



## 33 SAS PROC GLM and MIXED For Recovering Information in Augmented Designs

34

### 35 Abstract

36 The SAS GLM and MIXED procedures can be useful for experimenters desiring to analyze data from  
37 screening experiments using a member of the class of augmented experiment designs. Since application  
38 of the procedures is typically not straightforward for these designs, several programs of possible interest  
39 are described. We show how to recover interblocking and intervareity information when the blocking and  
40 varieties are random effects, how to arrange varietal responses in descending order, and a number of other  
41 options.

42

### 43 Abbreviations:

44 AED: augmented experiment design

45 ARCBD: augmented randomized complete block design

46 AIBD: augmented incomplete block design

47 ANOVA: analysis of variance

48 EBLUP: empirical best linear unbiased predictor

49 REML: restricted maximum likelihood

50

### 51 1. Introduction

52

53 The class of augmented designs was introduced by Federer (1956, 1961, 1991) as an alternative to the  
54 systematic arrangement of a single check variety in every  $k$ th plot. This latter method does not provide  
55 for an estimate of experimental error and is inefficient in that too much space is devoted to check plots  
56 (Yates, 1936). Also, the random nature of genotypes in the early stages of a selection program needs to be  
57 taken into account (Federer, 1996). Cullis *et al.* (1989) provide an analysis for the systematic check  
58 arrangement and consider the genotypes as random effects, but the over-abundance of check plots and the  
59 single check deficiencies remain. An augmented experiment design (AED) is obtained by selecting an

60 experiment design for checks and enlarging the blocks (rows and columns) to accommodate new  
61 genotypes (treatments) which usually only appear once in an experiment. Our purpose here is to present  
62 PROC GLM (SAS Institute Inc., 1989) and PROC MIXED (SAS Institute Inc., 1996) programs for  
63 obtaining the analysis for this class of designs in much the same manner as described by Federer and  
64 Wolfinger (1996). Describing codes for various statistical analyses make procedures readily available to  
65 an experimenter in order that more efficient statistical analyses and better use of resources may be made.  
66 Also, some other related program ideas are described.

67

## 68 2. Materials and Methods

69

70 The SAS software package with the PROC GLM (general linear models) and PROC MIXED (mixed  
71 model of random and fixed effects) was used to develop codes or programs for recovering interblocking  
72 and intervariety or intergenotype information associated with the random variables in the experimnet.  
73 Such variables as complete blocks, incomplete blocks, rows columns, and/or genotypes quality for  
74 consideration as random effects. Programs in PROC GLM consider all variables as fixed effects as is  
75 done in regression analyses. In this procedure some variables may be designated as random for the  
76 purpose of obtaining expected values of mean squares only. PROC MIXED uses REML (restricted  
77 maximum likelihood) solutions for the variance components for random effects but other options are  
78 available. Ordinary textbook analyses use analysis of variance (ANOVA) solutions for the variance  
79 components used for recovering intereffect information. Another procedure, PROC IML (interactive  
80 matrix language), was used to generate orthogonal polynomial regression coefficients for statistical  
81 analyses.

82

83 To develop codes for specific analyses, the desired statistical analysis was determined. Then, a matrix  
84 software package such as GAUSS was used to obtain the numerical values for an example. The example  
85 was constructed from known fixed effect parameters for all variables, allowing for a check on the GAUSS  
86 solutions obtained. Knowing the numerical values desired, various SAS commands were investigated to

87 obtain the desired results. When the commands giving the desired results were obtained, the program for  
88 the analysis was finalized.

89

### 90 3. Results and Discussion

91

92 PROC GLM and PROC MIXED codes were developed for an augmented balanced incomplete block  
93 design using  $n = 6$  new genotypes and  $c = 4$  checks as an illustrative example. A balanced incomplete  
94 block design for  $c = 4$  checks in incomplete blocks of size two in  $r = 3$  complete blocks or replicates was  
95 augmented by including one new treatment in each of the six incomplete blocks. The  $n + c = v = 10$   
96 entries need to be divided into a set (new) which represents the random set and a set (checks) which  
97 represent the entries considered to be fixed effects. The outputs for the PROC GLM and PROC MIXED  
98 codes are given in the Appendix. These same procedures plus PROC IML are used to develop a program  
99 for the analysis of a 15 row by 12 column augmented design with  $n = 120$  new and  $c = 2$  checks replicated  
100 30 times each. Since this row-column design is not connected (i.e., not all row, column, and entry effects  
101 have solutions under the usual restrictions that effects of a variable sum to zero), orthogonal polynomial  
102 regression functions of row and columns were used in the analysis. In addition, owing to the nature of  
103 spatial variation, interactions of the row and column regressions were needed to account for the particular  
104 type of variation encountered.

105

#### 106 PROC GLM For AEDs With One-Way Blocking

107

108 With respect to an augmented randomized complete block (ARCBD) or an augmented incomplete block  
109 (AIBD) design, every blocked design is incomplete with respect to the new treatments. Hence, it is  
110 desirable to recover interblock information even for ARCBDs. However, for completeness, we begin by  
111 showing how to use SAS PROC GLM to obtain only intrablock and intravariety analyses for an AIBD as  
112 follows:

113

```
114 data augbibd;  
115     infile 'augbibd.dat';  
116     input yield rep block treat;  
117 proc glm data=augbibd;  
118     class rep block treat;  
119     model yield = rep block(rep) treat / solution;  
120     random rep block(rep);  
121     lsmeans treat;  
122 run;
```

123

124 The program starts with a SAS DATA step inputting a raw data file named "augbibd.dat" containing data  
125 from an AIBD with  $n$  new treatments and  $c$  check treatments. The input variables are YIELD (the  
126 response), REP (the replicate), BLOCK (the block), and TREAT (the treatment).

127

128 The DATA= option in the PROC GLM statement reads in the newly created SAS data set. The CLASS  
129 statement declares REP, BLOCK, and TREAT to be classification (qualitative) variables. The MODEL  
130 statement lists the dependent variable YIELD and the effects to be used in the analysis. Since REP,  
131 BLOCK, and TREAT are all classification variables, the effects involving them are constructed using 0-1  
132 indicator variables.

133

134 Owing to the nature of PROC GLM's model parameterization, we assume the levels of TREAT are 1 to  $n$   
135 for the  $n$  new genotypes and  $n + 1$ ,  $n + 2$ , ...,  $n + c$  for the  $c$  checks. The SOLUTION option in the  
136 MODEL statement prints out estimates of all of these levels, and since there is an overparameterization,  
137 PROC GLM sets the last treatment effect equal to zero. Therefore the new and the other check effects will  
138 all have the last treatment effect subtracted from them. The standard error listed with the solution is a  
139 standard error of a difference of the two effects, and the highest numbered check should be the one of most  
140 interest.

141

142 The RANDOM statement declares the REP and BLOCK(REP) effects to be random effects, leaving  
143 TREAT as the lone fixed effect in this analysis. In spite of being declared random, PROC GLM will still  
144 consider REP and BLOCK(REP) to be fixed effects during the model fit, but it will compute the expected  
145 mean squares for the replicate (complete block for checks) and blocks within replicates mean squares.  
146 Note that the sum of the new treatment effects are not required to sum to zero when using the constraint  
147 that the highest numbered treatment effect is equal to zero.

148

149 The LSMEANS statement computes estimated population marginal means for TREAT with equal weights  
150 applied to each of the treatment levels. The resulting estimates are the same as if using the constraint that  
151 the sum of all the treatment effects is zero. If it is desired to sort the lsmeans from the highest to the  
152 lowest, change the LSMEANS statement to the following:

153

```
154 lsmeans treat / out=lsmeans noprint;
```

155

156 and then add

157

```
158 proc sort data=lsmeans;
```

```
159     by descending lsmean;
```

```
160 proc print;
```

```
161 run;
```

162

163 The NOPRINT option prevents the lsmeans from being printed during the PROC GLM invocation. Since  
164 most AEDs have large numbers of new treatments, the above ordering is a desirable feature for the  
165 experimenter who wishes to select the top performers and to discard poor performers.

166

167 To run an analysis on check yields only for  $n = 6$  new treatments and  $c = 4$  checks, add the following  
168 statement to the end of the DATA step:

169

170 `if treat > 6 and treat < 11 then check = treat;`

171

172 To obtain additional sums of squares, the following code can be useful:

173

174 `data augbibd;`

175 `infile 'augbibd.dat';`

176 `input yield rep block treat;`

177 `if (treat > 6) then new = 0;`

178 `else new = 1;`

179 `if (new) then treatn = 999;`

180 `else treatn = treat;`

181 `proc glm data=augbibd;`

182 `class rep block treat treatn;`

183 `model yield = rep block(rep) treatn treat*new;`

184 `random rep block(rep);`

185 `lsmeans treatn;`

186 `run;`

187

188

189 PROC MIXED for Recovering Interblocking Information

190

191 Although the preceding PROC GLM code can provide a fairly complete analysis of data from an  
192 augmented design, it can be a difficult chore sorting through the various sums of squares and constructing

193 appropriate tests. We therefore recommend PROC MIXED be used for most augmented design problems.  
194 The output is much more straightforward and direct account is made of random effects.

195

196 Our first analysis using PROC MIXED considers both check and new treatments to be fixed effects and  
197 replicates and blocks to be random effects:

198

```
199 data augbibd;  
200     infile 'augbibd.dat';  
201     input yield rep block treat;  
202 proc mixed data = augbibd;  
203     class rep block treat;  
204     model yield = treat;  
205     random rep block(rep);  
206     lsmeans treat;  
207 run;
```

208

209 Note the syntax for PROC MIXED is nearly identical to that of PROC GLM, with one important  
210 exception: only fixed effects are listed in PROC MIXED's MODEL statement, whereas both fixed and  
211 random effects are listed in PROC GLM's MODEL statement. As noted before, this is not really an  
212 inconsistency because PROC GLM considers all effects to be fixed when it fits the linear model. On the  
213 other hand, PROC MIXED handles random effects directly by estimating their variance components using  
214 Gaussian restricted maximum likelihood.

215

216 The treatment means resulting from the LSMEANS statement are adjusted for interreplicate (for the new  
217 treatments) and interblock information, and associated t-statistics take into account all estimated variance  
218 components.

219

220

221 PROC MIXED for Recovering Both Interblocking and Intervariety Information

222

223 Since the SAS system is not designed to partition a variable such as TREAT into a set which is fixed (the  
224 checks) and a set which is random (the new treatments), it is necessary to construct some auxiliary  
225 variables in order to accomplish this. The following program considers the checks as fixed effects and the  
226 other effects as random, still assuming  $n = 6$  new treatments and  $c = 4$  checks:

227

228 data augbibd;

229 infile 'augbibd.dat';

230 input yield rep block treat;

231 if (treat &gt; 6) then new = 0;

232 else new = 1;

233 if (new) then treatn = 999;

234 else treatn = treat;

235 proc mixed data=augbibd;

236 class rep block treat treatn;

237 model yield = treatn;

238 random rep block(rep) treat\*new / solution;

239 lsmeans treatn;

240 make 'solutionr' out=sr noprint;

241 run;

242 proc sort data=sr;

243 by descending est;

244 proc print;

245 run;

246

247 The DATA step creates two new auxiliary variables: NEW and TREATN. NEW indicates whether or not  
248 a treatment is a new treatment, and is subsequently used to construct the random effect corresponding to  
249 the new treatments. TREATN equals TREAT for all of the check treatments but has a constant level for  
250 all of the new treatments. It is used as a fixed effect to model different means for each of the check  
251 treatments and a common mean for the new treatments. The new treatments are thus assumed to vary  
252 randomly about a common mean, and note this mean is free to fall anywhere in relation to the check  
253 means. The MODEL statement thus lists TREATN as the sole fixed effect, and the subsequent  
254 LSMEANS statement uses TREATN to construct mean estimates.

255

256 The RANDOM statement again lists REP and BLOCK(REP) as random effects along with a new one:  
257 TREAT\*NEW. This last effect equals 0 for all of the check treatments (note that NEW is not a CLASS  
258 variable) and has a different level for all of the new treatments. The SOLUTION option in the RANDOM  
259 statement requests empirical best linear unbiased predictors (EBLUPs) of the random effects.

260

261 The MAKE statement is PROC MIXED's mechanism for creating output data sets, and the one listed here  
262 creates a data set named SR from the EBLUP table printed by the SOLUTION option in the RANDOM  
263 statement. The 'solutionr' string is the label for this table and is a necessary part of the MAKE statement.  
264 All such labels as well as detailed information on every statement can be found in the PROC MIXED  
265 documentation (SAS Institute Inc., 1996).

266

267 The final lines of the program sort and print the EBLUPs. Note that the EBLUPs for REP, BLOCK(REP),  
268 TREAT\*NEW will all be intermixed in this printout, and one may wish to extract just those for  
269 TREAT\*NEW in a different analysis. The sorted EBLUPs for TREAT\*NEW provide a means for  
270 comparing the new treatments.

271

272

273 Other Augmented Designs

274

275 The above ideas for recovering interblock as well as intervariety information is easily extended to other  
276 augmented designs. The analysis described by Federer (1996) provides a useful example, and others can  
277 be found in Federer *et al.* (1975) and Federer and Raghavarao (1975). This example uses an input data  
278 file augmerc1.dat, and there are  $c = 2$  checks repeated 30 times each in 15 rows and 12 columns with  $n =$   
279 120 new genotypes.

280

281 PROC GLM and PROC MIXED programs are now presented for obtaining Type III sums of squares,  
282 intrarow-column (fixed effects) least squares means, check treatment means adjusted for interrow and  
283 intercolumn information, and new treatment means adjusted for interrow, intercolumn, and intervariety  
284 information. A PROC GLM program for obtaining some relevant sums of squares is as follows:

285

```
286 data augmerc1;
287     infile 'augmerc1.dat';
288     input site col row treat gw c1 c2 c3 c4 r1 r2 r3 r4;
289     if (treat > 120) then new = 0;
290     else new = 1;
291     if (new) then treatn = 999;
292     else treatn = treat;
293     ll = r1*c1;
294     lq = r1*c2;
295 proc glm data=augmerc1;
296     class row col treat treatn;
297     model gw = r1 r2 r3 r4 c1 c2 c3 c4 ll lq treatn treat*new;
298     random row col treat*new;
299 run;
```

300

301 The dependent variable GW is grain weight.  $R_i$  and  $C_i$  ( $i = 1, 2, 3, 4$ ) are orthogonal polynomial  
302 regressions for row and column numbers. These variables can be created using the ORPOL function in  
303 SAS/IML. For example, the following program creates a data set OPN15 containing variables ROW and  
304 R1-R4. This data set can then be match-merged with the original data set.

305

```
306 proc iml;
307   opn15 = orpol(1:15,4);
308   opn15[,1] = (1:15)';
309   op15 = opn15;
310   create opn15 from opn15[colname={'ROW' 'R1' 'R2' 'R3' 'R4'}];
311   append from opn15;
312   close opn15;
313 run;
```

314

315 The term columnname refers to the column in the created data set and not to the column of the  
316 experiment design. In the preceding DATA step LL and LQ are created to represent interactions of row  
317 and column regressions. These are then specified along with the other variables in the MODEL statement  
318 of PROC GLM. This particular model is used here because the design is not connected. Other regression  
319 terms may be added to the model if deemed necessary to explain the experimental variation.

320

321 Note that the row, column, and new treatment effects are considered random. The RANDOM statement is  
322 used to obtain the expected values of mean squares in the event ANOVA solutions for the row, column,  
323 and new variance components are required.

324

325 The following code constructs the fixed effects means and arranges them in order from highest to lowest:

326

```
327 proc glm data=augmerc1;
```

```
328   class row col treat treatn;  
329   model gw = r1 r2 r3 r4 c1 c2 c3 c4 || lq treat;  
330   lsmeans treat / out = lsmeans noprint;  
331 run;  
332 proc sort data=lsmeans;  
333   by descending lsmean;  
334 proc print;  
335 run;
```

336

337 While the preceding PROC GLM code can be used to obtain various partitions of the sums of squares, a  
338 more straightforward analysis can be obtained with PROC MIXED. This analysis adjusts the check  
339 means for interrow and intercolumn information and adjusts the new effects for interrow, intercolumn,  
340 and intervariety information:

341

```
342 proc mixed data=augmerc1;  
343   class row col treat treatn;  
344   model gw = r1 r2 r3 r4 c1 c2 c3 c4 || lq treatn / solution;  
345   random r1 r2 r3 r4 c1 c2 c3 c4 || lq treat*new / solution;  
346   lsmeans treatn;  
347   make 'solutionr' out=sr noprint;  
348 run;  
349 proc sort data=sr;  
350   by descending est;  
351 proc print;  
352 run;
```

353

354 If there were no row and column regression interactions, the  $R_i$  and  $C_i$  in the random statement could be  
355 replaced with ROW and COL as REML makes use of the normal distribution theory and the design need  
356 not be connected for the PROC MIXED procedure. Using the row and column designation for random  
357 corrects for all row and column effects and not just the  $R_i$  and  $C_i$  used in the regression model. Note that  
358 TREATN is used in the MODEL statement, and it has a distinct level for all of the check treatments and a  
359 common level (999) for all of the new treatments. The effect TREAT\*NEW in the RANDOM statement  
360 models all of the new treatments as random effects varying about the common mean modeled by the 999  
361 level of TREATN. For some computer set-ups, it may be necessary to use the command \_EST\_ in place  
362 of EST.

363 A report describing the use of the above programs has been prepared by Federer and Wolfinger  
364 (1996a). A small numerical example has been used and the computer outputs of the programs have been  
365 annotated with descriptions of the results.

366

#### 367 4. Conclusions

368

369 Present software literature is inadequate for experimenters to obtain programs for the analyses described  
370 herein. Considerable computer and package expertise and several trial and error runs were required to  
371 obtain the final programs. These programs are in a readily usable form for experimenters who desire  
372 statistical analyses of augmented experiment designs and to recover interblock, interrow, intercolumn,  
373 interregression, and/or intervariety information. Recovery of the information associated with the random  
374 effects leads to more efficient analyses of experimental data, and hence more efficient use of experimental  
375 resources.

376

#### 377 Acknowledgement

378 Appreciation is expressed for the constructive comments of an associate editor and referees. These were  
379 useful in clarifying presentation of results.

380

381 References

382

383 Cullis, B. R., J. L. Warwick, J. A. Fisher, and B. J. Read (1989). A new procedure for the analysis of  
384 early generation variety trials. *Applied Statistics* 38:361-375.

385

386 Federer, W. T. (1956). Augmented (or hoonuiaku) designs. *Hawaiian Planters' Record* LV(2):191-208.

387

388 Federer, W. T. (1961). Augmented designs with one-way elimination of heterogeneity. *Biometrics*  
389 17:447-473.

390

391 Federer, W. T. (1991). *Statistics and Society*, 2nd edition. Marcel Dekker, Inc., New York, Section 7.11.

392

393 Federer, W. T. (1996). Recovery of interblock, intergradient, and intervariety information in incomplete  
394 block and lattice rectangle designed experiments. *Biometrics* (in press).

395

396 Federer, W. T. , R. C. Nair, and D. Raghavarao (1975). Some augmented row-column designs.  
397 *Biometrics* 31:361-373.

398

399 Federer, W. T. and D. Raghavarao (1975). On augmented designs. *Biometrics* 31:29-35.

400

401 Federer, W.T. and R. D. Wolfinger (1996). SAS PROC GLM and PROC MIXED Code for Recovering  
402 Inter-Factor Information. BU-1330-M in the Technical report Series of the Biometrics Unit, Cornell  
403 University, Ithaca, New York.

404

405 Federer, W. T. and R. D. Wolfinger (1996a). GAUSS and SAS for recovering interblock and intervariety  
406 information. BU-1384-M in the Technical Report Series of the Biometrics Unit, Cornell University,  
407 Ithaca, New York.

- 408 SAS Institute Inc. (1989). SAS/STAT Users Guide, Version 6, Fourth Edition, Volume 2. SAS Institute  
409 Inc., Cary, NC.  
410
- 411 SAS Institute Inc. (1996). SAS/STAT Software: Changes and Enhancements through Release 6.11. SAS  
412 Institute Inc., Cary, NC.  
413
- 414 Yates, F. (1936). A new method of arranging variety trials involving a large number of varieties. Journal  
415 of Agricultural Science 26:424-455.