

A comparative study of methods for identifying epigenetic aging moderators

Colin Farrell¹, Kalsuda Lapborisuth¹, Chanyue Hu¹, Kyle Pu¹, Sagi Snir², and Matteo Pellegrini^{1,3}

¹Dept. of Molecular, Cell and Developmental Biology;
University of California, Los Angeles, CA 90095, USA;;

²Dept. of Evolutionary Biology, University of Haifa, Israel;

³Corresponding Author, matteop@mcdb.ucla.edu

1

2 Epigenetic clocks, DNA methylation based chronological age prediction
3 models, are commonly employed to study age related biology. The error
4 between the predicted and observed age is often interpreted as a form of
5 biological age acceleration and many studies have measured the impact
6 of environmental and other factors on epigenetic age. Epigenetic clocks
7 are fit using approaches that minimize the error between the predicted
8 and observed chronological age and as a result they reduce the impact of
9 factors that may moderate the relationship between actual and epigenetic
10 age. Here we compare the standard methods used to construct epigenetic
11 clocks to an evolutionary framework of epigenetic aging, the epigenetic
12 pacemaker (EPM) that directly models DNA methylation as a function
13 of a time dependent epigenetic state. We show that the EPM is more
14 sensitive than epigenetic clocks for the detection of factors that moderate
15 the relationship between actual age and epigenetic state (or epigenetic
16 age). Specifically, we show that the EPM is more sensitive at detecting
17 sex and cell type effects in a large aggregate dataset and in an example
18 case study is more sensitive at detecting age-related methylation
19 changes associated with polybrominated biphenyl exposure. Thus we find
20 that the pacemaker provides a more robust framework for the study of
21 factors that impact epigenetic age acceleration than traditional clocks based
22 on linear regression models.

23

24 1 Introduction

25 Epigenetic clocks, accurate age prediction models made using DNA methylation,
26 are promising tools for the study of aging and age related biology.
27 Beyond predicting the age of an individual to within a couple of years,
28 multiple studies have shown that the difference between the observed and
29 expected epigenetic age can be interpreted as a measure of biological age
30 acceleration [1]. Age acceleration observed using the first generation of epi-
31 genetic clocks [2, 3] has been associated with a variety of health outcomes
32 including mortality risk[4, 5], cancer risk [6], cardiovascular disease[7] and
33 other negative health outcomes[8–10]. However, as epigenetic clocks be-
34 come more accurate, epigenetic age acceleration is no longer associated
35 with mortality [11].

Epigenetic clocks are generally trained using a regularized regression model. Given an elastic net model of the form $y = \beta X$ the goal of penalized regression is to maximize the likelihood by reducing the prediction error of the model, $L(\lambda_1, \lambda_2, \beta) = |y - X\beta|^2 + \lambda_2|\beta|^2 + |\lambda_1\beta|$. In the case of epigenetic clocks, the likelihood is maximized by minimizing the difference between the observed and predicted age subject to the elastic net penalty, λ_1 and λ_2 . . Methylation sites that increase modeled error but contain biologically meaningful information may be discarded during model fitting. This problem is magnified in the case of epigenetic clocks where the relationship between methylation and time is nonlinear[12].

An alternative and complementary approach to studying epigenetic aging is to model how methylation changes for a predetermined collection of sites with respect to time. To this end, we have developed the epigenetic pacemaker (EPM) [13, 14] to model methylation changes with age. Given j individuals and i methylation sites, under the EPM an individual methylation site can be modeled as $\hat{m}_{ij} = m_i^0 + r_i s_j + \epsilon_{ij}$ where \hat{m}_{ij} is the observed methylation value, m_i^0 is the initial methylation value, r_i is the rate of change, s_j is the epigenetic state, and ϵ_{ij} is a normally distributed error term. The r_i and m_i^0 are characteristic of the sites across all individuals and the epigenetic state of an individual s_j is set using information from all modeled sites. Given an input matrix $\tilde{M} = [\hat{m}_{i,j}]$ the EPM utilizes a fast conditional expectation maximization algorithm to find the optimal values of m_i^0 , r_i , and s_j to minimize the error between the observed and predicted methylation values across a set of sites. This is accomplished by first fitting a linear model per site using age as the initial s_j . The s_j of the modeled samples is then updated to minimize the error between the observed and predicted methylation values. This process is performed iteratively until the reduction in error is below a specified threshold or the maximum number of iterations is reached. Under the EPM, the epigenetic state has a linear relationship with the modeled methylation data, but not necessarily with chronological age. This allows for nonlinear relationships between time and methylation to be modeled without prior knowledge of the underlying form. Every modeled methylation site has a characteristic m_i^0 and r_i that describes the site in relation to other modeled sites and the output epigenetic states. In the current work, we ask whether the EPM formalism can be utilized for the identification of moderators that impact the association between age and epigenetic state (i.e factors that accelerate or decelerate the changes in epigenetic states with time). To this end we extend the EPM model to simulate methylation matrices associated with age and age accelerating phenotypes. We then evaluate the ability of regularized regression and EPM models to detect age acceleration traits that have linear and nonlinear associations with age. Utilizing a large aggregate dataset we validate the simulation results and in one illustrative example further assess the ability of both approaches to detect age related methylation changes associated with PBB exposure.

81 **2 Results**

82 **2.1 Simulation of Trait Associated Methylation Matrices**

83 Under the EPM the epigenetic state for individual j , S_j , can be interpreted
 84 as a form of biological age that represents a weighted sum of aging
 85 associated phenotypes $S_j = \sum_{k=1}^n \alpha_k p_{k,j} + \dots + \alpha_k p_{k,j}$. Under this model α_k
 86 is the weight for phenotype k and $p_{k,j}$ is the value of phenotype k . Phenotypes
 87 may contribute to increased or decreased aging respectively and when considered as a whole contribute to the overall aging rate observed
 88 for an individual.

89 As shown in our previous work[12], the relationship between $p_{k,j}$ and
 90 time is not necessarily linear. When simulating age associated phenotypes,
 91 each phenotype can be represented as $p_{k,j} = Age_j^{\gamma_k} q_{k,j}$, where γ_k is
 92 a phenotype specific parameter shared among all individuals and $q_{j,k}$ repre-
 93 sents the magnitude of exposure for a simulated trait and is personal
 94 to an individual. The observed phenotype is modeled as an interaction
 95 between age and an exposure of varying magnitude among individuals.
 96 This formulation is flexible as non-age dependent traits can be easily sim-
 97 ulated by setting $\gamma_k = 0$, $p_{k,j} = q_{k,j} = Age_j^0 q_{k,j}$. Individual sites can be
 98 described as a linear model where $\hat{m}_{i,j} = m_i^0 + r_i P_{i,j} + \epsilon_{i,j}$. $P_{i,j}$ is a weighted
 99 sum of phenotypes influencing the methylation status of an individual site,
 100 $P_{i,j} = \sum_{k=1}^n v_k p_{k,j} + \dots + v_k p_{k,j}$.

101 To assess the sensitivity of the EPM and penalized regression approaches
 102 at detecting moderator of epigenetic aging we simulated a methylation ma-
 103 trix containing linear and nonlinear age associated traits of form $p_{k,j} =$
 104 $Age_j^{\mathcal{N}(0.5,0.01)} q_{k,j}$ and $p_{k,j} = Age_j^{\mathcal{N}(1,0.01)} q_{k,j}$. The trait γ parameter was gener-
 105 ated by sampling from a normal distribution $\mathcal{N}(0.5,0.01)$ to generate traits
 106 with varying relationships with time (Figure 1). Samples were simulated
 107 by assigning an age from a uniform distribution, $\mathcal{U}(0,100)$ and setting sam-
 108 ple health by sampling from a normal distribution. Sample health is a
 109 sample specific metric that influences the magnitude and direction of the
 110 simulated age accelerating trait. Simulated traits included a binary pheno-
 111 type ($P = 0.5$), continuous phenotypes influenced by only age, or by age and
 112 sample health (Table 1). The effect, q , of a binary trait was varied from
 113 0.995 to 1.0 over 5 equally spaced intervals. Given a binary trait form of
 114 $p_{k,j} = Age_j^{0.5} q_{k,j}$ a 0.001 decrease in q corresponds to a 1 percent decrease
 115 in epigenetic state by age 100 relative to samples not assigned the binary
 116 trait. Within each interval the standard deviation of the sample health
 117 sampling distribution was varied from 0.0 to 0.01 over 5 equally spaced
 118 intervals. The simulation was repeated 50 times for each binary, continu-
 119 ous trait combination with 500 simulated samples within each simulation.
 120 Additionally, at a binary q of 0.995 the range of continuous traits was ex-
 121 panded over a broader range to assess the model sensitivity for detecting
 122 the continuous trait. Five methylation sites for all continuous traits were
 123 then simulated and 50 methylation sites for the binary trait. An additional
 124 50 sites were simulated that were equally influenced by a mixture of four
 125 continuous traits and the simulated binary trait. The resulting simulation
 126 matrix contains 450 methylation sites.

127 Given a simulation dataset, the samples were split randomly in half for

model training and testing. EPM and penalized regression models were fit for each simulation training set and epigenetic state and age predictions were made for the testing set. e then fit a regression model where the epigenetic age or state is dependent on the age, square-root of the age, the health status, and binary trait status of the sample ($S_j = Age + \sqrt{Age} + health_j + binary_j$). The square-root of the age is included in the regression model to account for the nonlinear relationship between the simulated age and methylation data.

As the exposure size of the binary trait is decreased from 1.00 to 0.995 the ability to detect the influence of the trait on the epigenetic state and age is improved (Figure 2A and B). At an effect size of 0.995 the estimated effect of the binary trait on the epigenetic state is significant ($\mu = 0.035, \sigma = 0.089$) while the effect on the epigenetic age it is not ($\mu = 0.269, \sigma = 0.282$). At an exposure size of 1.0, where the simulated binary trait has no effect, the distribution of p values forEPM and linear models is randomly distributed. The ability to observe the health effect of the simulated continuous traits improves in both the linear and EPM models as the standard deviation of the sample health sampling distribution is increased (Figure 2 C and D). At an exposure size of 0.002 and 0.0025 the average EPM model is significant ($\mu = 0.0194, \sigma = 0.0436$) while the average linear model is not ($\mu = 0.0607, \sigma = 0.128$). At a continuous trait standard deviation above 0.005 both models produce significant results.

2.2 Universal Blood EPM and Penalized Regression Models

We validated the simulation results using a large aggregate dataset composed of Illumina 450k array data[15–27] deposited in the Gene Expression Omnibus[28] (GEO). All methylation array datasets were processed using a unified pipeline from raw array intensity data (IDAT) files using minfi (Aryee et al., 2014). Sex and blood cell type abundance predictions were made for each processed as previously described[29, 30]. The aggregate dataset contains 6,251 whole blood tissue samples representing 16 GEO series.

We trained EPM and penalized regression models using data assembled from four GEO series[31–34] ($n = 1605$) with samples spanning a wide age range (0.01 - 94.0 years). The training set was split by predicted sex, within each sex we used stratified sampling by age to select 95% of the samples for model traning. The selected samples from each sex were combined ($n = 1524$) and the remaining samples ($n = 81$) left out for model evaluation. Methylation values for all samples were quantile normalized by probe type[2] using the median site methylation values across all training samples for each methylation site. Principal component analysis (PCA) was performed on the cell type abundance estimates using the training data. The trained PCA model was used to predict the cell type PCs for the testing and validation datasets.

We fit a penalized regression model to the training matrix as follows. The normalized training methylation matrix was first filtered to remove sites with a variance below 0.001, resulting in a training matrix with

176 183,114 sites. A cross validated ($cv = 5$) elastic net model was trained
177 against training sample ages using the filtered methylation matrix. The
178 trained model performed well on the training ($R^2 = 0.981$) and testing
179 ($R^2 = 0.940$) datasets (S.Figure 2).

180 In contrast to penalized regression based approaches, site selection for
181 the EPM model is performed outside of model fitting. Methylation sites
182 were selected for model training if the absolute Pearson correlation coeffi-
183 cient between methylation values and age was greater than 0.4 ($n = 16,880$).
184 A per site regression model was fit using the observed methylation value
185 as the dependent variable and age as the explanatory variable. Sites with
186 a mean absolute error (MAE) less than 0.025 between the predicted and
187 observed methylation values were retained for further analysis ($n = 7,013$).
188 An EPM model was fit using these sites (Figure 3A). We then sought to
189 identify subsets of sites that had functionally similar forms between age and
190 methylation. This was done to filter sites that were associated with age by
191 chance and to select clusters of sites with low prediction error. Subsets of
192 sites with similar functional form were identified by clustering sites using
193 affinity propagation [35]) by the euclidean distance between the single site
194 regression model residuals. Cross validated EPM and penalized regression
195 models were trained for all clusters with greater than ten sites ($n = 55$).
196 The cluster EPM models show varying associations between the epigenetic
197 state and age relative to the EPM model fit with all sites initially selected
198 by absolute PCC(Figure 3B). Clusters with an observed EPM and penali-
199 zed regression MAE less than 6 ($n = 5$) were combined to fit final EPM
200 and penalized regression models. This resembles the simulated methy-
201 lation matrices where sites with differing functional forms are modeled
202 collectively. The combined cluster EPM and combined cluster regression
203 model performed well on the training and testing datasets (S.Figure 1).

204 We evaluated the combined cluster EPM, combined cluster penalized
205 regression, and the full penalized regression models against a validation
206 data set consisting of 14 GEO series experiments representing 4,600 whole
207 blood tissue samples. Each model accurately predicted the epigenetic state
208 or epigenetic age of the validation samples (Figure 4). We then fit an
209 ordinary least squares regression model for every validation experiment
210 individually to predict the observed epigenetic age or state using the sample
211 age, the square root of age, cell type PCs, and predicted sex ($S_j = Age +$
212 $\sqrt{Age} + PC1 + PC2 + PC3 + Sex + Intercept$). If the proportion of female samples
213 to the total number of samples was greater than 0.7 the sex term was
214 dropped from the regression model. Significant cell type PC2 coefficients
215 were observed for all EPM models and the majority of the cluster and
216 full penalized regression models (Figure 5A). Significant cell type PC1 and
217 PC3 coefficients were observed for the majority of the EPM models but not
218 for the cluster or full penalized regression models. Significant sex effects
219 ($p < 0.0038$) were observed for 9, 4, 0 out of 15 models for the EPM, cluster
220 penalized regression, and full penalized regression respectively (Figure 5B).

221 2.3 Polybrominated Biphenyls Exposure

222 Polybrominated biphenyls (PBB) were widely used throughout the United
223 States in the 1960's and 1970's for a variety of industrial applications.

224 Widespread PBB exposure occurred in the state of Michigan from the
225 summer of 1973 to later spring of 1974 when an industrial PBB mixture
226 was incorrectly substituted for a nutritional supplement used in livestock
227 feed[36]. PBB is biologically stable and has a slow biological half life; indi-
228 viduals exposed during the initial 1973 - 1974 incident still have detectable
229 PBB in their blood[37]. PBB is an endocrine-disrupting compound and ex-
230 posure has been linked to numerous adverse health outcomes in Michigan
231 residents such as thyroid dysfunction[38, 39] and various cancers[40, 41].
232 A study by Curtis et al. showed total PBB exposure is associated with
233 altered DNA methylation at CpG sites enriched for an association with
234 endocrine-related autoimmune disease [42]. Utilizing the publicly available
235 Illumina Methylation EPIC array [43] profiles ($n = 679$), that covered a wide
236 age range (23 - 88 years), we sought to compare the ability of penalized
237 regression and the EPM to detect epigenetic age acceleration associated
238 with PBB exposure.

239 Briefly, 50% of samples ($n = 339$) were selected for model training using
240 stratified cross validation by age. A cross validated elastic net model was
241 trained using all methylation sites with a site variance above 0.001, ($n =$
242 529,703). The trained model performed well on the training and testing
243 datasets ($R^2 = 1.00, R^2 = 0.740, S.Figure2A - B$). EPM sites were selected
244 and models fit as described with the universal blood EPM. Four EPM
245 clusters ($MAE < 6$) were merged for a combined EPM model built using
246 413 CpG sites. The combined EPM model performed well in training and
247 testing datasets ($R^2 = 0.794, R^2 = 0.812, S.Figure2C - D$). Epigenetic age and
248 epigenetic state predictions were then made for the testing samples using
249 the penalized regression and EPM models.

250 We then fit an OLS regression model to predict the epigenetic age or
251 state dependent on PBB-153 exposure, h age, the square root of age, cell
252 type PCs, and predicted sex ($S_j = Age + \sqrt{Age} + PC1 + PC2 + PC3 + Sex +$
253 $PBB - 153 + Intercept$). PBB-153 exposure was highly significant in the
254 EPM regression model ($p = 5.9e - 10$) but not the penalized regression
255 model ($p = 0.141$).

256 3 Discussion

257 A long standing question in the field of epigenetics was whether biomarkers
258 could be trained to predict various traits using methylation measurements.
259 The most successful biomarkers to date have been epigenetic clocks that
260 can accurately predict the age of an individual based on their methylation
261 pattern. These have been shown to be useful for human studies of aging,
262 as well as for animal studies, including mice[44] and dogs[45]. DNA methy-
263 lation biomarkers are typically constructed using penalized regression ap-
264 proaches. Given a large enough matrix, penalized regression will select
265 sites that minimize the prediction error given a modeled trait. Epigenetic
266 clocks are examples of such models. Beyond predicting actual ages, these
267 models have also been used to measure the influence of external factors on
268 the rates of aging, and multiple studies have shown that the resulting age
269 accelerations (i.e the differences between actual and predicted ages) are sig-
270 nificantly associated with multiple factors such as cardiovascular disease[7]

271 and mortality risk[4, 5].

272 While epigenetic clocks have proven to be useful they have significant
273 limitations. Because they are based on linear models, it may be difficult
274 to model aging when the underlying methylation changes are non-linear in
275 time. Moreover, epigenetic clocks are prone to overfitting, and while cross
276 validation schemes are often used to construct robust clocks, they often do
277 not yield accurate estimates for other datasets. Finally, epigenetic clocks
278 are not very interpretable, and highly degenerate, so that it is difficult to
279 extract biological insights from the weights of the models.

280 To overcome some of these limitations, we have previously developed
281 the epigenetic pacemaker formalism. In this approach, rather than building
282 a model for the age, we construct a model for the observed methylation data
283 that depends on age. The advantage of this approach is that this formalism
284 allows us to identify non-linear associations between methylation and age
285 across a lifespan. Moreover, these models tend to be robust to training as
286 they are fit to large methylation matrices rather than age vectors. Finally,
287 the model describes the change in methylation at each site with respect
288 to a time dependent epigenetic state, and therefore all parameters of the
289 model are directly interpretable as either initial values of methylation or
290 rates of change of methylation.

291 Depending on the context, epigenetic clocks are both more and less
292 effective than the EPM. The penalized regression models provide more
293 accurate age predictions ($R^2 = 0.875, 0.911$) than the EPM model ($R^2 =$
294 0.821), and the model output can be directly compared to the age of a
295 sample. However, because these models are optimized to reduce the error
296 between actual and predicted age, they tend to minimize the effect of
297 extraneous factors on the predicted age. As such, epigenetic clocks are
298 not optimal for identifying external factors that moderate the relations
299 between actual and predicted age. By contrast, the EPM models are not
300 optimized to minimize the difference between predicted and actual age,
301 but rather try to minimize the difference between observed and modeled
302 methylation values. As such, they retain many of the effects that other
303 factors may have on the association between methylation and epigenetic
304 states.

305 In this study we find that while the penalized regression models were
306 more accurate for predicting age, the epigenetic state generated by the
307 EPM is significantly impacted by cell type and sex effects in both simula-
308 tions and real data. We also found that The EPM model generated for
309 individuals exposed to PBB was sensitive to e PBB exposure, which has
310 been linked to negative health outcomes, while the penalized regression
311 epigenetic aging model was not. Additionally, the sensitivity of the EPM
312 to moderators of epigenetic aging has been supported by two two recent
313 studies investigating epigenetic aging in marmots[46] and zebras[47]. In the
314 first of these studies, EPM models showed an association between hiberna-
315 tion and slowed epigenetic aging in marmots and in the second an increased
316 epigenetic age associated with zebra inbreeding; no such associations were
317 observed with penalized regression epigenetic age models.

318 Most studies of human epigenetic aging are not motivated by the de-
319 velopment of accurate age predictors, since ages are nearly always known
320 in studies, but rather by the discovery of biological aging moderators. The

321 EPM is a more sensitive approach than epigenetic clocks for the detection
322 of factors other than age that influence the epigenome and therefore
323 potentially more useful for discovering moderators of biological aging.

324 4 Methods

325 4.1 Simulation

326 We implemented the simulation framework as a python package with numpy($\geq v1.16.3$)[48]
327 and scikit-learn(v0.24)[49] as dependencies. A simulation run generates a
328 trait-associated methylation matrix and samples are tied to the simulated
329 traits. The simulation procedure is implemented as follows:

- 330 • Traits are initialized that contain the information about the trait re-
331 lationship with age and a simulated sample phenotype. Given the
332 structure $p_{k,j} = Age_j^{\gamma_k} q_{k,j}$, and k samples and j traits γ is character-
333 teristic of the trait. When a sample is passed to a trait, a value of q
334 is generated for the sample by sampling from a normal distribution
335 with a variance characteristic of the simulation trait. Additionally,
336 each trait can be optionally influenced by a characteristic measure
337 of sample health, h_j . Given, a normally distributed trait $\mathcal{N}(\mu, \sigma^2)$
338 and a health effect h_j , the sampled distribution for individual j is
339 $\mathcal{N}(\mu + h_j, \sigma^2)$. Continuous and binary traits can be simulated. If a
340 binary trait is simulated, a q other than 1 is assigned at a specified
341 probability.
- 342 • Samples are simulated by setting the age by sampling from a uniform
343 distribution over a specified range and by setting a sample health
344 metric h by sampling from a normal distribution centered on zero
345 with a specified variance. Traits passed to a sample simulation object
346 are then set according to the age and health of the sample. Simulated
347 samples retain all the set phenotype information for downstream ref-
348 erence.
- 349 • Methylation sites are simulated by randomly setting the initial methy-
350 lation value, maximum observable methylation value, the rate of
351 change at the site, and the error observed at each site. Sites are then
352 assigned traits that influence the methylation values at each site.
- 353 • Methylation values are simulated for each site for every individual
354 given the simulated phenotypes with a specified amount of random
355 noise.

356 4.2 Simulation EPM and Penalized Regression Models

357 Simulation data was randomly split in half into training and testing sets.
358 EPM models were fit using the simulated methylation matrix against age.
359 Penalized regression models were fit using scikit-learn(v0.24)[49] Elastic-
360 Net ($\text{alpha}=1$, $\text{l1_ratio}=0.75$, and $\text{selection}=\text{random}$). All other parame-
361 ters were set to their default values. Ordinary least squares regression as
362 implemented in statsmodels (0.11.1)[50] was utilized to describe the epi-
363 genetic state or age with the following form ($S_j = Age + \sqrt{Age} + health_j +$

364 $binary_j$). Full analysis is found in the EPMSimulation.ipynb supplementary
365 file.

366 4.3 Methylation Array Processing

367 Metadata for Illumina methylation 450K Beadchip methylation array ex-
368 periments deposited in the Gene Expression Omnibus (GEO) [28] with
369 more than 50 samples were parsed using a custom python toolset. Experi-
370 ments that were missing methylation beadchip array intensity data (IDAT)
371 files, made repeated measurements of the same samples, utilized cultured
372 cells, or assayed cancerous tissues were excluded from further processing.
373 IDAT files were processed using minfi[30] (v1.34.0). Sample IDAT files
374 were processed in batches according to GEO series and Beadchip identi-
375 fication. Methylation values within each batch were normal-exponential
376 normalized using out-of-band probes[51]. Blood cell types counts were
377 estimated using a regression calibration approach[29] and sex predictions
378 were made using the median intensity measurements of the X and Y chro-
379 mosomes as implemented in minfi[30]. Whole blood array samples were
380 used for downstream analysis if the sample median methylation probe in-
381 tensity was greater than 10.5 and the difference between the observed and
382 expected median unmethylation probe intensity is less than 0.4, where the
383 expected median unmethylated signal is described by ($y = 0.66x + 3.718$).

384 4.4 Blood EPM and Penalized Regression Models

385 Methylation sites with an absolute Pearson correlation coefficient between
386 methylation values and age greater than 0.40 and 0.45 for the unified whole
387 blood and PBB datasets respectively were initially selected for EPM model
388 training. A linear model was generated using numpy polyfit [48] with age
389 and the independent variable and methylation values as the dependent
390 variable. Mean absolute error (MAE) was calculated as the mean absolute
391 difference between the observed and predicted meth values according to
392 the site linear models. A vector of residuals generated using these models
393 were utilized for clustering by affinity propagation[35]) as implemented in
394 scikit-learn (v0.24)[49] with a random state of 1 and a cluster preference
395 of -2.5. Cross-validated EPM, and penalized regression models for the
396 universal blood analysis, were trained for all clusters containing greater
397 than ten sites. Clusters with an observed EPM and penalized regression
398 MAE less than 6.0 were combined to fit final EPM and regression models.

399 Penalized regression models were fit using scikit-learn(v0.24)[49] Elas-
400 ticNetCV (cv=5 alpha=1, l1_ratio=0.75, and selection=random). All other
401 parameters were set to their default values. Principal Component Anal-
402 ysis as implemented in scikit-learn was utilized with default parameters
403 to perform PCA on training sample cell type adundances. The trained
404 PCA was utilized to calculate cell type PCs for the testing and validation
405 samples. Ordinary least squares regression as implemented in statsmodels
406 (0.11.1)[50] was utilized describe the epigenetic state or age with the fol-
407 lowing form ($S_j = Age + \sqrt{Age} + CellTypePC1 + CellTypePC2 + CellTypePC3 +$
408 $Sex + Intercept$). Full analysis is found in the EPMUniversalClock.ipynb
409 supplementary file.

410 **4.5 Analysis Environment**

411 412 413 Analysis was carried out in a Jupyter[52] analysis environment. Joblib[53],
SciPy[54], Matplotlib[55], Seaborn[56], Pandas[57] and TQDM[58] packages were utilized during analysis.

414 **4.6 Supplementary Information**

415 416 Analysis code and notebooks can be found at <https://github.com/NuttyLogic/EPM-ModeratorsOfAging>.

References

1. Horvath, S. & Raj, K. DNA methylation-based biomarkers and the epigenetic clock theory of ageing. en. *Nat. Rev. Genet.* **19**, 371–384 (June 2018).
2. Horvath, S. *DNA methylation age of human tissues and cell types* 2013.
3. Hannum, G. *et al.* Genome-wide methylation profiles reveal quantitative views of human aging rates. en. *Mol. Cell* **49**, 359–367 (Jan. 2013).
4. Marioni, R. E. *et al.* DNA methylation age of blood predicts all-cause mortality in later life. en. *Genome Biol.* **16**, 25 (Jan. 2015).
5. Perna, L. *et al.* Epigenetic age acceleration predicts cancer, cardiovascular, and all-cause mortality in a German case cohort 2016.
6. Dugué, P.-A. *et al.* DNA methylation-based biological aging and cancer risk and survival: Pooled analysis of seven prospective studies. en. *Int. J. Cancer* **142**, 1611–1619 (Apr. 2018).
7. Huang, R.-C. *et al.* Epigenetic Age Acceleration in Adolescence Associates With BMI, Inflammation, and Risk Score for Middle Age Cardiovascular Disease. en. *J. Clin. Endocrinol. Metab.* **104**, 3012–3024 (July 2019).
8. Armstrong, N. J. *et al.* Aging, exceptional longevity and comparisons of the Hannum and Horvath epigenetic clocks. en. *Epigenomics* **9**, 689–700 (May 2017).
9. Horvath, S. *et al.* Decreased epigenetic age of PBMCs from Italian semi-supercentenarians and their offspring. en. *Aging* **7**, 1159–1170 (Dec. 2015).
10. Horvath, S. *et al.* Obesity accelerates epigenetic aging of human liver. en. *Proc. Natl. Acad. Sci. U. S. A.* **111**, 15538–15543 (Oct. 2014).
11. Zhang, Q. *et al.* Improved precision of epigenetic clock estimates across tissues and its implication for biological ageing 2019.
12. Snir, S., Farrell, C. & Pellegrini, M. Human epigenetic ageing is logarithmic with time across the entire lifespan. en. *Epigenetics* **14**, 912–926 (Sept. 2019).

13. Snir, S., vonHoldt, B. M. & Pellegrini, M. A Statistical Framework to Identify Deviation from Time Linearity in Epigenetic Aging. en. *PLoS Comput. Biol.* **12**, e1005183 (Nov. 2016).
14. Farrell, C., Snir, S. & Pellegrini, M. The Epigenetic Pacemaker: modeling epigenetic states under an evolutionary framework. en. *Bioinformatics* **36**, 4662–4663 (Nov. 2020).
15. Marabita, F. *et al.* Author Correction: Smoking induces DNA methylation changes in Multiple Sclerosis patients with exposure-response relationship. en. *Sci. Rep.* **8**, 4340 (Mar. 2018).
16. Ventham, N. T. *et al.* Integrative epigenome-wide analysis demonstrates that DNA methylation may mediate genetic risk in inflammatory bowel disease. en. *Nat. Commun.* **7**, 13507 (Nov. 2016).
17. Tan, Q. *et al.* Epigenetic signature of birth weight discordance in adult twins. en. *BMC Genomics* **15**, 1062 (Dec. 2014).
18. Johnson, R. K. *et al.* Longitudinal DNA methylation differences precede type 1 diabetes. en. *Sci. Rep.* **10**, 3721 (Feb. 2020).
19. Voisin, S. *et al.* Many obesity-associated SNPs strongly associate with DNA methylation changes at proximal promoters and enhancers 2015.
20. Soriano-Tárraga, C. *et al.* Epigenome-wide association study identifies TXNIP gene associated with type 2 diabetes mellitus and sustained hyperglycemia. en. *Hum. Mol. Genet.* **25**, 609–619 (Feb. 2016).
21. Dabin, L. *et al.* Altered DNA methylation profiles in blood from patients with sporadic Creutzfeldt-Jakob disease
22. Horvath, S. & Ritz, B. R. Increased epigenetic age and granulocyte counts in the blood of Parkinson’s disease patients. en. *Aging* **7**, 1130–1142 (Dec. 2015).
23. Kurushima, Y. *et al.* Epigenetic findings in periodontitis in UK twins: a cross-sectional study 2019.
24. Zannas, A. S. *et al.* Epigenetic upregulation of FKBP5 by aging and stress contributes to NF- κ B–driven inflammation and cardiovascular risk. en. *Proc. Natl. Acad. Sci. U. S. A.* **116**, 11370–11379 (June 2019).
25. Braun, P. R. *et al.* Genome-wide DNA methylation comparison between live human brain and peripheral tissues within individuals. en. *Transl. Psychiatry* **9**, 47 (Jan. 2019).
26. Demetriou, C. A. *et al.* Methylome analysis and epigenetic changes associated with menarcheal age. en. *PLoS One* **8**, e79391 (Nov. 2013).
27. Tserel, L. *et al.* Age-related profiling of DNA methylation in CD8+ T cells reveals changes in immune response and transcriptional regulator genes. en. *Sci. Rep.* **5**, 13107 (Aug. 2015).
28. Barrett, T. *et al.* NCBI GEO: archive for functional genomics data sets—update. en. *Nucleic Acids Res.* **41**, D991–D995 (Nov. 2012).

29. Houseman, E. A. *et al.* DNA methylation arrays as surrogate measures of cell mixture distribution. en. *BMC Bioinformatics* **13**, 86 (May 2012).
30. Aryee, M. J. *et al.* Minfi: a flexible and comprehensive Bioconductor package for the analysis of Infinium DNA methylation microarrays. en. *Bioinformatics* **30**, 1363–1369 (May 2014).
31. Johansson, A., Enroth, S. & Gyllensten, U. Continuous Aging of the Human DNA Methylome Throughout the Human Lifespan. en. *PLoS One* **8**, e67378 (June 2013).
32. Liu, Y. *et al.* Epigenome-wide association data implicate DNA methylation as an intermediary of genetic risk in rheumatoid arthritis 2013.
33. Butcher, D. T. *et al.* CHARGE and Kabuki Syndromes: Gene-Specific DNA Methylation Signatures Identify Epigenetic Mechanisms Linking These Clinically Overlapping Conditions. en. *Am. J. Hum. Genet.* **100**, 773–788 (May 2017).
34. Dámaso, E. *et al.* Comprehensive Constitutional Genetic and Epigenetic Characterization of Lynch-Like Individuals. en. *Cancers* **12** (July 2020).
35. Frey, B. J. & Dueck, D. Clustering by passing messages between data points. en. *Science* **315**, 972–976 (Feb. 2007).
36. Fries, G. F. The PBB episode in Michigan: an overall appraisal. en. *Crit. Rev. Toxicol.* **16**, 105–156 (1985).
37. Safe, S. Polychlorinated biphenyls (PCBs) and polybrominated biphenyls (PBBs): biochemistry, toxicology, and mechanism of action. en. *Crit. Rev. Toxicol.* **13**, 319–395 (1984).
38. Jacobson, M. H. *et al.* Serum Polybrominated Biphenyls (PBBs) and Polychlorinated Biphenyls (PCBs) and Thyroid Function among Michigan Adults Several Decades after the 1973–1974 PBB Contamination of Livestock Feed 2017.
39. Curtis, S. W. *et al.* Thyroid hormone levels associate with exposure to polychlorinated biphenyls and polybrominated biphenyls in adults exposed as children. en. *Environ. Health* **18**, 75 (Aug. 2019).
40. Terrell, M. L., Rosenblatt, K. A., Wirth, J., Cameron, L. L. & Marcus, M. Breast cancer among women in Michigan following exposure to brominated flame retardants. en. *Occup. Environ. Med.* **73**, 564–567 (Aug. 2016).
41. Hoque, A. *et al.* Cancer among a Michigan cohort exposed to polybrominated biphenyls in 1973. en. *Epidemiology* **9**, 373–378 (July 1998).
42. Curtis, S. W. *et al.* Exposure to polybrominated biphenyl (PBB) associates with genome-wide DNA methylation differences in peripheral blood. en. *Epigenetics* **14**, 52–66 (Jan. 2019).
43. Pidsley, R. *et al.* Critical evaluation of the Illumina MethylationEPIC BeadChip microarray for whole-genome DNA methylation profiling. en. *Genome Biol.* **17**, 208 (Oct. 2016).

44. Thompson, M. J. *et al.* A multi-tissue full lifespan epigenetic clock for mice. en. *Aging* **10**, 2832–2854 (Oct. 2018).
45. Thompson, M. J., vonHoldt, B., Horvath, S. & Pellegrini, M. An epigenetic aging clock for dogs and wolves. en. *Aging* **9**, 1055–1068 (Mar. 2017).
46. Pinho, G. M. *et al.* *Hibernation slows epigenetic aging in yellow-bellied marmots* en. Mar. 2021.
47. Larison, B. *et al.* *Epigenetic models predict age and aging in plains zebras and other equids* en. Mar. 2021.
48. Harris, C. R. *et al.* Array programming with NumPy. en. *Nature* **585**, 357–362 (Sept. 2020).
49. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.* **12**, 2825–2830 (2011).
50. Seabold, S. & Perktold, J. *Statsmodels: Econometric and statistical modeling with python* in *Proceedings of the 9th Python in Science Conference* **57** (2010), 61.
51. Triche Jr, T. J., Weisenberger, D. J., Van Den Berg, D., Laird, P. W. & Siegmund, K. D. Low-level processing of Illumina Infinium DNA Methylation BeadArrays. en. *Nucleic Acids Res.* **41**, e90 (Apr. 2013).
52. Basu, A. *Reproducible research with jupyter notebooks*
53. Varoquaux, G. & Grisel, O. Joblib: running python function as pipeline jobs. *packages. python. org/joblib* (2009).
54. Virtanen, P. *et al.* SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat. Methods* (Feb. 2020).
55. Hunter, J. D. *Matplotlib: A 2D Graphics Environment* 2007.
56. Waskom, M. *seaborn: statistical data visualization*. *J. Open Source Softw.* **6**, 3021 (Apr. 2021).
57. McKinney, W. *Python for Data Analysis: Data Wrangling with Pandas, NumPy, and IPython* en (“O’Reilly Media, Inc.”, Oct. 2012).
58. Da Costa-Luis, C. O. *tqdm: A Fast, Extensible Progress Meter for Python and CLI*. *JOSS* **4**, 1277 (May 2019).

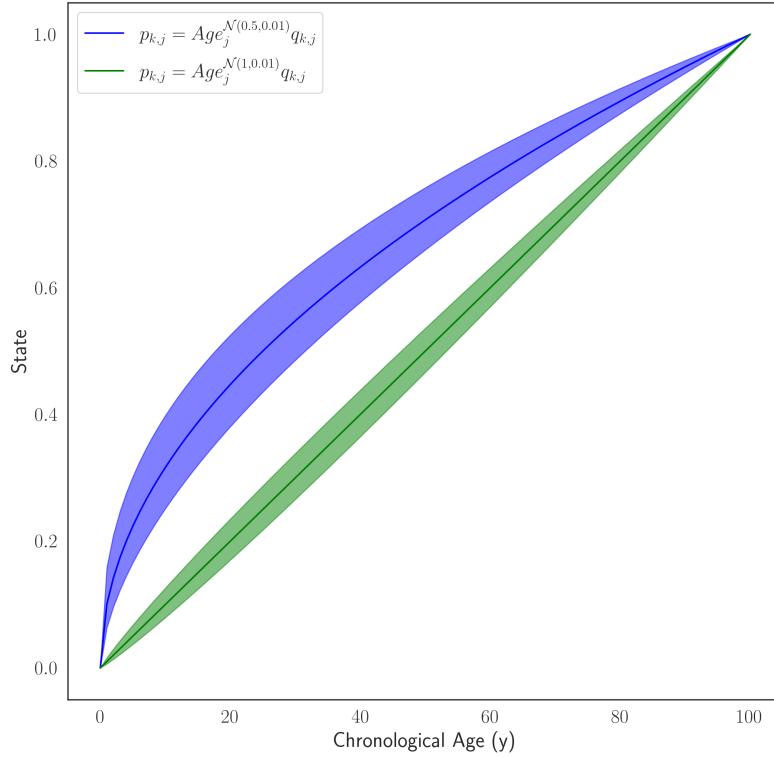


Figure1: Simulated trait forms where the shaded area represent one standard deviation away from the mean γ , given $p_{k,j} = Age_j^{\gamma_k} q_{k,j}$.

Table 1: Simulated Trait Conditions

Trait Form	Beta	Gamma	Gamma Std. Dev.	Sample Effect	Age Only	Generated Phenotypes
Continuous	0.1	$\mathcal{N}(0.5, 0.01)$	0.05	Yes	No	10
Continuous	0.1	$\mathcal{N}(1.0, 0.01)$	0.05	Yes	No	10
Continuous	0.1	$\mathcal{N}(0.5, 0.01)$	0.05	No	Yes	20
Continuous	0.1	$\mathcal{N}(1.0, 0.01)$	0.05	No	Yes	20
Binary ($Pr = 0.5$)	0.1	0.5	0	Yes	No	1

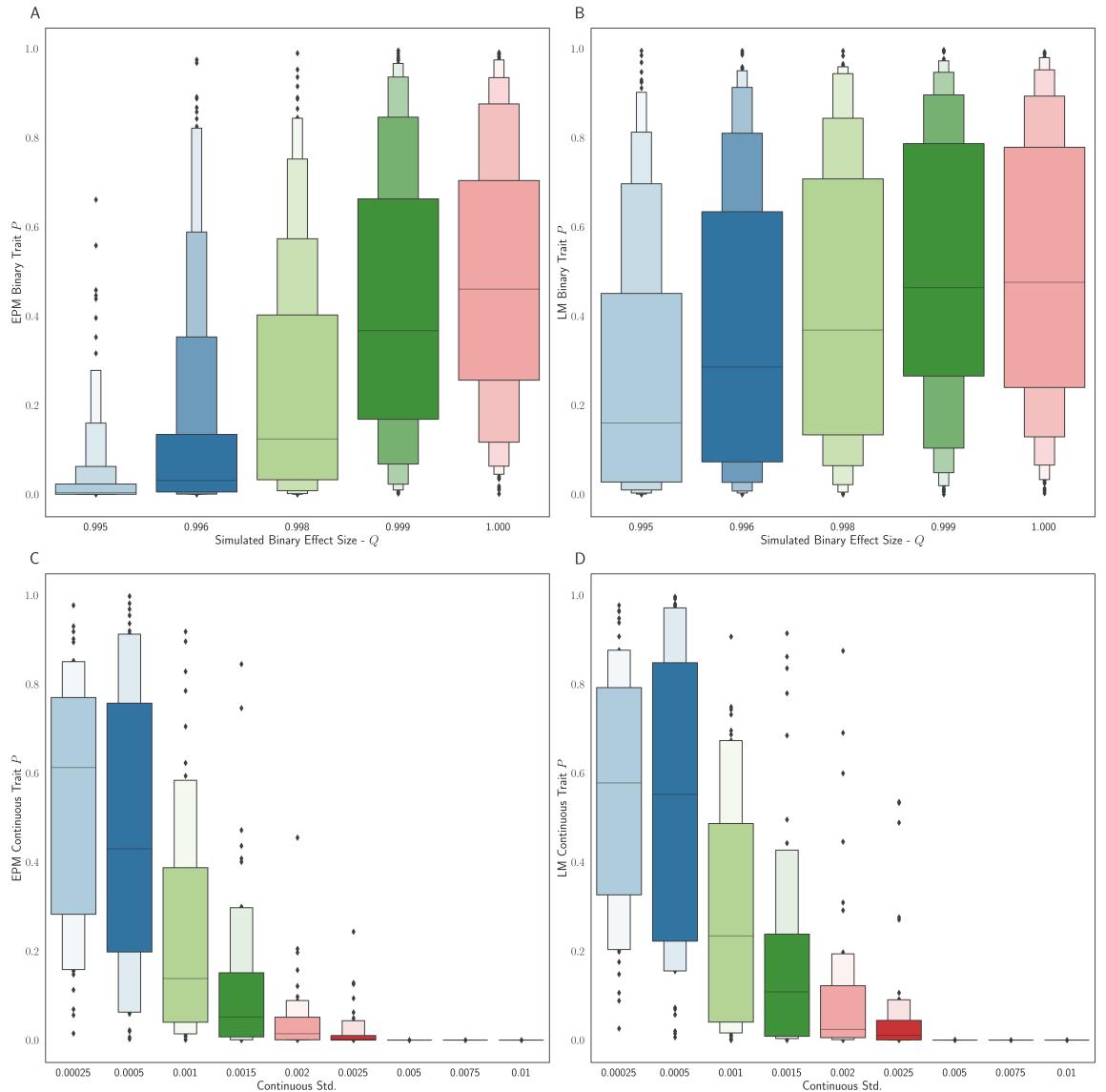


Figure 2: The distribution binary coefficient p-values for **A** EPM and **B** penalized regression models. The distribution of p-values given a simulation health standard deviation for **C** EPM and **D** penalized regression models.

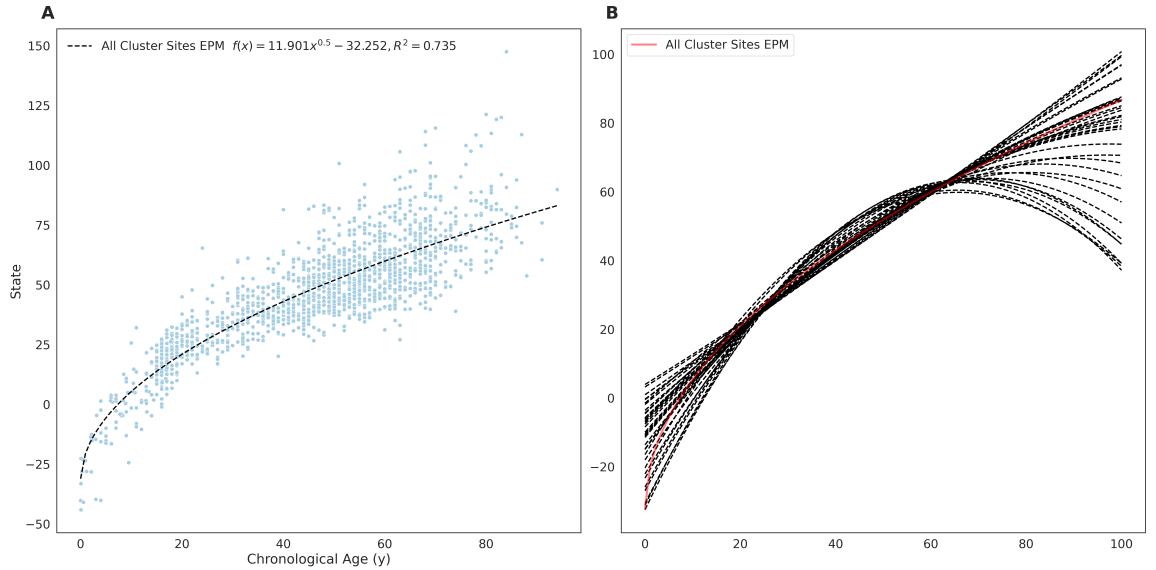


Figure3: **A** EPM model fit with 3832 methylation sites with a MAE below 0.025. **B** The fit trend line for EPM clusters with more than 10 sites and an $R^2 \geq 0.4$.

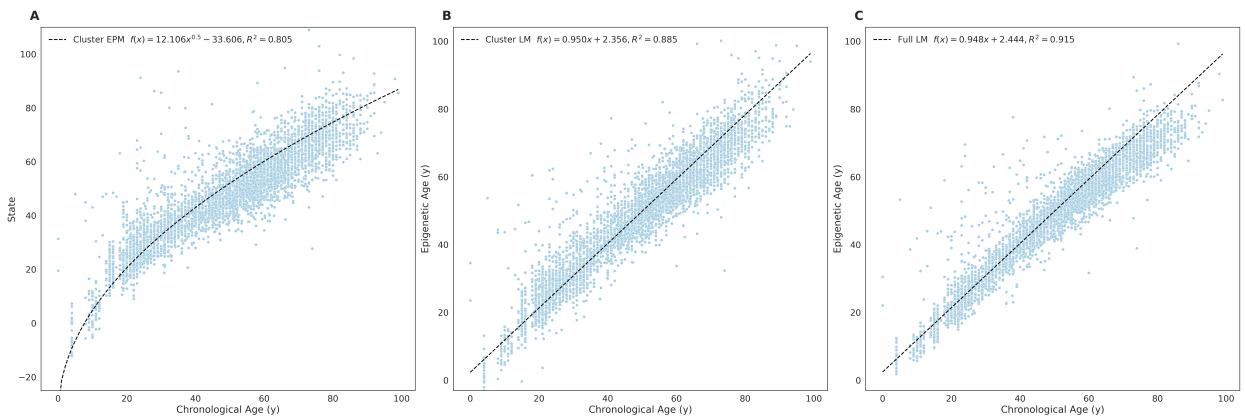


Figure4: Whole blood tissue validation **A** EPM, **B** cluster penalized regression and **C** full penalized regression models.

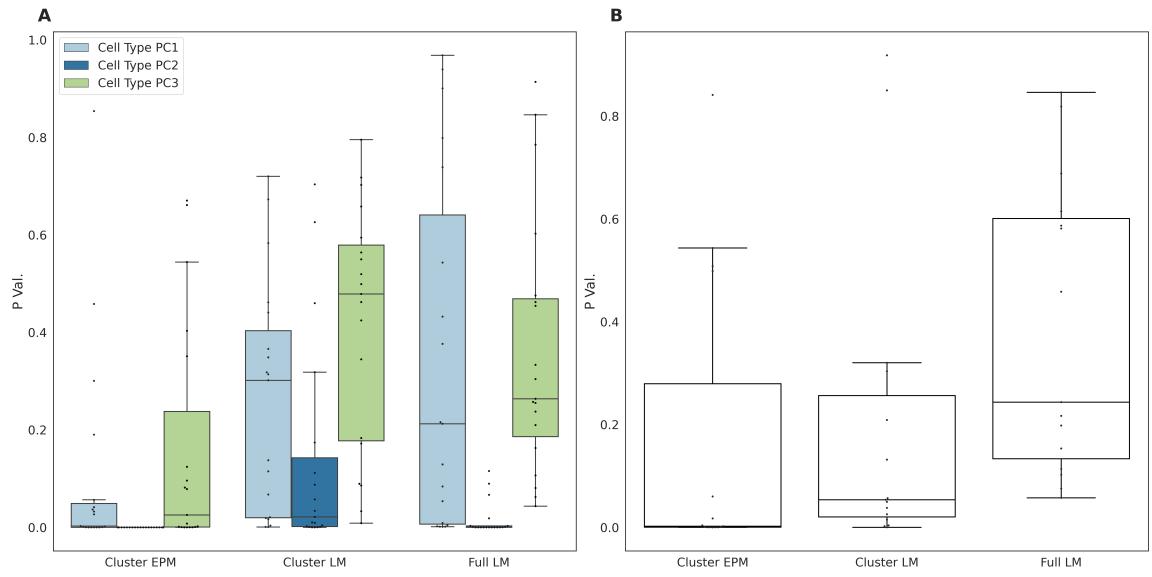
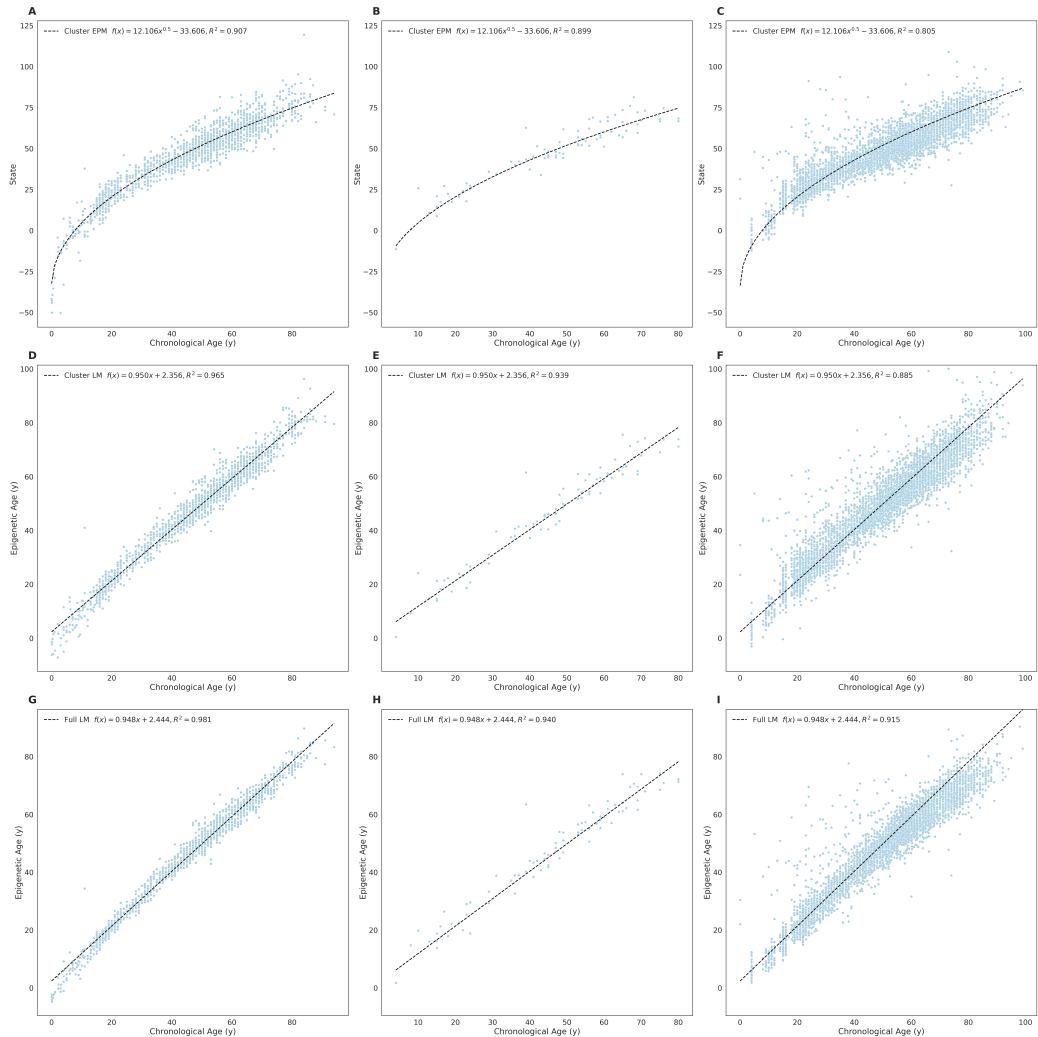
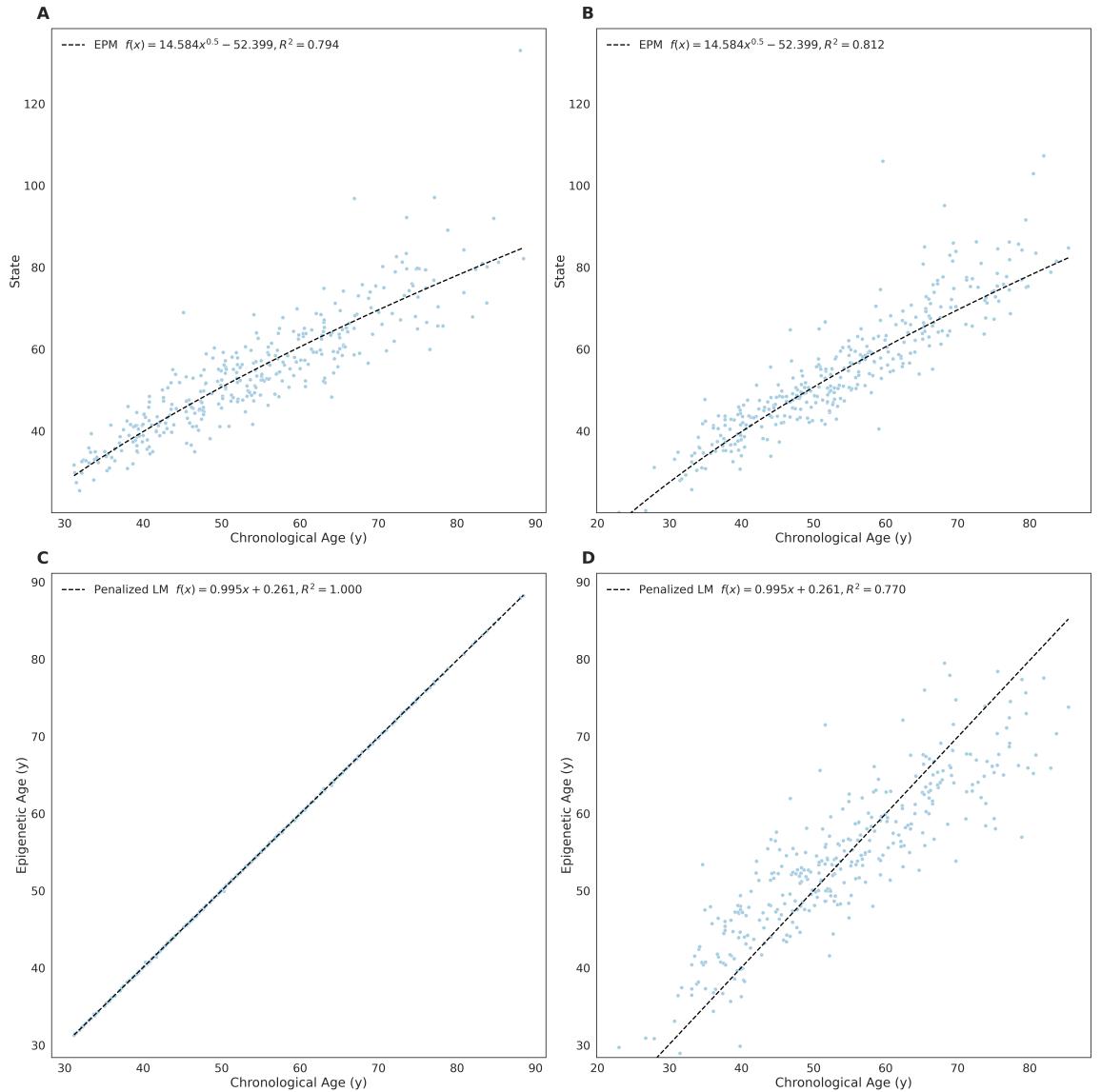


Figure 5: **A** Cell type principal component and **B** predicted sex regression coefficient p-values.



S.Figure1: Universal blood EPM and regression models. **A - C** Train, testing, and validation EPM model. **D-E** Train, testing, and validation cluster penalized regression model. **G-J** Train, testing, and validation full penalized regression model.



S.Figure2: PBB EPM and regression models. **A - B** Train and testing EPM model.
C-D Train and testing penalized regression model.