

# Mixed models and multivariate analysis for selection of superior maize genotypes

Gustavo H.F. Oliveira<sup>1\*</sup>, Camila B. Amaral<sup>1</sup>, Flávia A.M. Silva<sup>1</sup>, Sophia M.F. Dutra<sup>1</sup>, Marcela B. Marconato<sup>1</sup>, and Gustavo V. Mõro<sup>1</sup>



## ABSTRACT

Selections via the mixed model and the multivariate analysis approach can be powerful tools for selecting cultivars in plant breeding programs. Therefore, this study aimed to compare the use of mixed models, multivariate analysis and traditional phenotypic selection to identify superior maize (*Zea mays* L.) genotypes. Seventy-one (71) maize Topcrosses and three commercial cultivars were evaluated using these three methods. Plant height, ear height, ear placement, stalk lodging and breakage, and grain yield were evaluated. There was a difference between selection methods, as the selection with mixed models and the selection based on the average phenotypic afforded the inclusion of genotypes with high productivity, which did not occur for the multivariate analysis. The selection by multivariate analysis allowed the inclusion of genotypes with better agronomic and other desirable traits, not only those with highest productivity, in a maize breeding program.

**Key words:** Blup, K-means, *Zea mays*.

## INTRODUCTION

The main objective of plant breeding programs is the development of marketable cultivars, and success depends on applying appropriate selection methods to available genotypes. For this process, selections by mixed models focusing on multivariate analyses are powerful tools for selecting cultivars.

For mixed models, the Best Linear Unbiased Predictor (BLUP) is near the phenotypic average observed for the true genotypic value of the individual, which is a property of an accurate estimator (Piepho et al., 2008). BLUP can be used for both the selection of superior individuals and the estimation of future generations' gains. This allows selection for precocity, thus making possible crossings of only individuals that are promising for the characteristic of interest (Cruz Baldissera et al., 2012). On the other hand, truncated selection, where characteristics are selected for one at a time, may compromise the flexibility of the breeding program. Thus, multivariate analyses, besides combining various data from a single experiment (Gonçalves et al., 2014), can contribute to the robustness of the selection (Magalhães Bertini et al., 2010).

Multivariate analysis refers to a broad category of methods used when different variables are measured in a single set of experimental data (Yeater et al., 2014). Among multivariate analyses is the cluster analysis, which seek to minimize the differences within groups and maximize the differences between clusters. Therefore, the objective of this study was to compare the application of mixed models, multivariate analysis and traditional phenotypic selection to identify superior genotypes of maize.

## MATERIALS AND METHODS

The study was conducted over the harvests of 2012-2013 and 2013-2014, with the evaluation of 71 maize Topcrosses and three commercial cultivars at the experimental farm of the College of Agriculture and Veterinary Sciences of Universidade Estadual Paulista (UNESP), Jaboticabal (21°15'17" S, 48°19'20" W; 605 m a.s.l.), São Paulo, Brazil.

The design was randomized blocks with two replicates. The plots consisted of two 5-m-long rows, spaced 0.50 m apart, with 18 plants in each row, for an ideal stand of 36 plants per plot. The evaluated characteristics were: ear height (EH, cm) measured as the distance from the ground to the insertion point of the main ear; plant height (PH, cm) measured as the distance from the ground to the flag leaf; ear position (EP) as the ratio between ear

<sup>1</sup>Universidade Estadual Paulista 'Júlio de Mesquita Filho' (UNESP), Faculdade de Ciências Agrárias e Veterinárias, Via de acesso Prof. Paulo Donato Castellane s/n, 14884-900, Jaboticabal, São Paulo, Brasil.

\*Corresponding author (gustavo.genetica@posgrad.fcav.unesp.br).

Received: 3 March 2016.

Accepted: 16 September 2016.

doi:10.4067/S0718-58392016000400005



height and plant height; grain yield (GY, kg ha<sup>-1</sup>) obtained from the average weight of the useful grains; and stalk lodging and breakage (SLB) estimated by the number of broken plants (plants with breakage in the stem below the ear), fallen plants, and those plants growing at less than a 45° angle. Grain production was corrected to 13% moisture and adjusted by the average stand for covariance (Silva et al., 2014), converting it to kg ha<sup>-1</sup>.

Phenotypic selection was performed using the average obtained from the joint variance of the 2 yr evaluation, selecting the 20 best genotypes for productivity. For the mixed model selection, both environments were considered in joint analysis for extracting the BLUP's of all genotypes, which were then ordered by the 20 best genetic values (BLUP) for productivity (GY).

The linear mixed model used was:

$$y = X\beta + Z_1g + Z_2w + \varepsilon$$

where  $y$  is the vector of phenotypic observations;  $\beta$  is the vector of fixed effects due to the blocks, local, and general average;  $g$  is the vector of the effects of the assumed random genotypes;  $w$  is the vector of the effects of the interaction of Genotype  $\times$  Environment (random);  $X$ ,  $Z_1$ , and  $Z_2$  are incidence matrices of these effects, and  $\varepsilon$  is the vector of random residues.

Whereas the BLUE model for  $\beta$  and BLUP for the random effects  $g$  and  $w$  are given in the following mixed model equation:

$$\begin{bmatrix} X'R^{-1}X & X'R^{-1}Z_1 & X'R^{-1}Z_2 \\ Z_1'R^{-1}X & Z_1'R^{-1}Z_1 + G^{-1} & Z_1'R^{-1}Z_2 \\ Z_2'R^{-1}X & Z_2'R^{-1}Z_1 & Z_2'R^{-1}Z_2 + W^{-1} \end{bmatrix} \begin{bmatrix} \beta \\ g \\ w \end{bmatrix} = \begin{bmatrix} X'R^{-1}y \\ Z_1'R^{-1}y \\ Z_2'R^{-1}y \end{bmatrix}$$

The estimates of variance components necessary to obtain the subjects' genetic values (BLUP) were determined using the restricted maximum likelihood (REML) method. The confidence intervals for the averages were estimated by:

$$IC = u \pm t \cdot sdaf$$

where  $u$  is the average of the family,  $t$  is the tabulated value of the "t" student distribution at 5% probability, and  $sdaf$  is the standard deviation of the average family.

The non-hierarchical, multivariate K-means cluster analysis was carried out to find the shortest distance between the genotypes and select the superior genotypes closest to the commercial cultivars. Because of the high

number of genotypes,  $K = 5$  was stipulated, the number of groups was chosen *a priori* according to the assumption for the separation of clusters (Kanungo et al., 2002; Nazeer and Sebastian, 2009).

This tool selected genotypes from the commercial hybrids and/or genotypes that showed good agronomic attributes according to the analyzed variables. All analyses were performed using SAS (SAS Institute, Cary, North Carolina, USA).

## RESULTS AND DISCUSSION

The genotypes differed for all traits and environments (Table 1), indicating that, although grown in the same place, cultivation in different crops was sufficient to influence the genotypes. The overall average grain yield in both environments was 8208.65 kg ha<sup>-1</sup>. Therefore, the genotypes presented high productivity estimates, since average productivity in Brazil is about 4.2 t ha<sup>-1</sup> (Lyra et al., 2014). The variation coefficients for all variables were within the normal range, revealing good experimental precision (Hallauer et al., 2010; Fritsche-Neto et al., 2012).

There was Genotypes  $\times$  Environments interaction for all traits except ear height and ear placement. Thus, these results indicate that there is variability between individuals, environments were contrasting, and the selection made considering both genotypes and environments together will provide gains generated from contrasting environments with the genotypes becoming more adapted and stable.

For the selection, it is important to consider that according to the amplitudes, there was high variability among individuals (Table 2). The top 20 genotypes selected by joint analysis were the same as those selected via mixed models (Table 2). This indicates that with respect to the true genetic value of individuals, the averages estimated by the least squares method were sufficiently precise, even when considering the effect of genotypes as fixed and making changes in the ranking of genotypes. However, in the mixed model analysis, there was a reduction in the variation range of the genotypes (2660.94 kg ha<sup>-1</sup>). This difference is due to the shrinkage effect inherent in BLUP, which eliminates residual effects embedded in phenotypic data, thereby providing data for genetic and non-

**Table 1. Summary of the analysis of joint variance of five agronomic traits in 74 evaluated maize genotype.**

Variation sources	DF	MS				
		GY (kg ha <sup>-1</sup> )	PH (cm)	EH (cm)	+EP	SLB
Genotype (G)	73	3112.32**	242.82**	122.98**	1.58**	0.64**
Environment (E)	1	2614147.48**	28346.39**	1352.22**	65.40**	9.48**
G $\times$ E	73	1473.91*	124.89**	63.46 <sup>ns</sup>	0.28 <sup>ns</sup>	0.37*
Residual	145	1014.57	78.83	51.55	0.38	0.26
CV, %		12.27	3.78	5.31	3.40	37.39
Average		8208.65	234.58	134.99	0.57	1.36
Amplitude		5047.74	43.07	32.34	0.09	1.52

MS: Mean squares, DF: degree of freedom, GY: Grain yield, PH: plant height, EH: ear height, EP: ear placement, SLB: stalk lodging and breakage, CV: experimental coefficient of variation, <sup>ns</sup>: nonsignificant by F test.

\*, \*\*Significant at the 0.05 and 0.01 probability levels, respectively by F test.

\*Value multiplied by 1000.

**Table 2. Phenotypic and mixed model selection of 20 individuals from 74 genotypes by productivity.**

Genotype	Phenotypic selection			Genotype	Selection by mixed models		
	GY	ULCI	LLCI		GY	ULCI	LLCI
2B707	11034.30	14502.76	7575.88	2B707	9702.91	10959.44	8446.38
DKB 390	10914.10	14382.54	7445.67	DKB 390	9639.53	10896.07	8383.01
AG7000	10144.40	13612.84	6675.97	AG7000	9233.78	10490.32	7977.25
G-54	9494.34	12962.78	6025.97	G-54	8891.10	10147.64	7634.57
G-19	9296.78	12765.21	5828.34	G-19	8786.95	10043.49	7530.42
G-61	9211.94	12680.37	5743.50	G-61	8742.23	9998.76	7485.69
G-64	9208.22	12676.66	5739.78	G-64	8740.27	9996.81	7483.74
G-55	9204.36	12672.80	5737.92	G-55	8738.24	9994.77	7481.70
G-71	9132.10	12600.53	5663.66	G-71	8700.14	9956.67	7443.61
G-14	9107.20	12575.64	5638.76	G-14	8687.02	9943.55	7430.48
G-26	9012.46	12480.90	5544.03	G-26	8637.07	9893.61	7380.54
G-15	8964.03	12432.46	5495.59	G-15	8611.54	9868.07	7355.01
G-05	8949.33	12417.76	5480.89	G-05	8603.79	9860.33	7347.26
G-13	8929.33	12397.80	5460.93	G-13	8593.27	9849.80	7336.74
G-57	8859.88	12328.31	5391.44	G-57	8556.64	9813.17	7300.11
G-47	8816.15	12284.59	5347.72	G-47	8533.59	9790.12	7277.06
G-56	8759.79	12228.22	5291.35	G-56	8503.87	9760.41	7247.34
G-17	8686.51	12154.95	5218.07	G-17	8465.25	9721.78	7208.72
G-09	8610.64	12097.44	5160.57	G-09	8434.93	9691.47	7178.41
G-50	8610.64	12079.08	5142.21	G-50	8425.25	9681.78	7168.72

GY: Grain yield (kg ha<sup>-1</sup>); ULCI and LLCI is the upper and lower limit (kg ha<sup>-1</sup>) of the confidence interval at 5% probability.

phenotypic selection similar to that which occurs in a least-squares analysis (Resende, 2007). In addition, it confirms the efficiency of least squares with respect to mixed models for balanced information.

The prediction of genotypic values for grain yield, based on the effect of the families for the 20 best evaluated genotypes via BLUP are listed in Table 2. According to Arnhold et al. (2009), it is better to use the mixed model methodology for breeding programs and BLUP is preferable to BLUE, as BLUP can be used to estimate the breeding values of individuals.

Out of the top 20 genotypes, the three commercial hybrids outperformed the general average with the highest yields (Table 2). However, when the overlap of the confidence intervals of the mean yield is considered (Table 1), it was possible to obtain genotypes with higher productivity than the general average and similar to the commercial controls. This reveals the existence of genotypes with high yielding potential for inclusion in breeding programs, but assessments in different locations are needed for additional studies of adaptability and stability.

Rivas and Barriga (2002) discuss the importance of choosing the parents for superior combinations in breeding programs. The choice of these parents can be through the genotypes' combining ability, which provides

information for obtaining crosses for improving the various characteristics of interest (Arnhold et al., 2009).

In the multivariate analysis selection, it was possible to characterize the genotypes into five groups through the K-means method (Table 3, Figure 1).

Group 1 contained approximately 8% of genotypes and had the highest average plant and ear height (247.5 and 145.78 cm, respectively) (Table 3). Theoretically, these genotypes will have an increased tendency for stalk lodging and breakage because the ears are higher on the stalk (Figure 1). It is important to remember that in breeding experiments, the harvesting is often manually done, with broken and lodged plants being harvested, and thus, it is possible to account for the productivity of these genotypes. However, this would not occur with mechanical harvesting, despite the positive correlation between productivity and plant height (Yin et al., 2011).

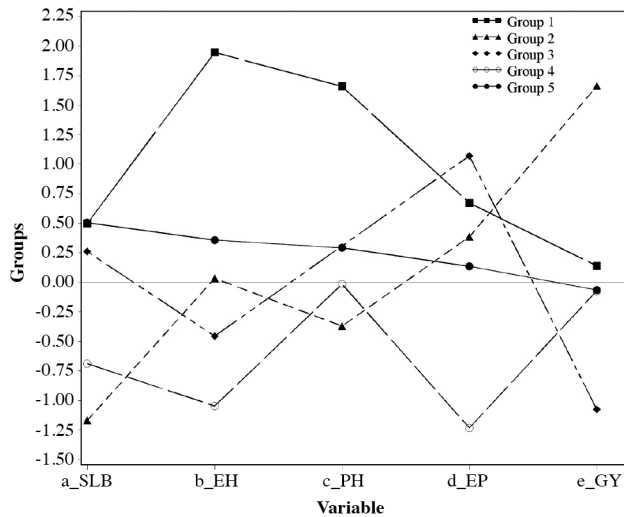
Group 2 was made up of the commercial cultivars and an additional five genotypes (approximately 10% of genotypes). It had the highest grain yields and less lodging, revealing smaller sized candidates for genotype breeding programs (Table 3). The genotypes in group 2 should be evaluated in more places and different growing seasons, because they performed almost as well as the commercial cultivars and therefore have increased agronomic potential (Figure 1).

**Table 3. Group means using the multivariate method K-means.**

Group	Variable					Genotype
	PH	EH	EP	SLB	GY	
1	247.50	145.78	0.58	1.56	8329.73	G-10, G-16, G-35, G-43, G-47, G-71
2	231.69	135.17	0.58	0.89	9677.76	DKB, AG7000, 2B707, G-19, G-22, G-26, G-61, G-64
3	222.05	132.46	0.59	1.46	7252.94	G-11, G-12, G-21, G-27, G-36, G-37, G-38, G-39, G-52
4	234.47	129.18	0.55	1.08	8141.33	G-18, G-29, G-46, G-48, G-49, G-55, G-56, G-57, G-58, G-59, G-60, G-61, G-63, G-66, G-67, G-68, G-69, G-70
5	236.86	136.97	0.57	1.56	8149.45	G-07, G-08, G-09, G-13, G-14, G-15, G-17, G-23, G-24, G-25, G-28, G-30, G-31, G-32, H-33, G-34, G-01, G-40, G-41, G-42, G-44, G-45, G-02, G-50, G-51, G-53, G-54, G-03, G-65, G-07, G-05, G-06

PH: Plant height, EH: ear height, EP: ear placement, SLB: stalk lodging and breakage, GY: grain yield.

**Figure 1. Multivariate non-hierarchical cluster analysis obtained by the K-means method with the standardized averages of the effects of grain yield (GY), plant height (PH), ear height (EH), ear placement (EP), and stalk lodging and breakage (SLB).**



The genotypes G-19, G-61, G-64, G-26 and G-22 are present in this group and were also selected for in the mixed and average phenotypic model methods. They may be incorporated into breeding programs for low breakage and stalk lodging with high yields. Except for G-22, these genotypes are on the list of 20 selected genotypes found in Table 2. This shows that the selection made with multivariate models encompassed not only higher yielding genotypes, but those which had better agronomic traits (taking into account all features together) and can be considered as one of the functions of the selection index (Gonçalves et al., 2014). Thus, one can have significant productivity gains combined with good agronomic traits with the selection of these genotypes.

Group 3, comprising approximately 12% of the genotypes, was the least productive. This is possibly due to the high incidence of genotypes with relative ear height above the mean, as well as the presence of plants with high lodging and breakage index (1.46). This may be because when a genotype has ears high on the stalk (ear height) (Figure 1), the possibility of stalk lodging and breakage is greater, and even though harvesting is by hand, the yield is negatively impacted by contact with the ground and the increase in diseases associated with that contact.

These values should only be used as exploratory data for the genotypes' agronomic characteristics, as the correlation between the ear position and productivity are relatively low (Alvi et al., 2003; Bello et al., 2010; Toebe and Cargnelutti Filho, 2013). This group's genotypes were not selected when using the mixed model and average phenotypic methodology, probably because they do not exhibit high productivity. This shows the effectiveness of using contrasts when applied to the multivariate methodology by K-means (Kanungo et al., 2002).

Group 4, comprising approximately 25% of genotypes, were shorter with low lodging rates, which are desirable traits in breeding programs. Of the genotypes in this group, only three (G-55, G-56, G57) were also among the 20 best genotypes selected according to their genetic values. These genotypes deserve special attention because they have good qualities found in both analyses.

Finally, in Group 5, where approximately 44% of all genotypes were allocated, the agronomic traits are not favorable, since their above-average characteristics (SLB, EP, EH, PH) are all undesirable and their average yields were below the overall average (Figure 1). On the other hand, there are eight genotypes (G-05, G-09, G-13, G-14, G-15, G-17, G-50, G-54) in this group that were selected via the mixed model considering BLUP for productivity, revealing that even within groups with below-average performance, promising genotypes may be found. This shows that the analysis via the multivariate approach using k-means can be a tool for selection based on mixed models, because despite these genotypes having been selected for productivity, they do not meet appropriate agronomic standards.

In the multivariate approach, the genotypes in groups 1 and 2 were selected because they present high productivity averages and agronomic characteristics similar to those found in the commercial cultivars. However, the mixed model methodology found five genotypes (G-10, G-16, G-22, G-35, G-43) which were not selected in the multivariate analysis, emphasizing the complementarity of the two approaches. In addition, in the mixed model selection via BLUP, selection favors the most productive individuals, while the multivariate analysis reveals individuals with better agronomic standards.

The line drawn on the zero axis "y" of "Groups" is the overall average of the standardized variables where  $\mu = 0$  and  $\sigma = 1$ . Group 2 is represented by the red line, and this line is noticeably different from the others and its superiority is apparent. This shows the average of the evaluated variables that were within the acceptable standard for maize genetic improvement programs, when the goal is to reduce all evaluated traits except productivity (Figure 1).

## CONCLUSION

It was possible to accurately select genotypes with high productivity and good agronomic traits. There was a significant difference between selection methods. Furthermore, the mixed model and the average phenotypic-based methods selected the same genotypes with high productivity. This did not occur in the multivariate analysis.

Selection with multivariate analysis allowed the breeding improvement program to include genotypes with better agronomic and other desirable traits, instead of only those with the highest productivity.

The selection based on BLUP via mixed models can be supplemented with multivariate analysis using k-means, contributing to the accuracy and robustness of the selection, and thereby providing superior genotypes in less time.

## ACKNOWLEDGEMENTS

The authors thank the Foundation for Cordenação de Aperfeiçoamento de Pessoal de Nível Superior (CAPES), for granting a scholarship.

## LITERATURE CITED

- Alvi, M.B., M. Rafique, M.S. Tariq, A. Hussain, T. Mahmood, and M. Sarwar. 2003. Character association and path coefficient analysis of grain yield and yield components maize (*Zea mays* L.) Pakistan Journal of Biological Sciences 6:136-138. doi:10.3923/pjbs.2003.136.138.
- Arnhold, E., F. Mora, R.G. Silva, P.I. Good-God, and M.A. Rodovalho. 2009. Evaluation of top-cross popcorn hybrids using mixed linear model methodology. Chilean Journal of Agricultural Research 69:46-53. <http://dx.doi.org/10.4067/S0718-58392012000100026>.
- Bello, O., S. Abdulmalik, M. Afolabi, and S. Ige. 2010. Correlation and path coefficient analysis of yield and agronomic characters among open pollinated maize varieties and their F1 hybrids in a diallel cross. African Journal of Biotechnology 9:2633-2639. doi:10.5897/AJAR2015.9613.
- Cruz Baldissera, J.N., J.G. Bertoldo, G. Valentini, M.M.D. Coan, D.S. Rozeto, A.F. Guidolin, et al. 2012. Uso do melhor preditor linear não-viesado (BLUP) na predição de híbridos em feijão. Bioscience Journal 28:395-403.
- Fritsche-Neto, R., R.A. Vieira, C.A. Scapim, G.V. Miranda, and L.M. Rezende. 2012. Updating the ranking of the coefficients of variation from maize experiments. Acta Scientiarum. Agronomy 34:99-101.
- Gonçalves, A., L. Simões, S.D.P. Freitas Júnior, A.T. Amaral Júnior, C.A. Scapim, R. Rodrigues, et al. 2014. Estimating combining ability in popcorn lines using multivariate analysis. Chilean Journal of Agricultural Research 74:10-15. doi:10.5897/AJAR2015.9613.
- Hallauer, A.R., M.J. Carena, and J.D. Miranda Filho. 2010. Quantitative genetics in maize breeding. Handbook of Plant Breeding Vol 6. 664 p. Springer-Verlag, New York, USA.
- Kanungo, T., D.M. Mount, N.S. Netanyahu, C.D. Piatko, R. Silverman, and A.Y. Wu. 2002. An efficient *k*-means clustering algorithm: Analysis and implementation. IEEE Transactions on Pattern Analysis and Machine Intelligence 24:881-892. doi:10.1109/TPAMI.2002.1017616.
- Lyra, G.B., A.E.Q. Rocha, G.B. Lyra, J.L. Souza, e I. Teodoro. 2014. Crescimento e produtividade do milho, submetido a doses de nitrogênio nos Tabuleiros Costeiros de Alagoas. Revista Ceres 61:578-586. <http://dx.doi.org/10.1590/0034-737X201461040019>.
- Magalhães Bertini, C.H.C., W.S. Almeida, A.P.M. Silva, J.W. Lima, e E.M. Teófilo. 2010. Análise multivariada e índice de seleção na identificação de genótipos superiores de feijão-caupi. Acta Scientiarum. Agronomy 32:613-619. doi:10.4025/actasciagron.v32i4.4631.
- Nazeer, K.A., and M. Sebastian. 2009. Improving the accuracy and efficiency of the *k*-means clustering algorithm. p. 1-3. In Proceedings of the World Congress on Engineering Vol. 1, London. 1-3 July. International Association of Engineers (IAENG), London, UK.
- Piepho, H., J. Möhring, A. Melchinger, and A. Büchse. 2008. BLUP for phenotypic selection in plant breeding and variety testing. Euphytica 161:209-228. doi:10.1007/s10681-007-9449-8.
- Resende, M.D.V. 2007. Matemática e estatística na análise de experimentos e no melhoramento genético. Embrapa Florestas, Colombo, Paraná, Brasil. Available at [http://livraria.sct.embrapa.br/liv\\_resumos/pdf/00083145.pdf](http://livraria.sct.embrapa.br/liv_resumos/pdf/00083145.pdf) (accessed 2 April 2016).
- Rivas, P., y B. Barriga. 2002. Capacidad combinatoria para rendimiento de grano y caracteres de calidad maltera en cebada (*Hordeum vulgare* L.) Agricultura Técnica 62:347-356. <http://dx.doi.org/10.4067/S0365-28072002000300001>.
- Silva, K.J., C.B. Menezes, F.D. Tardin, V.F. Souza, e C.V. Santos. 2014. Comparação de métodos de correção de estande para estimar a produtividade de sorgo granífero. Pesquisa Agropecuária Tropical 44:175-181. <http://dx.doi.org/10.1590/S1983-40632014000200005>.
- Toebe, M., e A. Cargnelutti Filho. 2013. Não normalidade multivariada e multicolinearidade na análise de trilha em milho. Pesquisa Agropecuária Brasileira 48:466-477. <http://dx.doi.org/10.1590/S0100-204X2012000900002>.
- Yeater, K.M., S.E. Duke, and W.E. Riedell. 2014. Multivariate analysis: greater insights into complex systems. Agronomy Journal 107:799. doi:10.2134/agronj14.0017.
- Yin, X., M.A. McClure, N. Jaja, D.D. Tyler, and R.M. Hayes. 2011. In-season prediction of corn yield using plant height under major production systems. Agronomy Journal 103:923-929. doi:10.2134/agronj2010.0450.