# Regional importance of crop yield constraints: Linking simulation models and geostatistics to interpret spatial patterns

*David B. Lobell* [a,b,*], *J. Ivan Ortiz-Monasterio* [c]

[a] *Department of Global Ecology, Carnegie Institution of Washington, Stanford, CA 94305, USA*
[b] *Department of Geological and Environmental Science, Stanford University, Stanford, CA 94305, USA*
[c] *International Maize and Wheat Improvement Center (CIMMYT), Wheat Program, Apdo. Postal 6-641, 06600 Mexico D.F., Mexico*

## ARTICLE INFO

## ABSTRACT

Over the next few decades, a central goal of agricultural research and policy will be to increase average regional crop yields in the face of diminished gains in genetic yield potential, likely climatic changes, decreased resource availability, and stricter environmental standards. Fundamental to the pursuit of effective investment strategies is an ability to quantify tradeoffs associated with potential policy and management changes. However, the data needed to predict regional yield responses to change, namely observations of yields and climatic, soil, and management conditions in farmers' fields, are often difficult to obtain. In this paper, we investigate the value of data on the spatial distribution of yields for understanding causes of landscape yield variability. Stochastic simulation models, which employ the CERES model to simulate crop yields across a landscape, are used to translate assumed spatial patterns of soil and management conditions into spatial pattern of yields. Monte Carlo simulation is then used to repeat this process for many different realizations of conditions, resulting in a modeled relationship between yield patterns and the relative importance of soil and management yield constraints, both of which can be computed in the controlled simulation environment. The derived relationship then allows one to infer from observed yield patterns the true proportion of yield variability explained by soil and management.

This procedure was tested for wheat in the Yaqui Valley, an intensive agricultural region in Sonora, Mexico, where yield patterns have been previously estimated with remote sensing. Comparison of simulated and observed yield patterns indicated that roughly 80% of spatial yield variance in 2001–2002 was attributable to management variations. The ability of simulation models to aid interpretation of landscape patterns is potentially invaluable for understanding yield constraints in many agricultural regions, where direct observations of soil and management variables are infeasible.

* *Corresponding author.* Present address: Energy and Environment Directorate, Lawrence Livermore National Laboratory, P.O. Box 808, L-103 Livermore, CA 94550, USA. Tel.: +1 925 422 4148.
E-mail address: dlobell@llnl.gov (D.B. Lobell).

## 1.    Introduction

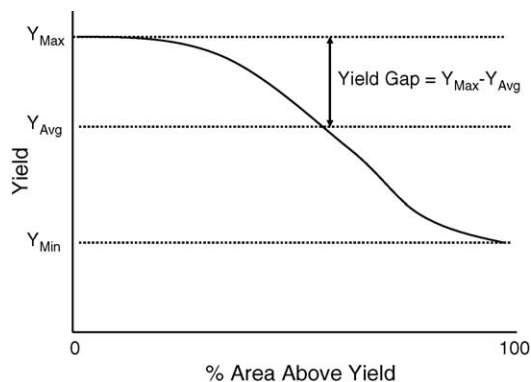Public investment in agricultural regions must consider trade-offs between multiple policy objectives, such as increased food production, improved environmental quality, decreased resource use, and higher farmer income. Fundamental to the pursuit of appropriate investment strategies is an ability to quantify the impact of policy decisions on each objective of interest. In this context, a key goal of agricultural research and policy is to close the gap between potential and average yields achieved in farmers' fields. Closing this yield gap will play an important role in increasing food production in the face of diminished gains in genetic yield potential, possible climatic changes, decreased resource availability, and stricter environmental standards (Cassman et al., 2003).

Identification of strategies to reduce the yield gap requires an understanding of its causes. This understanding may also be useful in other contexts, such as in assessing regional land quality and its change through time (Bindraban et al., 2000; Dumanski and Pieri, 2000). Here, we define the yield gap as the difference between attainable yield potential and average regional yields, where attainable yield is measured as yield on the highest yielding farmers' fields (Fig. 1). This definition contrasts with some authors who use more standard measures of yield potential, namely yields from experimental stations or outputs of crop simulation models (Evans, 1993; Penning de Vries et al., 1997; Cassman, 1999), which will often be larger than maximum economically viable yields. However, we utilize a definition based on farmers' yields for three main reasons. First, in some cases farmer yields may be more readily available than experimental or modeled data. Second, in intensively managed irrigated systems, the difference between true yield potential and maximum achieved yields is likely to be small. Third, in cases where the difference is substantial, the gap between average and maximum yields still provides an important measure of the potential to improve average yields under current technological and economic conditions (Mosher, 1978).



**Fig. 1 – Graphical representation of the spatial distribution of crop yields in a region. Spatial variability may arise from many factors, such as differences in land quality and various management practices. The yield gap is defined here as the difference between maximum and average yields.**

Yields achieved in a farmer's field can be considered a function of environmental conditions:

$$y(\mathbf{u}) = f[e(\mathbf{u})] \tag{1}$$

where $y(\mathbf{u})$ is the yield at location $\mathbf{u}$, $e(\mathbf{u})$ represents the conditions experienced by the crop at location $\mathbf{u}$, and $f$ describes the relationship between conditions and yield. Both $\mathbf{u}$ and $e(\mathbf{u})$ can be considered multivariate vectors, for instance, $\mathbf{u}$ may consist of latitude and longtitude, and $e(\mathbf{u})$ may describe the many soil, climate, and management variables affecting crop growth.

To evaluate the impact on the yield gap of changes in environmental conditions, such as those that might be achieved through policy instruments, an understanding of what drives variability in $y(\mathbf{u})$ is needed. Typically, this means obtaining observations of $e(\mathbf{u})$ in farmers fields, which can then be combined with pre-defined models of $f$ (e.g., Aggarwal and Kalra, 1994; Singh et al., 1994; Matthews et al., 2002) or with statistical models relying on joint observations of $y(\mathbf{u})$ to determine the contribution of different factors to yield variability (Calvino and Sadras, 2002; Lobell et al., 2004).

While these studies have provided new insights into the yield gap in specific regions, detailed understanding of the causes and often even the magnitudes of yield gaps are poorly known in many regions. This paucity of information reflects the difficulty of acquiring extensive information on both yields and environmental conditions in fields. For instance, information on the spatial heterogeneity of soil properties and management practices between fields are poorly known in most, if not all, regions (Hansen and Jones, 2000).

Given the difficulty of measuring conditions on farmers' fields, an approach to understanding the yield gap that utilizes only readily available information would be of considerable value. For instance, data on crop yields and their spatial locations are increasingly available through GIS technologies, such as combine-mounted yield monitors and remote sensing (e.g., Dobermann et al., 2003; Lobell et al., 2005). What can be learned from this information alone? In this paper, we investigate the use of spatially referenced yield data to infer causes of the yield gap, based on the concept that different potential limiting factors can be associated with different spatial patterns of variability.

The notion that yield patterns reflect causes of variability is common in precision agriculture studies aimed at understanding within-field yield gradients (Lotz, 1997; Plant, 2001). Underlying these studies are assumptions about the spatial patterns of yield controls within fields (e.g., management practices are uniform along straight lines; Lotz, 1997), which are often based on farmer experience and can be readily tested with field measurements. At the landscape scale, knowledge of the spatial structure of soil, climate, and management conditions is currently more limited. However, many datasets that exist for other purposes can potentially be used to analyze spatial patterns (see example below), and the collection and storage of georeferenced data is likely to become more affordable and common in the future.

There are many ways to quantitatively describe spatial patterns (see, e.g., Cressie, 1991; Getis and Ord, 1992; Gustafson, 1998; Haining, 2003). Here, we adopt a geostatistical approach, where the spatial distribution of a variable $Z(\mathbf{u})$ is character-

ized by its auto-covariance function $C(\mathbf{h})$, which defines the degree of association between two values separated by a vector $\mathbf{h}$.

$$C(\mathbf{h}) = \text{Cov}(\mathbf{u}, \mathbf{u} + \mathbf{h}) = E[Z(\mathbf{u})Z(\mathbf{u} + \mathbf{h})] - E[Z(\mathbf{u})]^2 \qquad (2)$$

or, alternatively, by its semi-variance function $\gamma(\mathbf{h})$:

$$2\gamma(\mathbf{h}) = \text{Var}[Z(\mathbf{u} + \mathbf{h}) - Z(\mathbf{u})] = E[(Z(\mathbf{u} + \mathbf{h}) - Z(\mathbf{u}))^2] \qquad (3)$$

The spatial auto-covariance or semi-variance is a statistical summary of the relationship between values of $Z(\mathbf{u})$ separated by $\mathbf{h}$ (for more discussion see, e.g., Deutsch and Journel, 1992). Empirical determination of $C(\mathbf{h})$ or $\gamma(\mathbf{h})$ requires systematic sampling of the variable(s) of interest (i.e., $e(\mathbf{u})$) across the landscape. Unlike obtaining joint observations of $y(\mathbf{u})$ and $e(\mathbf{u})$ within farmers fields, however, $C(\mathbf{h})$ or $\gamma(\mathbf{h})$ can be defined independently for each variable (provided that covariance between the variables can be neglected). For example, it is possible to use datasets collected using different sample locations and at different times.

To summarize, information on the spatial pattern (e.g., $\gamma(\mathbf{h})$) of both $y(\mathbf{u})$ and $e(\mathbf{u})$) may be easy to obtain relative to simultaneous measurement of yield and conditions at locations throughout a landscape. The challenge then becomes using these spatial patterns alone to infer causes of the yield gap. Similar attempts in ecological research to relate measurable landscape patterns to underlying processes have benefited from the use of stochastic simulation models, which can quantitatively explore the sensitivity of landscape pattern to various factors (e.g., Gardner et al., 1987; Turner et al., 2001). Effectively, such models relate underlying processes to resulting patterns, and thereby provide a way to test hypotheses about the causes of observed landscape patterns.

The goal of this study was to develop and test a modelling framework for combining information on patterns of $y(\mathbf{u})$ and $e(\mathbf{u})$ to identify the relative importance of soil and management constraints at the regional scale. This approach was tested in an agricultural region in Northwest Mexico, where data on the spatial patterns of conditions and yields were available.

## 2. Methods

The proposed methodology consists of using a landscape simulation model to train a simple regression (i.e., meta-model) between yield patterns and the importance of different yield-controlling factors. This regression can then used to quantify factor importance based on observed yield patterns. The procedure entails five main steps, as illustrated in Fig. 2 and outlined below.

(1) The first stage is to select a simulation model capable of predicting yield under a specified set of conditions. This model should include, as much as possible, all processes operating in the landscape. In this case, we have chosen the CERES-wheat crop model (Tsuji et al., 1994), which has proven capable of simulating responses to climatic variations, and to a lesser degree water and nutrient stress (Jones et al., 2003).

(2) The next step is to define the spatial structure of variability (e.g., $\gamma(\mathbf{h})$) for all conditions under consideration. The choice of which conditions to consider will be dictated by the required inputs for the crop yield model (defined in step 1). In this process, one should also account for any spatial covariance between input variables, as discussed in the example below. The definition of variability can also be extended in a straightforward way to include temporal variation. While not done in the example below, this would be useful for analyzing the complete spatio-temporal pattern of yield variability.

(3) The third stage is a Monte Carlo simulation of landscape yield patterns. The structures defined in step 2 are used to randomly simulate spatial distributions of conditions, for example through sequential Gaussian simulation with
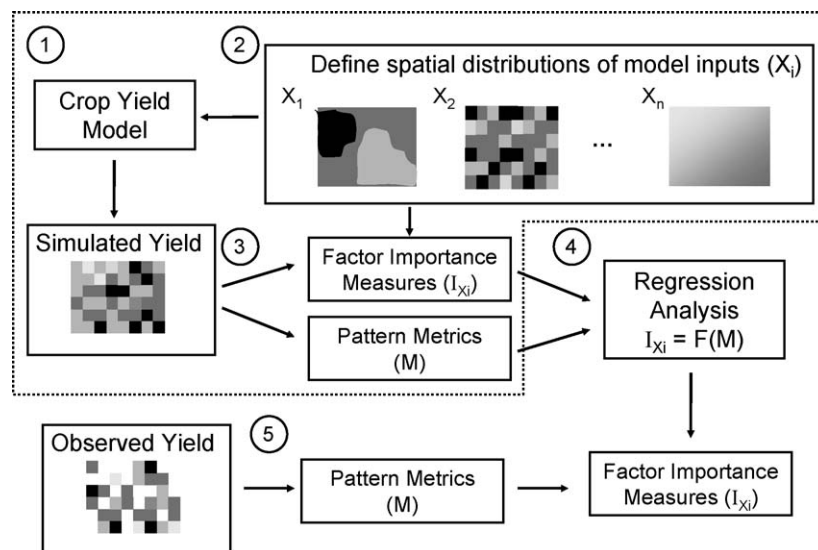


**Fig. 2 – Schematic representation of simulation approach to interpreting yield patterns. Numbers indicate steps described in text.**

simple kriging (Deutsch and Journel, 1992). The yield model is then applied at each point on the landscape, resulting in an array of simulated yields (see Fig. 2). Two independent analyses are then performed on these arrays. First, the outputs are regressed against the known inputs to determine the proportion of yield variability explained by each input. In this case, we use simple linear regression although non-linear techniques may also be used. Second, the spatial pattern of yield is measured, using one or many quantitative metrics of landscape variability. The choice of a metrics is arbitrary, and could incorporate several individual metrics. In this study, we employ as a metric the ratio of $\gamma(\mathbf{h})$ for yield at 350 m and 2 km lag distance, which represents a measure of short versus long-range variability (see below).

This sequential process of simulating inputs and then outputs, and finally analyzing sources of variability and yield patterns, is repeated a large number of times, in this case 100, in order to quantify the sensitivity of the simulated pattern to the underlying processes. The result of this Monte Carlo simulation is a value of percent variance explained by (or relative importance of) each input for each run, and an array of pattern metrics for each run.
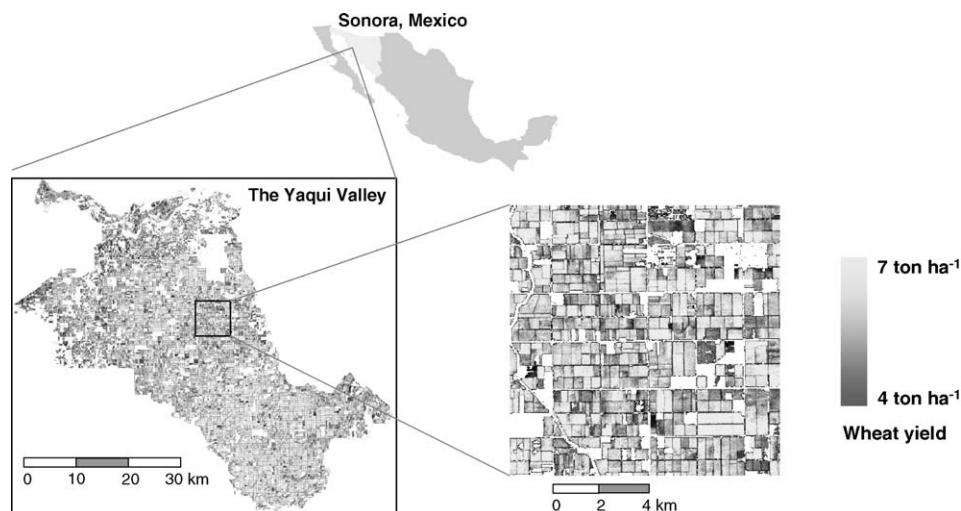
This step involves explicit definition of the model spatial scale, which includes both the grain (also commonly called support or resolution) and extent. The grain should be small enough that model inputs are relatively homogeneous (Hansen and Jones, 2000), but cannot be so small that simulating a proper extent is computationally infeasible. The specification of model extent also involves a tradeoff, in that it should be large enough to allow significant variation of inputs within the landscape, but small enough to ensure different degrees of variability in different simulation runs. Essentially, we seek a simulation extent such that some simulations have low spatial variance in a given input and some simulations have high spatial variance. In the example below, we define the model grain as 1 ha (100 m × 100 m) and model extent as 9 km² (3 km × 3 km). The former was chosen to

be smaller than individual field sizes, which are roughly 10 ha in the example. The extent of 9 km² was selected because it exceeds the maximum correlation length of soil properties (see below), but was small enough to produce substantial variations between simulations.

(4) The fourth step is to statistically compare the outputs of the simulation experiment, specifically the importance of each input and the yield pattern metrics. The ability of a statistical model to relate the two quantifies the degree to which the observable pattern metric(s) is sensitive to (and therefore, can be used to predict) the desired quantity. As mentioned in step 3, ensuring sufficient contrast between simulation runs is critical for training the statistical model. This stage may also involve comparing the predictive power of several candidate pattern metrics.

(5) Finally, the fifth step is to apply the selected metric(s) to an actual dataset of spatial yield distributions. Such a dataset may be available, for instance, through field surveys or remote sensing analyses. The value of the derived metric(s) can then be related to the relative importance of underlying causes based on the simulation analysis in the previous steps.

### 2.1. Data used in this study

The procedure outlined above was tested in the Yaqui Valley (YV), an irrigated region comprising 225,000 ha along the western coast of mainland Sonora, Mexico (27°N 110°W). In YV, spring wheat is sown on average to ~60% of the irrigation district in November–December and harvested in April–May. Previous work with Landsat satellite imagery has furnished well-validated datasets on the spatial distributions of wheat yields in this region (Lobell et al., 2003, 2005; see Fig. 3). Specifically, a method that combines multiple Landsat images with a temperature-based model of crop growth provides yield estimates at the 30 m × 30 m resolution of Landsat with root mean square errors of 0.5 t ha⁻¹ (~10%). For this study, we used yield data from the 2001 to 2002 growing season.



**Fig. 3 – Wheat yields in the Yaqui Valley study region of Northwest Mexico for the 2001–2002 growing season, estimated from Landsat data.**

In addition to the availability of yield estimates, the YV was selected because previous and ongoing studies of the yield gap provide an independent means of evaluating the performance of the simulation approach. Specifically, water and nutrient management practices in YV have been found to explain at least 50% of the spatial yield variability (Lobell et al., 2005). Remaining variability is likely due to additional management practices, with at least ∼10% due to soil properties (Lobell et al., 2002).

To characterize the spatial distribution of soil conditions, measurements of soil texture at three depth intervals (0–30, 30–60, and 60–100 cm) were obtained from the National Water Commision (CNA) for 2594 georeferenced points collected throughout the valley over the past decade. Unfortunately, the accuracy of this dataset was not well-documented and significant errors may be present. For instance, clay values appeared consistently lower than values measured in our previous work. However, since we use this dataset mainly to specify the spatial structure of soil variability, any systematic errors should not greatly affect the analysis.

Since CERES requires soil inputs in terms of soil lower limit (SLL), drained upper limit (DUL), and saturated upper limit (SAT), the texture values were converted based on the following pedo-transfer functions (Rawls et al., 1982):

$$SLL\,(\%) = 2.6 + 0.5 \times \%\,clay + 1.58 \times OM \qquad (4)$$

$$DUL\,(\%) = 25.76 + 0.36 \times \%\,clay - 0.20 \times \%\,sand + 2.99 \times OM \qquad (5)$$

$$SAT\,(\%) = 78.99 - 0.37 \times \%\,sand + 1.0 \times OM - 13.15 \times BD \qquad (6)$$

The equations for SLL, DUL, and SAT were obtained from the equations in Rawls et al. (1982) for −1500, −33, and −4 kPa, respectively, which were derived from multiple linear regression analysis of 1323 soils. Soil organic matter (OM) and bulk density (BD) were set to 0.5% and 1.5 g cm$^{-3}$, respectively, which are representative value for soils in this region. Because pedo-transfer functions can exhibit large errors when applied to new datasets (Gijsman et al., 2002), we first tested these equations using a database of 60 local measurements of soil

texture, SLL, and DUL, which ranged from 24 to 49% clay, 10–30% silt, and 24–66% sand (Ortiz-Monasterio, unpublished data). The root mean square error (rmse) and bias of predictions were 4.06 and 3.36%, respectively, for SLL and 4.89 and −0.61% for DUL. Thus, the equations were deemed appropriate for the purposes of the simulation, namely to characterize the spatial structure of soil water limits. While SAT was not directly validated, repeated runs of CERES with different values of SAT revealed a negligible sensitivity of modeled yields to percent saturation.
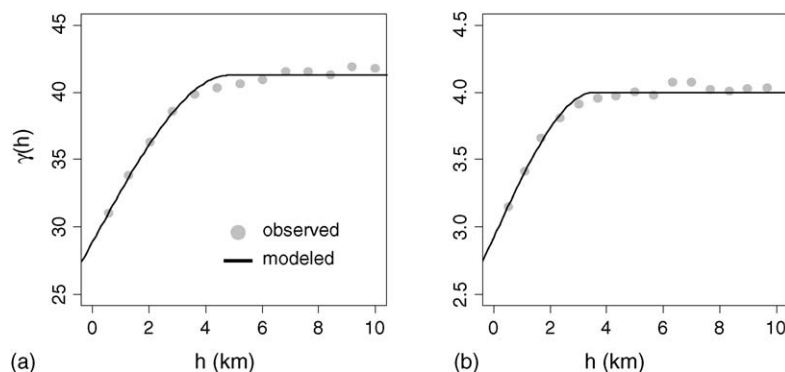
Data on management practices were taken from a survey of 74 wheat farmers conducted in 2002–2003 (Lobell et al., 2005). While this survey measured numerous aspects of management, for the purpose of the simulation we considered only the following main variables: planting date, fertilizer rate, seeding densities, number of irrigations, and timing of irrigations (Table 1). All soil and management inputs required by CERES other than those mentioned above were set to representative values and held constant for all simulations. Variability of model outputs were therefore due entirely to variability in the soil water limits and management practices discussed above. Data on temperature, rainfall, and solar radiation were obtained from a local meteorological station for the 2001–2002 growing season (http://www.pieaes.org.mx/datos.htm).

### 2.2. Spatial simulation of conditions

#### 2.2.1. Soil
The soil conditions measured in YV exhibited spatial autocorrelation up to approximately 8 km. Fig. 4a illustrates the sample variogram for field capacity in the upper 30 cm in YV, where here and throughout this paper spatial variability is assumed to be isotropic (only the magnitude, and not direction, of **h** is considered). Importantly, as the separation distance approaches zero, $\gamma(\mathbf{h})$ tends to a non-zero value. Commonly called a nugget effect, this value reflects the variability between two samples located infinitesimally close to each other, which can result from fine scale variability and/or measurement error.

Because soil properties are highly correlated with each other, it is important to account for this correlation when simulating spatial distributions (Wackernagel, 2003). We, therefore, performed a principal component analysis (PCA) on the nine soil properties (permanent wilting point, water holding



Fig. 4 – Empirical semi-variogram for (a) soil 0–30 cm field capacity and (b) first principal component of soil properties (see text for details). Solid lines shows best-fit model semi-variogram.

Table 1 – Attributes of management variables used in this study (from Lobell et al., 2005)

| Name | Units | Mean | Standard deviation | Minimum | Maximum | Spatial auto-correlation (Moran's I) | p-Value of Moran's I |
|---|---|---|---|---|---|---|---|
| Planting day | Days after November 1 | 34.1 | 13.0 | 8 | 62 | 0.08 | 0.13 |
| N fertilizer rate | kg ha$^{-1}$ | 250.8 | 50.3 | 66 | 368 | 0.01 | 0.32 |
| Seeding density | kg ha$^{-1}$ | 146.6 | 23.5 | 90 | 200 | −0.09 | 0.93 |
| Number of irrigations | # | 3.73 | 0.5 | 3 | 5 | 0.11 | 0.05 |
| Day of first irrigation | Days before planting | 20.8 | 9.1 | 0 | 48 | −0.04 | 0.61 |
| Day of second irrigation | Days after planting | 54.7 | 8.8 | 26 | 80 | 0.06 | 0.21 |
| Day of third irrigation | Days after planting | 85.4 | 9.5 | 55 | 113 | 0.05 | 0.15 |
| Day of fourth irrigation | Days after planting | 102.4 | 7.0 | 77 | 115 | 0.24 | 0.01 |

capacity, and saturation point at 0–30, 30–60, and 60–100 cm). A spherical variogram model was then fit to the sample variogram of each of the derived principal components, which by definition were orthogonal (a spherical model provided the best-fit to the sample variogram; see Fig. 4b). For each simulation run, spatial distributions of the principal components were generated using unconditional sequential Gaussian simulation, implemented with the gstat library in the R software package (Pebesma, 2004). The sequential simulation used simple block kriging, with the block size defined as 100 m × 100 m (1 ha) and the neighborhood size limited to 15 points. In this procedure, grid points are randomly visited and values are generated from a Gaussian distribution defined by the existing values in neighboring cells (for details see, e.g., Deutsch and Journel, 1992). The simulated values of the principal components were finally back-transformed to the original units using the covariance matrix derived in the PCA. Fig. 5 displays the simulated distribution of soil properties for five simulation runs.

### 2.2.2. Management

Management practices in YV did not generally exhibit significant spatial correlation, with Moran's I statistic exceeding 5% significance only for the day of fourth irrigation (Table 1). We, therefore, considered management to be randomly distributed between fields. For the simulation, fields were defined as an area comprising 3 × 3 cells (9 ha), which is close to the size of field divisions in YV (each lot is 10 ha). For each field, one of the 74 observed management regimes in the survey was randomly selected as input to the CERES model. This method of using the empirical distribution was chosen because the fairly small sample size of the survey precluded an accurate model of the management distribution (e.g., multi-variate Gaussian).
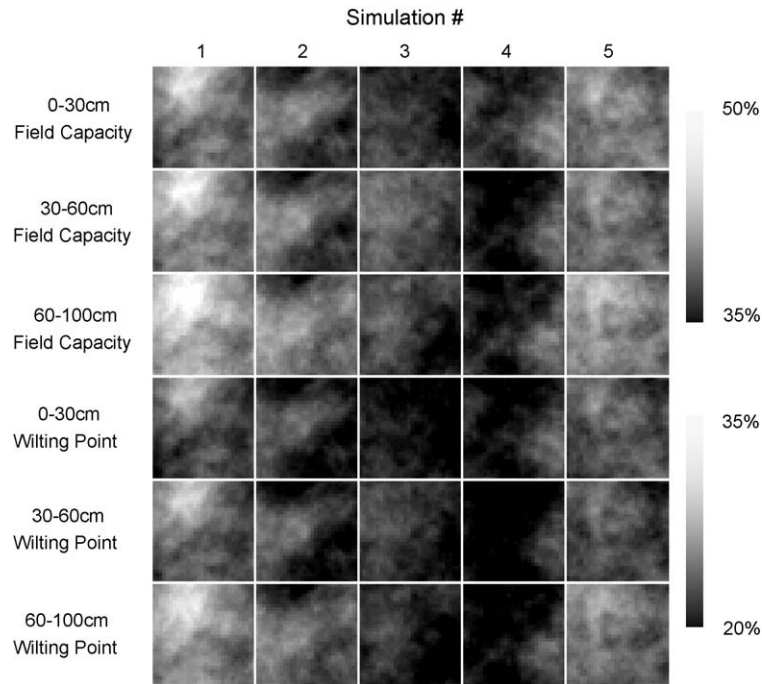
### 2.2.3. Climate

Climatic conditions were assumed to be spatially uniform in YV, which exhibits very little variation in elevation. While this simplification may introduce some error into the analysis, we currently do not understand spatial climatic variability enough to justify a non-constant model. We expect any errors due to climatic variability to be small, however, given the minimal importance of rainfall for irrigated wheat and the likely high spatial auto-correlation of temperature.

## 3. Results

### 3.1. Simulated yields, factor importance, and spatial patterns

Simulated yield distributions for four model runs are illustrated in Fig. 6. The 90,000 simulated values (30 ha × 30 ha × 100 runs) exhibited a reasonable distribution for YV, with a mean of 5.15 t ha$^{-1}$ and a standard deviation of 1.24 t ha$^{-1}$.
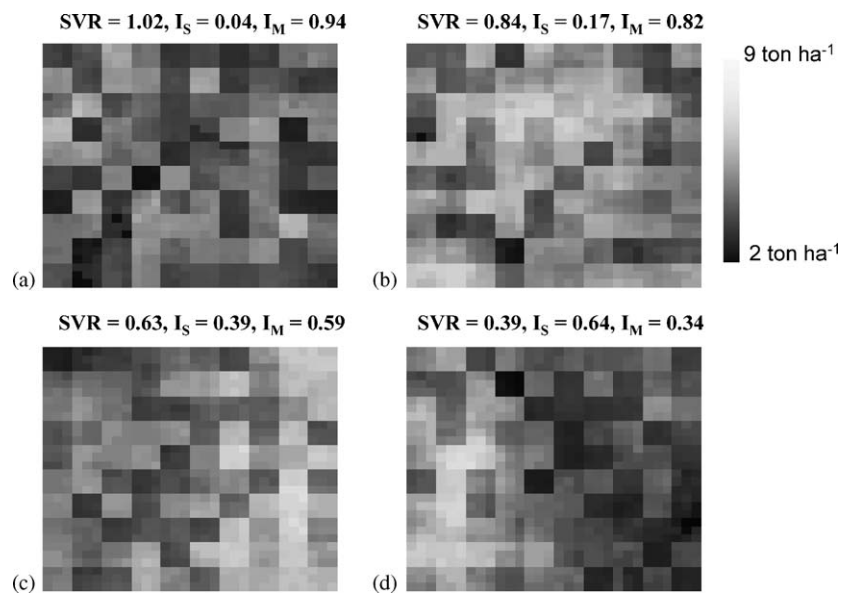
As discussed in step 3, measures of factor importance were derived from a linear regression relating simulated yields to soil and management inputs. Specifically, yields were first regressed on the six soil properties, with the importance of soil ($I_S$) quantified as the coefficient of determination ($R^2$) for
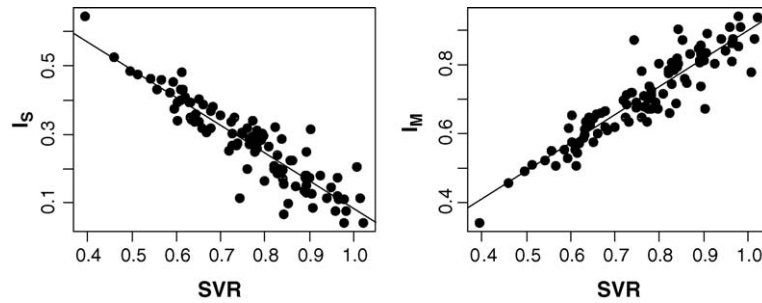
Fig. 5 – Simulated spatial distribution of soil field capacity and wilting point for five sample simulation runs, generated with unconditional sequential Gaussian simulation. The grain and extent of each simulation is 100 m × 100 m and 3 km × 3 km, respectively. These soil properties, along with management and climatic variables, were input into CERES-wheat for the spatial simulation of crop yields.

the linear regression. Management, which was represented as a categorical variable with 76 classes corresponding to the 76 individual management regimes, was then added as a predictor variable in the regression. The importance of management ($I_M$) was quantified as the difference in $R^2$ between this full model and the soil-only model.

Fig. 6 displays the values of $I_S$, $I_M$, and the ratio of semi-variance at 350 m and 2 km (hereafter called SVR for semi-variance ratio) for the corresponding simulation runs. As the field sizes in the simulation were 300 m × 300 m, 350 m represents the distance between two adjacent fields, while 2 km represents a distance close to the extent of the model



Fig. 6 – (a–d) Simulated yield patterns for four sample simulation runs, along with computed values of the ratio of semi-variance at 350 m to semi-variance at 2 km (SVR), importance of soil ($I_S$), and management ($I_M$). SVR is higher for those simulations with higher values of $I_M$, the proportion of overall yield variance explained by management.

**Fig. 7 – Plots of SVR vs. $I_S$ (left) and $I_M$ (right) for 100 simulation runs. Solid line indicates best-fit linear regression. $R^2$ of regression was 0.84 for both $I_S$ and $I_M$.**

(3 km × 3 km). The SVR was therefore expected to provide an indicator of how much more similar adjacent fields were than fields separated by relatively large distances. In agreement with expectation, simulations in which yield variations were predominantly driven by management (e.g., Fig. 6a) were characterized by a value of SVR near 1, indicating that points in adjacent fields (separated by 350 m) were no more similar to each other than points separated by 2 km. In contrast, simulations in which soil properties explained a substantial part of yield variation (e.g., Fig. 6d) exhibited significant correlation between yields on adjacent fields. As a result, the values of SVR for these simulations tended to be lower.

Using all 100 simulation runs, the statistical relationship between SVR and $I_S$ or $I_M$ was established (Fig. 7). This relationship was not heavily influenced by any single simulation, implying that 100 Monte Carlo simulations were sufficient in this particular example. As expected, higher values of SVR tended to reflect landscapes with a greater contribution of management to crop yield variance. While the association between yield patterns and factor importance was strong ($R^2 = 0.84$), a significant amount of scatter was also evident. This reflects the fact that soil exhibits a substantial amount of short-range variability because of the nugget effect, as well as the fact that even totally random variations in management can result in non-negligible long-range variability over the spatial extent of the simulation. The scatter would likely be reduced considerably with more accurate soil measurements, as a significant fraction of the nugget effect may be due to measurement error.

The relationship between factor importance and yield patterns may also be improved through the use of alternative pattern metrics (Gustafson, 1998; Haining, 2003). Despite the limitations of the current approach, however, it is clear that landscape yield patterns can considerably constrain the perceived importance of soil and management as underlying causes of yield variability. Such information can be used to rapidly assess the relevant constraints to regional production and guide more detailed studies of specific factors.

### 3.2. Application to remote sensing data

The metric SVR was computed for the yield dataset illustrated in Fig. 3. Since the spatial support of the remotely sensed yields (30 m × 30 m) differed from that of the simulation (100 m × 100 m), the metric derived from remote sensing

data was adjusted for the change in support. This adjustment, known as regularization, requires a new variogram $\gamma_V(\mathbf{h})$ to be computed from the sample variogram $\gamma(\mathbf{h})$, where $\gamma_V(\mathbf{h})$ represents the variogram for support V. $\gamma_V(\mathbf{h})$ can be expressed as (Armstrong, 1998):

$$\gamma_V(\mathbf{h}) = \overline{\gamma}(V, V_h) - \overline{\gamma}(V, V) \tag{7}$$

where V and $V_h$ are two regions of area (support) V whose midpoints are separated by $\mathbf{h}$, $\overline{\gamma}(V, V_h)$ the average value of the variogram between an arbitrary point in V and $V_h$, and $\overline{\gamma}(V, V)$ is the average value of the variogram between two arbitrary points in V. The two terms on the right hand side of Eq. (7) can be readily computed numerically if the support of $\gamma_V(\mathbf{h})$ is a multiple of the original support of $\gamma(\mathbf{h})$, by computing the variogram value for every possible combination of points in V (Armstrong, 1998). Indeed, the relative ease with which variograms can be transformed to different scales is one attractive feature of geostatistical pattern metrics.

Since 100 is not a multiple of 30, $\gamma_V(\mathbf{h})$ was computed for both 90 m × 90 m and 120 m × 120 m support. The corresponding values of SVR were computed as 0.90 and 0.95. Thus, the observed value of SVR for the remotely sensed yield data at the support of the simulation was roughly 0.9. As seen in Fig. 7, this indicates that management variations explain roughly 70–90% of the yield variability for this year in the Yaqui Valley. This finding agrees well with previous studies in YV that employed joint observations of soil type, management practices, and yields (Lobell et al., 2002, 2005), which concluded that the majority of yield variance is explained by management. The causes of the yield gap in YV obviously may differ from those in other regions, in particular rainfed systems. Further study would be needed to evaluate how the sources of the yield gap differ in other settings.

## 4. Discussion and conclusions

The results demonstrate that spatial patterns of yields possess substantial information on the relative importance of soil and management factors for yield variability. While previous studies have applied crop models to spatially explicit input datasets (Beaujouan et al., 2001; Priya and Shibasaki, 2001; Mo et al., 2005), to our knowledge this is the first study to explicitly evaluate the pattern of derived yields and its sensitivity

to underlying factors. As discussed in Section 1, the ability to directly interpret aspects of observed yield patterns could potentially reduce the amount of data and effort needed to assess yield constraints. Specifically, data on the actual spatial distribution of soil, climate, and management factors can be replaced by knowledge of only their geostatistical properties.

The simulation framework presented here was simple in many respects. For instance, only selected soil and management properties were considered in the model, and unmodeled factors, such as planting depth, water quality, pest dynamics, and myriad others may contribute to yield variability in reality. In addition, soil properties were simulated using simple geostatistical techniques whereas features, such as anisotropy or sharp transitions from one soil type to another were not modeled (Heuvelink and Webster, 2001). Spatial variability in climatic factors was also ignored though factors, such as frost or rainfall variations would likely be critical in many regions.

Thus, we emphasize that each step of the framework may be improved through more accurate measurements and more sophisticated modelling and pattern analysis techniques. Nonetheless, the modelling framework itself appears a promising avenue for interpreting readily available information (spatial distribution of crop yields) in terms of desired knowledge (sources of the yield gap). We also note that the presence of measurement error in the soil samples likely led to an overestimation of short-range variability in soil conditions, and therefore an underestimation of the ability of spatial patterns to discriminate soil from management controls.

Successful application of this approach relies on proper specification of the spatial distribution of environmental conditions. Thus, more data on the semi-variance structure of relevant soil, climate, and management variables is needed. It is likely that to some extent such information can be mined from existing databases collected by government agencies and extension services. Alternatively, it is possible to specify models for spatial dependence of conditions based not on hard data, but on assumptions about landscape patterns. For instance, one might assume that management is randomly distributed between fields, and that soil properties covary up to 8 km. The simulation framework can then be used to evaluate the implications of these assumptions for the interpretation of yield patterns. More generally, repeated analyses with different scenarios of landscape soil and management patterns can be used to test the sensitivity of interpretations to underlying assumptions.

The spatial range of correlation for soil properties in YV was similar to the values of 10–19 km for permanent wilting point and field capacity in the Netherlands presented in (Hoosbeek and Bouma, 1998). While a full synthesis of landscape soil variability is beyond the scope of this paper, this raises the question of whether general constraints can be placed on the range of soil semi-variograms when entering new regions. Similarly, it is not yet clear whether management is randomly distributed in most regions or whether this is unique to YV. A more complete understanding of the diversity in soil and management patterns among different regions, combined with a sensitivity study of various pattern metrics, may facilitate the development of robust metrics that can be used to rapidly assess limits to production even in regions with very scarce ground data.

Finally, we note that the approach presented here can be readily extended to look at additional factors, such as climatic variability and interactions between soil, climate, and management. Individual soil and management variables may also be discernible to the degree that their spatial pattern of variability is unique from other variables. In general, there are many dimensions to space-time patterns of yield that may contain useful information for scientists and policy makers, provided that we possess the imagination and quantitative tools to interpret them. While inverse models based on landscape pattern will likely never provide the detail of knowledge possible with ground-based studies, their minimal cost and time requirements and ability to use a growing wealth of spatial information may be of great value in regions with limited resources and time available for making important investment decisions.

## Acknowledgements

### REFERENCES

Aggarwal, P.K., Kalra, N., 1994. Analyzing the limitations set by climatic factors, genotype, and water and nitrogen availability on productivity of wheat. 2. Climatically potential yields and management strategies. Field Crops Res. 38, 93–103.

Armstrong, M., 1998. Basic Linear Geostatistics. Springer–Verlag, Germany, 153 pp.

Beaujouan, V., Durand, P., Ruiz, L., 2001. Modelling the effect of the spatial distribution of agricultural practices on nitrogen fluxes in rural catchments. Ecol. Model. 137, 93–105.

Bindraban, P.S., Stoorvogel, J.J., Jansen, D.M., Vlaming, J., Groot, J.J.R., 2000. Land quality indicators for sustainable land management: proposed method for yield gap and soil nutrient balance. Agric. Ecosyst. Environ. 81, 103–112.

Calvino, P., Sadras, V., 2002. On-farm assessment of constraints to wheat yield in the south-eastern Pampas. Field Crops Res. 74, 1–11.

Cassman, K.G., 1999. Ecological intensification of cereal production systems: yield potential, soil quality, and precision agriculture. Proc. Natl. Acad. Sci. 96, 5952–5959.

Cassman, K.G., Dobermann, A., Walters, D.T., Yang, H., 2003. Meeting cereal demand while protecting natural resources and improving environmental quality. Annu. Rev. Environ. Resour. 28, 315–358.

Cressie, N., 1991. Statistics for Spatial Data. Wiley, New York.

Deutsch, C.V., Journel, A.G., 1992. GSLIB: Geostatistical Software Library and User's Guide. Oxford University Press, New York.

Dobermann, A., Ping, J.L., Adamchuk, V.I., Simbahan, G.C., Ferguson, R.B., 2003. Classification of crop yield variability in irrigated production fields. Agron. J. 95, 1105–1120.

Dumanski, J., Pieri, C., 2000. Land quality indicators: research plan. Agric. Ecosyst. Environ. 81, 93–102.

Evans, L.T., 1993. Crop Evolution, Adaptation, and Yield. Cambridge University Press, New York, 500 pp.

Gardner, R.H., Milne, B.T., Turner, M.G., O'Neill, R.V., 1987. Neutral models for the analysis of broad-scale landscape pattern. Landscape Ecol. 1, 19–28.

Getis, A., Ord, J.K., 1992. The analysis of spatial association by use of distance statistics. Geogr. Anal. 24, 189–206.

Gijsman, A.J., Jagtap, S.S., Jones, J.W., 2002. Wading through a swamp of complete confusion: how to choose a method for estimating soil water retention parameters for crop models. Eur. J. Agron. 18, 75–105.

Gustafson, E.J., 1998. Quantifying landscape spatial pattern: what is the state of the art? Ecosystems 1, 143–156.

Haining, R., 2003. Spatial Data Analysis: Theory and Practice. Cambridge University Press, Cambridge, 432 pp.

Hansen, J.W., Jones, J.W., 2000. Scaling-up crop models for climate variability applications. Agric. Syst. 65, 43–72.

Heuvelink, G.B.M., Webster, R., 2001. Modelling soil variation: past, present, and future. Geoderma 100, 269–301.

Hoosbeek, M.R., Bouma, J., 1998. Obtaining soil and land quality indicators using research chains and geostatistical methods. Nutr. Cycl. Agroecosyst. 50, 35–50.

Jones, J.W., Hoogenboom, G., Porter, C.H., Boote, K.J., Batchelor, W.D., Hunt, L.A., Wilkens, P.W., Singh, U., Gijsman, A.J., Ritchie, J.T., 2003. The DSSAT cropping system model. Eur. J. Agron. 18, 235–265.

Lobell, D., Ortiz-Monasterio, J., Addams, C., Asner, G., 2002. Soil, climate, and management impacts on regional wheat productivity in Mexico from remote sensing. Agric. Forest Meteorol. 114, 31–43.

Lobell, D.B., Asner, G.P., Ortiz-Monasterio, J.I., Benning, T.L., 2003. Remote sensing of regional crop production in the Yaqui Valley, Mexico: estimates and uncertainties. Agric. Ecosyst. Environ. 94, 205–220.

Lobell, D.B., Ortiz-Monasterio, J.I., Asner, G.P., 2004. Relative importance of soil and climate variability for nitrogen management in irrigated wheat. Field Crops Res. 87, 155–165.

Lobell, D.B., Ortiz-Monasterio, J.I., Asner, G.P., Naylor, R.L., Falcon, W.P., 2005. Combining field surveys, remote sensing, and regression trees to understand yield variations in an irrigated wheat landscape. Agron. J. 97, 241–249.

Lotz, L., 1997. Yield Monitors and Maps: Making Decisions. Ohio State University, Columbus, OH.

Matthews, R., Stephens, W., Hess, T., Middleton, T., Graves, A., 2002. Applications of crop/soil simulation models in tropical agricultural systems. Adv. Agron. 76, 31–124.

Mo, X., Liu, S., Lin, Z., Xu, Y., Xiang, Y., McVicar, T.R., 2005. Prediction of crop yield, water consumption and water use efficiency with a SVAT-crop growth model using remotely sensed data on the North China Plain. Ecol. Model. 183, 301–322.

Mosher, A.T., 1978. An Introduction to Agricultural Extension. Agricultural Development Council, New York, 114 pp.

Pebesma, E.J., 2004. Multivariable geostatistics in S: the gstat package. Comput. Geosci. 30, 683–691.

Penning de Vries, W.T.P., Rabbinge, R., Groot, J.J.R., 1997. Potential and attainable food production and food security in different regions. Philos. Trans. R. Soc. Lond. Ser. B Biol. Sci. 352, 917–928.

Plant, R.E., 2001. Site-specific management: the application of information technology to crop production. Comput. Electron. Agric. 30, 9–29.

Priya, S., Shibasaki, R., 2001. National spatial crop yield simulation using GIS-based crop production model. Ecol. Model. 136, 113–129.

Rawls, W.J., Brakensiek, D.L., Saxton, K.E., 1982. Estimation of soil-water properties. Trans. ASAE 25, 1316.

Singh, P., Boote, K.J., Rao, A.Y., Iruthayaraj, M.R., Sheikh, A.M., Hundal, S.S., Narang, R.S., 1994. Evaluation of the groundnut model PNUTGRO for crop response to water availability, sowing dates, and seasons. Field Crops Res. 39, 147–162.

Tsuji, G., Uehara, G., Balas, S. (Eds.), 1994. DSSAT v3 Crop Simulation Software. University of Hawaii, Honolulu.

Turner, M.G., Gardner, R.H., O'Neill, R.V., 2001. Landscape Ecology in Theory and Practice: Pattern and Process. Springer–Verlag, New York, 401 pp.

Wackernagel, H., 2003. Multivariate Geostatistics: An Introduction with Applications. Springer, New York, 387 pp.