

1 **RESEARCH ARTICLE**

2 **Guidance for DNA methylation studies: Statistical insights from the Illumina EPIC array**

3

4 Georgina Mansell<sup>1</sup>, Tyler J Gorrie-Stone<sup>2</sup>, Yanchun Bao<sup>3</sup>, Meena Kumari<sup>3</sup>, Leonard S Schalkwyk<sup>2</sup>,  
5 Jonathan Mill<sup>1</sup>, Eilis Hannon<sup>1\*</sup>

6

7 <sup>1</sup>University of Exeter Medical School, University of Exeter, RD&E Hospital, Barrack Road, Exeter,  
8 Devon, EX2 5DW, United Kingdom

9 <sup>2</sup>School of Biological Sciences, University of Essex, Colchester, Essex, CO4 3SQ, United Kingdom

10 <sup>3</sup>Institute for Social and Economic Research, University of Essex, Colchester, Essex, CO3 3LG,  
11 United Kingdom

12

13 \* Correspondence to: Eilis Hannon, University of Exeter Medical School, RILD Building, Royal  
14 Devon and Exeter Hospital, Barrack Road, Exeter. EX2 5DW. UK. Tel.: +44 1392 408278. E-mail:  
15 [e.j.hannon@exeter.ac.uk](mailto:e.j.hannon@exeter.ac.uk).

16

17 Georgina Mansell G.Mansell@exeter.ac.uk

18 Tyler Gorrie-Stone tgorri@essex.ac.uk

19 Yanchun\_Bao ybaoa@essex.ac.uk

20 Meena Kumari mkumari@essex.ac.uk

21 Leonard Schalkwyk lschal@essex.ac.uk

22 Jonathan Mill J.Mill@exeter.ac.uk

23 Eilis Hannon E.J.Hannon@exeter.ac.uk

24

- 1 **Keywords:** DNA methylation, Epigenome-wide association study (EWAS), multiple testing, Illumina
- 2 EPIC array, power
- 3

1 **Abstract**

2

3 **Background:** There has been a steady increase in the number of studies aiming to identify DNA  
4 methylation differences associated with complex phenotypes. Many of the challenges of epigenetic  
5 epidemiology regarding study design and interpretation have been discussed in detail, however there  
6 are analytical concerns that are outstanding and require further exploration. In this study we seek to  
7 address three analytical issues. First, we quantify the multiple testing burden and propose a standard  
8 statistical significance threshold for identifying DNA methylation sites that are associated with an  
9 outcome. Second, we establish whether linear regression, the chosen statistical tool for the majority of  
10 studies, is appropriate and whether it is biased by the underlying distribution of DNA methylation  
11 data. Finally, we assess the sample size required for adequately powered DNA methylation  
12 association studies.

13

14 **Results:** We quantified DNA methylation in the Understanding Society cohort (n = 1,175), a large  
15 population based study, using the Illumina EPIC array to assess the statistical properties of DNA  
16 methylation association analyses. By simulating null DNA methylation studies, we generated the  
17 distribution of p-values expected by chance and calculated the 5% family-wise error for EPIC array  
18 studies to be  $9 \times 10^{-8}$ . Next, we tested whether the assumptions of linear regression are violated by  
19 DNA methylation data and found that the majority of sites do not satisfy the assumption of normal  
20 residuals. Nevertheless, we found no evidence that this bias influences analyses by increasing the  
21 likelihood of affected sites to be false positives. Finally, we performed power calculations for EPIC  
22 based DNA methylation studies, demonstrating that existing studies with data on ~1000 samples are  
23 adequately powered to detect small differences at the majority of sites.

24

25 **Conclusion:** We propose that a significance threshold of  $P < 9 \times 10^{-8}$  adequately controls the false  
26 positive rate for EPIC array DNA methylation studies. Moreover, our results indicate that linear  
27 regression is a valid statistical methodology for DNA methylation studies, despite the fact that the  
28 data do not always satisfy the assumptions of this test. These findings have implications for

1 epidemiological-based studies of DNA methylation and provide a framework for the interpretation of  
2 findings from current and future studies.

3

4

## 1 **Background**

2 There is increasing interest in the role of epigenetic processes in health and disease, with the primary  
3 focus of most epigenetic epidemiological studies being on DNA methylation (DNAm) [1]. Platforms  
4 such as the Illumina 450K Human Methylation microarray (450K array) and the Illumina EPIC  
5 Human Methylation microarray (EPIC array) have enabled the economical, high-throughput profiling  
6 of methylomic variation across large numbers of samples. In recent years a number of epigenome-  
7 wide association studies (EWAS), which aim to identify DNAm differences associated with  
8 environmental exposure and disease, have been reported for a range of complex phenotypes including  
9 cancer [2-4], autoimmune disorders [5-7], psychiatric illnesses [8, 9], neurodevelopmental disorders  
10 [10, 11] and dementia [12, 13].

11

12 Primarily due to the dynamic nature of the epigenome throughout development, across different cell  
13 types and in response to environmental exposures, much has previously been written regarding the  
14 specific nuances of performing an EWAS compared to a genome-wide association study (GWAS) of  
15 genetic variation [14-16]. However, this literature is mainly focused on study design and  
16 interpretation rather than specific analytical issues relating to the characteristics of the data. One  
17 concern that has merited some discussion relates to whether the distribution of DNAm data violates  
18 the assumptions of Gaussian linear regression [17, 18], the most commonly used analysis model as it  
19 allows for the inclusion of covariates relating to both biological and technical confounders. For each  
20 molecule of DNA in a single cell, DNAm is a binary entity, in that at any cytosine it is either present  
21 or absent. However, as almost all DNAm studies profile either bulk tissue - comprising multiple cell  
22 types - or a population of purified cells, the analyses are essentially measuring the proportion of cells  
23 (taking a value between 0 and 1) in a sample that are methylated at a specific genomic position [19].  
24 While across the sites profiled on Illumina arrays DNAm has a bimodal distribution with peaks of  
25 hypomethylation (i.e. unmethylated sites) and hypermethylation (i.e. methylated sites), there is a  
26 significant subset of sites exhibiting intermediate levels of DNAm (proportion of methylated alleles =  
27 ~0.5). As the presence/absence of DNAm primarily distinguishes different cell types and tissues, in  
28 studies of a single tissue, which the majority of epigenetic epidemiology studies are, it is unlikely that

1 the distribution at individual DNAm sites (the standard unit of analysis in an EWAS) will be bimodal.  
2 However, it is likely that the distributions will be variable and often non-normal, meaning that the  
3 assumption that the residuals of a linear regression fit are normally distributed may not hold.  
4 Furthermore, as DNAm levels are bounded by the limits of 0 and 1 it means that at the extreme ends  
5 of the distribution the variance is compressed. States of hypo and hypermethylation often define cell  
6 types and would not be expected vary biologically within a cell type, beyond any technical noise in  
7 the assay. This is exacerbated by the fact that the sensitivity of the microarray technology is less  
8 precise at these extremes of the distribution, and hence some measured variation is often present for  
9 these theoretically non-variable sites. This property of the data is called heteroskedasticity, defined as  
10 a relationship between the mean and variance of a dataset, and violates another assumption of linear  
11 regression. Although these concerns should be considered when it comes to deciding the statistical  
12 methodology, it is not currently known whether these violations are sufficient to bias analyses and  
13 introduce false positive or even false negative findings.

14  
15 Consistent with studies of other types of genomic variation, another challenge for EWAS is how to  
16 account for the multiple testing burden in a typical analysis; for example, the Illumina EPIC array  
17 assays DNAm at base pair resolution for > 850 000 sites across the genome. Currently, a range of  
18 approaches are used to establish an appropriate significance threshold and there is no standard  
19 significance threshold as is used in GWAS. A common approach is a Bonferroni correction for the  
20 number of probes on the array [20-23] although this is often presumed to be too conservative as  
21 DNAm values at neighboring probes are known to be correlated [24], and many sites on the array are  
22 non-variable. An alternative, potentially more powerful, approach sets a permissible false discovery  
23 rate (FDR), and identifies the top associated sites that satisfy this criterion [25]. While FDR can be  
24 calculated by generating the empirical null distribution of test statistics [26], it is most commonly  
25 applied using the approach introduced by Benjamini and Hochberg [27]. This makes the assumption  
26 that under the null hypothesis the p-values across individual sites are uniformly distributed [28],  
27 which is not necessarily true. In EWAS it is not uncommon to see inflated test statistics [29, 30], even  
28 in the scenario of no true associations [31], indicating a skewed p-value distribution and perhaps

1 reflecting unaccounted confounders such as differences in cellular composition, or certain  
2 environmental exposures such as smoking. This variation in the distribution of p-values across studies  
3 means that the FDR approach often demonstrates variable behaviour making it challenging to  
4 compare results across studies. A better approach would be to estimate the number of independent  
5 tests performed in a EWAS and make the appropriate adjustment to the significance level. Saffari and  
6 colleagues have previously applied the methodology successfully used for GWAS to DNAm data  
7 profiled on the Illumina 450K array [32] in an attempt to establish a standard multiple testing  
8 threshold, however this is yet to be repeated for the EPIC array.

9

10 In this study, we used a large population based study, Understanding Society (n = 1,175), where DNA  
11 extracted from whole blood was profiled using the EPIC array[33, 34] to investigate potential  
12 statistical biases of DNAm association analyses, with the goal of providing recommendations for  
13 future epigenetic epidemiology studies. First, we used a permutation procedure to establish an  
14 appropriate significance threshold that accounts for the multiple testing burden of the EPIC array.  
15 Second, we investigated whether the assumptions of linear regression are satisfied when measuring  
16 DNAm as beta-values and whether any violations bias the results of DNAm studies. Although  
17 transformations of beta-values (e.g. conversion to M-values[18]) have been proposed in order to better  
18 satisfy the assumptions of linear regression, these approaches have not been unanimously adopted by  
19 the community therefore we seek to determine the validity of studies that analysed beta-values.  
20 Finally, we used the significance threshold derived from our simulations to explore the statistical  
21 power of DNAm studies across various scenarios. These results of our analyses will inform the  
22 optimal approach to designing and analysing DNAm data.

23

## 24 **Results**

25

### 26 *Estimating a multiple testing corrected significance threshold for the EPIC array*

27 After a stringent quality control (QC) pipeline (see **Methods**) and the exclusion of DNAm sites that  
28 may be technically biased by either the presence of genetic variants or cross-hybridisation to multiple

1 genomic loci, our final dataset included DNAm estimates for 804,826 sites across the autosomes and  
2 X chromosome derived from 1,175 individuals. Applying the Bonferroni correction formula for  
3 multiple testing, the significance threshold for hypothesis testing would be set to  $P < 6.21 \times 10^{-8}$   
4 ( $0.05/804,826$ ). In order to establish a significance threshold for EPIC array DNAm studies that  
5 controls for the number of independent tests (as opposed to the total number of sites tested), we used a  
6 permutation approach previously applied to GWAS [35] and 450K array DNAm studies [32]. This  
7 method preserves the correlation structure between sites and simulates null association studies by  
8 randomly assigning case control status. Repeating this process 100 times generates the distribution of  
9 p-values obtained by chance. From this distribution we calculated the 5% family-wise error rate  
10 (FWER) to be  $9.42 \times 10^{-8}$  (**Additional File 1: Figure S1**). Using the inverse of the Bonferroni  
11 correction formula this is equivalent to correcting for 530,639 independent tests ( $0.05/9.42 \times 10^{-8}$ ), a  
12 reduction of 34.1% compared to the total number of sites included in the analysis.

13

#### 14 *Estimating multiple testing corrected significance threshold for a genome-wide DNAm study*

15 As DNAm microarrays only profile a subset of the ~28 million potentially methylated sites in the  
16 human genome, the threshold calculated above is specific to an EPIC array-based experiment and  
17 hence we will refer to it as an “experiment-wide significance threshold”. Next, we were interested in  
18 using our permutations to extrapolate from this experiment-wide threshold to a significance threshold  
19 that accounts for all variation in DNAm across the genome. Given the correlation in DNAm between  
20 proximal DNAm sites, the content of the EPIC array provides some information about neighboring  
21 sites that are not directly profiled. Continuing to increase the genomic coverage of the microarray  
22 should, therefore, have diminishing returns in terms of novel association tests as we can use the sites  
23 present on the array to infer the status of other unmeasured neighboring sites. In order to model the  
24 information gain in terms of number of independent tests as the coverage of the microarray increases,  
25 we applied our permutation procedure to subsamples of DNAm sites at increasing densities ( $x_i = 5\%$ ,  
26  $15\%$ , ...,  $95\%$ ). For each density, we estimated the significance threshold 100 times and calculated the  
27 mean 5% FWER (denoted  $P_{T_i}$  for density  $i\%$ ). These estimated  $P_{T_i}$  values are plotted in **Figure 1a**,  
28 and demonstrate the expected monotonic, non-linear relationship where  $P_{T_i}$  becomes more significant



1 as the number of sites sampled increases. Each  $P_{Ti}$  value was then used to calculate the effective  
2 number of independent tests ( $m_i$ ) at density  $i\%$  using the inverse of Bonferroni formula ( $m_i =$   
3  $0.05/P_{Ti}$ ). Again, we observe a monotonic relationship where the effective number of tests increases as  
4 the proportion of sites sampled increases (**Figure 1b; Additional File 2 Table S1**). As the proportion  
5 of additional independent tests should decrease as the number of sites increases, this relationship is  
6 expected to be non-linear and converge to an asymptote which represents the total number of  
7 independent tests across the genome. These properties can be represented by the Monod function,  
8 which was originally proposed for the growth of microorganisms but is applicable to scenarios where  
9 subsequent growth is increasingly restricted over time. In this application, continually increasing the  
10 number of sites profiled in an experiment leads to smaller and smaller increments in the number of  
11 independent sites tested until all variation in DNAm is captured. This upper limit represents the total  
12 number of independent tests in the genome and is the value we want to estimate in order to determine  
13 the genome-wide multiple testing burden. We observe that this non-linear behaviour only starts to  
14 appear after ~600,000 sites. Fitting a Monod function to the subsampling results, we estimated the  
15 asymptote to be 5,803,067 (**Figure 2a**) reflecting the total number of independent tests across the  
16 DNA methylome. Compared to the total number of sites in the genome, this is a reduction of 79.3%.  
17 Calculating the Bonferroni corrected significance threshold based on this estimate gives a methylome-  
18 wide significance threshold of  $8.62 \times 10^{-9}$  ( $=0.05/5.80 \times 10^6$ ) (**Figure 2b**). Comparing this to a  
19 Bonferroni corrected significance threshold for all sites in the genome of  $1.79 \times 10^{-9}$  ( $0.05/2.8 \times 10^7$ ), our  
20 estimate is almost an order of magnitude smaller. The Monod function was also fitted to the  
21 subsample 95% confidence interval (CI) limits, estimating a 95% CI for the asymptote of  $1.69 \times 10^6$  to  
22  $3.36 \times 10^{13}$ , which equates to a 95% CI of  $2.97 \times 10^{-8}$  to  $1.49 \times 10^{-15}$  for the methylome-wide significance  
23 threshold.

24

### 25 *Testing the assumptions of linear regression for DNAm analyses*

26 To assess whether the assumptions of linear regression are satisfied, we performed an EWAS of age, a  
27 trait known to robustly co-vary with DNAm at multiple loci [21, 36]. The four assumptions of linear  
28 regression were assessed using four statistical tests implemented within the *gvlma* R package [37].

1 Specifically, these were tests for i) skewness, an asymmetrical distribution of the residuals, ii)  
2 kurtosis, a non-bell-shaped distribution of the residuals, iii) incorrect link function, a non-linear  
3 relationship between independent and dependent variables, and iv) heteroskedasticity, inconstant  
4 variance of the residuals (**Additional File 1: Figure S2**). In addition, a global test was performed  
5 providing an omnibus test of the four individual statistical tests. QQ plots of all five tests  
6 demonstrated dramatic inflation of p-values smaller than expected by chance (**Additional File 1:  
7 Figure S3**), indicating that the null hypothesis that DNAm data meets the assumptions of linear  
8 regression can be rejected for a large number of DNAm sites. Based on the experiment-wide  
9 significance threshold we previously derived for the EPIC array (i.e.  $P < 9.42 \times 10^{-8}$ ), 71.8% of sites  
10 rejected the null hypothesis for at least one assumption, with the majority of sites having non-normal  
11 residuals that exhibited evidence of excess skewness (41.3%) or excess kurtosis (67.6%) (**Table 1**).  
12 Furthermore, the specific DNAm sites whose residuals were skewed overlapped with the sites whose  
13 residuals were kurtotic (i.e. either highly or shallowly peaked) (**Figure 3**). A much smaller percentage  
14 of sites reject the null hypothesis in favour of a non-linear model (7.4%) or heteroskedasticity (4.3%).  
15

	<b>Global</b>	<b>Skewness</b>	<b>Kurtosis</b>	<b>Link Function</b>	<b>Heteroskedasticity</b>
<b>N reject null hypothesis</b>	577,919	332,457	544,460	59,572	35,001
<b>% reject null hypothesis</b>	71.8	41.3	67.6	7.4	4.3

16 **Table 1:** Summary of DNA methylation sites significantly rejecting the assumptions of linear  
17 regression. For each of the 5 tests performed by the *gvlma* package the number and percentage of  
18 DNA methylation sites with significant p-values ( $P < 9.42 \times 10^{-8}$ ) are reported.

19

### 20 *Characterising DNAm sites that infringe the assumptions of linear regression*

21 In order to propose guidelines for future EWAS studies, we were interested in whether DNAm sites  
22 that performed poorly in the *gvlma* tests could be characterized by common features such as DNAm  
23 level or variability. First, we considered the level of DNAm at each site, hypothesising that sites  
24 which are located at the extremes of the distribution would be more likely to violate the assumptions  
25 of the tests. We observed that the sites with the most significant p-values in the *gvlma* tests (i.e. those

1 with the largest  $-\log_{10}$  p-values) are generally either hypo- or hypermethylated (**Additional File 1:**  
2 **Figure S4**). Furthermore, by grouping sites based on their mean DNAm level we can pinpoint where  
3 in the distribution of DNAm values the assumptions are typically not satisfied. We observe a U-  
4 shaped relationship whereby sites with DNAm levels at the extremes (i.e. approaching 0 or 1), are  
5 more likely to violate the assumptions compared to sites with intermediate levels of DNAm (**Figure**  
6 **4; Additional File 2: Table S2**). This pattern generally holds for all four tests, but is most apparent  
7 for tests of skewness and kurtosis. Of interest, the relationship is not symmetrical, with the first two  
8 bins on the left of the distribution (containing sites with means of between 0 and 0.2) but only one bin  
9 on the far right of the distribution (containing sites with means of between 0.9 and 1.0) showing  
10 elevated mean  $-\log_{10}$  p-value compared to the middle seven bins. Second, we considered site  
11 variability, hypothesising that sites with low levels of variation would be more likely to violate the  
12 test assumptions. Using the standard deviation to index variability, we observed that sites with lower  
13 standard deviations had larger  $-\log_{10}$  p-values when testing the assumptions of linear regression  
14 (**Additional File 1: Figure S5**). This was most evident for the tests of skewness, kurtosis and  
15 heteroskedasticity, in particular for sites with a standard deviation  $< 0.02$  (**Figure 5; Additional File**  
16 **2: Table S3**). A more complex pattern was seen for the link function test, where the most variable  
17 probes and the second group of least variable probes had the highest  $-\log_{10}$  p-values. Using an  
18 alternative non-parametric method to characterize sites as ‘variable’ (range of middle 80% of values  
19  $>5\%$ ) or ‘non-variable’, we observed a similar pattern of results (**Additional File 1: Figure S6;**  
20 **Additional File 2: Table S4**) where non-variable sites were more likely to reject the assumptions of  
21 linear regression compared to variable sites. Taken together, these findings suggest that sites with  
22 extreme DNAm levels or low variation are most likely not to satisfy the assumptions of linear  
23 regression. These characteristics are not unrelated because sites with low levels of variation are  
24 typically located at the boundaries of the distribution of DNAm (**Additional File 1: Figure S7**).

25

26 Recently, M-values have been proposed as an alternative to beta-values in EWAS analyses of traits  
27 and exposures due to their more desirable statistical properties [18]. Although a direct comparison of  
28 beta-values and M-values is beyond the scope of this manuscript, we repeated our analyses on M-

1 values to further interpret the results presented above. Using our experiment-wide significance  
2 threshold, 70.1% of DNAm sites demonstrated significant bias of at least one assumption when using  
3 M-values; that is just 1.09% less than the original analysis based on beta-values (**Additional File 2:**  
4 **Table S5**). Furthermore, 85.9% of DNAm sites that are considered statistically inappropriate based on  
5 beta-values were also classed as statistically inappropriate when analysed as M-values. As with the  
6 beta-value analysis, the primary assumption violated by M-values related to the shape of the  
7 distribution of residuals. In fact, a comparable number of sites demonstrated excess kurtosis  
8 regardless of whether beta-values (67.6%) or M-values (66.7%) were used. Furthermore, albeit more  
9 subtly, DNAm sites with methylation levels at the extreme ends of the distribution were more likely to  
10 fail the statistical tests (**Additional File 1: Figure S8**), consistent with the results of the analysis using  
11 beta-values.

12

13 *Evaluating the impact on DNAm studies of sites that do not meet the assumptions of linear regression*  
14 The primary concern about using an invalid analytical model is the risk of either reporting false  
15 positive or false negative findings in tests of association. As linear regression is considered robust to  
16 violations of the assumptions, we next explored whether sites that violated an assumption were more  
17 likely to be significant in a DNAm analysis using a linear regression model. Using our simulated null  
18 association studies, DNAm sites were ranked by their association p-value to calculate the mean rank  
19 across the simulations. In a scenario where all sites are equally likely to be associated and there is no  
20 bias in the analysis, the distribution of these mean ranks should be symmetrical and unimodal with a  
21 mean of 402,413.5. Any skew in the distribution, or the presence of outliers and/or multiple peaks,  
22 would indicate an underlying bias in which DNAm sites are often identified as significant or not. We  
23 found that the distribution of the mean rank was normally distributed with a mean of 402,446  
24 (**Additional File 1: Figure S9**), similar to the expected value. We observed no association between p-  
25 values from the *gvlma* tests and a DNAm site's mean rank indicating that even highly significant  
26 rejections of the assumptions of linear regression do not bias EWAS results in terms of either false  
27 positives or false negatives (**Figure 6; Additional File 1: Figure S10; Additional File 2: Table S6**).

28

1 *Estimating the power of an EPIC array DNAm study*

2 The power of a test is defined as the probability that it correctly rejects the null hypothesis when the  
3 alternative hypothesis is true. As with other types of genomic analyses, large sample sizes are required  
4 for EWAS in order to obtain the statistical power required to identify a significant non-zero effect  
5 with a p-value that survives the adjustment for multiple testing. Having derived an appropriate  
6 multiple testing corrected significance threshold for the EPIC array, we investigated the typical  
7 sample sizes required for a DNAm study using this platform. In order to estimate power we need to  
8 know the sample size, multiple testing threshold, expected effect size and variance. While the first  
9 three of these parameters will remain constant for a particular study, the variance of DNAm will vary  
10 across sites. This means that a single power calculation, perhaps based on an average probe, provides  
11 limited information about the overall power of a DNAm study. We therefore performed a power  
12 calculation for each individual site on the EPIC array and then established the proportion of sites that  
13 surpass a specific power threshold. The estimated power for a single association test across a range of  
14 standard deviations and sample sizes for a binary phenotype (as would be tested in a disease case-  
15 control study) are shown in **Table 2**. For example, to detect a mean difference of 2% with 500 cases  
16 and 500 controls (total N = 1000), we have 100% power at sites with a standard deviation  $\leq 0.03$ .  
17 Performing separate power calculations tailored by the variance of each site, we plotted power curves  
18 for a range of typical DNAm studies (**Figure 7**). This analysis demonstrates that when N=200 (100  
19 cases and 100 controls), 85% of sites have >80% power to detect an effect of 5% (yellow line in  
20 **Figure 7b**), and when N=1000 (500 cases and 500 controls), 81% of probes have >80% power to  
21 detect an effect of 2% (light blue line in **Figure 7a**). While these examples provide a general  
22 overview of power for EPIC array studies, the results are also available for browsing in an interactive  
23 web application (<https://epigenetics.essex.ac.uk/shiny/EPICDNAmPowerCalcs/>) where the  
24 parameters can be adjusted in order to generate bespoke power calculations allowing researchers to  
25 assess the power of their individual study.

26

27

Sample Size	Standard Deviation					
	0.01	0.03	0.05	0.07	0.09	0.15
<b>Mean Difference = 2%</b>						
<b>100</b>	100%	1.26%	0.03%	0.00%	0.00%	0.00%
<b>200</b>	100%	21.4%	0.45%	0.04%	0.01%	0.00%
<b>500</b>	100%	97.8%	17.6%	1.43%	0.19%	0.01%
<b>1000</b>	100%	100%	82.7%	19.7%	3.22%	0.06%
<b>Mean Difference = 5%</b>						
<b>100</b>	100%	99.1%	24.4%	2.18%	0.29%	0.01%
<b>200</b>	100%	100%	93.0%	32.0%	6.07%	0.11%
<b>500</b>	100%	100%	100%	99.4%	78.4%	4.81%
<b>1000</b>	100%	100%	100%	100%	100%	45.8%

1 **Table 2: Summary of statistical power to significantly detect differential methylation between**  
2 **cases and controls.** Presented are example power calculations for a range of scenarios, varying effect  
3 size, sample size and variance for a binary phenotype. Power calculations are for a two-sided, two-  
4 sample t-test with a significance threshold of  $P < 9.42 \times 10^{-8}$ . The sample size is the total number of  
5 samples with a 50:50 split between groups.

## 7 Discussion

8 This study used a large DNAm dataset generated using the Illumina EPIC array to assess the  
9 statistical properties that influence the analytical design for hypothesis testing in epigenome-wide  
10 association studies. We estimated that there are 530,639 independent tests in a whole blood EPIC  
11 array DNAm study, which equates to a corrected significance threshold of  $9.42 \times 10^{-8}$ . For ease, we  
12 propose  $9 \times 10^{-8}$  would be an appropriate EPIC array experiment-wide significance threshold that  
13 should be adopted by the field to minimize the reporting of false positives. Although this EPIC array  
14 experiment-wide threshold is not substantially different to a Bonferroni correction for the actual  
15 number of tests, our estimate is comparable to that proposed using a similar methodology to data from

1 the older 450K array [32], which includes approximately half the number of sites ( $P = 2.4 \times 10^{-7}$ ) that  
2 were converted to M-values. Our results indicate that the correlation in DNAm across sites included  
3 on the Illumina EPIC array is relatively small and does not encompass large genomic regions; Saffari  
4 and colleagues also observed that strong correlations between neighboring sites were typically only  
5 observed within 1 kilobase [32], consistent with the minimal reduction from number of actual tests to  
6 number of independent tests we report. This challenges the argument that a Bonferroni correction is  
7 too conservative and therefore a more relaxed multiple testing threshold can be applied. Existing and  
8 future studies which report results at a more lenient threshold, particularly those with small sample  
9 sizes and lower statistical power should be interpreted with caution.

10

11 We attempted to extrapolate from the experiment-wide threshold for the EPIC array to estimate an  
12 appropriate threshold for all potential tests across the genome, including those not currently profiled  
13 by the EPIC array, by using simulations to profile how the number of independent tests changes as the  
14 coverage of the microarray increases. At sufficient density, the number of independent tests should  
15 plateau; however this behaviour was not really evident across the range of densities we were able to  
16 simulate, suggesting that the EPIC array does not interrogate a large part of the variation in DNAm  
17 across the genome. Therefore, our estimate of the number of independent tests in the genome is likely  
18 to be imprecise. Moreover, given the wide confidence interval around the estimated genome-wide  
19 multiple testing burden, we recommend this result is taken with some caution. Future large population  
20 based studies that include more DNAm sites across the genome would be required to address this  
21 question. We propose that our experiment-wide significance threshold should be adopted for all future  
22 EPIC array EWAS. The use of a standardized significance threshold would benefit the field by  
23 providing a common standard for reporting associations and facilitate the comparison of results across  
24 different studies. While the threshold has been determined to minimize the reporting of false  
25 positives, it does not prevent them entirely; prudent study design and effective control of confounders  
26 are still required for high quality EWAS studies. Furthermore, replication of associations across  
27 independent datasets is still required to validate robust associations.

28

1 We also tested the assumptions of linear regression, the most commonly used tool for identifying  
2 associations between differential DNAm and a trait, when measuring DNAm using beta-values (i.e. as  
3 a proportion) and conclude that the majority of sites do not satisfy the assumption of normally  
4 distributed errors. This was particularly the case for DNAm sites that have low levels of variation or  
5 are located at the extreme ends of the distribution. While we use our experiment-wide p-value  
6 threshold to quantify the number of probes not satisfying these assumptions in order to gauge the  
7 pattern of results, we caution against using this threshold to classify sites as passing or failing these  
8 assumptions. As the statistical evidence required to reject the null hypothesis in these tests is unlikely  
9 to equate to the degree of violation of the assumption needed to influence the results of the regression  
10 analysis, it may not follow that sites that fail these tests will lead to incorrect conclusions if a linear  
11 regression model is used. As these assumptions were tested on an EWAS of chronological age, it is  
12 possible that our results are specific to this particular analysis. Furthermore, we used a European adult  
13 whole blood cohort as a basis for our assessment, which may mean that the results are not applicable  
14 to studies of other tissues, cell-types, ages or ethnicities. It is also likely that these violations of these  
15 assumptions will be more important for studies based on smaller sample sizes. For these reasons,  
16 rather than report a list of DNAm sites that do not satisfy the assumptions, we focused on  
17 characterising these sites in order to provide general guidelines. Although the specific sites that not do  
18 vary within a sample may differ between studies, we predict that it is always the non-variable sites  
19 that fail the tests of the assumptions. Some studies remove non-variable sites prior to hypothesis  
20 testing[38, 39] [40] and our results support such a filtering step. However, as we found no evidence  
21 that the lack of normal residuals, an incorrectly specified link function, or heteroskedasticity leads to  
22 either false positive or false negative associations, our data also suggests that this is not strictly  
23 necessary. A number of studies have used transformations of beta-values, for example using log ratios  
24 of methylation percentage referred to as M-values in order to obtain a normal distribution[3, 18, 41]  
25 or regression based on an alternative distribution (e.g. beta regression[42]); our results show that the  
26 use of linear regression with beta values in DNAm studies, even if the data do not satisfy the standard  
27 assumptions of this test, does not appear to lead to biased results. Despite considering the four key  
28 assumptions of linear regression, we did not specifically investigate the effect of outlier DNAm



1 values, which may arise due to either technical or biological artefacts (e.g. rare SNP effects). The  
2 presence of outliers can introduce false positive associations as linear regression estimates are derived  
3 by minimising the sum of the residuals, therefore extreme values, which would lead to large residuals,  
4 can lead to larger, and therefore significant, estimated slope coefficients.

5

6 Finally, we performed power calculations to ascertain the sample size required for EPIC array studies  
7 using our proposed experiment-wide significance threshold. Most complex phenotypes are expected  
8 to be associated with small effects (typically < 5% difference between cases and controls), and our  
9 calculations indicate that with a sample size of 500 cases and 500 controls, 81% of sites have >80%  
10 power to detect an effect of 2%. This estimate should be reassuring to the epigenetic community, as  
11 there are an increasing number of studies approaching or surpassing this sample size [9, 43-46]. Our  
12 approach advances previous efforts[47] by taking into account the individual properties of each  
13 DNAm site and uses an empirically derived significance threshold to provide an overview of power  
14 across the EPIC array. Finally, we have developed an online tool  
15 (<https://epigenetics.essex.ac.uk/shiny/EPICDNAmPowerCalcs/> ) where users can perform their own  
16 bespoke calculations to quantify the power of their specific study for individual DNAm sites; we are  
17 currently extending this power calculation application for use with quantitative trait variables, and  
18 will implement an updated version in the near future.

19

## 20 **Conclusions**

21 We show that linear regression is a valid statistical methodology for DNAm studies, despite the fact  
22 that the data do not always satisfy the assumptions of the test. Additionally, we propose that a  
23 significance threshold of  $P < 9 \times 10^{-8}$  should be adopted to adequately control the false positive rate for  
24 EPIC array based analyses and should be accepted as a standard for reporting results. These findings  
25 have implications for epidemiological-based DNAm studies and provide a framework for the  
26 interpretation of findings from current and future studies.

27

## 28 **Methods**

1 All analyses were performed using the statistical language R [48].

2

### 3 *Genomic-wide profiling of DNAm in Understanding Society*

4 The DNAm dataset generated as part of the Understand Society study has been analysed in two  
5 previously published studies [34, 49] and a detailed description of the sample, DNAm data generation  
6 and data preprocessing can be found in the original publication [34]. Briefly, Understanding Society  
7 (<https://www.understandingsociety.ac.uk>) is a longitudinal panel survey of 40,000 UK households  
8 which has collected sociodemographic information, biomedical measures and blood samples from  
9 participants. DNA was extracted from whole blood samples to facilitate genomic profiling including  
10 DNAm.

11

### 12 *DNAm data preprocessing*

13 DNAm was profiled for a subset of 1,193 individuals from the Understanding Society study using the  
14 Illumina Infinium HumanMethylationEPIC BeadChip. Raw signal intensity data were processed from  
15 idat files through a standard pipeline using the bigmelon [49] and wateRmelon [50] packages in R. A  
16 number of quality control steps were performed to these data prior to normalization. First, outlier  
17 samples were identified using principal component analysis and mahalanobis distance equivalents,  
18 second, successful bisulphite conversion was confirmed using control probes, third the ages of the  
19 samples were estimated using the Horvath Epigenetic Clock algorithm [51] and compared to reported  
20 age at sampling, and fourth visualisation of principal components. These data were then normalized  
21 using the *dasen* method [50], which performs background adjustment and between-sample quantile  
22 normalization of methylated (M) and unmethylated (U) intensities separately for Type I and Type II  
23 probes. A second round of sample filtering was then performed excluding samples that were either  
24 dramatically altered as a result of normalisation or samples that had  $> 1\%$  of sites with detection p-  
25 value  $> 0.05$ . DNAm sites were filtered to exclude those with a bead count  $< 3$  or  $> 1\%$  of samples  
26 with detection p-value  $> 0.05$ . The raw DNAm data of the final sample set was then re-normalized  
27 with the *dasen* method. Prior to data analysis, SNP probes, probes with non-specific binding, probes  
28 affected by common SNPs [52], and 65 probes annotated to the Y chromosome were additionally

1 removed. The final dataset contained 1,175 individuals and 804,826 DNAm sites (787,400 annotated  
2 to autosomes, and 17,426 annotated to the X chromosome).

3

#### 4 *Estimating a significance threshold for DNAm studies using the EPIC array*

5 To estimate an experiment-wide significance threshold for the EPIC array, we applied the permutation  
6 procedure previously described by Dudbridge and Gusnanto [35]. For each permutation, 50% of our  
7 1,175 samples ( $n = 557$ ) were randomly assigned as “cases” and 50% ( $n = 558$ ) as “controls” to  
8 simulate a null EWAS (i.e. no differences between cases and controls). Each of the 804,826 sites was  
9 then tested for association with this simulated phenotype using a linear regression model including  
10 sex, age, and six estimated cellular composition variables (B cells, CD8 T cells, CD4 T cells,  
11 monocytes, granulocytes, natural killer T cells) [53, 54] as covariates. We repeated this procedure  
12 1000 times recording the smallest p-value (i.e. the most significant) from each permutation. The  
13 EPIC array significance threshold was estimated by taking the 5<sup>th</sup> percentile point of these 1000  
14 minimum p-values representing the 5% family-wise error rate (FWER).

15

#### 16 *Estimating a genome-wide significance threshold for DNAm studies*

17 In order to extrapolate from our experiment-wide significance thresholds to one appropriate for  
18 genome wide DNAm association studies, we implemented the subsampling procedure also  
19 implemented by Dudbridge and Gusnanto [35]. Briefly, to simulate experiments with a reduced  
20 number of sites that capture a smaller proportion of genome-wide variation, sites were randomly  
21 subsampled at a range of densities ( $x_i=5\%, 15\%, \dots, 95\%; i=1, 2, \dots, 10$ ). From each permutation, the  
22 smallest p-value across the subset of sites was extracted and the 5<sup>th</sup> percentile point across all 1000  
23 minimum p-values was recorded. This subsampling was repeated 100 times and the mean, 2.5 and  
24 97.5 percentile points were calculated to set the significance threshold ( $P_{T_i}$ ) and confidence intervals  
25 for density  $i$ . At low densities, where the coverage is sparse, it is assumed that all included DNAm  
26 sites will be independent and a Bonferroni correction for multiple testing is appropriate. As coverage  
27 increases, correlations between neighboring sites mean that the number of additional independent tests  
28 decreases. In other words, continually increasing the number of sites studied has diminishing returns

1 in terms of the increase in additional variation captured. Therefore, as the number of sites profiled in  
2 an experiment increases, the effective number of independent tests converges to an asymptote. To  
3 estimate the value of this asymptote, we fitted a Monod function across the site densities and their  
4 estimated number of independent tests. For each of the site densities ( $x_i$ ), the effective number of  
5 independent tests ( $m_i$ ), was calculated by using the inverse of the Bonferroni correction for multiple  
6 testing ( $m_i = 0.05/P_{Ti}$ ). A Monod function, originally a mathematical model for bacterial population  
7 growth with limited resources, takes the form:

$$8 \quad m = f(x, u, k) = \frac{ux}{k + x}$$

9 where  $u$  is the limit and  $k$  is the half-saturation parameter, their values given by:

$$10 \quad f(k) = \frac{u}{2}$$

$$11 \quad f(\infty) = u$$

12 This function was fitted using a least squares approach in R to find the value of  $u$ , which represents  
13 the number of independent tests in the entire DNA methylome. To calculate the methylome-wide  
14 significance threshold we applied the Bonferroni correction using this estimate ( $P_{\text{genome}} = 0.05/u$ ).

15

### 16 *Testing the assumptions of linear regression models used in DNAm studies*

17 To assess the validity of linear regression models in studies of DNAm, an EWAS of age was  
18 performed including sex, processing chip and six estimated cellular composition variables (B cells,  
19 CD8 T cells, CD4 T cells, monocytes, granulocytes, natural killer T cells) [53, 54] as covariates. For  
20 each of the 804,826 models (one per DNAm site) we tested for violations of the assumptions of linear  
21 regression using the *gvlma* (Global Validation of Linear Model Assumptions) R package [37]. This  
22 package performs four tests to test the performance of the model fit with regards to the four  
23 assumptions of a linear regression: linearity, homoskedasticity, uncorrelatedness and normality of the  
24 residuals (**Additional File 1: Figure S2**). The *gvlma* package provides a numerical measure of  
25 violation through significance testing for skewness, kurtosis, link function, and heteroskedasticity.  
26 Briefly, the package calculates a directional test statistic for each assumption using the standardized

1 residuals from the fitted linear model. These test statistics are each compared to a 1 degree-of-  
2 freedom chi-square distribution to calculate a p-value for hypothesis testing. In addition to obtaining a  
3 p-value for each of these four tests, the software also generates a “global” p-value, which is an  
4 omnibus test of the four others. The global test statistic is the sum of the four components (one for  
5 each assumption) and compared to a 4 degree-of-freedom chi-square distribution. The formula for  
6 each component and further details can be found in the original manuscript proposing the method  
7 [37]. The null hypothesis for the global test is that all four assumptions hold, and the alternative  
8 hypothesis is that at least one does not (i.e. a significant p-value indicates that a linear model is not  
9 appropriate). In order to assess how DNAm sites on the EPIC array performed across these five tests  
10 we plotted Quantile-Quantile (QQ) plots of the observed vs expected p-values. To characterize sites  
11 which perform poorly in these tests we visualized correlations between the p-values from the five  
12 *gvlma* tests and both the mean level of DNAm and two measures of variance (standard deviation and  
13 range of the middle 80% of values). For the purpose of assessing which assumptions are most  
14 commonly violated, and which are most commonly violated simultaneously, we applied the  
15 experiment wide p-value threshold derived in the previous sections ( $P < 9.42 \times 10^{-8}$ ), to identify sites  
16 that reject the assumptions of linear regression. Finally to investigate the impact of violating the  
17 assumptions of linear regression, we calculated the mean rank across the 1000 null EWAS  
18 permutations as an indicator of how likely a site was to be associated by chance and any bias in  
19 association analyses. These mean ranks were then compared with the p-values of the *gvlma* tests.

20

### 21 *Estimating statistical power for EPIC array studies*

22 Power calculations were performed for each of the 804,826 sites in the dataset using the function  
23 *pwr.t.test* from the R package *pwr* [55]. We consider the scenario with a binary outcome (i.e. case  
24 control study), using a two-sample t-test to compare the means of the two groups where the null  
25 hypothesis of each test is that the means of the two groups are equal. To calculate power, the  
26 parameters sample size, effect size and significance level were provided. The significance level was  
27 set as our previously calculated experiment-wide threshold of  $9.42 \times 10^{-8}$ . The effect size was provided  
28 as Cohen’s d, which is the expected difference between the two group means divided by their pooled

1 standard deviation [56]. In order to get a power estimate for the overall study, calculations were  
2 performed for every site individually using that site's variance estimated from the Understanding  
3 Society dataset, for two different mean differences (2% and 5%). Power calculations were also  
4 performed for a range of total sample sizes ( $n = 100, 200, 500, 1000, 2000$  and  $5000$ ) consisting of  
5 equal numbers of cases and controls. For each combination of parameters (sample size and mean  
6 difference), we calculated the percentage of sites that had sufficient statistical power across the full  
7 range of possible values (0-100%). While we only present results for a subset of the possible scenarios  
8 as a guide to the power of a typical EWAS study, we have also developed an R shiny app [7] to allow  
9 users to perform bespoke power calculations  
10 (<https://epigenetics.essex.ac.uk/shiny/EPICDNAmPowerCalcs/>). In this app, the user can specify  
11 sample size and mean difference to generate a summary results table and downloadable figure. As  
12 performing  $>800,000$  power calculations is time consuming, the app uses a binning method, grouping  
13 sites with similar variances and plotting a smoothed curve, to speed up the calculation. For more  
14 accurate results the user can increase the number of bins, or chose to calculate the power for all sites  
15 individually. There is also an option to search for a specific DNAm site of interest and calculate its  
16 power under user defined parameters.

17

## 18 **Figure Legends**

19

### 20 **Figure 1: Subsampling sites on the EPIC array to estimate a genome-wide significance**

21 **threshold.** Line graphs depicting the relationship between the number of EPIC array DNA  
22 methylation sites (x-axis) and A) the 5% family-wise error rate (FWER) ( $-\log_{10}(\text{p-values})$ ; y-axis) and  
23 B) the mean effective number of tests (y-axis) estimated from 1000 simulated null association studies.  
24 Error bars present the 95% confidence intervals from 1000 simulations. The final point includes all  
25 DNA methylation sites on the EPIC array and therefore could not be resampled to generate a  
26 confidence interval.

27

1 **Figure 2: Extrapolation to a genome-wide significance threshold.** Line graphs depicting the  
2 relationship between the number of DNA methylation sites (x-axis) and **A**) the effective number of  
3 independent tests (y-axis) and **B**) the multiple testing corrected threshold ( $-\log_{10}(\text{p-value})$ ; y-axis)  
4 estimated after fitting a Monod function to the observed data presented in Figure 1B. The observed  
5 values are plotted as the solid black line, and the estimated Monod model is plotted as a dashed line.  
6 The grey shaded region represents the 95% CI created by fitting a Monod model to the 95% CI of the  
7 subsampled data. The blue horizontal line represents the estimated asymptote of the Monod model of  
8 5,803,067 independent tests equivalent to a genome-wide significance threshold of  $8.62 \times 10^{-9}$ .

9  
10 **Figure 3: Overlap of significant violations of linear regression assumptions.** Venn diagram  
11 depicting the overlap of DNA methylation sites significant for each test of a linear assumption ( $P <$   
12  $9.42 \times 10^{-8}$ ). Presented are the number of overlapping DNA methylation sites along with the percentage  
13 of all tested sites shown in brackets.

14  
15 **Figure 4: Comparison of tests of linear regression assumptions across the distribution of DNA**  
16 **methylation levels.** Boxplots of  $-\log_{10}(\text{p-value})$  for each of the 5 tests (a) global (b) skewness (c)  
17 kurtosis (d) link function and (e) heteroskedasticity for groups of DNA methylation sites binned by  
18 their mean DNA methylation level. The boxes are coloured by their mean  $-\log_{10}(\text{p-value})$  from light  
19 yellow (low) to red (high).

20  
21 **Figure 5: Comparison of tests of linear regression assumptions against DNA methylation**  
22 **variability.** Boxplots of  $-\log_{10}(\text{p-value})$  for each of the 5 tests (a) global (b) skewness (c) kurtosis (d)  
23 link function and (e) heteroscedasticity for groups of DNA methylation sites binned by their standard  
24 deviation. The boxes are coloured by their mean  $-\log_{10}(\text{p-value})$  from light yellow (low) to red (high).

25

1 **Figure 6: Comparison of tests of linear regression assumptions with bias in DNA methylation**  
2 **association studies.** Scatterplots of  $-\log_{10}(\text{p-value})$  (y-axis) from the (a) global (b) skewness (c)  
3 kurtosis (d) link function and (e) heteroskedasticity tests performed in the R *gvlma* package against  
4 average (mean) ranking from 1000 simulated null association studies (x-axis) for all DNA  
5 methylation sites. Each point represents a single site, and the color represents the density of points  
6 plotted at that position (low density in grey to high density in yellow).

7

8 **Figure 7: Power curves of typical DNA methylation studies.** Line graphs depicting the proportion  
9 of sites on the EPIC array (y-axis) with sufficient power (x-axis) to detect a mean difference in DNA  
10 methylation between two groups of (a) 2% and (b) 5%. The different coloured lines represent  
11 different sample sizes where the value of N the total sample size set to be a 50:50 split between  
12 groups.

13

14



1 **Supplemental Data**

2

3 **Additional File 1:** Supplementary figures S1-S10 (pdf)

4 **Additional File 2:** Supplementary tables S1-S6 (pdf)

5

6 **Declarations**

7 *Ethics approval and consent to participate*

8 Ethical approval for the Understanding Society nurse visit was obtained from the National Research  
9 Ethics Service (Reference: 10/H0604/2). Participants gave written consent for blood sampling.

10 *Consent for publication*

11 Not applicable

12 *Availability of data and material*

13 Individual level DNA methylation are available on application through the European Genome-  
14 phenome Archive under accession EGAS00001001232 (<https://www.ebi.ac.uk/ega/home>). Specific  
15 details can be found here (<https://www.understandingsociety.ac.uk/about/health/data>). Phenotype  
16 linked to DNA methylation data are available through application to the METADAC  
17 ([www.metadac.ac.uk](http://www.metadac.ac.uk)). Analysis scripts used in this manuscript are available on  
18 <https://github.com/ejh243/EPICStatsPaper>.

19 *Competing interests*

20 The authors declare that they have no competing interests

21 *Funding*

22 DNA methylation data generation in UKHLS was funded through enhancements to the Economic and  
23 Social Research council (ESRC) grants ES/K005146/1 and ES/N00812X/1. MK is supported by the

1 University of Essex and ESRC (RES-596-28-0001). EH, JM and LS time on this project was  
2 supported by MRC grant K013807. The funding bodies played no role in the design of the study, data  
3 collection, analysis, and interpretation or in the writing of the manuscript.

#### 4 *Authors' contributions*

5 EH designed the study. GM undertook primary statistical analyses with input from TG, YB, MK, LS,  
6 JM and EH. GM and EH drafted the manuscript. All authors read and approved the final manuscript.

#### 7 *Acknowledgements*

8 Analysis was facilitated by access to the Genome high performance computing cluster at the  
9 University Of Essex School Of Biological Sciences.

10

#### 11 **References**

- 12 1. Murphy TM, Mill J: **Epigenetics in health and disease: heralding the EWAS era.** *Lancet* 2014,  
13 **383**(9933):1952-1954.
- 14 2. Heyn H, Carmona FJ, Gomez A, Ferreira HJ, Bell JT, Sayols S, Ward K, Stefansson OA, Moran  
15 S, Sandoval J *et al*: **DNA methylation profiling in breast cancer discordant identical twins**  
16 **identifies DOK7 as novel epigenetic biomarker.** *Carcinogenesis* 2013, **34**(1):102-108.
- 17 3. Irizarry RA, Ladd-Acosta C, Wen B, Wu Z, Montano C, Onyango P, Cui H, Gabo K, Rongione M,  
18 Webster M *et al*: **The human colon cancer methylome shows similar hypo- and**  
19 **hypermethylation at conserved tissue-specific CpG island shores.** *Nat Genet* 2009,  
20 **41**(2):178-186.
- 21 4. Lange CP, Campan M, Hinoue T, Schmitz RF, van der Meulen-de Jong AE, Slingerland H, Kok  
22 PJ, van Dijk CM, Weisenberger DJ, Shen H *et al*: **Genome-scale discovery of DNA-**  
23 **methylation biomarkers for blood-based detection of colorectal cancer.** *PLoS One* 2012,  
24 **7**(11):e50266.
- 25 5. Liu Y, Aryee MJ, Padyukov L, Fallin MD, Hesselberg E, Runarsson A, Reinius L, Acevedo N,  
26 Taub M, Ronninger M *et al*: **Epigenome-wide association data implicate DNA methylation**  
27 **as an intermediary of genetic risk in rheumatoid arthritis.** *Nat Biotechnol* 2013, **31**(2):142-  
28 147.
- 29 6. Rakyán VK, Beyan H, Down TA, Hawa MI, Maslau S, Aden D, Daunay A, Busato F, Mein CA,  
30 Manfras B *et al*: **Identification of type 1 diabetes-associated DNA methylation variable**  
31 **positions that precede disease diagnosis.** *PLoS Genet* 2011, **7**(9):e1002300.
- 32 7. Murphy TM, Wong CC, Arseneault L, Burrage J, Macdonald R, Hannon E, Fisher HL, Ambler A,  
33 Moffitt TE, Caspi A *et al*: **Methylomic markers of persistent childhood asthma: a**  
34 **longitudinal study of asthma-discordant monozygotic twins.** *Clin Epigenetics* 2015, **7**:130.
- 35 8. Pidsley R, Viana J, Hannon E, Spiers HH, Troakes C, Al-Saraj S, Mechawar N, Turecki G,  
36 Schalkwyk LC, Bray NJ *et al*: **Methylomic profiling of human brain tissue supports a**  
37 **neurodevelopmental origin for schizophrenia.** *Genome Biol* 2014, **15**(10):483.

- 1 9. Hannon E, Dempster E, Viana J, Burrage J, Smith AR, Macdonald R, St Clair D, Mustard C,  
2 Breen G, Therman S *et al*: **An integrated genetic-epigenetic analysis of schizophrenia:  
3 evidence for co-localization of genetic associations and differential DNA methylation.**  
4 *Genome Biol* 2016, **17**(1):176.
- 5 10. Ladd-Acosta C, Hansen KD, Briem E, Fallin MD, Kaufmann WE, Feinberg AP: **Common DNA  
6 methylation alterations in multiple brain regions in autism.** *Mol Psychiatry* 2014, **19**(8):862-  
7 871.
- 8 11. Berko ER, Suzuki M, Beren F, Lemetre C, Alaimo CM, Calder RB, Ballaban-Gil K, Gounder B,  
9 Kampf K, Kirschen J *et al*: **Mosaic epigenetic dysregulation of ectodermal cells in autism  
10 spectrum disorder.** *Plos Genet* 2014, **10**(5):e1004402.
- 11 12. De Jager PL, Srivastava G, Lunnon K, Burgess J, Schalkwyk LC, Yu L, Eaton ML, Keenan BT,  
12 Ernst J, McCabe C *et al*: **Alzheimer's disease: early alterations in brain DNA methylation at  
13 ANK1, BIN1, RHBDL2 and other loci.** *Nat Neurosci* 2014, **17**(9):1156-1163.
- 14 13. Lunnon K, Smith R, Hannon E, De Jager PL, Srivastava G, Volta M, Troakes C, Al-Sarraj S,  
15 Burrage J, Macdonald R *et al*: **Methylomic profiling implicates cortical deregulation of ANK1  
16 in Alzheimer's disease.** *Nat Neurosci* 2014, **17**(9):1164-1170.
- 17 14. Mill J, Heijmans BT: **From promises to practical strategies in epigenetic epidemiology.** *Nat  
18 Rev Genet* 2013, **14**(8):585-594.
- 19 15. Relton CL, Davey Smith G: **Epigenetic epidemiology of common complex disease: prospects  
20 for prediction, prevention, and treatment.** *PLoS Med* 2010, **7**(10):e1000356.
- 21 16. Rakyan VK, Down TA, Balding DJ, Beck S: **Epigenome-wide association studies for common  
22 human diseases.** *Nat Rev Genet* 2011, **12**(8):529-541.
- 23 17. Laird PW: **Principles and challenges of genomewide DNA methylation analysis.** *Nat Rev  
24 Genet* 2010, **11**(3):191-203.
- 25 18. Du P, Zhang X, Huang CC, Jafari N, Kibbe WA, Hou L, Lin SM: **Comparison of Beta-value and  
26 M-value methods for quantifying methylation levels by microarray analysis.** *BMC  
27 Bioinformatics* 2010, **11**:587.
- 28 19. Birney E, Smith GD, Grealis JM: **Epigenome-wide Association Studies and the Interpretation  
29 of Disease -Omics.** *PLoS Genet* 2016, **12**(6):e1006105.
- 30 20. Panni T, Mehta AJ, Schwartz JD, Baccarelli AA, Just AC, Wolf K, Wahl S, Cyrus J, Kunze S,  
31 Strauch K *et al*: **A Genome-Wide Analysis of DNA Methylation and Fine Particulate Matter  
32 Air Pollution in Three Study Populations: KORA F3, KORA F4, and the Normative Aging  
33 Study.** *Environ Health Perspect* 2016.
- 34 21. Spiers H, Hannon E, Schalkwyk LC, Smith R, Wong CC, O'Donovan MC, Bray NJ, Mill J:  
35 **Methylomic trajectories across human fetal brain development.** *Genome Res* 2015,  
36 **25**(3):338-352.
- 37 22. Cardenas A, Houseman EA, Baccarelli AA, Quamruzzaman Q, Rahman M, Mostofa G, Wright  
38 RO, Christiani DC, Kile ML: **In utero arsenic exposure and epigenome-wide associations in  
39 placenta, umbilical artery, and human umbilical vein endothelial cells.** *Epigenetics* 2015,  
40 **10**(11):1054-1063.
- 41 23. Elliott HR, Tillin T, McArdle WL, Ho K, Duggirala A, Frayling TM, Davey Smith G, Hughes AD,  
42 Chaturvedi N, Relton CL: **Differences in smoking associated DNA methylation patterns in  
43 South Asians and Europeans.** *Clin Epigenetics* 2014, **6**(1):4.
- 44 24. Ong ML, Holbrook JD: **Novel region discovery method for Infinium 450K DNA methylation  
45 data reveals changes associated with aging in muscle and neuronal pathways.** *Aging Cell*  
46 2014, **13**(1):142-155.
- 47 25. Heyn H, Moran S, Hernando-Herraez I, Sayols S, Gomez A, Sandoval J, Monk D, Hata K,  
48 Marques-Bonet T, Wang L *et al*: **DNA methylation contributes to natural human variation.**  
49 *Genome Res* 2013, **23**(9):1363-1372.
- 50 26. Noble WS: **How does multiple testing correction work?** *Nat Biotechnol* 2009, **27**(12):1135-  
51 1137.

- 1 27. Benjamini Y, Hochberg Y: **Controlling the False Discovery Rate: A Practical and Powerful**  
2 **Approach to Multiple Testing.** *Journal of the Royal Statistical Society* 1995, **57**(1):289-300.
- 3 28. Moskvina V, Schmidt KM: **On multiple-testing correction in genome-wide association**  
4 **studies.** *Genet Epidemiol* 2008, **32**(6):567-573.
- 5 29. Zaghlool SB, Al-Shafai M, Al Muftah WA, Kumar P, Falchi M, Suhre K: **Association of DNA**  
6 **methylation with age, gender, and smoking in an Arab population.** *Clin Epigenetics* 2015,  
7 **7**(1):6.
- 8 30. Absher DM, Li X, Waite LL, Gibson A, Roberts K, Edberg J, Chatham WW, Kimberly RP:  
9 **Genome-wide DNA methylation analysis of systemic lupus erythematosus reveals**  
10 **persistent hypomethylation of interferon genes and compositional changes to CD4+ T-cell**  
11 **populations.** *Plos Genet* 2013, **9**(8):e1003678.
- 12 31. Lehne B, Drong AW, Loh M, Zhang W, Scott WR, Tan ST, Afzal U, Scott J, Jarvelin MR, Elliott P  
13 *et al*: **A coherent approach for analysis of the Illumina HumanMethylation450 BeadChip**  
14 **improves data quality and performance in epigenome-wide association studies.** *Genome*  
15 *Biol* 2015, **16**:37.
- 16 32. Saffari A, Silver MJ, Zavattari P, Moi L, Columbano A, Meaburn EL, Dudbridge F: **Estimation**  
17 **of a significance threshold for epigenome-wide association studies.** *Genet Epidemiol* 2018,  
18 **42**(1):20-33.
- 19 33. Gorrie-Stone TJ, Smart MC, Saffari A, Malki K, Hannon E, Burrage J, Mill J, Kumari M,  
20 Schalkwyk LC: **Bigmelon: Tools for analysing large DNA methylation datasets.**  
21 *Bioinformatics* 2018:bty713-bty713.
- 22 34. Hannon E, Gorrie-Stone TJ, Smart MC, Burrage J, Hughes A, Bao Y, Kumari M, Schalkwyk LC,  
23 Mill J: **Leveraging DNA-Methylation Quantitative-Trait Loci to Characterize the Relationship**  
24 **between Methylomic Variation, Gene Expression, and Complex Traits.** *Am J Hum Genet*  
25 2018, **103**(5):654-665.
- 26 35. Dudbridge F, Gusnanto A: **Estimation of significance thresholds for genomewide**  
27 **association scans.** *Genet Epidemiol* 2008, **32**(3):227-234.
- 28 36. Bell JT, Tsai PC, Yang TP, Pidsley R, Nisbet J, Glass D, Mangino M, Zhai G, Zhang F, Valdes A *et*  
29 *al*: **Epigenome-wide scans identify differentially methylated regions for age and age-**  
30 **related phenotypes in a healthy ageing population.** *PLoS Genet* 2012, **8**(4):e1002629.
- 31 37. Peña EA, Slate EH: **Global Validation of Linear Model Assumptions.** *J Am Stat Assoc* 2006,  
32 **101**(473):341.
- 33 38. Glossop JR, Nixon NB, Emes RD, Haworth KE, Packham JC, Dawes PT, Fryer AA, Matthey DL,  
34 Farrell WE: **Epigenome-wide profiling identifies significant differences in DNA methylation**  
35 **between matched-pairs of T- and B-lymphocytes from healthy individuals.** *Epigenetics*  
36 2013, **8**(11):1188-1197.
- 37 39. Fryer AA, Emes RD, Ismail KM, Haworth KE, Mein C, Carroll WD, Farrell WE: **Quantitative,**  
38 **high-resolution epigenetic profiling of CpG loci identifies associations with cord blood**  
39 **plasma homocysteine and birth weight in humans.** *Epigenetics* 2011, **6**(1):86-94.
- 40 40. Gao Z, Fu HJ, Zhao LB, Sun ZY, Yang YF, Zhu HY: **Aberrant DNA methylation associated with**  
41 **Alzheimer's disease in the superior temporal gyrus.** *Exp Ther Med* 2018, **15**(1):103-108.
- 42 41. Ladd-Acosta C, Hansen K, Briem E, Fallin M, Kaufmann W, Feinberg A: **Common DNA**  
43 **methylation alterations in multiple brain regions in autism.** *Molecular Psychiatry* 2014,  
44 **19**(8):862-871.
- 45 42. Triche TJ, Laird PW, Siegmund KD: **Beta regression improves the detection of differential**  
46 **DNA methylation for epigenetic epidemiology.** *bioRxiv* 2016:054643.
- 47 43. Hannon E, Schendel D, Ladd-Acosta C, Grove J, Hansen CS, Andrews SV, Hougaard DM,  
48 Bresnahan M, Mors O, Hollegaard MV *et al*: **Elevated polygenic burden for autism is**  
49 **associated with differential DNA methylation at birth.** *Genome Med* 2018, **10**(1):19.

- 1 44. Wahl S, Drong A, Lehne B, Loh M, Scott WR, Kunze S, Tsai PC, Ried JS, Zhang W, Yang Y *et al*:  
2 **Epigenome-wide association study of body mass index, and the adverse outcomes of**  
3 **adiposity**. *Nature* 2017, **541**(7635):81-86.
- 4 45. Joehanes R, Just AC, Marioni RE, Pilling LC, Reynolds LM, Mandaviya PR, Guan W, Xu T, Elks  
5 CE, Aslibekyan S *et al*: **Epigenetic Signatures of Cigarette Smoking**. *Circ Cardiovasc Genet*  
6 2016, **9**(5):436-447.
- 7 46. Joubert BR, Felix JF, Yousefi P, Bakulski KM, Just AC, Breton C, Reese SE, Markunas CA,  
8 Richmond RC, Xu CJ *et al*: **DNA Methylation in Newborns and Maternal Smoking in**  
9 **Pregnancy: Genome-wide Consortium Meta-analysis**. *Am J Hum Genet* 2016, **98**(4):680-696.
- 10 47. Tsai PC, Bell JT: **Power and sample size estimation for epigenome-wide association scans to**  
11 **detect differential DNA methylation**. *Int J Epidemiol* 2015.
- 12 48. R Development Core Team: **R: A Language and Environment for Statistical Computing**. In.  
13 Vienna, Austria: R Foundation for Statistical Computing; 2008.
- 14 49. Gorrie-Stone TJ, Smart MC, Saffari A, Malki K, Hannon E, Burrage J, Mill J, Kumari M,  
15 Schalkwyk LC: **Bigmelon: tools for analysing large DNA methylation datasets**. *Bioinformatics*  
16 2019, **35**(6):981-986.
- 17 50. Pidsley R, Wong CCY, Volta M, Lunnon K, Mill J, Schalkwyk LC: **A data-driven approach to**  
18 **preprocessing Illumina 450K methylation array data**. *Bmc Genomics* 2013, **14**.
- 19 51. Horvath S: **DNA methylation age of human tissues and cell types**. *Genome Biol* 2013,  
20 **14**(10):R115.
- 21 52. McCartney DL, Walker RM, Morris SW, M. MA, J. PD, L. EK: **Identification of polymorphic**  
22 **and off-target probe binding sites on the Illumina Infinium MethylationEPIC BeadChip**.  
23 *Genomics Data* 2016, **9**(September):22-24.
- 24 53. Houseman EA, Accomando WP, Koestler DC, Christensen BC, Marsit CJ, Nelson HH, Wiencke  
25 JK, Kelsey KT: **DNA methylation arrays as surrogate measures of cell mixture distribution**.  
26 *BMC Bioinformatics* 2012, **13**:86.
- 27 54. Koestler DC, Christensen B, Karagas MR, Marsit CJ, Langevin SM, Kelsey KT, Wiencke JK,  
28 Houseman EA: **Blood-based profiles of DNA methylation predict the underlying**  
29 **distribution of cell types: a validation analysis**. *Epigenetics* 2013, **8**(8):816-826.
- 30 55. Champely S: **pwr: Basic Functions for Power Analysis**. In., 1.2-2 edn. [https://CRAN.R-](https://CRAN.R-project.org/package=pwr)  
31 [project.org/package=pwr](https://CRAN.R-project.org/package=pwr); 2018.
- 32 56. Cohen J: **Statistical power analysis for the behavioral sciences** . , Second edn. Hillsdale, N.J.:  
33 Lawrence Erlbaum Associates; 1988.

34