

Supplementary material

Table S1. Specific relevant information associated to each TF property tagged in sentences of manual summaries

Property	Specific information manually tagged	Tag
ACT	Growth condition in negative regulation Growth condition in positive regulation Effector Regulation of the TF activity Active conformation of the TF	ACTCONDN ACTCONDP ACTEFFECTE ACTREG ACTCONF
EVO	Domain position and percentage of domain identity Percentage of TF identity with other TFs	EVPIDT EVPIT
SIT	Symmetry Size	SSM SSZ
TU	Regulation of the TU Organization of the TU Localization of the TU	TUR TURO TUL

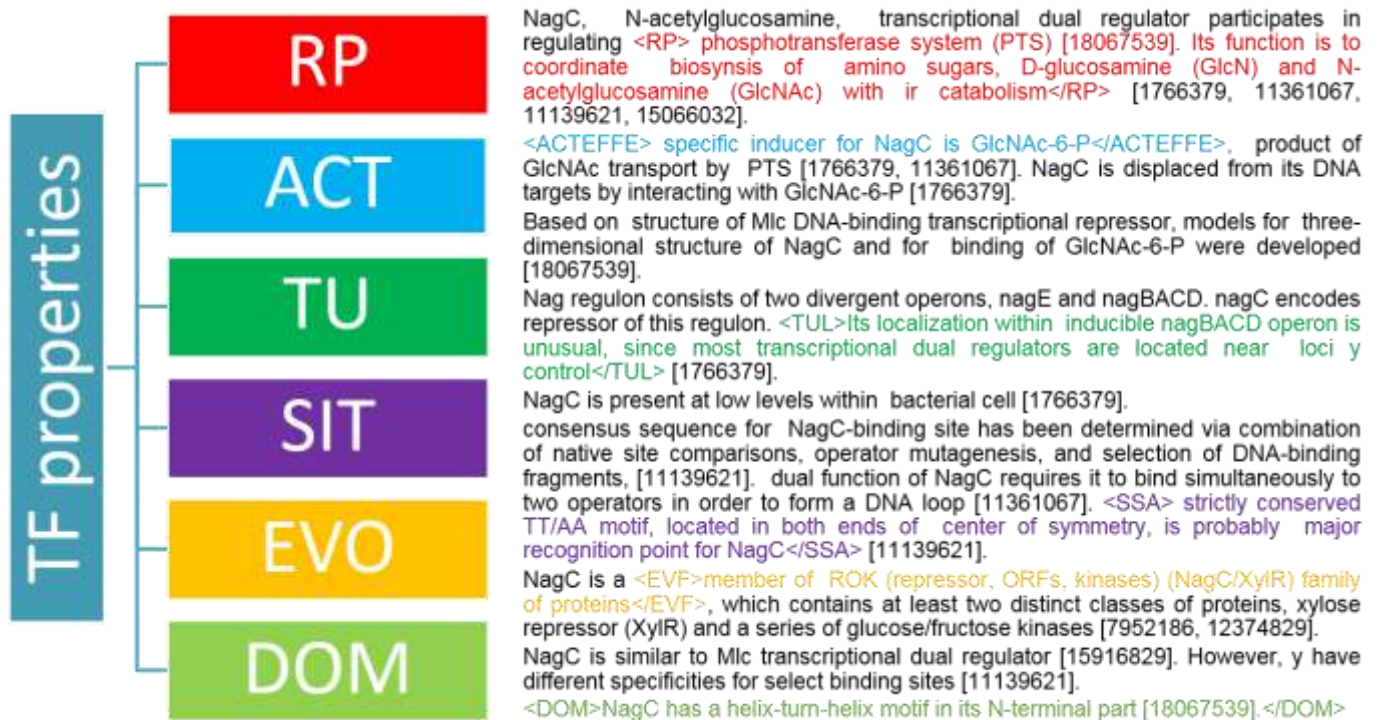


Figure S1. Example of tagged manual summary with specific relevant information.

Table S2. Biological information automatically tagged in manual summaries to enrich features for supervised learning. This information was also used to tag sentences from articles of *E. coli* and *Salmonella*. For *E. coli*, we indicate the source, and for *Salmonella* we indicate if the source was the same, the new source, or if the tag was not used

Property	Biological information	Source for <i>E. coli</i>	Source for <i>Salmonella</i>	Tag
ACT	Dictionary of TFs Dictionary of growth conditions	RegulonDB RegulonDB	NO SÉ The same	ACTTF ACTCOND

	Dictionary of effectors Keywords of effectors Dictionary of regulatory verbs Keywords of regulation Dictionary of conformations Keywords of conformations	RegulonDB Manually collected Manually collected Manually collected RegulonDB Manually collected	The same The same The same The same Not used The same	ACTEFFE ACTEFFE ACTREG ACTREG ACTCONF ACTCONF
DOM	Dictionary of structural domain families Dictionary of molecular functions Dictionary of structural motifs Dictionary of TFs Keywords of structural domains	DBD: Transcription Factor Prediction Database GO's OBO file Interpro RegulonDB Frequent words	The same The same The same NO SÉ The same	DFAM MF DMOT TF FWDOM
EVO	Dictionary of evolutionary families Keywords of percentage of TF identity with other TFs and percentage of domain identity Dictionary of structural domain position	RegulonDB Manually collected Manually collected	Not used The same The same	EVF EVPI EVDOM
RP	Dictionary of biological processes Keywords of regulated processes	GO Frequent words	The same The same	PRO FWRP
SIT	Keywords of symmetry Keywords of size Spatial arrangement	Manually collected Manually collected Regular expression	The same The same The same	SSM SSZ SSA
TU	Dictionary of genes Dictionary of transcription units Keywords of organization of the TU Keywords of regulation of the TU Keywords of localization of the TU	RegulonDB RegulonDB Manually collected Manually collected Manually collected	NO SÉ NO SÉ The same The same The same	TURO TURO TURO TUR TUL

Table S3. A general description of the experimental setup, including the tested values of the different aspects employed for training the six classifiers

Aspect	Values
Combination of features	lemma POS, lemma tag, tag for lemma
Vectorizer	Binary, TF-IDF, TF-IDF binary
N-grams	1, 1-2, 1-3
Dimensionality reduction with SVD	300, 200 dimensions
Feature selection with χ^2	1000, 800, 500 features
Under-sampling technique	RandomUS, Tomek, IHT, and OSS
SVM kernel	rbf, lineal, poly
Class weighting	true, false

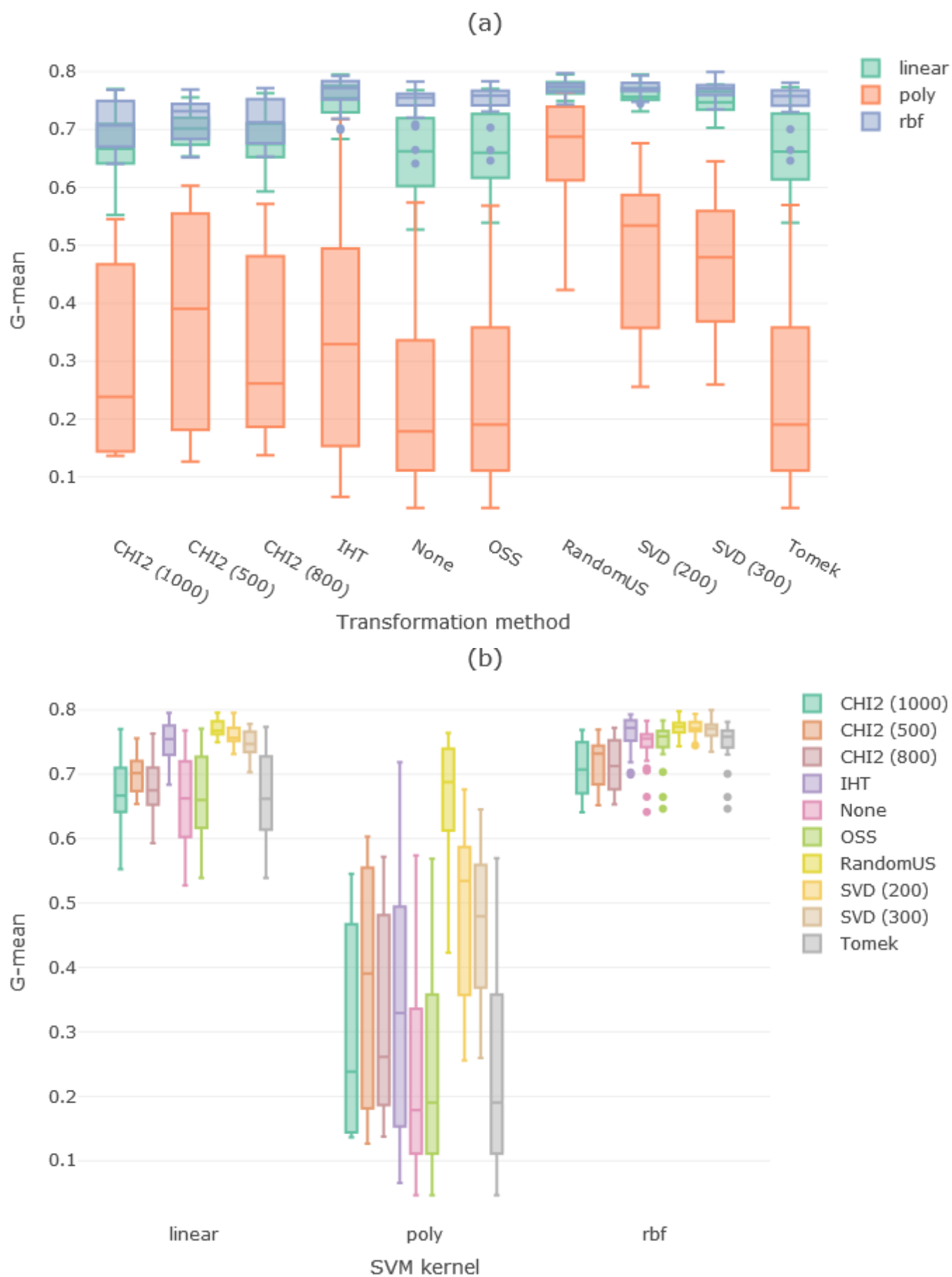


Figure S2. Distribution of performance in cross-validation of all trained predictive models for ACT property. (a) shows transformation methods in horizontal axis, CHI2 (χ^2) and SVD include number of final dimensions/features. (b) shows SVM kernel in horizontal axis.

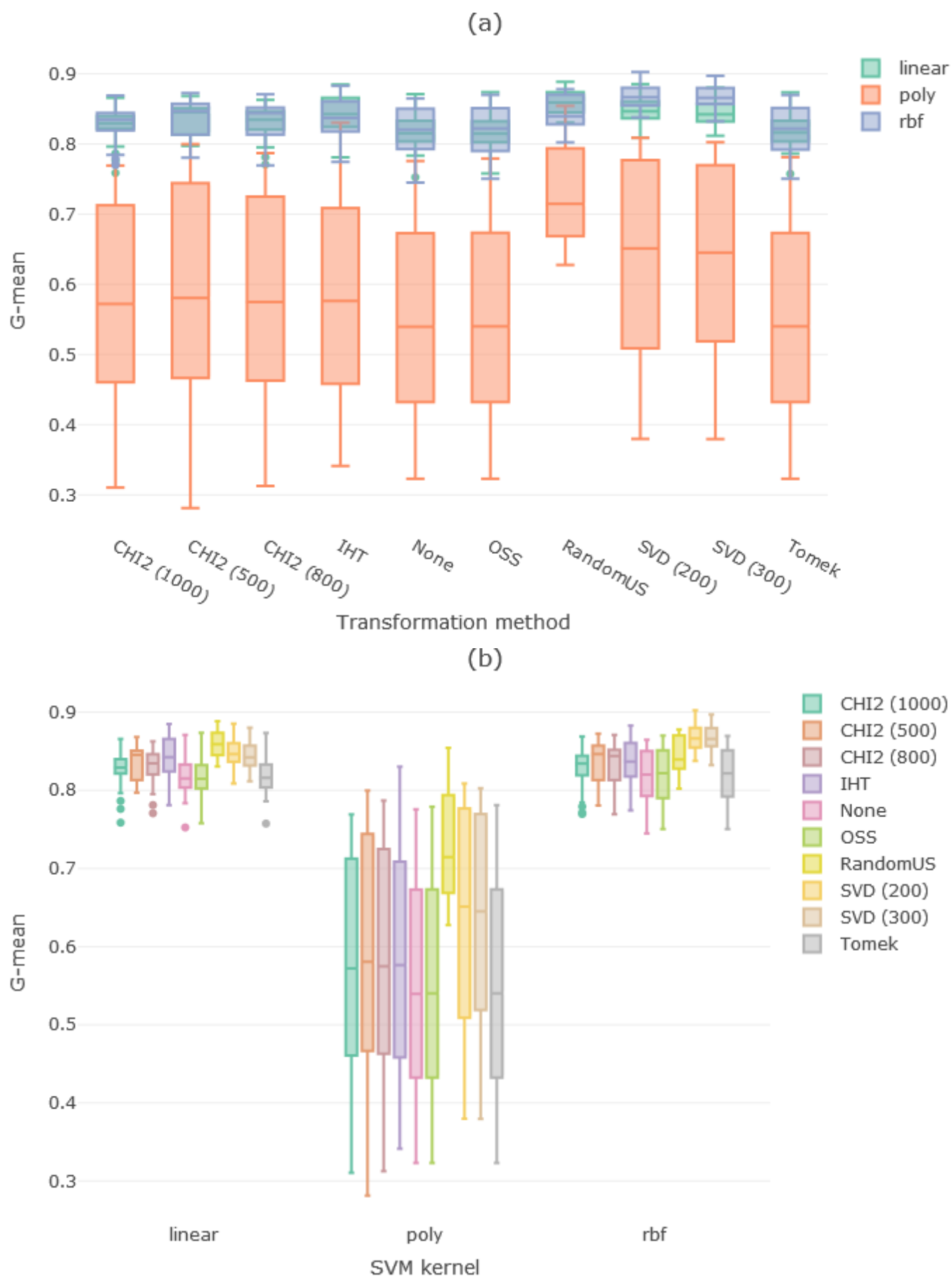


Figure S3. Distribution of performance in cross-validation of all trained predictive models for DOM property. (a) shows transformation methods in horizontal axis, CHI2 (χ^2) and SVD include number of final dimensions/features. (b) shows SVM kernel in horizontal axis.

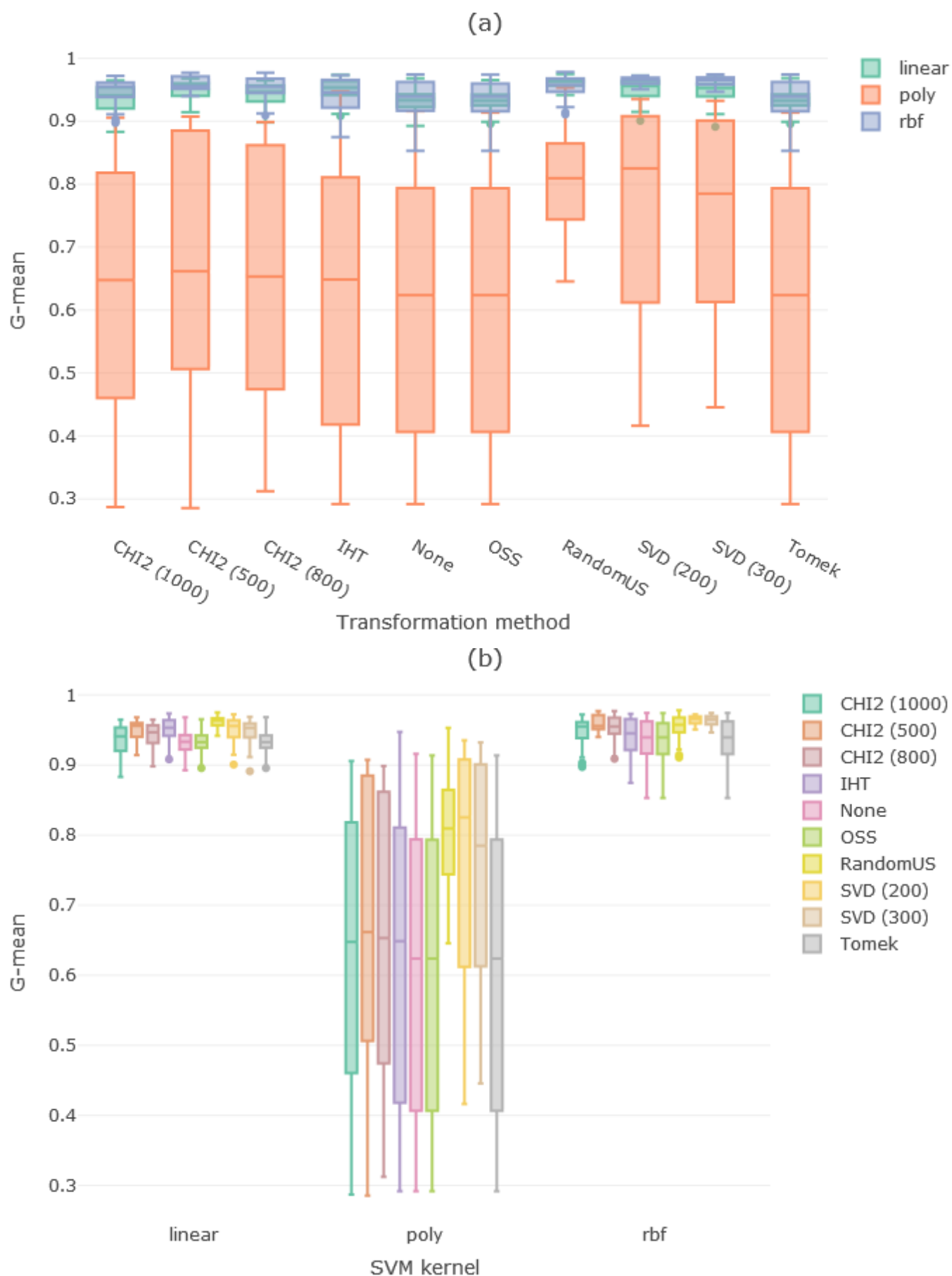


Figure S4. Distribution of performance in cross-validation of all trained predictive models for EVO property. (a) shows transformation methods in horizontal axis, CHI2 (χ^2) and SVD include number of final dimensions/features. (b) shows SVM kernel in horizontal axis.

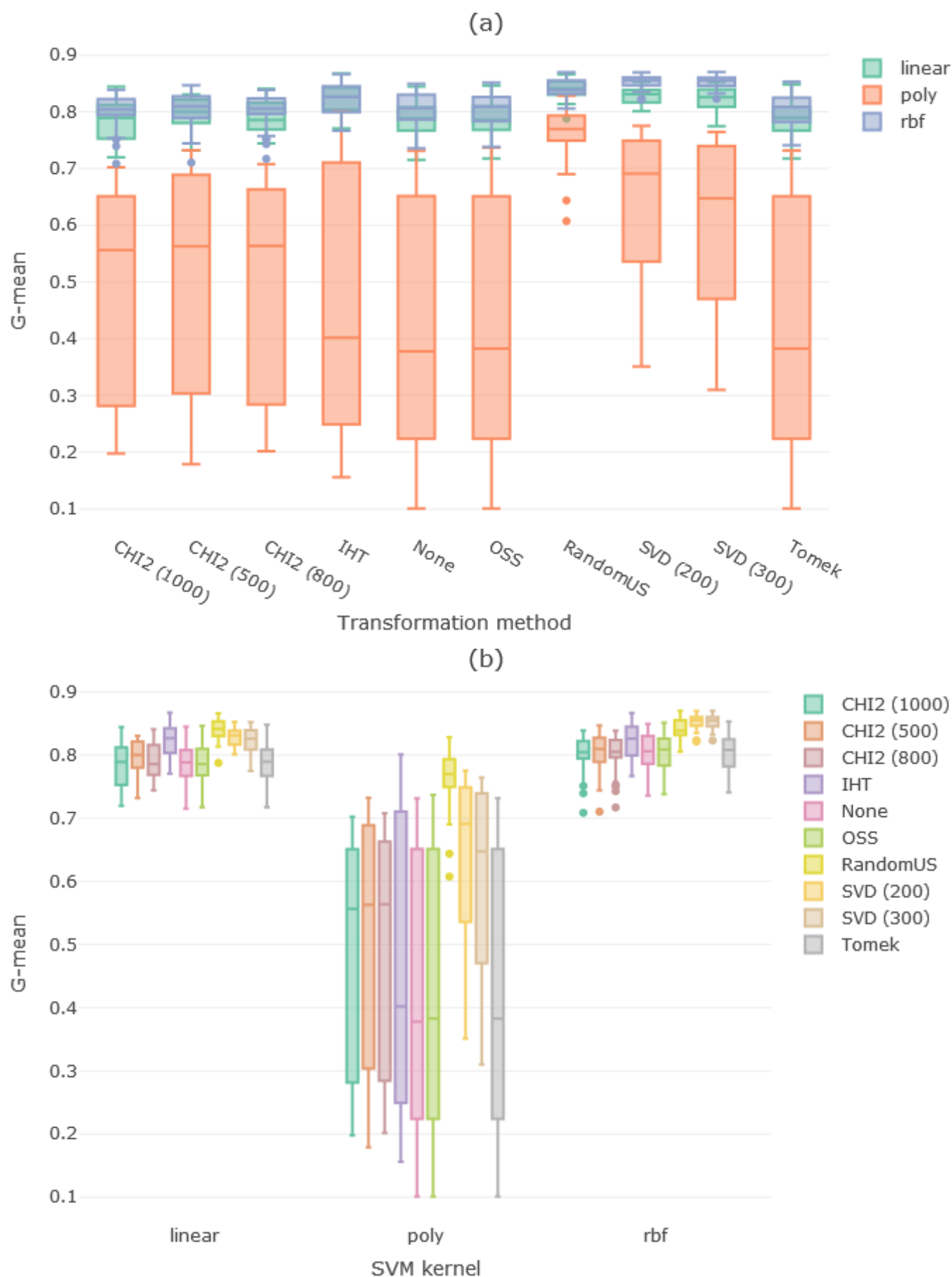


Figure S5. Distribution of performance in cross-validation of all trained predictive models for RP property. (a) shows transformation methods in horizontal axis, CHI2 (χ^2) and SVD include number of final dimensions/features. (b) shows SVM kernel in horizontal axis.

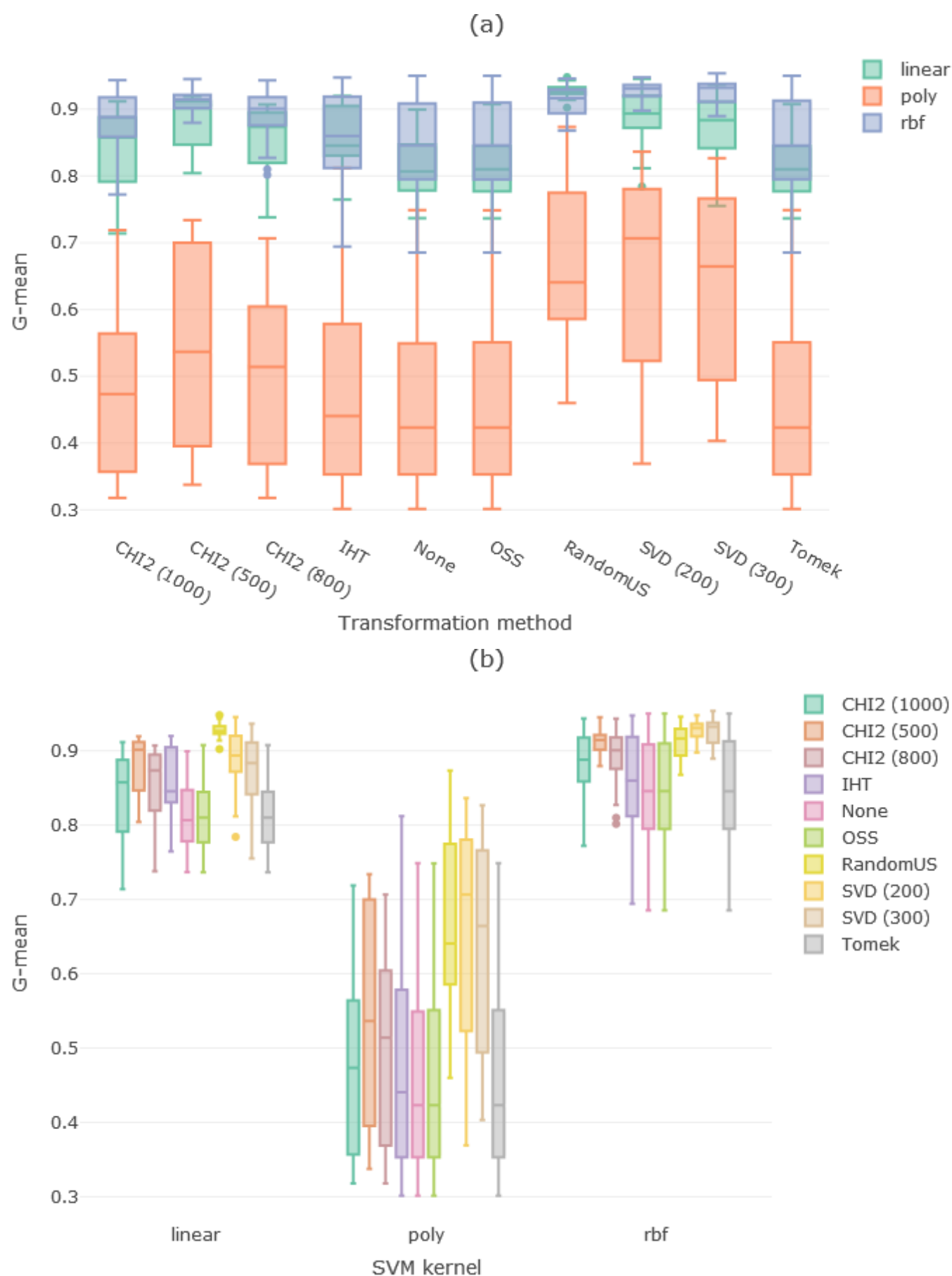


Figure S6. Distribution of performance in cross-validation of all trained predictive models for SIT property. (a) shows transformation methods in horizontal axis, CHI2 (χ^2) and SVD include number of final dimensions/features. (b) shows SVM kernel in horizontal axis.

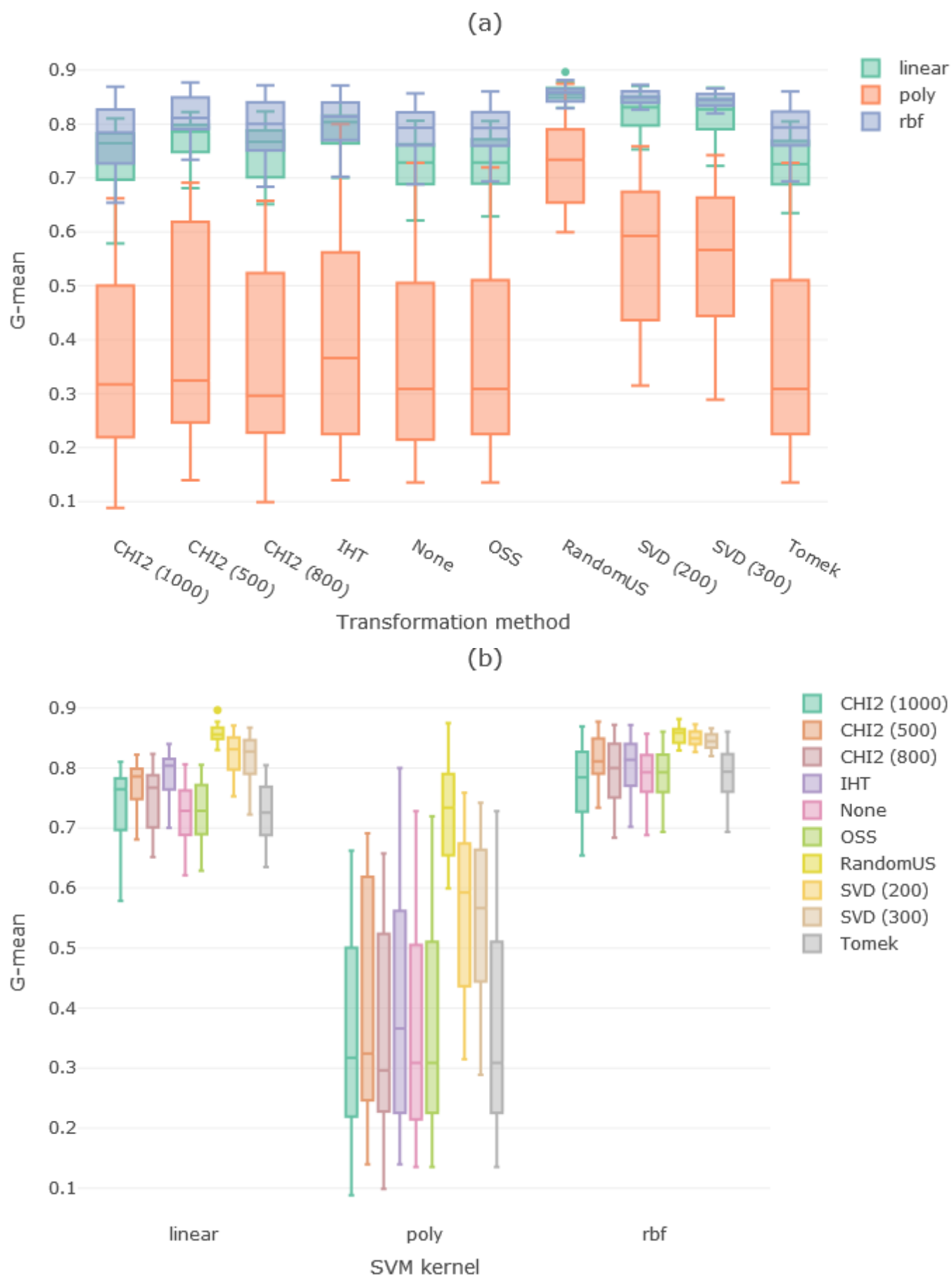


Figure S7. Distribution of performance in cross-validation of all trained predictive models for TU property. (a) shows transformation methods in horizontal axis, CHI2 (χ^2) and SVD include number of final dimensions/features. (b) shows SVM kernel in horizontal axis.

Table S4. Detailed characteristics of the best predictive model per TF property

Property	Vectorization	N-grams	Transformation	Final dimensions	SVM hyperparameters				G-mean score
					Kernel	C	Gamma	Class weight	
ACT	TF-IDF	1,2	RandomUS	--	RBF	3.0	1.0	Balanced	0.80
DOM	TF-IDF	1	SVD	200	RBF	1.0	1.0	Balanced	0.90
EVO	TF-IDF binary	1	SVD	200	RBF	0.5	1.0	Balanced	0.97
RP	TF-IDF binary	1,2	SVD	200	RBF	1.0	1.0	Balanced	0.87
SIT	TF-IDF binary	1	SVD	200	RBF	0.5	1.0	Balanced	0.95
TU	TF-IDF binary	1	SVD	200	RBF	3.0	0.1	Balanced	0.87