Research Article

The last imprint: Differential methylation assay for forensic identification of brain tissue and death determination

Supplementary file

The score for each genomic position is calculated using the following Equation (S1):

Score =
$$\sum_{1}^{n} \left[\sum_{1}^{k} b \log(|d - Mt[n][k]| + p) + \sum_{1}^{w} \log(|(1 - d) - Mb[n][k]| + p) \right]$$
 (S1)

where Mt_1 ... Mt_k refers to n vectors in length k, describing the methylation state in k non-brain tissue samples in n genomic positions; Mb_1 ... Mb_w refers to n vectors in length k, describing the methylation state in w brain tissue samples in n genomic positions; p is the penalty parameter; p is the "brain preference" parameter, representing the relative importance of methylation values in brain samples; and p is the ideal methylation state for each tissue, for each score type. The value of p is defined based on the sample type and the corresponding score classification. For brain samples, the hypo score was assigned a value of p0, while the hyper score was assigned 1. In contrast, for non-brain samples, the hypo score was assigned 1, and the hyper score was assigned 0.

Potential genomic regions containing at least five CpG sites within 300 bp were ranked based on CpG count and the mean informativeness score of the top 50% most informative sites using Equation (S2):

$$Score = \frac{\min(amplicon_len, \max_rewarding_size)}{\min(\max_len, \max_rewarding_size)} * |\max_score| * length_weight + \max_bias * \max_score + \frac{top_half_sum}{amplicon_len}$$
(S2)

The parameters and their corresponding final values for Equation (S2) are defined as follows. The maximum length (max_len) was set to 300. A value of 1 was used to represent the bias toward the best position within the amplicon range (max_bias). The maximum amplicon size rewarded (max_rewarding_size) was set at 150. The bias toward amplicon length was given a weight of 0.0 (length_weight), while 0.5 was assigned to represent the fraction of top positions given higher importance (top_half_fraction).

For information regarding potential genomic position ranges, the sum of the best top_half_fraction positions was recorded as Top_half_sum. The score of the best position within the range was referred to as top_score, while the distance between the first and last CpG sites was defined as amplicon_len. In addition, the maximum and minimum scores within the range were denoted as max_score and min_score, respectively.