# Sample Design in the Bayesian Analysis of Radiocarbon Dates in Paleoecology

J. Andrés Christen     Maarten Blaauw

Centro de Investigación en Matemáticas (CIMAT)
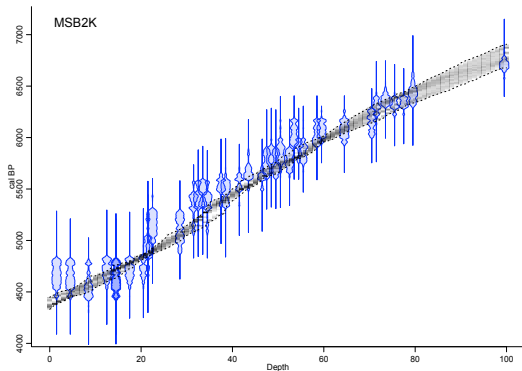Guanajuato, Mexico.

U. de Belfast, UK.
email: *jac@cimat.mx*, web page: http://www.cimat.mx/~jac

PBW, San Miguel A., Aug. 2010

# Paleoecology



Figure: Sampling from a peat core, recent goelogical past (less than 40,000 years to the present). Samples are collected at several depths and some samples are radiocarbon dated.

Figure: Chronology for the MSB2K peat profile, sampled from a peat bog in Meerstalblok, Neatherlands. 40 radiocarbon datings; the posterior depth-age relationship is obtained using a self adjusted MCMC.

We have a set of radiocarbon datings $y_j \pm \sigma_j; j = 1, 2, \ldots, m$, of samples taken along the peat core, at depths $d_j$.

We use a semiparametric model, to establish the (unknown) relationship between depth $d$ and age $G$ of the profile

$$G(d, \theta, x) = \theta + \sum_{j=1}^{i} x_j \Delta c + x_{i+1}(d - c_i);$$

where $c_i \leq d < c_{i+1}, i < K$, and $c_0 < c_1 < \cdots < c_K$ are depths uniformly spaced along the core every $\Delta c$ cm and $x = (x_1, x_2, \ldots, x_K)$.

# Radiocarbon dating, Bayesian chronology and MCMC

That is, the core profile is divided into $K$ equidistant sections and $x_j$ is the peat accumulation rate (yr/cm) in section $j$.

Every radiocarbon determination has an asociated standard (measuring) error $\sigma_j$. We use a new robust t model (fallowing Christen and Pérez E., 2009), with $E[y_j \mid d_j, x] = \mu(G(d_j, \theta, x))$ where $\mu$ is the *Internationally Agreed Radiocarbon Calibration Curve, INTCAL09*, that relates (calibrates) calendar years to radiocarbon ages.

We perform a MCMC analysis (the prior information on the accumulation rates is very important) to obtain our posterior distributions. We use a robust and self adjusting MCMC sampler: the t-walk (Christen and Fox, 2010).

# Where to date? Costs?

But, how should we decide on the depths where it is wise to radiocarbon date and how many dates should we consider?

Considering the fact that each radiocarbon dating will cost around **€500** each.

# Brute force (formal) approach

Christen and Buck (1998) present a brute force design approach, in a similar radiocarbon dating problem:

1. We define an Utility Function $U$, measuring the relative gain of radiocarbon dates against their cost and start with a (trial) set of already dated samples $\boldsymbol{y}_N$.

2. We select a test set $M$ (a candidate design) from the set of all samples available $A$.

3. Using our MCMC we simulate, from the predictive distribution of samples in $M$, a set of radiocarbon results. For each of these simulated samples, a second MCMC is used to recalculate the chronology.

4. The corresponding utility is evaluated each run, and over many sampled data, the corresponding expected utility of the trail set, $U^*(M)$, may be approximated.

5. The same is repeated for all possible designs (all $M \subset A$) to obtain an optimal design $M^*$.

## Brute force (formal) approach

Christen and Buck (1998) present a brute force design approach, in a similar radiocarbon dating problem:

1. We define an Utility Function $U$, measuring the relative gain of radiocarbon dates against their cost and start with a (trial) set of already dated samples $y_N$.

2. We select a test set $M$ (a candidate design) from the set of all samples available $A$.

3. Using our MCMC we simulate, from the predictive distribution of samples in $M$, a set of radiocarbon results. For each of these simulated samples, a second MCMC is used to recalculate the chronology.

4. The corresponding utility is evaluated each run, and over many sampled data, the corresponding expected utility of the trail set, $U^*(M)$, may be approximated.

5. The same is repeated for all possible designs (all $M \subset A$) to obtain an optimal design $M^*$.

## Brute force (formal) approach

Christen and Buck (1998) present a brute force design approach, in a similar radiocarbon dating problem:

1. We define an Utility Function $U$, measuring the relative gain of radiocarbon dates against their cost and start with a (trial) set of already dated samples $\boldsymbol{y}_N$.
2. We select a test set $M$ (a candidate design) from the set of all samples available $A$.
3. Using our MCMC we simulate, from the predictive distribution of samples in $M$, a set of radiocarbon results. For each of these simulated samples, a second MCMC is used to recalculate the chronology.
4. The corresponding utility is evaluated each run, and over many sampled data, the corresponding expected utility of the trail set, $U^*(M)$, may be approximated.
5. The same is repeated for all possible designs (all $M \subset A$) to obtain an optimal design $M^*$.

## Brute force (formal) approach

Christen and Buck (1998) present a brute force design approach, in a similar radiocarbon dating problem:

1. We define an Utility Function $U$, measuring the relative gain of radiocarbon dates against their cost and start with a (trial) set of already dated samples $y_N$.

2. We select a test set $M$ (a candidate design) from the set of all samples available $A$.

3. Using our MCMC we simulate, from the predictive distribution of samples in $M$, a set of radiocarbon results. For each of these simulated samples, a second MCMC is used to recalculate the chronology.

4. The corresponding utility is evaluated each run, and over many sampled data, the corresponding expected utility of the trail set, $U^*(M)$, may be approximated.

5. The same is repeated for all possible designs (all $M \subset A$) to obtain an optimal design $M^*$.

# Brute force (formal) approach

Christen and Buck (1998) present a brute force design approach, in a similar radiocarbon dating problem:

1. We define an Utility Function $U$, measuring the relative gain of radiocarbon dates against their cost and start with a (trial) set of already dated samples $\mathbf{y}_N$.

2. We select a test set $M$ (a candidate design) from the set of all samples available $A$.

3. Using our MCMC we simulate, from the predictive distribution of samples in $M$, a set of radiocarbon results. For each of these simulated samples, a second MCMC is used to recalculate the chronology.

4. The corresponding utility is evaluated each run, and over many sampled data, the corresponding expected utility of the trail set, $U^*(M)$, may be approximated.

5. The same is repeated for all possible designs (all $M \subset A$) to obtain an optimal design $M^*$.

# Brute force (formal) approach

As may be seen, Christen and Buck (1998) is the standard Bayesian design solution (maximize posterior predictive expected utility over all possible designs), but in a MCMC setting is simply too computationally demanding and unfeasible to run in many cases.

# An heuristic alternative: Active Learning

Christen and Sansó (2010) propose an index to sequentially select sample design points in the context of Gaussian Processes (for the statistical analysis of computer experiments).

## The index A

Christen and Sansó (2010) propose an index based on an apprximation to the *Active Learning* ideas of Cohon, that propose to select a design point that contributes the most to the reduction in variance of the fitted model.

Christen and Sansó (2010) index is

$$A(d_{N+1}|\boldsymbol{y}_N) = \frac{1 - ||r(d_{N+1})||}{C^2} \frac{1}{m} \sum_{j=1}^{m} c(d_j, d_{N+1})^2 \tag{1}$$

where $||r(d_{N+1})||^2 = \sum_{i=1}^{N} c(d_{N+1}, d_i)^2$, $C = \max_{j=1,2,...,m} V(d_j)$, $c(\cdot, \cdot)$ is the covariance and $V(d) = c(d, d)$, the variance, at depth $d$. Using a renormalization such that $\min_{j=1,2,...,m} V(d_j) = 1$.

## The index A

Christen and Sansó (2010) propose an index based on an apprximation to the *Active Learning* ideas of Cohon, that propose to select a design point that contributes the most to the reduction in variance of the fitted model.

Christen and Sansó (2010) index is

$$A(d_{N+1}|\boldsymbol{y}_N) = \frac{1 - ||r(d_{N+1})||}{C^2} \frac{1}{m} \sum_{j=1}^{m} c(d_j, d_{N+1})^2 \tag{1}$$

where $||r(d_{N+1})||^2 = \sum_{i=1}^{N} c(d_{N+1}, d_i)^2$, $C = \max_{j=1,2,...,m} V(d_j)$, $c(\cdot, \cdot)$ is the covariance and $V(d) = c(d, d)$, the variance, at depth $d$. Using a renormalization such that $\min_{j=1,2,...,m} V(d_j) = 1$.

# The index A

Our index may be rewritten as

$$\frac{1 - ||r(d_{N+1})||}{C^2} \frac{1}{m} \left( V(d_{N+1})^2 + \sum_{j \neq N+1} c(d_j, d_{N+1})^2 \right).$$

- Our index prefers points (depths) with high variance (high variance in the fitted chronology),

- but also design points that are correlated with other points (potentially, we may obtain information about those points by sampling at one location only).

- Moreover, our index prefers points far (not correlated) with already sample points, given the factor $||r(d_{N+1})||$.

# The index A

Our index may be rewritten as

$$\frac{1 - ||r(d_{N+1})||}{C^2} \frac{1}{m} \left( V(d_{N+1})^2 + \sum_{j \neq N+1} c(d_j, d_{N+1})^2 \right).$$

- Our index prefers points (depths) with high variance (high variance in the fitted chronology),
- but also design points that are correlated with other points (potentially, we may obtain information about those points by sampling at one location only).
- Moreover, our index prefers points far (not correlated) with already sample points, given the factor $||r(d_{N+1})||$.

## The index A

Our index may be rewritten as

$$\frac{1 - ||r(d_{N+1})||}{C^2} \frac{1}{m} \left( V(d_{N+1})^2 + \sum_{j \neq N+1} c(d_j, d_{N+1})^2 \right).$$

- Our index prefers points (depths) with high variance (high variance in the fitted chronology),
- but also design points that are correlated with other points (potentially, we may obtain information about those points by sampling at one location only).
- Moreover, our index prefers points far (not correlated) with already sample points, given the factor $||r(d_{N+1})||$.

1. Given a trial of already radiocarbon dated samples $\boldsymbol{y}_N$, we have our MCMC simulation for $\theta_0^{(t)}, x^{(t)}, w^{(t)}, t = 1, 2, \ldots, T$ and for each simulation $t$ we are able to calculate

$$G(c_1, \theta_0^{(t)}, x^{(t)}), G(c_2, \theta_0^{(t)}, x^{(t)}), \ldots, G(c_k, \theta_0^{(t)}, x^{(t)}).$$

2. We may then calculated (estimate) the covariance of any pair of depths, for our fitted chronology, $c(d_i, d_j)$.

3. Using this estimated covariance $c(d_i, d_j)$, we calculate the index $A(d_{N+1}|\boldsymbol{y}_N)$, for each depth and that depth $d^*$ with the maximum $A$ index is the next sample considered for dating.

4. Send the sample at depth $d^*$ for a radiocarbon analysis and wait for the result. Alternatively, depending on logistics etc., we may impute a new date at depth $d^*$ with, for example

$$y_{N+1} = \frac{1}{T} \sum_{t=1}^{T} G(d^*, \theta_0^{(t)}, x^{(t)}).$$

JA Christen (CIMAT)  Design in Paleo Chronologies  PBW, San Miguel A., Aug. 2010  30

# Sequential sampling procedure for chronology building

1. Given a trial of already radiocarbon dated samples $y_N$, we have our MCMC simulation for $\theta_0^{(t)}, x^{(t)}, w^{(t)}, t = 1, 2, \ldots, T$ and for each simulation $t$ we are able to calculate

$$G(c_1, \theta_0^{(t)}, x^{(t)}), G(c_2, \theta_0^{(t)}, x^{(t)}), \ldots, G(c_k, \theta_0^{(t)}, x^{(t)}).$$

2. We may then calculated (estimate) the covariance of any pair of depths, for our fitted chronology, $c(d_i, d_j)$.

3. Using this estimated covariance $c(d_i, d_j)$, we calculate the index $A(d_{N+1}|y_N)$, for each depth and that depth $d^*$ with the maximum $A$ index is the next sample considered for dating.

4. **Send the sample at depth $d^*$ for a radiocarbon analysis and wait for the result**. Alternatively, depending on logistics etc., we may **impute** a new date at depth $d^*$ with, for example

$$y_{N+1} = \frac{1}{T} \sum_{t=1}^{T} G(d^*, \theta_0^{(t)}, x^{(t)}).$$

1. Given a trial of already radiocarbon dated samples $y_N$, we have our MCMC simulation for $\theta_0^{(t)}, x^{(t)}, w^{(t)}, t = 1, 2, \ldots, T$ and for each simulation $t$ we are able to calculate

$$G(c_1, \theta_0^{(t)}, x^{(t)}), G(c_2, \theta_0^{(t)}, x^{(t)}), \ldots, G(c_k, \theta_0^{(t)}, x^{(t)}).$$

2. We may then calculated (estimate) the covariance of any pair of depths, for our fitted chronology, $c(d_i, d_j)$.

3. Using this estimated covariance $c(d_i, d_j)$, we calculate the index $A(d_{N+1}|y_N)$, for each depth and that depth $d^*$ with the maximum $A$ index is the next sample considered for dating.

4. **Send the sample at depth $d^*$ for a radiocarbon analysis and wait for the result**. Alternatively, depending on logistics etc., we may **impute** a new date at depth $d^*$ with, for example

$$y_{N+1} = \frac{1}{T} \sum_{t=1}^{T} G(d^*, \theta_0^{(t)}, x^{(t)}).$$

# Sequential sampling procedure for chronology building

1. Given a trial of already radiocarbon dated samples $\boldsymbol{y}_N$, we have our MCMC simulation for $\theta_0^{(t)}, x^{(t)}, w^{(t)}, t = 1, 2, \ldots, T$ and for each simulation $t$ we are able to calculate

   $$G(c_1, \theta_0^{(t)}, x^{(t)}), G(c_2, \theta_0^{(t)}, x^{(t)}), \ldots, G(c_k, \theta_0^{(t)}, x^{(t)}).$$

2. We may then calculated (estimate) the covariance of any pair of depths, for our fitted chronology, $c(d_i, d_j)$.

3. Using this estimated covariance $c(d_i, d_j)$, we calculate the index $A(d_{N+1}|\boldsymbol{y}_N)$, for each depth and that depth $d^*$ with the maximum $A$ index is the next sample considered for dating.

4. **Send the sample at depth $d^*$ for a radiocarbon analysis and wait for the result**. Alternatively, depending on logistics etc., we may **impute** a new date at depth $d^*$ with, for example

   $$y_{N+1} = \frac{1}{T} \sum_{t=1}^{T} G(d^*, \theta_0^{(t)}, x^{(t)}).$$

# Sequential sampling procedure for chronology building

1. Given a trial of already radiocarbon dated samples $\boldsymbol{y}_N$, we have our MCMC simulation for $\theta_0^{(t)}, x^{(t)}, w^{(t)}, t = 1, 2, \ldots, T$ and for each simulation $t$ we are able to calculate

   $$G(c_1, \theta_0^{(t)}, x^{(t)}), G(c_2, \theta_0^{(t)}, x^{(t)}), \ldots, G(c_k, \theta_0^{(t)}, x^{(t)}).$$

2. We may then calculated (estimate) the covariance of any pair of depths, for our fitted chronology, $c(d_i, d_j)$.

3. Using this estimated covariance $c(d_i, d_j)$, we calculate the index $A(d_{N+1}|\boldsymbol{y}_N)$, for each depth and that depth $d^*$ with the maximum $A$ index is the next sample considered for dating.

4. **Send the sample at depth $d^*$ for a radiocarbon analysis and wait for the result**. Alternatively, depending on logistics etc., we may **impute** a new date at depth $d^*$ with, for example

   $$y_{N+1} = \frac{1}{T} \sum_{t=1}^{T} G(d^*, \theta_0^{(t)}, x^{(t)}).$$

## Sequential sampling procedure for chronology building

1. Given a trial of already radiocarbon dated samples $\boldsymbol{y}_N$, we have our MCMC simulation for $\theta_0^{(t)}, x^{(t)}, w^{(t)}, t = 1, 2, \ldots, T$ and for each simulation $t$ we are able to calculate

$$G(c_1, \theta_0^{(t)}, x^{(t)}), G(c_2, \theta_0^{(t)}, x^{(t)}), \ldots, G(c_k, \theta_0^{(t)}, x^{(t)}).$$

2. We may then calculated (estimate) the covariance of any pair of depths, for our fitted chronology, $c(d_i, d_j)$.

3. Using this estimated covariance $c(d_i, d_j)$, we calculate the index $A(d_{N+1}|\boldsymbol{y}_N)$, for each depth and that depth $d^*$ with the maximum $A$ index is the next sample considered for dating.

4. **Send the sample at depth $d^*$ for a radiocarbon analysis and wait for the result**. Alternatively, depending on logistics etc., we may **impute** a new date at depth $d^*$ with, for example

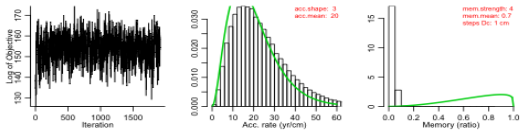$$y_{N+1} = \frac{1}{T} \sum_{t=1}^{T} G(d^*, \theta_0^{(t)}, x^{(t)}).$$

Figure: We take a trail sample of 3 radiocarbon dates only. Black dots: index *A* plotted in arbitraty scale.

Figure: We take a trail sample of 3 radiocarbon dates only. Black dots: index *A* plotted in arbitraty scale.

# Example with the MSB2K core, using our index *A*



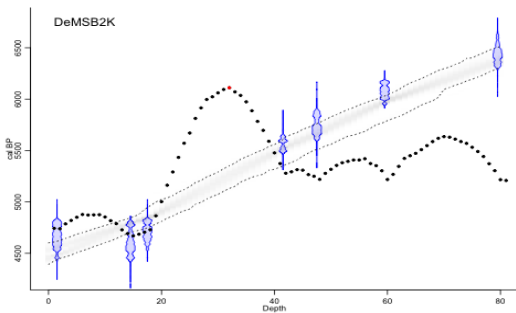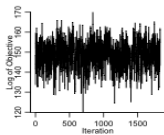Figure: We take a trail sample of 3 radiocarbon dates only. Black dots: index *A* plotted in arbitraty scale.

- We concentrate on the high density sampled section of 0 to 80 cm depth
- We calculate the index *A* for every possible depth to be radiocarbon dated.
- We plot the index *A*. The red dot is the maximum.
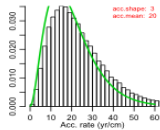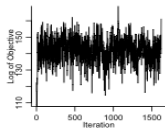- We radiocarbon date at the selected point, recalculate the chronology and start again.

Figure: We take a trail sample of 3 radiocarbon dates only. Black dots: index *A* plotted in arbitraty scale.

- We concentrate on the high density sampled section of 0 to 80 cm depth
- We calculate the index *A* for every possible depth to be radiocarbon dated.
- We plot the index *A*. The red dot is the maximum.
- We radiocarbon date at the selected point, recalculate the chronology and start again.

Figure: We take a trail sample of 3 radiocarbon dates only. Black dots: index *A* plotted in arbitraty scale.

- We concentrate on the high density sampled section of 0 to 80 cm depth
- We calculate the index *A* for every possible depth to be radiocarbon dated.
- We plot the index *A*. The red dot is the maximum.
- We radiocarbon date at the selected point, recalculate the chronology and start again.

# Example with the MSB2K core, using our index *A*



Figure: We take a trail sample of 3 radiocarbon dates only. Black dots: index *A* plotted in arbitraty scale.

- We concentrate on the high density sampled section of 0 to 80 cm depth
- We calculate the index *A* for every possible depth to be radiocarbon dated.
- We plot the index *A*. The red dot is the maximum.
- We radiocarbon date at the selected point, recalculate the chronology and start again.

Figure: We take a trail sample of 3 radiocarbon dates only. Black dots: index *A* plotted in arbitraty scale.

- We concentrate on the high density sampled section of 0 to 80 cm depth
- We calculate the index *A* for every possible depth to be radiocarbon dated.
- We plot the index *A*. The red dot is the maximum.
- We radiocarbon date at the selected point, recalculate the chronology and start again.

# Averange variance along the core



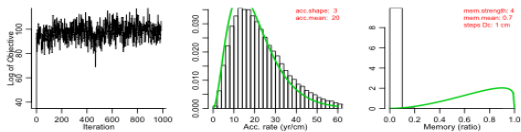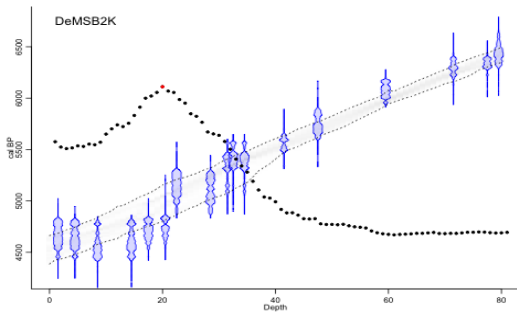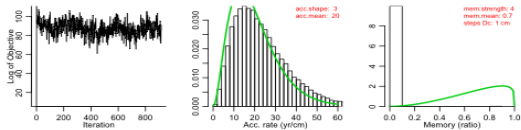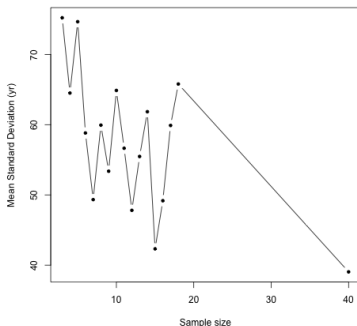Figure: Average variance for our chronology with increasing number of samples.

1. The use of our index *A* in this context seems like a good, feasible, alternative to the a formal and complex procedure.

2. It is however, not clear weather radiocarbon labs will be able to date samples sequentially, or could paleoecologists will be willing to wait for the results. We need also to consider sampling in batches, and experiment with imputing dates for candidate samples.

3. A stopping rule is also needed, to balance out the added gain in chronology precision with the cost of sampling.

# Discusión

1. The use of our index *A* in this context seems like a good, feasible, alternative to the a formal and complex procedure.

2. It is however, not clear weather radiocarbon labs will be able to date samples sequentially, or could paleoecologists will be willing to wait for the results. We need also to consider sampling in batches, and experiment with imputing dates for candidate samples.

3. A stopping rule is also needed, to balance out the added gain in chronology precision with the cost of sampling.

# Discusión

1. The use of our index *A* in this context seems like a good, feasible, alternative to the a formal and complex procedure.

2. It is however, not clear weather radiocarbon labs will be able to date samples sequentially, or could paleoecologists will be willing to wait for the results. We need also to consider sampling in batches, and experiment with imputing dates for candidate samples.

3. A stopping rule is also needed, to balance out the added gain in chronology precision with the cost of sampling.

# References

Christen, J.A. and Fox, C. (2010), "A General Purpose Sampling Algorithm for Continuous Distributions (the t-walk)", *Bayesian Analysis*, **4**(2, June), .

Christen, J.A and Prez E., S. (2009). "A New Robust Statistical Model for Radiocarbon Data", *Radiocarbon*, **51**(3), 1047–1059.

Christen, J.A. and Buck, C.E. (1998), "Sample selection in radiocarbon dating", *Applied Statistics*, **47**(3), 543–557.