

Probabilistic modeling methods for cell-free DNA methylation based cancer classification

Viivi Halla-aho¹ and Harri Lähdesmäki¹

¹Department of Computer Science, Aalto University, Espoo, Finland

Abstract

Background: cfMeDIP-seq is a low-cost method for determining the DNA methylation status of cell-free DNA and it has been successfully combined with statistical methods for accurate cancer diagnostics. We investigate the diagnostic classification aspect by applying statistical tests and dimension reduction techniques for feature selection and probabilistic modeling for the cancer type classification, and we also study the effect of sequencing depth.

Methods: We experiment with a variety of statistical methods that use different feature selection and feature extraction methods as well as probabilistic classifiers for diagnostic decision making. We test the (moderated) t-tests and the Fisher’s exact test for feature selection, principal component analysis (PCA) as well as iterative supervised PCA (ISPCA) for feature generation, and GLMnet and logistic regression methods with sparsity promoting priors for classification. Probabilistic programming language Stan is used to implement Bayesian inference for the probabilistic models.

Results and conclusions: We compare overlaps of differentially methylated genomic regions as chosen by different feature selection methods, and evaluate probabilistic classifiers by evaluating the area under the receiver operating characteristic (AUROC) scores on discovery and validation cohorts. While we observe that many methods perform equally well as, and occasionally considerably better than, GLMnet that was originally proposed for cfMeDIP-seq based cancer classification, we also observed that performance of different methods vary across sequencing depths, cancer types and study cohorts. Overall, methods that seem robust and promising include Fisher’s exact test and ISPCA for feature selection as well as a simple logistic regression model with the number of hyper and hypomethylated regions as features.

Background

In recent years the interest in the possibilities of utilizing circulating cell-free DNA (cfDNA) for cancer diagnostics has grown, enhanced by the development of next-generation sequencing (NGS) technologies. Cell-free DNA refers to DNA

fragments that are not associated with cells and is considered to origin from cell apoptosis and necrosis [1, 2]. In the case of a presence of a tumor, part of cfDNA can be of tumor origin, and can be called circulating tumor DNA (ctDNA). Cell-free DNA can be extracted in a minimally non-invasive manner from a bodily fluid sample, such as blood, to identify and detect cancer type specific biomarkers [3].

ctDNA is believed to represent the tumor burden and to carry the same genomic and epigenetic properties as the tumor of origin [3]. Therefore multiple types of cancer biomarkers can be identified and detected from cfDNA, including mutations, epigenetic modifications and copy-number alterations [3]. The DNA fragmentation profiles of cfDNA can also be used to classify cancer types [4]. As quantification of somatic mutations from sequencing data necessarily requires a high sequencing coverage, assays that use methylation biomarkers can provide a significant cost reduction. Consequently, in this work we concentrate on cancer classification which is based on DNA methylation biomarkers. The most common way to measure DNA methylome is bisulfite sequencing (BS-seq), and tools such as CancerLocator [5] utilize BS-seq data to learn machine learning models to classify different cancer types. However, bisulfite conversion step of the BS-seq method leads to high degradation of the input DNA [6], making it unsuitable for cfDNA analysis where the amount of sequencing material is small.

Cell-free methylated DNA immunoprecipitation and high-throughput sequencing (cfMeDIP-seq) is a protocol for measuring the methylation status of cell-free DNA [7]. cfMeDIP-seq is a version of the MeDIP-seq method that takes into account specific needs of cfDNA sequencing. The amount of cfDNA material available for sequencing is often very small, so filler DNA from *Enterobacteria phage λ* is used in cfMeDIP-seq to increase the amount of DNA material [7]. Compared to bisulfite sequencing, cfMeDIP-seq is even more cost-effective, as only the methylated reads are sequenced in the immunoprecipitation-based approach [8]. While bisulfite sequencing provides information of the methylation status in a single-base resolution, cfMeDIP-seq can give information on the methylation status of genomic regions of length around 100bp or more [8].

Along with the cfMeDIP-seq protocol, statistical methods for finding differentially methylated regions (DMRs) and machine learning methods for classification of the cancer types were proposed in [7]. In brief, DMRs are found for each cancer type (i.e., a class) using a moderated t-test that separately compares that one cancer type vs. all other cancer types (including healthy controls), one at a time. Then GLMnet [9] binary classifiers are trained for each of the classes, using the DMRs found in the previous step. The results of these methods presented in [7] show high accuracy both in the discovery and validation data cohorts. The results for renal cell carcinoma (RCC) class were further validated in [10], where the same methods were applied to classify RCC patients from healthy controls. The cfMeDIP-seq assays and analysis steps were performed not only for plasma cfDNA, but also for cfDNA of urine origin. Both resulted in high AUROC scores, although the plasma-based classifier performed slightly better. In [11], the cfMeDIP-seq data set from [7] was extended with

samples from intracranial tumor patients. There are also other cases where usage of cfMeDIP-seq or MeDIP-seq measurements of cfDNA for cancer classification has been reported. Similar to [7], GLMnet models were utilized to classify pancreatic cancer patients and healthy controls in [12], but the model features were based on both cfMeDIP-seq and cell-free 5hmC sequencing data. Peak calling was performed for both cfMeDIP-seq and cell-free 5hmC signals with MACS2 tool [13], and differential peaks between the cancer samples and healthy controls were then determined with t-test. Both 5mC and 5hmC peaks used separately as model features gave high prediction accuracies, but using both peak types together reached even better results. In [14], the performances of detecting metastatic renal cell carcinoma (mRCC) using cfMeDIP-seq based cfDNA methylation analysis and cfDNA variant analysis were compared. The TMM-normalized cfMeDIP-seq count data was used to find DMRs with limma-voom [15] and the DMRs were then utilised as features in a GLMnet model, similar to the approach in [7]. The comparison showed that the classification method based on cfMeDIP-seq data had considerably higher sensitivity than cfDNA variant analysis. MeDIP-seq has also been applied to small number of cfDNA samples to find DMRs between cancer patients and healthy individuals, in particular for lung cancer [16] and pancreatic cancer [17]. In these cases the DMRs were found using MEDIPS tool [18] and then the found DMRs were further used for analysis of methylation data of tissue origin. MEDIPS is a tool for quality control and analysis of immunoprecipitation sequencing data, and it does differential coverage analysis using negative binomial model from edgeR package [19].

The results in [7] showed that methylation-based cfDNA biomarkers have great potential in cancer classification and that cfMeDIP-seq is a sensitive yet low-cost method for measuring the methylome. However, if the cfMeDIP-seq method was to be applied in clinical use, we hypothesize that for enhancing cost-efficiency the sequencing depth would have to be lower than in the demonstrative data set shown in [7]. But how well would the classification methods presented in [7] cope with lower sequencing depth? In this work we attempt to simulate a situation where the sequencing depth would be considerably lower. Additionally, we present statistical methods for improving the feature selection and probabilistic modeling to improve the classification of the cancer types. We compare our approaches to the machine learning methods presented in [7]. For feature selection, we experimented with classical principal component analysis (PCA) and iterative supervised PCA (ISPCA) [20], which can utilize the class information for finding the optimal principal components for separating classes from each other. We also tested Fisher's exact test for DMR finding, as a simpler statistical test could be more robust when the sequencing depth is lower. For the classification methods, we experimented with logistic regression with regularized horseshoe (RHS) prior [21] and logistic regression with DMR count variables, both implemented with probabilistic programming language Stan [22].

Table 1: Number of samples in each class in discovery and validation cohorts.

Class	Discovery cohort	Validation cohort
Healthy controls (Normal)	24	62
Renal cancer (RCC)	20	-
Pancreatic ductal adenocarcinoma (PDAC)	24	47
Colorectal cancer (CRC)	23	-
Bladder cancer (BLCA)	20	-
Breast cancer (BRCA)	25	-
Lung cancer (LUC)	25	55
Acute myeloid leukemia (AML)	28	35
Total	189	199

Methods

Aim of the study

The aim of this study was to design and test various feature selection and probabilistic classification methods and compare them to the methods presented in [7] on cfMeDIP-seq data across varying sequencing depths and cancer types.

Description of materials

The cfMeDIP-seq data set used in this work was received from the authors of [7] by request. In this work we use read count data, where the sequencing read counts have been determined for genomic windows of length 300bp. The details of the data processing can be found from [7]. The discovery and validation cohorts consisted of 189 and 199 samples respectively. The number of samples in each of the eight classes (corresponding to the healthy controls and 7 cancer types) is presented in Table 1.

Workflow

The workflow of the feature selection, model training and evaluation of the models followed the one presented in [7] but some modifications, such as sub-sampling of the cfMeDIP-seq count data, were done. First, both discovery and validation cohort data set were subsampled and 100 data splits of the discovery cohort were generated. In each data split, the discovery data was divided with 80%-20% ratio into balanced training and test sets using `caret` R package [23]. We utilized the scripts from [7] for generating the data splits. Then features were selected with different feature selection methods for each data split using the corresponding training data set, with a data transformation applied when applicable.

Using the found features, the probabilistic classifiers were trained using the training data. This resulted in 100×8 binary classifiers for each classification method. The classifiers were then used to classify the corresponding test data

sets to evaluate classifier performances. Finally, the trained models for each data split were applied to the validation cohort and the classifier performances were again evaluated.

Data preprocessing

Data subsampling

The data subsampling was done to simulate a lower sequencing depth than in the original data. The total read count in the discovery cohort data set, calculated from the preprocessed read count data, varied between 10659729 and 67228099, before extracting the 505027 genomic windows of interest. The thinning was done by sampling the original reads from all genomic windows without replacement. The probability of obtaining a read from a genomic window is the number of reads in that window in the original data divided by the total read count. The total read count of the thinned sample is the number of reads sampled from the original. The thinned total read counts per sample used in this work were 10^4 , 10^5 and 10^6 , where the highest value is already a magnitude lower than in the original, non-thinned data. After thinning the 505027 genomic windows can be extracted.

Data transformation

Depending on the classification model, the count data could be used as it is or transformed. We used logarithmic transformations as proposed in [7]

$$X_T = \log_2(c \cdot X + s). \quad (1)$$

The transformation used in [7] is obtained with $c = 0.3$ and $s = 10^{-6}$, but we also experimented with a modified version where $s = 0.5$. The difference of these transformations is best visible for the zero-count genomic windows. The original transformation maps the zero counts far away from the nonzero counts, while with the modified version the gap between the zero counts and nonzero counts is more moderate.

Scaling and normalization

Before fitting probabilistic classification methods, the count data was normalized based on the total read count in the 505027 genomic windows. This was done by dividing the read counts of these genomic windows with the sum of the read counts and multiplying with the mean of the read count sums over the discovery cohort. The read count normalization accounts for possible differences in sequencing depth between samples. In the case of the subsampled data, where the total read counts per sample have been made the same, this step should only have a small effect, but with the non-thinned data the total read counts can vary greatly. Also, the data was standardized to have zero mean and standard deviation of 1. The scaling was done based on the training data of each data

split. The scaling was done to standardize the different features, so that the same prior mean and scale can be used for all features during the probabilistic classification step.

Methods for feature selection

Moderated t-statistic

To generate results with the same methods as in [7], we used the same moderated t-statistic as implemented in the `limma` R-package [24] to find DMRs. 150 hypo- and hypermethylated DMRs were picked for each one class versus one other class comparison, totaling in $7 \times 300 = 2100$ DMRs per class. The 2100 DMRs are necessarily not unique, so the number of unique DMRs can be lower. This was repeated for each class and each 100 data split. Two versions of the modified t-statistic based DMRs were produced: one with data that was transformed with the original data transformation and one with data that was transformed with the modified version of the data transformation.

Fisher’s exact test

Fisher’s exact test was performed to the count data as an alternative method to the moderated t-statistic based DMR finding. For each genomic window, a contingency table of one class versus one other class comparison setting was formed on the training data and p-value was calculated with `fisher.test` function from the R package `stats` [25]. Then 150 hyper- and hypomethylated genomic windows with smallest p-values were picked as DMRs, totaling in $7 \times 300 = 2100$ DMRs per class. The 2100 DMRs are necessarily not unique, so the number of unique DMRs can be lower. This was repeated for each class and each 100 data split. Genomic windows with zero counts for all samples in the training data set were removed before conducting the tests.

PCA

We utilized the 2100-dimensional¹ DMRs as the input vectors for principal component analysis (PCA) and used the projections of the DMR vectors on the principal components as features for the binary classifiers. PCA was conducted for each data split and each class separately, using `prcomp` function from R package `stats` [25]. The data was shifted to be zero-mean and scaled to have variance of 1 by using `center` and `scale` options of the `prcomp` function. The found components were standardized by dividing them with their standard deviations.

ISPCA

Iterative supervised principal component analysis (ISPCA) [20] is a method for finding features that are most relevant for predicting the target value. Following

¹In case of overlapping DMRs, the dimension is less than 2100.

the notation of the description of the algorithm in [20], let us call the matrix of size $N_{\text{samples}} \times N_{\text{features}}$ containing the original features as \mathbf{X} and the target value vector of length N_{samples} as \mathbf{y} . The process of finding supervised components iterates three steps. First, scores $S(\mathbf{x}_j, \mathbf{y})$, $j = 1, \dots, N_{\text{features}}$ which tell how relevant each feature \mathbf{x}_j is for predicting target value \mathbf{y} are calculated and the features that have a score higher than threshold γ are chosen. From these features, \mathbf{X}_γ , the first principal component \mathbf{v}_γ is calculated. The threshold γ should be chosen so that the score $S(\mathbf{z}_\gamma, \mathbf{y})$, where \mathbf{z}_γ is the projection of \mathbf{X}_γ onto the principal component, is maximized. Finally, the variation explained by the found feature is subtracted from \mathbf{X} , so that a modified feature matrix \mathbf{X}' is retrieved. \mathbf{X}' is then used as the starting point for the next iteration.

For our case with eight classes, both binary and multiclass approaches are possible. In the case of multiclass setting, features maximizing the score for one class versus other classes are found for each class separately, and the feature with maximum score is picked, thus maximizing the relevance of distinguishing one of the classes from the others. With the multiclass approach ISPCA needs to be run only once per data split. For the binary approach, the multiclass labels have to be transformed into one-class-versus-other-classes-type of binary labels before running ISPCA. ISPCA is then run for each class separately.

As finding too many supervised components might lead to overfitting, the ISPCA method includes a permutation test approach to calculate the p-value of there being relevant information left in \mathbf{X}' . This test can be conducted after each iteration of finding supervised components and when the p-value exceeds a desired threshold, no more supervised components are searched. \mathbf{X}' can then be used for performing non-supervised PCA to retrieve up to $N_{\text{samples}} - 1$ components in total. We used the default threshold for the p-value, which is 0.1.

We used the implementation from R package `dimreduce` [26], namely function `ispca`. We gave the read counts for all 505027 genomic windows as input data to ISPCA. Before running ISPCA, the data was normalized for total read counts and transformed with the new version of the transformation. We used the `center` and `scale` options of the `ispca` function to make the data zero-mean and to have unit variance. We also used the option `normalize` so that the extracted features would have standard deviation of 1.

Classification methods

GLMnet

We use the GLMnet method as described in [7] and utilized the provided R scripts for reproducing results. Briefly, the model training utilities from R package `glmnet` [9] are used to learn a binomial GLMnet model to classify one class from the other classes using the DMRs found with moderated t-statistic as features. The binomial GLMnet model corresponds to the logistic regression that uses elastic net regularization on the model coefficients [9]. The model is fitted

by maximizing the penalized log-likelihood

$$\max_{(\beta_0, \boldsymbol{\beta}) \in \mathbb{R}^{p+1}} \frac{1}{N} \sum_{i=1}^N y_i (\beta_0 + \mathbf{x}_i^T \boldsymbol{\beta}) - \log \left(1 + e^{(\beta_0 + \mathbf{x}_i^T \boldsymbol{\beta})} \right) - \lambda P_\alpha(\boldsymbol{\beta}), \quad (2)$$

where y_i is the binary response variable, \mathbf{x}_i is the corresponding feature vector, N is the number of observations, β_0 and $\boldsymbol{\beta}$ are the intercept and coefficient parameters of the model and $P_\alpha(\boldsymbol{\beta})$ is the penalty term multiplied with penalty parameter λ . Elastic net penalty is the sum of ridge-regression and Lasso penalties

$$P_\alpha(\boldsymbol{\beta}) = (1 - \alpha) \frac{1}{2} \|\boldsymbol{\beta}\|_{\ell_2}^2 + \alpha \|\boldsymbol{\beta}\|_{\ell_1}, \quad (3)$$

where mixing parameter α controls the proportions of the two penalty terms. If $\alpha = 0$, elastic net simplifies into ridge regression and if $\alpha = 1$ the penalty term becomes the same as in Lasso. As in [7], the parameters λ and α are optimized using three iterations of 10-fold cross validation for grid values $\lambda = \{0, 0.01, 0.02, 0.03, 0.04, 0.05\}$ and $\alpha = \{0, 0.2, 0.5, 0.8, 1\}$. Before training the binary classifiers, data transformation was applied. Training data from each data split was used for training the models, which were then used to predict the class of each sample in the test data set. We also experimented with GLMnet model that uses the DMRs found with Fisher's exact test or with the moderated t-test using the modified version of the transformation.

Logistic regression with regularized horseshoe prior

In logistic regression model, each of the elements in the target vector $\mathbf{y} = (y_1, \dots, y_N)$ containing binary outcomes is assumed Bernoulli distributed with parameter p_i , $i = 1, \dots, N$, where N is the number of samples. A linear model can then be connected to parameter p_i with inverse-logit function

$$p_i = \frac{1}{1 + \exp^{-(\beta_0 + \mathbf{x}_i^T \boldsymbol{\beta})}}, \quad (4)$$

where β_0 is an intercept term, $\boldsymbol{\beta}$ is a coefficient vector and \mathbf{x}_i is a vector containing the values for the chosen features for sample i . After estimating the model coefficients, the class of a new sample can be predicted by calculating \hat{p} and using 0.5 (or some other value) as decision boundary. To classify the eight different classes, a logistic regression model is learned for each class separately by first transforming the class labels into binary labels using one-class-versus-the-other-classes approach.

The regularized horseshoe prior [21] is a technique to achieve sparsity in a regression model when the number of features is large and only few of them are expected to be relevant, and thus should have nonzero regression coefficient. The regularized horseshoe prior enforces sparsity to the regression coefficients by defining the scale of the coefficients to be a product of local and global terms, where the global term pulls all coefficients towards zero, while the local

term allows the relevant features to have nonzero coefficients. The prior for the regression coefficients β_j can be expressed more formally as

$$\beta_j | \lambda_j, \tau, c \sim N(0, \tau^2 \tilde{\lambda}_j^2), j = 1, \dots, D, \quad (5)$$

where D is the number of features and the modified local scale parameter $\tilde{\lambda}_j^2$ is defined as

$$\tilde{\lambda}_j^2 = \frac{c^2 \lambda_j^2}{c^2 + \tau^2 \lambda_j^2}, \quad (6)$$

where the local scale parameter λ_j is given a half-Cauchy prior

$$\lambda_j \sim C^+(0, 1). \quad (7)$$

The modified local scale parameter makes sure that all coefficients are shrunk at least a little, including the relevant coefficients as well. The parameter c controls the magnitude of the largest coefficients and is given an inverse-Gamma prior

$$c \sim \text{Inv-Gamma}(\nu/2, \nu s^2/2),$$

where ν and s were set to 4 and 2 respectively, following the default values of the regularized horseshoe prior implementation in R package **brms** [27]. The global scale parameter τ has also a half-Cauchy prior

$$\tau \sim C^+(1, \tau_0) \quad (8)$$

with scale τ_0 defined as

$$\tau_0 = \frac{p_0}{D - p_0} \frac{\sigma}{\sqrt{n}}, \quad (9)$$

where D is the number of features, n is the number of training samples, p_0 is the expected number of nonzero coefficients and σ is a pseudo standard deviation. The RHS prior enables using the knowledge on the expected number of nonzero coefficients to define the global scale parameter hyperprior. In our case p_0 was defined to be 300 after experimenting with a few different options and comparing the resulting coefficient posteriors. The pseudo variance for a model with binomial data and logit link function proposed in [21] is

$$\sigma^2 = \mu^{-1}(1 - \mu)^{-1}, \quad (10)$$

where μ is replaced with sample mean \bar{y} . The intercept term β_0 is handled separately and is given a Gaussian prior with mean 0 and standard deviation of 10.

The logistic regression model with horseshoe prior was fitted with different kinds of features: moderated t-statistic DMRs with the original data transformation, the moderated t-statistic DMRs with the modified data transformation, Fisher's exact test DMRs, PCA coordinates and ISPCA coordinates. The original data transformation was applied before model fitting on the corresponding moderated t-statistic DMRs and Fisher's exact test DMRs, while the modified data transformation was applied on the corresponding moderated statistic

DMRs. Before logistic regression we normalized the features for total read counts, standardized each feature to have zero mean and variance of 1.

The model was implemented with R version of the probabilistic programming language Stan [28]. For the logistic regression with horseshoe prior, we adopted the example code presented in [21], which used model parametrization proposed in [29].

The predictions for test and validation data sets were made by calculating the posterior predictive probabilities. For the i^{th} test/validation sample (y_i, \mathbf{x}_i) we compute the probability of belonging to the class of interest ($y_i = 1$) using posterior samples of β_0 and β

$$p(y_i = 1 | \mathbf{x}_i, \mathbf{y}, X) = \int p(y_i = 1 | \mathbf{x}_i, \theta) p(\theta | \mathbf{y}, X) d\theta \quad (11)$$

$$\approx \frac{1}{N_S} \sum_{k=1}^{N_S} p(y_i = 1 | \mathbf{x}_i, \theta_k) \quad (12)$$

$$= \frac{1}{N_S} \sum_{k=1}^{N_S} \text{logit}^{-1}(\beta_{0,k} + \mathbf{x}_i^T \beta_k), \quad (13)$$

where N_S is the number of samples retrieved from the posterior distribution and $\theta_k = (\beta_{0,k}, \beta_k)$ denotes the k^{th} parameter sample.

Logistic regression with DMR count variables

We also formulated a simpler logistic regression model for binary classification, where the model includes only an intercept term and two features. These two features are the numbers of hypermethylated and hypomethylated DMRs with nonzero read counts. The DMRs can be found either with the moderated t-statistic or Fisher's exact test. The two features were normalized based on total read counts. The features were scaled to have zero mean and standard deviation of 0.5 and the model intercept and coefficients were given Cauchy priors with scale parameters 10 and 2.5, respectively, as recommended in [30]. The model was implemented with Stan.

The classification of the test and validation sets are done in the same way as for the logistic regression model with RHS prior.

Sampling from posterior distributions with Stan

Stan uses MCMC sampling to retrieve samples from the posterior distribution, specifically the no-U-turn sampler (NUTS) algorithm, which is a variant of the Hamiltonian Monte Carlo algorithm [22]. The parameters for sampling the user can define include the number of MCMC chains, number of samples per chain, maximum tree depth and target acceptance rate. The sampling parameters for our models are presented in Table 2.

Table 2: Sampling parameters for the models implemented with Stan. Values marked with * are the default values in Stan.

Sampling parameter	LR RHS	LR RHS ISPCA/PCA	LR DMR counts
Number of chains	4*	4*	4*
Number of samples per chain	2000*	4000	3500
Max. tree depth	10*	10*	10*
Target acceptance rate	0.8*	0.99	0.9

Evaluation of the classification methods

The evaluation of the classification methods was performed in similar manner as in [7] and the distributed code for the publication was utilized when implementing the methods. For each method, each of the eight binary classifiers for each data split were used to classify the corresponding test data set and the validation data. The performance of the classifiers was evaluated by calculating class-specific area under receiver-operating characteristics curve (AUROC) statistics. The distribution of the AUROC values for each class over data splits could then be described by calculating median and quantile statistics or by plotting boxplots. The former could further be presented as barplots or scatterplots. For the validation cohort, we also calculated a mean of the class predictions over the data splits and plotted a receiver-operating characteristics curve (ROC) and calculated corresponding AUROC values for each of the four classes in the cohort.

Results

Feature selection

Comparison of the DMR finding methods

We tested three different methods to find differentially methylated regions that can be utilized as features in the classification: the moderated t-test method used in [7], moderated t-test with new data transformation and Fisher’s exact test. Each of these methods is used to pick 300 DMRs per a one class vs. other classes comparison using the training data, totaling in $300 \times 7 = 2100$ DMRs for each class. This is repeated for every data split. To compare the DMRs found with the different methods, we plotted Venn diagrams to find their overlaps. In Fig. 1, for each method, all of the DMRs for the 100 data splits were first combined and duplicates were removed to keep each DMR in the set only once. This was done for each class separately. Then Venn diagrams were plotted to show the overlaps between the methods. One Venn diagram is plotted for each class, and this is repeated for the three thinning versions. In Supplementary Fig. S1 we did the same, but only DMRs that were found in 50 or more data

splits of the total 100 were kept for the comparison.

Comparing Fig. 1 and Supplementary Fig. S1, we can see that the number of DMRs is overall higher in Fig. 1, where the DMRs were not filtered. This suggests, that many of the DMRs are only found in less than half of the data splits. In Fig. 1 we can also notice that overall the number of DMRs is smaller when the total read count is smaller, i.e. data has been thinned more. This might mean, that when the total read count is smaller, the found DMRs are more consistent between different data splits, resulting in a smaller number of DMRs. Supplementary Fig. S1 supports this, as the number of filtered DMRs is overall higher in Supplementary Fig. S1A than in Supplementary Fig. S1B and C. However, the overlap between all three methods is low in Supplementary Fig. S1A, indicating the different DMR finding methods work rather inconsistently in the case where the total read count is low.

When comparing the different methods, in most cases it seems that a big fraction of the DMRs is shared with all of the three methods. The overlap between the Fisher's exact test and t-test with the new data transformation is often quite high, as is the number of DMRs unique to the original t-test method. The numbers of DMRs unique to either Fisher's exact test or the t-test with the new data transformation are often low in comparison. The overlaps between the original t-test method and the two other methods separately are quite modest compared to the overlap between Fisher's exact test and t-test with new transformation. Altogether, it seems that a large part of the DMRs is shared between all of the methods, which suggests that DMRs can indeed be identified reliably from cfMeDIP-seq data. On the other hand, there are also DMRs that are not shared by all three methods. These DMRs may partly explain the differences in the performances of the classifiers utilizing these DMR sets.

Comparison to the RRBS-seq based DMCs

In [7], DMRs identified from cfMeDIP-seq data were compared to differentially methylated cytosines (DMCs) identified from reduced representation bisulfite sequencing (RRBS-seq) data that was obtained from solid samples. Shen et al. (2018) presented two sets of RRBS-seq DMCs: from comparisons between normal and tumor tissues as well as between tumor tissue and normal peripheral blood mononuclear cells (PBMCs). The comparison of cfMeDIP-seq DMRs and RRBS-seq DMCs presented in [7] showed that there was significant enrichment in concordantly hypermethylated and hypomethylated cfMeDIP-seq DMR and RRBS-seq DMC pairs. To see if there is overlap still after data subsampling, we made a simple comparison between the DMRs found from the thinned cfMeDIP-seq data and the two types of DMCs provided in [7]. The comparison was carried out by finding the cfMeDIP-seq DMRs from PDAC class vs. normal class comparison which had one or more overlapping DMCs. The direction of differential methylation was required to be the same in the RRBS-seq DMCs and cfMeDIP-seq DMRs when finding the overlaps. The number of such cfMeDIP-seq DMRs was calculated for each data split for each of the three subsampling versions.

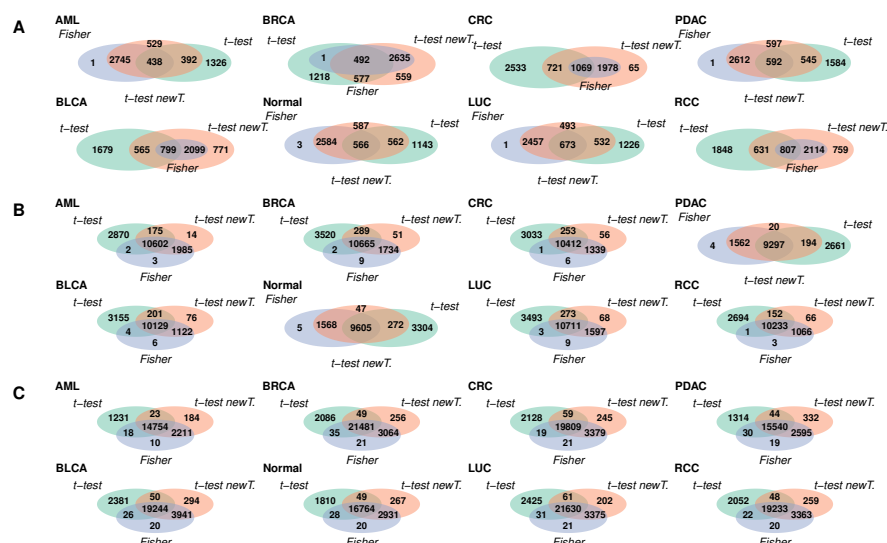


Figure 1: Number of overlapping DMRs between the DMRs found with Fisher's exact test, moderated t-test with the original data transformation and moderated t-test with new data transformation. DMRs from the 100 data splits have been combined before the comparison. The results are presented for thinning with total read counts 10^4 , 10^5 and 10^6 in figures A, B and C, respectively.

This was repeated for the three DMR finding methods: the original moderated t-test approach, moderated t-test approach with the new data transformation and Fisher's exact test.

The results of the comparison in Fig. 2 show that the number of cfMeDIP-seq DMRs with overlapping DMCs in a data split is overall quite low, ranging from 0 to 14. The overlap with the DMCs from the comparison of tumor tissue and normal tissue is generally lower than with the DMCs from the comparison of tumor tissue and PBMC. The severity of the thinning does not seem to affect the number of cfMeDIP-seq DMRs with RRBS DMC overlap, suggesting that the most significant DMRs can still be identified from the subsampled data as well. Overall the level of overlaps with the RRBS-seq DMCs seems to be the same for all three DMR finding methods, with the exception of the thinning with total read count 10^4 , where the moderated t-test with new data transformation has considerably more DMRs with DMC overlap than the two other methods.

PCA and ISPCA

Fig. 3 demonstrates the retrieved features from the three different dimension reduction approaches, PCA, binary ISPCA and multiclass ISPCA, for one of the data splits in the case of the subsampled data with total read count 10^6 .

In Fig. 3A the six first principal components from the standard PCA have

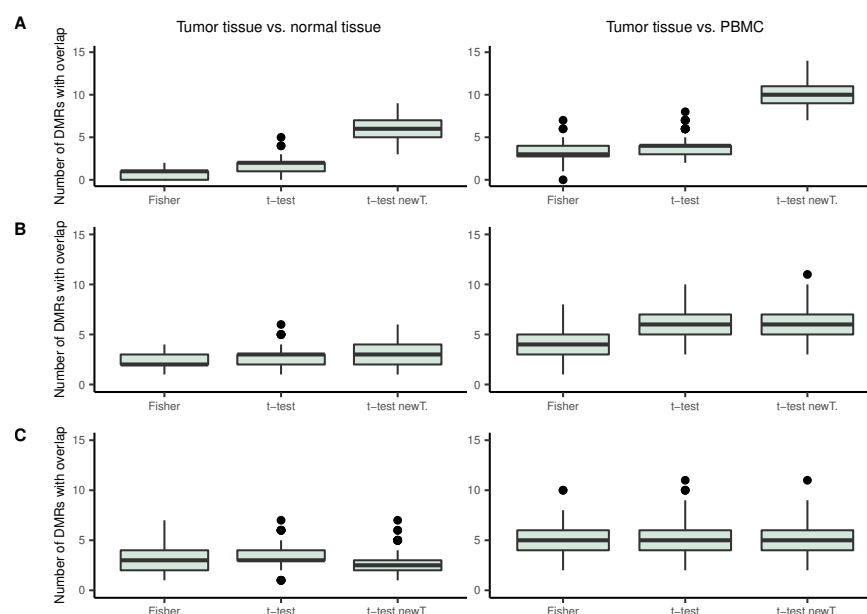


Figure 2: Boxplots of the number of cfMeDIP-seq DMRs having overlap with RRBS-seq DMCs in one data split, for each of the DMR finding approaches. The left side column shows the results for the RRBS-seq DMCs from comparing PDAC tumor tissue to normal tissue, while the right side column shows the results for RRBS-seq DMCs from comparing PDAC tumor tissue to PBMC. The results are presented for thinning with total read counts 10^4 , 10^5 and 10^6 in figures **A**, **B** and **C**, respectively.

been plotted, two at a time. The input data to the analysis was the DMR set for the AML class, and based on the plot, the AML class can be well separated from the other classes using PC1. The training and test set samples from AML class seem to mix to some extent, which indicates that the principal component generalizes well to the test samples. The other plotted components seem not to separate AML class from the others.

Similar plots were generated for the binary and multiclass ISPCA approaches, presented in Fig. 3B and Fig. 3C, respectively. For the binary ISPCA we set the samples from the AML class to have label 1 and for the other classes the label was set to 0. We gave data from all 505027 genomic windows as input to the ISPCA. This results again in the first component separating the AML class from the other classes quite well, while the test samples blend in with the training samples. The rest of the plotted principal components do not separate the AML class from the other classes. The first six components from multiclass ISPCA seem to each separate one class from the others. The only class for which the test samples blend in with the plotted training samples is AML, which could

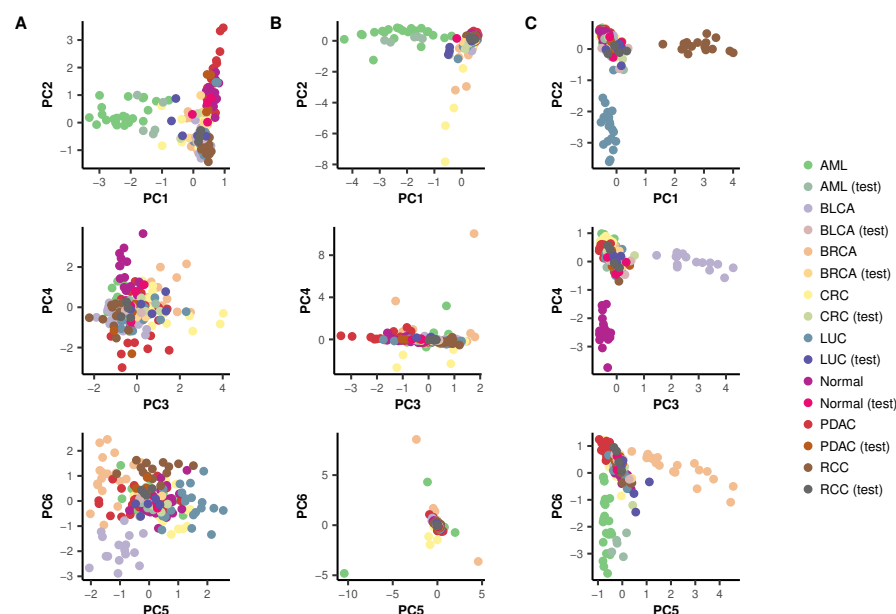


Figure 3: An example of the six first principal components from PCA and two ISPCA approaches for one of the data splits. The training and test samples from the discovery cohort have been plotted with different colors to demonstrate how well the test samples will blend in with the training samples. **A**: Components from PCA, using AML class DMRs as input to the analysis. **B**: Components from binary ISPCA, where class labels have been set to 1 for the samples in AML class and to 0 for the samples from other classes. **C**: Components from multiclass ISPCA.

indicate that making good predictions with a model that uses these components as features could be hard.

With the two other subsampled data versions, with total read counts 10^4 and 10^5 , the ISPCA approaches occasionally end up in a situation where no supervised components are found. Before finding next supervised component, ISPCA method performs a permutation test to see whether there is enough information in the data to find relevant supervised components that would not overfit to the training data. If the test fails, the data is used for computing standard PCA components instead. In the case of more heavily thinned data, there seems sometimes not to be enough information to find any supervised components, and the method returns standard PCA components only. As we give all the 505027 genomic windows as input to the method, this might not always lead to components which would separate the classes well. On the other hand, the standard PCA approach which utilizes the class-specific DMRs does not thrive in such a situation either.

Discovery cohort

After finding the model features, the different classification models were fitted using the training data sets of the discovery cohort. The data was partitioned to training and test data sets 100 times, and each of the trained models was evaluated with corresponding test set. This resulted in 100 AUROC values per each of the eight classes, from which we can calculate median AUROC values which are presented in Supplementary Fig. S2. Based on this figure, the overall trend is that the lower the total read count is after the data subsampling, the lower the median AUROC values are. This is expected, as the more the data is thinned the less there is information for us to use for the classification task. When the total read count is 10^4 , the median AUROC levels are very similar for all methods and for all classes. When the total read count is higher, some classes such as AML and PDAC begin to stand out with higher median AUROC values, while some classes such as BRCA, CRC and LUC have lower median AUROC values even with higher total read counts.

To better compare each of the methods to the original approach, we computed the differences between median AUROCs (Fig. 4). The differences were calculated for each class separately, but a mean over classes is also presented for overall performance comparison purposes. Looking at the means of the AUROC median differences over classes, all GLMnet-based approaches seem to perform overall equally well. However, there is some variance in the class-specific median AUROC differences for the GLMnet with Fisher's exact test DMRs as features when the subsampled data had total read count 10^4 or 10^5 . Similarly, the two logistic regression models with RHS prior using the original moderated t-test DMRs or the moderated t-test with new data transformation DMRs seem to work equally as well as the original GLMnet method in all of the three thinning versions. The LR model with RHS prior and Fisher's exact test DMRs, on the other hand, has a positive mean AUROC difference value when the sequencing depth is low. The logistic regression models with DMR count variables seem to perform overall as well as the original GLMnet method or slightly worse.

Continuing the comparisons with the LR models that use dimension reduction features, the logistic regression model with RHS prior that uses PCA components as features works approximately equally well as the original GLMnet method for all subsampled data versions. But perhaps the most promising methods in this comparison are the logistic regression models with RHS priors that use ISPCA components as features. Both the models using multiclass ISPCA and binary ISPCA components got higher median AUROC values than the original GLMnet method for most of the classes in Fig. 4 **A** and **B**. Surprisingly, when the total read count after thinning is set to 10^6 (Fig. 4 **C**), the model using multiclass ISPCA components performs overall considerably worse than the original GLMnet method. However, the mean median AUROC differences over all classes is still slightly on the positive side for the model using binary ISPCA components.

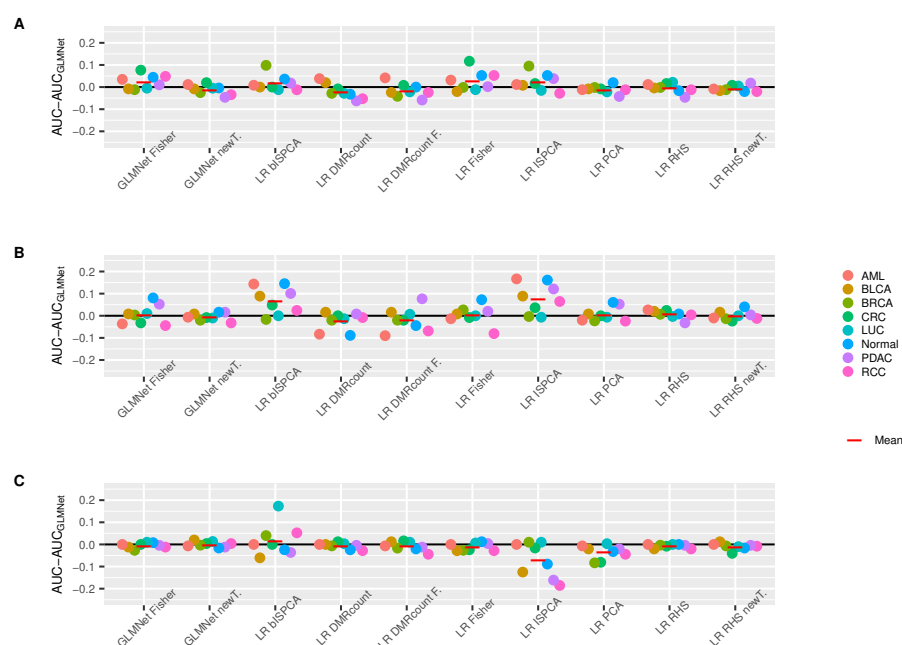


Figure 4: The differences between the test set AUROC medians over the 100 data splits for the original GLMnet method and other approaches. The AUROC has been calculated for the test data sets in the discovery cohort. The results are presented for each class separately. Positive values indicate, that the compared (new) method has higher AUROC median than the original method and negative values indicate that the original method has higher AUROC median. The results are presented for thinning with total read counts 10^4 , 10^5 and 10^6 in figures **A**, **B** and **C**, respectively. The red lines indicate the means over the eight classes.

Validation cohort

The models trained with the discovery cohort training sets were used to predict the class of samples in validation data set. The validation data was also thinned in the same manner as the discovery cohort data. The predictions were made with the models for each of the 100 data splits, and AUROC values were calculated for each set of predictions. This resulted in 100 AUROC values for each of the four classes in the validation set.

The AUROC medians over the 100 values were calculated and they can be found from Supplementary Fig. S3. As for the discovery cohort in Supplementary Fig. S2, the median AUROC values seem to be overall higher when the total read count is higher. When compared to the discovery cohort results, there seems to be more variance in the performance of different methods for the validation data.

We also calculated the median AUROC differences between the original GLMnet method and the new methods (Fig. 5). The logistic regression models with RHS prior that use the components from multiclass or binary ISPCA did well in the discovery cohort performance comparison. Conversely, these two models do not perform very well with the validation data and at best perform equally well as the original GLMnet method. Logistic regression with RHS prior using PCA components seems to be doing slightly better than the ISPCA approaches, the mean of median AUROC differences reaching positive value for the case with total read count 10^5 . The overall performance of the logistic regression models with DMR count variables is usually as good or slightly weaker than for the original approach. However, for the subsampled data with total read count 10^4 , the approach using Fisher’s exact test DMRs has a considerably positive mean difference value. The GLMnet methods with the alternative DMR choices both perform approximately equally well as the original GLMnet method with all of the three subsampling versions.

All in all, compared to the discovery cohort results, there seems to be a lot more variation in the AUROC differences for the validation data set. When the total read count is 10^6 , i.e. closest to the original read counts, none of the methods have positive mean difference over the classes. But when the thinning becomes more severe, there are some cases where other methods outperform the original method.

As in [7], we also calculated the mean of the validation set predictions over the 100 data splits and used the mean predictions to produce ROC curves for each class. In Fig. 6, we can see a general trend of the area under the curves getting higher the larger the total read count is. The ROC curves confirm the findings we did based on the median AUROC differences. When the total read count is 10^6 , several methods perform approximately equally well as the original GLMnet method, but rarely outperform it. When the total read count is 10^5 or 10^4 , some methods have their ROC curves surmounting the original method’s curve for some of the classes.

Non-thinned data

Until now, the results we have presented have been for the thinned data set. To assess how well the classifiers presented in this work perform with the original non-thinned data set, we ran the feature selection methods and classifier training to the data without thinning it first. We calculated the median AUROC values, the median AUROC differences between the original GLMnet method and the other classifiers for both the discovery and validation cohorts and produced ROC curves of the mean predictions for the validation cohort. The results are presented in Supplementary Figs. S4 and S5.

The median AUROC values over the 100 data splits are presented in Supplementary Fig. S4 for both discovery and validation cohorts. We can see that overall the AUROC medians are higher on the more deeply sequenced, non-subsampled data, both for the discovery and validation cohorts, than for the thinned data sets. The differences between the methods are moderate for the

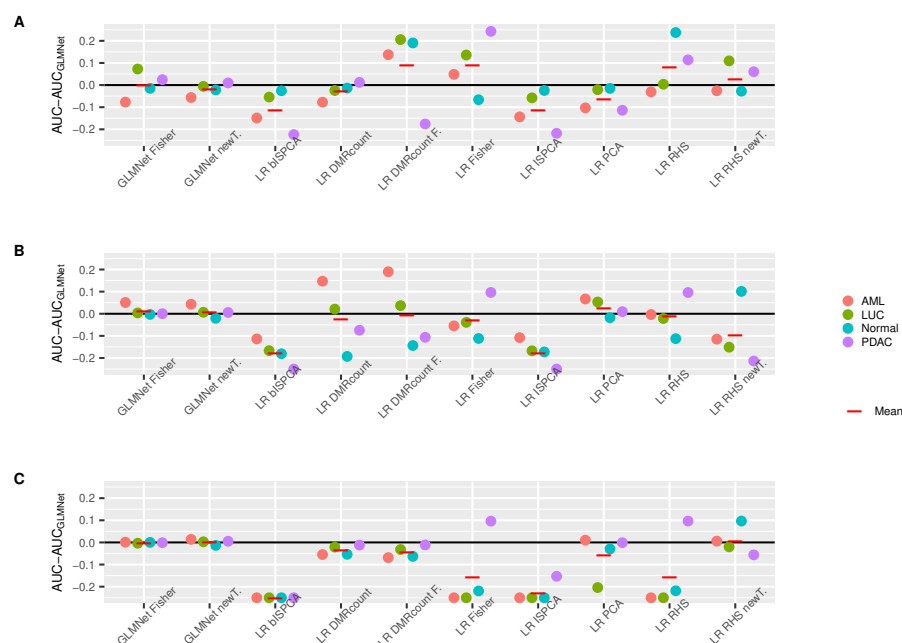


Figure 5: The differences between the validation cohort AUROC medians over the 100 data splits for the original GLMnet method and other approaches. The results are presented for each class separately. Positive values indicate, that the compared method has higher AUROC median than the original method and negative values indicate that the original method had lower AUROC median. The results are presented for thinning with total read counts 10⁴, 10⁵ and 10⁶ in figures **A**, **B** and **C**, respectively. The negative values have been truncated to -0.25. The red lines indicate the means over the four classes.

discovery cohort, but for the validation cohort there are greater differences between the performance of the methods. There are also differences between the classes, indicating that some classes are easier to distinguish from the other classes with the used methods. For example, considering the discovery cohort, all methods reach median AUROC values close to 1 for the AML class, while for LUC class the values are at best around 0.875.

In Supplementary Fig. S5A we present the median AUROC value differences between the original GLMnet method and the other methods for the discovery cohort. For most of the methods the mean differences are close to 0, meaning that overall there was no difference in the median AUROC values for the original GLMnet method and other approaches. However, the performance of logistic regression model with RHS prior and multiclass ISPCA components and logistic regression models with DMR count features, the mean difference values are

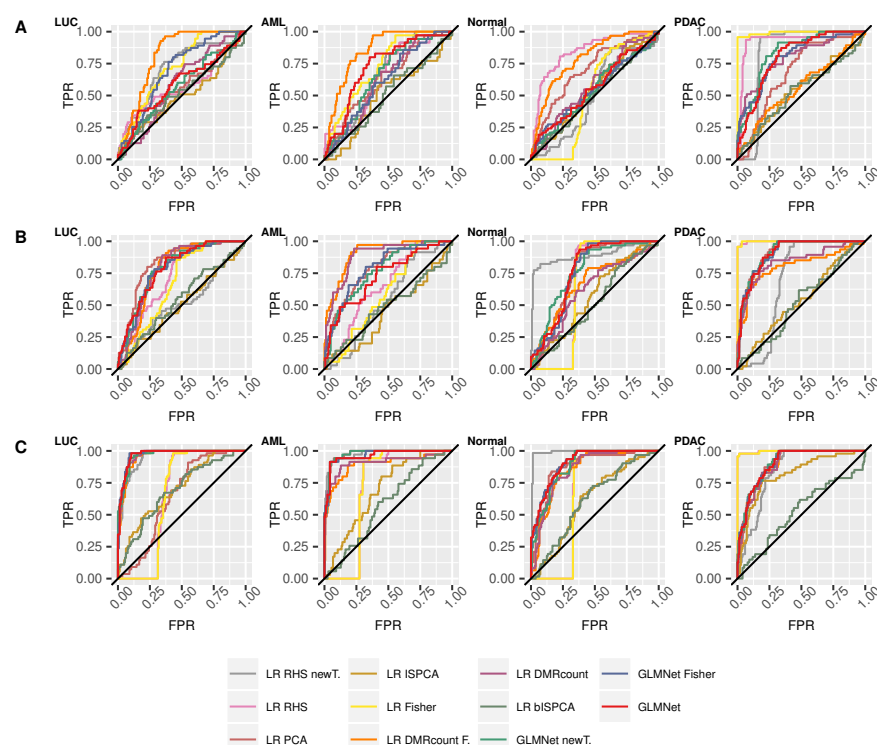


Figure 6: Validation cohort ROCs calculated with prediction means over 100 data splits and corresponding fitted models. The results are presented for each class separately. The results are presented for thinning with total read counts 10^4 , 10^5 and 10^6 in figures **A**, **B** and **C**, respectively.

negative.

Supplementary Fig. S5B presents the median AUROC differences between the original method and the other classifiers for the validation data set. Based on this comparison, the overall performance of the other classifiers is again in most cases equally good or even considerably weaker than for the original GLMnet method. There are few exceptions in the case of the PDAC class, for which the logistic regression models with RHS prior using either of the two moderated t-test or the Fisher's exact test DMRs have all approximately 0.1 better median AUROC value than the original method. Also, the modified GLMnet approaches both have approximately equally good results as the original method.

From the validation cohort ROC curves in Supplementary Fig. S5C we can see, that the original GLMnet method is performing similarly as in [7], with almost equally good performance for LUC, AML and Normal classes and slightly weaker performance for PDAC class. The corresponding AUROC values in [7] were 0.971, 0.980, 0.969 and 0.918 for classes LUC, AML, Normal and PDAC

respectively. The ROC curves for the other two GLMnet approaches behave quite similarly. For the rest of the methods, the ROC curves support the findings based on the median AUROC differences in Supplementary Fig. S5B.

Discussion

There seems to be a lot of variation in how well the compared methods performed with different classes, thinning versions and also between discovery and validation cohorts. For example, the LR model with RHS prior using binary and multiclass ISPCA features seemed to perform better than the original method when looking at the discovery cohort results for the lowest sequencing depth. But looking at the validation cohort performance, these two approaches had considerably weaker median AUROC values than the original GLMnet method. One possible cause to the differences between the three subsampled data versions is that the subsampling is done by taking random subsets of the original data, and that could cause differences even if the probabilities of obtaining reads from each of the genomic window were the same for all three thinning versions. The differences between the discovery and validation cohort data sets are visible with the non-thinned data set too, and this indicates that there are perhaps differences between the two cohorts, making it harder for the classifiers to perform well. However, the GLMnet-based methods seem to all cope quite well with the validation cohort, especially if the data is non-thinned or the data has not been thinned very much. It could also be, that not all DMR finding methods and classifiers are suited to all classes. It could be considered, that different classifiers would be used with different classes to obtain optimal performance.

The overlaps between the DMRs found with the three different approaches showed that there are indeed differences between the methods, even if many of the DMRs were shared between all methods. We noticed, that the Fisher's exact test and the moderated t-test with new data transformation often shared DMRs which were not found by the original moderated t-test method, while the moderated t-test found DMRs that were not found by other methods. The numbers of unique DMRs to Fisher's exact test or the moderated t-test with the new data transformation were often low in comparison. The DMRs that were not found by all methods could be a source of differences in the classification results.

The dimension reduction techniques, PCA and two versions of ISPCA, showed varying performance. While logistic regression models using ISPCA components did not perform well with the validation data cohort, their AUROC values were promising for the discovery cohorts when total read count was 10^4 and 10^5 . We also tested using DMRs as input data for multiclass and binary ISPCA. The classification results with LR model and RHS prior using the resulting principal components as features showed stable performance, but these approaches did not outperform the original GLMnet method neither in the discovery or validation cohorts.

Both logistic regression with RHS prior and GLMnet methods implement

the logistic regression model and enable sparsity of the feature coefficients. The difference between these two models is the prior enabling sparsity and how the model is fitted. In our approach we use the regularized horseshoe prior and GLMnet utilizes elastic net regularization. On the model fitting side, we utilized probabilistic programming language Stan to obtain posterior samples of the model, while GLMnet model is fitted with cyclical coordinate descent approach [9]. The GLMnet model fitting is very efficient, but while sampling with Stan requires more computational resources, we obtain samples that inform us of the whole posterior distribution. If we look at the performance of the LR RHS method with the moderated t-test DMRs as features, we notice that with the thinned data sets the performance is approximately equally good as for the original GLMnet method. The same applies to the LR RHS method with moderated t-test with new data transformation. For the validation cohort, there are both classes for which the performance is either considerably better or considerably weaker. Based on this comparison, our implementation of the logistic regression method with enabled model sparsity seems to have potential to give better results than GLMnet method, but it is perhaps not as stable as GLMnet. Combining promising feature selection methods such as ISPCA with the LR RHS model could further enhance the classification.

Based on the promising results in TCR analysis application [31], we expected the simple logistic regression model with two DMR count variables to perform well with the lowest sequencing coverage due to model robustness. The results with the thinned validation cohort with total read count 10^4 somewhat support this hypothesis, but the discovery cohort results were not as impressive.

Conclusions

We performed a method comparison to investigate if we the classification accuracy of the classifier based on cfMeDIP-seq data could be improved in a case where the sequencing depth of a cfMeDIP-seq experiments is low. To simulate lower sequencing depths, we thinned a cfMeDIP-seq data set by randomly sub-sampling the reads of each of the samples. We obtained three data sets with total read counts of 10^4 , 10^5 and 10^6 . Then we tested three different DMR finding methods and three different dimension reduction methods to find the features to be used in the classification. After finding the features, three different types of classifiers were trained, using the found features and performance was evaluated for discovery and validation cohorts. These steps were also performed for a non-thinned data set.

Based on the comparisons between the performances of the classification methods, there seems to be no one method that would consistently perform better with all thinning versions, all classes and for both discovery and validation cohorts. But there are cases, where the different feature selection and classifying methods seem to give advantage when the data has been thinned. Such methods include the Fisher's exact test, binary and multiclass ISPCA for DMR finding and feature selection and logistic regression model with DMR count variables.

Acknowledgements

The calculations presented above were performed using computer resources within the Aalto University School of Science “Science-IT” project.

Funding

This work was supported by the Academy of Finland (292660, 311584, 335436).

Abbreviations

AML: acute myeloid leukemia
AUROC: area under receiver-operating characteristics curve
BLCA: bladder cancer
BRCA: breast cancer
BS-seq: bisulfite sequencing
cfDNA: cell-free DNA
cfMeDIP-seq: cell-free methylated DNA immunoprecipitation and high-throughput sequencing
CRC: colorectal cancer
ctDNA: circulating tumor DNA
DMC: differentially methylated cytosine
DMR: differentially methylated region
GO: gene ontology
GP: Gaussian process
ISPCA: iterative supervised principal component analysis
Lasso: least absolute shrinkage and selection operator
LR: logistic regression
LUC: lung cancer
MCMC: Markov chain Monte Carlo
NGS: next-generation sequencing
PCA: principal component analysis
PDAC: pancreatic ductal adenocarcinoma
RCC: renal cell carcinoma
RHS: regularized horseshoe
RRBS-seq: reduced representation bisulfite sequencing

Availability of data and materials

Access to the cfMeDIP-seq data set was requested from the authors of [7]. Lists of the differentially methylated RRBS-seq based DMCs are available as Supplementary information of [7]. The scripts used to produce the results presented in this work are available at <https://github.com/hallav/cfMeDIP-seq>.

Authors' contributions

VH and HL developed methods. VH performed the computational analysis and wrote the manuscript. HL participated in the revision of the manuscript. Both authors have read and accepted the manuscript.

References

- [1] Yanni Y.N. Lui, Ki-Wai Chik, Rossa W.K. Chiu, Cheong-Yip Ho, Christopher W.K. Lam, and Y.M. Dennis Lo. Predominant Hematopoietic Origin of Cell-free DNA in Plasma and Serum after Sex-mismatched Bone Marrow Transplantation. *Clinical Chemistry*, 48(3):421–427, 03 2002. ISSN 0009-9147. doi: 10.1093/clinchem/48.3.421.
- [2] Frank Diehl, Meng Li, Devin Dressman, Yiping He, Dong Shen, Steve Szabo, Luis A. Diaz, Steven N. Goodman, Kerstin A. David, Hartmut Juhl, Kenneth W. Kinzler, and Bert Vogelstein. Detection and quantification of mutations in the plasma of patients with colorectal tumors. 102(45):16368–16373, 2005. ISSN 0027-8424. doi: 10.1073/pnas.0507904102.
- [3] Ellen Heitzer, Imran S. Haque, Charles E.S. Roberts, and Michael R. Speicher. Current and future perspectives of liquid biopsies in genomics-driven oncology. *Nature Reviews Genetics*, 20(2):71–88, 2019.
- [4] Stephen Cristiano, Alessandro Leal, Jillian Phallen, Jacob Fiksel, Vilmos Adleff, Daniel C. Bruhm, Sarah Østrup Jensen, Jamie E. Medina, Carolyn Hruban, James R. White, Doreen N. Palsgrove, Noushin Niknafs, Valsamo Anagnostou, Patrick Forde, Jarushka Naidoo, Kristen Marrone, Julie Brahmer, Brian D. Woodward, Hatim Husain, Karlijn L. van Rooijen, Mai-Britt Worm Ørntoft, Anders Husted Madsen, Cornelis J. H. van de Velde, Marcel Verheij, Annemieke Cats, Cornelis J. A. Punt, Geraldine R. Vink, Nicole C. T. van Grieken, Miriam Koopman, Remond J. A. Fijneman, Julia S. Johansen, Hans Jørgen Nielsen, Gerrit A. Meijer, Claus Lindbjerg Andersen, Robert B. Scharpf, and Victor E. Velculescu. Genome-wide cell-free dna fragmentation in patients with cancer. *Nature*, 570(7761):385–389, 2019.
- [5] Shuli Kang, Qingjiao Li, Quan Chen, Yonggang Zhou, Stacy Park, Gina Lee, Brandon Grimes, Kostyantyn Krysan, Min Yu, Wei Wang, Frank Alber, Fengzhu Sun, Steven M. Dubinett, Wenyuan Li, and Xianghong Jasmine Zhou. Cancerlocator: non-invasive cancer diagnosis and tissue-of-origin prediction using methylation profiles of cell-free dna. *Genome biology*, 18(1):1–12, 2017.
- [6] C. Grunau, S. J. Clark, and A. Rosenthal. Bisulfite genomic sequencing: systematic investigation of critical experimental parameters. *Nucleic Acids Research*, 29(13):e65–e65, 07 2001. ISSN 0305-1048. doi: 10.1093/nar/29.13.e65. URL <https://doi.org/10.1093/nar/29.13.e65>.

- [7] Shu Yi Shen, Rajat Singhania, Gordon Fehrer, Ankur Chakravarthy, Michael H. A. Roehrl, Dianne Chadwick, Philip C. Zuzarte, Ayelet Borgida, Ting Ting Wang, Tiantian Li, Olena Kis, Zhen Zhao, Anna Spreafico, Tiago da Silva Medina, Yadon Wang, David Roulois, Ilias Ettayebi, Zhuo Chen, Signy Chow, Tracy Murphy, Andrea Arruda, Grainne M. O’Kane, Jessica Liu, Mark Mansour, John D. McPherson, Catherine O’Brien, Natasha Leighl, Philippe L. Bedard, Neil Fleshner, Geoffrey Liu, Mark D. Minden, Steven Gallinger, Anna Goldenberg, Trevor J. Pugh, Michael M. Hoffman, Scott V. Bratman, Rayjean J. Hung, and Daniel D. De Carvalho. Sensitive tumour detection and classification using plasma cell-free dna methylomes. *Nature*, 563(7732):579–583, 2018.
- [8] Shu Yi Shen, Justin M. Burgener, Scott V. Bratman, and Daniel D. De Carvalho. Preparation of cfmedip-seq libraries for methylome profiling of plasma cell-free dna. *Nature protocols*, 14(10):2749–2780, 2019.
- [9] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software*, 33(1):1–22, 2010. URL <http://www.jstatsoft.org/v33/i01/>.
- [10] Pier Vitale Nuzzo, Jacob E. Berchuck, Keegan Korthauer, Sandor Spisak, Amin H. Nassar, Sarah Abou Alaiwi, Ankur Chakravarthy, Shu Yi Shen, Ziad Bakouny, Francesco Boccardo, John Steinharter, Gabrielle Bouchard, Catherine R. Curran, Wenting Pan, Sylvan C. Baca, Ji-Heui Seo, Gwo-Shu Mary Lee, M. Dror Michaelson, Steven L. Chang, Sushrut S. Waikar, Guru Sonpavde, Rafael A. Irizarry, Mark Pomerantz, Daniel D. De Carvalho, Toni K. Choueiri, and Matthew L. Freedman. Detection of renal cell carcinoma using plasma and urine cell-free dna methylomes. *Nature medicine*, 26(7):1041–1043, 2020.
- [11] Farshad Nassiri, Ankur Chakravarthy, Shengrui Feng, Shu Yi Shen, Romina Nejad, Jeffrey A. Zuccato, Mathew R. Voisin, Vikas Patil, Craig Horbinski, Kenneth Aldape, Gelareh Zadeh, and Daniel D. De Carvalho. Detection and discrimination of intracranial tumors using plasma cell-free dna methylomes. *Nature Medicine*, 26(7):1044–1047, 2020.
- [12] Feng Cao, Ailin Wei, Xinlei Hu, Yijing He, Jun Zhang, Lin Xia, Kailing Tu, Jue Yuan, Ziheng Guo, Hongying Liu, Dan Xie, and Ang Li. Integrated epigenetic biomarkers in circulating cell-free dna as a robust classifier for pancreatic cancer. *Clinical epigenetics*, 12(1):1–14, 2020.
- [13] Yong Zhang, Tao Liu, Clifford A. Meyer, Jérôme Eeckhoutte, David S. Johnson, Bradley E. Bernstein, Chad Nusbaum, Richard M. Myers, Myles Brown, Wei Li, and X. Shirley Liu. Model-based analysis of chip-seq (macs). *Genome biology*, 9(9):1–9, 2008.
- [14] Kathryn Lasseter, Amin H. Nassar, Lana Hamieh, Jacob E. Berchuck, Pier Vitale Nuzzo, Keegan Korthauer, Atul B. Shinagare, Barbara Ogorek,

- Rana McKay, Aaron R. Thorner, Gwo-Shu Mary Lee, David A. Braun, Rupal S. Bhatt, Matthew Freedman, Toni K. Choueiri, and David J. Kwiatkowski. Plasma cell-free dna variant analysis compared with methylated dna analysis in renal cell carcinoma. *Genetics in Medicine*, pages 1–8, 2020.
- [15] Charity W Law, Yunshun Chen, Wei Shi, and Gordon K Smyth. voom: Precision weights unlock linear model analysis tools for rna-seq read counts. *Genome biology*, 15(2):1–17, 2014.
- [16] Wei Xu, Jun Lu, Qiang Zhao, Jun Wu, Jieli Sun, Baohui Han, Xiaodong Zhao, and Yani Kang. Genome-wide plasma cell-free dna methylation profiling identifies potential biomarkers for lung cancer. *Disease markers*, 2019, 2019.
- [17] Shengyue Li, Lei Wang, Qiang Zhao, Zhihao Wang, Shuxian Lu, Yani Kang, Gang Jin, and Jing Tian. Genome-wide analysis of cell-free dna methylation profiling for the early diagnosis of pancreatic cancer. *Frontiers in genetics*, 11, 2020.
- [18] Matthias Lienhard, Christina Grimm, Markus Morkel, Ralf Herwig, and Lukas Chavez. Medips: genome-wide differential coverage analysis of sequencing data derived from dna enrichment experiments. *Bioinformatics*, 30(2):284–286, 2014.
- [19] Mark D. Robinson, Davis J. McCarthy, and Gordon K. Smyth. edgeR: a bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics*, 26(1):139–140, 2010.
- [20] Juho Piironen and Aki Vehtari. Iterative supervised principal components. In *International Conference on Artificial Intelligence and Statistics*, pages 106–114. PMLR, 2018.
- [21] Juho Piironen and Aki Vehtari. Sparsity information and regularization in the horseshoe and other shrinkage priors. *Electronic Journal of Statistics*, 11(2):5018–5051, 2017.
- [22] Stan Development Team. Stan modeling language users guide and reference manual, 2020. URL <https://mc-stan.org/>. 2.26.
- [23] Max Kuhn. *caret: Classification and Regression Training*, 2020. URL <https://CRAN.R-project.org/package=caret>. R package version 6.0-85.
- [24] Matthew E. Ritchie, Belinda Phipson, Di Wu, Yifang Hu, Charity W. Law, Wei Shi, and Gordon K. Smyth. limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Research*, 43(7):e47, 2015. doi: 10.1093/nar/gkv007.

- [25] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2020. URL <https://www.R-project.org/>.
- [26] Juho Piironen. *dimreduce: Supervised Dimension Reduction*, 2020. R package version 0.2.1.
- [27] Paul-Christian Bürkner. brms: An R package for Bayesian multilevel models using Stan. *Journal of Statistical Software*, 80(1):1–28, 2017. doi: 10.18637/jss.v080.i01.
- [28] Stan Development Team. RStan: the R interface to Stan, 2020. URL <http://mc-stan.org/>. R package version 2.19.3.
- [29] Tomi Peltola, Aki S. Havulinna, Veikko Salomaa, and Aki Vehtari. Hierarchical bayesian survival analysis and projective covariate selection in cardiovascular event risk prediction. In *BMA@ UAI*, pages 79–88. Cite-seer, 2014.
- [30] Andrew Gelman, Aleks Jakulin, Maria Grazia Pittau, and Yu-Sung Su. A weakly informative default prior distribution for logistic and other regression models. *The annals of applied statistics*, 2(4):1360–1383, 2008.
- [31] Ryan O. Emerson, William S. DeWitt, Marissa Vignali, Jenna Gravley, Joyce K. Hu, Edward J. Osborne, Cindy Desmarais, Mark Klinger, Christopher S. Carlson, John A. Hansen, Mark Rieder, and Harlan S. Robins. Immunosequencing identifies signatures of cytomegalovirus exposure history and hla-mediated effects on the t cell repertoire. *Nature genetics*, 49(5): 659–665, 2017.