#### Simulated Annealing Model Search (SAMS) for Subset Selection in Screening Experiments



Mark Wolters Statistical and Actuarial Sciences University of Western Ontario



Derek Bingham Statistics and Actuarial Science Simon Fraser University

ASQ Fall Technical Conference

October 14, 2011

# Outline

- Screening Experiments and Nonregular Designs
- Example: Plackett-Burman designs
- Model Selection from Ensembles of Oversized Models
  - Raster plot
  - Link plot
- SAMS for Generating Model Sets
- Demonstrations
- Conclusions

#### **Screening Experiments and Nonregular Designs**

# **Screening Experiments**

- Goal of screening: select the few important predictors from the many available ones.
  - Small number of runs
  - Many variables
  - Not very interested in inference or prediction
- Designs of interest: nonregular factorial designs.

  - nonregular designs **–** Exhibit **complex aliasing**.
- Benefit: can consider main effects and interactions.
- Cost: huge model set, model selection problem.

## **Principles of Analysis**

- Effect sparsity: Only a few effects are important.
  - Justifies the use of a screening experiment (few runs).
- Effect heredity: If interactions are active, then main effects should be active too.
  - Different ways to specify.
  - Heredity specification defines "interpretable" models.
  - Usual default: weak heredity. AB can only be active if A or B are active too.
- Effect hierarchy: Higher-order effects are less likely to be active than lower-order ones.

#### Model Selection is Difficult in Screening Experiments

- Key aspects of problem:
  - Small n, large k
  - Huge model set
  - Complex aliasing
  - Heredity requirement
- Why is model selection difficult in this situation?
  - High model selection uncertainty (aka model aliasing).
  - Exhaustive search may not be possible.
  - Strong tendency toward overfitting.
  - Common search methods don't respect heredity.
- Opinion: model aliasing is not given enough attention in this context. Usually multiple competing models should be identified.

**Example: Plackett-Burman Designs** 

#### **Plackett-Burman Designs**

- Design matrix is orthogonal array
- Add 2-way interactions, lose orthogonality
- Smallest cases: 12-run and 20-run
  - PB<sub>12</sub>: 11 main effects, 55 interactions
  - PB<sub>20</sub>: 19 main effects, 171 interactions

k = 66 variables

k = 190 variables

	1	A	в	С	D	Е	F	G	н	I	J	к	ΑB	AC	A D	AE	 IJ	IK	JK		
	[1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	 1	1	1		
	1	-1	1	-1	1	1	1	-1	-1	-1	1	-1	-1	1	-1	-1	 -1	1	-1		
	1	-1	-1	1	-1	1	1	1	-1	-1	-1	1	1	-1	1	-1	 1	-1	-1		
	1	1	-1	-1	1	-1	1	1	1	-1	-1	-1	-1	-1	1	-1	 1	1	1		
	1	-1	1	-1	-1	1	-1	1	1	1	-1	-1	-1	1	1	-1	 -1	-1	1		
7 –	1	-1	-1	1	-1	-1	1	-1	1	1	1	-1	1	-1	1	1	 1	-1	-1		
	1	-1	-1	-1	1	-1	-1	1	-1	1	1	1	1	1	-1	1	 1	1	1	> 12	runs
	1	1	-1	-1	-1	1	-1	-1	1	-1	1	1	-1	-1	-1	1	 -1	-1	1		
	1	1	1	-1	-1	-1	1	-1	-1	1	-1	1	1	-1	-1	-1	 -1	1	-1		
	1	1	1	1	-1	-1	-1	1	-1	-1	1	-1	1	1	-1	-1	 -1	1	-1		
	1	-1	1	1	1	-1	-1	-1	1	-1	-1	1	-1	-1	-1	1	 1	-1	-1		
	1	1	-1	1	1	1	-1	-1	-1	1	-1	-1	-1	1	1	1	 -1	-1	1	J	
	-											)									
Design matrix X: 11 Main Effects 55 Interactions																					

#### **Plackett-Burman Designs**

• With interactions included, *finding active variables becomes a model selection problem*.

	PI	$B_{12}$	$PB_{20}$				
$\mathbf{p}$	All Models	Hereditary	All Models	Hereditary			
1	66	11	190	19			
2	2145	165	$1.796 { imes} 10^4$	513			
3	$4.576 \times 10^{4}$	1705	$1.125 { imes} 10^{6}$	9861			
4	$7.207 \times 10^{5}$	$1.551{ imes}10^4$	$5.260 \times 10^{7}$	$1.705{ imes}10^{5}$			
5	$8.937 \times 10^{6}$	$1.252{ imes}10^5$	$1.957{ imes}10^9$	$2.680 \times 10^{6}$			
6	$9.086 \times 10^{7}$	$9.026{ imes}10^5$	$6.033{ imes}10^{10}$	$3.857{ imes}10^7$			
7	$7.788 \times 10^{8}$	$5.894 \times 10^{6}$	$1.586{ imes}10^{12}$	$5.147 \times 10^{8}$			

## **Summary of Available Methods**

	Pros	Cons
		-No heredity
Stepwise and	-Easy to use	-Greedy searches
variants		-Not theoretically sound
		-Returns 1 model
		-No heredity
Model selection	-Optimality in some sense	-Multiple criteria
criteria		-Doesn't adjust for size of model set
Boyosian	-Handles heredity	-Sensitive to priors
Methods <sup>1</sup>	-Returns distribution of models	-Interpretation of output
		-Returns 1 model
Groupwise	-Handles heredity	-Need stopping rule
LARJ-		

**Model Selection from Ensembles of Oversized Models** 

#### **Small-Scale Example**

**<u>Goal</u>**: introduce the idea and the key graphs. Avoid the model search problem for now to illustrate ideas.

- PB<sub>12</sub> design
- Generate  $Y = 1 + 2C + 1.5CD + \epsilon$ ,  $\epsilon \sim N(0,1)$
- Do exhaustive search of models size p = 5 (125202 models)
- Sort all models by residual sum of squares (RSS).
- Expect most of the top models to contain (C, CD)
  - There are 2904 such models
- Plot coefficients of top models... can the truth be detected?

#### **Raster Plot**



# Link Plot



- -Each effect represented by a point.
- -Each pair of variables represented by a line.
- -Line color and thickness represents occurrence frequency in the goodmodel set.
- -Color and thickness scaled relative to most frequent pair.
- -More frequent variables have points larger and more spaced out.

**SAMS for Generating Model Sets** 

#### **Summary of the Method**

# Search for lots of good, oversized models, then extract most common variable combinations from them.

- 1. Choose maximum truth size,  $s_{max}$ , and large model size,  $p > s_{max}$ .
- 2. Use a search heuristic to get 5000-10000 well-fitting models of size p.
  - Use RSS as goodness measure.
- 3. Use plots to visualize the model set and extract good small models.

The search heuristic used is based on simulated annealing, with two important changes.

# **Simulated Annealing**



#### **SA Innovation 1: Hereditary Move**

#### Hereditary move

- 1. Randomly choose one variable in M<sub>old</sub>.
- 2. Drop the chosen variable.
- 3. Drop any other variables that no longer respect heredity.
- 4. Build the model back up to size p:
  - a) Make a list of admissible variables.
  - b) Randomly choose one and add to the model.
  - c) repeat a) and b) until size = p.

Move		M	odel	
1	A	В	C	BD
2	A	B	$\overline{AB}$	<u>BD</u>
3	A	С	AB	AE
4	A	С	AJ	AE
5	A	$\overline{AK}$	$\underline{AJ}$	$\underline{AE}$
	$\overline{D}$	H	J	DG

## **SA Innovation 2: Adaptive Temperature Control**

- Premise: reduce T on every accepted move, increase T on every rejected move
- Improving moves always accepted; this rule influences uphill moves.
  - Accept: set T :=  $\rho$ T (0 <  $\rho$  < 1)  $\Rightarrow$  makes it harder to accept bad moves
  - Reject: set T :=  $\alpha$ T ( $\rho < \alpha < 1$ ) makes it easier to accept bad moves
- Control by setting  $\rho$  and  $\kappa,$  where  $\rho/\alpha^{\kappa}$  = 1
  - Call  $\kappa$  the "search depth" approx.  $\kappa$  rejections per acceptance.
  - Larger  $\kappa$ , more thorough local search.
- In the code, also put a lower bound on P(accept). Call this  $P_{min}$ .
- Main idea: this control makes the search *non-convergent*. Keep generating models indefinitely.

#### **Simulated Annealing Model Search**

- · Call the modified algorithm SAMS.
- Run until n<sub>gen</sub> models are accepted.
- Default parameters:  $\rho = 0.95$ ,  $\kappa = 4$ ,  $P_{min}=0.01$ ,  $n_{gen}=10000$ .

- Performs well across wide range of problems.

- Characteristics:
  - Will generate good models indefinitely
  - May re-visit models
  - Fast (~20sec for  $n_{gen}$ = 10000 with PB<sub>20</sub>)

**Demonstrations** 

#### **DEMO 1: Another PB<sub>12</sub> Example**

- Choose true model  $Y = 0.5 + 1.1B + 1.5E + 0.9BG + 1EI + \epsilon$ .
  - Response 1:  $\varepsilon \sim N(0, 0.25) \Rightarrow$  Low model selection uncertainty
  - Response 2:  $\epsilon \sim N(0,0.75) \Rightarrow$  High model selection uncertainty
- Set  $\tau = 4$ , p = 7.
- Set parameters:  $\rho = 0.95$ ,  $\kappa = 4$ ,  $P_{min}=0.001$ ,  $n_{qen}=10000$
- Run SAMS and view raster and link plots.

#### **Low Model Selection Uncertainty**



to emphasize heredity relationships

#### **High Model Selection Uncertainty**



#### **Clustered Raster Plot**

 Perform K-means clustering on rows to make common models more visible:



## **DEMO 2: The Blood Glucose Study**

- Mixed 2- and 3-level design, 18 runs.
- Factor A: 2 levels.
- Factors B-H: 3 levels. Use linear & quadratic terms.
- Include all valid 2-way interactions: 113 factorial effects.
- More complicated heredity specification due to quadratic terms, e.g.
  - B<sup>2</sup>C has parents BC and B<sup>2</sup>

#### **DEMO 2: The Blood Glucose Study**



- SAMS results agree with previous studies
- Model selection uncertainty particularly clear with SAMS
- Only graphical analysis required.

#### **DEMO 3: The Ozone Data**

- Response: ozone concentration (**n = 330**).
- Predictors: eight meteorological variables.
- Include all interactions and squared terms (k =44).
- Observational data → this is a *regression problem*.
- Large n, could consider larger models.
- Can our methodology choose important variable combinations?

Variable	Quantity	Label
Y	ozone concentration (ppm)	_
$X_1$	500 millibar height (m)	$\operatorname{Ht}$
$X_2$	wind speed (mph)	Wd
$X_3$	relative humidity $(\%)$	$\mathbf{R}\mathbf{H}$
$X_4$	surface temperature (°F)	Т
$X_5$	inversion height (ft)	iHt
$X_6$	pressure gradient (mmHg)	Р
$X_7$	inversion temperature (°F)	iT
$X_8$	visibility (mi)	V

#### **DEMO 3: The Ozone Data**

- Run SAMS with 13-variable models.
- A particular 8-variable model stands out.



Conclusions

# Conclusions

- Summary
  - Developed a heredity-respecting, non-convergent search heuristic
  - Developed graphical displays for model selection
  - Can find good small models from an ensemble of good large ones
- Advantages
  - Good performance
  - Easy to use, graphical approach
  - Heredity built in
  - Reduces the connection between GOF and choice of model size
- The SAMS code is available as a supplement to the paper.

#### **Future Work**

- More investigation of regression problems
  - Effect of collinearity, variable scaling, etc. on performance.
- Handle a wider range of heredity specifications
  - Currently each effect can have any other effect(s) as parents.
  - No strong heredity
  - Can't handle grouped variables (if A is in model, B must be as well)
  - Can handle exclusions (if A is in model, B must not be)
- Data-driven choice of m (size of ensemble of good models)

**Supporting Slides** 

## **Auto-Extraction of Best-Guess Model**

- Problem: can't use graphs in a simulation study
- Steps:
  - Find the 5 most frequent models of each size 1,...,τ in the good model set (use *branch and bound* method)
  - 2. Calculate the entropy criterion for each.
    - ➡ For a model with a main effects and b interactions:

Occurrence frequency of M in good set = fPopulation count of models containing M  $= N_p(a, b)$ Total population size  $= N_p$ Population proportion models containing M  $= \pi = \frac{N_p(a, b)}{N_p}$ Entropy criterion:  $H(M) = f \log_2(\frac{f}{\pi}) + (1 - f) \log_2(\frac{1 - f}{1 - \pi})$ 

3. Take maximum-entropy model as best

**A Simulation Study** 

# **Simulation Setup**

- Compare performance of methods that return a single model.
- Methods:

**Oracle.** Assume true model known; include each effect only if hypothesis test is significant.

Modified stepwise selection. Method of Wu and Hamada, 2000

 $AIC_c$  criterion. Select best model of size  $s_{max}$  or smaller.

SAMS + entropy criterion.

- Test cases:
  - Use both  $PB_{12}$  and  $PB_{20}$  designs ( $s_{max} = 4$  and 6, resp.).
  - Generate 5000 true models with random active factors, random coefficients.
  - Set coefficients large enough to give the model some chance to be detected, but small enough to keep R<sup>2</sup> to realistic levels.

## Results

 Any selected model will be in one of five categories: The truth (T). All variables correctly included; no extras. Underfitted (U). Only correctly included variables, but some omitted. Overfitted (O). All true variables included, with some extras. Partial-truth (P). Some (not all) true variables, some extras. Wrong (W). None of the true variables included.

			$PB_{12}$ design		$PB_{20}$ design					
Method	$\mathcal{T}$	U	$\mathcal{O}$	${\cal P}$	$\mathcal{W}$	Т	U	O	${\cal P}$	$\mathcal{W}$
Oracle	62.8	37.2	0.0	0.0	0.0	61.4	38.6	0.0	0.0	0.0
Step	11.1	3.7	41.2	17.6	26.4	2.0	0.6	42.3	37.2	17.9
SAMS	43.3	16.2	15.8	15.0	9.7	35.7	30.4	10.0	17.6	6.3
$AIC_c$	7.2	0.7	53.8	25.5	12.8	0.5	0.0	52.7	41.0	5.8

#### Partitioning of selected models (%)

#### Results

- Alternatively, count the total number of errors made in a model selection:
  - Excluding a truly-active variable.
  - Including a spurious variable.

		<i>PB</i> <sub>12</sub>	design		PB <sub>20</sub> design				
Size	oracle	step	$AIC_c$	SAMS	oracle	step	$AIC_c$	SAMS	
1	0.01	1.91	2.84	0.69	0.00	4.13	5.01	0.70	
2	0.15	2.09	2.04	0.46	0.05	3.87	4.19	0.35	
3	0.72	3.54	2.69	1.60	0.22	4.47	3.83	0.73	
4	1.99	5.31	4.97	3.97	0.61	5.89	4.23	1.99	
5					1.31	7.81	6.17	3.96	
6					2.32	8.86	7.69	5.47	

#### Average number of errors made