

## The PM<sub>10-2.5</sub> DQO Model

This is a companion document to the report titled: *Sensitivity of the PM<sub>10-2.5</sub> Data Quality Objective to Spatially Related Uncertainties* provided to the Clean Air Scientific Committee (CASAC) for the meeting scheduled for September 21-22, 2005.

### 1.0 THE SIMULATION MODEL

The simulation model used for establishing the PM<sub>10-2.5</sub> DQOs and equivalency method requirements is fairly complex. It was developed as a generalization of the model used for the PM<sub>2.5</sub> DQO simulation model. The basic model was first implemented in MathCad and then reprogrammed in C<sup>++</sup> with corresponding functions to increase the speed of execution. The C<sup>++</sup> functions have been successfully tested against the MathCad functions. Consequently, it suffices to describe the MathCad functions that are much easier to read. The MathCad version is described below.

### 2.0 CONCEPTUAL MODELS

There are several layers to the conceptual model. The PM<sub>10-2.5</sub> and PM<sub>2.5</sub> components are viewed separately. Since the PM<sub>2.5</sub> data generation is simpler and motivates the rest, it is described first. Within the code, however, it is easier to generate it last.

#### 2.1 The PM<sub>2.5</sub> Conceptual Model

The PM<sub>2.5</sub> component is thought of as a spatially uniform (within the scale of the monitor) process with a sinusoidal mean having one period per year. The basic sine curve has its modulus and vertical shift related so that models vary multiplicatively. This is a very important point for the implementation, as all the parameters have a multiplicative interpretation. This means that a site is associated with a fixed ratio between the highest point on the sine curve and its lowest point. Consequently, the sine wave is parameterized as:

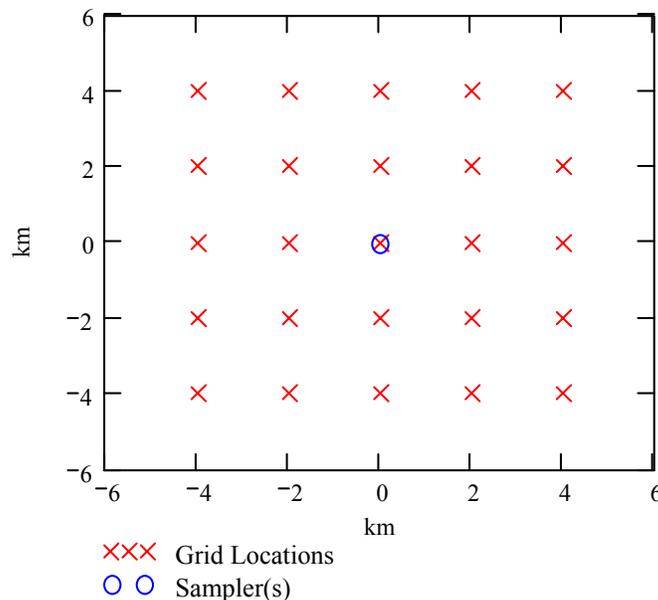
$$A \cdot \left( 1 + \frac{R-1}{R+1} \sin \left( \frac{2\pi \cdot \text{day}}{365} \right) \right)$$

where  $A$  is a multiplicative factor,  $R$  is the seasonal ratio, and  $\text{day}$  is the number of days into the simulation period (day ranges from 1 to  $3 \times 365$ ). Since three full years are modeled, there is no need for a phase shift in the PM<sub>2.5</sub> component.

From this mean the true day-to-day concentrations are generated as lognormally distributed random deviations from the mean with a constant coefficient of variation (CV). These deviations are allowed to have autocorrelation that is modeled as the exponential of a simple autoregressive process.

## 2.2 The PM<sub>10-2.5</sub> Conceptual Model

The coarse portion of the PM is simulated as a random field that varies with correlation both spatially and temporally. The spatial domain is simulated on a grid that is supposed to represent a neighborhood scale monitor with grid points located every 2 km in a square grid that is 8 km on a side. The temporal variation also includes a seasonal pattern throughout the three-year simulation periods. For each grid point, separate three-year design values are computed. Then the grid's design values (one for the annual standard and one for the daily standard) are computed by averaging over the grid-point specific design values. Finally, this design value is compared to the design value for the simulated observations corresponding to the center point. This whole process is repeated 1,000 times to determine the probability that the observed-design value correctly predicts whether or not the grid's true-design value is above the standard.



**Figure A-1. Simulation grid points.**

The temporal pattern for the coarse fraction is also based on random deviations from a sinusoidal pattern. The pattern can have either one or two periods per year and can be phase shifted from the fine PM portion. The random deviations can have autocorrelation and be correlated with the random deviations of the fine fraction.

## 3.0 PROGRAMMING DETAILS

The details are made somewhat complex by the need to pre-specify the various correlations, variances, and other model parameters. To achieve the desired results, the simulation is done with a series of functions/modules with separate tasks. Each of these is described below.

Notes:

1. The indexing always starts with 0.
2. The parameter vector is referred to as *parms*.
3. Within the code,  $A^{<n>}$  refers to the  $n^{\text{th}}$  column of the matrix  $A$ .
4. The code shown is the MathCad code that was used to develop the model. The C<sup>++</sup> code is tested against the MathCad code. To accommodate that testing, the format of the final output has been modified from what is currently shown.

### 3.1 Scaled Design Value Distribution

Input: The vector *parms* of parameter values.

```
Gen_scaled_dist (parms) := | Seed(1233547)
                           | bchol ← BChol(parms)
                           | for k ∈ 0..reps - 1
                           |   | in ← Gen_truth(bchol, parms)
                           |   | X ← submatrix(in, 0, rows(in) - 2, 0, cols(in) - 1)
                           |   | tm ← inrows(in)-1, 0
                           |   | Y ← Gen_obs(X, parms)
                           |   | Z<k> ← Cal_pobs(Y, parms, tm)
                           | ZT
```

This function forms the main simulation loop. For each repetition of the loop, a three-year set of scaled true concentrations are generated for each point in the grid. The output of the *Gen\_truth*, however, only contains the simulation values for the center grid point and a summary appended to the end. The next steps in the code separate the simulation values and summary for use later in the code. The fourth step in the loop generates the observations corresponding to the scaled true concentrations. The final step in the loop summarizes the observed-design values for this three-year simulation. The entire loop is repeated as often as is practical and is controlled by the value of the global variable *reps*.

The output of the function is a collection of scaled design values that all correspond to having the mean over the grid points of the design values of the true concentrations equal to 1. The scaling is such that multiplying the simulated observed-design values by 20 results in a set of simulated observed-design values. These values all correspond to examples with simulated true concentrations such that the mean of the true-design values over the grid points is exactly 20. This is true for both the daily standard design values and the annual standard design values. (Note that 20 is *not* the value of the long-term mean for the hypothetical site.)

### 3.2 The Cholesky Decomposition Matrix

Input: *parms*, the vector of simulation parameters.

```
BChol(parms) := if parms22 > 0
  for i ∈ 0..4
    for j ∈ 0..4
      for s ∈ 0..4
        for t ∈ 0..4
          x5·i+j, 5·s+t ← parms22 · exp  $\left[ - \left( \frac{|i-s| + |j-t|}{parms_{23}} \right) \cdot \frac{2}{parms_{23}} \right]$ 
        out ← cholesky(x)
      out24,24 ← 0 if parms22 = 0
    out
```

#### Note:

The matrix multiplication of the Cholesky decomposition of a matrix **A** and a vector/matrix of values independently drawn from a standard normal distribution,  $N(1,0)$ , results in a vector/matrix with the elements that have a variance-covariance matrix equal to **A**. We use this to create the user-specified spatial correlation structure. Specifically, this function creates the Cholesky decomposition of the spatial correlation matrix, except when the sill is set to 0. When the sill is set to zero, the matrix is all zeros. In either case, a separate temporal pattern is added to reflect the fact that the temporal variance is larger than the spatial variance. An exponential spatial correlation is assumed. The factor of 2 in the exponential is to reflect the fact that the grid points are 2 km apart.

### 3.3 The Coarse Mass Deviations Generator

Input: The matrix *bchol* from the function BChol and *parms*, the vector of simulation parameters.

```

Gen_CM_deviations (bchol ,parms) :=
  tmp<0> ← rnorm(26, 0, 1)
  tmp<0> ← tmp<0> - mean(tmp)
  for i ∈ 1.. 365·3 + 10
    tmp<i> ← rnorm(26, 0, 1)·√(1 - (parms13)2) + tmp<i-1>·parms13
  tmp2 ← bchol·submatrix(tmp, 0, 24, 11, 3·365 + 10)
  for i ∈ 0.. 365·3 - 1
    tmp2<i> ← tmp2<i> + √((parms16)2 - parms22)·tmp25, i+10
  tmp2T

```

#### Notes:

This function creates the (multiplicative) temporal deviations of the long-term trend for the PM<sub>10-2.5</sub> component. The first loop creates an autocorrelated temporal sequence for all of the grid points and an extra sequence. The first 25 sequences are then given the desired spatial correlation structure by multiplying by the Cholesky decomposition. However, the spatial variation can be less than temporal variance. This is fixed by adding a spatially uniform temporal sequence to all of the points in the grid based on the “extra” sequence generated in the first loop. Since everything has the same temporal autocorrelation, then the end result has the A1 structure ( $X_{n+1} = \alpha X_n + Y_{n+1}$ , where the  $Y$ s are independent with an appropriate variance).

Also, note that there is a very short “burn in.” The temporal sequences are generated for a slightly longer period than needed by the rest of the program. This is needed to adjust for the fact that the A1 structure does not hold initially.

### 3.4 The Fine Mass Deviations Generator

Input: The matrix  $CM\_dev$  from the function  $Gen\_CM\_deviations$ , and  $parms$ , the vector of simulation parameters.

```

Gen_FM_deviations (CM_dev , parms) :=
  x ← parms14 ·  $\frac{parms_{15} \cdot (1 - parms_{12} \cdot parms_{13})}{parms_{16} \cdot [1 - (parms_{13})^2 \cdot (parms_{14})^2]}$ 
  a ← (parms14)2
  b ← (parms12)2
  c ← (parms13)2 · (parms14)2
  d ← 2 · (parms14)2 · parms13 · parms12
  y ←  $\sqrt{\frac{1 - a - b - c + d}{1 - (parms_{13})^2 \cdot (parms_{14})^2}}$ 
  k ←  $\frac{parms_{12} - (parms_{14})^2 \cdot parms_{13}}{1 - (parms_{13})^2 \cdot (parms_{14})^2}$ 
  tmp ← rnorm(3·365, 0, parms15)
  tmp0 ← CM_dev0, 12 · parms14 ·  $\frac{parms_{15}}{parms_{16}}$  + tmp0 ·  $\sqrt{1 - (parms_{14})^2}$ 
  for i ∈ 1.. 3·365 - 1
    tmpi ← x·CM_devi, 12 + k·tmpi-1 + y·tmpi
  tmp

```

Note:

As mentioned earlier, it turns out to be easier to generate the fine mass deviations after generating the coarse mass deviations. The reason is that, in general, we want sequences that have a distinct autocorrelations and are correlated with each other either spatially or just through time. This function takes as an input the coarse mass deviations (although it only uses the center point). The constants  $x$ ,  $k$ , and  $y$  are adjusted to yield (approximately) the desired correlations.

### 3.5 The True Concentrations Generator

Input: The matrix  $bchol$ , from the function BChol, and  $parms$ , the vector of simulation parameters.

```

Gen_truth(bchol, parms) :=
  X ← Gen_CM_deviations (bchol, parms)
  z ← Gen_FM_deviations (X, parms)
  f25 ← e $\frac{-\left[\sqrt{\ln\left(\text{parms}_{15}\right)^2+1}\right]^2}{2}$ 
  f10 ← e $\frac{-\left[\sqrt{\ln\left(\text{parms}_{16}\right)^2+1}\right]^2}{2}$ 
  B25 ←  $\frac{-1 + \text{parms}_{18}}{1 + \text{parms}_{18}}$ 
  B10 ←  $\frac{-1 + \text{parms}_{19}}{1 + \text{parms}_{19}}$ 
  for j ∈ 0.. 3·365 - 1
    Truth25j ←  $\left(1 + B25 \sin\left(\frac{j \cdot 2 \cdot \pi}{365}\right)\right) \cdot \exp(z_j) \cdot f25$ 
    TruthC(j) ←  $\left[\left(1 + B10 \sin\left(\frac{j \cdot 2 \cdot \pi \cdot \text{parms}_{21}}{365} + \frac{\text{parms}_{20}}{12} \cdot 2 \cdot \pi\right)\right) \cdot \exp\left[\left(X^T\right)^{(j)}\right] \cdot f10\right]$ 
  trans_tC ← TruthCT
  for i ∈ 0.. 24
    tmp1i ← Cal_pt(trans_tC(i), parms)
  tm ← mean(tmp1)
  out ← augment $\left(\text{Truth25} \cdot \text{parms}_{17}, \text{trans\_tC}^{(12)} + \text{parms}_{17} \cdot \text{Truth25}\right) \cdot \frac{1}{\text{tm}}$ 
  outrows(out), 0 ←  $\frac{\text{mean}(\text{trans\_tC})}{\text{tm}}$ 
  out

```

#### Notes:

This function calls the previous two functions and combines them to create the sequences of true concentrations with the appropriate long-term seasonal trends, except that

everything is scaled to have a daily true-design value of 1. The output includes the annual true-design value corresponding to the scaled sequences.

Note that the output only includes the temporal sequences for the center grid and a summary value.

### 3.6 The Observed Values Generator

Input: The matrix of (scaled) true concentrations from the function Gen\_truth and *parms*, the vector of simulation parameters.

```

Gen_obs(truth, parms) :=
  v1_0 ← 0
  v2_0 ← 0
  for ii ∈ 0..2
    for j ∈ 0..3
      min ← floor( (j·365) / (4·parms_5) ) · parms_5 + ii·365
      max ← floor( (j+1)·365 / (4·parms_5) ) · parms_5 + ii·365
      num_sampled25 ← ceil( (max - min) / parms_5 · parms_6 )
      num_sampled10 ← ceil( (max - min) / parms_5 · parms_7 )
      v1 ← stack [ v1, parms_5 · Select( (max - min) / parms_5, num_sampled25 ) + min ]
      v2 ← stack [ v2, parms_5 · Select( (max - min) / parms_5, num_sampled10 ) + min ]
    k ← 0
    for ii ∈ 1..rows(v1) - 1
      for jj ∈ 1..rows(v2) - 1
        if v1_ii = v2_jj
          v_k ← v1_ii
          k ← k + 1
      total_sampled ← rows(v)
      e ← augment(rnorm(total_sampled, 0, 1), rnorm(total_sampled, 0, 1))
      for i ∈ 0..total_sampled - 1
        obs_i,1 ← [ (truth_v_i,1) · (1 + parms_11) · (1 + e_i,0 · parms_9) ] - truth_v_i,0 · (1 - parms_10) · (1 + e_i,1 · parms_8)
        obs_i,0 ← [ truth_v_i,1 · (1 - parms_11) · (1 + e_i,0 · parms_9) ] - truth_v_i,0 · (1 + parms_10) · (1 + e_i,1 · parms_8)
        obs_i,0 ← 0 if obs_i,0 < 0
        obs_i,1 ← 0 if obs_i,1 < 0
      augment(obs, v)

```

Note:

This function takes in the true concentrations for the center grid, selects the days throughout the 3-year period that will have corresponding sample values, and then simulates the observed values. Two versions of the observed values are generated: one corresponding to the maximum positive bias and one corresponding to the most extreme negative bias. The last column gives the number of days into the simulation to which the observed values correspond.

### 3.7 The Design Values Calculator for the Observations

Input: The matrix  $X$  of (scaled) observed concentrations from the function `Gen_obs`,  $parms$ , the vector of simulation parameters, and the scaling factor,  $tm$ , for switching between the daily and annual standard.

```

Cal_pobs (X,parms,tm) :=
  j ← 0
  k ← 0
  for i ∈ 0..rows(X) - 1
    | j ← j + (Xi,2 < 365)
    | k ← k + (Xi,2 < 365·2)
  r0 ← floor(j·parms4)
  r1 ← floor[(k - j)·parms4]
  r2 ← floor[(rows(X) - k)·parms4]
  for i ∈ 0..1
    | p0,i ← sort(submatrix(X,0,j - 1,i,i))(r0)
    | am0,i ← mean(submatrix(X,0,j - 1,i,i))
  for i ∈ 0..1
    | p1,i ← sort(submatrix(X,j,k - 1,i,i))(r1)
    | am1,i ← mean(submatrix(X,j,k - 1,i,i))
  for i ∈ 0..1
    | p2,i ← sort(submatrix(X,k,rows(X) - 1,i,i))(r2)
    | am2,i ← mean(submatrix(X,k,rows(X) - 1,i,i))
  for i ∈ 0..1
    | mi ← mean(p⟨i⟩)
    | mi+2 ← mean(am⟨i⟩) · 1/tm
  m

```

Note:

Since the number of days with both PM<sub>10</sub> and PM<sub>2.5</sub> data can vary, the function first determines to which year the data values correspond. Then, as with real data, annual values are evaluated for each of the three years. Finally, these are averaged over the three years to obtain the two scaled design values for the simulation for the positive bias case and two scaled design values for the negative bias case.

### 3.8 The Percentile Standard Calculator for the Truth

Input: A vector  $X$  of (scaled) true concentrations and the vector  $parms$ , of the input parameters.

$$\text{Cal\_pt}(X, parms) := \left| \begin{array}{l} r \leftarrow \text{floor}\left(\text{floor}\left(\frac{\text{rows}(X)}{3}\right) \cdot \text{parms}_4\right) \\ \text{for } i \in 0..2 \\ \quad p_i \leftarrow \text{sort}\left[\text{submatrix}\left[X, i \cdot \text{floor}\left(\frac{\text{rows}(X)}{3}\right), (i+1) \cdot \text{floor}\left(\frac{\text{rows}(X)}{3}\right) - 1, 0, 0\right]\right]_r \\ \text{mean}(p) \end{array} \right|$$

Note:

This function is less complicated than the version for the observations since the data are always complete with the same number of data points per year. Moreover, this completeness property implies that the annual design value is just the mean of all of the values and is not computed by this function.

### 3.9 Random Selection of Sample Days

Input: Integers  $n$  and  $k$ .

$$\text{Select}(n, k) := \left| \begin{array}{l} j \leftarrow 0 \\ \text{for } i \in 0..n-1 \\ \quad \text{if } [\text{rnd}(1) \cdot (n-i) < (k-j)] \\ \quad \quad \left| \begin{array}{l} v_j \leftarrow i \\ j \leftarrow j+1 \end{array} \right. \\ v \end{array} \right|$$

Note:

This function randomly selects  $k$  integers between 0 and  $n-1$ . This is an auxiliary function used in the generation of the observations to randomly pick the days that will have  $PM_{10}$  data and which will have  $PM_{2.5}$  data (the selections are made independent of each other).

### 3.10 Decision Performance Calculation

Input: A scalar  $x$  equal to the desired observed-design value target, an indicator  $std\_ind$  to indicate whether the daily or annual standard is being considered, the matrix output  $tt$  from the function `Gen_scaled_dist`, and the vector  $parms$  of input parameters.

```
P_obdv_ge_stnd (x, std_ind , tt, parms) :=
  if std_ind = 0
  | std ← parms2
  | round ← 0.5
  otherwise
  | std ← parms3
  | round ← 0.05
  mm ← tt·x
  rr0 ← 0
  rr1 ← 0
  for k ∈ 0..reps - 1
  | rr0 ← [mmk, 0+std_ind > (std + round)] + rr0
  | rr1 ← [mmk, 1+std_ind > (std + round)] + rr1
  rr ← rr ·  $\frac{1}{reps}$ 
  rr
```

Note:

This function computes the probability that the observed-design value exceeds the standard given that the design value of the true concentrations is equal to  $x$ .

### 3.11 Gray Zone and Performance Curve Computation

Input: The vector *parms* of input parameters.

```
Gray_zone(parms) := | tt ← Gen_scaled_dist (parms)
                    | out_day ← Gray_zone_day (parms,tt)
                    | out_ann ← Gray_zone_ann (parms,tt)
                    | out ← augment(out_day ,out_ann )
                    | out
```

Note:

In addition to computing the gray zones, this function computes the values for plotting the performance curves.

The final output consists of a matrix with 102 rows and 6 columns. The last row has the gray zones for the daily and annual standard. The other rows contain the values for plotting the decision performance curves. The first two columns are the probabilities of observing design values greater than the daily standard for the positive and negative bias cases. The third column consists of the daily true-design values for the grid. The fourth and fifth columns are the probabilities of observing design values greater than the annual standard for the positive and negative bias cases. The sixth column consists of the annual true-design values for the grid.

### 3.11.1 Gray Zone and Performance Curve Computation for the Daily Standard

Input: The vector *parms* of input parameters and the output of the Gen\_scaled\_dist function.

```

Gray_zone_day (parms, tt) :=
  min2 ← min(tt<sup>0</sup>)
  max2 ← max(tt<sup>1</sup>)
  start_day ← floor( (parms<sub>2</sub> - .0) / max2 )
  end_day ← ceil( (parms<sub>2</sub> + .75) / (min2 + .001) )
  delta2 ← (end_day - start_day) / 100
  for k ∈ 0..100
    ttest<sup>k</sup> ← P_obdv_ge_stnd (start_day + k·delta2, 0, tt, parms)
    v<sub>k</sub> ← start_day + k·delta2
  L ← 0
  u ← 0
  for i ∈ 0..cols(ttest) - 2
    L ← i if ttest<sub>1,i</sub> < parms<sub>0</sub>
    u ← i if ttest<sub>0,i</sub> < (1 - parms<sub>1</sub>)
  out_day ← (start_day end_day) if (ttest<sub>1,L+1</sub> - ttest<sub>1,L</sub>) · (ttest<sub>0,u+1</sub> - ttest<sub>0,u</sub>) = 0
  out_day ← start_day + [ L + ( (parms<sub>0</sub> - ttest<sub>1,L</sub>) / (ttest<sub>1,L+1</sub> - ttest<sub>1,L</sub>) ) u + [ (1 - parms<sub>1</sub>) - ttest<sub>0,u</sub> / (ttest<sub>0,u+1</sub> - ttest<sub>0,u</sub>) ] ] · delta2 otherwise
  plot_data ← augment(ttest<sup>T</sup>, v)
  out_day<sub>0,2</sub> ← 0
  stack(plot_data, out_day)

```

### 3.11.2 Gray Zone and Performance Curve Computation for the Annual Standard

Input: The vector *parms* of input parameters and the output of the Gen\_scaled\_dist function.

```

Gray_zone_ann(parms, tt) :=
  min1 ← min(tt<sup>2</sup>)
  max1 ← max(tt<sup>3</sup>)
  start_ann ← floor( $\frac{\text{parms}_3 - .0}{\text{max1}}$ )
  end_ann ← ceil( $\frac{\text{parms}_3 + .5}{\text{min1} + .001}$ )
  delta ←  $\frac{\text{end\_ann} - \text{start\_ann}}{100}$ 
  for k ∈ 0..100
    test<sup>k</sup> ← P_obdv_ge_stdn (start_ann + k·delta, 2, tt, parms)
    v_k ← start_ann + k·delta
  L ← 0
  u ← 0
  for i ∈ 0..cols(test) - 2
    L ← i if test1,i < parms0
    u ← i if test0,i < (1 - parms1)
  out_ann ← (start_ann end_ann) if (test1,L+1 - test1,L) · (test0,u+1 - test0,u) = 0
  out_ann ← start_ann +  $\left[ L + \frac{\text{parms}_0 - \text{test}_{1,L}}{\text{test}_{1,L+1} - \text{test}_{1,L}} \right] u + \left[ \frac{(1 - \text{parms}_1) - \text{test}_{0,u}}{\text{test}_{0,u+1} - \text{test}_{0,u}} \right] \cdot \text{delta}$  otherwise
  out_ann0,2 ← 0
  plot_data ← augment(testT, v)
  stack(plot_data, out_ann)

```