

Supplementary Information

A telescope GWAS analysis strategy, based on SNPs-Genes-Pathways ensemble and on multivariate algorithms, to characterize Late Onset Alzheimer's Disease

Margherita Squillario^{1,**}, Giulia Abate², Federico Tomasi¹, Veronica Tozzo¹, Annalisa Barla¹ and Daniela Uberti² “for the Alzheimer's Disease Neuroimaging Initiative”

* **Correspondence:** margherita.squillario@unige.it

Method Details

Data matrix transformation

Machine learning methods are state-of-the-art methods for a number of high-throughput data but they are still not very popular in the field of SNP analysis. Two are the main reasons: first, SNP data are usually very high-dimensional making some algorithms difficult to use on a regular workstation and second SNPs are usually represented by categorical variables (0, 1 or 2), a data type which may lead to inconsistent results with some machine learning methods.

Therefore for the SNP and pathway analyses we decided to use a different representation of the data that we defined “reinforced” and that uses weights calculated with a function of the Sequence Kernel Association Test (SKAT) method (see Supplementary Information), employed in the gene analysis performed in this work. Specifically we divided the samples of each datamatrix in two equal halves: one of these halves was used to calculate a weight for each SNP and the other half was used to address the actual classification task. Before this latter analysis we modified the data matrix multiplying the usual additive *code* associated to each SNP (*i.e.*, 0,1, or 2) to its respective weight, defined *beta*, as shown in this formula:

$$\text{SNP}_{(i)} = \text{SNP}_{\text{code}(i)} * \text{beta}_{(i)}$$

The few missing values present in the data matrix were replaced with the weight value associated to that SNP. Using this procedure we obtained datamatrices of continuous values, more suitable to be analyzed by many machine learning methods.

The SKAT function used to calculate these SNPs weights is called “Get Logistic Weights” and its formula is:

$$\text{weights} = \frac{e^{(\text{par1} - \text{MAF})\text{par2}}}{1 + e^{(\text{par1} - \text{MAF})\text{par2}}}$$

The SNPs weights are obtained considering the minor allele frequencies (MAF) that refer to the frequency at which the second most common allele occurs in a given population. It is an index used in population genetics to distinguish between common and rare variants. *Par1* and *par2* are two parameters, for the common and rare variants respectively, that a user can set up in order to give more or less relevance to rare and common variants. In this study we used the values suggest in ¹ in order to give a value different from 0 but very low (*i.e.*, *par1* = 1) to common variants and an high value to rare variants (*i.e.*, *par1* = 25), which are thought to have a major role in complex diseases like AD.

SNP analysis

When dealing with high-dimensional data, a natural problem arises in relation to the low number of available samples *n* with respect to the dimensionality of the problem that concerns the number of features *p* (*i.e.*, $n \ll p$). This issue happens frequently when dealing with biological data, and it usually referred to as “curse of dimensionality”. In this setting, usual statistical guarantees are lost, since the problem is over-determined, *i.e.*, there is no unique solution to the problem. A way to overcome this difficulty is to incorporate prior knowledge into the problem at hand, for example by employing sparse techniques, such as Elastic-Net or l_1l_2 feature selection ($l_1l_2\text{FS}$)².

$l_1l_2\text{FS}$ validate the robustness of the method using a model assessment framework, for which the model we selected as the “best” on our data is trained on a portion of the data (*learning set*) and tested on other unseen data (*test set*), iteratively (Figure S1A). Following this procedure it is possible to obtain a performance score for each selected model, through which we ensure the generalization properties our model achieves on unseen data.

In the present work we chose $l_1l_2\text{FS}$ within PALLADIO³ a machine learning python framework that can be customized to consider various combinations of feature selections and classification methods. Independently of the chosen methods, this tool ensures the reliability of the results performing two sets of experiments, *regular* and *permutation* batches (Figure S1B).

For each experiment, we resamples different learning and test sets a large number of times, in order to estimate the performance score distribution for both batches. The regular batch performs experiments on the given dataset, while the permutation batch performs experiments where the relationship between input and output is destroyed by shuffling the labels in the learning set, following what is referred to as *permutation test*. The two distributions are then compared by testing the null hypothesis H_0 by means of a Two-sample Kolmogorov–Smirnov test ⁴, a principled way to measure the statistical robustness of the obtained result. Then, we can reject the null hypothesis when the computed p-value is smaller than the confidence interval. Rejecting H_0 implies a clear difference between the two distributions and the sample size is large enough to describe the relationship between data and labels. The final outcome of a classification process is the prediction of the labels associated with a set of input samples. In order to assess the performance of l_1l_2FS , PALLADIO computes, among other performance metrics, the *balanced accuracy score* and the *Matthews correlation coefficient* (MCC).

The balanced accuracy score is the ratio of correctly predicted labels, adjusted for unbalanced problems. In particular, a predictor that always returns the label of the most represented class would yield a score of 50%, independently of the proportion of the labels in the dataset. MCC, also, is a comprehensive measure shown to be particularly useful for unbalanced problems, since it is always defined to be +1 for perfect match between predicted labels and ground truth, -1 for total disagreement and 0 for random prediction.

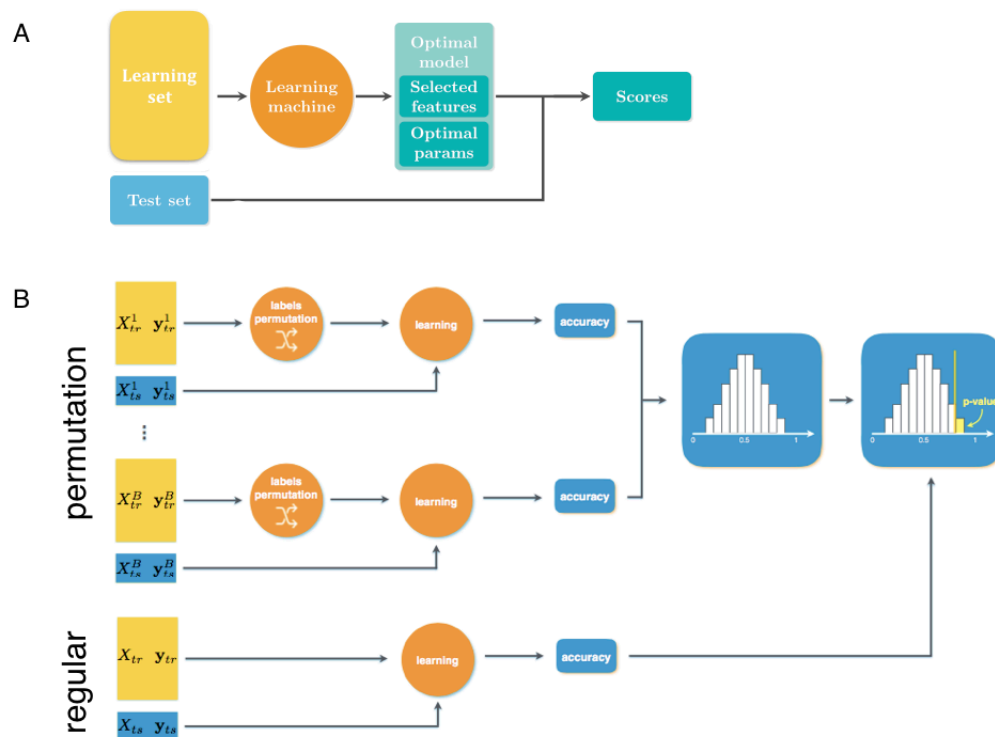


Figure S1. l_1l_2 and PALLADIO schemes.

(A) It shows the model assessment framework of l_1l_2FS .

(B) It shows the approach that PALLADIO adopts to ensure the reliability of the results.

Gene analysis

SKAT ¹ uses a multiple regression model to directly regress the phenotype on genetic variants in a region and on covariates allowing different variants to have different directions (i.e., protective or causal) and magnitude of effects, including no effects. To collapse the information of the

variants in a region, SKAT does not need a threshold because it uses a variance-component score test that is a kernel association test. In the formula of this test there are those weights that a SKAT user can chose to improve the power of the analysis (see “Data matrix transformation” section in Supplementary Information).

In the present work we chose to utilize the *SKATBinary* package, more suitable with PLINK formatted files, that encloses a function that computes p-values for Burden, SKAT and SKAT-O test for binary traits (in our case cases@controls and APOEε4 tasks) using asymptotic and resampling methods.

Burden test is suitable in case when a large proportion of variants are causals and their effect are in the same direction (*i.e.*, all protective or causal). SKAT test is suitable when only a small proportion of variants are causal or their effects have mixed directions (*i.e.*, some protective and some causal). SKAT-O test is suitable whenever we have a genetic scenario that is a mixture between that one suitable for Burden and that one suitable for SKAT.

For ADNI-1 and ADNI-2 we applied the following conservative thresholds: $0.05/36,000 = 1.37 \times 10^{-6}$ and $0.05/29,484 = 1.70 \times 10^{-6}$, where 0.05 is the level of significance, and the denominator indicates the known genes and intergenic regions in which the platforms have been subdivided.

In silico SNP characterization

The functional characterization of the gene lists derived from the SNPs signatures identified with PALLADIO was performed through enrichment analysis using the online toolkit WebGestalt ⁵. This tool takes as input a list of relevant genes/probesets and performs an enrichment analysis based on a hypergeometric test, providing several methods to correct for multiple hypothesis and using several databases (*e.g.*, the Kyoto Encyclopedia of Genes and Genomes (KEGG), Gene Ontology (GO)) for identifying the most relevant pathways and ontologies in each signature. In other words, given a KEGG pathway and a reference set (such as the entire human genome or the list of genes in a microarray platform), the enrichment is based on the comparison between the fraction of lists genes in the pathway and the fraction of KEGG pathway genes in the reference set. The gene list is enriched in that specific KEGG pathway if the former is larger than the latter fraction.

In the present work we enriched the two longest gene lists deriving from the SNP signatures of ADNI-1 (APOEε4 task), and of ADNI-2 (cases@controls task) considering KEGG ⁶ a database of pathway. In this analysis we considered the human genome as reference, 0.05 as level of significance, Benjamini-Hochberg as test for multiple hypothesis correction and 3 as the minimum number of genes in a KEGG pathway.

In order to understand which genes of the identified SNP signatures are already known to be involved in AD and which genes are not, we utilized Phenopedia ⁷, two web-based applications that explore the literature in a gene-centric and disease-centric way. We obtained the list of genes associated to AD from Phenopedia and we compared it to our lists of genes derived from our SNP signatures. The genes highlighted in red in Tables S2 and S3 are those genes known to be involved in AD.

Supplementary tables

	ADNI-1	ADNI-2
AD cases	179	126
Controls	214	155
1/2 APOEε4 carriers	178	120
0 APOEε4 carriers	215	161

Table S1. Sample size of the classes compared in the two classification tasks addressed in ADNI-1 and ADNI-2 datasets. AD cases vs. controls define the cases@controls task while 1/2 APOEε4 carriers (high risk) vs. 0 APOEε4 carriers (low risk) define the APOEε4 task.

	ADNI-1		
	SNP	# Chr	Gene Symbol
Cases@Controls	rs12205042	6	<i>HIVEP2</i>
	rs1940890	6	<i>LOC101928911 SPACA1</i>
	rs16881241	6	<i>RNGTT LOC100131124</i>
	rs6914160	6	<i>KLHL31 LRRC1</i>
	rs9465982	6	<i>CDKAL1</i>
	rs543049	6	<i>LOC101928911 SPACA1</i>
	rs6923298	6	<i>HIVEP2</i>
	rs677120	6	<i>LOC101928911 SPACA1</i>
	rs4896228	6	<i>IL20RA</i>
	rs6053572	20	<i>GPCPD1 SHLD1</i>
	rs708925	20	<i>PLCB1</i>
	rs2983626	20	<i>CST9L CST9</i>
	rs2247337	20	<i>CST9L</i>
	rs13040567	20	<i>CST9L</i>
APOEε4	rs11260977	1	<i>IGSF21</i>
	rs4839223	1	<i>SLC6A17</i>
	rs4839225	1	<i>SLC6A17</i>
	rs10749753	1	<i>LEPROT LEPR</i>
	rs1338138	1	<i>LINC01781 MTND2P30</i>
	rs11102933	1	<i>NGF LOC112268234</i>
	rs2147085	1	<i>LINC01781 MTND2P30</i>
	rs10789215	1	<i>SGIP1</i>
	rs12487324	3	<i>IQSEC1 NUP210</i>
	rs8180086	3	<i>GPR87 P2RY13</i>
	rs7651843	3	<i>RBMS3</i>
	rs7641352	3	<i>TBL1XR1 KCNMB2</i>
	rs276117	3	<i>LOC728290 GBE1</i>
	rs9875152	3	<i>P2RY12 MED12L</i>
	rs7625229	3	<i>BFSP2 LOC391578</i>
rs661798	3	<i>KALRN</i>	

rs12491760	3	<i>PTPRG</i>
rs9311976	3	<i>GRM7</i>
rs9310917	3	<i>RBMS3</i>
rs6784803	3	<i>MECOM</i>
rs7653603	3	<i>P2RY14</i>
rs276125	3	<i>LOC728290 GBE1</i>
rs13283389	9	<i>GRIN3A CYLC2</i>
rs11139921	9	<i>RASEF FRMD3</i>
rs10820215	9	<i>GRIN3A CYLC2</i>
rs7046513	9	<i>GRIN3A CYLC2</i>
rs10119403	9	<i>LINC01505</i>
rs1891999	9	<i>OLFM1 C9orf62</i>
rs7041138	9	<i>SLC25A6P5 LINC01505</i>
rs1350996	9	<i>FLJ35282 LOC101929563</i>
rs2075650	19	<i>TOMM40</i>
rs8106922	19	<i>TOMM40</i>
rs439401	19	<i>LOC100129500 APOC1</i>
rs236137	20	<i>SHDL1 CHGB</i>
rs6041265	20	<i>BTBD3 PA2G4P2</i>
rs1287032	20	<i>SHDL1 CHGB</i>
rs6041271	20	<i>BTBD3 PA2G4P2</i>
rs2294575	20	<i>CHD6</i>
rs2057291	20	<i>GNAS</i>

Table S2. ADNI-1 SNP signatures identified in the cases@control and in the APOEε4 tasks. “# Chr” indicates the chromosome number and “|” in the Gene Symbol column indicates that a specific SNP is located in a intergenic region between two genes. In red color are highlighted those genes that are already known to be associated to AD. In bold black color are highlighted those genes/SNPs mentioned in the main manuscript.

ADNI-2			
	SNP	# Chr	Gene Symbol
Cases@Control s	rs640688	1	<i>HIVEP3</i>
	rs2093933	1	<i>KIAA1324</i>
	rs3913318	1	<i>MIR4471 LOC100287877</i>
	rs10493973	1	<i>OLFM3</i>

rs2843130	1	<i>MORN1</i>
rs724309	1	<i>ATPAF1</i>
rs6739882	2	<i>MIR4431 / ASB3</i>
rs13014133	2	<i>RBMS1 / LOC100131736</i>
rs7562244	2	<i>AGAP1</i>
rs2160782	2	<i>LINC01800 / LINC02245</i>
rs13389584	2	<i>GPR39</i>
rs266410	3	<i>MRPS35P1 / GRM7-AS3</i>
rs341981	3	<i>EDEM1</i>
rs1598915	3	<i>EPHA6</i>
rs9683798	4	<i>ZEB2P1 / LDB2</i>
rs7660498	4	<i>HAUS3 / MXD4</i>
rs224489	4	<i>HS3ST1</i> / <i>LOC101929019</i>
rs4689726	4	<i>SORCS2</i>
rs12507259	4	<i>STK32B</i>
rs2203758	4	<i>LCORL / RPL21P46</i>
rs4648016	4	<i>NFKB1</i>
rs10043779	5	<i>KIAA0825</i>
rs1494699	5	<i>MSNP1 / LOC100131678</i>
rs6882967	5	<i>MSNP1 / LOC100131678</i>
rs25754	5	<i>ADAMTS12</i>
rs17156151	5	<i>NUDT12 / RAB9BP1</i>
rs261747	5	<i>FYB</i>
rs6892938	5	<i>ARL15</i>
rs2028269	5	<i>MTX3 / LOC100500934</i>
rs6865330	5	<i>FGF10-AS1 / LOC100506674</i>
rs9647537	5	<i>PGBD3P2 / HPRTP2</i>
rs17136076	7	<i>RNA5SP230 / MYL7</i>
rs17166226	7	<i>SCIN</i>
rs10121110	9	<i>ENG</i>
rs9792690	9	<i>TRPM3</i>
rs10976614	9	<i>C9orf123 / PTPRD</i>
rs4740366	9	<i>ABL1</i>
rs9408761	9	<i>PTPRD</i>
rs7854386	9	<i>LOC401557 / C9orf62</i>
rs36100013	9	<i>JAK2</i>
rs10817547	9	<i>ZNF618</i>

rs10819687	9	NAMA / LOC101928438
rs7031871	9	ARL2BPP7 / LOC100127962
s11244450	10	CHST15 / OAT
rs4980929	12	IQSEC3
rs11609462	12	ERC1
rs3759347	12	LEPREL2
rs3217933	12	CCND2
rs4766200	12	PARP11 / HSPA8P5
rs9552886	13	SGCG
rs7996072	13	CYSLTR2
rs17085790	13	LNK2
rs12867878	13	RNA5SP30 / LOC101926897
rs11841581	13	TEX26 / WDR95P
rs10507296	13	MIPEPP3 LINC00539
rs9564566	13	SNRPF3 / SRSF1P1
rs1373904	13	LACC1 / DGKZP1
rs2389229	13	ABCC4
rs7999070	13	TPTE2P1
rs12861751	13	LINC00378 / MIR3169
rs1935179	13	RPL7L1P1 / PEX12P1
rs2407249	13	CYSLTR2 / PSME2P2
rs12894732	14	LOC100418768 LINC01800
rs10136784	14	LINC00639
rs7143462	14	ESRRB / CYCSP1
rs2748144	14	LOC101927598 / GNG2
rs3825604	14	GNG2
rs11156929	14	SLC25A21
rs8010556	14	C14orf132
rs6497287	15	HERC2
rs2672680	15	FAM189A1
rs870185	15	ZFAND6
rs1883005	15	SNORD115-21 SNORD115-15
rs12908255	15	PSTPIP1
rs11634439	15	ARHGAP11A
rs2010459	15	TMED3
rs4965785	15	LRRK1 / CHSY1
rs7167588	15	GABRG3

rs34261044	15	<i>HERC2</i>
rs2239307	16	<i>ADCY9</i>
rs8061043	16	<i>CLEC16A</i>
rs12598337	16	<i>GRIN2A / ATF7IP2</i>
rs11643000	16	<i>GRIN2A</i>
rs6497898	16	<i>HS3ST4</i>
rs36474	16	<i>MYLK3</i>
rs9933735	16	<i>RBFOX1</i>
rs1124018	16	<i>RBFOX1 / LOC100131080</i>
rs4390571	16	<i>RBFOX1 / LOC100131080</i>
rs9940785	16	<i>RBFOX1</i>
rs2075158	16	<i>RSL1D1</i>
rs3116150	16	<i>SLC5A2</i>
rs13330742	16	<i>WWOX</i>
rs7189472	16	<i>XPO6</i>
rs2079268	17	<i>ALOX15P1 / SLC13A5</i>
rs9891398	17	<i>NF1</i>
rs16950363	17	<i>CA10</i>
rs9900961	17	<i>RPL17P41 / BPTF</i>
rs6504840	17	<i>LOC100419014 / RPS2P48</i>
rs8066872	17	<i>LINC00673</i>
rs12947685	17	<i>COX11</i>
rs12938347	17	<i>LINC01483 / LINC01028</i>
rs2007530	17	<i>ARHGAP27P1</i>
rs740642	17	<i>NTN1</i>
rs7236390	18	<i>PIEZO2</i>
rs630285	18	<i>AQP4-AS1</i>
rs605961	18	<i>MPPE1</i>
rs678570	18	<i>LAMA1</i>
rs8091074	18	<i>LINC01387</i>
rs12984574	19	<i>ZNF627</i>
rs2288867	19	<i>ATP13A1</i>
rs4807347	19	<i>ZNF555</i>
rs17639568	19	<i>NFIX</i>
rs367209	19	<i>LOC101928063</i>
rs7252291	19	<i>CELF5 / NFIC</i>
rs760629	20	<i>PPIAP21 / EIF4EBP2P</i>

	rs13041524	20	<i>PLCB4</i>
	rs236114	20	<i>MCM8</i>
	rs6075924	20	<i>LOC284744 LINC00261</i>
	rs6082789	20	<i>LNCNEF KRT18P3</i>
	rs6014017	20	<i>PFDN4 DOK5</i>
	rs8119892	20	<i>PPIAP21 EIF4EBP2P</i>
	rs9679935	20	<i>VAPB</i>
	rs12152036	21	<i>MIR548XHG PPIAP22</i>
	rs13048883	21	<i>C1QBPP FDPSP6</i>
	rs4816257	21	<i>MRPL39</i>
	rs2257008	21	<i>MIR5009</i>
	rs7283527	21	<i>LOC101927869 LINC01692</i>
	rs12484854	22	<i>LOC102724653</i>
	rs9612352	22	<i>ZDHHC8P1 LINC01659</i>
	rs11703440	22	<i>LOC284898 LINC02554</i>
	rs2298372	22	<i>DRICH1</i>
	rs5764804	22	<i>FBLN1</i>
	rs5760912	22	<i>CRYBB2</i>
	rs5752839	22	<i>ZNRF3</i>
	rs7288379	22	<i>SHISAL1 LINC01656</i>
rs9614616	22	<i>NUP50 KIAA0930</i>	
rs11703546	22	<i>CPSF1P1 RFPL3</i>	
APOEε4	rs367209	19	<i>LOC101928063</i>
	rs383133	19	<i>ZNF221</i>
	rs365745	19	<i>ZNF221</i>
	rs415499	19	<i>ZNF155</i>

Table S3. ADNI-2 SNP signatures identified in the cases@control and in the APOEε4 tasks. “# Chr” indicates the chromosome number and “|” in the Gene Symbol column indicates that specific SNP is located in an intergenic region between two genes. In red color are highlighted those genes that are already known to be associated to AD. In bold black color are highlighted those genes/SNPs mentioned in the main manuscript.

GROUP	PATHWAYS
1a	Caspase activation via extrinsic apoptotic signaling pathway, intrinsic pathway for apoptosis, apoptosis execution phase, regulated necrosis, transmission across chemical synapse, amyloid fiber formation, deregulated CDK5 triggers multiple neurodegenerative pathways.
1b	Macroautophagy, cellular response to hypoxia, cellular response to heat stress, cellular senescence, detoxification of reactive oxygen species, potassium channels.
1c	Cellular senescence, detoxification of reactive oxygen species, PIP3 activates AKT signaling.
2	Metabolism of nitric oxide, mitochondrial protein import, mitochondrial iron-sulfur cluster biogenesis, the citric acid (TCA) cycle and respiratory electron transport, cellular senescence, detoxification of reactive oxygen species, mitochondrial translation, mitochondrial calcium ion transport.
3	Caspase activation via extrinsic apoptotic signaling pathway, intrinsic pathway for apoptosis, apoptosis execution phase, regulated necrosis, death receptor signaling
4	Clathrin-mediated endocytosis, translocation of GLUT4 to the plasma membrane, trans-golgi network vesicle budding, mitochondrial calcium ion transport, ABC-family proteins mediated transport, cellular hexose transport
5a	Amyloid fiber formation, unfolded protein response, regulation of insulin-like growth factor (IGF), mitochondrial protein import, chaperoning-mediated protein folding, post-chaperoning tubuling folding pathway, asparagin N-linked glycosilation, gamma carboxylation, carboxyterminal post-translation.
5b	Post-translation protein phosphorylation, neddylation, protein ubiquitination, deubiquitination, O-linked glycosilation, post-translational modification: synthesis of GPI-anchored proteins.
6a	Biological oxidation, the citric acid (TCA) cycle and respiratory electron transport, regulation of insulin secretion, glucagon signaling in metabolism regulation, metabolism of carbohydrates, digestion.
6b	Metabolism of nitric oxide, metabolism of lipids.
7	Cellular response to hypoxia, Cellular response to heat stress, detoxification of reactive oxygen species, cellular senescence, HSP90 chaperone cycle for steroid hormone receptors (SHR), cell junction organization, macroautophagy.
8	mTOR signaling, death receptor signaling, PIP3 activates AKT signaling, MAPK1/MAPK3 signaling, MAPK6/MAPK4 signaling, integrin signaling by leptin, integrin signaling by hippo, WNT ligand biogenesis and trafficking, degradation of beta-catenin by destruction complex, TCF dependent signaling in response to WNT, beta-catenin independent WNT signaling.
9a	GPCR ligand binding, GPCR downstream signaling, GASTRIN-CREB signaling, pre-NOTCH expression and processing, signaling by NOTCH1, signaling by NOTCH2, signaling by NOTCH3, signaling by NOTCH4.
9b	GPCR downstream signaling, GASTRIN-CREB signaling.
9c2	Signaling by TGF-beta family members.
9c3	Signaling by receptor tyrosine kinases.

Table S4. Groups of pathways selected in REACTOME and analyzed with Group Lasso with overlap in ADNI-1.

GROUP	PATHWAYS
1a	Caspase activation, intrinsic pathway for apoptosis, apoptosis execution phase, regulated necrosis, transmission across chemical synapse, amyloid fiber formation, deregulated CDK5 triggers multiple neurodegenerative pathways.
1b1	Cellular response to hypoxia, cellular response to heat stress, potassium channels.
1b2	Macroautophagy, detoxification of reactive oxygen species, cellular senescence.
1c	Cellular senescence, detoxification of reactive oxygen species, PIP3 activates AKT signaling.
2	Metabolism of nitric oxide, mitochondrial protein import, mitochondrial iron-sulfur cluster biogenesis, the citric acid (TCA) cycle and respiratory electron transport, cellular senescence, detoxification of reactive oxygen species, mitochondrial translation, mitochondrial calcium ion transport.
3	Caspase activation via extrinsic apoptotic signaling pathway, intrinsic pathway for apoptosis, apoptosis execution phase, regulated necrosis, death receptor signaling
4	Clathrin-mediated endocytosis, translocation of GLUT4 to the plasma membrane, trans-golgi network vesicle budding, mitochondrial calcium ion transport, ABC-family proteins mediated transport, cellular hexose transport
5a	Amyloid fiber formation, unfolded protein response, regulation of insulin-like growth factor (IGF), mitochondrial protein import, chaperoning-mediated protein folding, post-chaperoning tubuling folding pathway, asparagin N-linked glycosilation, gamma carboxylation, carboxyterminal post-translation.
5b	Post-translation protein phosphorylation, neddylation, protein ubiquitination, deubiquitination, O-linked glycosilation, post-translational modification: synthesis of GPI-anchored proteins.
6a	Biological oxidation, the citric acid (TCA) cycle and respiratory electron transport, regulation of insulin secretion, glucagon signaling in metabolism regulation, metabolism of carbohydrates, digestion.
6b	Metabolism of nitric oxide, metabolism of lipids.
7	Cellular response to hypoxia, Cellular response to heat stress, detoxification of reactive oxygen species, cellular senescence, HSP90 chaperone cycle for steroid hormone receptors (SHR), cell junction organization, macroautophagy.
8	mTOR signaling, death receptor signaling, PIP3 activates AKT signaling, MAPK1/MAPK3 signaling, MAPK6/MAPK4 signaling, integrin signaling by leptin, integrin signaling by hippo, WNT ligand biogenesis and trafficking, degradation of beta-catenin by destruction complex, TCF dependent signaling in response to WNT, beta-catenin independent WNT signaling.
9a	GPCR ligand binding, GPCR downstream signaling, GASTRIN-CREB signaling, pre-NOTCH expression and processing, signaling by NOTCH1, signaling by NOTCH2, signaling by NOTCH3, signaling by NOTCH4.
9b	GPCR downstream signaling, GASTRIN-CREB signaling.
9c	Signaling by TGF-beta family members, signaling by receptor tyrosine kinases.

Table S5 Groups of pathways selected in REACTOME and analyzed with Group Lasso with overlap in ADNI-2.

Supplemental Results

Validation of the SNP signatures

In order to verify the robustness of the identified SNPs signatures, a validation procedure was performed considering the dataset (ADNI-1 or ADNI-2) left available. Two steps characterize the validation: the first one consists in mapping the SNPs of the signature in another independent dataset, and the second one consists in analyzing the data matrix (having all the subjects and just the selection of SNPs of the identified list) evaluating the classification performance of the signature. When we started the validation procedure of the SNP signature identified in ADNI-1 cases@controls task, we tried to map all the 14 SNPs in ADNI-2 but just 9 SNPs were found. Despite this issue, we build the data matrix, having these 9 SNPs and 281 subjects, and we analyzed it using Regularized Least Square (RLS) classifier inside PALLADIO. Figure S2 shows that we could not validate this signature. A possible reason resides in the failure of the complete mapping of all the SNPs of our signature in ADNI-2.

In the validation procedure of the other ADNI-1 SNP signature, considering APOEε4 task, we encountered the same SNP mapping issue: in ADNI-2 we found just 24 SNPs over 39 total SNPs. Nonetheless this incomplete SNP mapping was sufficient to validate the signature in ADNI-2 (Figure S2).

When we tried to validate the SNP signatures identified in ADNI-2 in ADNI-1, we encountered the same SNP mapping issue explained before: for the cases@controls signature we found just 46 over 138 total SNPs and for the APOEε4 signature we found 2 over 4 total SNPs. Both these signatures did not pass the validation procedure (Figure S2).

Even if the validation procedure truly succeeded just for the APOEε4 signature identified in ADNI-1, we cannot state with certainty that the other three SNP signatures failed the validation because we were unable to map all the SNPs of these signatures in the validation dataset. A possible reason why we could validate just APOEε4 signature identified in ADNI-1, even if also in this case the SNP mapping was incomplete, could be found in the different “value” or “weight” of the SNPs of a signature. Since the SNPs are characterized by different weights, as confirmed the different risk of developing AD based on the number of copies of APOEε4, the probability that a SNP signature passes the validation will increase proportionally to the number and the weights of the mapped SNPs.

Furthermore the method we chose to perform the SNP-based analysis (*i.e.*, $l_{1/2FS}$) is designed to identify a list of features discriminant but also correlated. This last characteristic means that if the mapping procedure does not comprehends the most correlated SNPs, high is the probability that the validation of the signature will not succeed.

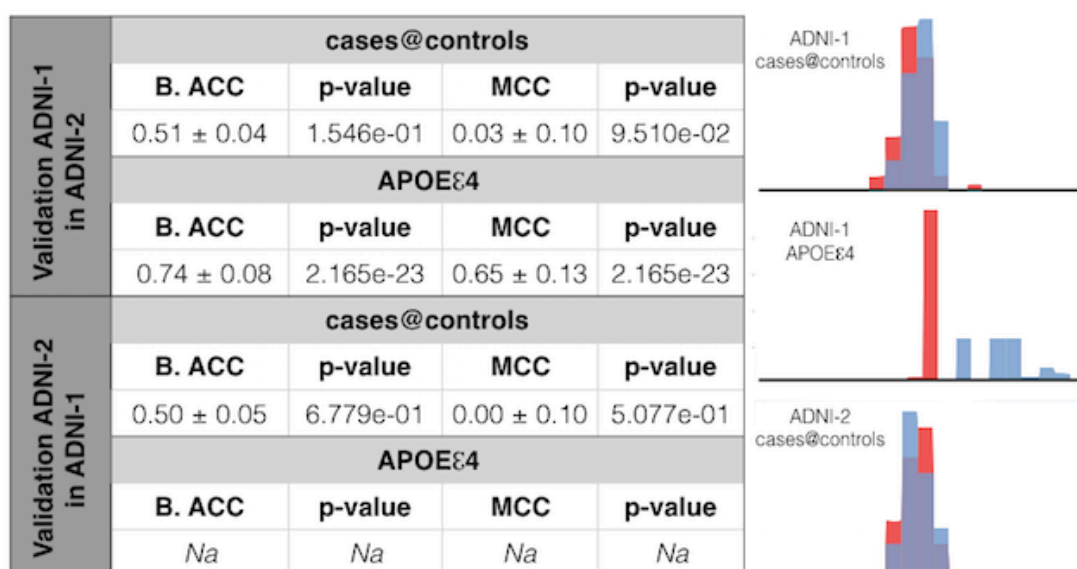


Figure S2. The validation results of the SNP signatures identifies in ADNI-1 and ADNI-2 dataset. In the validation procedure we consider the same two classification tasks: AD vs. healthy controls (cases@controls) and 1/2 APOEε4 vs. 0 APOEε4 carriers (APOEε4 task). B. ACC, Balanced Accuracy; MCC, Matthews Correlation Coefficient.

Supplementary References

1. Wu, M. C. *et al.* Rare-variant association testing for sequencing data with the sequence kernel association test. *Am. J. Hum. Genet.* **89**, 82–93 (2011).
2. Hastie, T., Tibshirani, R. & Wainwright, M. *Statistical learning with Sparsity: the lasso and generalizations.*
3. Barbieri, M., Fiorini, S., Tomasi, F. & Barla, A. PALLADIO: a parallel framework for robust variable selection in high-dimensional data. *Proc. 6th Work. Python High-Performance Sci. Comput.* 19–26 (2016). doi:10.1109/pyhpc.2016.13
4. Daniel, W. W. 'Kolmogorov-Smirnov one-sample test' in Applied Nonparametric statistics. in 635 (PWS-KENT Pub, 1990).
5. Wang, J., Duncan, D., Shi, Z. & Zhang, B. WEB-based GEne SeT AnaLysis Toolkit (WebGestalt): update 2013. *Nucleic Acids Res.* **41**, 77–83 (2013).
6. Kanehisa, M. & Goto, S. KEGG: Kyoto encyclopedia of genes and genomes. *Nucleic Acids Res.* **27**, 29–34 (1999).
7. Yu, W., Clyne, M., Khoury, M. J. & Gwinn, M. Phenopedia and genopedia: Disease-centered and gene-centered views of the evolving knowledge of human genetic associations. *Bioinformatics* **26**, 145–146 (2009).