

Supporting Information

Kim and Kim 10.1073/pnas.1318383110

SI Text

For classifying individuals into one of nine phenotypes denoted by $T = \{B, C, H, K, O, L, R, U, N\}$, where each phenotype is labeled by the first letter of the full initials of each trait (class) (Table 2). We use Bayesian inference of prediction results of four methods that are composed of two classification algorithms of the support vector machine (SVM) and k -nearest neighbor (k NN) analysis applied to two different descriptors of SNP and SNP syntax (SNP-S). The method names are abbreviated to k NN/SNP-S, k NN/SNP, SVM/SNP-S, and SVM/SNP, and they are mathematically denoted by m^1 , m^2 , m^3 , and m^4 , respectively. Each method requires training of its own parameters, which attempts to identify the best-performing parameter compositions. Once all methods are optimally fitted on the dataset, for each test individual i , we select a trait of having highest posterior probability conditioned on prediction results from trained methods, which can be formulated as $P(s_i|M_i^1, M_i^2, M_i^3, M_i^4)$, where s_i denotes the predicted trait of individual i , and M_i^j denotes the trait of individual i predicted by method m^j . By Bayes theorem, thus, we write

$$t_{max} = \operatorname{argmax}_{t \in T} P(s_i = t | M_i^1, M_i^2, M_i^3, M_i^4) \\ = \operatorname{argmax}_{t \in T} \frac{P(M_i^1, M_i^2, M_i^3, M_i^4 | s_i = t) \times P(s_i = t)}{P(M_i^1, M_i^2, M_i^3, M_i^4)},$$

where the denominator $P(M_i^1, M_i^2, M_i^3, M_i^4)$ is a normalizing constant. Because the predictive decisions of each method are inherently independent from each other and applying the chain rule (31)

$$= \operatorname{argmax}_{t \in T} \prod_{j=1}^4 P(M_i^j | s_i = t) \times P(s_i = t),$$

where $P(M_i^j | s_i = t)$ and $P(s_i = t)$ can be empirically inferred from the observations during the training phase of each of four methods by maximum likelihood estimation. For example, $P(M_i^1 = C | s_i = B)$ can be estimated by identifying a fraction of true breast invasive carcinoma (BRCA) individuals who were predicted to belong to colon adenocarcinoma (COAD) class by k NN/SNP-S method among entire BRCA samples in the training set. For $P(s_i = t)$, it corresponds to a fraction of samples of trait t of all training individuals, which is identical for each of nine traits, because the same sample size was used for each trait.

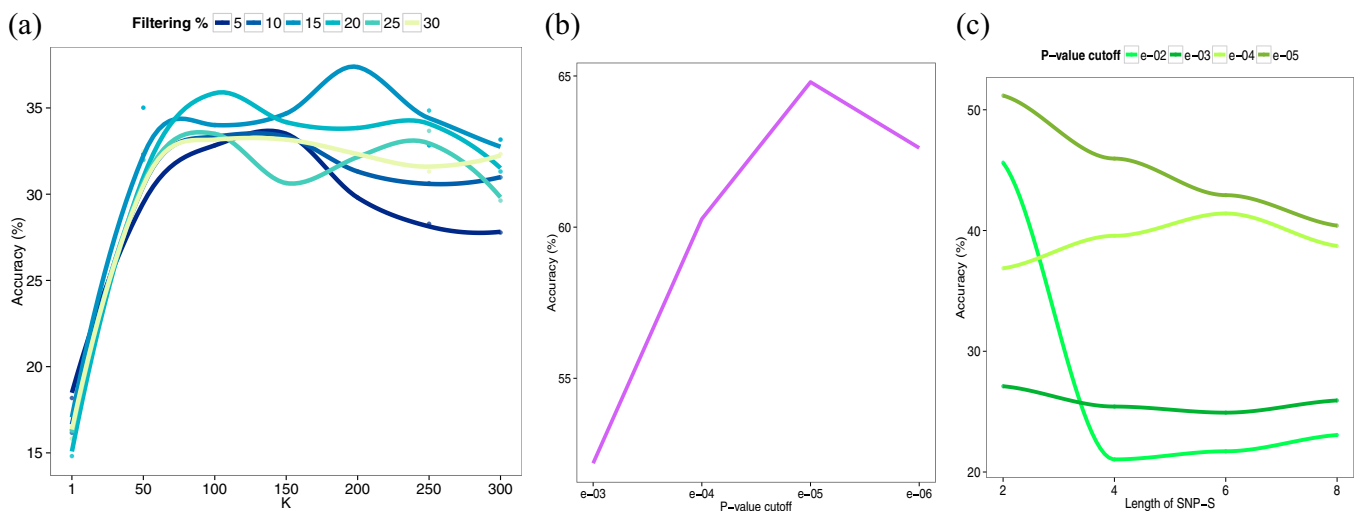


Fig. S1. (A) Optimization of parameters for the process of applying the k NN algorithm to the profiles of SNPs. The k NN/SNP method has two parameters: (i) filtering percentage for selecting rare features below specified frequency threshold, (e.g., for 1% filtering, the features below 1% frequency are selected for analysis) and (ii) k for selecting number of nearest neighbors of a test individual. (B) Optimization of parameters for the process of applying SVM algorithm to the profiles of SNPs. The SVM/SNP method has one parameter of the P value threshold for selecting the features with their P values below a specified value. (C) Optimization of parameters for the process of applying SVM algorithm to the profiles of SNP-Ss. The SVM/SNP-S method has two parameters: (i) P value threshold for filtering out features whose P values are greater than a specified value and (ii) the length of SNP-S.

Table S3. Training performance of kNN algorithm applied to profiles of SNPs

Actual trait	Predicted trait									Sample size	Accuracy (%)
	BRCA	COAD	HNSC	KIRC	LGG	OV	READ	UCEC	CEU		
BRCA	2	13	10	2	4	12	6	8	9	66	3.0%
COAD	0	27	12	0	6	11	4	3	3	66	40.9%
HNSC	0	15	36	1	2	1	3	2	6	66	54.5%
KIRC	1	20	9	7	4	6	8	4	7	66	10.6%
LGG	0	19	17	4	9	4	3	2	8	66	13.6%
OV	2	4	5	1	3	45	2	2	2	66	68.2%
READ	2	20	12	6	4	4	13	0	5	66	19.7%
UCEC	4	8	9	1	2	21	2	15	4	66	22.7%
CEU	0	0	1	1	0	0	0	0	66	66	100%
										Sum 594	Overall 37.0%

For the abbreviations, refer to Table S2 legend.

Table S4. Training performance of the SVM algorithm applied to profiles of SNPs

Actual trait	Predicted trait									Sample size	Accuracy (%)
	BRCA	COAD	HNSC	KIRC	LGG	OV	READ	UCEC	CEU		
BRCA	35	0	3	3	0	0	0	24	1	66	53.0%
COAD	8	33	1	12	2	0	4	5	1	66	50.0%
HNSC	1	0	52	9	1	0	0	1	2	66	78.8%
KIRC	5	1	4	45	1	0	0	9	1	66	68.2%
LGG	1	0	4	1	57	0	0	3	0	66	86.4%
OV	14	0	0	2	0	24	0	26	0	66	36.4%
READ	9	0	1	17	4	0	28	6	1	66	42.4%
UCEC	17	0	1	3	2	0	0	41	2	66	62.1%
CEU	0	0	0	0	0	0	0	0	66	66	100%
										Sum 594	Overall 64.1.0%

For the abbreviations, refer to Table S2 legend.

Table S5. Training performance of SVM algorithm applied to profiles of SNP-Ss

Actual trait	Predicted trait									Sample size	Accuracy (%)
	BRCA	COAD	HNSC	KIRC	LGG	OV	READ	UCEC	CEU		
BRCA	31	1	2	13	0	3	2	3	11	66	47.0%
COAD	5	11	1	28	1	2	3	1	14	66	16.7%
HNSC	1	2	39	7	2	0	7	0	8	66	59.1%
KIRC	5	10	1	31	2	3	3	4	7	66	47.0%
LGG	1	3	2	7	37	0	4	1	11	66	56.1%
OV	13	2	0	8	1	31	0	4	7	66	47.0%
READ	7	4	1	17	1	2	23	1	10	66	34.8%
UCEC	6	5	1	10	3	1	0	34	6	66	51.5%
CEU	0	0	0	0	0	0	0	0	66	66	100%
										Sum 594	Overall 51.1.0%

For the abbreviations, refer to Table S2 legend.