

EpiSmoker: An Epigenetic Smoking Status Estimator

Classification of individuals based on self-reported smoking status is prone to errors due to under-reporting and poor recall of long-term smoking history. Traditional biomarkers, such as cotinine, only measure short-term exposure. Cigarette smoking strongly influences DNA methylation, with **differential methylation patterns** detected among current, former and never smokers. We developed, **EpiSmoker** a robust prediction tool to infer the smoking status based on the whole-blood DNA methylation profiles.

S BOLLEPALLI^{1,2}, T KORHONEN^{1,3,4}, J KAPRIO^{1,2}, M OLLIKAINEN^{1,2} and S ANDERS^{1,5}
¹Institute for Molecular Medicine (FIMM), University of Helsinki, Finland; ²Department of Public Health, University of Helsinki, Finland; ³National Institute for Health and Welfare (THL), Finland; ⁴Institute of Public Health and Clinical Nutrition, University of Eastern Finland, Finland; ⁵Centre for Molecular Biology (ZMBH), University of Heidelberg, Germany.

Contact: sailalitha.bollepalli@helsinki.fi

METHODOLOGY

- DNA methylation data from DILGOM (N=517 and Age range: 25-74) a sub-cohort of FINRISK was used as the training dataset.
- After stringent QC, quantile normalization was performed. **Quantiles** were later used to fit the distributions of validation datasets.
- Multinomial LASSO regression was performed to select a parsimonious set of CpG sites predictive of smoking status (Figure 1).
- Tuning parameter, **optimal lambda** was chosen through a tenfold internal cross-validation from a sequence of 100 lambda values.
- 121** methylation sites and the corresponding training coefficients selected from this regression were used to predict the smoking statuses in the three test datasets.
- Three independent and external test datasets were used to assess the performance.

RESULTS

- EpiSmoker** showed high accuracy for predicting current and never smokers in all the 3 test datasets.
- Figures 2 - 4 illustrate results from the Finnish Twin Cohort.
- Precise classification of self-reported former smokers from current and never smokers depends on the extent and duration of smoking prior to smoking cessation. **Longer cessation time leads to reversal of methylation profiles in former smokers making them indistinguishable from never smokers** (Figures 3 & 4).

ADVANTAGES of EpiSmoker

- Considers three smoking statuses
- Implementation of quantile normalization on test datasets using quantiles from discovery dataset has allowed transferring prediction model from one study to another. Thereby allowing and improving cross-study performance.
- Serves as a robust predictor of biological smoking status based on methylation data.
- Globally applicable to all populations.
- Available on GitHub: <https://github.com/sailalithabollepalli/EpiSmoker>
- Output provided as Smoking status labels, in HTML and CSV file formats.

1

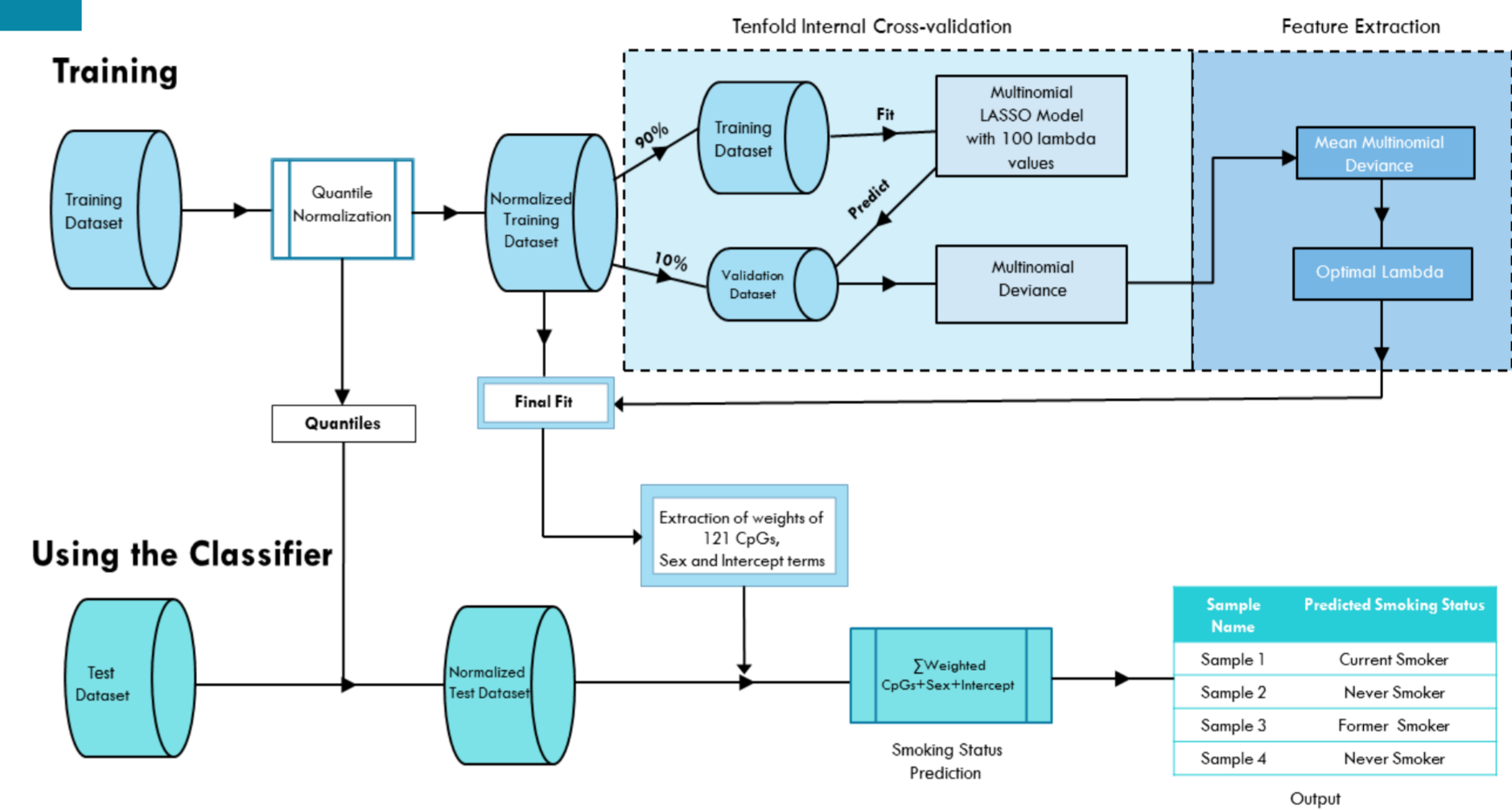
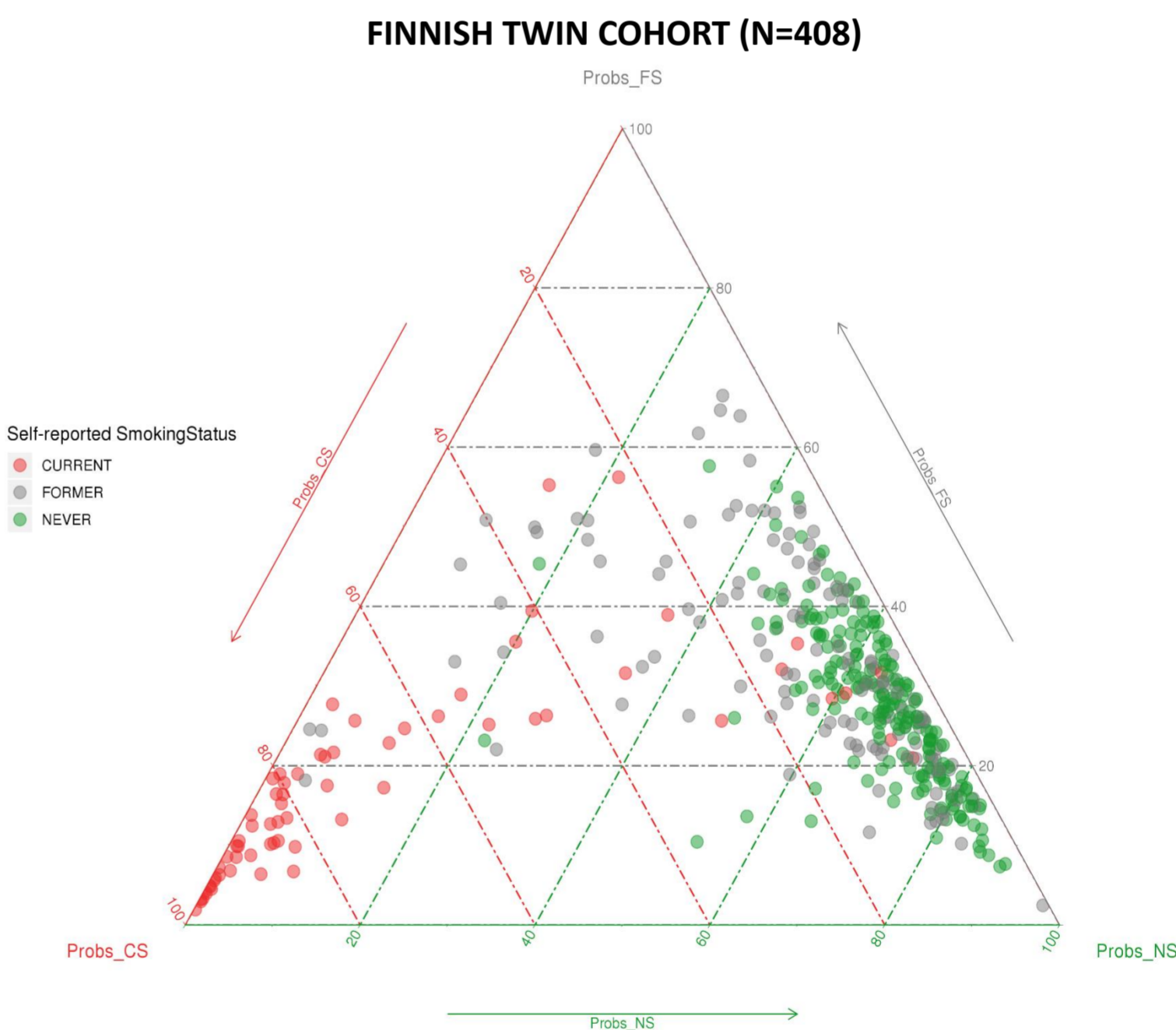


Figure 1: Schematic Representation of the Workflow of the EpiSmoker

2



Ternary plot based on the predicted probabilities given by our classifier in the FTC dataset. The ternary plot is an equilateral triangle with each vertex corresponding to a smoking status. Here color of the points and grid lines corresponds to self-reported smoking status. Each point in the ternary plot represents an individual with the corresponding triplet of probabilities that adds up to 100%. So higher the probability for a class closer is the point to the corresponding corner

3

FINNISH TWIN COHORT (N=408)				
PREDICTED	NEVER	9	101	192
	FORMER	3	31	7
	CURRENT	55	9	1
		CURRENT	FORMER	NEVER
		ACTUAL		

Figure 3: Confusion matrix comparing self-reported versus predicted smoking statuses

Self-reported smoking status is compared with the predicted smoking status determined by the classifier. Numbers along the anti-diagonal (7) indicate correctly predicted cases.

4

