

Closed-Form Estimation of Multiple Change-Points Models (Supplemental Material)

Greg Jensen
Columbia University

Abstract

Supplemental material to accompany “Closed-Form Estimation of Multiple Change-Points Models” by Greg Jensen. A variety of non-psychological datasets are analyzed to demonstrate broad applications of the *Conjugate Partitioned Recursion* (CPR) algorithm. A method for performing a sensitivity analysis is demonstrated. A weakness in the algorithm is identified in cases of large-scale stationary distributions, and two modifications to the standard procedure are proposed to circumvent that weakness: The ‘dicing’ operation (whose function is strictly exploratory) and the ‘forward-retrospective’ algorithm, which makes assessments sequentially rather than recursively. Finally, the mathematical basis for the conjugate priors invoked in the main article and the formulas for empirical Bayes ‘rule-of-thumb’ priors are provided.

Keywords: bayesian statistics, change-point analysis, marginal likelihood, time series analysis

Additional Examples

The examples provided in the main text showcase several experimental applications of the CPR algorithm. Although change-point algorithms remain uncommon in many experimental domains, they have a long history of use in actuarial, industrial, and econometric applications. In the interest of providing a bridge between that tradition, several further examples are provided here. These examples also demonstrate some of the probability distributions not discussed in the main text. In all cases, a decision criterion of $\tau = 10$ is used, and all analyses were performed on a dual-core 2.93GHz laptop.

British Coal-Mining Disasters

One of the signature datasets in the statistical analysis of change-points reports British coal-mining disasters (involving at least 10 men) over the period from March 15, 1851 to March 22, 1962. These data have been reported both as the interval between disasters in days (Jarrett, 1979) and as the number of disasters per year (Carlin et al., 1992).

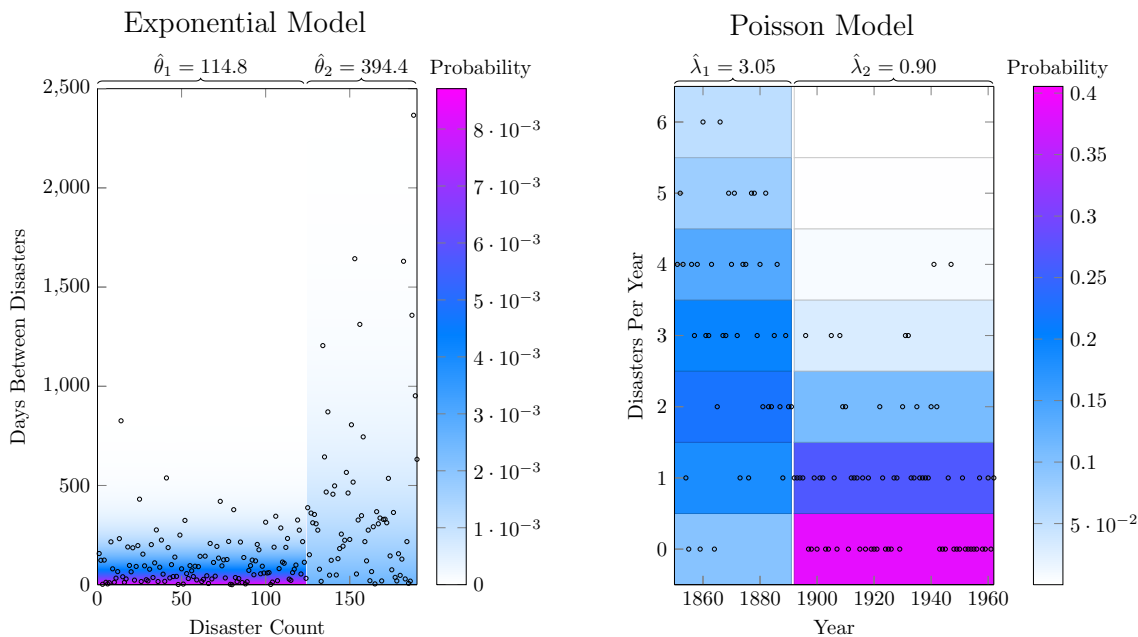


Figure 1. British coal-mining disasters between 1851 and 1962, represented as days between disasters (left) and number of disasters per year (right). In both cases, the CPR algorithm, using the appropriate rule-of-thumb empirical prior, finds a single change-point. Each plot shows its distribution’s probability function as a color map, based on the posterior parameter estimates.

1 Because the probability of a major coal mining disaster is low, and because disasters
 2 at different mines are presumably independent of one another, the intervals between dis-
 3 asters follow an exponential distribution, while the number of disasters per year (or any
 4 other interval) is Poisson-distributed. This intimate relationship between the exponential
 5 distribution and the Poisson distribution is clear from their conjugate relationship: The
 6 MML for the exponential distribution is $\text{Gamma}(n, \sum x)$, while the MML for the Poisson
 7 distribution is $\text{Gamma}(\sum x, n)$.

8 Figure 1 shows the results of the CPR algorithm for the coal-mining disaster
 9 data, represented both as an exponential model of intervals measured in days (prior
 10 $= \text{Pr}(\theta) = \text{Gamma}(1, 114)$) and as a Poisson model of disasters per year (prior $=$
 11 $\text{Pr}(\lambda) = \text{Gamma}(1, 1)$). Here, the exponential parameter θ estimates the distribution of
 12 ‘days between disasters’ and the Poisson parameter λ estimates the distribution of ‘disasters
 13 per year.’ In both cases, a single change-point is identified. The exponential model places
 14 the change between 124th and 125th disasters, between March 10th, 1890 and April 2nd,
 15 1891, whereas the Poisson model places the change between 1892 and 1893. These changes
 16 coincide with a period of regulatory reform in British mining, particularly the Coal Mines
 17 Act of 1887 (Anderson, 1911). Note that the parameters estimates are very similar, but
 18 not identical: $[\hat{\theta}_1 = 114.8 \approx \frac{365}{\lambda_1} = 119.7]$; $[\hat{\theta}_2 = 394.4 \approx \frac{365}{\lambda_2} = 405.6]$. Given the small size
 19 of the dataset (190 intervals over 112 years), computing the position of this change-point
 20 required less than a tenth of a second.

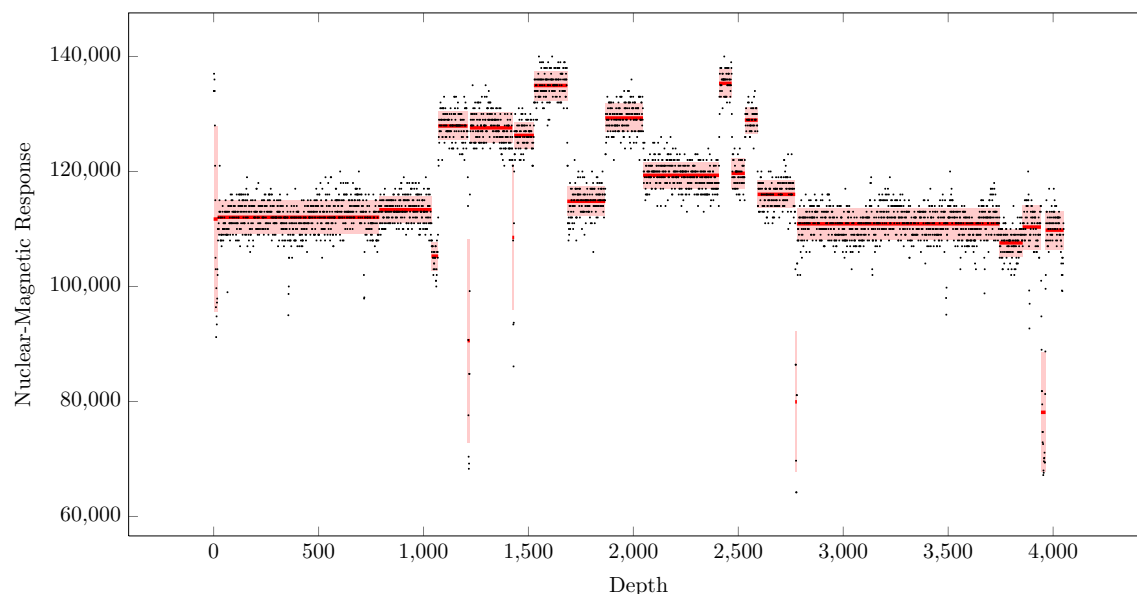


Figure 2. Nuclear-magnetic response of rock strata in an oil well bore-hole. Points correspond to individual observations. The red line indicates the estimated mean (with one standard deviation shown in light red) within segments identified by the CPR algorithm.

1 Given a choice between representing the data using either distribution, using the
 2 Poisson representation results in a loss of some information (dependent on the size of the
 3 intervals), so an exponential representation is preferred when precise timestamps are known
 4 for all events. Consequently, the exponential estimate is likely to be a better estimate. In
 5 many cases, however, counts within regular intervals are the only data available from the
 6 records.

7 Well-Log Data

8 Another well-studied dataset is a series of 4,050 nuclear-magnetic response measure-
 9 ments from a geological assay. Measurements were taken at regular intervals as a probe
 10 was lowered into a bore-hole, with changes in the response implying strata with different
 11 properties. Originally introduced by Ó Ruanaidh & Fitzgerald (1996), this dataset pro-
 12 vides a robust test of change-point algorithms because it only weakly conforms to standard
 13 distributional assumptions, both in terms of outliers (which could reflect either distinctive
 14 strata or mere measurement noise) and substantial and inconsistent rounding error.

15 Assuming that the data followed a Gaussian distribution with unknown parameters,
 16 Figure 2 shows the change-points identified by the CPR algorithm, given a Normal-Gamma
 17 prior based on the median and median absolute deviation of the data ($a = 114000$, $b = 0.1$,
 18 $c = 1$, $d = 2965.2$). This prior remains relatively weak with respect to the mean (as $b = 0.1$
 19 corresponds to only 1/10 of a hypothetical observation), but is somewhat stronger with
 20 respect to the variance (as $c = 1$ corresponds to one full hypothetical observation). The
 21 prior standard deviation d was estimated from the median absolute deviation of the the first
 22 differences (that is, the difference between consecutive observations), a robust econometric
 23 technique for assessing volatility in the presence of omitted variables (Wooldridge, 2002).

1 The resulting change-point model consists of 22 change-points, taking approximately 14
2 seconds to calculate.

3 One can get a sense of the CPR algorithm’s rapidity by comparing it to the “posterior
4 simulation” analyses performed on the dataset by Fearnhead (2006). These analyses
5 use recursive sampling to effect a numerical integration of the posterior odds. When their
6 analysis was performed on a 3.4GHz PC, a “piecewise-constant” posterior simulation algo-
7 rithm performed 10,000 simulations that required 26 seconds to run, identifying between
8 45 and 60 change-points. Fearnhead also ran a more robust “random walk” model, which
9 took approximately 19 minutes and only identified between 12 and 21 change-points; these
10 models were subsequently compared to a brute-force Markov-chain Monte Carlo (MCMC)
11 method.

12 This result is impressive, with a runtime that outperformed MCMC by about two
13 orders of magnitude. Nevertheless, implementation requires that the analyst make sub-
14 stantial assumptions (such as consistent variance over time), and required that outliers be
15 removed from the data prior to analysis. Additionally, because the method relies on random
16 resampling, slightly different results are obtained with each simulation. The CPR algorithm
17 compares very favorably, as it identifies close to the optimal number of change-points in
18 less time than the faster algorithm reported by Fearnhead, produces identical results every
19 time it is run, permits the variance to change from one segment to the next, and does not
20 require the pruning of outliers.

21 **Treasury Bill Rates**

22 In econometrics, change-point analysis is often demonstrated using the nominal rate
23 on three-month U.S. Treasury bills (or T-bills). Because these data are freely available to
24 the public, and are regularly updated, it has been used as a test case in both retrospective
25 (Bai, 1997) and prospective (Pesaran et al., 2006) change-point analysis.

26 Figure 3 shows two different implementation of the CPR algorithm to monthly T-bill
27 data over the period from June 1947 to February 2013. The top plot shows the posterior
28 estimates for a linear regression model. This approach identifies 28 change-points, and is
29 highly sensitive to big shifts. The bottom plot shows the ‘first differences’ (Wooldridge,
30 2002), which emphasizes the volatility in the data over the current rate (Meligkotsidou
31 & Dellaportas, 2011). Here, as in Figure 2, both the mean and one standard deviation
32 are shown. This approach identifies only 15 change-points, because segments with similar
33 volatility (e.g. throughout the 90s) are likely to be grouped together.

34 The contrasting benefits of these two methods depends on the purpose of the analysis.
35 The regression analysis is more complex, devoting three parameters to each segment (slope,
36 intercept, and the free cells of the estimated coefficient covariance matrix), and provides a
37 much more accurate historical description. On the other hand, the first differences provide
38 a more useful summary statistic for examining volatility. A comparable experimental case
39 might be the study of motor impairment: Studying of the distance moved over each time
40 interval is likely to be more informative than absolute position when tremors are a symptom
41 of interest.

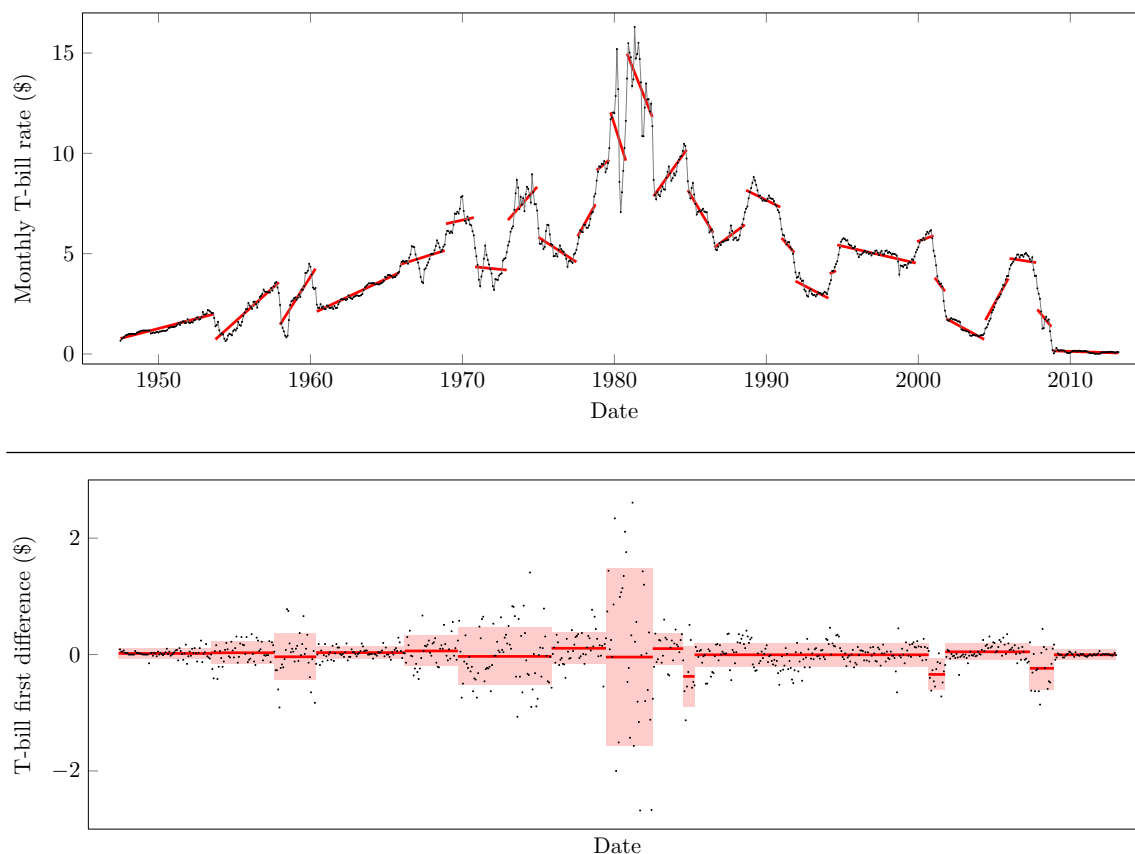


Figure 3. Rates on three-month U.S. Treasury bills, calculated monthly from June 1947 to February 2013. (Top) Nominal rates, with red lines displaying the model inferred by a linear regression change-point analysis. (Bottom) First differences of nominal rates $(x_i - x_{i-1})$, with the means (red) and one standard deviation (light red) from a Gaussian change-point analysis.

1

Sensitivity Analysis

2 An essential consideration in using the CPR algorithm is the selection of the decision
 3 criterion τ . In the examples presented, a criterion of $\tau = 10$ has been used because it
 4 provides a good balance between description and parsimony. ‘Good’ is a subjective term,
 5 however, and there is not always a good correspondence between results that are statistically
 6 vs. theoretically significant. This is further complicated by the relationship between a
 7 criterion’s efficacy and subtle differences in the assumptions underlying different models.

8 A benefit of the CPR algorithm is that, given its reliance on closed-form solutions, it
 9 is trivial to perform a post-hoc calculation of the likelihood of the data once the parameters
 10 have been estimated. The overall efficacy of the entire model (across all segments) can then
 11 be compared to other such models using the Schwarz-Bayes Information Criterion (SBIC;
 12 Schwarz, 1978). As noted in the main text, SBIC arises from an approximation of the

1 marginal model likelihood $m(x, M)$:

$$m(x, M) \tilde{\propto} f(x|\theta, M) \left(\frac{1}{n}\right)^{p/2} \quad (1)$$

2 The conventional practice is to convert this approximation into a score measured on a log
3 scale:

$$SBIC(x) = -2 \cdot \log(f(x|\theta, M)) + p \cdot \log(n) + C_x \quad (2)$$

4 Here, given the observations x , a value is calculated that depends on twice the negative log
5 of the likelihood function $f(x|\theta, M)$, on the number of free parameters p , and on the total
6 number of observations n . For data x , there is also a constant C_x , which prevents the direct
7 comparison of two datasets x and y (as they might have different constants). However, the
8 SBIC can be used to compare how well different models fit the data x because the constant
9 C_x will cancel out of such comparisons. The model with the *lowest* SBIC is considered
10 the best, because its overall marginal model likelihood is approximately the highest (Yang,
11 2005).

12 Although the SBIC was originally intended to be used in cases where observations in-
13 dependent and identically distributed, its derivation has been generalized to non-identically
14 distributed cases as well, having particular success in comparing time-series models (Ca-
15 vanaugh & Neath, 1999). This makes it a suitable metric for performing a sensitivity
16 analysis.

17 In order to examine the effects of the decision criterion τ on the resulting change-
18 point model, a sensitivity analysis consists of running the analysis multiple times using
19 different criteria but keeping other parameters constant. For each criterion that is tested,
20 the resulting change-point model \mathbf{M} is then used to calculate an SBIC using the following
21 equation:

$$SBIC(x, \mathbf{M}) = -2 \left[\sum_{i \in \mathbf{M}} \log(f(x_i|\theta_i, M)) \right] + (\text{length}(\mathbf{M})(p_M + 1) - 1) \cdot \log(n) + C_x \quad (3)$$

22 Here, the log likelihood is summed in each of the segments of \mathbf{M} , and the free parameters
23 consist of the one set of model parameters p_M for each segment, plus a free parameter for the
24 position of each change-point (Yao, 1988). In the analyses performed below, the parameters
25 θ_i were estimated post-hoc, using frequentist methods, to minimize the relationship between
26 the prior assumptions and the subsequent sensitivity analysis.

27 Figure 4 demonstrates the results of such a sensitivity analysis for the datasets dis-
28 cussed in the main article and in the supplement. Although each dataset displays its own
29 idiosyncrasies, several themes are evident from their comparison.

30 First, it is clear that models vary in their sensitivity to the decision criterion. For
31 example, linear regression models are quite resistant to changes in τ . The regression model
32 describing T-bill rates differs only by a few change-points over the range $1 \leq \tau \leq 100$
33 (Figure 4, lower left), and the reaction time data from Palmeri (1997) are not included in
34 Figure 4 because they returned the same result regardless of the decision criterion used. On
35 the other hand, the binomial data describing SimChain performance (Figure 4, upper left,
36 Jensen et al., 2013), and the Gaussian distribution of T-bill rate first differences (Figure 4,

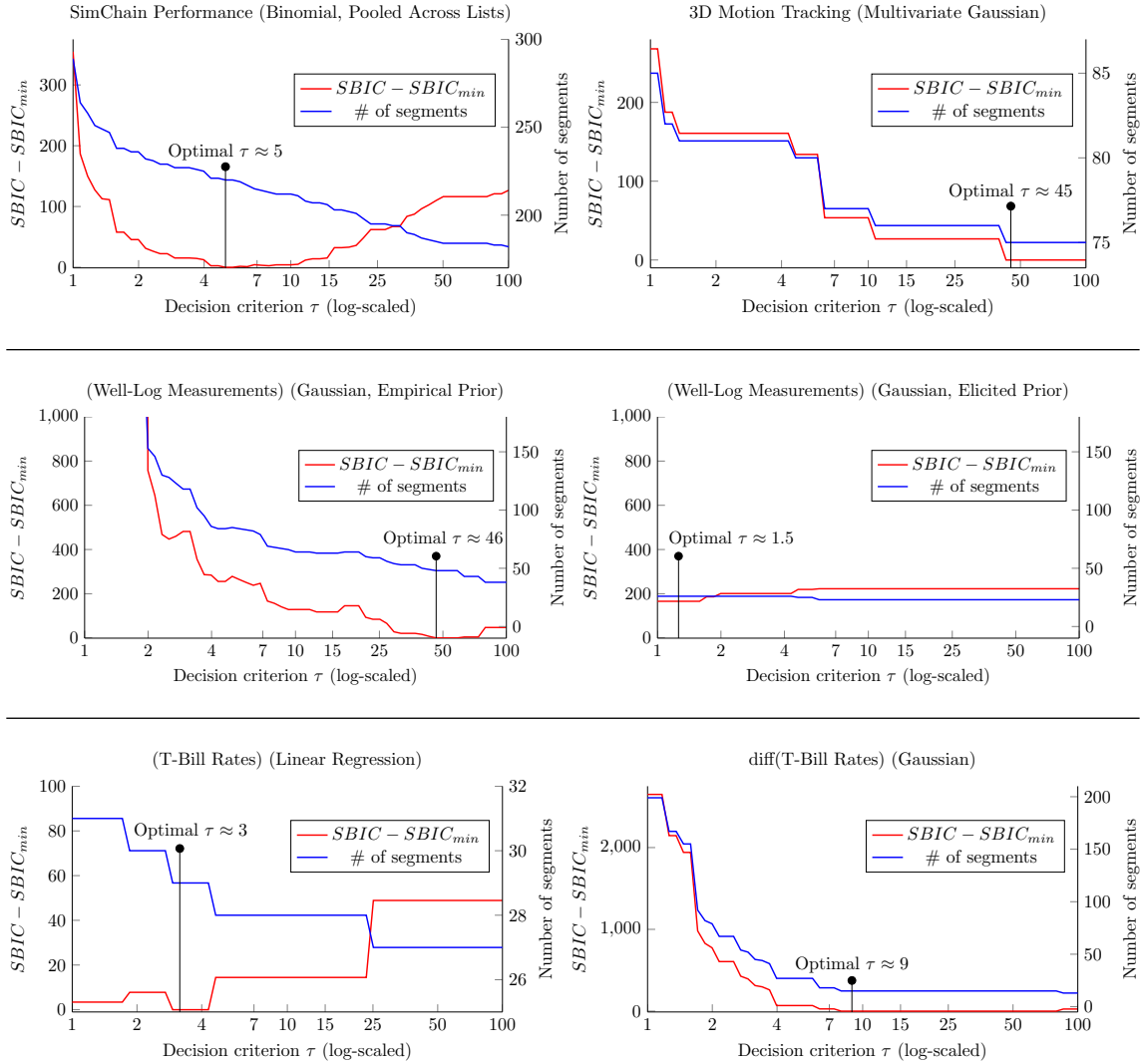


Figure 4. Sensitivity analysis performed on two datasets from the main body text and three from the supplement. The decision criterion τ was varied between 1 and 100, and the SBIC was calculated for each resulting model, using Equation 3. Relative SBIC among the models is plotted in red (with the best model set at zero), and the number of data segments is plotted in blue.

1 lower right) both demonstrate the risk of overfitting when the decision criterion is too lax.
 2 Additionally, the SimChain data show a risk of underfitting if the criterion is too harsh.

3 A different issue is raised by the sensitivity analysis performed on the well-log mea-
 4 surements. Figure 4 (center left) shows the results of the sensitivity analysis when the
 5 empirical ‘rule-of-thumb’ prior is used (hyperparameters $a = 114000$, $b = 0.1$, $c = 0.1$,
 6 $d = 5930.4$). However, these data do not conform well to the assumptions of a Gaussian
 7 distribution (due to device imprecision), and a different prior was used in the analysis de-
 8 scribed above (hyperparameters $a = 114000$, $b = 0.1$, $c = 1$, $d = 2965.2$). Figure 4 (center
 9 right) shows the results of using this ‘elicited’ prior (so-called because its selection depends
 10 on the analyst’s discretion about the meaning of the data). When using the rule-of-thumb
 11 prior (which is weakly subjective), a great deal of sensitivity to the decision criterion was
 12 observed. Contrastingly, when a prior with a stronger prior assumption about the disper-
 13 sion of the data was used, the influence of the decision criterion was much smaller. The
 14 sensitivity of posterior results to the prior is a well-established problem, and one of the
 15 central pillars of the argument in favor of objective Bayesian methods (Samaniego, 2012).
 16 Responsible use of subjective priors should include a sensitivity analysis with respect to the
 17 prior (Gelman et al., 2003).

18 Across these six sensitivity analyses, $\tau = 10$ performs reasonably well, neither result-
 19 ing in dramatic over- or under-fitting of the data. Although performing sensitivity analyses
 20 is an important element of statistical due diligence, $\tau = 10$ is a reasonable default value for
 21 the decision criterion.

22 **A Failure Condition: CPR Insensitivity When Data Are Stationary**

23 Although intended for large datasets, the sensitivity of the CPR algorithm to change-
 24 points is reduced in cases where changes occur frequently and values oscillate around some
 25 central value. Closed-form, conjugate solutions for $m(x, M)$ depend only on the prior hy-
 26 perparameters and the data’s sufficient statistics, both of which usually ignore the temporal
 27 structure of values within each segment (that is, they assume ‘exchangeability’). Thus, seg-
 28 ments containing many changes that oscillate around a stable central value are likely to be
 29 interpreted as a stationary process with an inflated variance.

30 This problem can be demonstrated using a simulated dataset. A “Gaussian iterated
 31 map” (or “mouse map”) is a nonlinear iterated function with two parameters α and β :

$$x_{(i+1)} = \beta + \exp\left(-\alpha x_{(i)}^2\right) \tag{4}$$

This function’s behavior is chaotic given certain parameters (Hilborn, 2001), making it
 useful for generating replicable datasets whose “true” parameters can be compared to those
 estimated. We defined two mouse maps:

$$\begin{aligned} x_{(i+1)} &= -0.6 + \exp\left(-2\pi x_{(i)}^2\right); x_{(1)} = 0 \\ y_{(i+1)} &= -0.4 + \exp\left(-e^2 y_{(i)}^2\right); y_{(1)} = 1 \end{aligned}$$

These, in turn, were used to specify a series of parameters:

$$\begin{aligned} \mu_{(i)} &= 4x_{(i)} \\ \sigma_{(i)} &= \exp\left(y_{(i)}\right) \end{aligned}$$

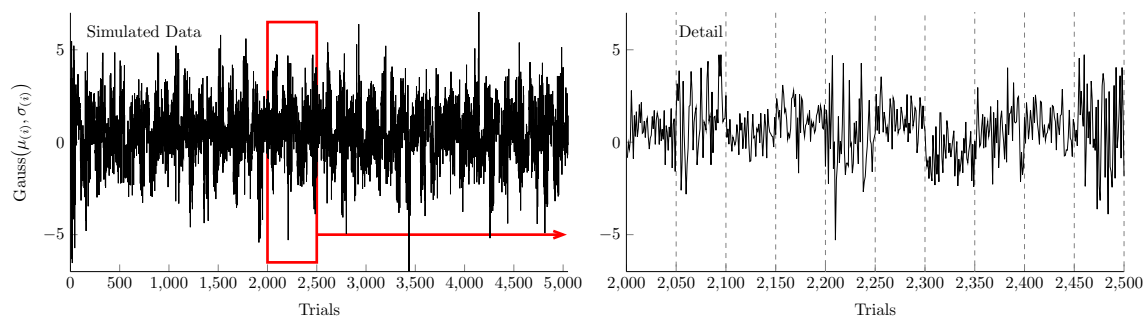


Figure 5. Simulated data generated using a chaotic distribution of Gaussian parameters derived from Equation 4. (Left) The full dataset D , consisting of 5050 simulated observations that were generated in blocks of 50. The default CPR algorithm will not identify most of these change-points. (Right) A subset of D , with the true change-points denoted by dashed gray lines. The CPR algorithm will reliably identify these change-points if it is run only on this subset, despite failing to detect them when run on the complete dataset.

1 The dataset D consisted of 101 segments, where the i th segment was generated using a
 2 Gaussian distribution with parameters $(\mu_{(i)}, \sigma_{(i)})$. A characteristic example is presented
 3 in Figure 5, which presents both the full data (Left) and a detail view of a section of the
 4 data (Right). When the true change-points are clearly labeled, their subdivisions appear
 5 highly plausible, but the CPR algorithm is not omniscient, and these relatively subtle
 6 discontinuities are subsumed by the overall variance of the data when too many of them
 7 occur within a particular segment.

8 The Dicing Operation

9 A strictly exploratory solution to this problem is to introduce the preliminary step
 10 of ‘dicing’ the data into arbitrary segments and examining each segment for change-points,
 11 outlined in Algorithm 1.

12 When invoking the dicing operation, the analyst must specify some integer value $d > 1$
 13 that specifies how many times to subdivide the data, doing so into equal parts. Each segment
 14 is checked for a single change-point, and any change-points that are identified are added to
 15 the initial model \mathbf{M} . After performing the dicing operation, the CPR Algorithm proceeds
 16 with the new model \mathbf{M} , rather than its default value of $\langle 0, n \rangle$. This effectively bypasses
 17 the first few partitioning operations (during which the sensitivity might be impaired) and
 18 applies the CPR algorithm to an already-subdivided dataset.

19 Figure 6 (top) demonstrates the efficacy of introducing the dicing operation to the
 20 data in Figure 5. Given one subdivision (the default for the CPR algorithm), no change-
 21 points are identified between $t = 150$ and $t = 5000$. However, with only three subdivisions,
 22 the number of change-points identified rises to 24. Given four or more subdivisions, the
 23 change-points identified were highly consistent, with difficulty only the range from 4400 to
 24 5000. Figure 6 (center) shows the total number of times different points were identified over
 25 the thirty different dicing operations. Only 284 of the 5049 possible points were identified
 26 at least once. Of those, 27 points were specified at least twenty-five times, and 62 were
 27 specified at least twenty times, without even considering ‘near misses.’

Algorithm 1: The dicing operation, which identifies the single best change-point in each arbitrary segment of data.

Data: events $x_{(1:n)}$, times $t_{(0:n)}$, model C_0 , decision criterion τ , diced segment count d

Result: diced model \mathbf{M}

```

begin
   $\mathbf{S} \leftarrow 0, \frac{n}{d}, \frac{2n}{d}, \dots, \frac{dn}{d}$  /* get segment fenceposts */
   $\mathbf{M} \leftarrow \langle 0, n \rangle; p_c = \frac{1}{n-1}$  /* empty "new CPs" array and prior odds of a
  change */
  for  $s = 1$  to  $\text{length}(\mathbf{S}) - 1$  do
     $i \leftarrow \mathbf{S}(s) + 1; j \leftarrow \mathbf{S}(s + 1); |\mathbf{K}| \leftarrow j - i + 1$  /* assign indices; allocate
    BF array */
    for  $c = i$  to  $j$  do
       $k_c \leftarrow \frac{m(x_{(i:c-1)}, C_0) m(x_{(c:j)}, C_0)}{m(x_{(i:j)}, C_0)}; \mathbf{K}_{(c)} \leftarrow \frac{k_c \cdot (t_{(c)} - t_{(c-1)})}{(t_{(j)} - t_{(i-1)}) \cdot \exp(SB_{(c)})}$  /* Bayes
      factors */
    if  $\text{sum}(\mathbf{K}) \cdot p_c \cdot \frac{1}{j-i+1} > \tau$  then
       $\hat{c} \leftarrow \text{index}(\max(\mathbf{K})); \text{Push}(\mathbf{M}, \hat{c})$  /* insert  $\hat{c}$  into  $\mathbf{M}$  if  $\tau$  permits
      */
  Sort( $\mathbf{M}$ ); return  $\mathbf{M}$ 

```

1 Although the dicing operation has practical use when detecting stationary meta-
 2 stable periods within a more generally stable process, it should not be necessary in the vast
 3 majority of cases. Data should not be diced into segments so small that no more than one
 4 or two change-points can reasonably be expected to exist.

5 Forward-Retrospective Change-Point Detection

6 Another approach to the problem of oscillating data is to identify each change-point in
 7 chronological order. Rather than approaching the analysis as a batch process, change-point
 8 evaluations can take place within a more constrained frame.

9 Sequential change-point detection algorithms have a long history, with the work of
 10 Page (1954) being among the earliest (see also Venkatraman, 1992; Chen & Gopalakrishnan,
 11 1998; Gallistel et al., 2004). Beginning with the indices $i = 1$ and $j = 2$, the data range
 12 $x_{(i:j)}$ is tested for a change-point repeatedly as j is incrementally increased. Once a change
 13 is detected at \hat{c} (where $i < \hat{c} < j$), the data prior to \hat{c} are removed from consideration and
 14 $i = \hat{c}$. This process continues until all observations have been considered.

15 Thus, when a change is detected, the range of data subsequently considered is trun-
 16 cated. Although a deliberate feature in the interests of reducing the algorithm's compu-
 17 tational load, this creates problems when 'false alarms' are observed. Over a very long
 18 sequence of Gaussian observations, for example, a handful of sizable outliers are reasonable;
 19 when only a small segment of data is being observed, however, estimates of variance are
 20 likely to be conservative and a random outlier may be confused for a true change in the
 21 distribution. This problem is particularly salient in data that suffer from some degree of

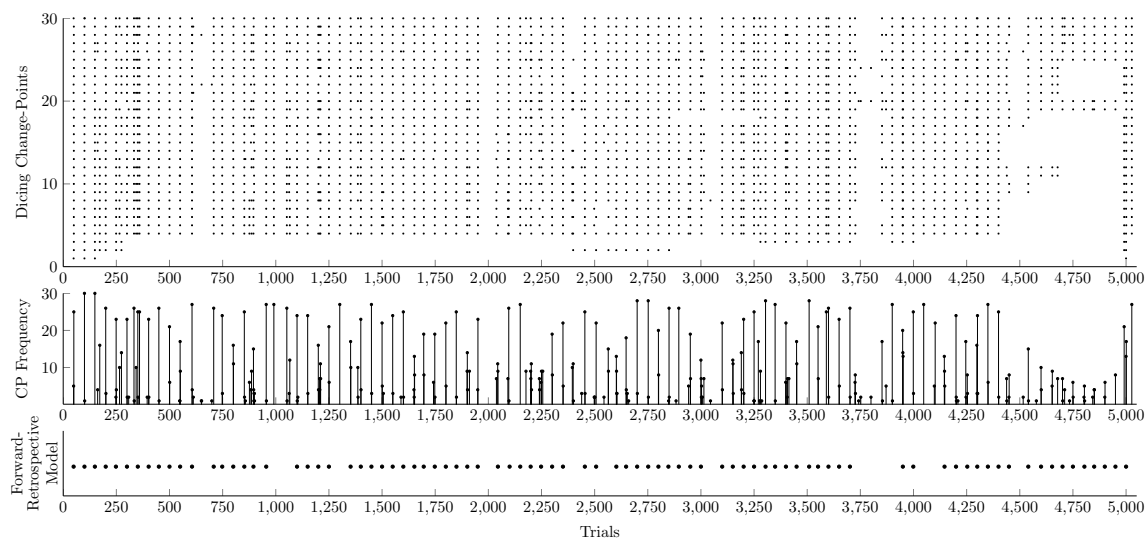


Figure 6. Change-points identified in simulated data. (Top) Trials for which change-points were identified, as a function of the number of subdivisions examined by the dicing operation. (Center) A histogram of how often, across 30 values for the dicing parameter, a precise change-point was identified. (Bottom) Trials for which change-points were identified using the retrospective sequential algorithm.

1 out-of-sample ‘contamination.’ Furthermore, even in cases where the data are well-behaved,
 2 sequential analyses are vulnerable to mistaken inference because of multiple comparisons.
 3 When a true change is identified, there is a risk that its position will be estimated prema-
 4 turely.

5 In order to combat this, a substantial improvement to traditional sequential analyses
 6 is proposed by Gallistel et al. (Submitted). After the first change-point is identified, each
 7 additional change-point’s identification is immediately followed by a retrospective test of
 8 the preceding change-point. For example, once $\hat{c}_{(3)}$ is identified, the interval $x_{\hat{c}_{(1)}:\hat{c}_{(3)}}$ is
 9 reanalyzed. If $\hat{c}_{(2)}$ was the result of a false alarm, it will probably not be detected again,
 10 and it can be removed from the model. If, on the other hand, the evidence still supports a
 11 change between $\hat{c}_{(1)}$ and $\hat{c}_{(3)}$, this re-analysis can be used to update the position of $\hat{c}_{(2)}$.

12 The forward-retrospective change-point analysis proposed by Gallistel et al. (Sub-
 13 mitted) is combined with the CPR algorithm in Algorithm 2. Here, rather than making
 14 use of recursive binary partitioning, the CPR algorithm is only invoked to find a single
 15 change-point at a time.

16 Figure 6 (bottom) shows the change-point model reported by the retrospective se-
 17 quential algorithm when analyzing the simulated data from Figure 5. Whereas the CPR
 18 algorithm displays considerable insensitivity in the absence of the dicing operation, the
 19 forward-retrospective analysis identifies nearly all changes, doing so close to their true times
 20 of occurrence.

21 Although an effective alternative in many circumstances, the forward-retrospective
 22 approach also has shortcomings. Both the standard CPR algorithm and its sequential
 23 counterpart have expected runtimes that scale as a function both of the number of data

Algorithm 2: The forward-retrospective change-point algorithm, which scans forward through the data in order to identify each change-point in chronological order

Data: events $x_{(1:n)}$, parameters \mathbf{P}

Result: change-point model \mathbf{M}

```

begin
   $\mathbf{M} \leftarrow \langle 0, n \rangle; s \leftarrow 1$            /* empty "new CPs" array, segment index */
  for  $i = 2$  to  $n$  do
     $\mathbf{N} \leftarrow \text{CPRBayes}(x_{(\mathbf{M}(s)+1:i)}, \mathbf{P})$    /* check for single cp */
    if  $\text{length}(\mathbf{N}) > 2$  then
       $\text{Push}(\mathbf{M}, \mathbf{M}(s) + \mathbf{N}(2))$            /* insert cp into  $\mathbf{M}$  */
       $\text{Sort}(\mathbf{M}); s = s + 1$            /* sort  $\mathbf{M}$ , update  $s$  */
      if  $s > 2$  then
         $\mathbf{O} \leftarrow \text{CPRBayes}(x_{(\mathbf{M}(s-2)+1:\mathbf{M}(s))}, \mathbf{P})$  /* check previous cp */
        if  $\text{length}(\mathbf{O}) > 2$  then
           $\mathbf{M}_{(s-1)} = \mathbf{M}_{(s-2)} + \mathbf{O}(2)$  /* refine previous cp */
        else
           $\text{Pull}(\mathbf{M}_{(s-1)}); s = s - 1$  /* remove previous cp */
    return  $\mathbf{M}$ 

```

1 segments s and of the number of observations n . When using a binary partitioning strategy,
 2 the CPR algorithm’s runtime is expected to fall between $O(n \cdot s)$ and $O(n \cdot \log s)$ (depend-
 3 ing on the order in which the points are identified). However, the sequential alternative has
 4 a runtime between $O((n - s)^2)$ and $O(\frac{n^2}{s})$, depending on how evenly the change-points
 5 are spaced. When changes are infrequent, the forward-retrospective algorithm is potentially
 6 much slower.

7 Table 1 gives a sense of these contrasting relationships. When change-points occur
 8 relatively frequently (as in the 3D motion tracking or the T-Bill rate datasets), the sequential
 9 approach only performs somewhat less well. However, when there are very few changes
 10 relative to the number of observations (as is particularly the case for the reaction times data),
 11 the increase in runtime can be dramatic. The quadratic runtime renders the sequential
 12 approach unusable for large datasets with large gaps between change-points.

13 One way to mitigate the quadratic runtime is to increment the index i by more than
 14 one observation at a time (Chen & Gopalakrishnan, 1998). If change-points are consis-
 15 tently widely spaced, then the value of i could conceivably increase by tens or hundreds
 16 of observations with each step. In practice, however, consistently wide spacing renders the
 17 sequential approach unnecessary, while inconsistent spacing (in which some change-points
 18 are spaced much closer together than others) introduces the risk that large increments of i
 19 will overlook some changes.

20 Another difficulty with the forward-retrospective strategy is its vulnerability to the
 21 stopping problem. In principle, this is mitigated to some extent by the retrospective oper-

Dataset	Binary Partitioned	Forward-Retrospective
SimChain	2.84 s	24.80 s
Reaction Time	21.23 s	3756.04 s
3D Motion Tracking	317.77 s	527.57 s
Mining Disasters (Exponential)	0.35 s	4.79 s
Mining Disasters (Poisson)	0.14 s	1.62 s
Well-Log Measurements	15.08 s	169.43 s
T-Bill Rates	4.45 s	15.70 s
diff(T-Bill Rates)	1.78 s	9.25 s
Simulation Data	(diced) 23.15 s	71.72 s

Table 1

Processing time for the CPR algorithm in each of the datasets described in the main text and supplement, performed using a dual-core 2.93GHz laptop.

1 ation, but the forward-retrospective strategy will, all things being equal, identify a larger
 2 number of change-points than the CPR strategy, because it is more likely to be impacted
 3 by small subsets of the data. As a counterweight against this, the forward-retrospective
 4 algorithm uses a default criterion of $\tau = 20$, twice the default recommended for the CPR
 5 algorithm.

6 Finally, the forward-retrospective algorithm yields different results examining the data
 7 forward and backward in time. This lack of symmetry arises from several factors, but the
 8 most substantive is the process of updating p_c as additional change-points are detected.
 9 This results in a growing propensity to identify change-points over time.

10 On the basis of these limitations, there are two scenarios in which the sequential
 11 approach should consistently be favored. The first is when changes are expected to occur
 12 relatively frequently, particularly if they occur in the oscillating manner characterized by the
 13 simulated data in Figure 5. The second is when temporal order is of primary importance,
 14 such as when the earliest change can be said to have occurred (e.g. the first trial in which
 15 an animal engaged in conditioned responding to a stimulus).

16 Closed Form Solutions for Conjugate Priors

17 This section provides closed-form solutions for the conjugate priors and marginal
 18 model likelihoods of various distributions that are likely to be of interest. These are used
 19 to compute the values of $m(x, M)$ that appear in the change-point algorithm. This list is
 20 not meant to be exhaustive. Unless otherwise noted, the conjugacy of these distributions is
 21 demonstrated by DeGroot (2005).

22 As detailed in the body text, the prior probability distribution (and the corresponding
 23 hyperparameters that describe its shape) can have a determining effect in Bayesian analysis,
 24 and must be selected in a principled manner. Whenever possible, the hyperparameters
 25 should reflect an analyst’s principled assumptions, and they should be explicitly reported
 26 along with a justification for their use. For example, if data were collected using equipment
 27 with a known degree of imprecision, this should be reflected by the prior. For each of the
 28 distributions below, the interpretation of the hyperparameters is provided to facilitate this
 29 process.

1 In some cases, an analyst’s intuitions must arise from the data themselves. In such ex-
 2 ploratory analyses, setting an inappropriate prior will not only result in misleading marginal
 3 likelihoods, but can also severely impair the sensitivity of the CPR algorithm. To the extent
 4 possible, the “rule-of-thumb” priors provided for each distribution are intended to facilitate
 5 the identification of change-points in keeping with an empirical Bayes approach (Casella,
 6 1985; Carlin & Louis, 2000). Insofar as their use deviates from standard Bayesian inference,
 7 an analyst choosing to use them should explicitly note cases where the prior hyperparame-
 8 ters are derived from the data.

9 By preference, the rule-of-thumb priors make use of robust estimation techniques
 10 (Huber & Ronchetti, 2009). Time series with substantial discontinuities will, almost by
 11 definition, appear to contain ‘outliers’ when viewed without regard for time, so robust
 12 estimators (such as the median and median absolute deviation) are less likely to suggest
 13 parameters that are not characteristic of the preponderance of the data.

	Notation	Formula	Interpretation
	$\Gamma(x)$	$(x-1)!$	Gamma Function
14	Beta (x, y)	$\frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)}$	Beta Function
	SumSq (x, y)	$(x-y)(x-y)^\top$	Sum of squared differences
	MAD (x)	$\text{median}(x - \text{median}(x))$	Median absolute deviation

15 **Discrete Data (Univariate)**

16 **The Binomial Distribution.**

- 17 • Support:

$$\{x_1, \dots, x_n | x_i = 0 \text{ or } 1\}$$

- 18 • Sufficient statistics:

$$n \text{ observations} \quad S_x = \sum_{i=1}^n x_i$$

- 19 • With unknown probability of success p :

20 Prior Hyperparameters

- 21 $\{a | a > 0\}$ Number of successes-1
 $\{b | b > 0\}$ Number of failures-1

22 Conjugate Prior

$$\Pr(p | a, b) = \frac{p^{a-1} (1-p)^{b-1}}{\text{Beta}(a, b)}$$

23 Normalizing Constant

$$m(a, b) = \text{Beta}(a, b) \tag{5}$$

24 Posterior Hyperparameters

$$a' = a + S_x \quad b' = b + n - S_x$$

25 Posterior Predictive

$$\hat{p} = \frac{a'}{a' + b'}$$

1 'Rule-of-Thumb' Empirical Prior

$$a = 0.5 \quad b = 0.5$$

2 **The Geometric Distribution.**

3 • Support:

$$\{x_1, \dots, x_n | x_i \in (0, 1, 2, 3, \dots)\}$$

4 • Sufficient statistics:

$$n \text{ observations} \quad S_x = \sum_{i=1}^n x_i$$

5 • With unknown probability of success p :

6 Prior Hyperparameters

7 $\{a | a > 0\}$ Number of observations
 $\{b | b > 0\}$ Sum of observations

8 Conjugate Prior

$$\Pr(p | a, b) = \frac{p^{a-1} (1-p)^{b-1}}{\text{Beta}(a, b)}$$

9 Normalizing Constant

$$m(a, b) = \text{Beta}(a, b) \tag{6}$$

10 Posterior Hyperparameters

$$a' = a + n \quad b' = b + S_x$$

11 Posterior Predictive

$$\hat{p} = \frac{a' + n}{a' + b' + n + S_x}$$

12 'Rule-of-Thumb' Empirical Prior

$$a = 0.5 \quad b = 0.5$$

13 **The Poisson Distribution.**

14 • Support:

$$\{x_1, \dots, x_n | x_i \in (0, 1, 2, 3, \dots)\}$$

15 • Sufficient statistics:

$$n \text{ observations} \quad S_x = \sum_{i=1}^n x_i$$

16 • With unknown rate λ :

17 Prior Hyperparameters

18 $\{a | a > 0\}$ Total occurrences
 $\{b | b > 0\}$ Number of intervals

1 Conjugate Prior

$$\Pr(\lambda|a, b) = \frac{b^a}{\Gamma(a)} \frac{\lambda^{a-1}}{\exp(b\lambda)}$$

2 Normalizing Constant

$$m(a, b) = \frac{\Gamma(a)}{b^a} \tag{7}$$

3 Posterior Hyperparameters

$$a' = a + S_x \quad b' = b + n$$

4 Posterior Predictive

$$\hat{\lambda} = \frac{a'}{b'}$$

5 ‘Rule-of-Thumb’ Empirical Prior

$$a = \text{median}(x) \quad b = 1$$

6 **Continuous Data (Univariate)**

7 **The Exponential Distribution.**

8 • Support:

$$\{x_1, \dots, x_n | x_i \in \mathbb{R}^{>0}\}$$

9 • Sufficient statistics:

$$n \text{ observations} \quad S_x = \sum_{i=1}^n x_i$$

10 • With unknown θ :

11 Prior Hyperparameters

12 $\{a|a > 0\}$ Number of observations
 $\{b|b > 0\}$ Sum of observations

13 Conjugate Prior

$$\Pr(\theta|a, b) = \frac{b^a}{\Gamma(a)} \frac{\theta^{a-1}}{\exp(b\theta)}$$

14 Normalizing Constant

$$m(a, b) = \frac{\Gamma(a)}{b^a} \tag{8}$$

15 Posterior Hyperparameters

$$a' = a + n \quad b' = b + S_x$$

16 Posterior Predictive

$$\hat{\theta} = \frac{b'}{a'}$$

17 ‘Rule-of-Thumb’ Empirical Prior

$$a = 1 \quad b = \text{median}(x)$$

1 **The Gaussian Distribution.**

2 The conjugate analysis in this section, originally outlined by DeGroot (2005), is based on
 3 the excellent treatment of the subject by Murphy (2007).

4 • Support:

$$\{x_1, \dots, x_n | x_i \in \mathbb{R}\}$$

5 • Sufficient statistics:

$$n \text{ observations} \quad S_x = \sum_{i=1}^n x_i \quad \bar{x} = \frac{S_x}{n}$$

6 • With unknown mean μ and known precision λ :

7 Prior Hyperparameters

$$\begin{aligned} \{a | a \in \mathbb{R}\} & \quad \text{Prior mean} \\ \{b | b > 0\} & \quad \text{Prior total precision} \end{aligned}$$

9 Conjugate Prior

$$\Pr(\mu | a, b) = \sqrt{\frac{b}{2\pi}} \exp\left(-\frac{b}{2}(\mu - a)^2\right)$$

10 Normalizing Constant

$$m(a, b) = \sqrt{\frac{2\pi}{b}} \exp\left(\frac{a^2 b}{2}\right) \tag{9}$$

11 Posterior Hyperparameters

$$a' = \frac{ab + \lambda S_x}{b + \lambda n} \quad b' = b + \lambda n$$

12 Posterior Predictive

$$\hat{\mu} = a'$$

13 ‘Rule-of-Thumb’ Empirical Prior

$$a = \text{median}(x) \quad b = 0.1$$

14 • With known mean μ and unknown precision λ :

15 Prior Hyperparameters

$$\begin{aligned} \{a | a > 0\} & \quad (\text{Number of observations})/2 \\ \{b | b > 0\} & \quad (\text{Sum of squared deviations from } \mu)/2 \end{aligned}$$

17 Conjugate Prior

$$\Pr(\lambda | a, b) = \frac{b^a}{\Gamma(a)} \lambda^{a-1} \exp(-b\lambda)$$

18 Normalizing Constant

$$m(a, b) = \frac{\Gamma(a)}{b^a} \tag{10}$$

1 Posterior Hyperparameters

$$a' = a + \frac{n}{2} \quad b' = b + \frac{\text{SumSq}(x, \mu)}{2}$$

2 Posterior Predictive

$$\hat{\lambda} = \sqrt{\frac{b'}{a'}}$$

3 ‘Rule-of-Thumb’ Empirical Prior

$$a = 0.1 \quad b = \text{MAD}(x) \cdot 1.4826$$

4 • With unknown mean μ and unknown precision λ :

5 Prior Hyperparameters

- 6
- $\{a|a \in \mathbb{R}\}$ Prior mean
 - $\{b|b > 0\}$ Observations supporting the prior mean
 - $\{c|c > 0\}$ (Observations supporting the prior precision)/2
 - $\{d|d > 0\}$ (Sum of squared deviations)/2

7 Conjugate Prior

$$\Pr(\mu, \lambda|a, b, c, d) = \frac{d^c \lambda^{c-1}}{\Gamma(c) \exp(\lambda d)} \cdot \sqrt{\frac{\lambda b}{2\pi}} \cdot \exp\left(-\frac{\lambda b}{2}(\mu - a)^2\right)$$

8 Normalizing Constant

$$m(a, b, c, d) = \frac{\Gamma(c)}{d^c} \sqrt{\frac{2\pi}{b}} \tag{11}$$

9 Posterior Hyperparameters

$$a' = \frac{ab + S_x}{b + n} \quad b' = b + n \quad c' = c + \frac{n}{2}$$

10

$$d' = d + \frac{\text{SumSq}(x, \bar{x})}{2} + \frac{bn(\bar{x} - a)^2}{2(b + n)}$$

11 Posterior Predictive

$$\hat{\mu} = a' \hat{\lambda} = \sqrt{\frac{d' \cdot (b' + 1)}{b' \cdot c'}}$$

12 ‘Rule-of-Thumb’ Empirical Prior

$$a = \text{median}(x) \quad b = 0.1$$

13

$$c = 0.1 \quad d = \text{MAD}(x) \cdot 1.4826$$

Basic Linear Regression.

Linear regression presents an analyst with a hybrid problem, because the dependent observations y are being summarized in terms of a multivariate vector of regression coefficients β and a univariate variance parameter σ^2 . As with the Gaussian distribution, these parameters can each be understood in terms of their respective priors distributions. The implementation below uses the multivariate Gaussian distribution as the conjugate prior for the coefficients β , and the inverse gamma distribution as the conjugate prior for the variance.

This implementation is also distinctive in that the most arithmetically tractable form for the normalizing constant incorporates both the prior and posterior hyperparameters.

• Support:

$$\{y_1, \dots, y_n | y_i \in \mathbb{R}\} X = \begin{bmatrix} X_{1,1} & \dots & X_{1,k} \\ \vdots & \ddots & \vdots \\ X_{n,1} & \dots & X_{n,k} \end{bmatrix}, X_{j,i} \in \mathbb{R}$$

• With unknown regression coefficients β and unknown variance σ^2 :

Prior Hyperparameters

model = $m(\mathbf{m}, a, b, \mathbf{\Lambda})$	Gaussian error model
$\mathbf{m} = \{m_1, \dots, m_k m_i \in \mathbb{R}\}$	Prior regression coefficients
$\mathbf{\Lambda} = \mathbf{\Sigma}^{-1}$	Prior coefficient precision matrix
$\{a a > 0\}$	(Observations supporting the error)/2
$\{b b > 1\}$	(Sum of squared deviations)/2

Conjugate Prior

$$\Pr(\beta, \sigma^2 | y, X, \text{model}) = \frac{\Pr(\beta, \sigma^2 | \text{model}) \cdot \Pr(y | X, \beta, \sigma^2, \text{model})}{\Pr(y | \text{model})}$$

$$\text{Given the model, } \Pr(\beta, \sigma^2 | y, X) \propto \Pr(\beta | \sigma^2, y, X) \cdot \Pr(\sigma^2 | y, X)$$

$$\Pr(\beta | \sigma^2, y, X) = \mathcal{N}(\mathbf{m}, \sigma^2 \mathbf{\Lambda}^{-1}) \Pr(\sigma^2 | y, X) = \text{InvGamma}(a, b)$$

Normalizing Constant

$$\Pr(y | \text{model}) = \frac{1}{(2\pi)^{n/2}} \cdot \sqrt{\frac{|\mathbf{\Lambda}|}{|\mathbf{\Lambda}'|}} \cdot \frac{b^a}{b^{a'}} \cdot \frac{\Gamma(a')}{\Gamma(a)} \tag{12}$$

Posterior Hyperparameters

$$\mathbf{m}' = (X^\top X + \mathbf{\Lambda})^{-1} \cdot (\mathbf{\Lambda} \mathbf{m} + X^\top y) \quad a' = a + \frac{n}{2}$$

$$\mathbf{\Lambda}' = X^\top X + \mathbf{\Lambda} \quad b' = b + \frac{1}{2} (y^\top y + \mathbf{m}^\top \mathbf{\Lambda} \mathbf{m} - \mathbf{m}'^\top \mathbf{\Lambda}' \mathbf{m}')$$

Posterior Predictive

$$\hat{\beta} = \mathbf{m}' \quad \hat{y} = X^\top \hat{\beta}$$

1

$$\hat{\sigma}^2 = \frac{\text{SumSq}(y, \hat{y})}{n}$$

2

‘Rule-of-Thumb’ Empirical Prior

$$a = \text{regress}(x) \quad b = 1$$

3

$$c = 1 \quad d = \frac{X^\top \cdot X}{n^2}$$

4

The Uniform Distribution.

5 Applying the CPR algorithm to the continuous uniform distribution is a difficult because of
 6 how it responds to virtual observations. The sufficient statistics for a uniform distribution
 7 are the number of observations n , the minimum x_{mn} , and the maximum x_{mx} . Given n
 8 observations from a uniform distribution, each additional observation x_{n+1} always increases
 9 n , but otherwise it *only* impacts the shape of the distribution when it falls outside the range
 10 of $[x_{\text{mn}}, x_{\text{mx}}]$ and changes the position of one of the endpoints. Thus, a “virtual” observation
 11 provided by a prior has very little effect on the distribution when $x_{\text{mn}} < x_{n+1} < x_{\text{mx}}$ but
 12 has a very powerful effect otherwise. This confounds the principle that one should use a
 13 ‘weak’ prior, because setting prior minimum and maximum values for the data either has a
 14 negligible effect, or a massive one, with no intermediary scenario.

15 When both the minimum λ_α and maximum λ_ω are unknown, setting a prior for the
 16 corresponding hyperparameters b_1 and b_2 has a determining effect on the viable range of
 17 values. This is reasonable in special cases where the system being modeled is very well
 18 defined, but in almost all empirical scenarios, the role played by the prior is to make an
 19 inference about the *support* of the distribution, which corresponds to a uniform step whose
 20 width is $(\lambda_\omega - \lambda_\alpha)$ and whose position is arbitrary. As such, an alternate parameterization
 21 of the conjugate prior only requires that the analyst specify prior observations a and the
 22 minimum allowable range d , which describes minimum variance (i.e. the maximum pre-
 23 cision) of the measurement. For example, if observations are rounded to the nearest tens
 24 place, then $d = 10$ is an appropriate value.

25 In the more straightforward case that the minimum is known *a priori* to be 0.0, the
 26 prior maximum b and the allowable range d are identical. The likelihoods in any case
 27 when either the minimum or maximum is known can be reduced to this simple case by
 28 repositioning the origin to the known parameter and (in the case of a known maximum)
 29 reversing the resulting signs.

- 30 • Support:

$$\{x_1, \dots, x_n | x_i \in \mathbb{R}\}$$

- 31 • Sufficient statistics:

$$n \text{ observations} \quad x_{\text{mn}} = \min(x) \quad x_{\text{mx}} = \max(x)$$

- 32 • With known minimum 0.0 and unknown maximum λ_ω :

33

Prior Hyperparameters

34

$$\begin{aligned} \{a | a > 0\} & \text{ Number of observations} \\ \{b | b > 0\} & \text{ Minimum width of the interval} \end{aligned}$$

1 Conjugate Prior

$$\Pr(\lambda_\omega | a, b) = \frac{ab^a}{(\lambda_\omega)^{a+1}}$$

2 Normalizing Constant

$$m(a, b) = \frac{1}{ab^a} \tag{13}$$

3 Posterior Hyperparameters

$$a' = a + n \quad b' = \max(b, x_{\text{mx}})$$

4 Posterior Predictive

$$\hat{\lambda}_\omega = \max(b, x_{\text{mx}})$$

5 ‘Rule-of-Thumb’ Empirical Prior

$$a = 1 \quad b = \text{minimum nonzero value of diff}(\text{sort}(x))$$

6 • With unknown minimum λ_α and unknown maximum λ_ω :

7 Prior Hyperparameters

$$\begin{array}{ll} \{a|a > 0\} & \text{Number of observations} \\ \{b_1, b_2|b_1, b_2 \in \mathbb{R}\} & \text{Minimum and maximum values} \end{array}$$

8 OR

$$\begin{array}{ll} \{a|a > 0\} & \text{Number of observations} \\ \{d|d > 0\} & \text{Minimum non-zero value for } \lambda_\omega - \lambda_\alpha \end{array}$$

9 Conjugate Prior

$$\Pr(\lambda_\alpha, \lambda_\omega | a, b_1, b_2) = \frac{a(a+1)(b_2 - b_1)^a}{(\lambda_\omega - \lambda_\alpha)^{a+2}} = \frac{a(a+1)(d)^a}{(\lambda_\omega - \lambda_\alpha)^{a+2}}$$

10 Normalizing Constant

$$m(a, b_1, b_2) = \frac{1}{a(a+1)(b_2 - b_1)^a} = \frac{1}{a(a+1)(d)^a} \tag{14}$$

11 Posterior Hyperparameters

$$\begin{array}{lll} a' = a + n & b'_1 = \min(b_1, x_{\text{mn}}) & b'_2 = \max(b_2, x_{\text{mx}}) \\ & d' = \max(d, x_{\text{mx}} - x_{\text{mn}}) & \end{array}$$

13 Posterior Predictive

$$\begin{array}{ll} \hat{\lambda}_\alpha = b'_1 & \hat{\lambda}_\omega = b'_2 \\ \hat{\lambda}_\alpha = \min\left(\left(\frac{x_{\text{mn}} + x_{\text{mx}} - d'}{2}\right), x_{\text{mn}}\right) & \hat{\lambda}_\omega = \max\left(\left(\frac{x_{\text{mn}} + x_{\text{mx}} + d'}{2}\right), x_{\text{mx}}\right) \end{array}$$

15 ‘Rule-of-Thumb’ Empirical Prior

$$a = 1 \quad d = \text{minimum nonzero value of diff}(\text{sort}(x))$$

$$b_1 = N/A \quad b_2 = N/A$$

1 **Discrete Data (Multivariate)**

2 **The Multinomial Distribution.**

3 • Support:

$$\mathbf{x} = \begin{bmatrix} x_{1,1} & \dots & x_{1,k} \\ \vdots & \ddots & \vdots \\ x_{n,1} & \dots & x_{n,k} \end{bmatrix}, x_{j,i} = 0 \text{ or } 1$$

4 • With unknown vector \mathbf{w} :

5 Prior Hyperparameters

6 $\mathbf{a} = \{a_1, \dots, a_k | a_i > 0\}$ Occurences per category–1

7 Conjugate Prior

$$\Pr(\mathbf{w}|\mathbf{a}) = \frac{\Gamma(\sum a_i)}{\prod(\Gamma(a_i))} \prod w_i^{a_i-1}$$

8 Normalizing Constant

$$m(\mathbf{a}) = \frac{\prod(\Gamma(a_i))}{\Gamma(\sum a_i)} \tag{15}$$

9 Posterior Hyperparameters

$$\mathbf{a}' = \{a'_1, \dots, a'_k\} \text{ given that } a'_i = a_i + \sum_{j=1}^n x_{j,i}$$

10 Posterior Predictive

$$\hat{w}_i = \frac{a'_i}{\sum \mathbf{a}'}$$

11 ‘Rule-of-Thumb’ Empirical Prior

$$a_i = 1 \text{ for all } i$$

12 **Continuous Data (Multivariate)**

13 **The Multivariate Normal Distribution.**

14 As in the univariate case above, the support for this section is based on the demonstrations
 15 provided by by DeGroot (2005) and Murphy (2007).

16 Because covariance matrices are highly sensitive to outliers, robust covariance esti-
 17 mation is highly complex (Huber & Ronchetti, 2009). The rule-of-thumb prior provided in
 18 this section is based on the estimates for the mean and covariance provided by Campbell
 19 (1980), which was selected because its implementation is straightforward and closed-form.
 20 However, this method’s robustness against outliers depends on the dimensionality of the
 21 data, tolerating at most $\frac{1}{k+1}$ outliers across the data’s k dimensions (Rousseeuw & van
 22 Zomeren, 1990). Consequently, this estimator should be used with caution when the data
 23 display very high dimensionality.

24 • Support:

$$\mathbf{x} = \begin{bmatrix} x_{1,1} & \dots & x_{1,k} \\ \vdots & \ddots & \vdots \\ x_{n,1} & \dots & x_{n,k} \end{bmatrix}, x_{j,i} \in \mathbb{R}$$

- 1 • Sufficient statistics:

$$n \text{ observations} \quad \bar{\mathbf{x}} = \left\{ \frac{1}{n} \sum_{j=1}^n x_{j,1}, \dots, \frac{1}{n} \sum_{j=1}^n x_{j,k} \right\}$$

2 $I_n = [n \times n]$ identity matrix $J_n = [n \times n]$ matrix of ones

- 3 • With known vector $\boldsymbol{\mu}$ and unknown covariance matrix $\boldsymbol{\Sigma}$:

4 Prior Hyperparameters

5 $\mathbf{m} = \{m_1, \dots, m_k | m_i \in \mathbb{R}\}$ Prior mean vector
 $\{p | p > 0\}$ Observations supporting the prior mean
 $\{a | a > k - 1\}$ Observations supporting the prior precision
 $\boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1}$ Prior precision matrix

6 Conjugate Prior

$$\Pr(\boldsymbol{\mu}, \boldsymbol{\Sigma} | \mathbf{m}, p, a, \boldsymbol{\Lambda}) = \frac{p^{k/2} |\boldsymbol{\Lambda}|^{a/2} \exp\left(-\frac{1}{2} \text{tr}\left(\boldsymbol{\Lambda} \boldsymbol{\Sigma}^{-1}\right) - \frac{p}{2} (\boldsymbol{\mu} - \mathbf{m})^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} - \mathbf{m})\right)}{(2\pi)^{k/2} |\boldsymbol{\Sigma}|^{(a+k)/2+1} 2^{ak/2} \pi^{k(k-1)/4} \prod_{j=1}^k \Gamma\left(\frac{a+1-j}{2}\right)}$$

7 Normalizing Constant

$$m(\mathbf{m}, p, a, \boldsymbol{\Lambda}) = \left(\frac{2\pi}{p}\right)^{k/2} \cdot \frac{2^{ak/2} \pi^{k(k-1)/4} \prod_{j=1}^k \Gamma\left(\frac{a+1-j}{2}\right)}{|\boldsymbol{\Lambda}|^{a/2}} \quad (16)$$

8 Posterior Hyperparameters

$$\mathbf{m}' = \frac{p}{p+n} \mathbf{m} + \frac{n}{p+n} \bar{\mathbf{x}} \quad p' = p+n \quad a' = a+n$$

9 $\boldsymbol{\Lambda}' = \boldsymbol{\Lambda} + \mathbf{x}^\top \left(I_n - \frac{1}{n} J_n\right) \mathbf{x} + \frac{pn}{p+n} \text{SumSq}(\bar{\mathbf{x}}, \mathbf{m})$

10 Posterior Predictive

$$\hat{\boldsymbol{\mu}} = \mathbf{m}' \quad \hat{\boldsymbol{\Sigma}} = \frac{\boldsymbol{\Lambda}'}{a' - k - 1}$$

11 ‘Rule-of-Thumb’ Empirical Prior

$$a = \text{robustmean}(x) \quad b = 1$$

12 $c = 1 \quad d = \text{inv}(\text{robustcov}(x))$

13 References

14 Anderson, A. (1911). Labour legislation. In H. Chisholm (Ed.), *Encyclopædia britannica* (11th ed.,
 15 Vol. 16, pp. 7–28). Encyclopedia Britannica Company.

16 Bai, J. (1997). Estimation of a change point in multiple regression models. *Review of Economics*
 17 *and Statistics*, 79, 551–563.

- 1 Campbell, N. A. (1980). Robust procedures in multivariate analysis i: Robust covariance estimation.
2 *Journal of the Royal Statistical Society, Series C (Applied Statistics)*, 29, 231–237.
- 3 Carlin, B. P., Gelfand, A. E., & Smith, A. F. M. (1992). Hierarchical bayesian analysis of changepoint
4 problems. *Applied Statistics*, 41, 389–405.
- 5 Carlin, B. P., & Louis, T. A. (2000). Empirical bayes: Past, present, and future. *Journal of the*
6 *American Statistical Association*, 95, 1286–1289.
- 7 Casella, G. (1985). An introduction to empirical bayes data analysis. *The American Statistician*,
8 39, 83–87.
- 9 Cavanaugh, J. E., & Neath, A. A. (1999). Generalizing the derivation of the schwarz information
10 criterion. *Communications in Statistics - Theory and Methods*, 28, 49–66.
- 11 Chen, S. S., & Gopalakrishnan, P. S. (1998). Speaker, environment, and channel change detection
12 and clustering via the bayesian information criterion. *Proceedings of the DARPA Broadcast News*
13 *Transcription and Understanding Workshop*, 8-13.
- 14 DeGroot, M. H. (2005). *Optimal statistical decisions*. Wiley-Interscience.
- 15 Fearnhead, P. (2006). Exact and efficient bayesian inference for multiple changepoint problems.
16 *Statistics and Computing*, 16, 203–213.
- 17 Gallistel, C. R., Fairhurst, S., & Balsam, P. (2004). The learning curve: Implications of a quantitative
18 analysis. *Proceedings of the National Academy of Sciences of the United States of America*, 101,
19 13124–13131.
- 20 Gallistel, C. R., Krishan, M., Liu, Y., Miller, R., & Latham, P. E. (Submitted). The perception of
21 probability.
- 22 Gelman, A., Carlin, J. B., Stern, H. S., & Rubin, D. B. (2003). *Bayesian data analysis*. Chapman
23 & Hall/CRC.
- 24 Hilborn, R. C. (2001). *Chaos and nonlinear dynamics* (2nd ed.). Oxford University Press.
- 25 Huber, P. J., & Ronchetti, E. M. (2009). *Robust statistics* (2nd ed.). John Wiley & Sons.
- 26 Jarrett, R. G. (1979). A note on the intervals between coal-mining disasters. *Biometrika*, 66,
27 191–193.
- 28 Jensen, G., Altschul, D., Danly, E., & Terrace, H. S. (2013). Transfer of a spatial representation of
29 two distinct serial tasks by rhesus macaques. *PLOS ONE*, 8, e70825.
- 30 Meligkotsidou, L., & Dellaportas, P. (2011). Forecasting with non-homogeneous hidden markov
31 models. *Statistics and Computing*, 21, 439–449.
- 32 Murphy, K. P. (2007). *Conjugate bayesian analysis of the gaussian distribution* (Tech. Rep.).
33 University of British Columbia.
- 34 Ó Ruanaidh, J. J. K., & Fitzgerald, W. J. (1996). *Numerical bayesian methods applied to signal*
35 *processing*. Springer.
- 36 Page, E. S. (1954). Continuous inspection schemes. *Biometrika*, 41, 100–115.
- 37 Palmeri, T. J. (1997). Exemplar similarity and the development of automaticity. *Journal of*
38 *Experimental Psychology: Learning, Memory, and Cognition*, 23, 324–354.

- 1 Pesaran, M. H., Pettenuzzo, D., & Timmermann, A. (2006). Forecasting time series subject to
2 multiple structural breaks. *Review of Economic Studies*, 73, 1057–1084.
- 3 Rousseeuw, P. J., & van Zomeren, B. C. (1990). Unmasking multivariate outliers and leverage
4 points. *Journal of the American Statistical Association*, 85, 633–639.
- 5 Samaniego, F. J. (2012). *A comparison of the bayesian and frequentist approaches to estimation*.
6 Springer New York.
- 7 Schwarz, G. (1978). Estimating the dimension of a model. *Annals of Statistics*, 6, 461–464.
- 8 Venkatraman, E. S. (1992). *Consistency results in multiple change-point problems* (Tech. Rep.
9 No. 24). Department of Statistics, Stanford University.
- 10 Wooldridge, J. M. (2002). *Econometric analysis of cross section and panel data*. MIT Press.
- 11 Yang, Y. (2005). Can the strengths of aic and bic be shared? a conflict between model identification
12 and regression estimation. *Biometrika*, 92, 937–950.
- 13 Yao, Y. C. (1988). Estimating the number of change-points vis schwarz' criterion. *Statistics and*
14 *Probability Letters*, 6, 181–189.