# Probabilistic soil strata delineation using DPT data and Bayesian changepoint detection

Stephen K. Suryasentana[1], Ph.D.

Myles Lawler[2], Ph.D.

Brian B. Sheil[3], Ph.D.

Barry M. Lehane[4], Ph.D.

**Affiliations**

[1] Lecturer, Department of Civil and Environmental Engineering, University of Strathclyde, 75 Montrose St, Glasgow G1 1XJ, UK.

[2] Independent Geotechnical Consultant, Ireland.

[3] RAEng Research Fellow, Department of Engineering Science, University of Oxford, Parks Road, Oxford OX1 3PJ, UK.

[4] Winthrop Professor, Department of Civil, Environmental and Mining Engineering, University of Western Australia, 35 Stirling Hwy, Crawley WA 6009, Australia.

**Full contact details of corresponding author**

Stephen K. Suryasentana

stephen.suryasentana@strath.ac.uk

Main text word count: 2517

Figures: 5

Tables: 1

Mar 26, 2022

**Abstract**

Soil strata delineation is a fundamental step for any geotechnical engineering design. The dynamic penetration test (DPT) is a fast, low cost in-situ test that is commonly used to locate boundaries between strata of differing density and driving resistance. However, DPT data are often noisy and typically require time-consuming, manual interpretation. This paper investigates a probabilistic method that enables delineation of dissimilar soil strata (where each stratum is deemed to belong to different soil groups based on their particle size distribution) by processing DPT data with Bayesian changepoint detection methods. The accuracy of the proposed method is evaluated using DPT data from a real-world case study, which highlights the potential of the proposed method. This study provides a methodology for faster DPT-based soil strata delineation, which paves the way for more cost-effective geotechnical designs.

**Keywords**

**List of notation**

| | |
|---|---|
| $N$ | DPT no. of blows |
| $r_z$ | 'Run length' random variable |
| $x_{1:n}$ | Set of data $\{x_1, x_2, \ldots, x_{n-1}, x_n\}$ |
| $\alpha, \beta$ | Parameters of inverse gamma distribution |
| $p_{\text{cp}}$ | Changepoint probability threshold |

1    **Introduction**

2    Soil strata delineation is a fundamental step for any geotechnical engineering design.

3    Delineation divides the soil volume into separate layers of geological material deemed to belong

4    to the same group. This process typically requires a time-consuming, manual interpretation of a

5    combination of borehole data and associated in-situ and laboratory test results (Parry et al.

6    2014). It is highly desirable to develop a rapid approach that can delineate the soil strata

7    automatically.

8    The cone penetration test (CPT) (Lunne et al. 1997) is an in-situ ground investigation method

9    that is widely used for soil delineation by applying soil behaviour type (SBT) classification rules

10   (e.g. Robertson 1990; Jefferies and Davies 1993; Schneider et al. 2008) to the measured CPT

11   data. Other delineation approaches include fuzzy analysis (Zhang and Tumay 1999), clustering

12   analysis (Hegazy and Mayne 2002; Depina et al. 2016), signal processing analysis (Ching et al.

13   2015) and statistical/Bayesian analysis (Wickremesinghe and Campanella 1991; Phoon et al.

14   2003; Wang et al. 2013, 2019, 2020; Li et al. 2016; Cao et al. 2019). Bayesian analysis has the

15   advantages of being robust to noisy data and allowing quantification of uncertainty, although it

16   tends to be computationally intensive.

17   The dynamic probing/penetration test (DPT) is a fast and low cost in-situ ground investigation

18   method (BS 2005), which bears some similarities to both CPT and the standard penetration test

19   (SPT). Like CPT, DPT uses a cylindrical steel cone penetrometer. However, DPT drives the

20   cone into the ground using a hammer, and the measured result is the number of blows $N$ for a

21   given penetration (e.g. 100mm). The primary advantage of DPT over CPT is lower costs, faster

22   speed of operation and applicability in terrains with poor accessibility. However, there are limited

23   methods to interpret DPT results for soil strata delineation.

24   This paper aims to develop a method that enables fast soil strata delineation using DPT data.

25   The proposed method uses Bayesian changepoint detection (BCPD) methods to detect abrupt

26   changes in the soil data trends indicative of transitions between different soil strata. Unlike most

27   Bayesian approaches, the proposed method is computationally efficient. Two BCPD methods

28   are explored: (i) 'online', where each data point is processed as it becomes available and

29   inferences are made without knowledge of future measurements (e.g. Fearnhead and Liu, 2007;

30    Adams and MacKay, 2007); and (ii) 'offline', where the entire DPT dataset is required before

31    making inference (e.g. Barry and Hartigan, 1993; Stephens, 1994; Fearnhead, 2005, 2006).

32    The proposed method divides the soil profile up into three dissimilar soil categories: (i)

33    predominantly fine-grained soils (e.g. clay, silt), (ii) predominantly sand, and (iii) predominantly

34    gravel. These soil categories have very different permeability, stiffness and strength properties

35    such that poor identification will have a negative impact on optimal geotechnical design. The

36    proposed method bears some similarities to that of Zhang and Tumay (1999), who applied fuzzy

37    analysis to CPT data to identify three soil categories, although the methodology and nature of

38    the data are different. The performance of the proposed BCPD methods are evaluated using a

39    real-world case study.

40

41    **Methodology**

42    Changepoints are abrupt changes in data, which typically represent transitions between states,

43    as shown in Fig. 1. Given a sequence of data, these changepoints split the data into a set of

44    non-overlapping partitions, where it is assumed that the data within a partition are generated by

45    the same model. While many changepoint detection methods are available (Reeves et al. 2007;

46    Aminikhanghahi and Cook 2017; Truong et al. 2020), this paper focuses on Bayesian

47    changepoint detection (BCPD) methods.

48    *Online Bayesian changepoint detection*

49    The first method investigated in this paper is an online BCPD method (Adams and Mackay

50    2007), denoted 'BCPD-ON'. In the following exposition, the notation $x_{1:n}$ refers to the set of data

51    $\{x_1, x_2, \dots, x_{n-1}, x_n\}$. BCPD-ON estimates the probability of a changepoint at a given depth based

52    only on data processed up to that depth. It does so by computing the probability distribution of a

53    random variable called the 'run length' $r_z$, which represents the length of the current data

54    partition. Each new data point either (a) comes from the same distribution, in which case the

55    parameter estimates of the current distribution is updated using Bayes' theorem and $r_z$

56    increases by one, or (b) it belongs to a new distribution which means a changepoint occurs and

57    the new distribution will reset back to the prior distribution and $r_z$ resets to zero. When the most

58    probable value of $r_z$ is zero, it is likely that there is a changepoint at depth $z$, the probability of

59    which is equivalent to the posterior probability of $r_z = 0$:

$$p(\text{changepoint at } z| \, x_{1:z}) = p(r_z = 0|x_{1:z}) \tag{1}$$

60    The posterior distribution of the run length i.e. $p(r_z|x_{1:z})$ can be calculated as:

$$p(r_z|x_{1:z}) = \frac{p(r_z, x_{1:z})}{p(x_{1:z})} \tag{2}$$

61    where $p(x_{1:z}) = \sum_{r_z} p(r_z, x_{1:z})$. The joint distribution $p(r_z, x_{1:z})$ can be calculated using the

62    following recursive relationship:

$$p(r_z, x_{1:z}) = \sum_{r_{z-1}} p(r_z, x_z, |r_{z-1}, x_{1:z-1}) p(r_{z-1}, x_{1:z-1})$$

$$= \sum_{r_{z-1}} p(r_z|r_{z-1}) \, p(x_z|r_{z-1}, \boldsymbol{x}_z^r) \, p(r_{z-1}, x_{1:z-1}) \tag{3}$$

63    where $\boldsymbol{x}_z^r$ is the set of data associated with the run length $r_z$. $p(r_{z-1}, x_{1:z-1})$ is a recursive term,

64    which represents the previous iteration of Eq. 3 at depth $z - 1$. $p(r_z|r_{z-1})$ is the conditional

65    distribution of the run length. Finally, $p(x_z|r_{z-1}, \boldsymbol{x}_z^r)$ is the posterior predictive distribution and it

66    can be calculated analytically by assuming that the data point $x_z$ comes from some probability

67    distribution (e.g. Gaussian) and by adopting conjugate priors. More details about these

68    calculations can be found in Adams and Mackay (2007).

69    *Offline Bayesian changepoint detection*

70    The second method investigated in this paper is an offline BCPD method (Fearnhead 2005,

71    2006) denoted 'BCPD-OFF', which was previously employed by Houlsby and Houlsby (2013) for

72    clay layer delineation using undrained shear strength data. BCPD-OFF is based on a recursive

73    algorithm that computes the posterior probability distribution exactly over the location of

74    changepoints. This is significantly more efficient than previous Markov Chain Monte Carlo

75    (MCMC) approaches for computing the posterior (e.g. Punskaya et al. 2002).

76    In this case, the data within each partition are modelled by some probability distribution, with

77    distribution parameters independent of those determined for other partitions.  Let $c_j$ represent

78     the $j$th changepoint. The posterior distribution of $c_j$ is $p(c_j|x_{1:n})$. The probability of a

79     changepoint occurring at depth $z$ can be calculated as:

$$p(\text{changepoint at } z|x_{1:n}) = \sum_{j=1}^{z} p(c_j = z|x_{1:n}) \qquad (4)$$

80     where all possible scenarios of 1 to $z$ changepoints thus far are considered. This approach

81     differs from that of Houlsby and Houlsby (2013), which first identifies the *maximum a posteriori*

82     (MAP) number of changepoints and then the conditional MAP locations of the changepoints.

83     This modification makes the outputs of BCPD-OFF and BCPD-ON identical, thereby allowing

84     direct comparisons.

85     $p(c_j|x_{1:n})$ in Eq. 4 is obtained by marginalising out the previous changepoints:

$$p(c_j|x_{1:n}) = \int p(c_j, \dots, c_1|x_{1:n}) \, dc_{j-1} \dots dc_1 \qquad (5)$$

86     As the probability of a changepoint is assumed to be dependent only on the previous

87     changepoint, the integrand in Eq. 5 can be calculated as:

$$p(c_j, \dots, c_1|x_{1:n}) = p(c_j|c_{j-1}, x_{1:n})p(c_{j-1}|c_{j-2}, x_{1:n}) \dots p(c_2|c_1, x_{1:n})p(c_1|x_{1:n}) \qquad (6)$$

88     Each of the terms on the right hand side of Eq. 6 can be calculated exactly and efficiently using

89     the recursive algorithm described in Fearnhead (2005, 2006).

90

91     **Case study**

92     The proposed BCPD methods are evaluated using a case-study involving multi-layered alluvial

93     deposits, consisting of sands, silts, clays, and gravels. This case study is based on the

94     Deutsche Bahn AG (German Rail) 'DB46/2' project, which is an expansion line from Emmerich

95     to Oberhausen in Germany. A complex three-dimensional (3D) ground model for this project

96     has been documented in Prinz (2019). This paper considers 26 DPT tests from the case study:

97     20 (approximately 77% of the dataset) are randomly selected for calibration of the priors and

98     hyperparameters for BCPD-OFF and BCPD-ON; the remaining 6 DPT locations (labelled 'T1' to

99     'T6') are used for testing to evaluate the performance of the calibrated methods. A plan map of

100    the DPT calibration and test locations is shown in Fig. 2.

101   Expert predictions are also made for each DPT location, where the soil strata are identified

102   among the three soil categories defined in the introduction. These expert predictions were

103   extracted from the 3D ground model that was developed separately for the case study (Prinz

104   2019). This ground model was based on careful, manual interpretation of both the DPT data and

105   the borehole data in an integrated manner, ensuring no conflicts between the interpretation of

106   the soil layering boundaries based on both types of data (e.g. the soil stratification interpreted

107   from the DPT data should be consistent with that observed from a neighbouring borehole). Fig.

108   3 shows a typical DPT profile from one of the DPT locations and its corresponding expert

109   prediction of the soil strata. The proposed BCPD methods will be applied to DPT data only.

110

**Calibration**

112   For both BCPD-OFF and BCPD-ON, the data in each partition are assumed to be normally

113   distributed with unknown mean $\mu$ and variance $\sigma^2$. Therefore, the DPT data were preprocessed

114   using a Freeman-Tukey transformation (Freeman and Tukey 1950): $N_{\text{transformed}} = \sqrt{N} + \sqrt{N+1}$

115   where $N$ represents the raw DPT blowcount data. This transformation is typically used to make

116   discrete count data better approximate a normal distribution (Mosteller and Youtz 2006; Lin and

117   Xu 2020). To test for normality of the transformed data, the Shapiro-Wilk test (Shapiro and Wilk

118   1965) was applied to the transformed data in each soil layer at DPT locations where

119   neighbouring borehole data is available to determine the approximate locations of the soil layer

120   boundaries. The p-values obtained are greater than 0.05 and thus the null hypothesis that the

121   transformed data is normally distributed is not rejected. Following Houlsby and Houlsby (2013),

122   the variance $\sigma^2$ is assumed to follow an inverse gamma distribution and the distribution

123   parameters $\alpha$ = 1.8 and $\beta$ = 0.38 are obtained by curve-fitting the cumulative distribution of the

124   variance for the DPT calibration dataset, as shown in Fig. 4.

125   Outputs of interest for both BCPD-ON and BCPD-OFF are the probabilities of a changepoint

126   occurrence at each depth (i.e. using Eq. 1 and Eq. 4, respectively). When the changepoint

127   probability exceeds a predefined threshold $p_{\text{cp}}$, the soil is considered to have changed category

128   at this depth. The optimal value of $p_{\text{cp}}$ is dependent on the method adopted (BCPD-ON or

129   BCPD-OFF) and is calibrated as a hyperparameter. For each method, a grid search is

130    implemented within the set of trial $p_{cp} = \{0.1, 0.15, 0.2, 0.25, 0.3, \ldots, 0.8, 0.85, 0.9\}$ to identify the

131    value of $p_{cp}$ that achieve the best match with the expert predictions for the soil stratification at

132    each DPT calibration location. To quantify the match with expert predictions, the accuracy

133    measure, F1 score, is adopted,

$$\text{F1 score} = 2(\text{Precision} * \text{Sensitivity})/(\text{Precision} + \text{Sensitivity}) \tag{7}$$

134    where Precision = True Positive/(True Positive + False Positive) and Sensitivity = True

135    Positive/(True Positive + False Negative). True Positive (TP) is the number of times an expert

136    prediction for soil layer boundary has been correctly identified, while False Positive (FP) is the

137    number of times an expert prediction for soil layer boundary has been incorrectly identified.

138    False Negative (FN) is the number of times an expert prediction for soil layer boundary has not

139    been identified. A higher F1 value indicates a better match with the expert predictions. As the

140    predicted boundaries based on the DPT data are not expected to exactly match the expert

141    predictions, this paper considers a soil layer boundary to be correctly identified if the DPT-

142    predicted boundary is within a distance of 1m from the expert prediction for a boundary. The

143    grid search exercise gives the optimal values of $p_{cp}$ = 0.45 and 0.4 for BCPD-OFF and BCPD-

144    ON respectively.

145

146    **Results**

147    Fig. 5 shows the soil strata predictions determined using BCPD-OFF and BCPD-ON for the 6

148    DPT test locations. The BCPD changepoint probability predictions are shown in the figure as

149    grey lines and a soil strata boundary is identified when these predictions exceed $p_{cp}$.

150    From this figure, it can be observed that both BCPD-OFF and BCPD-ON perform well for most

151    locations, where the predicted soil strata boundaries are similar to the expert predictions. The

152    exception to this is Location T3, where the expert prediction for the soil strata is very complex,

153    and both BCPD methods only detect some of the soil strata boundaries. Nevertheless, the

154    overall performance is encouraging as the BCPD predictions agree well with the expert

155    predictions, despite using information only from the local DPT data. Some of the soil strata

156    boundary detections are noteworthy (e.g. see Fig. 5d), as they are not obvious from manual

157    inspection of the noisy DPT data alone.

158    Comparing the two BCPD methods, it is evident that BCPD-OFF is the more sensitive of the

159    two, as it can detect more soil strata boundaries (e.g. at locations T3 and T4), despite having a

160    higher $p_{cp}$ than BCPD-ON. However, this increased sensitivity comes with the drawback of

161    producing more false positives (see Figs. 5b, c). To quantify the accuracy of both methods, their

162    F1 scores are calculated based on Eq. 7, as detailed in Table 1. BCPD-ON has a slightly higher

163    F1 score than BCPD-OFF, indicating that BCPD-ON has a slightly better balance of precision

164    and sensitivity. In terms of computational efficiency, BCPD-ON has the advantage of being

165    much faster than BCPD-OFF (on average, BCPD-ON takes approximately 0.03 seconds to

166    process each DPT location, while BCPD-OFF takes approximately 5 seconds).

167    A key highlight is that both BCPD-OFF and BCPD-ON could detect soil strata boundaries

168    quickly and automatically without manual intervention. This makes them helpful to industry

169    practitioners for extracting additional insights from the DPT data to complement their current

170    workflow for identifying soil strata. A useful application of the approach could be, for example, to

171    assist the design of large-scale foundation projects such as solar farms. Engineers could be

172    faced with up to 1000 DPT locations in one project, and this approach provides a consistent,

173    automated and rapid way to interpret the soil stratigraphy.

174    When applying these BCPD methods to a new site, a calibration process should be carried out

175    to obtain site-specific values for both the priors and the $p_{cp}$ hyperparameter; this should provide

176    improved soil layer boundary detection results. Site-specific calibration should not be an issue

177    as DPT tests are typically carried out in conjunction with borehole tests. However, if calibration

178    data is not available at the new site, the calibrated parameters in this paper may be used for

179    preliminary analysis, using the BCPD methods to highlight potential locations of soil layer

180    boundaries through the 'spikes' in the changepoint probability. However, caution is advised as a

181    non-site specific calibration of $p_{cp}$ and the priors will affect the precision of the soil layer

182    boundary detections. To investigate the sensitivity of the calibration to the number of DPT tests,

183    the calibration results (i.e. the calibrated values for $\alpha, \beta, p_{cp}$) were determined using random

184    selections of 3, 4, 5, 6, 7, 8, 9, 10, 15, 20 DPT tests. The analysis indicates that when 5 or more

185    DPT tests are used for calibration, the calibrated $p_{\mathrm{cp}}$ values are the same and the calibrated

186    values of $\alpha, \beta$ change by less than 4% from the values used in the current study. However,

187    caution should be advised against taking this as a general rule as these results may be specific

188    to the dataset used in the current study. Furthermore, in this study, each DPT test location is

189    near a calibration location. The effect of the distance between the calibration and test locations

190    on the predictive accuracy of the BCPD methods has not been evaluated in this study. Further

191    research is required with a comprehensive study, involving a larger database of DPT data from

192    a wider range of sites, to provide more definitive answers to the above questions and to obtain

193    values of the priors and hyperparameter more suited to general use across different sites.

194

195    **Conclusion**

196    This paper proposes a fast, automatic Bayesian approach for soil strata delineation using DPT

197    data. The proposed approach is based on the concept of offline and online Bayesian

198    changepoint detection, which allows both retrospective and real-time soil strata delineation. Its

199    reliability and utility have been evaluated using DPT data from a real-world case study. The

200    proposed approach is very fast to run and provides additional insights from the DPT data for a

201    more robust soil strata identification solution.

202

214

215 **Data availability statement**

216 Some or all data, models, or code that support the findings of this study are available from the

217 corresponding author upon reasonable request.

**References**

Adams, R. P., and MacKay, D. J. (2007). Bayesian online changepoint detection. arXiv preprint arXiv:0710.3742.

Aminikhanghahi, S., and Cook, D. J. (2017). A survey of methods for time series change point detection. Knowledge and information systems, 51(2), 339-367.

Barry, D., and Hartigan, J. A. (1993). A Bayesian analysis for change point problems. Journal of the American Statistical Association, 88(421), 309-319.

BS (2005). Geotechnical Investigation and Testing–Field Testing–Part 2; Dynamic Probing. BS EN ISO 22476-2, BSI, London, UK.

Cao, Z. J., Zheng, S., Li, D. Q., and Phoon, K. K. (2019). Bayesian identification of soil stratigraphy based on soil behaviour type index. Canadian Geotechnical Journal, 56(4), 570-586.

Ching, J., Wang, J.-S., Juang, C.H., and Ku, C.-S. (2015). Cone penetration test (CPT)-based stratigraphic profiling using the wavelet transform modulus maxima method. Canadian geotechnical journal, 52(12): 1993-2007.

Depina, I., Le, T.M.H., Eiksund, G., and Strøm, P. (2016). Cone penetration data classification with Bayesian Mixture Analysis. Georisk: Assessment and management of risk for engineered systems and geohazards, 10(1): 27-41.

Fearnhead, P. (2005). Exact Bayesian curve fitting and signal segmentation. IEEE Transactions on Signal Processing, 53(6), 2160-2166.3, 747–758.

Fearnhead, P. (2006). Exact and efficient Bayesian inference for multiple changepoint problems. Statistics and computing, 16(2), 203-213.

Fearnhead, P., and Liu, Z. (2007). On-line inference for multiple changepoint problems. Journal of the Royal Statistical Society: Series B (Statistical Methodology), 69(4), 589-605.

Freeman, M. F., and Tukey, J. W. (1950). Transformations related to the angular and the square root. The Annals of Mathematical Statistics, 607-611.

Hegazy, Y.A., and Mayne, P.W. (2002). Objective site characterization using clustering of piezocone data. Journal of Geotechnical and Geoenvironmental Engineering, 128(12): 986-996.

Houlsby, N. M. T., and Houlsby, G. T. (2013). Statistical fitting of undrained strength data. Géotechnique, 63(14), 1253-1263.

Jefferies, M. G. and Davies, M. P. (1993). Use of CPTU to estimate equivalent SPT N60. Geotech. Test. J. 16, No. 4, 458–468.

Li, J., Cassidy, M. J., Huang, J., Zhang, L., and Kelly, R. (2016). Probabilistic identification of soil stratification. Géotechnique, 66(1), 16-26.

Lin, L., and Xu, C. (2020). Arcsine-based transformations for meta-analysis of proportions: Pros, cons, and alternatives. Health Science Reports, 3(3), e178.

Lunne, T., Robertson, P. K. and Powell, J. J. M. (1997). Cone penetration testing in geotechnical practice. London, UK: Blackie Academic and Professional.

Mosteller, F., and Youtz, C. (2006). Tables of the Freeman-Tukey transformations for the binomial and Poisson distributions. In Selected Papers of Frederick Mosteller (pp. 337-347). Springer, New York, NY.

Parry, S., Baynes, F.J., Culshaw, M.G., Eggers, M., Keaton, J.F., Lentfer, K., Novotny, J. and Paul, D. (2014). Engineering geological models - an introduction: IAEG commission 25. Bull. Engng. Geol. Environ., 73, 689–706.

Phoon, K.-K., Quek, S.-T., and An, P. (2003). Identification of statistically homogeneous soil layers using modified Bartlett statistics. Journal of Geotechnical and Geoenvironmental Engineering, 129(7): 649-659.

Prinz, I., (2019). Digitale Baugrundmodelle: BIM in der Geotechnik Erfahrungen und Ableitungen aus dem Projekt Ausbaustrecke Emmerich – Oberhausen (ABS 46/2), einem BIM-Piloten der Deutschen Bahn. Geotechnik 22.

Punskaya, E., Andrieu, C., Doucet, A., and Fitzgerald, W. J. (2002). Bayesian curve fitting using MCMC with applications to signal segmentation. IEEE Transactions on signal processing, 50(3), 747-758.

Reeves, J., Chen, J., Wang, X. L., Lund, R., and Lu, Q. Q. (2007). A review and comparison of changepoint detection techniques for climate data. Journal of applied meteorology and climatology, 46(6), 900-915.
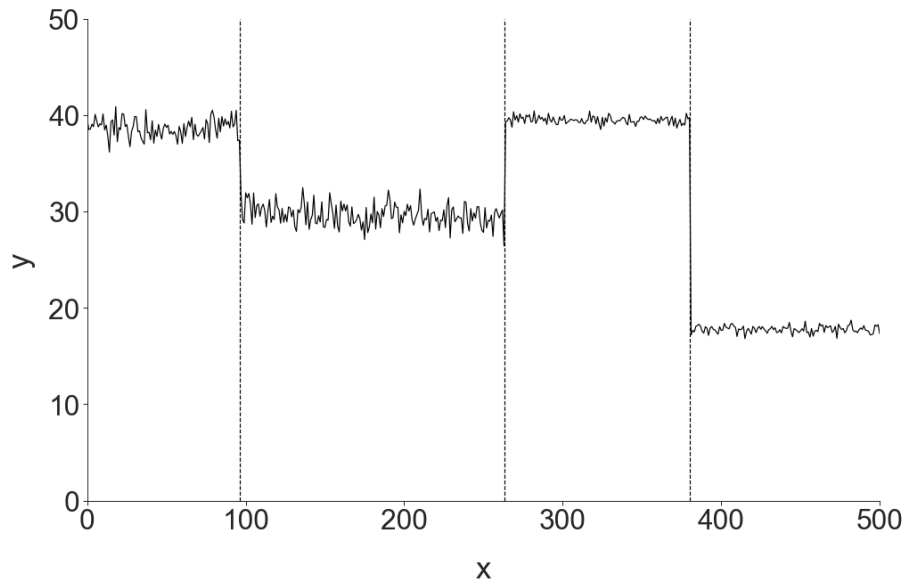
275   Robertson, P. K. (1990). Soil classification using the cone penetration test. Can. Geotech. J. 27,
276        No. 1, 151–158.
277   Schneider, J. A., Randolph, M. F., Mayne, P. W. and Ramsey, N. R. (2008). Analysis of factors
278        influencing soil classification using normalized piezocone tip resistance and pore pressure
279        parameters. J. Geotech. Geoenvir. Engng 134, No. 11, 1569–1586.
280   Shapiro, S. S., & Wilk, M. B. (1965). An analysis of variance test for normality (complete
281        samples). Biometrika, 52(3/4), 591-611.
282   Stephens, D. A. (1994). Bayesian retrospective multiple-changepoint identification. Journal of
283        the Royal Statistical Society: Series C (Applied Statistics), 43(1), 159-178.
284   Truong, C., Oudre, L., and Vayatis, N. (2020). Selective review of offline change point detection
285        methods. Signal Processing, 167, 107299.
286   Wang, Y., Huang, K., and Cao, Z. (2013). Probabilistic identification of underground soil
287        stratification using cone penetration tests. Canadian Geotechnical Journal, 50(7), 766–776.
288   Wang, Y., Hu, Y., & Zhao, T. (2020). Cone penetration test (CPT)-based subsurface soil
289        classification and zonation in two-dimensional vertical cross section using Bayesian
290        compressive sampling. Canadian Geotechnical Journal, 57(7), 947-958.
291   Wang, H., Wang, X., Wellmann, J. F., and Liang, R. Y. (2019). A Bayesian unsupervised
292        learning approach for identifying soil stratification using cone penetration data. Canadian
293        Geotechnical Journal, 56(8), 1184-1205.
294   Wickremesinghe, D., and Campanella, R. (1991). Statistical methods for soil layer boundary
295        location using the cone penetration test. Proc. ICASP6, Mexico City, 2: 636-643.
296   Zhang, Z., and Tumay, M.T. (1999). Statistical to fuzzy approach toward CPT soil classification.
297        Journal of Geotechnical and Geoenvironmental Engineering, 125(3): 179-186.
298
299
300

301 **Table 1** Accuracy calculations for the BCPD-OFF and BCPD-ON soil layer boundary predictions

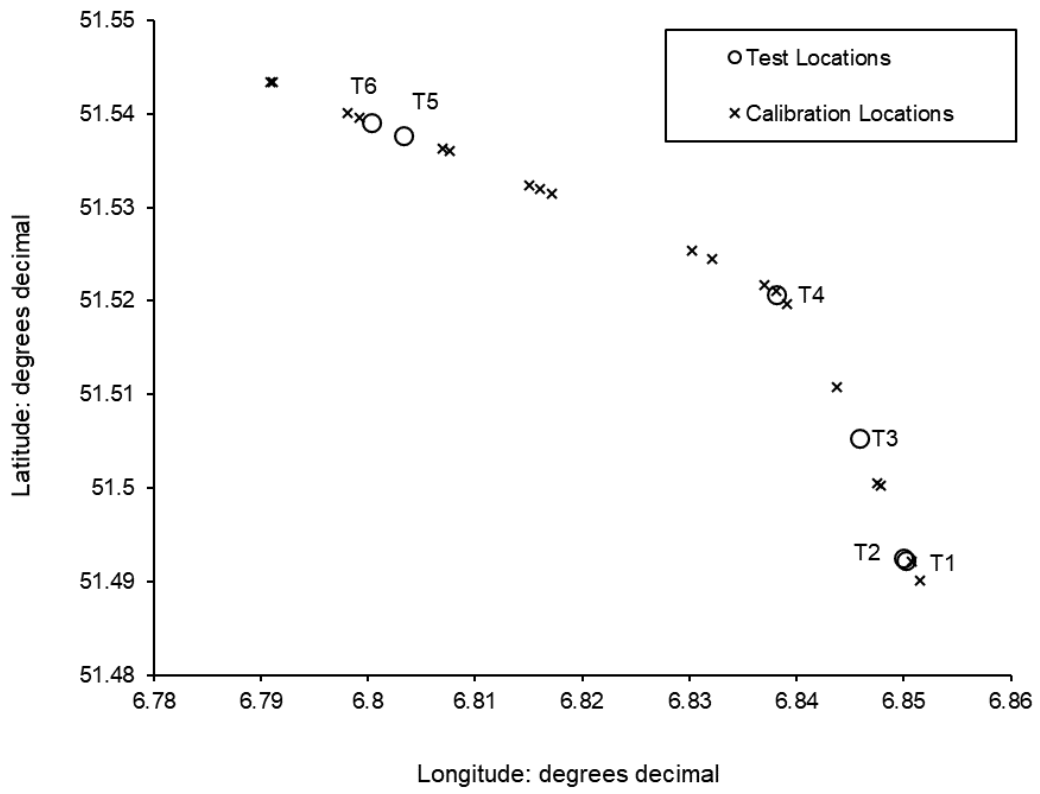|  | TP | FP | FN | Precision | Sensitivity | F1 score |
|---|---|---|---|---|---|---|
| BCPD-OFF | 12 | 3 | 2 | 0.80 | 0.857 | 0.827 |
| BCPD-ON | 10 | 0 | 4 | 1.00 | 0.714 | 0.833 |

302
303
304
305

306 **Figures**
307
308
309



310

311 **Fig. 1** Illustration of a sequence of data with abrupt changes, where y is the measured quantity
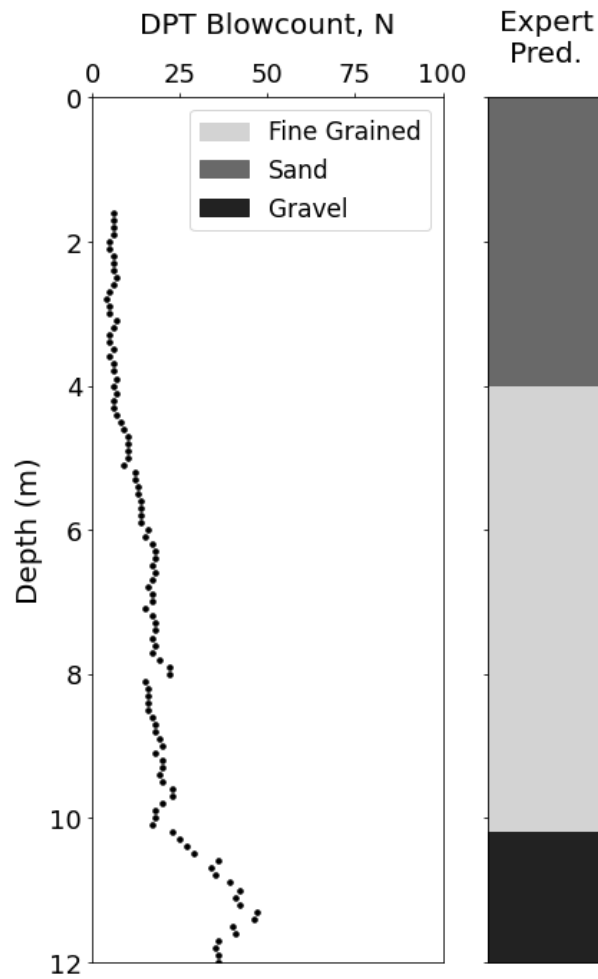312 and x is the index. The dashed lines represent the locations of the changepoints.
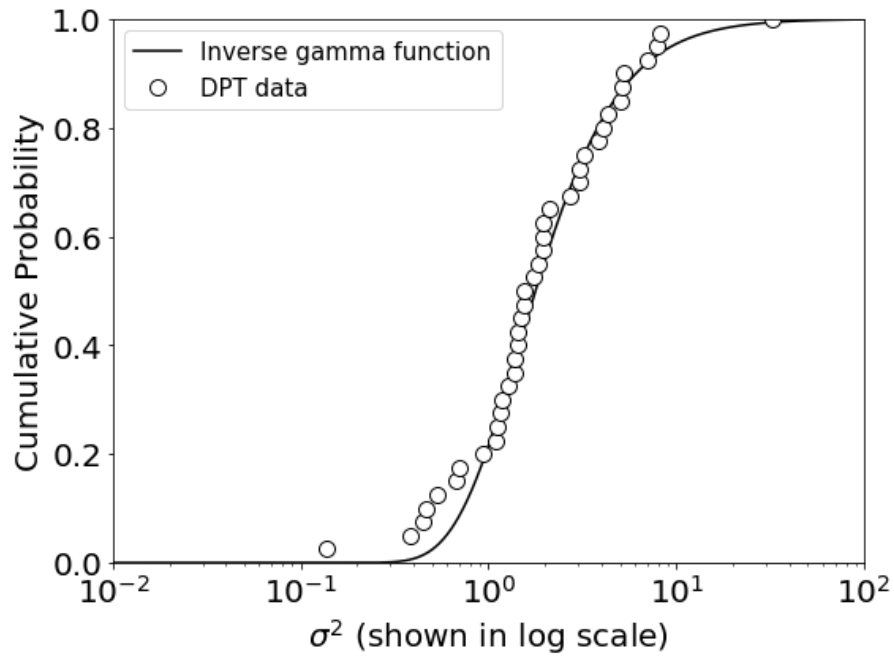
313

314

315

316 **Fig. 2** Locations of DPT dataset used for calibration and testing of the BCPD methods.
317

318

319    **Fig. 3** Exemplar DPT profile showing the development of the DPT blowcount, $N$, with depth.

320    The expert prediction for the soil strata at this location is also shown, where the soil categories
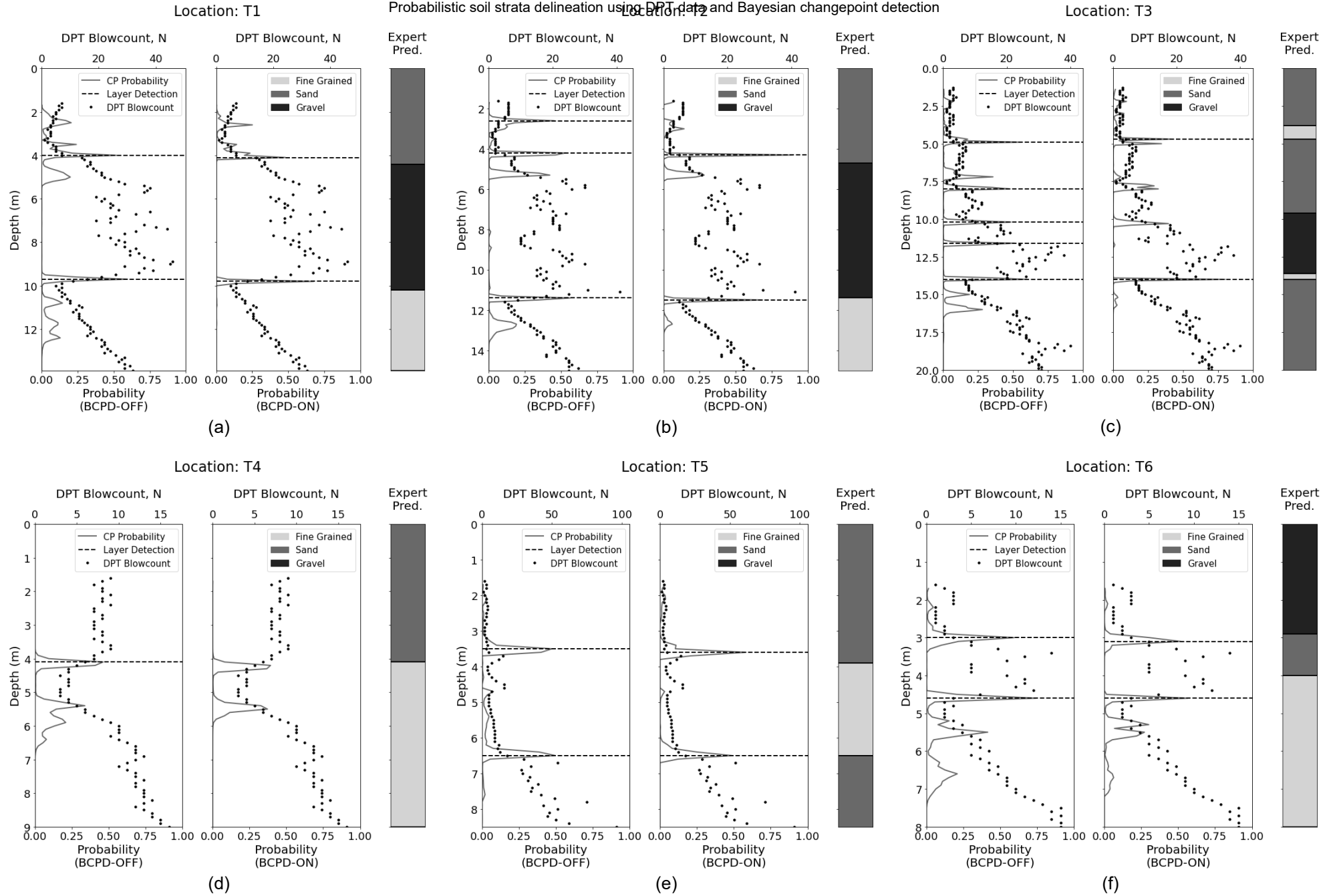
321    are shown in the legend.

322
323
324
325
326
327
328
329

330
331
332    **Fig. 4** Cumulative distribution of the variance of the transformed $N$ data within each soil strata

333        identified in the DPT calibration dataset, compared with the inverse gamma cumulative

334                        distribution with $\alpha = 1.8, \beta = 0.38$.

**Fig. 5** Comparison of soil strata boundaries (shown as horizontal black lines) predicted by BCPD-OFF and BCPD-ON, with the expert predictions, at locations T1 to T6.