconds to months

Models need to be developed & updated rapidly

Most of models operate in real time

Kordon, Smits & Kotanchek

GECCO 2007                                                                 6

## Intelligent Systems in Industrial Data Analysis: Lessons From the Past

inside

pentium PROCESSOR

**The Expert Systems campaign (late 80s) "We'll put engineers in the box"**

• static rule-based models not linked to numerical world

• the politics of knowledge acquisition

• the efforts of knowledge acquisition

**The Neural Networks campaign (early 90s) "We'll turn data into gold"**

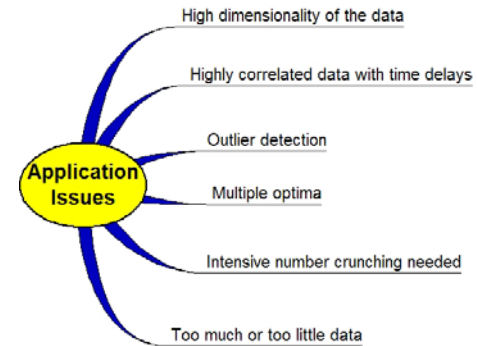• black-box models with inefficient structure

• fragile models and model validation

• maintenance nightmare

Kordon, Smits & Kotanchek

GECCO 2007                                                                 7

## Industrial Data Modeling Issues

High dimensionality of the data

Highly correlated data with time delays

Outlier detection

**Application Issues**

Multiple optima

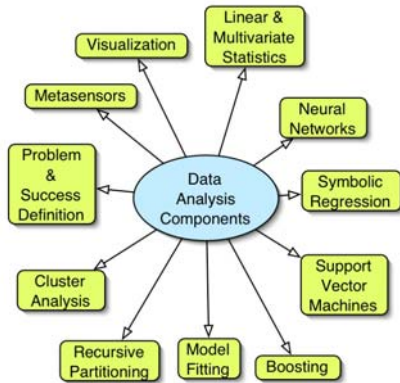Intensive number crunching needed

Too much or too little data

"The most exciting phrase to hear in science, the one that heralds new discoveries, is not 'Eureka!' (I found it!) but 'That's funny …'" — Isaac Asimov (1920 - 1992)

Kordon, Smits & Kotanchek

GECCO 2007                                                                 8

## Industrial data analysis components



The role of evolutionary computing (symbolic regression) is to …

– Facilitate physical/mechanism insight and **understanding**
– **Summarize** data behavior
– Identify data **transforms** and metasensors
– Perform **variable selection**
– Enable response surface **exploration and optimization**
– **Visualize** behavior in the form of a symbolic expression

The overall goal is to achieve speed, accuracy & efficiency. Symbolic regression is part of an integrated methodology.

Kordon, Smits & Kotanchek

GECCO 2007                                                                 9
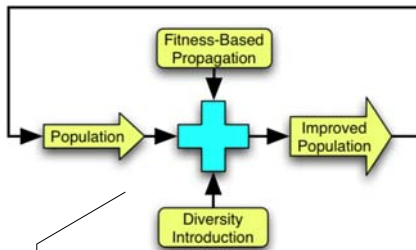
---

# Competing/Complementary Technologies

- Linear Models
  - Linear in coefficients, not necessarily linear in model
  - Often "good enough" and simple
  - Well developed criteria and foundations in linear statistical analysis
  - Typically easy and fast to develop (unless subtleties are involved)
- Neural networks
  - Often good performance but lots of "trust me"
  - A good reference for nonlinear modeling potential

- Support Vector Machines
  - Useful for data compression to match information content
  - Computationally demanding
  - Unique nonlinear outlier detection capability
- Fuzzy Rules/Recursive Partitioning
  - Human interpretability — if simple
  - Can handle categorical data

Kordon, Smits & Kotanchek

GECCO 2007                                                                 10

---
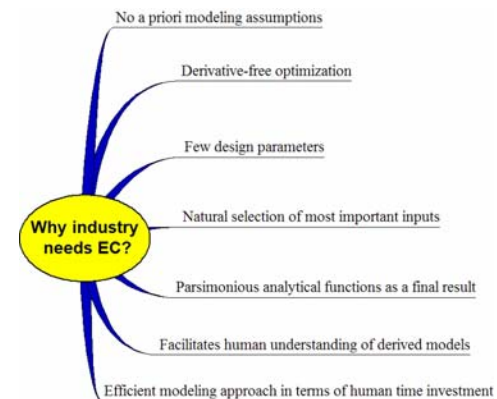
# Evolutionary Computing Theory



**It is this simple!**

Variants:
– Genetic Algorithms (GA)
– Evolutionary Strategies (ES)
– Evolutionary Programming (EP)
– Genetic Programming (GP)
– Particle Swarm Optimization (PSO)
– Gene Expression Programming (GEP)
– etc.

Genetic Programming
– Genome (genetic code) evolves
– Phenotype (realization) judged for fitness
– Goal is to evolve *programs* which solve problems
– The search space is *infinite*!
– Symbolic regression is one application of genetic programming

Symbolic Regression
– Goal is to identify expressions which summarize data
– NOT parameter fitting — discovery of both structure and parameters
– The search space is infinite!
– In practice, symbolic regression is part of an integrated methodology

Kordon, Smits & Kotanchek

GECCO 2007                                                                 11

---

## Why industry needs Evolutionary Computing?



No a priori modeling assumptions

Derivative-free optimization

Few design parameters

Natural selection of most important inputs

Parsimonious analytical functions as a final result

Facilitates human understanding of derived models

Efficient modeling approach in terms of human time investment

Why industry needs EC?

Kordon, Smits & Kotanchek

GECCO 2007                                                                 12

## Economic benefits from Evolutionary Computing



- Resolve complex optimization problems
  - Genetic Algorithms
  - Particle Swarm Optimization
  - Ant Colony Optimization
- Physical interpretation & insight (Symbolic regression)
  - Suggestions for profitable directions for R&D
  - Accelerate R&D
  - Higher credibility in comparison to black-boxes
- Reduce model development cost (significantly reduced development time relative to alternatives)
- Reduce model exploitation cost
  - Minimal model implementation cost (no need for specialized software)
  - Reduced maintenance cost (less frequent re-training)
- Reduce cost of industrial experiments (minimizes the number of additional experiments)

**Benefits from EC**

Kordon, Smits & Kotanchek

GECCO 2007                                                                 13

## Benefits of integrating Evolutionary Computing with other approaches



- Increased quality of generated models
  - Data with high information content
  - Model complexity measure
- Reduced model development time and cost
  - Condensed data sets
  - Faster model selection
- Broader support from different stakeholders
  - Final users
  - First-principle modelers
  - Statistical community
  - Machine learning community

**Benefits of integration**

Kordon, Smits & Kotanchek

GECCO 2007                                                                 14

## Application areas with impact



**Industrial Applications**

- Research Acceleration
  - Understand Variable Relationships
  - Cues to Physical Mechanisms
  - Explore Multivariate Relationships
- Inferential Sensors
  - Infer System States
  - Online Monitoring & Alarm
- Nonlinear DOE
  - Focus Data Gathering
  - Model Discrimination DOE
- Emulators
  - System Modeling
  - Coarse Optimization
  - Insight into System
- Variable Transforms
  - Meaningful Combinations
  - Convert into less nonlinear problem
- Variable Sensitivity
  - Identify Variables which drive system

Kordon, Smits & Kotanchek

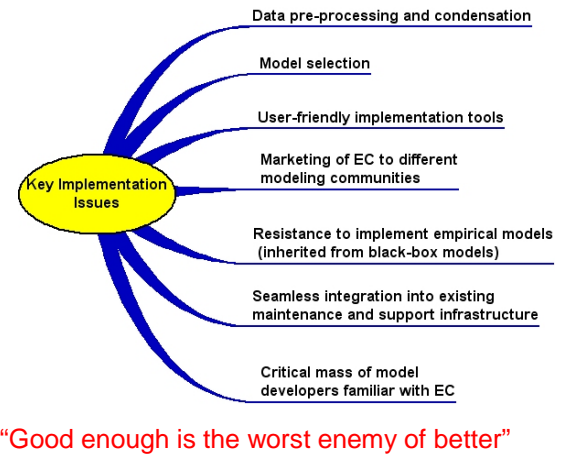GECCO 2007                                                                 15

## Implementation guidelines

- Requirements for successful empirical modeling
- Key issues to be overcome
- Implementation strategy
- Implementation tools

Kordon, Smits & Kotanchek

GECCO 2007                                                                 16
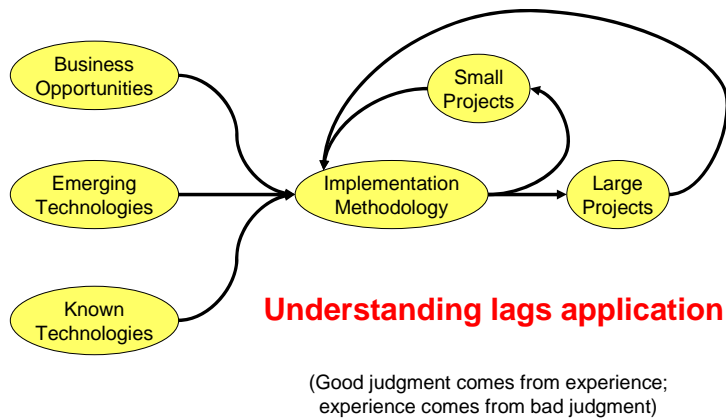
## Requirements for successful data-driven modeling

**Objective function:**
**Minimizing modeling cost and maximizing data analysis efficiency**
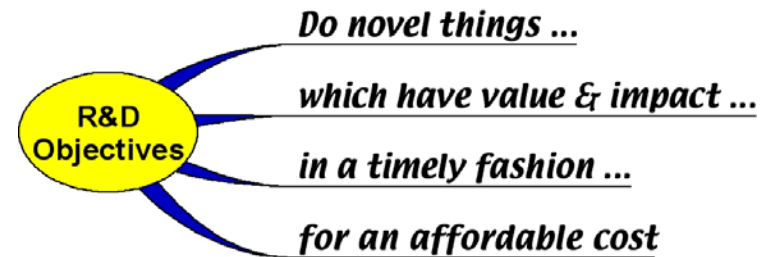**under broad range of operating conditions**
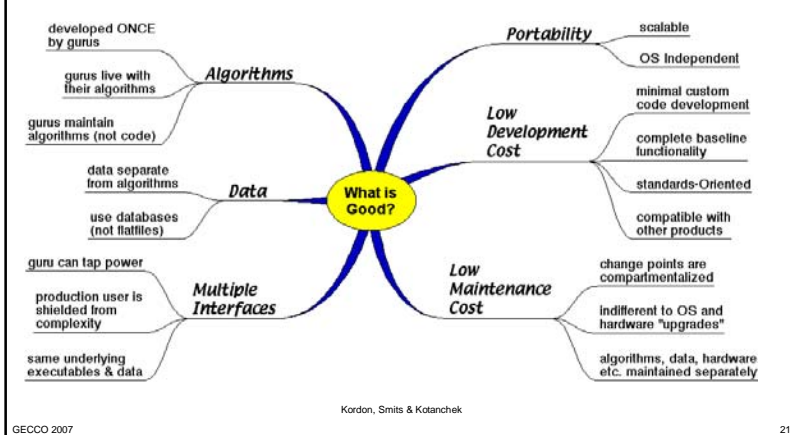


ability to withstand minor changes in targeted system — **Robustness**

**Self-Assessment** — ability to estimate quality of predictions

ability to operate outside training range — **Extrapolations**

**Good Model Aspects**

**Credibility** — the model matches the observed behavior

The total cost-of-ownership (development + operation + maintenance) is proper — **Cost-Effective**

**Interpretability** — humans are able to agree that the model is "reasonable"

Kordon, Smits & Kotanchek

GECCO 2007                    17

## Key issues to overcome



**Key Implementation Issues**

- Data pre-processing and condensation
- Model selection
- User-friendly implementation tools
- Marketing of EC to different modeling communities
- Resistance to implement empirical models (inherited from black-box models)
- Seamless integration into existing maintenance and support infrastructure
- Critical mass of model developers familiar with EC

"Good enough is the worst enemy of better"

Kordon, Smits & Kotanchek

GECCO 2007                    18

# Implementation Strategy



Business Opportunities

Emerging Technologies

Known Technologies

Implementation Methodology

Small Projects

Large Projects

**Understanding lags application**

(Good judgment comes from experience;
experience comes from bad judgment)

Kordon, Smits & Kotanchek

GECCO 2007                    19

# Corporate Research Objectives



**R&D Objectives**

*Do novel things …*

*which have value & impact …*

*in a timely fashion …*

*for an affordable cost*

Kordon, Smits & Kotanchek

GECCO 2007                    20

3301

## Characteristics of a "Good" Analysis System



developed ONCE by gurus

gurus live with their algorithms

gurus maintain algorithms (not code)

*Algorithms*

data separate from algorithms

use databases (not flatfiles)

*Data*

**What is Good?**

guru can tap power

production user is shielded from complexity

same underlying executables & data

*Multiple Interfaces*

*Portability*

scalable

OS Independent

*Low Development Cost*

minimal custom code development

complete baseline functionality

standards-Oriented

compatible with other products

*Low Maintenance Cost*

change points are compartmentalized

indifferent to OS and hardware "upgrades"

algorithms, data, hardware etc. maintained separately

Kordon, Smits & Kotanchek

## Implementation tools

- Mathematica (Dow & Evolved Analytics developed)
  - Symbolic regression package
  - AutoAnalysisTools
  - Analytic neural networks
  - Particle Swarm Optimization (PSO)
  - Group Methods of Data Handling (GMDH)
- MATLAB (Dow developed)
  - Genetic Algorithms (GA)
  - Genetic Programmimg (GP)
  - PSO (single objective and multi-objective)
  - Analytic neural networks
  - Support vector machines
- Tools for model deployment
  - Delphi
  - WebMathematica
  - Excel
  - Process control systems

*Using a commercial framework allows us to bring new concepts and technologies to bear while mitigating the development and long-term maintenance costs of exploiting those technologies.*
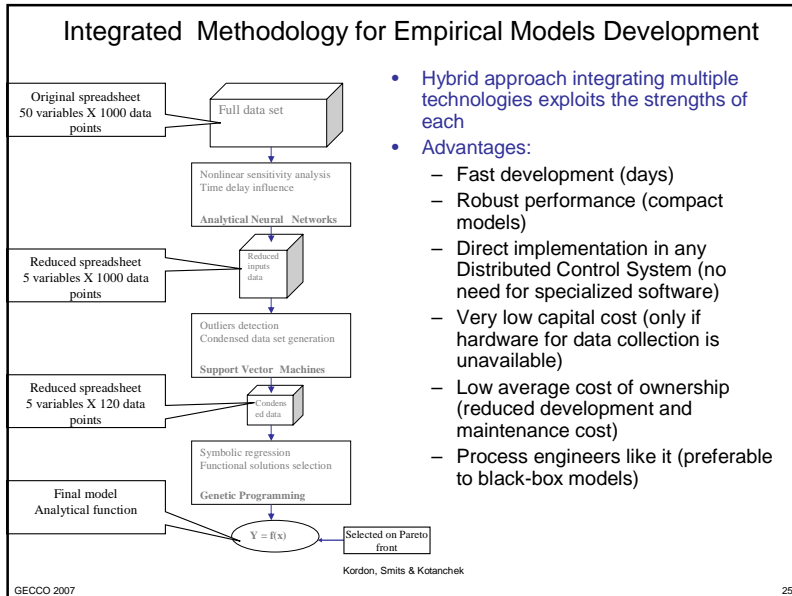
Kordon, Smits & Kotanchek

## Exploitation/Implementation Sequence of Computational Intelligence Approaches in Dow Chemical



Classical NN

GA/GP

Integrated methodology

Pareto GP

1990     1995     2000     2005

Analytic NN

SVM

PSO

Kordon, Smits & Kotanchek

## Integrate & Conquer



- Integrated methodology for successful EC implementation
- Related approaches
- A case study

Kordon, Smits & Kotanchek

## Integrated Methodology for Empirical Models Development



- Hybrid approach integrating multiple technologies exploits the strengths of each
- Advantages:
  - Fast development (days)
  - Robust performance (compact models)
  - Direct implementation in any Distributed Control System (no need for specialized software)
  - Very low capital cost (only if hardware for data collection is unavailable)
  - Low average cost of ownership (reduced development and maintenance cost)
  - Process engineers like it (preferable to black-box models)

Kordon, Smits & Kotanchek

GECCO 2007          25

## Structural Risk Minimization



GECCO 2007          26

## VC-dimension

- In general, VC-dimension does not coincide with the number of parameters (can be larger or smaller)
- VC-dimension of the set of functions is responsible for the generalization ability of learning machines
- Opens remarkable opportunities to overcome the "curse of dimensionality" (large number of parameters, but low VC-dimension)

Kordon, Smits & Kotanchek

GECCO 2007          27

## Two hidden nodes



Any complex surface can be approximated by combining simple surfaces corresponding to a single hidden node

Combination:

Kordon, Smits & Kotanchek

GECCO 2007          28

3303

## Slide 29

### Structural difference between classical and analytic neural networks

**Classical NN**

**Analytical NN**

An additional link between inputs $X_i$ and the output Y is introduced

Bias(1)

Hidden nodes calculation

$Z_1 = F_h(a_{10} + a_{11}X_1 + a_{12}X_2 + a_{13}X_3)$
$Z_2 = F_h(a_{20} + a_{21}X_1 + a_{22}X_2 + a_{23}X_3)$
$Z_3 = F_h(a_{30} + a_{31}X_1 + a_{32}X_2 + a_{33}X_3)$
$Z_4 = F_h(a_{40} + a_{41}X_1 + a_{42}X_2 + a_{43}X_3)$
$Y = F_o(b_0 + b_1Z_1 + b_2Z_2 + b_3Z_3 + b_4Z_4)$

$Z_1 = F_h(a_{10} + a_{11}X_1 + a_{12}X_2 + a_{13}X_3)$
$Z_2 = F_h(a_{20} + a_{21}X_1 + a_{22}X_2 + a_{23}X_3)$
$Z_3 = F_h(a_{30} + a_{31}X_1 + a_{32}X_2 + a_{33}X_3)$
$Z_4 = F_h(a_{40} + a_{41}X_1 + a_{42}X_2 + a_{43}X_3)$
$Y = F_o(b_0 + b_1Z_1 + b_2Z_2 + b_3Z_3 + b_4Z_4 + c_1X_1 + c_2X_2 + c_3X_3)$

Kordon, Smits & Kotanchek

GECCO 2007                                                                29

## Slide 30

### Analytic neural networks have a fixed Capacity

If input-to-hidden layer weights $a_{ij}$ are fixed, there is an analytical solution for the weights $b_i$ and $c_i$

Bias(1)

$$F_o^{-1}(Y) = [1\ X\ Z]* \begin{bmatrix} b_0 \\ c_i \\ b_j \end{bmatrix}$$

$Z_1 = F_h(a_{10} + a_{11}X_1 + a_{12}X_2 + a_{13}X_3)$
$Z_2 = F_h(a_{20} + a_{21}X_1 + a_{22}X_2 + a_{23}X_3)$
$Z_3 = F_h(a_{30} + a_{31}X_1 + a_{32}X_2 + a_{33}X_3)$
$Z_4 = F_h(a_{40} + a_{41}X_1 + a_{42}X_2 + a_{43}X_3)$
$Y = F_o(b_0 + b_1Z_1 + b_2Z_2 + b_3Z_3 + b_4Z_4 + c_1X_1 + c_2X_2 + c_3X_3)$

Standard linear regression problem
X – inputs data matrix (**known**)
Z – hidden layer values vector (**known**)
Unique least-squares solutions for $b_i$ and $c_i$

Kordon, Smits & Kotanchek

GECCO 2007                                                                30

## Slide 31

### Input-to-hidden layer initialization

f=Sig(x,1)

Hidden nodes have to be within the active region of the nonlinear function

The width of the active zone is defined by the steepness of the function or the "temperature"

The "temperature" depends also on the number of inputs to the hidden node

Empirical expression for a normalized "temperature" of a sigmoid function

Weights from the input-to-hidden layer are Sampled from a normal distribution

$$Tn = \eta.\frac{\log(2+\sqrt{3})}{\sqrt{ni - 0.5}}$$

Kordon, Smits & Kotanchek

GECCO 2007                                                                31

## Slide 32

### Analytic Neural Network Benefits

- **Robust** algorithm
  - No tunable parameters
  - One **global** optimum
- Very **fast**,
  - possible to use a whole range of cross-validation principles from statistics
  - No longer an NP-complete problem
- Strong **theoretical foundation**
  - statistical learning theory
  - Direct measure for the model capacity (VC-dimension)

Kordon, Smits & Kotanchek

GECCO 2007                                                                32

3304

## Stacked Analytic Neural Nets (SANN)



- Fast development
- Diverse subnet consensus indicator of model output quality
- Allows explicit calculations of input/output sensitivity
- Can handle time-delayed inputs by convolution functions
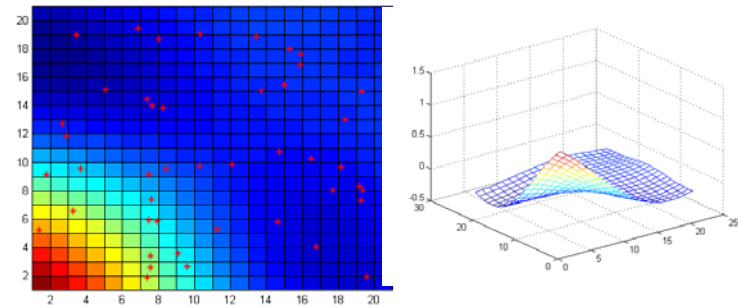- Gives more reliable estimates based on multiple models statistics

Internally developed in Dow Chemical

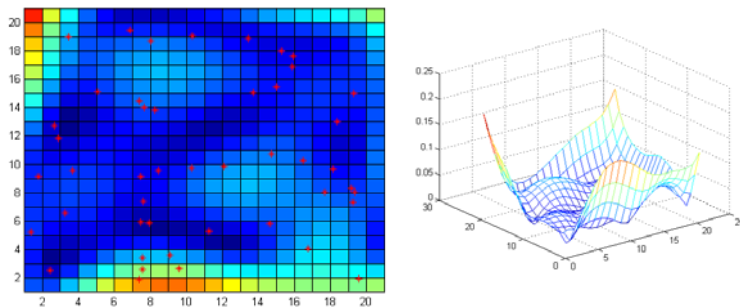Kordon, Smits & Kotanchek

GECCO 2007                                                                                 33

## Model Mismatch Indicator - 2D



Kordon, Smits & Kotanchek

GECCO 2007                                                                                 34

## Model Mismatch Indicator - 2D

Models tend to agree where there is data points and tend to disagree where there is no data.



Kordon, Smits & Kotanchek

GECCO 2007                                                                                 35

## Reduction of the number of input dimensions using Neural Networks

$$SI_j = \frac{\frac{1}{Np}\sum_{p=1}^{Np}\left[\left|\frac{\partial Y}{\partial X_j}\right|\right]_p}{\sqrt{(\mathbf{X'X})_{jj}^{-1}}}$$

$$\frac{\partial NN_n(\mathbf{X})}{\partial X_i} = w_i^n + \sum_{h=1}^{N_h^n} w_h^n a_h^n (1-a_h^n) t_h^n . w_{ih}^n \quad \text{where} \quad a_h^n = \text{Sig}(\sum_{i=0}^{N_i^n} w_{ih}^n X_i^n, t_h^n)$$

$$\frac{\partial CM(\mathbf{X})}{\partial X_i} = \sum_{1}^{N} w_n \frac{\partial NN_n(\mathbf{X})}{\partial X_i}$$

Kordon, Smits & Kotanchek

GECCO 2007                                                                                 36

3305

## An example of stacked analytic NN application - a model for catalyst efficiency

Sensitivity analysis of various process parameters on catalyst efficiency

NN model performance with model disagreement indicator

Model disagreement indicator

Kordon, Smits & Kotanchek

GECCO 2007                                                                 37

## Integrated Methodology for Empirical Models Development

- Hybrid approach integrating multiple technologies exploits the strengths of each
- Advantages:
  - Fast development (days)
  - Robust performance (compact models)
  - Direct implementation in any Distributed Control System (no need for specialized software)
  - Very low capital cost (only if hardware for data collection is unavailable)
  - Low average cost of ownership (reduced development and maintenance cost)
  - Process engineers like it (preferable to black-box models)

Original spreadsheet 50 variables X 1000 data points — Full data set

Nonlinear sensitivity analysis Time delay influence
**Analytical Neural Networks**

Reduced spreadsheet 5 variables X 1000 data points — Reduced inputs data

Outliers detection Condensed data set generation
**Support Vector Machines**

Reduced spreadsheet 5 variables X 120 data points — Condensed data

Symbolic regression Functional solutions selection
**Genetic Programming**

Final model Analytical function — $Y = f(x)$ — Selected on Pareto front

Kordon, Smits & Kotanchek

GECCO 2007                                                                 38

## Explicit Complexity Control in Support Vector Machines (SVM)

Kordon, Smits & Kotanchek

GECCO 2007                                                                 39

## Controlled Data Compression

Kordon, Smits & Kotanchek

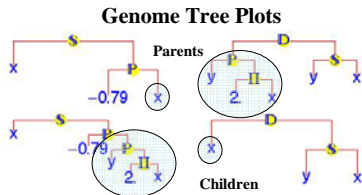GECCO 2007                                                                 40

3306

## Local Kernel



x=(0,0)  x=(0,0.25)  x=(0,0.5)  x=(0,0.75)  x=(0,1)

x=(0.25,0)  x=(0.25,0.25)  x=(0.25,0.5)  x=(0.25,0.75)  x=(0.25,1)

x=(0.5,0)  x=(0.5,0.25)  x=(0.5,0.5)  x=(0.5,0.75)  x=(0.5,1)

x=(0.75,0)  x=(0.75,0.25)  x=(0.75,0.5)  x=(0.75,0.75)  x=(0.75,1)

x=(1,0)  x=(1,0.25)  x=(1,0.5)  x=(1,0.75)  x=(1,1)

RBF Kernel with σ=0.2

Kordon, Smits & Kotanchek

GECCO 2007                                                                    41

## Interpolation/Extrapolation of Local Kernel

- Small widths of kernel interpolate better
- Outside input range, no local information is available and the kernel levels off – no extrapolation
- No single choice of width achieves both



Kordon, Smits & Kotanchek

GECCO 2007                                                                    42

## Global Kernel



x=(0,0)  x=(0,0.25)  x=(0,0.5)  x=(0,0.75)  x=(0,1)

x=(0.25,0)  x=(0.25,0.25)  x=(0.25,0.5)  x=(0.25,0.75)  x=(0.25,1)

x=(0.5,0)  x=(0.5,0.25)  x=(0.5,0.5)  x=(0.5,0.75)  x=(0.5,1)

x=(0.75,0)  x=(0.75,0.25)  x=(0.75,0.5)  x=(0.75,0.75)  x=(0.75,1)

x=(1,0)  x=(1,0.25)  x=(1,0.5)  x=(1,0.75)  x=(1,1)

Polynomial Kernel with degree=2

Kordon, Smits & Kotanchek

GECCO 2007                                                                    43

## Interpolation/Extrapolation of Global Kernel

- Lower order polynomials extrapolate better
- High order polynomials needed to interpolate
- No single choice of order achieves both



Kordon, Smits & Kotanchek

GECCO 2007                                                                    44

## Mix of Local and Global Kernel



Kordon, Smits & Kotanchek

## Interpolation/Extrapolation with Mixed kernels

- Mixture of first degree polynomial and RBF with $\sigma=0.01$
- RBF contribution makes interpolation possible
- Polynomial makes extrapolation possible
- Single choice of parameters achieves both



Kordon, Smits & Kotanchek

## Industrial Example: Polynomial Kernel



Kordon, Smits & Kotanchek

## Industrial Example: RBF Kernel



Kordon, Smits & Kotanchek

## Industrial Example: Mixed Kernel



ε–Insensitive SVM (Learning Set)

(ε = 5.000)

ε–Insensitive SVM (Test Set)

Kordon, Smits & Kotanchek

---

## Integrated Methodology for Empirical Models Development



Original spreadsheet 50 variables X 1000 data points → Full data set

Nonlinear sensitivity analysis
Time delay influence
**Analytical Neural Networks**

Reduced spreadsheet 5 variables X 1000 data points → Reduced inputs data

Outliers detection
Condensed data set generation
**Support Vector Machines**

Reduced spreadsheet 5 variables X 120 data points → Condensed data

Symbolic regression
Functional solutions selection
**Genetic Programming**

Final model Analytical function

$Y = f(x)$ → Selected on Pareto front

- Hybrid approach integrating multiple technologies exploits the strengths of each
- Advantages:
  - Fast development (days)
  - Robust performance (compact models)
  - Direct implementation in any Distributed Control System (no need for specialized software)
  - Very low capital cost (only if hardware for data collection is unavailable)
  - Low average cost of ownership (reduced development and maintenance cost)
  - Process engineers like it (preferable to black-box models)

Kordon, Smits & Kotanchek

---

## Genetic Programming

**Genome Tree Plots**



Parents

Children

**Example of Crossover Operation**

**Phenotypes (Expressions)**

Parents
$$-(-0.787701)^{x} + x \qquad \frac{y^2 x}{-x+y}$$

Children
$$-(-0.787701)^{y^2 x} + x \qquad \frac{x}{-x+y}$$

- Based on artificial evolution of millions of potential nonlinear functions => **survival of the fittest**
- **Many possible solutions** with different levels of complexity
- The final result is an **explicit (nonlinear) function**
- *Can* have better **generalization capabilities** than neural nets
- Low **implementation** requirements
- Issues include …
  - Time delays
  - Sensitivity analysis of large data sets
  - Relatively slow development (hours of computation time)

Kordon, Smits & Kotanchek

---

## Steps Based on Genetic Programming



Representative data collection

Data preprocessing and classification

Sensitivity analysis of all inputs

Convolution parameters' estimation

Outlier detection and data set condensation

GP function generation

Analytical function selection/verification

On-line implementation

Model maintenance

Kordon, Smits & Kotanchek

## Classic Problems with Genetic Programming

- 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 (e.g., neural nets) outperform in raw performance
- **Bloat**
  - Parsimony control requires user intervention and is problem dependent
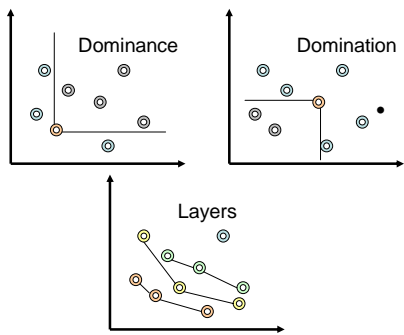
Kordon, Smits & Kotanchek

GECCO 2007　　　　　　　　　　　　　　　　　　53

## The Pareto Front



- Identifies trade-off surface between competing objectives
  - e.g., performance vs. complexity
- Pareto front solutions are the best "bang-for-the-buck"
- Introns are punished automatically
- How can we exploit?

Note that much evolutionary effort is spent exploring high complexity & high fitness regions

GECCO 2007　　　　　　　　　　　　　　　　　　54

## Pareto Performance



Dominance

Domination

Layers

- Characterizing Pareto Performance
  - Dominance
  - Domination
  - Layer
  - Combinations …
- Computational Issues
  - Brute force is $M N^2$
  - Can do $M N \log_{M-1}(N)$ or $M N \log_{M-2}(N)$ if clever
    - $M$ = # of objectives
    - $N$ = population size
  - Computation demands need to be considered in algorithm design

Kordon, Smits & Kotanchek

GECCO 2007　　　　　　　　　　　　　　　　　　55

## Genome Complexity



- What is complexity?
  - # of nodes?
  - Tree depth?
  - Included functions?
  - Number of variables?
  - Combinations?
- Chosen function is sum of sum of node counts
  - Provides more resolution at low end of complexity than simply using node count
  - Rewards fewer layers
- Real goal is to characterize the (relative) "smoothness" of the evolved function

Complexity = 36

$$\frac{1}{x} - 27\, x$$

Complexity = 17

Kordon, Smits & Kotanchek

GECCO 2007　　　　　　　　　　　　　　　　　　56

3310

# ParetoGP Algorithm



- Maintain archive based upon Pareto layers
- Each child results from one archive and one population parent
- Cascades …
  - Pareto archive maintained
  - Population wiped out (fresh genes!)
- Independent runs with independent archives for diversity
- This approach is intrinsically Pareto-aware

Kordon, Smits & Kotanchek

GECCO 2007                                                      57

# ClassicGP Algorithm



- ClassicGP can be Pareto-aware if a Pareto-aware selection scheme is used
- Most Pareto selection schemes are slow
- Finding the Pareto front can be relatively efficient
- Pareto Elite or Pareto Tourney may be viable selection schemes
  - Pareto tourney: select Pareto fronts from random subpopulations until desired number of models is reached
  - Pareto elite: select randomly from elite (defined using Pareto layers)

Kordon, Smits & Kotanchek

GECCO 2007                                                      58

# Symbolic Regression via GP

**Nuances…**

```
GenomeTreePlot[{parents,
  MutateSubtree[parents,
    MaximumTreeDepth → 3,
    MaximumArity → 2,
    DataVariables → {x, y}],
  Crossover[parents]}];
```



*Parent*

*Mutant*

*Child*

Introns are either overly complex or non-functional

choice of operators
- functional building blocks

parsimony pressure
- preference for simpler/smaller solutions

diversity operators
- modify fit solutions and the relative presence of each mechanism

fitness-based breeding rights
- proportional, ranking, elitist, tournament, random, etc.

evolution environment
- population size, number of generations, population interaction, fitness criteria, etc.

genetic modifications
- coefficient & structure optimization

automatically defined functions
- dynamically determined building blocks

metasensor definitions
- dynamically determined transforms and variable combinations

Kordon, Smits & Kotanchek

GECCO 2007                                                      59

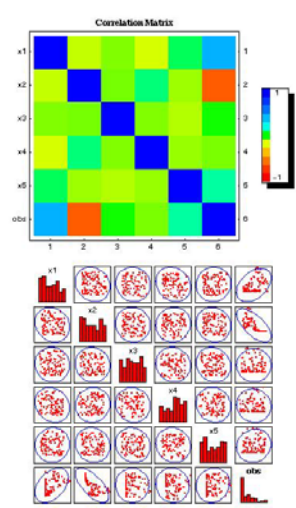# A Toy Problem for Illustration



- We sampled a function of two variables at 100 random points in the range [0,4]
- The data matrix has three random spurious variables in the range [0,4]
- Notice that the entire parameter space is not covered

$$\frac{e^{-(-1+b)^2}}{1.2 + (-2.5 + a)^2}$$

Kordon, Smits & Kotanchek

GECCO 2007                                                      60

3311

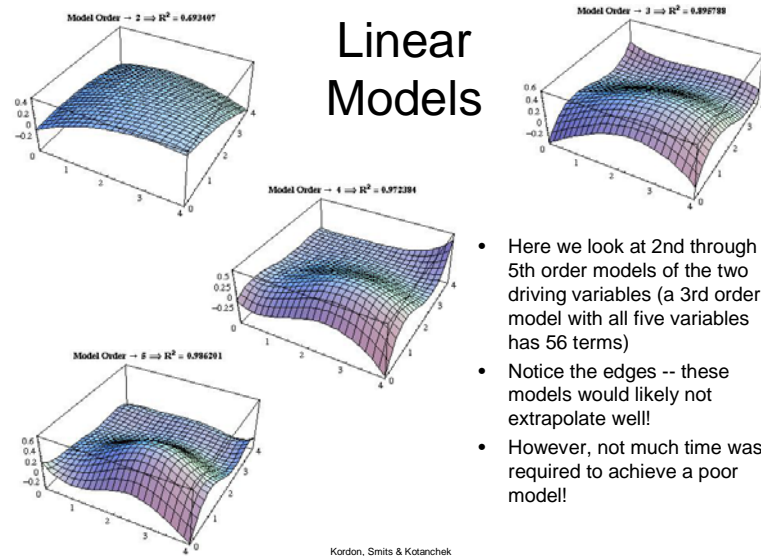## Slide 61



# Getting the Zen of the Data

- In this simple example, we could probably guess that only two variables were important for model building
- Correlated inputs can be a problem for some other modeling techniques
- However, lack of correlation to the response does not necessarily correspond to lack of importance

**Context-free analysis leads to confidently wrong answers!**
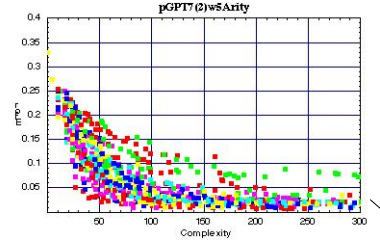
Kordon, Smits & Kotanchek

GECCO 2007    61

## Slide 62

# Linear Models



- Here we look at 2nd through 5th order models of the two driving variables (a 3rd order model with all five variables has 56 terms)
- Notice the edges -- these models would likely not extrapolate well!
- However, not much time was required to achieve a poor model!

Kordon, Smits & Kotanchek

GECCO 2007    62

## Slide 63

# The Pareto Front: Handling Competing Objectives

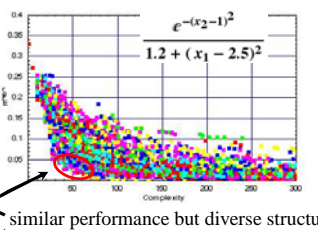No more things should be presumed to exist than are absolutely necessary — W. Occam [1280–1349]



- Identifies trade-off surface between competing objectives
  - e.g., performance vs. complexity
- Pareto front solutions are the best "bang-for-the-buck"
- Accuracy and simplicity are automatically rewarded
- Pareto Front Benefits
  - Avoids need for *a priori* combination of objectives into a single metric
  - The shape of the front gives us insight into the problem
  - Identifies multiple candidate solutions simultaneously

These are the error vs. complexity results of multiple independent symbolic regressions. Note that there is variability from run to run due to the random nature of the evolutionary process.

Kordon, Smits & Kotanchek

GECCO 2007    63

## Slide 64

# Evolved Models



- A run tends to fully explore a foundation structure
- Independent evolutions will result in different (but still fit) structures
- Cascading results from independent evolutions seems to be beneficial
- Note that we are not strictly restricted to the Pareto front in selecting models -- many models may be "good enough" and have the benefit of being structurally different and diverse

similar performance but diverse structure

Kordon, Smits & Kotanchek

GECCO 2007    64

3312

## Pareto Front Models

Truth

$$\frac{e^{-(x_2-1)^2}}{1.2 + (x_1 - 2.5)^2}$$

Explicit model complexity vs. accuracy control

## Parsimony & Extrapolation

Truth

- Note the pathologies at high complexity when extrapolating
- In general, we want to avoid over-modeling!

# 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
- Metavariables can also be used as **inputs into other modeling tools**

**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 high-dimensional data

**Rapid Modeling**
- Exploitation of the Pareto front has resulted in several orders-of-magnitude in the symbolic regression **performance** relative to more traditional GP. This greatly increases the range of possible applications.

There are many benefits to symbolic regression. These are enhanced when coupled with other analysis tools and techniques.

Kordon, Smits & Kotanchek

## Particle swarm optimization

Global best

Local best

Local best

An efficient technique to find the global optimum for model inversion and non-linear parameter estimation

At each time step $t$
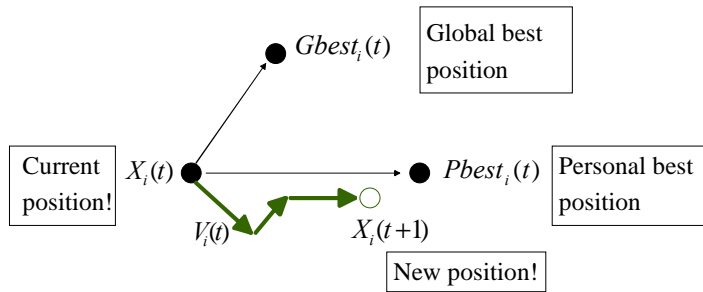
For each particle $i$

Update the position change (velocity)

$$V_i(t+1) = \chi \cdot (V_i(t) + c_1 \cdot rand(0,1) \cdot (P_i(t) - X_i(t))$$
$$+ c_2 \cdot rand(0,1) \cdot (P_g(t) - X_i(t))$$

Then move    $X_i(t+1) = X_i(t) + V_i(t+1)$

Note: - stochastic component
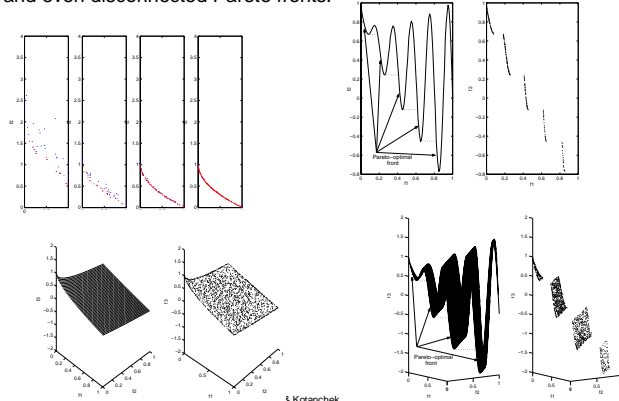- parameters $c_1, c_2, \chi$ default values (2.05, 2.05, 0.73)

Kordon, Smits & Kotanchek

## Particle's Movement – A Compromise

$Gbest_i(t)$ — Global best position

Current position! — $X_i(t)$

$Pbest_i(t)$ — Personal best position

$V_i(t)$   $X_i(t+1)$

New position!

## Multi-Objective PSO

Efficient technique to determine the Pareto front for problems with convex, non-convex and even disconnected Pareto fronts.

## Software tools

**DAP, Cave, IP21**

Representative data collection

**Excel, JMP, SIMCA, Mathematica**

Data preprocessing and classification

**MATLAB, Excel**

NN sensitivity analysis of all inputs

Convolution parameters' estimation

**MATLAB Toolbox**

Outlier detection and data set condensation

GP function generation

**MATLAB &MATHEMATICA Toolboxes**

Analytical function selection/verification

**G2, MOD, IP21, WebMathematica**

On-line implementation

Model maintenance

**DAP, Cave, IP21**

## Case Study: Inferential Sensors

**Key objective:**
To predict difficult-to-measure parameter (melt index) from easy-to-measure data (temperature, pressure, flow, etc.)

$$\rho C_v \frac{\partial T}{\partial t} + \nabla \cdot \left( -k\nabla T + \rho C_v T\mathbf{u} \right) = Q$$

Process

Process input → Process → Process Quality

Lab-test

Quality Prediction

**Training data**

**Inferential Sensors Development Software**

**Simple formulas**

$$y = a + b \cdot \left( e \cdot \left| \frac{x3}{x5} - d \right| \right)^c$$

**Easy On-Line implementation**

**Inferential Sensor**
An empirical model based on analytical equations with built-in self-assessment capability

3314

## Issues with neural net-based inferential sensors

**Issues with existing neural net-based inferential sensors:**
- High sensitivity to process changes
- Frequent re-training
- Complicated development & maintenance
- Low survival rate after 3 years in operation
- Engineers hate black-boxes

Black box

Analytical expression

VS.

$$Func2 = \frac{rate^2\left(vac + \frac{rate \cdot hopp \, wt}{temp}\right) pllt \, wt \, temp}{density \cdot temp^2}$$

Specialized run-time software

Directly coded into most on-line systems

Kordon, Smits & Kotanchek

GECCO 2007                                    73

---

## Inferential sensor for emission monitoring: A case study
## Data Collection

251 training data points

107 test data points (~40% outside training range)

143.0 ppm

Emission variable

Chemical Process        8 inputs        Design Of Experiments

Kordon, Smits & Kotanchek

GECCO 2007                                    74

---

## Inferential sensor for emission monitoring: A case study
## Sensitivity analysis by SANN

Input x3 removed after first sequence

A NN with 4 inputs: x2, x5, x6, and x8 is selected after discussion with process engineers

Input x7 removed after second sequence

Input x6 has the strongest sensitivity

Kordon, Smits & Kotanchek

GECCO 2007                                    75

---

## Inferential sensor for emission monitoring: A case study (SANN model performance)

Measured emission variable

Bad extrapolation (test data is 40% outside the range of training data)

Predicteded emission variable

Model based on 30 stacked NN with 10 neurons in hidden layer

Reduced number of inputs from 8 to 4

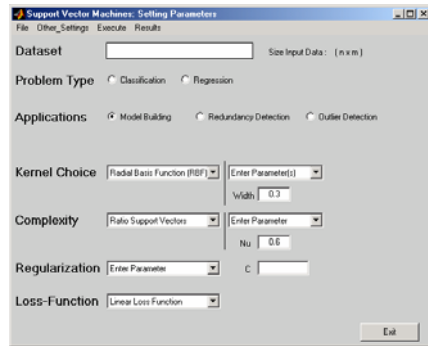**Fast test of the hypothesis about potential nonlinear relationship (in 20-30 min)**

Kordon, Smits & Kotanchek

GECCO 2007                                    76

## Slide 77

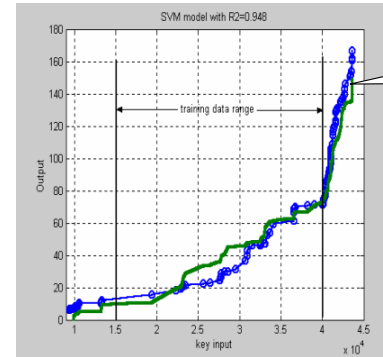### Inferential sensor for emission monitoring: A case study (SVM parameters)



Parameters:
% support vectors: 10
$C = 10^6$
Mixed Kernels: Polynomial and RBF
Range of Polynomial kernels: 1-3
Range of RBF kernel: 0.25-0.75
Range of ratio 0.5 – 0.99

Kordon, Smits & Kotanchek

GECCO 2007 — 77

## Slide 78

### Inferential sensor for emission monitoring: A case study (SVM model performance)



Impressive extrapolation (test data is 40% outside the range of training data)

Model based on a mixture of 2nd order polynomial global kernel and RBF local kernel with width of 0.5 and ratio of 0.95

Reduced number of training data points from 251 to 34 (based on support vectors)

Kordon, Smits & Kotanchek

GECCO 2007 — 78

## Slide 79

### Inferential sensor for emission monitoring: A case study (GP parameters)



Parameters for a GP simulated evolution

| | |
|---|---|
| Reference data | :34 |
| Random subset selection [%] | :100 |
| Number of runs | :20 |
| Population size | :500 |
| Number of generations | :100 |
| Probability for function as next node | :0.6 |
| Optimization function | :Corr. |
| Parsimony pressure | :0.1 |
| Prob. for random vs guided crossover | :0.5 |
| Probability for mutation of terminals | :0.3 |
| Probability for mutation of functions | :0.3 |

Kordon, Smits & Kotanchek

GECCO 2007 — 79

## Slide 80

### Inferential sensor for emission monitoring: A case study (Selected symbolic regression model)



Simple expression with acceptable performance (R2 = 0.87)

Response surface of model according to process physics

Selected model on Pareto front

Kordon, Smits & Kotanchek

GECCO 2007 — 80

3316

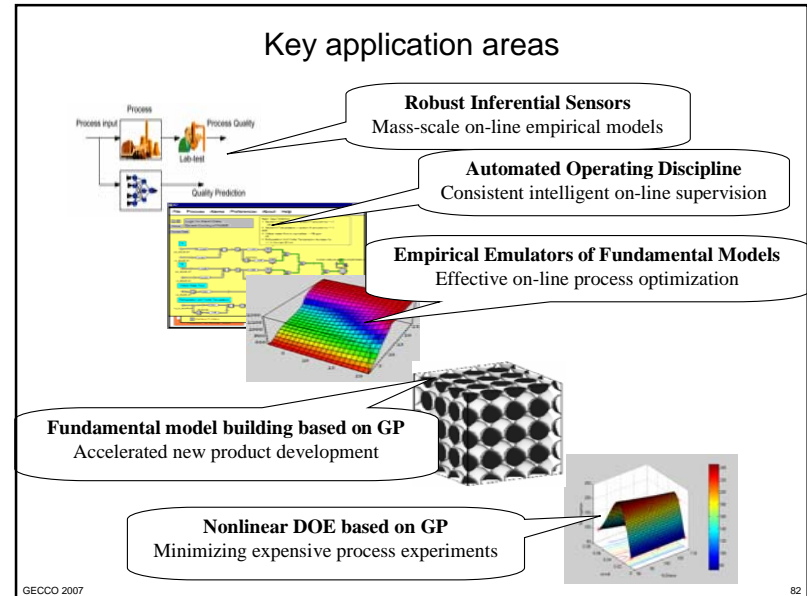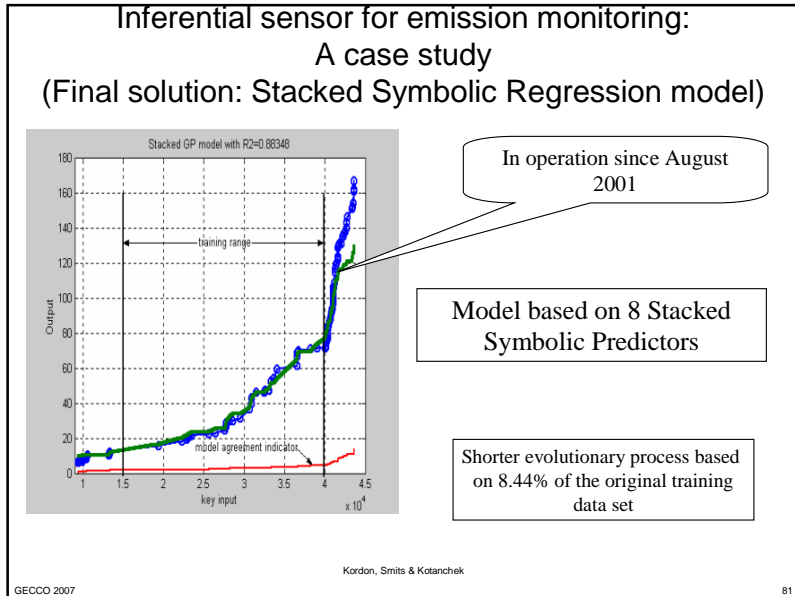## Inferential sensor for emission monitoring: A case study
### (Final solution: Stacked Symbolic Regression model)



Stacked GP model with R2=0.86348

In operation since August 2001

Model based on 8 Stacked Symbolic Predictors

Shorter evolutionary process based on 8.44% of the original training data set

Kordon, Smits & Kotanchek

GECCO 2007 — 81

## Key application areas



**Robust Inferential Sensors**
Mass-scale on-line empirical models

**Automated Operating Discipline**
Consistent intelligent on-line supervision

**Empirical Emulators of Fundamental Models**
Effective on-line process optimization

**Fundamental model building based on GP**
Accelerated new product development

**Nonlinear DOE based on GP**
Minimizing expensive process experiments

GECCO 2007 — 82

## EC Applications in Dow Chemical

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

Kordon, Smits & Kotanchek

GECCO 2007 — 83

## Automating Operating Discipline



- Heuristic rules defined verbally by process engineers/operators
- holdup predictor designed by stacked analytic NN and GP
- all decision blocks have fuzzy thresholds defined by membership functions
- simple empirical models and mass balances
- fundamental model predictions are used in the heuristic rules

- reduced major shutdowns
- reduced lab sampling

Kordon, Smits & Kotanchek

GECCO 2007 — 84

## Slide 85: Emulator for optimization of an industrial chemical process

Four levels DOE

10 inputs → Reactor Model 20-25 min/prediction → 12 outputs → Training Data set → Symbolic Regression Emulator 5 ms/prediction → **On-line process optimization**

Test Data set

Kordon, Smits & Kotanchek

GECCO 2007 — 85

## Slide 86: Fundamental Model Building Based on GP

1. Problem definition
2. Run symbolic regression
3. Identify key factors&transforms
4. Select GP generated models
5. Construct first principle models
6. Select&verify the final model solution
7. Validate the model

GP

Accelerated fundamental model building steps

Run simulated evolution before beginning fundamental modeling

GPfunction1

$y = a + b \left( e \cdot \left| \frac{x_3}{x_4} - d \right| \right)^c$

$S_k = \frac{3.13868 \times 10^{-17} e^{\sqrt{x_1}} \ln(x_3)^2 x_2}{x_4} + 1.00545$

Virtual modelers

**The evolutionary process identifies the key input variables as well as natural groupings & relationships. Combining this with a domain knowledge and first-principles insights is very powerful.**

Kordon, Smits & Kotanchek

GECCO 2007 — 86

## Slide 87: Approaches to accelerate fundamental model building process

AI approach

Reduce hypothesis search by GP

GP as automated invention machine

inside pentium PROCESSOR

out

Mimic the expert

Eliminate the expert

Maximize creativity of the expert

Kordon, Smits & Kotanchek

GECCO 2007 — 87

## Slide 88: The problem of structure-properties in fundamental modeling

Material structure

Properties:
- molecular weight
- particle size
- crystallinity
- volume fraction
- material morphology
- etc.

Key modeling effort for new product development

Modeling issues:
• nonlinear interaction
• large number of preliminary expensive experiments required
• large number of possible mechanisms
• slow fundamental model building
• insufficient data for training neural nets

Kordon, Smits & Kotanchek

GECCO 2007 — 88

3318

## Case Study with Structure-Property Relationships

**Fundamental Model Building**

Theoretical Analysis

$$\frac{dT}{dt} = a\frac{\partial^2 T}{\partial z^2} - \frac{DH}{C_p r}\frac{dc}{dt}$$

Hypothesis Search

Fundamental model

$$y = a + [b\,x_1 + c\,\log(x_2)]\,e^{kx_1} + d\,x_3$$

3 months

Structure-property data sets

**Fundamental Model Building + Symbolic Regression = Accelerated New Product Development**

**Symbolic Regression**

Sensitivity Analysis

Simulated Evolution

Symbolic Regression Model

$$y = a + b\sqrt{\frac{-x_3}{e^{\sqrt{\log(x_1,x_2)^2}}}} + \sqrt{x_1} + x_3 \quad (2)$$

10 hours

---

## GP and Design Of Experiments (DOE)
## Models Showing Lack of Fit

**Situations of Lack of Fit**

**1. Simple factorial DOE**
Enough experiments to fit first order model

$$y = \beta_o + \sum_{i=1}^{k}\beta_i x_i + \sum\sum_{i<j}\beta_{ij}x_i x_j$$

**Classical approach if LOF add experiments to fit second order model**

$$S_k = \beta_o + \sum_{i=1}^{k}\beta_i x_i + \sum\beta_{ii}x_i^2 + \sum\sum_{i<j}\beta_{ij}x_i x_j$$

More costly experiments

**2. A response surface DOE**
already had all experiments to fit second order model

$$S_k = \beta_o + \sum_{i=1}^{k}\beta_i x_i + \sum\beta_{ii}x_i^2 + \sum\sum_{i<j}\beta_{ij}x_i x_j$$

**Classical approach if LOF no alternative (use model as it is)**

**Suggested approach:
Use GP to transform inputs**

Kordon, Smits & Kotanchek

---

**1. Generate GP models**

$$S_k = \frac{3.13868\times10^{-17} e^{\sqrt{2x_1}}\ln\left[(x_3)^2\right]x_2}{x_4} + 1.00545 \quad (2)$$

Selected solution

**2. Generate input transforms**

**Variable transformations suggested by GP model**

| Original Variable | Transformed Variable |
|---|---|
| $x_1$ | $Z_1 = \exp\left(\sqrt{2x_1}\right)$ |
| $x_2$ | $Z_2 = x_2$ |
| $x_3$ | $Z_3 = \ln\left[(x_3)^2\right]$ |
| $x_4$ | $Z_4 = x_4^{-1}$ |

**3. Fit response surface model in transformed variables**

$$S_k = \beta_o + \sum_{i=1}^{4}\beta_i Z_i + \sum\sum_{i<j}\beta_{ij}Z_i Z_j + \sum_{i=1}^{4}\beta_{ii}Z_i^2$$

| Source | DF | Sum of Square | Mean Square | F Ratio |
|---|---|---|---|---|
| Lack of Fit | 2 | 0.00049190 | 0.000246 | 2.2554 |
| Pure Error | 2 | 0.00021810 | 0.000109 | **Prob > F** |
| Total Error | 2 | 0.00071000 | | 0.3072 |
| | | | | **Max RSq** |
| | | | | 0.9999 |

No Lack Of Fit (p=0.3037)

Kordon, Smits & Kotanchek

---

## PSO application: Optimizing color spectrum of plastics

Real-time optimization in 2-3 seconds

PSO and GA convergence

Multiple-objective PSO with 15 variables

ColourPro Formulation Optimization

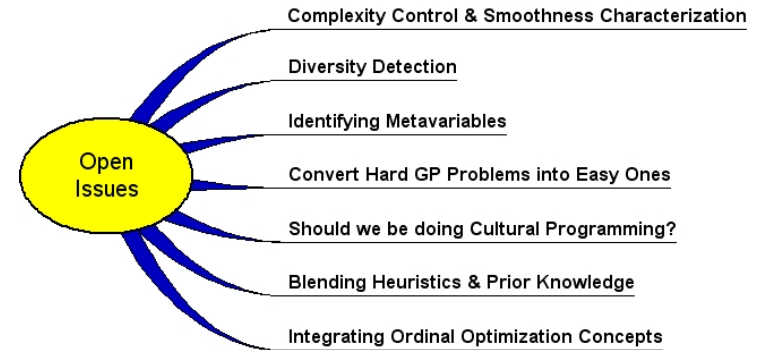Kordon, Smits & Kotanchek

## Other PSO applications

- Drug release predictor
  - 6 parameters
  - population size = 30
  - optimization time: ~ 30 seconds
- Foam acoustics performance predictor
  - 8 parameters
  - population size = 50
  - optimization time: ~ 5 seconds
- Crystallization kinetics predictor
  - 4 parameters
  - population size = 30
  - optimization time: ~ 2 seconds

Kordon, Smits & Kotanchek

GECCO 2007                                                                 93

## Open Issues & Current Research



Open Issues

Complexity Control & Smoothness Characterization

Diversity Detection

Identifying Metavariables

Convert Hard GP Problems into Easy Ones

Should we be doing Cultural Programming?

Blending Heuristics & Prior Knowledge

Integrating Ordinal Optimization Concepts

Kordon, Smits & Kotanchek

GECCO 2007                                                                 94

## Summary

- Evolutionary Computing can create significant value to industry by reducing model development time and model exploitation cost
- Integrating EC with Neural Networks, Support Vector Machines, and Statistics is recommended for successful industrial applications
- This strategy works for many real applications in the chemical industry
- The key application areas are:
  - Inferential sensors
  - Improved process monitoring and control
  - Accelerated new product development
  - Effective design of experiments
- And this is only the beginning …

Kordon, Smits & Kotanchek

GECCO 2007                                                                 95

## Acknowledgement

We would like to acknowledge the contribution of the following researchers from The Dow Chemical Company:

Alex Kalos

Kip Mercure

Flor Castillo

Elsa Jordaan

Leo Chiang

Irina Graf

Katya Vladislavleva – Tilburg University

Kordon, Smits & Kotanchek

GECCO 2007                                                                 96

## References

1. M. Kotanchek, G. Smits, and A. Kordon, *Industrial Strength Genetic Programming*, In *GP Theory and Practice* (R. Riolo and B. Worzel-Eds), Kluwer, 2003.
2. A. Kordon, G. Smits, A. Kalos, and E. Jordaan, *Robust Soft Sensor Development Using Genetic Programming*, In *Nature-Inspired Methods in Chemometrics*, (R. Leardi-Editor), Elsevier, 2003.
3. Kordon A.K, G.F. Smits, E. Jordaan and E. Rightor, *Robust Soft Sensors Based on Integration of Genetic Programming, Analytical Neural Networks, and Support Vector Machines*, Proceedings of WCCI 2002, Honolulu, pp. 896 – 901, 2002.
4. Kotanchek M., A. Kordon, G. Smits, F. Castillo, R. Pell, M.B. Seasholtz, L. Chiang, P. Margl, P.K. Mercure, A. Kalos, *Evolutionary Computing in Dow Chemical*, Proceedings of GECCO'2002, New York, volume Evolutionary Computation in Industry, pp. 101-110., 2002
5. Kordon A. K., H.T. Pham, C.P. Bosnyak, M.E. Kotanchek, and G. F. Smits, *Accelerating Industrial Fundamental Model Building with Symbolic Regression: A Case Study with Structure – Property Relationships*, Proceedings of GECCO'2002, New York, volume Evolutionary Computation in Industry, pp. 111-116, 2002

6. Castillo F., K. Marshall, J. Greens, and A. Kordon, *Symbolic Regression in Design of Experiments: A Case Study with Linearizing Transformations*, Proceedings of GECCO'2002, New York, pp. 1043-1048.
7. Kordon A., E. Jordaan, L. Chew, G. Smits, T. Bruck, K. Haney, and A. Jenings, *Biomass Inferential Sensor Based on Ensemble of Models Generated by Genetic Programming*, Proceedings of GECCO 2004, Seattle, WA, pp. 1078-1089, 2004
8. Smits G. and M. Kotanchek, *Pareto-Front Exploitation in Symbolic Regression*, In *GP Theory and Practice* (R. Riolo and B. Worzel-Eds), Kluwer, 2004.
9. Kordon A., A. Kalos, and B. Adams, *Empirical Emulators for Process Monitoring and Optimization*, Proceedings of the IEEE 11th Conference on Control and Automation MED'2003, Rhodes, Greece, pp.111, 2003.
10. Kordon A. and CT Lue, *Symbolic Regression Modeling of Blown Film Process Effects*, Proceedings of CEC 2004, Portland, OR, pp. 561-568, 2004.
11. Jordaan, E., A. Kordon, G. Smits, and L. Chiang, *Robust Inferential Sensors based on Ensemble of predictors generated by Genetic Programming*, Proceedings of PPSN 2004, Birmingham, UK, pp. 522-531, 2004

## References

12. Kordon A., E. Jordaan, F. Castillo, A. Kalos, G. Smits, and M. Kotanchek, *Competitive Advantages of Evolutionary Computation for Industrial Applications*, Proceedings of CEC 2005, Edinburgh, UK, pp. 166-173, 2005
13. F. Castillo, A. Kordon, J. Sweeney, and W. Zirk, *Using Genetic Programming in Industrial Statistical Model Building*, In: O'Raily U. , Yu T., Riolo, R. and Worzel, B. (eds): *Genetic Programming Theory and Practice II*. Springer, NY, New York, pp. 31 – 48, 2004.
14. A. Kordon, F. Castillo, G. Smits, and M. Kotanchek, *Application Issues of Genetic Programming in Industry*, In: Yu T., Riolo, R. and Worzel, B. (eds): *Genetic Programming Theory and Practice III*. Springer, NY, New York, pp. 241 - 258, 2006.
15. G. Smits, A. Kordon, E. Jordaan . C. Vladislavleva, and M. Kotanchek, *Variable Selection in Industrial Data Sets Using Pareto Genetic Programming*, In: Yu T., Riolo, R. and Worzel, B. (eds): *Genetic Programming Theory and Practice III*. Springer, NY, New York, pp. 79 - 92, 2006.

16. F. Castillo, A. Kordon, G. Smits, B. Christenson, and D. Dickerson, *Pareto Front Genetic Programming Parameter Selection Based on Design of Experiments and Industrial Data*, Proceedings of GECCO 2006, Seattle, WA, pp. 1613-1620, 2006.
17. Kordon A., G. Smits, E. Jordaan, A. Kalos, and L. Chiang, *Empirical Models with Self-Assessment Capabilities for On-Line Industrial Applications*, Proceedings of CEC 2006, Vancouver, pp. 10463-10470, 2006.
18. G. Smits, and K. Vladislavleva, *Ordinal Pareto Genetic Programming*, Proceedings of CEC 2006, Vancouver, pp. 10471-10477, 2006.
19. F. Castillo, A. Kordon, and G. Smits, *Robust Pareto Front Genetic Programming Parameter Selection Based on Design of Experiments and Industrial Data*, In: Riolo, R. and Worzel, B. (eds): *Genetic Programming Theory and Practice IV*, Springer, NY, New York, pp. 149 - 166, 2007.
20. M. Kotanchek, G. Smits, and K. Vladislavleva, *Pursuing the Pareto Paradigm: Tournaments, Algorithm Variations and Ordinal Optimization*, In: Riolo, R. and Worzel, B. (eds): *Genetic Programming Theory and Practice IV*, Springer, NY, New York, pp. 167 - 186, 2007.

3321