

Spatial Disaggregation of Agricultural Production Data by Maximum of Entropy[°]

Richard Howitt[^] and Arnaud Reynaud[†]

2/25/2002

Abstract

In this paper we develop a dynamic data-consistent way for estimating agricultural land use choices at a disaggregate level (district-level), using more aggregate data (regional-level). The disaggregation procedure requires two steps. The first step consists in specifying and estimating a dynamic model of land use at the regional-level. In the second step, we disaggregate outcomes of the aggregate model using maximum entropy (ME). The ME disaggregation procedure is applied to a sample of California data. The sample includes 6 districts located in Central Valley and 8 possible crops, namely: Alfalfa, Cotton, Field, Grain, Melons, Tomatoes, Vegetables and Subtropical. The disaggregation procedure enables the recovery of land use at the district-level with an out-sample prediction error of 16%. This result shows that the micro behavior, inferred from aggregate data with our disaggregation approach, seems to be consistent with observed behavior.

Keywords: Disaggregation, Bayesian method, Maximum entropy, Land use.

Code JEL: C11, C44, Q12.

[°] The authors wish to thank Guilherme Marques for valuable help.

[^] Department of Agricultural and Resource Economics, University of California at Davis.

[†] Department of Agricultural and Resource Economics, University of California at Davis and LEERNA-INRA, University of Social Sciences Toulouse.

1. Introduction

In this paper we develop a data-consistent way to estimate agricultural land use choices at a disaggregate level (district-level) using more aggregate data (regional-level). We argue that such a disaggregation method is of interest in agricultural production economics for three main reasons.

The first reason deals with data availability. As Just and Pope mention (1999-a), the most significant obstacle to progress in agricultural production is the lack of better and more detailed data. Therefore, applied economists often use aggregate data for estimating relationships, which are theoretically defined, at a more disaggregate level. The main criticism of using aggregate data deals with aggregation problems; namely the failure to consider heterogeneity across agricultural producers may result in misrepresenting technology, but it also may fail to support the regularity conditions needed to recover technology from estimated structures, Just and Pope (1999-a). It follows that a valid disaggregation method would partially bypass the lack of disaggregate data in agricultural economics.

The second argument for having a valid data disaggregation tool is the increasing demand for environmental and multidisciplinary policy models. Agricultural production models are being increasingly used in conjunction with biophysical process models¹. These latter models are often calibrated at a smaller scale. Disaggregation of economic models enables more effective interaction with physical process models. A good example of the scale problems involved in multidisciplinary studies is given by the Integrated Model to Predict European Land use (IMPEL). The IMPEL project is funded by the Commission of the European Communities under Framework IV Program (Climate and Environment – DGXII). IMPEL is a spatial model aiming to integrate physical and socio-economic modeling procedures to evaluate the impact of climate change on European land use at the regional scale. It includes five interrelated modules: climate, soil and crop, land degradation, socio-economics and hydrology. One of the key challenges for the IMPEL project is the successful integration of these modules defined at different scales. For example, the soil and crop modules operate at the scale of individual soil types, whereas the socio-economic module must operate at the scale of individual farms that include one or more soil types. This aggregation method addresses the issue of defining a compatible scale for the conjunctive use of these two modules. A valid

¹ As mentioned by Antle and Capalbo (2001), assessing environmental impacts of agriculture increasingly requires the use of linked disciplinary simulation models.

disaggregation method would allow the two modules to interact at the smaller scale level without information loss.

The third reason is based on efficient model use. Given the cost of disaggregated data collection and modeling, an aggregate economic model coupled with an efficient disaggregation procedure may provide a more cost-effective approach to annual policy modeling. For example, the CAPRI model for EEC wide agricultural policy² has 200 regional spatial units, mostly based on NUTS II definition³. A disaggregation procedure would allow the policy results to be disaggregated to the NUTS III more detailed spatial units at a low computational cost. In the empirical example used in this paper, we use an agricultural production and resource use model, CVPM⁴, defined over 21 production regions that are economically homogenous. However, one of the key uses of the model is water policy planning, and the hydrologic units used for this purpose are smaller than the economic regions and are termed Detailed Analysis Units (DAUs). There are 59 DAUs within the 21 economic regions. This clearly addresses the issue of a valid data disaggregation method from regional-level to DAU-level.

For these three reasons, we think that a disaggregation method is of interest for agricultural economists. The problem of data disaggregation should be related to the much wider econometrics literature on aggregation. A rapid survey of the aggregation literature shows that there has been, since the beginning of the seventies, a lot of work done on aggregation problems in econometrics. Two main lines of research have been particularly followed:

- *Aggregation Problems*: Identification of conditions under which aggregate models reflect and provide interpretable information on the underlying micro behavior.
- *Model selection problem*: Choice between different levels of aggregation specification when the objective is to predict some aggregate (macro) phenomena.

² See Heckelevi and Britz (2000).

³ The nomenclature of territorial units for statistics (NUTS) has been created by the European Office for Statistics (Eurostat) in order to create a single and coherent structure of territorial distribution. The current nomenclature subdivides the 15 countries of the European Union into 78 NUTS level 1 territorial units, 210 NUTS level 2 units and 1093 NUTS level 3 units.

⁴ The Central Valley California Model (CVPM) has been developed by the Bureau of Reclamation, US Department of the Interior. See USBR (1997) for a detailed presentation of this model.

The main conclusion of this literature is that, when (1) the disaggregate model is correctly specified and (2) the available data are free from measurement errors, then the investigator cannot improve on a disaggregate approach. Some arguments may however support use of aggregate data. First, the model specification may be subject to less error at the aggregate. Second, there are errors in variable measurement at the disaggregate level that may roughly cancel out at the aggregate level. Third, individual equations have unobserved influences that may cancel with aggregation. Finally, the use of aggregate data may simply result from data availability considerations.

In contrast to the aggregation literature, only a few papers explicitly address the disaggregation of economic model results. In macroeconomics, the linkage problem between an aggregate models and disaggregate sectoral models of the economy has been widely recognized, Barker and Pesaran (1989). It is, for example, a common practice to use a macro model to provide estimates and forecasts of national economic aggregates and then, to divide these up by various approaches to yield disaggregate results. Yet, little is known of the implication of such macro-micro linkages. In agricultural economics, Miller and Plantinga (1999) have proposed a maximum of entropy approach (ME) for estimating land use shares using aggregate data. They use ME to disaggregate land use shares from multi-county scale to county scale. They show that ME specification encompasses the traditional pooled logistic regression as a particular case. They apply ME approach for estimating land use in three Iowa counties and for predicting its impact on soil erosion. Our paper differs from Miller and Plantinga as we explicitly model cropping pattern choices as a dynamic process within an endogenous framework.

The question under study in this paper is the following. How can we combine in a dynamic framework partial information at disaggregate-level with complete information at an aggregated level to recover information at the disaggregate level? More precisely, we want to recover from year to year land use at a small-scale level using:

- observation of cropping patterns at a larger scale
- an initial allocation of land at the small-scale-level.

Our disaggregation procedure requires two steps. First we estimate a dynamic model of land use using aggregate data. In the second step, we disaggregate large-scale land use observations to a smaller scale by ME using first-step aggregate land use forecasts as priors. The remainder of the paper

is organized as follows. Section 2 presents the dynamic aggregate model of land use and the ME disaggregating approach. In section 3 we apply our model to a sample of Californian data.

2. The Disaggregation model

2.1 The Problem

Let us consider a region made of I sub-regions termed districts in the rest of the paper. Districts are indexed by i going from 1 to I . We assume that each year we observe *at the regional-level* the land that is allocated to each crop $S_k(t)$, where $k = 1, \dots, K$ and $t = 1, \dots, T$ respectively index crops and years. Let $Y_k(t)$ be the probability of producing crop k at date t . By definition:

$$Y_k(t) = \frac{S_k(t)}{\sum_k S_k(t)} \quad \forall k, t. \quad (1)$$

Moreover, we assume that the available information at the *district-level* is limited to the land allocated to each crop in each district $s_k^i(t)$ for the r first periods, $r < T$. We only have a partial information at the district-level limited to the *r first periods*. This partial information may come from detailed district-level surveys that are not conducted every year. Let $y_k^i(t)$ be the probability of producing crop k at date t in district i :

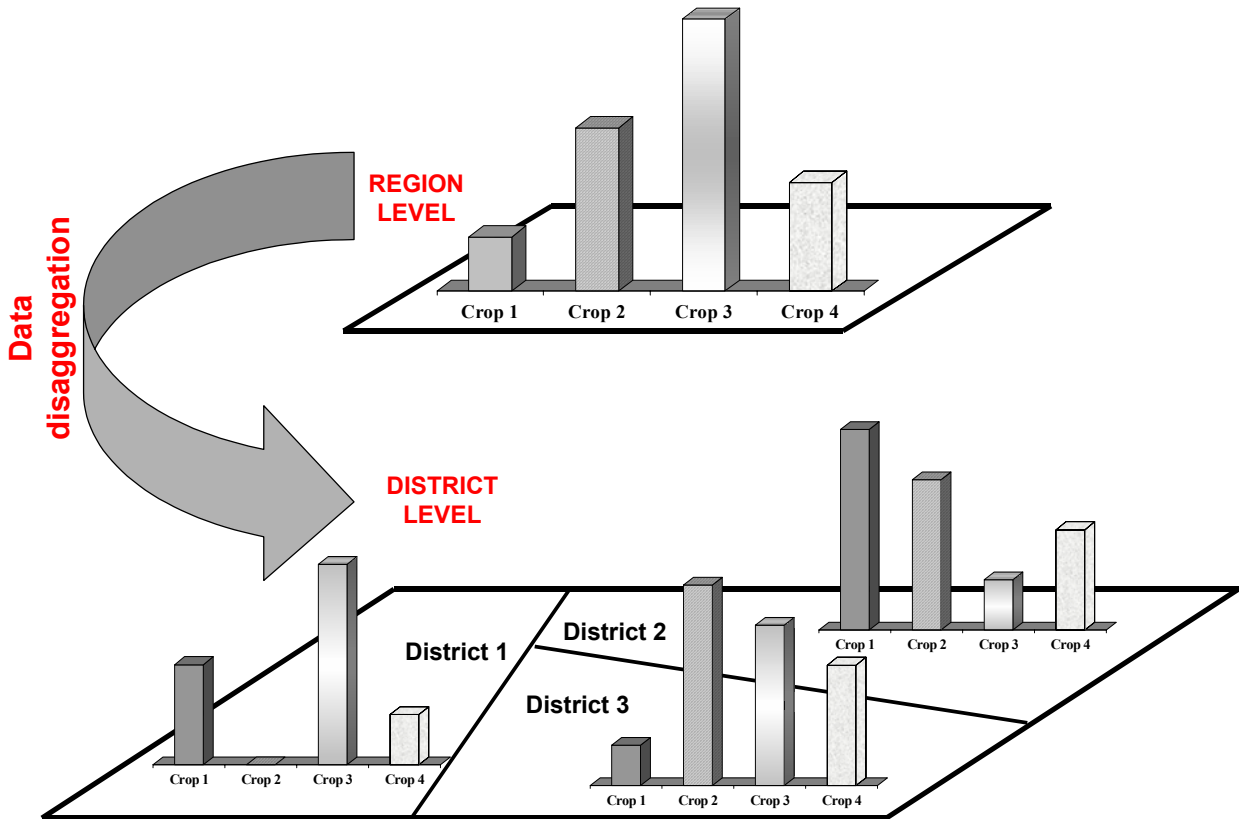
$$y_k^i(t) = \frac{s_k^i(t)}{\sum_k s_k^i(t)} \quad \forall k, t = 1, \dots, r. \quad (2)$$

The brief review of data availability in the EEC and US supports the view that we often get exhaustive data at aggregated level but only partial information at disaggregate-level. We want to combine the complete information at the regional-level for $t = 1, \dots, T$ with the partial information at the district-level for $t = 1, \dots, r < T$, in order to recover land use in each district for periods $r + 1, \dots, T$. In other words, we want to estimate $s_k^i(t)$ or in an equivalent way $y_k^i(t) \quad \forall k, i$ and $t = r + 1, \dots, T$. Notice that these estimates must satisfy the following data-compatibility constraint:

$$S_k(t) = \sum_{i=1}^I s_k^i(t) \quad \forall k, t = r + 1, \dots, T \quad (3)$$

The land allocated to a given crop in a given period must be equal to the total land allocated to this crop across all districts. Figure 1 presents the objective of the disaggregation method in the case of four crops and three districts.

Figure 1: Land use share disaggregation from regional level to district level



We want to go from land use, defined at the regional level, to a distribution of cropping patterns for each district.

2.2 The Model

We first present the dynamic land use at the district-level *we would have estimated* if data at this disaggregate level were available. Then, we turn to our disaggregation procedure that enables the recovery of cropland allocations in a dynamic framework at the district-level.

2.2.1 Dynamic land use at the district-level.

Let us first define the dynamic land-use model at the district-level. We assume that land use at a given period t , $t \in \{1, \dots, T\}$, only depends on the r previous periods. This assumption is based on the view that land use is a dynamic problem⁵ and that an agricultural producers' crop rotation horizon is finite.

Assumption 1: Land use at the district-level follows a finite non-stationary r -order Markov process.

A Markov process is a valid tool for estimating intertemporal relationships between economic variables when the current value of a variable only depends on the preceding values of the same variable. Moreover, a Markov process can be estimated even if we only observe aggregate data in the form of proportions, Lee *et al.* (1970). The assumption of non-stationarity means that we allow transition probabilities to be influenced by exogenous shocks (drought year, crop price changes...) but there is no systematic change in the dynamic relationships. This latter assumption can be easily relaxed if a richer data set is available to estimate the non-stationary Markov process. Specifying a r -order Markov process allows us to use all information at the district-level. Since we observe land use at the district-level for the r first periods, we can assign an initial probability to each Markov state for all districts.

As is well-known, any r -order Markov process may be rewritten as a more complicated *first-order* process by enlarging the space of possible states, Kijima (1997). Hence, a sequence of r -observed crops is characterized by a *first-order* Markov process. At each period, farmers choose among K possible crops at the district-level hence the Markov process is defined on the Markov space $\Omega^r = \{1, \dots, K^r\}$. There exist K^r states corresponding to the K^r possible r -tuplets. States are indexed by $j \in \{1, \dots, J\}$ with $J = K^r$. The probability associated with any state j in district i at time t is denoted $q_j^i(t)$. $q_j^i(t)$ is computed as the product of probabilities $y_k^i(t)$ corresponding to the crop

⁵ Farmer's choices are inherently dynamic. Four main types of intertemporal relationships between crops can be mentioned to justify the use of a dynamic process. First, crop rotation may be viewed as a way to reduce the loss of soil productivity due to erosion. Second, it may stabilize profits of risk-averse agricultural producers over time. Third, crop rotations may be used for breaking weed and disease cycles. Finally by reducing dependence on external inputs, crop rotation system offer the possibility of attenuating agriculture environmental impacts while maintaining profitability.

sequence indexed by j . For example, assuming a *second-order* Markov process, the probability of producing *alfalfa* in $t-1$ and *grain* in t is given by $y_{alfalfa}^i(t-1) \times y_{grain}^i(t)$. Now, let $T^i(t)$ be the $(K^r \times K^r)$ Markov transition matrix associated to land use in district i for period t . $T_{jj'}^i$ gives the probability of passing from any state $j \in \{1, \dots, K^r\}$ at date t to any state $j' \in \{1, \dots, K^r\}$ at date $t+1$. The transition probabilities satisfy the two following properties:

$$T_{jj'}^i \geq 0 \forall j, j' \quad (4)$$

$$\sum_{j'} T_{jj'}^i = 1 \forall j. \quad (5)$$

Given this notation, the probability of being in state j' in $t+1$ is given by:

$$q_{j'}^i(t+1) = \sum_{j=1}^J q_j^i(t) \cdot T_{jj'}^i(t), \quad \forall j' \in \{1, \dots, J\} \text{ and } \forall t \in \{r, \dots, T-1\}. \quad (6)$$

and the probability of producing crop k in $t+1$ is given is:

$$y_k^i(t+1) = \sum_{j=1}^J \sum_{j' \in \Psi(k)} q_j^i(t) \cdot T_{jj'}^i(t) \quad (7)$$

where $\Psi(k)$ is the set of Markov states for which crop k is produced at the last period.

We cannot directly estimate the *r-order* non-stationary transition Markov matrix at the district-level as, by assumption, we do not have data at this level⁶. Thus, the approach we follow is:

- Estimate a stationary *r-order* Markov process at the regional-level.
- Disaggregate land use at the district-level using a Generalized Maximum Entropy (GME) framework. Regional transition probability estimates are used as priors at the district-level. We then use the information in the regional level land allocations for a given year to calculate estimates of how the district-level land allocations must differ from the aggregate priors in order to be compatible with the regional-level land allocations for that year.

⁶ Estimating a *r-order* Markov process would require at least $r+1$ periods of observations and we assume we only have *r-periods* at the district-level.

2.2.2 Estimating a stationary r -order Markov process at the regional-level

We now define the dynamic land-use model at the regional-level. We assume that land use at the regional-level for a given period t depends only on the r previous periods in a stationary way.

Assumption 2: Land use at the regional-level follows a finite stationary r -order Markov process.

Two main reasons support the stationary assumption of the regional Markov process. The first reason is that, by aggregating over districts, we are losing some spatial heterogeneity. Aggregate data should be more stable than disaggregate data. The second reason is that we want to disaggregate data from regional-level to the district-level in a dynamic framework. Specifying a stationary Markov process at the regional-level allows us to omit the additional exogenous variables that may be needed to predict land use at the regional-level.

Keeping the same notation as used for district-level, the Markov process at the regional-level is defined on the Markov space Ω^r . States are indexed by $j \in \{1, \dots, J\}$ with $J = K^r$. The probability associated with any state j at time t is denoted $Q_j(t)$. It is computed as the product of probabilities $Y_k(t)$ corresponding to the crop sequence indexed by j . T is the $(K^r \times K^r)$ stationary Markov transition matrix associated with land use in district i for period t . The number of possible outcomes for any state is K and at most $K^r \times K$ transition probabilities are strictly positive. Moreover, as the sum of transition probabilities must be equal to one, $K^r \times (K - 1)$ transition probabilities have to be considered. As we have $T - r > 0$ periods of observation at the regional-level, $K^r \times (T - r)$ observations can be used.

When $T - r > K - 1$, T can be estimated using various classical statistical methods such as least chi-square, maximum likelihood and Bayesian methods. When $T - r > K - 1$ does not hold, there are more parameters to be estimated than available moment conditions and the problem is ill-posed. Using a maximum of entropy method (ME) allows a unique optimum solution to be achieved despite

this situation⁷. In the following section, we proceed by estimating the regional stationary r -order Markov transition matrix.

Let us first add an error term $e_{j'}(t)$ to equation (6) which is defined at the regional-level. Following the ME formalism we reparameterize parameters to be estimated, namely $T_{jj'}$ and $e_{j'}(t)$, in terms of unknown probability distributions.

- By definition, $T_{jj'}$ is between zero and one. It follows that we can define a set $\omega' = \{\omega_1, \dots, \omega_M\}$ of $M \geq 2$ points with $z_1 = 0$, $z_M = 1$ and a probability distribution $\{T_{jj'1}, \dots, T_{jj'M}\}$ such as

$$T_{jj'} = \sum_{m=1}^M \omega_m \cdot T_{jj'm}.$$

- The unknown disturbances $e_{j'}(t)$ may be treated in a similar way. By denoting the error support values $v' = \{v_1, \dots, v_N\}$ with $N \geq 2$ and defining $\{e_{j'1}(t), \dots, e_{j'N}(t)\}$ as the associated probabilities,

$$\text{we have: } e_{j'}(t) = \sum_{n=1}^N v_n \cdot e_{j'n}(t).$$

The problem of recovering the transition probabilities can be formulated in a standard generalized maximum entropy framework (GME). We want to estimate the probability distribution $\{T_{jj'1}, \dots, T_{jj'M}\}$ $\forall j, j'$ and $\{e_{j'1}(t), \dots, e_{j'N}(t)\}$, $\forall j', t$ solution of:

$$\text{Max}_{T, e} H(T, e) = \sum_{j=1}^J \sum_{j'=1}^J \sum_{m=1}^M T_{jj'm} \cdot \log(T_{jj'm}) + \sum_{j'=1}^J \sum_{n=1}^N \sum_{t=r}^{T-1} e_{j'n}(t) \cdot \log(e_{j'n}(t)) \quad (8)$$

subject to :

$$Q_{j'}(t+1) = \sum_{j=1}^J \left\{ Q_j(t) \cdot \sum_{m=1}^M \omega_m \cdot T_{jj'm} \right\} + \sum_{n=1}^N v_n \cdot e_{j'n}(t) \quad \forall j' \forall t \quad (9)$$

$$\sum_{j'=1}^J \sum_{m=1}^M \omega_m \cdot T_{jj'm} = 1 \quad \forall j \quad (10)$$

$$\sum_{m=1}^M T_{jj'm} = 1 \text{ and } T_{jj'm} \in [0, 1] \quad \forall j, j' \quad (11)$$

⁷ Maximum entropy is an effective tool for estimating a large number of parameters with limited data. Moreover, it eliminates problems associated with data endogeneity and collinearity. See Golan et al. (1997) for a complete description of maximum entropy methods and Howitt and Reynaud (2001) for estimating Markov transition metrics using ME.

$$\sum_{n=1}^N e_{j'n}(t) = 1 \text{ and } e_{j'n}(t) \in [0,1] \forall j', t \quad (12)$$

We seek to maximize the entropy of the probability distributions $\{T_{jj'1}, \dots, T_{jj'M}\} \forall j, j'$ and $\{e_{j'1}(t), \dots, e_{j'N}(t)\} \forall j', t$ under constraints (9)-(12). Constraint (9) defines land use at the regional-level as a stationary r -order Markov process. Constraints (11) and (12) ensure that the parameters $\{T_{jj'1}, \dots, T_{jj'M}\}$ and $\{e_{j'1}(t), \dots, e_{j'N}(t)\}$ to be estimated are defined over probability distributions. Finally, constraint (10) corresponds to the second property of transition probabilities, equation (5). It states that, for any initial Markov state, the sum of transition probabilities must be equal to 1.

The optimization program (9)-(12) constitutes a standard GME problem. As this program is convex, it has a unique solution. The interested reader may consult Golan, Judge and Miller (1996) for a complete and detailed derivation of this program's solution, $\hat{T}_{jj'm}$ and $\hat{e}_{j'n}(t)$. Point estimates both for transition probabilities and error term defined in equation (9) are recovered from the GME probability estimates $\hat{T}_{jj'm}$ and $\hat{e}_{j'n}(t)$. More formally we have:

$$\hat{T}_{jj'} = \sum_{m=1}^M \omega_m \cdot \hat{T}_{jj'm} = 1 \quad \forall j, j' \text{ and } \hat{e}_{j'}(t) = \sum_{n=1}^N \nu_n \cdot \hat{e}_{j'n}(t) \quad \forall j'. \quad (13)$$

At this point of the analysis, we have estimated the transition matrix of a Markov process using aggregate data.

2.2.3 Disaggregation at the district-level

Disaggregation of land use at the district-level requires two more steps. First, it requires estimating for each period a r -order non-stationary Markov metric at the district-level by a Generalized Cross-Entropy method (GCE). Then, using the transition probability estimates, we compute land use distribution at the district-level.

a- Estimation of the district-level non stationary r -order Markov process

At each period, the allocation of land between crop at the district-level must be compatible with the observed allocation of land at the regional-level. This data-compatibility constraint, first described by equation (3), can be rewritten as:

$$\sum_{i=1}^I \left(\sum_{j=1}^J \sum_{j' \in \Psi(k)} q_j^i(t) \cdot T_{jj'}^i(t) \right) \cdot s^i + e_k(t) = S_k(t+1) \quad \forall k = 1, \dots, K \quad (14)$$

The data-compatibility constraint states that at each period the total expected land predicted to produce crop k in all districts must be equal to the observed surface allocated at the regional-level to crop k , S_k plus an error term e_k . Notice that, given we only observe land use at the district-level for the r first periods, the probability of being in state j in district i at time t , $q_j^i(t)$, can initially only be computed for $t = r$. For $t = r + 1, \dots, T$, this probability is endogeneously computed from period to period.

Estimating a r -order non-stationary Markov matrix at the district-level can be formulated as a special generalized cross-entropy framework (GCE)⁸. For a given period⁹, we have to solve the following nonlinear optimization program:

$$\text{Min } H(T, e) = \sum_{i=1}^I \sum_{j=1}^J \sum_{j'=1}^J T_{jj'}^i \cdot \log(T_{jj'}^i / \hat{T}_{jj'}^i) + \sum_{k=1}^K \sum_{n=1}^N e_{kn} \cdot \log(e_{kn}) \quad (15)$$

$\{T^1, \dots, T^I\}, \{e_1, \dots, e_K\}$

subject to :

$$\sum_{i=1}^I \left(\sum_{j=1}^J \sum_{j' \in \Psi(k)} q_j^i \cdot T_{jj'}^i \right) \cdot s^i + \sum_{n=1}^N \zeta_n \cdot e_{kn} = S_k \quad \forall k \quad (16)$$

$$\sum_{j'=1}^J T_{jj'}^i = 1 \quad \forall i = 1, \dots, I \quad \text{and} \quad T_{jj'}^i \in [0, 1] \quad (17)$$

$$\sum_{n=1}^N e_{kn} = 1 \quad \forall k = 1, \dots, K \quad \text{and} \quad e_{kn} \in [0, 1] \quad (18)$$

where $\{\zeta_1, \dots, \zeta_N\}$ with $N \geq 2$ is the support associated with probabilities $\{e_{k1}, \dots, e_{kN}\}$ such as

$e_k = \sum_{n=1}^N \zeta_n \cdot e_{kn} \quad \forall k$. We seek to minimize the cross-entropy of the probability distribution for the

⁸ The cross entropy between two distributions p and q is: $I(q, p) = \sum_i p_i \cdot \log(p_i / q_i)$. It was first introduced by Kullback (1959) but was explicitly called cross-entropy by Good (1963). The cross-entropy measures the distance between two distributions. The cross-entropy is minimized when the two distributions are identical.

⁹ For reasons of simplicity we omit the time index in the following program.

Markov transition matrixes and the entropy for the error term. Constraints (18) ensure that $\{e_{k1}, \dots, e_{kN}\}$ is a probability distribution. Constraint (17) corresponds to the second property of transition probabilities, equation (5). Constraint (16) is the data-compatibility constraint.

The intuition of this program is the following. Let's first consider the problem without the data-compatibility constraint (16) holding. The solution of this relaxed program is $T_{jj'}^i = \hat{T}_{jj'}^i, \forall i, j, j'$ and $e_{kn} = 1/N \forall k, n$. Without any district heterogeneity, the estimated district-level transition probabilities are given by the stationary regional Markov matrix and the error distribution for each crop is uniform with an expected value of zero. Suppose we now impose condition (16). If there is some district heterogeneity, transition and errors probabilities must be changed from the previous solution to satisfy the data-compatibility constraint. In this case the optimal parameter estimates the tradeoff between deviations from the priors and information recovery. It can be shown, as previously, that this non-linear optimization program has a unique solution in (T, e) , see Golan, Judge and Miller (1996).

The result of this step is that we obtain for each district and for a given year, a Markov transition matrix associated with land use at the district-level and a point estimate of the error term in equation (12):

$$\hat{T}_{jj'}^i \text{ and } \hat{e}_k = \sum_{n=1}^N \zeta_n \cdot \hat{e}_{kn} \quad (19)$$

Finally we should mention that this ME approach provides an easy way to take into account out-of-sample information. Out-of-sample information may either consist of additional constraints in the optimization program or in particular priors for the Markov transition metrics. For example, some specific physical constraints (quality of soil, water availability, agronomic constraints) may prevent farmers in a given district from producing particular crops. This information may be added to the disaggregation program with additional constraints. In the same way, we could have out-of-the sample information on transition probabilities for a specific district. This information may be added to the model via a change of transition probability priors.

b- Land use at the district-level

From the previous estimates we can recover land use at the district-level. The probability of producing crop k in $t+1$ in district i and the expected land allocated to this crop are respectively given by:

$$\hat{y}_k^i(t+1) = \sum_{j=1}^J \sum_{j' \in \Psi(k)} q_j^i(t) \cdot \hat{T}_{jj'}^i(t) \text{ and } \hat{s}_k^i(t+1) = \hat{y}_k^i(t+1) \cdot s^i \quad (20)$$

c- Dynamic district-level land use

The disaggregation program is solved year by year. Since we have observed land use at the district-level for years 1 to r , we can compute the probability associated to each Markov state for year r : $q_j^i(r) \forall i, j$. The optimization program (15)-(18) can then be solved for year $r+1$ and the solution defined by (20) allows us to compute the probability associated with each Markov state in year $r+1$: $\hat{q}_j^i(r+1) \forall i, j$.

$$\hat{q}_{j'}^i(r+1) = \sum_{j=1}^J q_j^i(r) \cdot \hat{T}_{jj'}^i(r) \quad \forall j'. \quad (21)$$

For periods $r+1$ to T , a closed-form loop solution for the Markov state probabilities is obtained in the same way using the previous year's estimates. Hence for each period the program (15)-(18) is completely defined. The data disaggregation from the regional to district scale can therefore be performed year by year. In the next section, we apply this framework to a sample of Californian data.

3. An application of the disaggregation method to California

In California, the Department of Water Resources the US Bureau of Reclamation has developed a regional model of irrigated agricultural production that simulates the decisions of agricultural producers in the Central Valley of California. The Central Valley California Model (CVPM) is implemented as part of an integrated analysis with surface water hydrology, groundwater, agricultural economics land use and water transfer analysis. The model includes 21 production regions and 26 categories of crops. However, many water management issues require a smaller scale-level analysis. Often water issues are analyzed by the DWR at the Detailed Analysis Unit (DAU) level, which is generally defined by hydrologic features or boundaries of organized water service agencies.

In the major agricultural areas, a DAU typically includes 100,000 to 300,000 acres. A typical CVPM region is made of four to five DAUs.

Disaggregation of CVPM regions to the DAU level is of great interest as it would allow the agricultural production model and the hydrologic water models to interact more effectively.

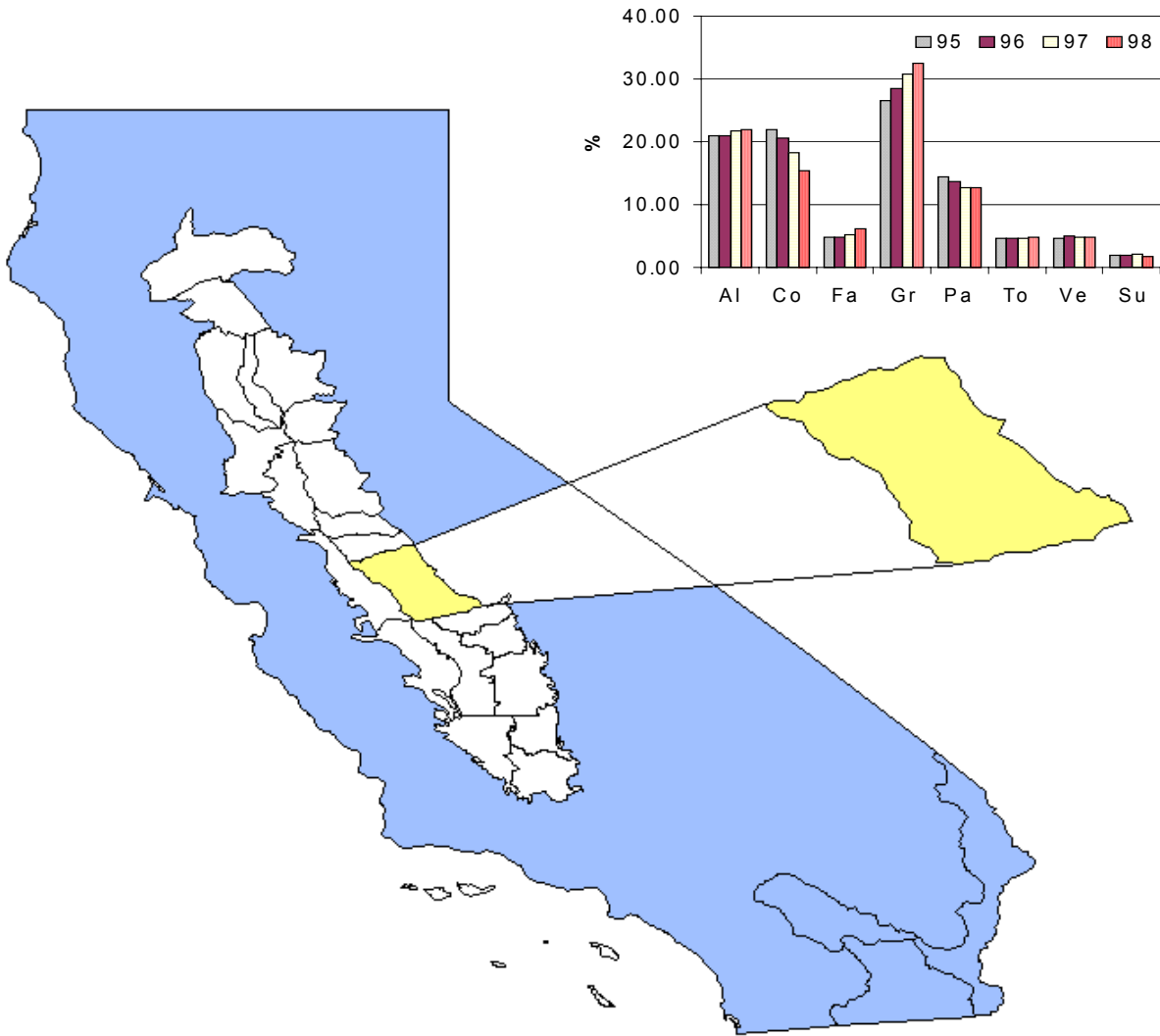
3.1 Data

We apply our disaggregation procedure to a set of Californian data. The area considered is CVPM region 13 located in Central Valley of California (See Map 1 and Table A.1 in Appendix A). This region is made of six Detailed Analysis Units (DAU): *Merced*, *Merced Stream Group*, *El Nido-Stevinson*, *Madera-Chowchilla*, *Adobe - Valley Eastside* and *Gravelly Ford*. We selected this region as it has more DAUs than others and because the DAUs are quite heterogeneous both in terms of size and cropping patterns. This makes the disaggregation procedure more interesting and also more difficult.

We consider eight possible crop groups, namely: *Alfalfa*, *Cotton*, *Field*, *Grain*, *Melon*, *Tomatoes*, *vegetables* and *subtropical*. In what follows, the first letter indexes each crop. Therefore, we have $k \in \{A, C, F, G, M, T, V, S\}$. Eleven years of data on land use at the regional-level (1988 to 1998) are available. In order to perform in-sample and out-of-sample estimates, we only use years 1988-94 to estimate the Markov transition matrix at the aggregated level. In terms of the notation in the previous section, we have $T=7$ and $K=8$. Eleven years of data on land use at the DAU level were also available. We use observations for year 1988 and 1989 to define the initial Markov probabilities at the DAU level and assume that land use at the DAU-level can be described by a *second-order* Markov process.

The objective of the disaggregation procedure is to recover land use at the DAU-level for years 1990 to 1998. Observation of the actual land use from 1990 to 1998 allows us to test the accuracy of the disaggregation procedure. See Table A.2 in appendix A for a complete presentation of DAU-level data.

Figure 2: Location and characteristics of CVPM region 13



3.2 Regional stationary second-order Markov matrix estimate

Data in Table A.1 in Appendix A are used to estimate the stationary *second-order* Markov process at the regional-level. A Markov state is a pair of crops observed during two consecutive years. As a consequence, there are 64 possible states each year. States are indexed in the following way: state $k_t k'_{t-1}$ observed in t means that crop k' was produced in $t-1$ and crop k in t for $k, k' \in \{A, C, F, G, M, T, V, S\}$. As we have $T - r \leq K - 1$, the problem is ill-posed and use of ME for estimating the Markov matrix is justified.

Table 1: Transition probability estimates \hat{T}_{jj}^i at the region-level

	$A_{t+1}A_t$	$C_{t+1}A_t$	$F_{t+1}A_t$	$G_{t+1}A_t$	$M_{t+1}A_t$	$T_{t+1}A_t$	$V_{t+1}A_t$	$S_{t+1}A_t$		$A_{t+1}C_t$	$C_{t+1}C_t$	$F_{t+1}C_t$	$G_{t+1}C_t$	$M_{t+1}C_t$	$T_{t+1}C_t$	$V_{t+1}C_t$	$S_{t+1}C_t$	
ALFALFA	A_tA_{t-1}	0.00	0.50	0.20	0.00	0.30	0.00	0.00	0.00	C_tA_{t-1}	0.50	0.00	0.19	0.00	0.16	0.00	0.15	0.00
	A_tC_{t-1}	0.00	0.26	0.00	0.50	0.11	0.00	0.13	0.00	C_tC_{t-1}	0.00	0.00	0.00	0.50	0.41	0.00	0.00	0.09
	A_tF_{t-1}	0.08	0.06	0.12	0.06	0.41	0.05	0.09	0.14	C_tF_{t-1}	0.05	0.12	0.11	0.05	0.52	0.05	0.08	0.03
	A_tG_{t-1}	0.00	0.23	0.00	0.60	0.17	0.00	0.00	0.00	C_tG_{t-1}	0.00	0.72	0.00	0.22	0.00	0.06	0.00	0.00
	A_tM_{t-1}	0.91	0.00	0.00	0.00	0.00	0.10	0.00	0.00	C_tM_{t-1}	0.43	0.00	0.00	0.57	0.00	0.00	0.00	0.00
	A_tT_{t-1}	0.20	0.12	0.10	0.11	0.08	0.12	0.11	0.16	C_tT_{t-1}	0.16	0.19	0.11	0.16	0.08	0.13	0.11	0.07
	A_tV_{t-1}	0.32	0.10	0.08	0.09	0.04	0.12	0.08	0.16	C_tV_{t-1}	0.24	0.18	0.09	0.21	0.03	0.13	0.08	0.04
	A_tS_{t-1}	0.17	0.13	0.10	0.12	0.09	0.12	0.11	0.16	C_tS_{t-1}	0.15	0.19	0.11	0.15	0.10	0.12	0.12	0.07
	$A_{t+1}F_t$	$C_{t+1}F_t$	$F_{t+1}F_t$	$G_{t+1}F_t$	$M_{t+1}F_t$	$T_{t+1}F_t$	$V_{t+1}F_t$	$S_{t+1}F_t$		$A_{t+1}G_t$	$C_{t+1}G_t$	$F_{t+1}G_t$	$G_{t+1}G_t$	$M_{t+1}G_t$	$T_{t+1}G_t$	$V_{t+1}G_t$	$S_{t+1}G_t$	
FIELD	F_tA_{t-1}	0.22	0.26	0.07	0.26	0.00	0.06	0.08	0.05	$G_{t-1}A$	0.50	0.00	0.00	0.00	0.50	0.00	0.00	0.00
	F_tC_{t-1}	0.22	0.25	0.06	0.30	0.05	0.00	0.07	0.05	G_tC_{t-1}	0.30	0.00	0.15	0.41	0.00	0.09	0.00	0.06
	F_tF_{t-1}	0.17	0.18	0.10	0.20	0.09	0.08	0.10	0.08	G_tF_{t-1}	0.06	0.03	0.18	0.05	0.50	0.03	0.09	0.07
	F_tG_{t-1}	0.09	0.14	0.00	0.28	0.50	0.00	0.00	0.00	G_tG_{t-1}	0.00	0.41	0.00	0.50	0.00	0.00	0.09	0.00
	F_tM_{t-1}	0.22	0.27	0.10	0.29	0.06	0.06	0.00	0.00	G_tM_{t-1}	0.10	0.50	0.00	0.24	0.16	0.00	0.00	0.00
	F_tT_{t-1}	0.15	0.15	0.11	0.15	0.11	0.11	0.12	0.10	G_tT_{t-1}	0.12	0.23	0.10	0.15	0.07	0.12	0.12	0.11
	F_tV_{t-1}	0.15	0.16	0.11	0.17	0.09	0.10	0.11	0.10	G_tV_{t-1}	0.12	0.32	0.07	0.18	0.03	0.12	0.08	0.07
	F_tS_{t-1}	0.15	0.15	0.12	0.15	0.11	0.11	0.12	0.10	G_tS_{t-1}	0.12	0.20	0.10	0.16	0.08	0.12	0.13	0.11
	$A_{t+1}M_t$	$C_{t+1}M_t$	$F_{t+1}M_t$	$G_{t+1}M_t$	$M_{t+1}M_t$	$T_{t+1}M_t$	$V_{t+1}M_t$	$S_{t+1}M_t$		$A_{t+1}T_t$	$C_{t+1}T_t$	$F_{t+1}T_t$	$G_{t+1}T_t$	$M_{t+1}T_t$	$T_{t+1}T_t$	$V_{t+1}T_t$	$S_{t+1}T_t$	
MELON	M_tA_{t-1}	0.00	0.00	0.00	0.89	0.00	0.11	0.00	0.00	T_tA_{t-1}	0.14	0.16	0.09	0.28	0.14	0.07	0.07	0.06
	M_tC_{t-1}	0.25	0.24	0.13	0.38	0.00	0.00	0.00	0.00	T_tC_{t-1}	0.15	0.29	0.04	0.34	0.12	0.06	0.00	0.00
	M_tF_{t-1}	0.12	0.09	0.15	0.09	0.34	0.04	0.10	0.07	T_tF_{t-1}	0.14	0.14	0.12	0.15	0.14	0.10	0.11	0.10
	M_tG_{t-1}	0.32	0.41	0.00	0.00	0.14	0.00	0.09	0.04	T_tG_{t-1}	0.24	0.22	0.00	0.28	0.16	0.00	0.05	0.05
	M_tM_{t-1}	0.12	0.22	0.05	0.11	0.50	0.00	0.00	0.00	T_tM_{t-1}	0.15	0.17	0.10	0.20	0.14	0.08	0.09	0.07
	M_tT_{t-1}	0.15	0.16	0.10	0.14	0.11	0.11	0.12	0.11	T_tT_{t-1}	0.13	0.14	0.12	0.14	0.13	0.12	0.11	0.12
	M_tV_{t-1}	0.15	0.17	0.08	0.20	0.08	0.10	0.11	0.11	T_tV_{t-1}	0.13	0.14	0.12	0.14	0.14	0.11	0.11	0.11
	M_tS_{t-1}	0.15	0.16	0.10	0.14	0.11	0.10	0.13	0.11	T_tS_{t-1}	0.13	0.14	0.12	0.14	0.13	0.12	0.11	0.12
	$A_{t+1}V_t$	$C_{t+1}V_t$	$F_{t+1}V_t$	$G_{t+1}V_t$	$M_{t+1}V_t$	$T_{t+1}V_t$	$V_{t+1}V_t$	$S_{t+1}V_t$		$A_{t+1}S_t$	$C_{t+1}S_t$	$F_{t+1}S_t$	$G_{t+1}S_t$	$M_{t+1}S_t$	$T_{t+1}S_t$	$V_{t+1}S_t$	$S_{t+1}S_t$	
VEGETABLES	V_tA_{t-1}	0.17	0.18	0.09	0.31	0.19	0.06	0.00	0.00	S_tA_{t-1}	0.15	0.17	0.09	0.21	0.14	0.08	0.10	0.08
	V_tC_{t-1}	0.16	0.31	0.00	0.41	0.12	0.00	0.00	0.00	S_tC_{t-1}	0.15	0.26	0.00	0.34	0.13	0.06	0.00	0.06
	V_tF_{t-1}	0.14	0.15	0.11	0.16	0.14	0.10	0.10	0.11	S_tF_{t-1}	0.14	0.15	0.11	0.16	0.13	0.10	0.11	0.10
	V_tG_{t-1}	0.22	0.24	0.04	0.27	0.16	0.00	0.08	0.00	S_tG_{t-1}	0.21	0.24	0.09	0.29	0.18	0.00	0.00	0.00
	V_tM_{t-1}	0.16	0.17	0.08	0.19	0.14	0.08	0.08	0.09	S_tM_{t-1}	0.16	0.16	0.09	0.20	0.14	0.08	0.10	0.07
	V_tT_{t-1}	0.14	0.13	0.12	0.14	0.13	0.11	0.12	0.12	S_tT_{t-1}	0.13	0.13	0.11	0.14	0.13	0.12	0.12	0.11
	V_tV_{t-1}	0.14	0.14	0.11	0.15	0.13	0.11	0.11	0.12	S_tV_{t-1}	0.13	0.14	0.11	0.15	0.13	0.11	0.12	0.11
	V_tS_{t-1}	0.13	0.13	0.12	0.14	0.13	0.12	0.11	0.12	S_tS_{t-1}	0.13	0.13	0.12	0.14	0.13	0.12	0.12	0.12

Note: Table 2 gives transition probabilities of passing from one Markov state to another. For example, the transition probability of producing Cotton in t after having produced Alfalfa in t and Grain int-1 is 0.23. It corresponds to the probability of passing from state AG to state CA.

Before solving the non-linear optimization program (8)-(12), we have to choose support values for the errors and parameters. The natural bounds for the $\omega = \{\omega_1, \dots, \omega_M\}$ terms are zero and one. Still there remains the choice of the number M of support values. Since previous studies have shown

that increasing the support space from 3 to 5 points has little effect on the estimates¹⁰, we fix M equal to 3. Then $\varpi' = \{0, 0.5, 1\}$. The choice of error support $v' = \{v_1, \dots, v_N\}$ is subject to more controversy and clearly depends on properties of errors e . By reference to the Chebyshev's inequality, some authors determine the bounds using a 3σ rule, Golan, Judge and Miller (1996). This is the rule followed here. The number N of values for error support is 3. The non-linear optimization program (8)-(12) is solved using GAMS.

Table 1 shows the point estimates of the Markov process transition probabilities. The probability of growing a crop differs according the crop patterns of the two previous years. For example, the average probability of producing COTTON in $t+1$ goes from 0.17, if SUBTROPICAL is produced in t , to 0.22 if GRAIN is produced in t . Moreover, given that ALFALFA has been produced in t , the probability of growing COTTON in $t+1$ varies from 0, if MELON is produced in $t-1$, to 0.5 if ALFALFA is produced in $t-1$. These results support our land use specification at the regional-level as, first a dynamic process and second, a second-order Markov process. In order to evaluate the Markov matrix, we compare the predicted crop shares from 1988 to 1998 with the observed aggregate shares. Table 2, presents closed-loop simulations of the Markov matrix¹¹. For out-sample simulations, we assume that we can observe the true distribution of land use at the DAU-level for years 1993 and 1994.

Table 2: Simulated land use shares per crop $\hat{Y}_k(t)$ at the regional-level (in %)

	90 ^(a)		91 ^(a)		92 ^(a)		93 ^(a)		94 ^(a)		95 ^(b)		96 ^(b)		97 ^(b)		98 ^(b)	
	Obs	Pre	Obs	Pre	Obs	Pre	Obs	Pre	Obs	Pre	Obs	Pre	Obs	Pre	Obs	Pre	Obs	Pre
A	18.2	18.1	18.6	18.7	17.9	18.0	18.8	18.5	19.6	17.9	20.9	17.8	21.0	18.0	21.7	18.1	22.0	18.2
C	20.7	22.1	23.3	22.0	22.5	22.6	23.9	21.9	22.5	22.4	22.0	22.4	20.6	22.0	18.2	22.2	15.3	22.3
F	6.7	5.3	4.7	5.3	4.4	5.3	4.5	5.3	5.1	5.3	4.8	5.2	4.8	5.2	5.2	5.1	6.1	5.2
G	27.6	28.1	27.0	27.6	30.0	28.5	29.0	28.5	29.6	28.6	26.6	29.0	28.4	28.9	30.8	28.8	32.4	28.5
M	18.4	16.4	16.9	16.4	15.6	15.2	14.2	15.3	13.4	15.1	14.4	15.1	13.6	15.3	12.7	15.1	12.7	15.1
T	2.4	3.1	2.7	3.2	2.6	3.3	3.4	3.4	3.6	3.4	4.6	3.3	4.7	3.3	4.6	3.3	4.8	3.4
V	3.6	4.2	4.5	4.2	4.6	4.3	3.7	4.3	3.7	4.4	4.6	4.3	4.9	4.4	4.9	4.4	4.8	4.3
S	2.5	2.7	2.3	2.6	2.4	2.8	2.5	2.8	2.4	2.9	2.0	2.8	2.0	2.8	2.0	2.9	1.8	2.9

Notes: Obs and Pre respectively give the observed and the predicted land use shares.
^(a) for in-sample estimates ^(b) for out-sample estimates

¹⁰ See for example Golan, Judge and Miller (1996). These authors have shown that passing from 2 points to 3 substantially decreases the mean-square-error of estimates. More increase in M is shown to only result in a smaller improvement.

¹¹ Given aggregate land use for 1988 and 1989, the Markov metric gives a prediction of land use shares for 1989. Then given observed land use share in 1988 and predicted in 1989, we estimate land use for 1990. Finally, land allocation in 1991 is based on predicted land use shares in 1989 and 1990...

The Markov metrics performs quite well, both in term of the level of prediction and a measure of prediction variation. Let us define the Percentage Absolute Predicted Error (PAPE) for a given crop k as:

$$PAPE_k = \left| \frac{Y_k - \hat{Y}_k}{Y_k} \right| * 100 \quad (22)$$

where Y_k and \hat{Y}_k respectively represent the observed and estimated probability of producing crop k .

The average crop Percentage Absolute Predicted Error (PAPE) for in-sample years 1990 to 1994 is 10.40%. For out-sample, the average PAPE is 19.04%. If we do not take into account subtropical crops, for which the PAPE is high but only represents a small proportion of total surface, the average out-of-sample PAPE is 14.46%. It increases from 9.80% in 1995 to 21.41% in 1998. Hence, the Markov metrics enable us to recover the aggregated surfaces allocated to crops in a precise way.

3.3 Disaggregation at the DAU level

In this section, we present the final results of the disaggregation method, namely the distribution of land use per DAU and per year.

As the total agricultural land use varies from year to year at the regional-level, Table A.1 Appendix 1, the data-compatibility constraint (14) must slightly be modified allowing the DAU size to vary from year to year:

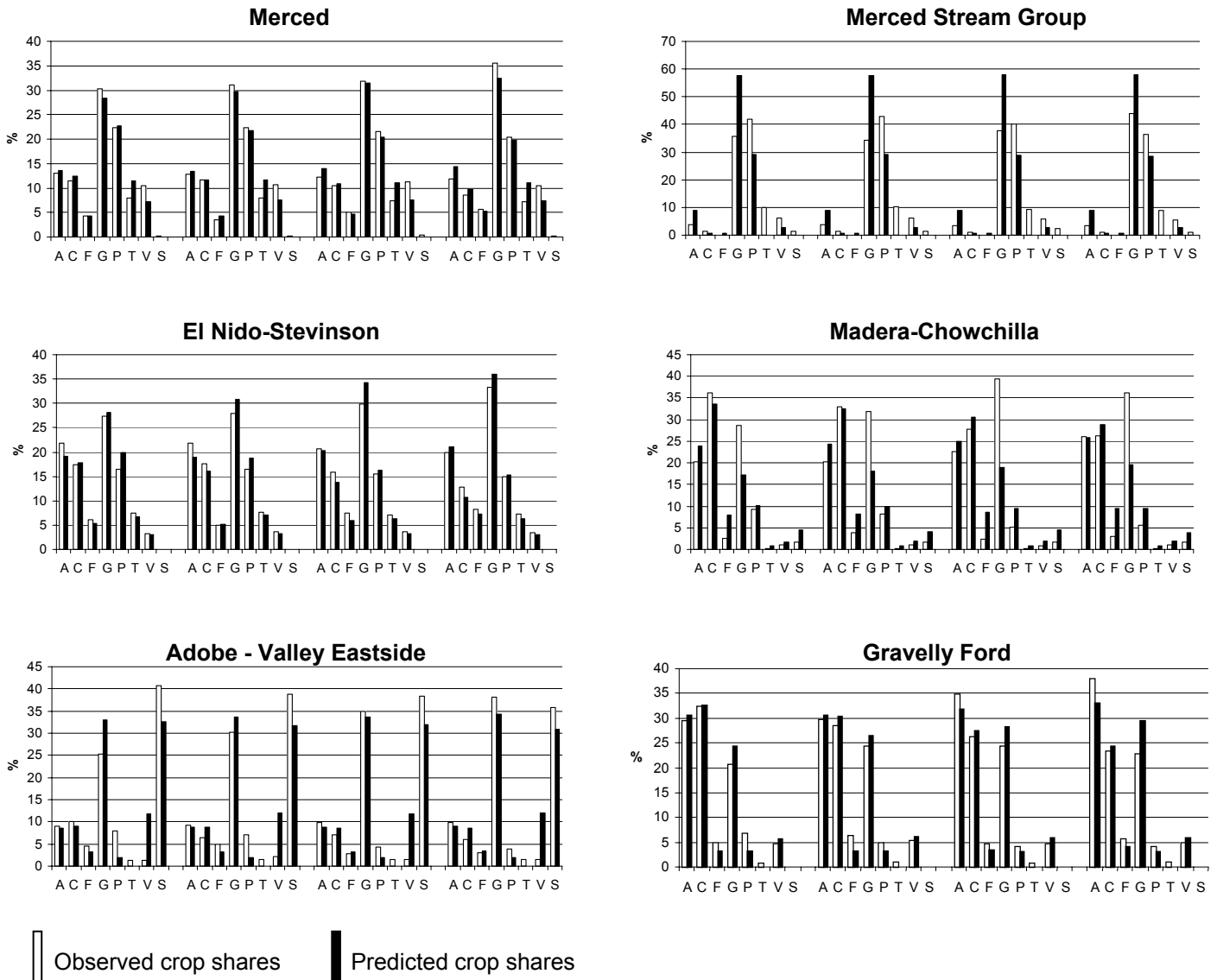
$$\sum_{i=1}^I \left(\sum_{j=1}^J \sum_{j' \in \Psi(k)} q_j^i(t) \cdot T_{jj'}^i(t) \right) \cdot s^i(t) + e_k(t) = S_k(t+1) \quad \forall k = 1, \dots, K \quad (23)$$

In equation (23), we add some more information at the DAU level, namely the total agricultural area of each DAU at each date, $s^i(t)$. In a more complex framework this variable could be endogenously predicted.

Since, in this empirical example, we also observe the true land distribution at the DAU level, we can evaluate the accuracy of the DAU disaggregation procedure. Figure 2 presents a comparison of land use shares resulting from the disaggregation procedure with the observed land use shares at the

DAU level for the out-of-sample years 1995 to 1998. The disaggregation framework is initialized assuming that we observe DAU land use in 1993 and 1994. Complete disaggregation results are also reported in Table B.1, Appendix B.

Figure 2: Observed and predicted crop shares $\hat{y}_k^i(t)$ at DAU level for year 1995-98



For each district, Figure 2 compares the out-of-sample observed and the predicted crop shares for years 1995, 96, 97 and 98. The disaggregated crop shares are on average close to the observed shares. This is especially true for four DAUs: Merced, El Nido-Stevinson, Adobe Valley-Eastside and Gravelly Ford. It is interesting to notice that for Merced, El Nido-Stevinson and Adobe Valley-Eastside, some long-term trends such as the grain share increase, are predicted well. The accuracy of the disaggregation for Merced Stream Group and Madera-Chowchilla seems to be less precise. The

main reason is that, for these two DAUs, there is an important exogenous change of grain land shares between 1994 and 1995¹². The stationary regional Markov process does not capture this change and the departures from the aggregate priors are not large enough to adjust the DAU transition probabilities. We should notice that specifying a non-stationary Markov process at the regional level would not improve predictions for these two DAUs as cropping patterns at the regional level do not exhibit drastic changes between 1994 and 1995. However, if the modeler is aware of such an exogenous change occurring, they can easily modify the transition probability priors.

Another useful measure of prediction errors is given by the Weighted Percentage Absolute Predicted Error (WPAPE) for each DAU and at the regional-level. For DAU i , the WPAPE is defined by:

$$WPAPE^i = \sum_{k=1}^K y_k^i \cdot \left| \frac{y_k^i - \hat{y}_k^i}{y_k^i} \right| \quad (24)$$

and at the regional-level by:

$$WPAPE = \sum_{i=1}^I \frac{S^i}{S} \cdot WPAPE^i \quad (25)$$

The DAU Weighted PAPE is the sum of crop PAPE weighted by the land allocated by each crop. The regional weighted PAPE is the sum of DAU PAPE weighted by the size of each DAU. Weighted PAPE results are presented in Table 3.

Table 3: DAU-level and regional weighted PAPE (in %)

	90 ^(a)	91 ^(a)	92 ^(a)	93 ^(a)	94 ^(a)	95 ^(b)	96 ^(b)	97 ^(b)	98 ^(b)
Merced	5.1	19.4	22.0	30.1	32.5	10.5	10.2	12.0	15.0
Merced Stream Group	4.9	7.6	4.1	16.3	28.3	54.7	58.1	52.3	40.4
El Nido-Stevinson	5.3	13.5	12.8	17.0	20.3	9.6	10.7	10.2	8.3
Madera-Chowchilla	8.7	6.3	16.8	14.2	27.3	28.2	28.0	41.1	33.2
Adobe – Valley Eastside	8.8	25.3	29.6	31.1	27.6	36.9	31.4	25.1	26.6
Gravelly Ford	8.6	7.9	8.6	8.2	9.3	12.0	11.4	12.3	17.7
CVPM Region 13	6.8	12.9	15.3	18.0	22.0	15.3	15.4	17.2	16.4

Notes: ^(a) for in-sample estimates ^(b) for out-sample estimates

Both for in-sample and out-of-sample data, the weighted PAPE values show a reasonable level of precision given the inherent difficulty of data disaggregation. The weighted PAPE increases from

¹² For Merced Stream Group the proportion of land allocated to grain goes down from 51.5% in 1993 and 57.7% in 1994 to 35.8% in 1995. For Madera-Chowchilla, it goes from 22.8% in 1993 and 18.2% in 1994 to 28.7% in 1995.

1990 to 1994 and from 1995 to 1998, as the predictions are done in closed-form loop. The high weighted PAPE for the Merced Stream Group DAU should be considered in context with its small size (less than 3% of the regional area in 1988).

3.3.2 Measuring information recovery

Finally we want to measure the information gains from the disaggregation procedure. We need to define a quantitative measure of information change due to disaggregation. This measure should have the following properties.

1. The measure of potential gain increases monotonically with the heterogeneity of the disaggregated sample.
2. The gain from disaggregating a uniform set of samples is zero.
3. The measure is invariant to changes in the number of disaggregated samples and the variability of the aggregated sample.
4. The measure has an information theoretic interpretation.

Let us define the cross-entropy between the aggregate observed land shares, y_k , and the true disaggregate land shares, y_k^i , as:

$$CE = \sum_i \sum_k y_k \cdot \ln \left(\frac{y_k}{y_k^i} \right) \quad (25)$$

and the cross-entropy between the disaggregate estimate land shares, \hat{y}_k^i , and the true disaggregate land shares, y_k^i , as:

$$C\hat{E} = \sum_i \sum_k \hat{y}_k^i \cdot \ln \left(\frac{\hat{y}_k^i}{y_k^i} \right) \quad (26)$$

First assume that we do not have any information at the district level. The disaggregation procedure would result in attributing the aggregate land share distribution y_k to each district. CE , which is an aggregate entropy-measure of the distance between the distributions y_k and y_k^i , measures how far we are from the actual district shares when we attribute the aggregate land use distribution to each district. Now assume that we use our disaggregation procedure to calculate the district land use distributions from the aggregated distribution. $C\hat{E}$ is an aggregate measure, in term of entropy, of how

far the posteriors \hat{y}_k^i are from the true distributions y_k^i . Hence, the Disaggregation Informational Gain (DIG) from the disaggregation procedure is defined by:

$$DIG = 1 - \frac{\sum_i \sum_k \hat{y}_k^i \cdot \ln\left(\frac{\hat{y}_k^i}{y_k^i}\right)}{\sum_i \sum_k y_k^i \cdot \ln\left(\frac{y_k^i}{y_k^i}\right)} = 1 - \frac{C\hat{E}}{CE} \quad (27)$$

The *DIG* is a measure of the proportion of district-level heterogeneity that is recovered. In case of a perfect disaggregation where $\hat{y}_k^i = y_k^i \forall k, i$, the *DIG* is equal to 1. In such a case, we are recovering 100% of the heterogeneity at the district level. In the case of no disaggregation procedure, we have $\hat{y}_k^i = y_k \forall k, i$ and the *DIG* is equal to 0, and we recover no information at the district level. In all other cases, the *DIG* is between 0 and 1. The *DIG* measure increases as the district posteriors get closer to the true district land use distributions y_k^i .

As an illustrative empirical example, we compute the Disaggregation Informational Gain for the out-of-sample years 1995 to 1998. The *DIG* are respectively equal to 56.34%, 69.03%, 62.08% and 65.54% for years 1995 to 1998. This means that the disaggregation procedure recovers a substantial part of the district heterogeneity (on average 63.75%). Moreover, it is interesting to notice that the proportion of information recovered does not decrease with time, as might be expected, since the disaggregation is calculated annually in a closed-loop form.

4. Conclusion

In this paper, we have addressed the issue of dynamic data disaggregation in agricultural economics. We have developed a data-consistent method to estimate cropping choices by farmers at a disaggregate level (district-level) using data from a more aggregate (regional-level) source. Our disaggregation procedure requires two steps. The first step consists of specifying a model of crop allocation and estimating it using aggregate data. In the second step, we disaggregate outcomes of the regional-model using maximum of entropy (ME). Two points should be noticed:

- First, we explicitly model aggregate cropping pattern choices as a dynamic process by using a Markov process. We believe that farmer's crop choices are dynamic *per se*
- Second, we use a ME approach for downscaling data. The ME approach gives an optimal solution using the Kullback-Leibler cross-entropy criterion in cases where traditional inversion methods do not result in identifying a set of parameters.

The resulting disaggregate data are consistent with priors, given by the Markov metrics, and with the data, given by the aggregate land use shares.

We have applied our disaggregation procedure to a sample of Californian data. The sample includes six districts for which we want to recover land use for eight possible crops, namely: Alfalfa, Cotton, Field, Grain, Melons, Tomatoes, Vegetables and Subtropical. Eleven years of cropping patterns are available, from 1988 to 1998. A second-order Markov process is specified as representing aggregate crop choices. The estimate of the aggregate Markov process is based on the years 1988 to 1994. This allows us to have in-sample crop predictions for years 1990 to 1994 and out-of-sample predictions for the rest of the periods. We have shown that the quality of predictions at the disaggregate level is relatively good. For out-sample estimates, the regional-level weighted PAPE is between 15.3% and 17.2% according to the year considered. These results show that, the district-level behavior inferred from aggregate data with our disaggregation approach, are consistent with the observed behavior.

The disaggregation approach partially bypasses one of the most significant obstacles to progress in agricultural production: the lack of *better and more detailed* data, Just and Pope (1999-b). Aggregate agricultural production data are now available in most of countries¹³. They can be disaggregated using this procedure. This is especially interesting as substantial site-specific data (soil-surveys, GIS data, satellite images) are becoming increasingly available. Disaggregation of economic data permits economic analysis at the most disaggregated level. It enables the combination of biophysical models, defined at this scale, with economic models. Moreover, the ME approach is

¹³ In the U.S., aggregate data may be found in the annual publication of the U.S. Department of Agriculture, *Agricultural Statistics*. There also are available in the *Census of Agriculture* published every 5 years. Most of country-level aggregate data are compiled by the Food and Administration Organization (FAO) and are easily available.

flexible enough to take into account out-of-sample information. Any specific out-of-sample information may be added to the disaggregation program via additional constraints. Any out-of-the sample information on transition probabilities may be added to the model via specification of priors.

Finally, as mentioned in the introduction of this paper, a valid disaggregation method is of interest in many other fields. The lack of high quality disaggregate data is a recurrent problem faced by many applied researchers.

5. References

Antle, J., and S. Capalbo. (2001). Econometric-Process Models for Integrated Assessment of Agricultural Production Systems. *American Journal of Agricultural Economics*, 83(2), May, 2001: 389-401.

Barker, T. and M.H. Pesaran. (1989). Disaggregation in Econometric Modeling – An Introduction. in *Disaggregation in Econometric Modeling*. T. Barker and M.H. Pesaran editors. Chapman and Hall, Inc. London.

Golan, A., G. G. Judge, and D. Miller (1996). Maximum Entropy Econometrics: Robust Estimation with Limited Data. John Wiley & Sons. 307 pages. New York.

Good, I.J. (1963). Maximum entropy for hypothesis formulation, especially for multidimensional contingency tables. *The Annals of Mathematical Statistics*, 34(3):911-934.

Heckeley, T., and W. Britz (2000). Concept and Explorative Application of an EU-wide Regional Agricultural Sector Model (CAPRI-Projekt). In: Heckeley, T., H.P. Witzke, and W. Henrichsmeyer (Eds.): *Agricultural Sector Modelling and Policy Information Systems*. Proceedings of the 65th EAAE Seminar, March 29-31, 2000 at Bonn University, Vauk Verlag Kiel, pp. 282-291.

Howitt, R.E., and A. Reynaud (2001). Dynamic land use - Disaggregation to field-level by maximum of entropy. Mimeo University of California at Davis. *To be presented at the 56th Econometric Society European Meeting, Lausanne 2001*.

Just, R., and R. Pope. (1999-a). Implications of Heterogeneity for Theory and Practice in Production Economics. *American Journal of Agricultural Economics*. 81(February): 711-718.

Just, R., and R. Pope. (1999-b). The agricultural producer: theory and statistical measurement. In *Handbook of Agricultural Economics*. B. Gardner and C. Rausser Eds.

- Kijima , M.** (1997). Markov processes for stochastic modeling. Chapman & Hall. London.
- Kullback, S.** (1959). Information Theory and Statistics. John Wiley and Sons. New York.
- Lee, T.C., G. G. Judge, and A. Zellner** (1970). Estimating the parameters of the Markov probability model from aggregate time series data. North-Holland Publishing Co., 254 pages. Amsterdam.
- Miller, D., and A. Plantinga.** (1999). Modeling land use decision with aggregate data. *American Journal of Agricultural Economics*, 81(February): 180-194.
- Plantinga, A.** (1996). The Effect of Agricultural Policies on Land Use and Environmental Quality. *American Journal of Agricultural Economics*, 78(November): 1033-1047.
- Stocker, T.M.,** ‘Empirical Approaches to the Problem of Aggregation over Individuals’, *Journal of Economic Literature*, 31(4), December, 1993: 1827-1874.
- United States Bureau of Reclamation (1997).** Central Valley Project Improvement Act. Technical Appendix, Volume 8. Sacramento.
- Wu, J., and K. Segerson.** ‘The Impact of policies and land Characteristics on Potential Groundwater Pollution in Wisconsin.’ *American Journal of Agricultural Economics*, 77(November 1995): 1033-1047.
- Wu, J., and R. Adams.** ‘Micro vs. Macro Acreage Response Models: Does Site-Specific Information Matter?’ mimeo. Oregon State University.

Appendixes

A. Data

Table A.1: Regional land use per year and per crop for CVPM 13

	88		89		90		91		92		93		94		95		96		97		98	
	Acres	%	Acres	%	Acres	%	Acres	%	Acres	%	Acres	%	Acres	%	Acres	%	Acres	%	Acres	%	Acres	%
A	60.2	17.5	64.3	19.7	59.7	18.2	63.5	18.6	61	17.9	63.1	18.8	65.8	19.6	67.74	20.9	69.2	21.0	71.4	21.7	69.9	22.0
C	73.3	21.3	61.8	19.0	68	20.7	79.6	23.3	76.4	22.5	80.4	23.9	75.5	22.5	71.37	22.0	68.1	20.6	59.9	18.2	48.6	15.3
F	30.7	8.9	25.8	7.9	22	6.7	16.2	4.7	15	4.4	15.3	4.5	17.1	5.1	15.57	4.8	15.8	4.8	17.3	5.2	19.4	6.1
G	96.5	28.1	89.2	27.4	90.7	27.6	92	27.0	101.9	30.0	97.6	29.0	99.3	29.6	86.26	26.6	93.8	28.4	101.4	30.8	103	32.4
P	61.9	18.0	61.4	18.9	60.5	18.4	57.7	16.9	53.2	15.6	47.6	14.2	44.9	13.4	46.74	14.4	44.9	13.6	41.9	12.7	40.5	12.7
T	6.4	1.9	7.1	2.2	7.8	2.4	9.1	2.7	8.8	2.6	11.4	3.4	11.9	3.6	15.01	4.6	15.5	4.7	15	4.6	15.4	4.8
V	8.61	2.5	8.3	2.5	11.8	3.6	15.3	4.5	15.7	4.6	12.4	3.7	12.5	3.7	15.06	4.6	16.3	4.9	16	4.9	15.4	4.8
S	6.1	1.8	7.8	2.4	8.2	2.5	7.7	2.3	8.2	2.4	8.5	2.5	8.2	2.4	6.51	2.0	6.5	2.0	6.7	2.0	5.6	1.8
Tot.	343.7		325.7		328.7		341.1		340.2		336.3		335.2		324.3		330.1		329.6		317.8	

Table A.2: Land use shares per year and crop at the DAU level

	88	89	90	91	92	93	94	95	96	97	98
<u>Merced</u>											
A	6.50	6.70	5.81	7.88	7.92	8.69	8.62	7.78	7.72	7.67	7.62
C	1.20	1.00	1.07	5.31	5.30	6.02	8.50	6.87	7.00	6.60	5.49
F	5.60	4.50	4.48	2.38	2.38	3.42	2.92	2.59	2.11	3.21	3.68
G	21.80	20.80	21.20	17.31	18.27	18.07	19.78	18.27	18.69	20.00	23.00
P	21.90	21.90	21.90	18.00	16.02	14.47	13.82	13.51	13.51	13.52	13.18
T	2.50	3.10	3.41	4.00	4.02	6.52	6.09	4.82	4.82	4.59	4.72
V	4.70	4.50	5.11	7.50	6.88	4.78	3.61	6.27	6.39	7.11	6.78
S	0	0.10	0.13	0.13	0.12	0.12	0.13	0.12	0.12	0.19	0.13
<u>Merced Stream Group</u>											
A	1.10	1.10	0.90	1.10	1.00	1.00	1.00	0.30	0.30	0.30	0.30
C	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10	0.10
F	0.30	0.20	0.20	0.30	0.20	0.10	0.10	0	0	0	0
G	4.30	3.80	4.00	3.90	4.21	5.00	6.41	2.90	2.70	3.20	4.00
P	3.40	3.40	3.40	3.90	3.70	3.20	3.20	3.40	3.40	3.40	3.30
T	0	0	0	0	0	0		0.80	0.80	0.80	0.80
V	0.30	0.30	0.30	0.40	0.40	0.30	0.30	0.50	0.50	0.50	0.50
S	0	0	0	0	0	0	0	0.10	0.10	0.20	0.10
<u>El Nido-Stevinson</u>											
A	23.30	24.00	20.03	21.99	21.09	21.39	21.34	25.25	25.09	25.11	24.74
C	21.00	18.20	20.79	21.99	22.94	22.95	22.85	20.06	20.14	19.17	15.96
F	12.40	11.00	9.59	9.63	9.02	6.68	6.73	7.15	5.76	8.98	10.14
G	28.00	26.20	28.54	25.38	28.05	32.75	37.00	31.59	32.23	36.15	41.32
P	22.10	22.10	22.08	22.97	20.54	20.50	20.53	18.91	18.88	18.92	18.56
T	3.20	4.00	4.42	5.03	4.67	4.68	5.45	8.65	8.75	8.73	8.91
V	2.00	2.00	2.26	2.41	2.50	2.34	2.20	3.81	4.14	4.25	4.21
S	0	0	0	0	0	0	0	0	0	0	0
<u>Madera-Chowchilla</u>											
A	8.00	8.80	9.02	8.77	8.42	10.01	11.40	9.54	9.89	10.60	10.30
C	19.40	19.00	21.02	23.93	21.51	19.92	16.81	17.06	16.02	13.09	10.42
F	4.40	3.30	1.21	2.02	2.28	3.69	4.02	1.27	1.90	1.08	1.19
G	15.70	11.50	12.00	12.48	17.18	12.02	8.72	13.49	15.49	18.60	14.30
P	5.90	5.70	5.29	5.01	4.72	4.69	4.60	4.37	3.99	2.40	2.22
T	0	0	0	0.11	0.11	0.21	0.38	0.05	0.10	0.09	0.12

V	0.40	0.40	1.11	1.09	1.02	0.90	0.72	0.47	0.49	0.38	0.40
S	0.60	0.50	0.81	1.09	1.71	1.32	1.29	0.80	0.78	0.80	0.71

Adobe - Valley Eastside

A	2.30	2.50	2.50	2.50	2.39	1.99	2.10	1.23	1.31	1.40	1.30
C	3.10	2.00	1.99	2.39	2.31	2.39	2.19	1.35	0.90	1.00	0.80
F	2.60	2.30	1.89	0.60	0.39	0.50	0.79	0.62	0.70	0.40	0.41
G	8.20	7.60	7.71	10.80	11.20	8.70	8.20	3.40	4.30	5.01	5.00
P	2.10	2.00	1.89	1.90	2.00	1.09	0.49	1.07	0.99	0.60	0.50
T	0.70	0	0	0	0	0	0	0.18	0.20	0.20	0.20
V	0.81	0.70	1.89	2.50	3.11	2.61	2.70	0.16	0.30	0.20	0.20
S	5.50	7.20	7.31	6.50	6.39	6.40	6.80	5.51	5.50	5.51	4.70

Gravelly Ford

A	19.00	21.20	21.51	21.16	20.22	20.02	21.32	23.71	24.92	26.27	25.73
C	28.50	21.50	22.97	25.91	24.26	28.99	25.06	25.87	23.91	19.93	15.91
F	5.40	4.50	4.62	1.32	0.69	0.92	2.57	4.01	5.29	3.62	3.93
G	18.50	19.30	17.32	22.10	22.97	21.02	19.18	16.58	20.39	18.50	15.37
P	6.50	6.30	5.87	5.91	6.33	3.61	2.28	5.53	4.11	3.10	2.78
T	0	0	0	0	0	0	0	0.56	0.84	0.60	0.68
V	0.40	0.40	1.10	1.40	1.83	1.53	3.01	3.85	4.53	3.47	3.32
S	0	0	0	0	0	0	0	0	0	0	0

B Disaggregation results

Table B.1: Simulated land use shares per crop $\hat{y}_k^j(t)$ versus observed at the DAU-level (in %)

	90 ^(a)		91 ^(a)		92 ^(a)		93 ^(a)		94 ^(a)		95 ^(b)		96 ^(b)		97 ^(b)		98 ^(b)	
	Obs	Pre	Obs	Pre	Obs	Pre	Obs	Pre	Obs	Pre	Obs	Pre	Obs	Pre	Obs	Pre	Obs	Pre
<u>Merced</u>																		
A	9.2	9.9	12.6	10.3	13.0	9.9	14.0	10.8	13.6	11.4	12.9	13.5	12.8	13.4	12.2	14.1	11.8	14.4
C	1.7	1.6	8.5	1.9	8.7	1.8	9.7	2.1	13.4	2.1	11.4	12.4	11.6	11.7	10.5	10.9	8.5	9.6
F	7.1	6.3	3.8	4.8	3.9	4.5	5.5	4.9	4.6	5.5	4.3	4.3	3.5	4.3	5.1	4.6	5.7	5.2
G	33.6	33.0	27.7	32.7	30.0	35.4	29.1	36.2	31.2	36.9	30.3	28.4	31.0	29.7	31.8	31.5	35.6	32.4
P	34.7	34.0	28.8	32.3	26.3	30.1	23.3	28.4	21.8	26.7	22.4	22.6	22.4	21.8	21.5	20.3	20.4	19.8
T	5.4	5.2	6.4	5.9	6.6	5.7	10.5	7.0	9.6	7.1	8.0	11.4	8.0	11.6	7.3	11.0	7.3	11.1
V	8.1	9.9	12.0	12.2	11.3	12.5	7.7	10.5	5.7	10.4	10.4	7.3	10.6	7.6	11.3	7.5	10.5	7.5
S	0.2	0.0	0.2	0.0	0.2	0.0	0.2	0.0	0.2	0.0	0.2	0.0	0.2	0.0	0.3	0.0	0.2	0.0
<u>Merced Stream Group</u>																		
A	10.1	12.3	11.3	12.5	10.4	12.4	10.3	12.6	9.0	12.7	3.7	8.9	3.8	8.9	3.5	9.0	3.3	9.0
C	1.1	1.1	1.0	1.1	1.0	1.1	1.0	1.2	0.9	1.2	1.2	0.8	1.3	0.8	1.2	0.8	1.1	0.8
F	2.2	2.0	3.1	2.0	2.1	2.0	1.0	2.0	0.9	2.1	0.0	0.8	0.0	0.8	0.0	0.8	0.0	0.8
G	44.9	42.8	40.2	42.8	43.8	43.3	51.5	43.4	57.7	43.5	35.8	57.6	34.2	57.7	37.6	58.0	44.0	58.0
P	38.2	38.2	40.2	37.9	38.5	37.4	33.0	37.1	28.8	36.7	42.0	29.3	43.0	29.1	40.0	28.7	36.3	28.5
T	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	9.9	0.0	10.1	0.0	9.4	0.0	8.8	0.0
V	3.4	3.6	4.1	3.7	4.2	3.8	3.1	3.7	2.7	3.8	6.2	2.7	6.3	2.7	5.9	2.7	5.5	2.8
S	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.2	0.0	1.3	0.0	2.4	0.0	1.1	0.0
<u>El Nido-Stevinson</u>																		
A	18.6	20.0	20.1	20.7	19.4	19.8	19.2	21.1	18.4	22.5	21.9	19.1	21.8	18.9	20.7	20.4	20.0	21.1
C	19.3	19.3	20.1	22.6	21.1	21.2	20.6	23.9	19.7	22.7	17.4	17.9	17.5	16.0	15.8	13.7	12.9	10.7
F	8.9	8.2	8.8	5.0	8.3	4.4	6.0	4.7	5.8	5.5	6.2	5.3	5.0	5.2	7.4	5.9	8.2	7.3
G	26.5	24.9	23.2	24.0	25.8	28.4	29.4	27.1	31.9	27.7	27.4	28.0	28.0	30.8	29.8	34.3	33.4	36.0
P	20.5	20.2	21.0	18.1	18.9	16.4	18.4	13.6	17.7	12.0	16.4	20.0	16.4	18.8	15.6	16.2	15.0	15.4
T	4.1	4.2	4.6	5.0	4.3	4.9	4.2	6.3	4.7	6.4	7.5	6.7	7.6	7.0	7.2	6.3	7.2	6.4
V	2.1	3.3	2.2	4.6	2.3	4.9	2.1	3.3	1.9	3.2	3.3	3.0	3.6	3.3	3.5	3.2	3.4	3.1
S	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
<u>Madera-Chowchilla</u>																		
A	17.9	16.9	16.1	17.0	14.8	16.7	19.0	16.8	23.8	17.4	20.3	24.0	20.3	24.3	22.5	25.0	26.0	25.7
C	41.7	40.6	43.9	43.4	37.8	42.1	37.8	43.6	35.1	42.7	36.3	33.6	32.9	32.5	27.8	30.7	26.3	28.8
F	2.4	5.9	3.7	4.6	4.0	4.4	7.0	4.4	8.4	4.8	2.7	8.1	3.9	8.1	2.3	8.6	3.0	9.4
G	23.8	23.4	22.9	22.7	30.2	24.9	22.8	23.9	18.2	24.2	28.7	17.2	31.8	18.2	39.5	19.0	36.1	19.7
P	10.5	11.3	9.2	10.6	8.3	10.1	8.9	9.1	9.6	8.7	9.3	10.1	8.2	10.0	5.1	9.5	5.6	9.5
T	0.0	0.0	0.2	0.0	0.2	0.0	0.4	0.0	0.8	0.0	0.1	0.8	0.2	0.8	0.2	0.8	0.3	0.9
V	2.2	0.9	2.0	1.0	1.8	1.1	1.7	0.9	1.5	0.9	1.0	1.8	1.0	1.9	0.8	2.0	1.0	2.0
S	1.6	1.0	2.0	0.7	3.0	0.8	2.5	1.3	2.7	1.3	1.7	4.5	1.6	4.1	1.7	4.5	1.8	4.0
<u>Adobe - Valley Eastside</u>																		
A	9.9	9.9	9.2	10.7	8.6	10.3	8.4	9.8	9.0	10.0	9.1	8.7	9.2	8.7	9.8	8.7	9.9	9.0
C	7.9	8.4	8.8	9.3	8.3	9.0	10.1	8.6	9.4	8.6	10.0	8.9	6.3	8.9	7.0	8.6	6.1	8.6
F	7.5	8.9	2.2	8.4	1.4	8.0	2.1	7.6	3.4	7.9	4.6	3.2	4.9	3.2	2.8	3.3	3.1	3.5
G	30.6	31.0	39.7	32.6	40.3	33.1	36.7	30.5	35.2	30.6	25.1	33.0	30.3	33.6	35.0	33.6	38.1	34.2
P	7.5	8.1	7.0	8.4	7.2	8.0	4.6	7.2	2.1	7.1	7.9	1.9	7.0	2.0	4.2	2.0	3.8	2.0
T	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.3	0.0	1.4	0.0	1.4	0.0	1.5	0.0
V	7.5	3.1	9.2	3.7	11.2	3.7	11.0	3.2	11.6	3.2	1.2	11.8	2.1	12.0	1.4	11.8	1.5	12.0
S	29.0	30.5	23.9	26.9	23.0	28.0	27.0	33.0	29.2	32.6	40.8	32.5	38.7	31.6	38.5	32.0	35.9	30.8
<u>Gravelly Ford</u>																		
A	29.3	26.9	27.2	27.0	26.5	26.1	26.3	26.7	29.0	27.8	29.6	30.5	29.7	30.5	34.8	31.8	38.0	33.0

C	31.3	32.1	33.3	35.2	31.8	33.6	38.1	35.7	34.1	34.4	32.3	32.7	28.5	30.4	26.4	27.4	23.5	24.4
F	6.3	5.2	1.7	3.6	0.9	3.4	1.2	3.5	3.5	3.9	5.0	3.2	6.3	3.2	4.8	3.5	5.8	4.0
G	23.6	26.7	28.4	25.6	30.1	28.7	27.6	27.2	26.1	27.5	20.7	24.5	24.3	26.5	24.5	28.2	22.7	29.5
P	8.0	8.4	7.6	7.7	8.3	7.2	4.7	6.2	3.1	5.7	6.9	3.4	4.9	3.3	4.1	3.0	4.1	3.0
T	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.7	0.0	1.0	0.0	0.8	0.0	1.0	0.0
V	1.5	0.7	1.8	0.9	2.4	0.9	2.0	0.7	4.1	0.7	4.8	5.7	5.4	6.1	4.6	6.0	4.9	6.0
S	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0

Notes: Obs and Pre respectively give the observed and the predicted land use shares.
 (a) for in-sample estimates (b) for out-sample estimates