# Best Linear Unbiased Prediction of Cultivar Effects for Subdivided Target Regions

H. P. Piepho* and J. Möhring

## ABSTRACT

Breeding for local adaptation may be economically viable providing there is sufficient genotype × subregion interaction. If the targeted subregion is part of a larger region covered by a testing network, information from neighboring subregions can be exploited to gain more precise estimates for the targeted subregion. For balanced data, the simplest approach is to use genotypic mean estimates for the whole target region, and this has often been shown to yield better predictions than simple means per subregion. A disadvantage of this approach is that it gives equal weight to all neighboring subregions and the targeted subregion, thus ignoring potential heterogeneity in information content. The objective of the present paper is to propose a method that allows a weighted combination of data from several subregions and to compare that method to other estimators. The proposed method is based on best linear unbiased prediction, which employs a weighted mean of subregion means. It follows from the theory of mixed models that the resulting estimator is optimal under the assumed model, minimizing prediction errors and maximizing the expected gain from selection. Using published variance component estimates, we found the resulting predictions to be superior to other approaches. We also show that the estimator is beneficial when selecting for global adaptation.

P LANT BREEDERS usually seek to develop broadly adapted varieties for a wider target region. If the target region is agroecologically diverse, it may be worthwhile to stratify the target into more homogeneous subregions. Stratification will allow more accurate overall performance estimates for candidate varieties in the target region, thus increasing gain from selection. Alternatively, plant breeders may opt to develop locally adapted varieties for specific subregions. Breeding for local adaptation will be worthwhile only when there is substantial genotype × subregion interaction. Moreover, division of a target region will be accompanied by a division of testing resources. Thus, despite presence of substantial genotype × subregion interaction, it may turn out to be more efficient to breed for broad adaptation, if resources are not sufficient for accurately detecting locally adapted genotypes. This has been lucidly demonstrated by Curnow (1988) and Atlin et al. (2000), who studied the response to selection in subdivided target regions. The authors considered genotypic means in a large target and constituent subregions as correlated traits. They showed that the correlated response to selection for overall performance may outperform the direct response to selection within subregions.

There are two opposing factors that will determine if breeding for local adaptation is worthwhile. On the one hand, subdividing the target tends to increase heritabilities (on a mean basis) within subregions, essentially because the genotype × subregion interaction variance becomes a genetic variance component. On the other hand, subdivision of resources will leave only a limited number of trials per subregion, thus decreasing heritabilities within subregions (Atlin et al., 2000). This is why selection based on subregion means alone is not necessarily a good strategy.

Regional trial networks designed to provide cultivar recommendations present breeders with a similar dilemma. If recommendations are based on overall performance in the target region, cultivars with good local adaptation in one or more subregions may go unnoticed, resulting in suboptimal cultivar recommendations. Conversely, an attempt to detect local adaptation by subdivision of the target may be compromised by a limited number of trials remaining per subregion.

The common view seems to be that there are basically two alternative approaches to estimation, depending on whether one strives for broad or local adaptation: (i) if the objective is broad adaptation, ignore subregions and use all data to make selections or recommendations based on overall performance in the target region; (ii) if the objective is local adaptation, use only data from the targeted subregion for selection or recommendation. An associated notion is that selecting or recommending for local adaptation requires almost the same amount of resources within each subregion as would be needed for assessing broad adaptation with the same accuracy. This notion assumes that subregions are not substantially differentiated and that local adaptation can be detected only on the basis of data from the targeted subregion (Comstock and Moll, 1963; Talbot, 1997; Atlin et al., 2000). More often than not, the result of this common view has been that global adaptation is favored over local adaptation. This is usually a reasonable choice if one considers only the two alternatives described above.

The objective of this paper is to show that accuracy of yield estimates can be increased both for global and for local adaptation, if a slightly modified route of analysis is followed. The suggestion is to (i) always contemplate a subdivision of the target and (ii) always use all the data, employing a suitable weighting scheme, based on genetic variances and covariances among subregions, no matter whether broad or local adaptation is the objective. We show how standard mixed model procedures (best linear unbiased prediction) can be used for this task. The method is developed by initially considering balanced data and a two-step approach. This simplifies the exposition and makes key features easier to appreciate. Subsequently, we will stress that a restricted maxi-

Bioinformatics Unit, Univ. of Hohenheim, Fruwirthstrasse 23, 70599 Stuttgart, Germany. Received 2 July 2004. Crop Breeding, Genetics & Cytology. *Corresponding author (piepho@uni-hohenheim.de).

**Abbreviations:** ANOVA, analysis of variance; BLUP, best linear unbiased prediction; MSEP, mean squared error of prediction; REML, restricted maximum likelihood.

mum likelihood (REML)-based mixed model analysis demonstrates its full power with unbalanced data and we will explain how replicate data, balanced or unbalanced, can be dealt with in a single-step analysis. Some easy-to-use SAS code, which also works for unbalanced data, is presented in an appendix.

## THEORY

### Subdivision for Global Adaptation

If a random sample of locations is used for yield testing in each subregion, the resulting data may be regarded as a random stratified sample, with strata corresponding to subregions. As is well known from the theory of survey sampling (Kish, 1965), stratification of a target region is beneficial providing there is heterogeneity between subregions, governed by environmental factors such as climate, soil type, or topography. Gain in accuracy is largest in the one extreme (but hypothetical) case that all heterogeneity occurs between subregions, while there is complete homogeneity within subregions. In this extreme case, a single location per subregion would suffice to assess the expected cultivar yield per subregion. Conversely, in the other extreme case, where all heterogeneity occurs within subregions and none among subregions, stratification is not beneficial.

In stratified samples, the overall mean is estimated by a weighted mean of the subregion means, with the growing area per subregion used as a weight. To see that stratification is beneficial, again consider the extreme case where there is no heterogeneity within subregions, but considerable heterogeneity between subregions. The weighted mean based on a stratified sample, with growing areas used as weights, will have a variance of zero, while the simple mean based on an unstratified sample will have a variance depending on the heterogeneity between subregions. The only additional prerequisite for a stratified estimate of overall performance is that the growing areas per subregion must be available.

### Estimation of Local Adaptation

Local adaptation in a subregion is usually assessed by analyzing data only from the targeted subregion (Talbot, 1997). It is often true, however, that some of the neighboring subregions are agroecologically similar to the subregion of interest. Thus, yield data from neighboring subregions may be exploited to improve yield estimates for the subregion of interest. A natural approach is to compute a weighted mean of mean yields in the targeted subregion and the neighboring subregions, with weights depending on the similarity between subregions and the number of trials per subregion. In fact, with a weighting approach, one could use data from all subregions for estimation in the targeted subregion, providing the availability of weights that are optimal or near-optimal in terms of the error of prediction for the targeted subregion. We will show, subsequently, that best linear unbiased prediction (BLUP; Searle et al., 1992) is the method of choice for this task.

### Mixed Model for Subdivided Target Regions

Our basic mixed model is

$$y_{irjk} = \mu_r + g_{ir} + e_{irjk}, \qquad [1]$$

where $\mu_r$ = expected value in the $r$th region, $g_{ir}$ = random genotypic value of $i$th genotype in $r$th subregion and $e_{irjk}$ = random environmental deviation of the $i$th genotype in the $k$th year and in the $j$th location within $r$th subregion.

The environmental deviation will be modeled by a standard partition of the form

$$e_{irjk} = Y_k + L(S)_{j(r)} + (SY)_{rk} + LY(S)_{jk(r)} + (GY)_{ik} +$$
$$GL(S)_{ij(r)} + (GSY)_{irk} + GLY(S)_{ijk(r)}, \qquad [2]$$

where $Y_k$ = main effect of $k$th year, $L(S)_{j(r)}$ = effect of $j$th location nested within $r$th subregion, $(SY)_{rk}$ = $rk$th subregion $\times$ year interaction, $LY(S)_{jk(r)}$ = $jk$th location $\times$ year interaction nested within $r$th subregion, $(GY)_{ik}$ = $ik$th genotype $\times$ year interaction, $GL(S)_{ij(r)}$ = $ij$th genotype $\times$ location interaction nested within $r$th subregion, $(GSY)_{irk}$ = $irk$th genotype $\times$ subregion $\times$ year interaction, $GLY(S)_{ijk(r)}$ = $ijk$th genotype $\times$ location $\times$ year interaction nested within $r$th subregion (includes error of a treatment mean). All effects appearing in $e_{irjk}$ are assumed to be independent homoscedastic normal deviates with zero mean.

The genetic effect may be further partitioned as

$$g_{ir} = G_i + (GS)_{ir}, \qquad [3]$$

where $G_i$ is a main effect for the $i$th genotype and $(GS)_{ir}$ is the $ir$th genotype $\times$ subregion interaction, assuming that $G_i$ and $(GS)_{ir}$ are independent homoscedastic normal deviates. This model implies that the variance-covariance model for $\boldsymbol{g}_i = (g_{i1}, g_{i2}, …, g_{im})'$, where $m$ is the number of subregions, has the compound symmetry structure, i.e.,

$$\mathrm{var}(\boldsymbol{g}_i) = \boldsymbol{\Sigma}_g = \boldsymbol{J}_m\sigma^2_G + \boldsymbol{I}_m\sigma^2_{GS}, \qquad [4]$$

where $\boldsymbol{J}_m$ is an $m \times m$ matrix of ones everywhere, $\boldsymbol{I}_m$ is an $m$-dimensional identity matrix, and $\sigma^2_G$ and $\sigma^2_{GS}$ are the variances of $G_i$ and $(GS)_{ir}$, respectively. Under the compound symmetry model, genetic variances are the same in each subregion, and genetic covariances (and correlations) are the same for each pair of subregions. While this assumption may be useful in simple settings, more diverse settings call for more refined modeling. Specifically, some pairs of subregions may be more alike than others, requiring heterogeneity of covariances to be allowed for. Also, there may be heterogeneity of genetic variance between subregions. Many extensions of the ANOVA-type model, Eq. [3], have been proposed, which can be used for modeling $\boldsymbol{\Sigma}_g$ (Piepho, 1998, 1999; Smith et al., 2001). Here, we will confine attention to the compound symmetry model in order to facilitate comparison to other methods. It is stressed, however, that quite frequently more complex variance–covariance structures are needed.

Analysis of regional yield trials should be based on replicate data, using the model described above. This allows exploiting regional subdivision for estimation of both local and global adaptation in an optimal way. Following this approach, estimates of $g_{ir}$ may be obtained by BLUP using standard procedures (Searle et al., 1992). Some sample code for SAS is given in Appendix A. This single-step analysis may be contrasted to a two-step analysis, in which genotypic means per subregion are estimated in the first step. These means are then subjected to a mixed model analysis to obtain BLUPs of $g_{ir}$ in the second step. In balanced settings, both procedures yield identical results, while with unbalanced data, results differ and the REML-based single-step analysis is to be preferred. To study the properties of the BLUP procedure, it is more convenient to use the two-step approach and restrict attention to the balanced set-up, and this will be done subsequently.

### Implied Model for Subregion Means

For demonstration purposes, we will assume here that the series of trials in each subregion is balanced in the following sense: On the basis of the mixed model described in the preced-

ing section, taking genotypic effects $g_{ir}$ fixed and environmental effects random, the means of all cultivars in a subregion are homoscedastic, i.e., they all have same variance. In addition, the means for all pairs of cultivars have the same covariance. This assumption is made here mainly to simplify the comparison of different estimators and to gain some insight into their statistical properties.

Let $y_{ir}$ denote the mean for the $i$th genotype in the $r$th subregion. We can assume the model

$$y_{ir} = \mu_r + g_{ir} + e_{ir}, \qquad [5]$$

where $e_{ir}$ is the error associated with $y_{ir}$. We find, conditioning on the genetic effects, that

$$\mathrm{E}(\boldsymbol{y}_i) = \boldsymbol{\mu} + \boldsymbol{g}_i,$$
$$\mathrm{E}(\boldsymbol{e}_i) = \boldsymbol{0}, \qquad [6]$$
$$\mathrm{var}(\boldsymbol{e}_i) = \boldsymbol{\Sigma}_{e1} + \boldsymbol{\Sigma}_{e2}, \text{ and}$$
$$\mathrm{cov}(\boldsymbol{e}_i, \boldsymbol{e}_{i'}) = \boldsymbol{\Sigma}_{e2} \ (i \neq i'),$$

where $\boldsymbol{y}_i = (y_{i1}, y_{i2}, ..., y_{im})'$, $\boldsymbol{\mu} = (\mu_1, \mu_2, ..., \mu_m)'$, $\boldsymbol{g}_i = (g_{i1}, g_{i2}, ..., g_{im})'$ and $\boldsymbol{e}_i = (e_{i1}, e_{i2}, ..., e_{im})'$. Now taking genetic effects random, we have $\mathrm{E}(\boldsymbol{g}_i) = \boldsymbol{0}$, $\mathrm{var}(\boldsymbol{g}_i) = \boldsymbol{\Sigma}_g$, and $\mathrm{cov}(\boldsymbol{g}_i, \boldsymbol{e}_i) = \mathrm{cov}(\boldsymbol{g}_i, \boldsymbol{e}_{i'}) = \boldsymbol{0}$. For $\boldsymbol{\Sigma}_g$, one may assume the compound symmetry structure, Eq. [4], or some more general model.

For illustration, consider the special case that the design is completely balanced, i.e., in each of the $m$ subregions the $n$ genotypes are tested in $l$ locations and $y$ years. In this case, $\boldsymbol{\Sigma}_{e1}$ has diagonal elements $\sigma_{GY}^2/y + \sigma_{GSY}^2/y + \sigma_{GSL}^2/l + \sigma_{GSLY}^2/(ly)$ and off-diagonal elements $\sigma_{GY}^2/y$, while $\boldsymbol{\Sigma}_{e2}$ has diagonal elements $\sigma_Y^2/y + \sigma_{SY}^2/y + \sigma_{SL}^2/l + \sigma_{SLY}^2/(ly)$ and off-diagonal elements $\sigma_Y^2/y$. It should be stressed, however, that Eq. [6] is more broadly applicable, e.g., when the number of locations differs among years or among subregions or both, when only some locations are used for several years, while others are exchanged every year, or when the number of years is not the same for each subregion. Under the assumed model in Eq. [1] and [2], the only prerequisite for the validity of Eq. [6] is that the same set of $n$ genotypes is tested in each trial.

## Estimation—Local and Global BLUP

On the basis of regional subdivision, the genotypic value in the target region can be expressed as

$$g_i = a_1 g_{i1} + a_2 g_{i2} + a_3 g_{i3} + ... + a_m g_{im} = \boldsymbol{a}'\boldsymbol{g}_i, \qquad [7]$$

where $a_r$ ($r = 1, ..., m$) are the relative growing areas in the $m$ subregions (expressed as proportions) and $\boldsymbol{a} = (a_1, a_2, ..., a_m)'$. We wish to either estimate $g_i$ (for assessing global adaptation) or $g_{ir}$ (for assessing local adaptation in the $r$th subregion). The estimable function of interest is of the general form

$$L_i = \boldsymbol{v}'\boldsymbol{g}_i \qquad [8]$$

with $\boldsymbol{v} = \boldsymbol{a}$ for $g_i$ (global adaptation) and $\boldsymbol{v} = \boldsymbol{u}_r$ for $g_{ir}$ (local adaptation), where $\boldsymbol{u}_r$ is a unit vector with $r$th element equal to unity and zeros elsewhere. For example, when there are five subregions, and an estimator is needed for the second region, we set $\boldsymbol{v} = \boldsymbol{u}_r = (0, 1, 0, 0, 0)'$. We consider estimators of $L_i$, which are of the form

$$\tilde{L}_i = \boldsymbol{w}'(\boldsymbol{y}_i - \boldsymbol{\mu}), \qquad [9]$$

where $\boldsymbol{w}$ are suitably chosen weights and the tilda indicates an estimator. Our preferred estimator is

$$\tilde{L}_i^{opt} = \boldsymbol{v}'\tilde{\boldsymbol{g}}_i \text{ with} \qquad [10]$$

$$\tilde{\boldsymbol{g}}_i = \boldsymbol{W}(\boldsymbol{y}_i - \boldsymbol{\mu}), \text{ where} \qquad [11]$$

$$\boldsymbol{W} = \boldsymbol{\Sigma}_g(\boldsymbol{\Sigma}_g + \boldsymbol{\Sigma}_{e1})^{-1}, \qquad [12]$$

because it involves an optimal weighting scheme, as detailed below. Estimator [10] is a special case of [9] with weights $\boldsymbol{w}' = \boldsymbol{v}'\boldsymbol{W}$. These weights are optimal and follow from BLUP theory (Searle et al., 1992). The estimator [11] is essentially the BLUP of genetic effects $\boldsymbol{g}_i$ (see Appendix B). When $\boldsymbol{v} = \boldsymbol{u}_r$, we will refer to Eq. [10] as local BLUP, while when $\boldsymbol{v} = \boldsymbol{a}$, we refer to Eq. [10] as global BLUP.

In practice, variance components in $\boldsymbol{W}$ are unknown and need to be estimated. Plugging in estimates for parameters yields empirical BLUP with added uncertainty because of estimated weights. Providing parameters are estimated by REML using, e.g., a Newton-Raphson or Fisher scoring algorithm, uncertainty can be accounted for when computing approximate standard errors of BLUPs using the asymptotic dispersion matrix (Kackar and Harville, 1984). Generally, when a REML-based mixed model package such as MIXED is employed, the user need not worry about computation of weights $\boldsymbol{W}$: these will be computed automatically for the BLUP of $\boldsymbol{g}_i$ on the basis of the fitted variance–covariance model. In the case of global adaptation, $L_i = \boldsymbol{v}'\boldsymbol{g}_i$ can be estimated from the BLUP for $\boldsymbol{g}_i$, which may require additional computation following the run of the mixed model routine. Note that the weighting matrix $\boldsymbol{W}$ is analogous to the broad-sense heritability in the case of an undivided target region. Generally, to estimate $g_{ir}$ for a particular subregion, information on the $i$th genotype is used from all subregions. The extent to which the information from neighboring subregions is exploited is determined by $\boldsymbol{W}$ and depends on the heritabilities in the neighboring subregions and on the genetic and environmental correlations with the subregion of interest, as determined by the structures of $\boldsymbol{\Sigma}_g$ and $\boldsymbol{\Sigma}_{e1}$, respectively.

## Variance-Bias Trade-Off

The BLUP of $\boldsymbol{g}_i$ is unbiased in the sense that, across all genotypes, BLUP has the same expected value as $\boldsymbol{g}_i$ itself, i.e., $\mathrm{E}(\boldsymbol{g}_i) = \mathrm{E}[\mathrm{BLUP}(\boldsymbol{g}_i)] = 0$. Clearly, this expectation is unconditional (Searle et al., 1992, p. 269). By contrast, there is a bias conditional on the genotype, i.e., $\mathrm{E}(\mathrm{BLUP}(\boldsymbol{g}_i)|\boldsymbol{g}_i) \neq \boldsymbol{g}_i$. The bias may increase when incorporating data from neighboring subregions. This may be counterbalanced, however, by the reduced estimation variance because of the use of more data. Thus, it is usually beneficial to exploit data from neighboring subregions.

The purpose of this section is to shed more light on our proposed procedure, explaining how the gain of efficiency comes about. The variance-bias trade-off can be most conveniently studied by considering pairwise differences among genotypic effect estimates. Note that the ranking of genotypes is fully determined by the set of all pairwise differences. The difference of two genotypes $i$ and $i'$ is given by

$$\delta_{ii'} = \boldsymbol{v}'(\boldsymbol{g}_i - \boldsymbol{g}_{i'}). \qquad [13]$$

For simplicity, we will drop the indices on $\delta$, i.e., we set $\delta \equiv \delta_{ii'}$. The difference $\delta$ is estimated by BLUP according to

$$\tilde{\delta} = \boldsymbol{v}'(\tilde{\boldsymbol{g}}_i - \tilde{\boldsymbol{g}}_{i'}). \qquad [14]$$

This estimator is biased, since for given genotypes $i$ and $i'$

$$E(\delta - \tilde{\delta}|\boldsymbol{g}_i, \boldsymbol{g}_{i'}) = Bias = (\boldsymbol{v} - \boldsymbol{w})'(\boldsymbol{g}_i - \boldsymbol{g}_{i'}), \qquad [15]$$

where $\boldsymbol{w}' = \boldsymbol{v}'\boldsymbol{W}$. Thus, exploiting information from neighboring subregions introduces a bias. This may be more than offset, however, by a reduction in the estimation variance.

The common criterion for combining bias and variance is the mean squared error of prediction (MSEP), which in the case at hand is

CROP SCIENCE, VOL. 45, MAY–JUNE 2005

**Table 1. Variance component estimates for three multienvironment trials (reproduced from Atlin et al., 2000), expressed as proportion of the phenotypic variance $\tilde{\sigma}_P^2 = \tilde{\sigma}_G^2 + \tilde{\sigma}_{GL}^2 + \tilde{\sigma}_{GY}^2 + \tilde{\sigma}_{GYL}^2 + \tilde{\sigma}_E^2$.**

| | | Variance component estimates | | | | |
|---|---|---|---|---|---|---|
| | | $\tilde{\sigma}_G^2$ | $\tilde{\sigma}_{GL}^2$ | $\tilde{\sigma}_{GY}^2$ | $\tilde{\sigma}_{GYL}^2$ | $\tilde{\sigma}_E^2$ |
| Crop | Region | Genotype | Genotype × location | Genotype × year | Genotype × year × location | Error |
| Winter wheat | Eastern Canada | 0.36 | 0.03 | 0.02 | 0.29 | 0.30 |
| Spring wheat | Western Canada | 0.29 | 0.11 | 0.02 | 0.27 | 0.31 |
| Spring wheat | Australia | 0.05 | 0.13 | 0.11 | 0.12 | 0.58 |

$$MSEP = E_g[(\delta - \tilde{\delta})^2 | \boldsymbol{g}_i, \boldsymbol{g}_{i'}] = E_g[Bias^2 | \boldsymbol{g}_i, \boldsymbol{g}_{i'}] + Var, \quad [16]$$

where

$$E_g[Bias^2 | \boldsymbol{g}_i, \boldsymbol{g}_{i'}] = 2(\boldsymbol{v} - \boldsymbol{w})' \Sigma_g (\boldsymbol{v} - \boldsymbol{w}) \quad [17]$$

and

$$Var = var(\delta | \boldsymbol{g}_i, \boldsymbol{g}_{i'}) = 2\boldsymbol{w}' \Sigma_{e1} \boldsymbol{w}. \quad [18]$$

The subscripted $g$ indicates that expectations are with respect to genotype pairs. It is seen from Eq. [16] that the MSEP depends on both bias and variance. The weights $\boldsymbol{W}$ in $\boldsymbol{w}' = \boldsymbol{v}'\boldsymbol{W}$ are chosen so as to minimize the MSEP. Thus, the optimal weights strike the best balance between bias and variance. Specifically, absence of bias is not a requirement: some bias can be tolerated, providing the MSEP is smaller than for an unbiased estimator. Not only does BLUP minimize the MSEP, but it also maximizes the response to selection in variance component models (Searle et al., 1992; but see Portnoy, 1982).

### Gain from Selection

We consider response to selection based on the estimator $\tilde{L}_i$. In the special case that $\boldsymbol{w}' = \boldsymbol{v}'\boldsymbol{W}$, where $L_i = \boldsymbol{v}\boldsymbol{g}_i$, this is our proposed optimal estimator $\tilde{L}_i^{opt}$. For generality, we regard $\tilde{L}_i$ as an indirect trait and formulate the selection response as a correlated response to selection (Falconer and Mackay, 2001; Atlin et al., 2000). Selection for a direct trait is included as a special case with genetic correlation equal to unity. We have

$$var(L_i) = \boldsymbol{v}'\Sigma_g \boldsymbol{v}, \quad [19]$$

$$var(\tilde{L}_i) = \boldsymbol{w}'(\Sigma_g + \Sigma_{e1})\boldsymbol{w}, \quad [20]$$

$$cov(\tilde{L}_i, L_i) = \boldsymbol{v}'\Sigma_g \boldsymbol{w}. \quad [21]$$

The genetic correlation between $L_i$ and $\tilde{L}_i$ is

$$\rho_g = \frac{\boldsymbol{v}'\Sigma_g \boldsymbol{w}}{\sqrt{(\boldsymbol{v}'\Sigma_g \boldsymbol{v})(\boldsymbol{w}'\Sigma_g \boldsymbol{w})}}. \quad [22]$$

The heritability of $\tilde{L}_i$ equals

$$h^2 = \frac{\boldsymbol{w}'\Sigma_g \boldsymbol{w}}{\boldsymbol{w}'(\Sigma_g + \Sigma_{e1})\boldsymbol{w}}. \quad [23]$$

The correlated response to selection is

$$R = i\rho_g h\sqrt{var(L_i)}, \quad [24]$$

where $i$ is the selection intensity (Falconer and Mackay, 2001). For evaluation of different methods it is convenient to consider the ratio of $R$-values, since the selection intensity as well as $\sqrt{[var(L_i)]}$ cancel out. Note that all results in this section are rather general in that they do not require a special structure for $\Sigma_g$ or $\Sigma_{e1}$, such as compound symmetric.

### Variance Component Estimates

To illustrate our procedure, we use published variance component estimates for three multienvironment trials with wheat (*Triticum aestivum* L.), which are reproduced in Table 1. These estimates were also used by Atlin et al. (2000) to demonstrate their approach for estimation in subdivided target regions. We used the variance components in Table 1 to assign values to the variance components in our mixed model (Eq. [1], [2], and [3]) (Table 2). Following Atlin et al. (2000), the number of replications per trial was three, the total number of locations in the trial network was set equal to 12, and locations were evenly split among subregions. For simplicity, the variance of $(GSY)_{irk}$ was set to zero. Since our model is based on trial means, we set the variance of the residual effect $GLY(S)_{ijk(r)}$ equal to $\tilde{\sigma}_{GLY}^2 + \tilde{\sigma}_E^2/3$, with variance components $\tilde{\sigma}_{GYL}^2$ and $\tilde{\sigma}_E^2$ taken from Table 1.

The compound symmetry structure was used for $\Sigma_g$ (Eq. [4]). Assuming an orthogonal design per subregion, the diagonal elements of $\Sigma_{e1}$ were equated to $\sigma_{GY}^2/y + \sigma_{GSY}^2/y + \sigma_{GSL}^2/l + \sigma_{GSLY}^2/(ly)$, where $l$ is the number of locations per subregion and $y$ is the number of years. The off-diagonal elements were set to $\sigma_{GY}^2/y$.

## RESULTS

### Estimation for Global Adaptation

To study the gain from stratification, we assumed that the target can be subdivided into subregions of equal size, i.e., the relative areas are $a_r = 1/m$. Estimation of the overall mean $(g_i)$ in the target region by global BLUP was performed by setting $\boldsymbol{v} = \boldsymbol{a} = (a_1, a_2)'$ and $\boldsymbol{w} = \boldsymbol{v}'\boldsymbol{W}$. The response to selection based on our method is denoted as $R_1$. For comparison, we computed the response to selection assuming that stratification is ignored and locations are a fully random sample from the target region ($R_0$). Ratios $R_1/R_0$ are reported in Table 3. The results show that one can only win by stratification and that the gain is largest when the genotype × subregion interaction variance is large relative to the variance of the genetic main effect. In the most favorable case reported in Table 3, stratification results in a 20% improvement in the response to selection.

We studied the effect of unequal subregion areas $a_r$ assuming that the target region is subdivided into two subregions with $a_1 = q$ and $a_2 = 1 - q$, where $q$ takes on

**Table 2. Variance components for mixed model given by Eq. [1], [2], and [3] as derived from estimates in Table 1. $p$ is the proportion of $\tilde{\sigma}_{GL}^2$ assigned to $\sigma_{GS}^2$.**

| Effect in model (2) | Variance | Value assigned from estimates in Table 1 |
|---|---|---|
| $G_i$ | $\sigma_G^2$ | $\tilde{\sigma}_G^2$ |
| $(GS)_{ir}$ | $\sigma_{GS}^2$ | $p\tilde{\sigma}_{GL}^2$ |
| $(GY)_{ik}$ | $\sigma_{GY}^2$ | $\tilde{\sigma}_{GY}^2$ |
| $GL(S)_{ij(r)}$ | $\sigma_{GSL}^2$ | $(1 - p)\tilde{\sigma}_{GL}^2$ |
| $(GSY)_{irk}$ | $\sigma_{GSY}^2$ | 0 |
| $GLY(S)_{ijk(r)}$ | $\sigma_{GSLY}^2$ | $\tilde{\sigma}_{GYL}^2 + \tilde{\sigma}_E^2/3$ |

**Table 3. Selection for global adaptation. Ratio $R_1/R_0$, where $R_0$ = response to selection for $g_i$ ignoring stratification and $R_1$ = response to selection for $g_i$ using global BLUP. Analysis based on three replications per trial. $p$ is the proportion of $\tilde{\sigma}^2_{GL}$ assigned to $\sigma^2_{GS}$ (see Table 2).**

| | Value of ratio $R_1/R_0$ | | | | | |
|---|---|---|---|---|---|---|
| | 1 yr | | | 2 yr | | |
| $p$ | $l = 2$ | $l = 4$ | $l = 6$ | $l = 2$ | $l = 4$ | $l = 6$ |
| Winter wheat in eastern Canada ($\tilde{\sigma}^2_{GL}/\tilde{\sigma}^2_G = 0.08$) | | | | | | |
| 0.1 | 1.000 | 1.000 | 1.001 | 1.000 | 1.000 | 1.000 |
| 0.3 | 1.001 | 1.001 | 1.002 | 1.001 | 1.001 | 1.001 |
| 0.5 | 1.002 | 1.002 | 1.003 | 1.002 | 1.002 | 1.002 |
| Spring wheat in western Canada ($\tilde{\sigma}^2_{GL}/\tilde{\sigma}^2_G = 0.37$) | | | | | | |
| 0.1 | 1.002 | 1.002 | 1.003 | 1.002 | 1.002 | 1.002 |
| 0.3 | 1.005 | 1.007 | 1.008 | 1.005 | 1.006 | 1.007 |
| 0.5 | 1.009 | 1.011 | 1.014 | 1.009 | 1.010 | 1.011 |
| Spring wheat in Australia ($\tilde{\sigma}^2_{GL}/\tilde{\sigma}^2_G = 2.67$) | | | | | | |
| 0.1 | 1.019 | 1.034 | 1.049 | 1.017 | 1.030 | 1.041 |
| 0.3 | 1.054 | 1.096 | 1.133 | 1.050 | 1.082 | 1.111 |
| 0.5 | 1.088 | 1.151 | 1.204 | 1.081 | 1.128 | 1.168 |

the values 0.3, 0.1, and 0.01. Estimation of the overall mean ($g_i$) in the target region by global BLUP was implemented by setting $\boldsymbol{v} = \boldsymbol{a} = (a_1, a_2)'$ and $\boldsymbol{w} = \boldsymbol{v}'\boldsymbol{W}$. The associated response to selection is denoted as $R_1$. For comparison, the response to selection based on the genotypic main effect $G_i$ using the mixed model in Eq. [1], [2], and [3] was computed as proposed by Atlin et al. (2000). This is denoted as $R_2$. The method corresponds to selection using a simple mean of yields in the two subregions, i.e., $\boldsymbol{w} = (0.5, 0.5)'$. It should be stressed that both estimators exploit stratification of the target region. Differences in performance are therefore solely due to the contrasting weighting schemes. Of course, when the relative areas are the same ($q = 0.5$), both methods yield identical results. The ratio $R_1/R_2$ is reported in Table 4. The more substantial the genotype × location interaction in the target region and the less equal the relative growing areas, the more pronounced is the gain from accounting for unequal subregion areas by global BLUP.

**Table 4. Selection for global adaptation. Ratio $R_1/R_2$, where $R_1$ = response to selection for $g_i$ using global BLUP and $R_2$ = response to selection for genotypic main effect $G_i$ as proposed by Atlin et al. (2000). Analysis assumes $l = 6$ locations per subregion, $s = 2$ subregions, and three replications per trial. Relative areas $a_1 = q$ and $a_2 = 1 - q$. $p$ is the proportion of $\tilde{\sigma}^2_{GL}$ assigned to $\sigma^2_{GS}$ (see Table 2).**

| | Value of ratio $R_1/R_2$ | | | | | |
|---|---|---|---|---|---|---|
| | 1 yr | | | 2 yr | | |
| $p$ | $q = 0.3$ | $q = 0.1$ | $q = 0.01$ | $q = 0.3$ | $q = 0.1$ | $q = 0.01$ |
| Winter wheat in eastern Canada ($\tilde{\sigma}^2_{GL}/\tilde{\sigma}^2_G = 0.08$) | | | | | | |
| 0.1 | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 |
| 0.3 | 1.000 | 1.001 | 1.001 | 1.000 | 1.001 | 1.001 |
| 0.5 | 1.000 | 1.001 | 1.002 | 1.001 | 1.002 | 1.003 |
| Spring wheat in western Canada ($\tilde{\sigma}^2_{GL}/\tilde{\sigma}^2_G = 0.37$) | | | | | | |
| 0.1 | 1.000 | 1.001 | 1.001 | 1.000 | 1.001 | 1.002 |
| 0.3 | 1.002 | 1.006 | 1.009 | 1.002 | 1.008 | 1.012 |
| 0.5 | 1.004 | 1.014 | 1.021 | 1.004 | 1.017 | 1.026 |
| Spring wheat in Australia ($\tilde{\sigma}^2_{GL}/\tilde{\sigma}^2_G = 2.67$) | | | | | | |
| 0.1 | 1.005 | 1.020 | 1.030 | 1.005 | 1.019 | 1.029 |
| 0.3 | 1.025 | 1.096 | 1.141 | 1.023 | 1.087 | 1.128 |
| 0.5 | 1.043 | 1.161 | 1.233 | 1.037 | 1.142 | 1.207 |

## Estimation for Local Adaptation

Local BLUPs of the subregion mean were obtained by setting $\boldsymbol{v} = \boldsymbol{u}_r$ and $\boldsymbol{w} = \boldsymbol{v}'\boldsymbol{W}$. The response to selection based on local BLUP is denoted as $R_3$. For comparison, selection using the subregion mean was implemented by setting $\boldsymbol{w} = \boldsymbol{v} = \boldsymbol{u}_r$. The response to selection by this method is denoted as $R_4$. The response to selection based on the unweighted global mean using the mixed model in Eq. [1], [2], and [3] was computed as proposed by Atlin et al. (2000). This is denoted as $R_2$. The ratios $R_2/R_4$ and $R_3/R_4$ are reported in Table 5.

The most important result is that local BLUP always does better than selection based on a subregion mean or than selection based on the global mean as proposed by Atlin et al. (2000). In some cases, the differences are quite marked; in others, differences are minor. When the genotype × subregion interaction is substantial (spring wheat in Australia), the global mean performs poorly, while local BLUP is slightly better than the subregion mean. When the genotype × subregion interaction is small (winter wheat in eastern Canada), the global mean almost always outperforms the subregion mean, but is itself outperformed by local BLUP.

Table 6 shows relative weights for the targeted subregion and the other subregions. The latter are equal for all subregions because of balancedness of the design and the compound symmetry model for $\boldsymbol{\Sigma}_g$. The larger the number of years and locations, the smaller will be the magnitude of $\boldsymbol{\Sigma}_{e1}$, thus increasing the weight for the targeted subregion relative to the other subregions. In the limit as $\boldsymbol{\Sigma}_{e1}$ tends to zero, all weight will be on the target subregion, no matter what is the assumed structure for $\boldsymbol{\Sigma}_g$. Also, the smaller the genetic correlation,

**Table 5. Selection for local adaptation. Ratios $R_2/R_4$ and $R_3/R_4$, where $R_2$ = response to selection for genotypic main effect $G_i$ as proposed by Atlin et al. (2000), $R_3$ = response to selection for $g_{ir}$ using local BLUP, and $R_4$ = response to selection based on subregion mean alone. Analysis based on three replications per trial. $p$ is the proportion of $\tilde{\sigma}^2_{GL}$ assigned to $\sigma^2_{GS}$ (see Table 2).**

| | | Value of ratio $R_2/R_4$ or $R_3/R_4$ | | | | | |
|---|---|---|---|---|---|---|---|
| | | 1 yr | | | 2 yr | | |
| $p$ | Ratio | $l = 2$ | $l = 4$ | $l = 6$ | $l = 2$ | $l = 4$ | $l = 6$ |
| Winter wheat in eastern Canada ($\tilde{\sigma}^2_{GL}/\tilde{\sigma}^2_G = 0.08$) | | | | | | | |
| 0.1 | $R_2/R_4$ | 1.184 | 1.074 | 1.034 | 1.106 | 1.040 | 1.017 |
| | $R_3/R_4$ | 1.185 | 1.077 | 1.038 | 1.108 | 1.043 | 1.021 |
| 0.3 | $R_2/R_4$ | 1.167 | 1.059 | 1.020 | 1.091 | 1.027 | 1.003 |
| | $R_3/R_4$ | 1.172 | 1.069 | 1.033 | 1.096 | 1.037 | 1.017 |
| 0.5 | $R_2/R_4$ | 1.150 | 1.045 | 1.006 | 1.076 | 1.014 | 0.990 |
| | $R_3/R_4$ | 1.160 | 1.062 | 1.029 | 1.086 | 1.031 | 1.014 |
| Spring wheat in western Canada ($\tilde{\sigma}^2_{GL}/\tilde{\sigma}^2_G = 0.37$) | | | | | | | |
| 0.1 | $R_2/R_4$ | 1.214 | 1.074 | 1.023 | 1.137 | 1.041 | 1.005 |
| | $R_3/R_4$ | 1.222 | 1.089 | 1.043 | 1.146 | 1.056 | 1.026 |
| 0.3 | $R_2/R_4$ | 1.139 | 1.012 | 0.962 | 1.070 | 0.983 | 0.947 |
| | $R_3/R_4$ | 1.169 | 1.061 | 1.027 | 1.101 | 1.034 | 1.014 |
| 0.5 | $R_2/R_4$ | 1.074 | 0.957 | 0.909 | 1.010 | 0.931 | 0.896 |
| | $R_3/R_4$ | 1.128 | 1.042 | 1.017 | 1.070 | 1.020 | 1.008 |
| Spring wheat in Australia ($\tilde{\sigma}^2_{GL}/\tilde{\sigma}^2_G = 2.67$) | | | | | | | |
| 0.1 | $R_2/R_4$ | 1.111 | 0.940 | 0.873 | 1.111 | 0.943 | 0.876 |
| | $R_3/R_4$ | 1.185 | 1.053 | 1.017 | 1.185 | 1.056 | 1.019 |
| 0.3 | $R_2/R_4$ | 0.795 | 0.677 | 0.628 | 0.800 | 0.685 | 0.635 |
| | $R_3/R_4$ | 1.036 | 1.001 | 1.001 | 1.039 | 1.003 | 1.000 |
| 0.5 | $R_2/R_4$ | 0.622 | 0.532 | 0.493 | 0.629 | 0.542 | 0.501 |
| | $R_3/R_4$ | 1.003 | 1.004 | 1.010 | 1.005 | 1.001 | 1.004 |

**Table 6. Selection for local adaptation. Standardized weights $w'$/($w'1$) for local BLUP of $g_{ir}$. Analysis based on three replications per trial. $p$ is the proportion of $\tilde{\sigma}^2_{GL}$ assigned to $\sigma^2_{GS}$ (see Table 2).**

| | | Standardized weights $w'$/($w'1$) | | | | | |
| | | 1 yr | | | 2 yr | | |
| $p$ | Subregion | $l = 2$ | $l = 4$ | $l = 6$ | $l = 2$ | $l = 4$ | $l = 6$ |
|---|---|---|---|---|---|---|---|
| | Winter wheat in eastern Canada ($\tilde{\sigma}^2_{GL}/\tilde{\sigma}^2_G = 0.08$) | | | | | | |
| 0.1 | Target | 0.180 | 0.355 | 0.524 | 0.190 | 0.370 | 0.540 |
| | Other | 0.164 | 0.323 | 0.476 | 0.160 | 0.315 | 0.460 |
| 0.3 | Target | 0.207 | 0.395 | 0.567 | 0.236 | 0.436 | 0.608 |
| | Other | 0.159 | 0.302 | 0.433 | 0.153 | 0.282 | 0.392 |
| 0.5 | Target | 0.233 | 0.432 | 0.604 | 0.279 | 0.493 | 0.661 |
| | Other | 0.153 | 0.284 | 0.396 | 0.144 | 0.254 | 0.339 |
| | Spring wheat in western Canada ($\tilde{\sigma}^2_{GL}/\tilde{\sigma}^2_G = 0.37$) | | | | | | |
| 0.1 | Target | 0.211 | 0.402 | 0.574 | 0.233 | 0.433 | 0.605 |
| | Other | 0.158 | 0.299 | 0.426 | 0.153 | 0.284 | 0.395 |
| 0.3 | Target | 0.294 | 0.513 | 0.681 | 0.352 | 0.579 | 0.737 |
| | Other | 0.141 | 0.243 | 0.319 | 0.130 | 0.210 | 0.263 |
| 0.5 | Target | 0.396 | 0.600 | 0.756 | 0.454 | 0.682 | 0.816 |
| | Other | 0.126 | 0.200 | 0.244 | 0.109 | 0.159 | 0.184 |
| | Spring wheat in Australia ($\tilde{\sigma}^2_{GL}/\tilde{\sigma}^2_G = 2.67$) | | | | | | |
| 0.1 | Target | 0.347 | 0.598 | 0.775 | 0.347 | 0.592 | 0.764 |
| | Other | 0.131 | 0.201 | 0.225 | 0.131 | 0.204 | 0.236 |
| 0.3 | Target | 0.644 | 0.942 | 1.062 | 0.633 | 0.900 | 1.007 |
| | Other | 0.071 | 0.029 | −0.062 | 0.073 | 0.050 | −0.007 |
| 0.5 | Target | 0.875 | 1.141 | 1.189 | 0.847 | 1.062 | 1.103 |
| | Other | 0.025 | −0.071 | −0.189 | 0.031 | −0.031 | −0.103 |

i.e., the larger the diagonal elements of $\Sigma_g$ in relation to the off-diagonal elements, the higher the weights for the target subregion. When the genetic correlation is very low, as in the Australian wheat data, negative weights may occur for nontarget subregions. This is a result of the genotype $\times$ year interaction component, which introduces a positive environmental covariance in $\Sigma_{e1}$. The negative weights are usually small in absolute value, and virtually all the weight lies on the target subregion.

## DISCUSSION

Plant breeders and extension service personnel frequently consider a subdivision of a target region into smaller subregions, which in themselves are more homogeneous than the overall target region. Subdivision does not necessarily imply, however, that data from one subregion are not informative regarding another targeted subregion. In this paper, we have described a weighting scheme, based on standard mixed model procedures (BLUP), which allows information from different subregions to be efficiently combined. The weighting approach has been shown to be beneficial for two different objectives, the one being estimation of performance in a specific subregion, while the other is estimation of the global mean in the target region. For estimating the mean in a targeted subregion, local BLUP combines yield data from all subregions in an optimal way. If there is little information in the neighboring subregions because of large genotype $\times$ subregion interaction, the neighbors will be assigned small weights. Conversely, if the information content is high, weights assigned to neighbors will be relatively high. The BLUP procedure will yield the optimal weights from the variance components in a quasi-automatic fashion. For estimating the global mean, accuracy can be improved by stratification

into subregions and using growing areas per subregion as weights to combine local BLUPs for subregions into a global BLUP.

Curnow (1988) and Atlin et al. (2000) have made a very important contribution in demonstrating that information from neighboring subregions may be informative for a targeted subregion. Atlin et al. (2000) considered two alternatives, i.e., use all data to estimate the global mean, giving equal weight to all subregions, or use only the data from the targeted subregion. They showed that the local mean might be outperformed by the global mean if genetic correlations are high between subregions. In these cases, they propose not to subdivide the target region. Our approach obviates the need to choose between these two simple alternatives. The method always uses all data, giving different weight to the targeted subregion and its neighbors, with weights depending on the objective of estimation (local or global adaptation) and on the information provided by each subregion. Also, as opposed to the procedures considered by Curnow (1988) and Atlin et al. (2000), our method will always benefit from a subdivision of the target, providing there is heterogeneity among subregions and variance components are known.

Under the assumed model and providing known variance components, our weighted estimator will be optimal in the sense that it minimizes the MSEP and maximizes the expected response to selection. Specifically, local BLUP will give the best estimate of $g_{ir}$ and global BLUP will give the best estimate of $g_i$. If sample estimates are plugged in (empirical BLUP), some loss of accuracy will result, and optimality can no longer be guaranteed, though near-optimality is usually achieved. It may be a worthwhile strategy to use long-term data to obtain more accurate variance component estimates. This gain in accuracy will need to be balanced, however, against a potential long-term shift in variance components due to breeding progress as well as advances in agronomic practices. It would be worthwhile to conduct a simulation study on the effect of estimation error in the variance parameters on the response to selection. This will be the subject of future research.

When considering the relative merits of local and global BLUP, it is useful to distinguish two different objectives: (i) geographical placement of cultivars and (ii) selection of entries in a breeding program. Local BLUP will maximize the gain from selection for local adaptation. Thus, placement of cultivars should be based on local BLUP, if a meaningful subdivision of the target region is available and data permit reliable estimates of all variance components. By contrast, it is not so straightforward to decide, whether local or global BLUP is preferable for breeding purposes. While selecting for global adaptation may reduce the selection gain in a particular subregion, it does have the advantage of addressing a larger growing area. Therefore, breeding for local adaptation will be worthwhile only if the superiority of local response to selection more than compensates the loss from a reduction in growing area where the cultivar can be successfully marketed.

We have shown that global BLUP outperforms the

unweighted mean across subregions, when areas per subregion are unequal. Assuming compound symmetry and balanced data, the genetic correlation between the unweighted mean and the true mean effect in the target $(g_i)$ is

$$\rho_g = \sqrt{\frac{m\sigma_G^2 + \sigma_{GS}^2}{m(\sigma_G^2 + a'a\sigma_{GS}^2)}}. \qquad [25]$$

This will equal unity only when all subregions are of equal size, so that $a'a = m^{-1}$. The more heterogeneous the growing areas, the smaller the genetic correlation. This shows why, for unequally sized subregions, the unweighted analysis is suboptimal.

Our conceptual framework assumes that there is a larger target region, which may be subdivided into agroecologically distinct subregions. The demarcation between target regions (broad, global) and subregions (narrow, local) will depend on the type and scale of breeding application. For example, some breeders are part of a multinational (company or IARC) effort, whereas others are part of a national, state or provincial level program. What is a subregion to one is a broad region to another. Thus, the breeder will have to decide on a clear definition of the target region and subregions. In so doing, one will need to balance the geographic and climatological context against the organizational context. For a given latitudinal zone, there may be a very high degree of similarity between distant locations, but across longitudinal or elevational distances, closer subregions may rapidly become quite dissimilar. Subdivision will be most useful if homogeneity within subregions is maximized, i.e., genotype × location interaction within subregions is minimized, while genotype × subregion interaction is maximized. In this case, data from the targeted subregion are very informative, while information from neighboring subregions will be relatively small. Conversely, if subdivision mainly follows administrative boundaries, the subregions may not be very distinct agroecologically, thus increasing the information content from neighboring subregions. Clearly, agroecological factors are usually better criteria for subdivision than administrative boundaries. In either case, optimal weighting of information from targeted and neighboring subregions is crucial, and this may be achieved in a convenient way by BLUP.

Performance of our methodology critically depends on good estimates of the variance–covariance structure for genetic effects. Efficiency will be compromised if subdivision leads to subregions with no more than one or two test locations, especially when heterogeneity of covariance needs to be accounted for. Thus, it is a good strategy to have a larger number of test locations per subregion. The optimal number of locations per subregion will depend on a number of factors, such as the magnitude of and relationship among of variance components, the degree of agroecological differentiation between subregions, the objective of estimation (global or local adaptation), the number of years available for analysis, the design used with individual trials, etc. These factors are best studied by comprehensive simulation, and this will be the subject of future work.

Mixed model analysis of multienvironment trials requires random locations to allow broad inferences with respect to the target region. The requirement of a truly random sample of locations from the whole target may not always be easy to satisfy in practice, particularly when locations are selected to be representative. The notion of representativeness of certain locations usually implies a stratification of the target region. Thus, instead of selecting representative locations, one may subdivide the target region into (representative) subregions and then take truly random samples of locations per subregion. Breeders may be more comfortable with this type of restricted random sampling from subregions than a fully random sample of locations from the whole target region.

Our approach treats genotypes as a random factor. Many trial systems are set up to test elite genotypes, which have undergone several cycles of selection at the later stages of a breeding program or released cultivars. This raises the question concerning the population of entries to which results should apply. One possible view is that the population is defined by the potential set of entries that could have been obtained by the same breeding programs that generated the entries under consideration. The entries in the trials can be considered as a random sample of this hypothetical population. Clearly, the entries are not a random sample from the genotypes available at the beginning of the breeding process. Under a random genotypes model, it is also possible to incorporate pedigree information, using a model allowing for genetic correlation among related genotypes (Piepho and Pillen, 2004). Such models may be useful for multi-environment testing in complex breeding programs. Generally, estimates of genetic effects are typically more efficient under a random genotypes model than under a fixed effects model, providing the genetic variance components can be accurately estimated (Piepho, 1998; Smith et al., 2001).

It is worth mentioning that Atlin et al. (2000) do not explicitly use BLUP, and their development is basically in terms of simple genotype means, although they regard genotypes as random. Under their assumed model and balanced data setting, simple means and BLUPs will be perfectly linearly related, so the selection decision will be the same with either estimator. With unbalanced data and other variance–covariance structures, however, the two estimators will differ, and BLUP will typically be more efficient than simple means in such circumstances (Piepho, 1998).

Our BLUP procedure has two salient statistical features: shrinkage and weighting. For balanced data, shrinkage toward the mean is the same for each genotype, so shrinkage will not affect ranking compared with alternative estimators such as the mean per subregion. Thus, the differences in performance we found in comparison to other estimators were not due to shrinkage. The differences were solely due to the weighting scheme. With unbalanced data, shrinkage will differ among genotypes and so may affect ranking.

In this paper, we mainly focused on balanced data per subregion because this facilitated studying the statistical properties of our procedure. For the same reason, we considered a two-step approach, in which genotypic means per subregion are estimated in the first step and are then submitted to a BLUP procedure in the second step. In practice, a single-step procedure is more convenient and, in fact, easier to implement for routine use. Moreover, data will often be unbalanced. A single-step analysis of possibly unbalanced data is straightforward when a good mixed model package is used. All that needs to be done is to fit the mixed model outlined in Eq. [1] and [2] and let the package compute the BLUPs of genetic effects and linear functions thereof, depending on the choice of the vector $\boldsymbol{v}$ (see Appendix A).

In the example, we have used the compound symmetry structure in Eq. [4] to model the genetic variance-covariance structure $\boldsymbol{\Sigma}_g$. The model was used to facilitate comparison with the method of Atlin et al. (2000), which is based on this same structure. It should be stressed, however, that the compound symmetry model is rather restrictive in that variances are assumed to be the same for all subregions and that covariances are the same for each pair of subregions. It is now relatively easy to fit more general structures including heteroscedastic and factor-analytic, and experience shows that such models often fit considerably better than compound symmetry (Piepho, 1998; Smith et al., 2001). When fitting more complex models, one needs to balance increased realism against the need to estimate more parameters. In the extreme case of an unstructured model, $\boldsymbol{\Sigma}_g$ has $m(m+1)/2$ parameters, where $m$ is the number of subregions, and sample size may not even permit fitting of such complex models. According to the principle of parsimony, it may be preferable to fit simpler models, particularly when the sample size is limited, and there are established procedures for striking the balance between overly simplistic and overly complex models (Piepho, 1999).

Departure from the compound symmetry structure implies heterogeneity in genetic correlations, which will affect the weights $\boldsymbol{W}$ in the BLUP equation. Specifically, for estimation of $g_{ir}$ in the targeted subregion, the weight of neighboring subregions increases with the genetic correlation, while subregions showing low genetic correlation with the targeted subregion are down-weighted. The dependence of weights on genetic correlations with neighboring subregions, especially in models with heterogeneity in the correlations, is both intuitively appealing and statistically desirable. It is our experience that the proposed method demonstrates its full power when models departing from compound symmetry fit the data well. This is likely to occur when heterogeneity among subregions is pronounced. A detailed account will be published elsewhere.

While modeling of $\boldsymbol{\Sigma}_g$ is certainly the most crucial model selection step, other model components require attention as well. For example, heterogeneity of variance components pertaining to $e_{ijkr}$ in Eq. [1] may be expected between subregions. Also, one may model replicate data instead of trial means. This opens several options for more refined modeling at the trial level, depending on experimental design, including fitting of incomplete block effects and spatial modeling (Gilmour et al., 1997; Smith et al., 2001). The optimal properties of BLUP and the near-optimality of empirical BLUP are retained with these more general and more flexible models.

In conclusion, it should be emphasized that our approach can only improve on the prediction of genotype × region interaction, while interactions of genotype × year and genotype × year × location essentially remain unpredictable for any farm in the target region. Unfortunately, these latter effects often dominate total variation. Thus, while the method proposed in this paper is a small step forward, the problem of unpredictable interactions related to years largely continues unabated.

## APPENDIX A

We present sample code in SAS for single-step local BLUP. Factors are assumed to be coded as follows: G = genotype, S = subregion, L = location, Y = year. The MIXED code for the response variable YIELD and the compound symmetry structure for $\boldsymbol{\Sigma}_g$ is as follows:

```
proc mixed;

class G S L Y;

model YIELD = S;

random Y S*L S*Y S*L*Y G*Y G*S*L G*S*Y;

random S/sub = G type = CS solution;

run;
```

Instead of the compound symmetry structure, other models such as factor-analytic of heteroscedasticity can be used for $\boldsymbol{\Sigma}_g$ by appropriately modifying the type = option in the second "random" statement (Piepho, 1999).

To reduce computation time, one may take all effects not crossed with genotypes as fixed. This will yield identical results for balanced data. With unbalanced data, this analysis will not make use of intertrial information. Since the intertrial information is often small (Patterson, 1997), sacrifice in efficiency will usually be marginal. The reduction in computing time results from taking genotypes as the "subject" for all effects (Littell et al., 1996).

```
proc mixed;

class G S L Y;

model YIELD = S Y S*L S*Y S*L*Y;

random Y S*L S*Y/subject = G;

random S/sub = G type = CS solution;

run;
```

If a genotype is missing in a particular subregion, a BLUP can still be computed. To do so, the input dataset needs to have at least one record for the subregion × genotype combination in question, coding the missing observation as a dot.

## APPENDIX B

Let $\boldsymbol{y}' = (\boldsymbol{y}_1', \boldsymbol{y}_2', ..., \boldsymbol{y}_n')$ and $\boldsymbol{g}' = (\boldsymbol{g}_1', \boldsymbol{g}_2', ..., \boldsymbol{g}_n')$, where $n$ is the number of genotypes. For balanced data (see Eq. [6]) we find

$$\mathrm{var}(\boldsymbol{g}) = \boldsymbol{I}_n \otimes \boldsymbol{\Sigma}_g,$$

$$\mathrm{var}(\boldsymbol{y}) = \boldsymbol{I}_n \otimes (\boldsymbol{\Sigma}_g + \boldsymbol{\Sigma}_{e1}) + \boldsymbol{J}_n \otimes \boldsymbol{\Sigma}_{e2},$$

$$\mathrm{cov}(\boldsymbol{g}, \boldsymbol{y}) = \mathrm{var}(\boldsymbol{g}),$$

$$\mathrm{E}(\boldsymbol{g}) = \boldsymbol{0}, \text{ and}$$

$$\mathrm{E}(\boldsymbol{y}) = \boldsymbol{1} \otimes \boldsymbol{\mu},$$

where $\boldsymbol{I}_n$ is an $n$-dimensional identity matrix, $\boldsymbol{J}_n$ is an $n \times n$ matrix of ones everywhere, and $\otimes$ denotes the Kronecker or direct product operator (Searle et al., 1992). The best linear unbiased predictor of $\boldsymbol{g}$ is (Searle et al. (1992), p. 269)

$$\mathrm{BLUP}(\boldsymbol{g}) = \tilde{\boldsymbol{g}}^o = \boldsymbol{W}^o\,[\boldsymbol{y} - \mathrm{E}(\boldsymbol{y})],$$

where

$$\boldsymbol{W}^o = \mathrm{var}(\boldsymbol{g})[\mathrm{var}(\boldsymbol{y})]^{-1}.$$

The selection decision will depend on ranking of genotypes, which in turn is fully determined by all pairwise differences among genotypes. Any estimator, which always yields the same pairwise differences as BLUP, can be considered essentially equivalent to BLUP with regard to the resulting selection decision. An estimate of the pairwise difference for an estimable function of interest can be obtained by multiplication of BLUP($\boldsymbol{g}$) with a contrast vector $\boldsymbol{c} \otimes \boldsymbol{v}$, where $\boldsymbol{c}$ is an $n$-dimensional pairwise contrast vector with $c_i = 1$ and $c_{i'} = -1$ for the two genotypes to be compared and zeros elsewhere, while $\boldsymbol{v}$ is a coefficient vector for the estimable function of interest. We find

$$(\boldsymbol{c}' \otimes \boldsymbol{v}')\tilde{\boldsymbol{g}}^o = (\boldsymbol{c}' \otimes \boldsymbol{v}')\boldsymbol{W}^o[\boldsymbol{y} - \mathrm{E}(\boldsymbol{y})] \text{ with}$$

$$\boldsymbol{W}^o = (\boldsymbol{I}_n \otimes \boldsymbol{\Sigma}_g)[\boldsymbol{I}_n \otimes (\boldsymbol{\Sigma}_g + \boldsymbol{\Sigma}_{e1}) + \boldsymbol{J}_n \otimes \boldsymbol{\Sigma}_{e2}]^{-1}.$$

It can be shown that

$$[\boldsymbol{I}_n \otimes (\boldsymbol{\Sigma}_g + \boldsymbol{\Sigma}_{e1}) + \boldsymbol{J}_n \otimes \boldsymbol{\Sigma}_{e2}]^{-1} = [\boldsymbol{I}_n \otimes \boldsymbol{I}_m - $$
$$[\boldsymbol{J}_n \otimes \boldsymbol{\Sigma}_{e2}(\boldsymbol{\Sigma}_g + \boldsymbol{\Sigma}_{e1} + n\boldsymbol{\Sigma}_{e2})^{-1}][\boldsymbol{I}_n \otimes (\boldsymbol{\Sigma}_g + \boldsymbol{\Sigma}_{e1})^{-1}]$$

This follows from the easily established fact that

$$(\boldsymbol{I}_n \otimes \boldsymbol{A} + \boldsymbol{J}_n \otimes \boldsymbol{B})^{-1} =$$
$$[\boldsymbol{I}_n \otimes \boldsymbol{I}_m - \boldsymbol{J}_n \otimes \boldsymbol{B}(\boldsymbol{A} + n\boldsymbol{B})^{-1}](\boldsymbol{I}_n \otimes \boldsymbol{A}^{-1}),$$

where $\boldsymbol{A}$ and $\boldsymbol{B}$ are $m \times m$ matrices. The proof is as follows:

$$(\boldsymbol{I}_n \otimes \boldsymbol{A} + \boldsymbol{J}_n \otimes \boldsymbol{B})^{-1}(\boldsymbol{I}_n \otimes \boldsymbol{A} + \boldsymbol{J}_n \otimes \boldsymbol{B}) =$$
$$[\boldsymbol{I}_n \otimes \boldsymbol{I}_m - \boldsymbol{J}_n \otimes \boldsymbol{B}(\boldsymbol{A} + n\boldsymbol{B})^{-1}]$$
$$(\boldsymbol{I}_n \otimes \boldsymbol{A}^{-1})(\boldsymbol{I}_n \otimes \boldsymbol{A} + \boldsymbol{J}_n \otimes \boldsymbol{B}) =$$
$$[\boldsymbol{I}_n \otimes \boldsymbol{I}_m - \boldsymbol{J}_n \otimes \boldsymbol{B}(\boldsymbol{A} + n\boldsymbol{B})^{-1}]$$
$$(\boldsymbol{I}_n \otimes \boldsymbol{I}_m + \boldsymbol{J}_n \otimes \boldsymbol{A}^{-1}\boldsymbol{B}) =$$
$$\boldsymbol{I}_n \otimes \boldsymbol{I}_m + \boldsymbol{J}_n \otimes \boldsymbol{A}^{-1}\boldsymbol{B} - \boldsymbol{J}_n \otimes \boldsymbol{B}(\boldsymbol{A} + n\boldsymbol{B})^{-1} -$$
$$n\boldsymbol{J}_n \otimes \boldsymbol{A}^{-1}\boldsymbol{B}^2(\boldsymbol{A} + n\boldsymbol{B})^{-1} =$$
$$\boldsymbol{I}_n \otimes \boldsymbol{I}_m + \boldsymbol{J}_n \otimes (\boldsymbol{A} + n\boldsymbol{B})^{-1}[\boldsymbol{A}^{-1}\boldsymbol{B}(\boldsymbol{A} + n\boldsymbol{B}) -$$
$$\boldsymbol{B} - n\boldsymbol{A}^{-1}\boldsymbol{B}^2] =$$
$$\boldsymbol{I}_n \otimes \boldsymbol{I}_m$$

Thus, we find

$$(\boldsymbol{c}' \otimes \boldsymbol{v}')\tilde{\boldsymbol{g}}^o = (\boldsymbol{c}' \otimes \boldsymbol{v}')(\boldsymbol{I}_n \otimes \boldsymbol{\Sigma}_g)[\boldsymbol{I}_n \otimes \boldsymbol{I}_m -$$
$$\boldsymbol{J}_n \otimes \boldsymbol{\Sigma}_{e2}(\boldsymbol{\Sigma}_g + \boldsymbol{\Sigma}_{e1} + n\boldsymbol{\Sigma}_{e2})^{-1}$$

$$[\boldsymbol{I}_n \otimes (\boldsymbol{\Sigma}_g + \boldsymbol{\Sigma}_{e1})^{-1}][\boldsymbol{y} - \mathrm{E}(\boldsymbol{y})]$$
$$= (\boldsymbol{c}' \otimes \boldsymbol{v}')(\boldsymbol{I}_n \otimes \boldsymbol{\Sigma}_g)[\boldsymbol{I}_n \otimes (\boldsymbol{\Sigma}_g + \boldsymbol{\Sigma}_{e1})^{-1}]$$
$$[\boldsymbol{y} - \mathrm{E}(\boldsymbol{y})]$$
$$= (\boldsymbol{c}' \otimes \boldsymbol{v}')[\boldsymbol{I}_n \otimes \boldsymbol{\Sigma}_g(\boldsymbol{\Sigma}_g + \boldsymbol{\Sigma}_{e1})^{-1}]$$
$$[\boldsymbol{y} - \mathrm{E}(\boldsymbol{y})]$$
$$= (\boldsymbol{c}' \otimes \boldsymbol{v}')[\boldsymbol{I}_n \otimes \boldsymbol{W}][\boldsymbol{y} - \mathrm{E}(\boldsymbol{y})]$$

with $\boldsymbol{W} = \boldsymbol{\Sigma}_g(\boldsymbol{\Sigma}_g + \boldsymbol{\Sigma}_{e1})^{-1}$. The key step in this derivation is to note that $\boldsymbol{c}'\boldsymbol{J}_n = \boldsymbol{0}_n$, where $\boldsymbol{0}_n$ is a null vector. It follows that, for selection purposes, it is sufficient to compute

$$\tilde{\boldsymbol{g}} = [\boldsymbol{I}_n \otimes \boldsymbol{W}][\boldsymbol{y} - \mathrm{E}(\boldsymbol{y})].$$

This is seen to be the estimator in [11].

## ACKNOWLEDGMENTS

## REFERENCES

Atlin, G.N., R.J. Baker, K.B. McRae, and X. Lu. 2000. Selection response in subdivided target regions. Crop Sci. 40:7–13.

Comstock, R.E., and R.H. Moll. 1963. Genotype-environment interaction. p. 164–194. *In* W.D. Hanson and H. F. Robinson (ed.) Statistical genetics and plant breeding. Publication 982. National Academy of Sciences-National Research Council, Washington, DC.

Curnow, R.N. 1988. The use of correlated information on treatment effects when selecting the best treatment. Biometrika 75:287–293.

Falconer, D.S., and T.F.C. Mackay. 2001. Introduction to quantitative genetics. Pearson Education, Harlow.

Gilmour, A.R., B.R. Cullis, A.P. Verbyla, and A. C. Gleeson. 1997. Accounting for natural and extraneous variation in the analysis of field experiments. J.Agric. Biol. Environ. Statist. 2:269–293.

Kackar, R.N., and D.A. Harville. 1984. Approximations for standard errors of estimators of fixed and random effects in mixed linear models. J. Am. Statist. Assoc. 79:853–862.

Kish, L. 1965. Survey sampling. John Wiley & Sons, New York.

Littell, R.C., P.R. Henry, and C.B. Ammerman. 1996. Statistical analysis of repeated measures data using SAS procedures. J. Anim. Sci. 76:1216–1231.

Patterson, H.D. 1997. Analysis of series of variety trials. p. 139–161. *In* R.A. Kempton, and P.N. Fox (ed.) Statistical methods for plant variety evaluation. Chapman and Hall, London.

Piepho, H. p. 1998. Empirical best linear unbiased prediction in cultivar trials using factor analytic variance-covariance structures. Theor. Appl. Genet, 97:195–201.

Piepho, H. p. 1999. Stability analysis using the SAS system. Agron. J. 91:154–160.

Piepho, H.P., and K. Pillen. 2004. Mixed modelling for QTL × environment interaction analysis. Euphytica 137:147–153.

Portnoy, S. 1982. Maximizing the probability of correctly ordering random variables using linear predictors. J. Multivariate Anal. 12: 256–269.

Searle, S.R., G. Casella, and C.E. McCulloch. 1992. Variance components. John Wiley & Sons, New York.

Smith, A., B.R. Cullis, and R. Thompson. 2001. Analyzing variety by environment data using multiplicative mixed models and adjustments for spatial field trend. Biometrics 57:1138–1147.

Talbot, M. 1997. Resource allocation for selection systems. p. 162–174. *In* R. A. Kempton and P. N. Fox (ed.) Statistical methods for plant variety evaluation. Chapman and Hall, London.