# Predictive Modeling of Key Process Variables in Granulation Processes based on Dynamic Partial Least Squares

**D. Ronen\*, C.F.W. Sanders\*, H.S. Tan\*\*, P. R. Mort\*\*\*, F.J. Doyle III\***

*\*Department of Chem. Eng., UCSB, Santa Barbara, CA 93106-5080 USA (e-mail: frank.doyle@isb.ucsb.edu)*
*\*\* P & G Technical Centre, Ltd. Whitley road, Longbenton, Newcastle Upon Tyne, NE12 9TS UK (e-mail: tan.h.9@pg.com)*
*\*\*\*Procter & Gamble Co. 5299 Spring Grove Ave, Cincinnati, OH 45217 USA (e-mail:mort.pr@pg.com)*

**Abstract:** Granulation is a multivariable process characterized by several physical attributes that are essential for product performance, such as granule size and size distribution. An optimally operated granulation process will yield, in a reproducible manner, product with tightly controlled performance attributes. In this paper predictive models of the dynamics of these key variables are developed using a dynamic partial least squares approach. The method, demonstrated here on process simulation as well as on an industrial mixer-granulator process, result in accurate predictions. These models motivate the development of model predictive controllers for these processes.

*Keywords:* Granulation, Process control, Dynamic modeling, Partial least squares

## 1. INTRODUCTION

Granulation is a complex process in which many input variables influence many product properties. As Iveson et al. describe in a review paper (2001), the understanding of the fundamental processes that control granulation behavior and product properties have increased in recent years. This knowledge can be used during process design, in choosing the right formulation and operating conditions, and it can also be used to improve process control. Although many variables are set constant during process design, variations during production in input variables occur due to the variable nature of the powder feed. Even if all granule properties, except for size, are ignored for process control, a one dimensional granule size distribution can be constructed by multiple discrete output variables, in order to represent the shape of the distribution (these can be mean sizes (with coefficients of variation), percentile sizes, moments or size bins). Model Predictive Control (MPC) is an effective method to control such multiple input, multiple output processes (García, et al., 1989). The majority of MPC applications in the chemical process industries utilize empirical models that are constructed from plant data. In this work, we explore the use of dynamic partial least squares to construct these empirical models.

## 2. METHODS

### 2.1 Partial Least Squares

Partial Least Squares (PLS) methods have been demonstrated as a useful tool for analysis of data and modeling of the systems from which the data are collected (Kaspar and Ray, 1993). Unlike related methods, such as Principal Component Analysis (PCA), which finds factors that capture the greatest amount of variance in the predictor (*X*) only, the PLS method attempts to find factors which both capture variance and achieve correlation. PLS handles this by projecting the information in high dimensional spaces (*X*,*Y*) down to low dimensional spaces defined by a small number of latent vectors ($t_1,t_2...t_a$). These new latent vectors summarize all the important information contained in the original data sets, by representing the scaled and mean-centered values of *X* and *Y* matrices as:

$$X = \sum_{i=1}^{a} t_i p_i^T + E$$

$$Y = \sum_{i=1}^{a} u_i q_i^T + F$$

$$(1)$$

where the $t_i$ are latent (score) vectors calculated sequentially for each dimension i=1,2,…a.

In the PLS method, the covariance between the linear combinations of *X* and the output measurement matrix *Y* is maximized at each iteration, using the vectors $p_i$ and $q_i$ which are the loading vectors whose elements express the contribution of each variable in *X* and *Y* toward defining the new latent vectors $t_i$ and $u_i$. *E* and *F* are residual matrices for *X* and *Y* blocks, respectively. The optimal number of latent vectors retained in the model is often determined by cross-validation (Dayal et al. 1994).

In an industrial environment, it is more often the case that many of the predictor variables (*X*) are highly correlated with one another and their covariance matrix is nearly singular,

which renders classical regression methods intractable. Reduced space methods such as PLS and PCA can overcome this problem (MacGregor and Kourti, 1995). PLS is also robust to measurement noise in the data and can be used in cases where there are random missing data and when the number of input variables is greater than the number of observations (Dayal et al. 1994). Various examples of the implementation of PLS analysis to industrial process modeling and control can be found in the literature (for example, Dayal et al., 1994, MacGregor and Kourti, 1995, and others).

Process dynamics can be incorporated into the PLS model by including columns of lagged outputs and/or inputs into the predictor block ($X$) (Dayal et al., 1994, Kaspar and Ray, 1993, Juricek et al. 2001). The resulting PLS model is actually an ARX type input-output model of the form:

$$A(q^{-1})y_i(k) = \sum_{j=1}^{nu} B_j(q^{-1})u_j(k - nk_j)$$

(2)

$$A(q) = 1 + a_1 q^{-1} + a_2 q^{-2} + ... + a_m q^{-m}$$

(3)

$$B_j(q) = b_{j,1}q^{-1} + b_{j,2}q^{-2} + ... + b_{j,m}q^{-m}$$

(4)

where $y$ denotes the output variable (e.g., median particle size, $d_{50}$), and $u$ denotes the manipulated variable (e.g., binder flow). The terms $A$ and $B$ contain the autoregressive and exogenous terms of the model, respectively. The autoregressive term captures dynamics through lagged terms of the output, and the exogenous term captures dynamics through lagged terms in the input.

Once the models have been calculated from the plant data, it is useful to evaluate their properties using several key statistical measures. Some of the useful statistics that are associated with reduced space models (Wise et al. 2006) are outlined below:

*Q residual* – is simply the sum of squares of each row of $E$ (from eq. 1), i.e. for the i[th] sample in $X$, $x_i$:

$$Q_i = e_i e_i^T$$

(5)

where $e_i$ is the i[th] row of $E$. The Q statistics is a measure of the difference between a sample and its projection into the $a$ principal components retained in the model.

*Hotelling $T^2$* is a measure of the variation in each sample within the model. Its value is the sum of normalized squared scores, defined by:

$$T_i^2 = t_i \lambda^{-1} t_i^T$$

(6)

where $t_i$ are the score vectors (eq. 1) and $\lambda$ is a diagonal matrix containing the eigenvalues corresponding to $a$ eigenvectors (principal components) retained in the model.

Together, the $T^2$ and Q residual statistics are useful in evaluating the fitness of a PLS model to specific data. It is possible to calculate statistically meaningful confidence limits for both cases.

### 2.2 Simulation studies

In our previous work, a nonlinear one dimensional population balance model (1D-PBM) was successfully used to model a laboratory continuous drum granulation process with fine particle recycle (Glaser et al., 2008). The same model is used here as a base for a process simulation (Figure 1) for a preliminary evaluation and sensitivity test of the applicability of the dynamic PLS modeling technique for granulation.
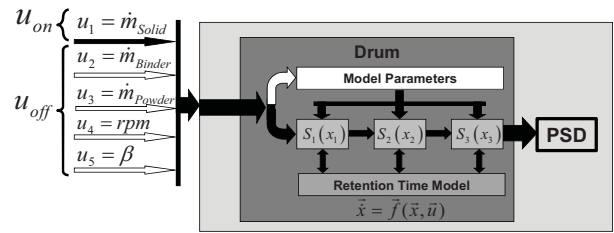


Fig. 1. Simulator structure: five inputs are included in the simulator: binder spray rate, fine powder feed-rate, drum rotation-rate and the drum inclination angle. The model is divided into three well mixed drum compartments, each described by an individual set of ODEs, a retention time model and a set of global parameters that influence the model behavior (taken from Glaser 2008).

Both particle median size ($d_{50}$) and, separately, particle size distribution width ($d_{84}/d_{16}$) were used as output variables for this study. The predictor ($X$) was constructed from 4 manipulated variables (solid feed flow rate, binder feed flow rate, drum rotation speed, recycle rate) and the computed recycle flow as an additional input variable. Process dynamics were incorporated into the $X$ block by including columns of lagged output variables. The lag time was estimated using an autocorrelation function. Delay times of each of the input variables were estimated using cross correlation function, and the predictor matrix was adjusted according to the obtained delay vector. During the simulation, the 4 manipulated variables were randomly perturbed around their nominal values at steady state sequentially, i.e. input variables were perturbed one after the other in fixed time gaps. The resulting PLS-based ARX model's short horizon predictive ability was tested by cross validation with a set of separately calculated simulation sequences with different excitation regimes. For each of these cross validation sequences, the root mean square error of the model based prediction (RMSEP), relative to the simulated plant measurements was calculated for a given short horizon period. In order to make a more representative quantification of the predicting ability of the model, the short horizon start point was moved along the time axis of the data one time step after another thus creating a set of RMSEP measures out of which an average and maximum RMSEP could be calculated. All variables were mean centered and scaled to unit variance prior to processing.

Based on this technique a sensitivity test was performed in order to estimate the required size of the data set needed for reliable process modeling. Figure 2 shows the convergence of RMSEP related to the length of data set used for the PLS model training. This plot is based on averaging 100 multiple simulations and modeling runs for each training length. Sample rate was set to 2 minutes and the simulated process step response time ($\tau$) was set to 4.5 minutes. The prediction horizon was set to 8 samples (i.e., 16 minutes). From this figure one can note that most of the dynamic features are captured by the model in the first 200 minutes of training data, as the mean RMSEP converges to low values. However using training data of up to 600 minutes would improve the model predictions. Notice that these results are not so sensitive to the excitations rate used in the modeling data set (i.e., time between two successive input variable perturbations), as long as this time is in the order of magnitude of the expected variations in process variables. Figure 3 depicts the response of the process model obtained with 520 minutes of training data to the process simulation for a step response in one of the input variables.

All of these models use one lagged output variable (granules median size) in the predictor block and 2 latent variables in the PLS model, which are linear combinations of the time-lagged values of the output variable and the delayed values of the five process variables. Figure 4 shows the prediction abilities of three PLS models obtained using different lengths of data sets for training (120, 220 and 520 minutes long) from a single simulation data, with input variable excitation every 15 minutes. In this example, the predictions of these models are cross-validated using data from a separate simulation run with randomly timed excitations of the inputs. Considering that PLS models captures covariance in X and Y, it is possible to calculated the percentage of variance captured by each of their latent variables by dividing the variance predicted by the latent variables to the total variance in the original data. The percentage of variance captured by the abovementioned 3 PLS models (from the training data) is detailed in Table 1. It is noticeable that the longer the training set used, more fine details of the process dynamics are captured by the models, confirming Figure 2 results. Notice that high values of explained variance do not guarantee good prediction of validation data by these models. The robustness of the PLS based models to measurement noise is demonstrated in Figure 5 and Table 2, where the simulated process was subject to 5% white noise on the output and input variable measurements.
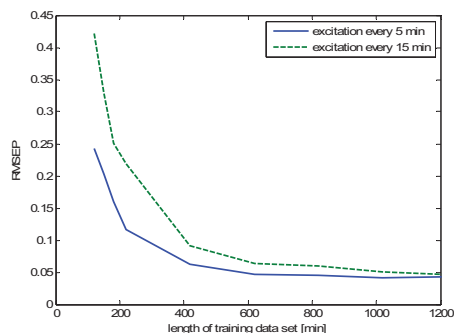


Fig. 2. Root Mean Square Error of Prediction (in validation simulation) versus length of training data set (based on modeling simulation), at different excitation rates.
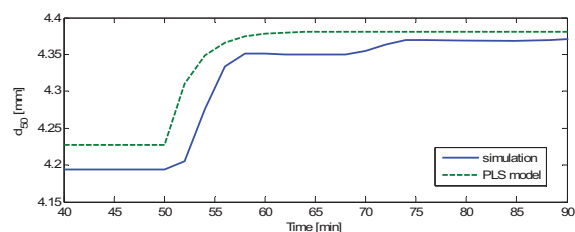


Fig. 3. Step response to a 1.2% step change in Binder feed flow - PLS model based on 520 minutes training data vs. process simulation.
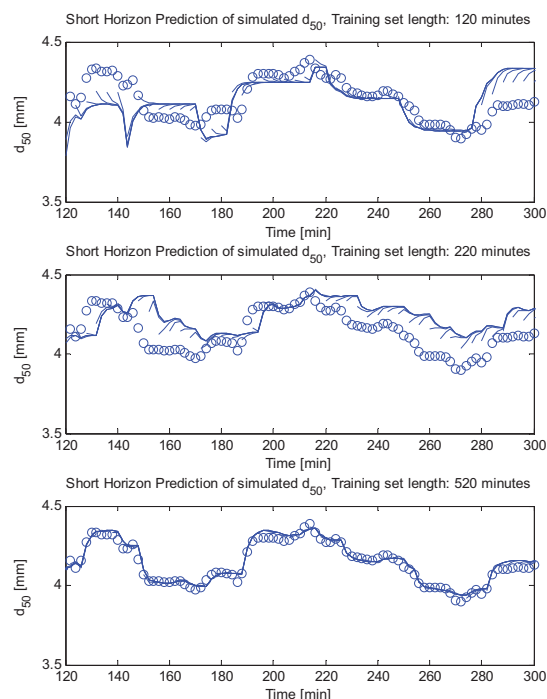


Fig. 4. PLS based dynamic model validation for different training set length. Circles represents simulation results, lines represent 8 point horizon prediction.

Table 1: Percent variance captured by PLS models based on different lengths of training data sets.

| Training set length | LVs | X Block | | Y Block | |
|---|---|---|---|---|---|
| | | This | Total | This | Total |
| 120min | 1 | 76.04 | 76.04 | 82.44 | 82.44 |
| | 2 | 7.66 | 83.7 | 11.42 | 93.85 |
| 220min | 1 | 59.45 | 59.45 | 93.05 | 93.05 |
| | 2 | 14.21 | 73.66 | 2.8 | 95.85 |
| 520min | 1 | 41.35 | 41.35 | 97.88 | 97.88 |
| | 2 | 17.29 | 58.64 | 1.04 | 98.92 |



Fig. 5. PLS based dynamic model validation, simulated process with 5% white measurement noise. Circles represents simulation results, lines represent 8 point horizon prediction.

Table 2: Percent variance captured by PLS model, simulated process with 5% white measurements noise.

| Training set length | LVs | X Block | | Y Block | |
|---|---|---|---|---|---|
| | | This | Total | This | Total |
| 520 min | 1 | 33.93 | 33.93 | 82.75 | 82.75 |
| | 2 | 33.10 | 67.03 | 1.61 | 84.37 |

## 3. DYNAMIC PLS MODELLING OF AN INDUSTRIAL PROCESS PLANT

### 3.1 Process plant description

In this section, we report on some preliminary studies on granulation model identification for a Procter & Gamble (P&G) industrial granulation process using normalized process data.

A complex industrial granulation process, such as the flowsheet shown in Fig. 6, was subjected to a series of (designed) random perturbations in a number of input parameters. This plant is equipped with an on-line granule size measurement system that measures particle size based on image analysis of 2-D camera images. The analysis constructed size distributions on the basis of the measured cross sectional area of the 2-D images. Granule size data along with all other plant variables were then sent to the UCSB team for modeling.

There are notable distinctions between the P&G study and the one reported in Section 2.2. The P&G process study is not meant to be used as a direct comparison (or validation) for the process simulation studies in Section 2.2. Rather, we present both as separate case studies to demonstrate the feasibility of using PLS methodology as an empirical modeling tool for granulation process control.

The low-shear drum-granulation pilot plant that was used to design the simulator in Section 2.2 produced particles with $d_{50}$'s of several mm; on the other hand, the medium-high shear mixer-granulation process shown in Fig 6 typically produced particles with $d_{50}$'s less than 1 mm. While the underlying physical mechanisms of growth and consolidation may be similar, the flow and shear fields are very different for the two processes (the drum granulator is relatively low shear, compared to the medium-high shear mixer-granulation process), the process layouts and control handles are different, as are the material properties. As such, the choices of process variables (manipulated and measured) are unique for each process.
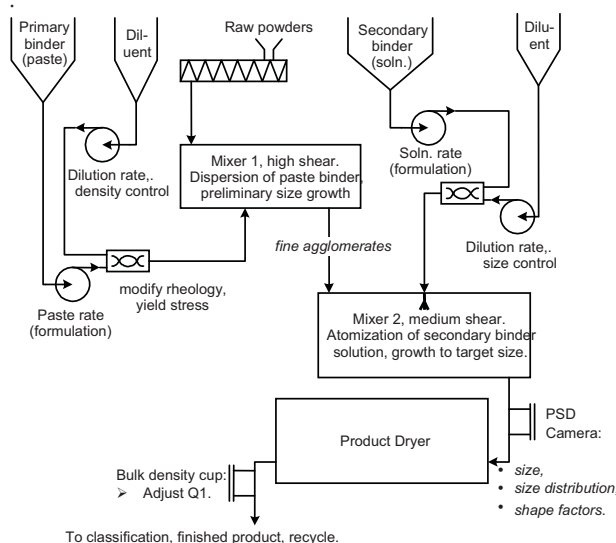


Figure 6. Representative P&G process flow diagram for mixer-granulator (Mort et al. 2001). For simplicity, this diagram omits the usual operations for classification and recycle.

### 3.2 PLS modeling of real plant data- Case I

Figure 7 describe the dynamic PLS model fitting obtained for the granules median size. The data set obtained from the plant originally contained 147 sampling points, each consist of 81 process variables, together with granules size measurements. Sampling time was 0.4 times the process characteristic time $\tau$. During this time period 4 manipulating variables were subjected to random perturbations around their nominal values at steady states (Fig. 8) in a similar way to the simulation work described earlier (adjusted to the process $\tau$), while other adjustments were continuously made to other plant variables (i.e. normal plant operations). Granules median size ($d_{50}$) was selected as the output variable. Nine out of the 81 process variables were chosen as predictor variables for the PLS model, based on engineering judgment, GA based variable selection (PLS Toolbox 5.0 by Eigenvector research incorporated), and trial and error. The

output lag time and process variables delay times were evaluated using the auto and cross correlation functions, respectively, as described in section 2.2. The process model uses two latent variables, which are linear combinations of the time-lagged values of the output variable and the delayed values of the nine process measurements. For an independent cross validation of the model, the above data was divided to two sections – the first half was used to train the PLS model, and the second used to test model predictions, yielding RMSEP value of 0.26. These results, as shown in Fig. 9. and Table 3, are very similar to the fitting obtained for the simulation data of the same training length to τ ratio (Fig. 2 and Fig 4. upper plot).
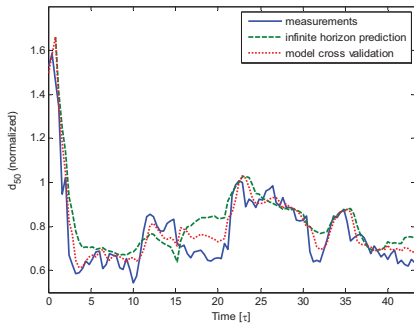


Fig. 7. Dynamic PLS model fitting to plant $d_{50}$ data.

Table 3: Percent variance captured by PLS model, plants' d50 data.

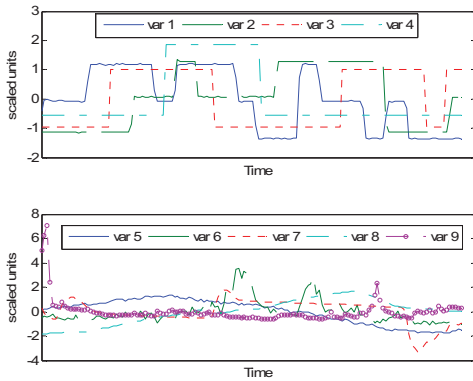| Training set length | LVs | X Block | | Y Block | |
|---|---|---|---|---|---|
| | | LVi | Total | LVi | Total |
| 60τ | 1 | 24.97 | 24.97 | 79.34 | 79.34 |
| | 2 | 25.09 | 50.06 | 9.04 | 88.38 |
| 36τ | 1 | 39.46 | 39.46 | 72.13 | 72.13 |
| | 2 | 23.31 | 62.77 | 16.8 | 88.92 |



Fig. 8. Values of the 9 predictor variables used in Case I PLS model – 4 manipulated variables (top) and 5 additional process variables (bottom)
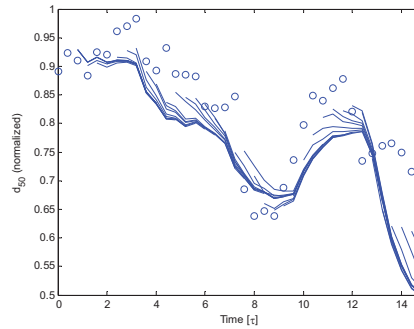


Fig. 9. Cross validation of the dynamic PLS model for plant $d_{50}$ data. Circles represents measurements, lines represent 8 point horizon prediction.

### 3.3 PLS modeling of real plant data- Case II

In a separate test, the granules distribution width as a function of selected process variables was modeled. This analysis was performed on two limited sets of data, each from a different operating day. A series of step tests were performed on one of the manipulating variables. As in the previous case, other adjustments were continuously made to other plant variables to maintain normal plant operation. The standard deviation of the granules measured area was used as the output variable to be modeled. A dynamic PLS model was built using 3 input variables and one lagged output variable, based on the first data set, and then validated using the second data set, resulting in an excellent fit (Figure 10). In this case, as well, the process model uses two latent variables.
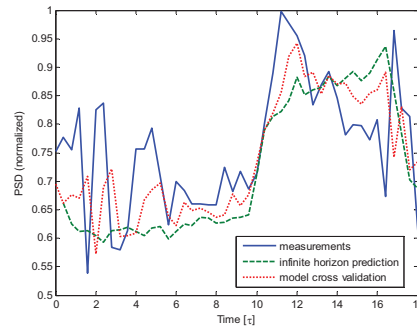


Fig. 10. Actual process data (granules area standard deviation): Cross validation of model based on first data set, tested on second data set

On the scores plot (Figure 11) it is clear, however, that these two sets represent different and distinct operating conditions. If one further examines the contribution of each variable for these two sets (Figures 12) we can see that the main difference is that on the validation set, much lower values of variable 2 were used, compared to the modeling set. It is also interesting to note that the only outlier of the modeling set also exhibits the same low value on variable 2.
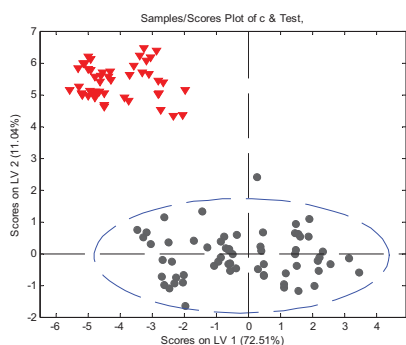
Fig. 11. Score Plots for the PSD width model. Circles are samples from the modeling set; Triangles are samples from the validation set. The ellipse marks the 95% confidence limit for the model.
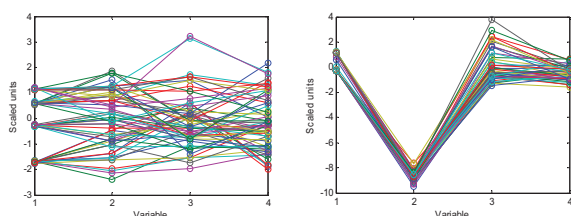


Figure 12: Values of input variables used in the modeling set (left) and in the validation set (right).

Although these results looks promising with respect to the ability to analyze and predict the granulation process variables, a quick look at the high values of Q residuals and Hotelling $T^2$ (Figure 13) indicates that this model is far from describing the whole complexity of the process, and many more measurements should be done in order to characterize the different operating regimes of this process.
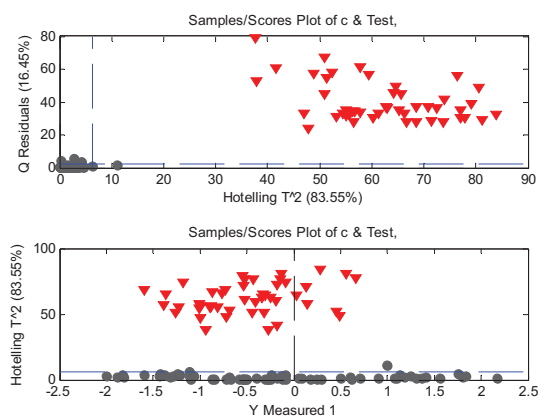


Figure 13: Q Residuals and Hotelling $T^2$ Values for the PSD width model

## 4. CONCLUSIONS

Dynamic PLS modeling was proven to be an effective tool in modeling key process variables in an industrial granulation process. Our future work will explore methods to capture the additional dynamics that remain in the plant data. We are also planning longer plant runs with larger input variable excitation to improve the model identification. Longer term goals are to develop a model-based controller for plant testing.

## REFERENCES

Dayal, B.S., MacGregor, J.F., Taylor, P.A., Kildaw, R. and Marcikic, S. (1994). Application of Feedforward Neural Networks and Partial Least Squares Regression for Modeling Kappa Number in a Continuous Kamyr Digester, *Pulp and Paper Canada*, 95, 26-32.

García, C.E., Prett, D.M. and Morari, M. (1989), Model Predictive Control: Theory and Practice—a Survey, *Automatica*, 25, 335–348.

Glaser, T., Sanders, C.F.W., Wang, F.Y., Cameron, I.T., Litster, J.D., Poon, J.M.H., Ramachandran, R., Immanuel, C.D. and Doyle, F.J. (2008), Model Predictive Control of Continuous Drum Granulation, *Journal of Process Control*, in press.

Iveson, S.M., Litster, J.D., Hapgood, K. and Ennis, B.J. (2001), Nucleation, Growth and Breakage Phenomena in Agitated Wet Granulation Processes: A Review, *Powder Technology*, 117, 3-39.

Juricek, B.C., Seborg, D.E., and Larimore, W. E. (2001), Identification of the Tennessee Eastman Challenge Process with Subspace Methods, *Control Engineering Practice*, 9, 1337-1351.

Kasper, M.H. and Ray, W.H. (1993), Dynamic PLS Modelling for Process Control, *Chemical Engineering Science*, 20, 3447-3461.

MacGregor, J.F. and Kourti, T. (1995), Statistical Process Control of Multivariate Processes, *Control Engineering Practice*, 3, 403-414.

Mort, P.R., Capeci, S.W. and Holder J.W. (2001), Control of Agglomerate Attributes in a Continuous Binder-Agglomeration Process, *Powder Technology*, 117, 173-176.

Wise, B.M., Gallagher, N.B., Bro, R., Shaver, J.M., Windig W. and Koch, R.S. (2006), *PLS_Toolbox 4.0*, Eigenvector Research Incorporated, Wenatchee WA.