

Solution proposal

Goal of generative model

Utility:

- Accuracy on learning task
- Global statistical properties

Privacy:

- Membership Inference Attack
- Attribute Inference Attack
- Reconstruction



TabDDPM



Usefulness of TabDDPM

_	AB (R2)	AD (F1)	BU (F1)	CA (<i>R</i> 2)	$\operatorname{CAR}(F1)$	$\operatorname{CH}(F1)$	DE (F1)	DI (F1)
TVAE	$0.433 {\pm .008}$	$0.781 {\scriptstyle \pm .002}$	$0.864 {\scriptstyle \pm .005}$	$0.752 {\scriptstyle \pm .001}$	$0.717 {\pm .001}$	$0.732 {\pm .006}$	$0.656 {\scriptstyle \pm .007}$	$0.714 {\scriptstyle \pm .039}$
CTABGAN	-	$0.783 {\scriptstyle \pm .002}$	$0.855 \pm .005$	<u> </u>	$0.717 {\scriptstyle \pm .001}$	$0.688 {\pm}.006$	$0.644 {\scriptstyle \pm .011}$	$0.731 {\scriptstyle \pm .022}$
CTABGAN+	$0.467 {\scriptstyle \pm .004}$	$0.772 {\pm .003}$	$0.884 {\pm} .005$	$0.525 {\scriptstyle \pm .004}$	$0.733 {\pm .001}$	$0.702 {\scriptstyle \pm .012}$	$0.686 {\pm} .004$	$0.734 {\scriptstyle \pm .020}$
SMOTE	$0.549 {\scriptstyle \pm .005}$	$0.791 {\scriptstyle \pm .002}$	$0.891 {\pm} .003$	$0.840 {\scriptstyle \pm .001}$	$0.732 {\scriptstyle \pm .001}$	$0.743 {\scriptstyle \pm .005}$	$0.693 {\scriptstyle \pm .003}$	$0.683 {\scriptstyle \pm .037}$
TabDDPM	$0.550 {\scriptstyle \pm .010}$	$0.795 {\scriptstyle \pm.001}$	$0.906 {\pm .003}$	$0.836 {\scriptstyle \pm .002}$	$0.737 {\scriptstyle \pm .001}$	$0.755 {\scriptstyle \pm .006}$	$0.691 {\scriptstyle \pm .004}$	$0.740 {\scriptstyle \pm .020}$
Real	$0.556 {\pm} .004$	$0.815 {\scriptstyle \pm .002}$	$0.906 {\scriptstyle \pm .002}$	$0.857 {\scriptstyle \pm .001}$	$0.738 {\scriptstyle \pm .001}$	$0.740 {\pm} .009$	$0.688 {\pm} .003$	$0.785 {\scriptstyle \pm .013}$
	FB (R2)	GE (F1)	$\operatorname{HI}(F1)$	HO(R2)	IN (R2)	KI (R2)	MI(F1)	WI(F1)
TVAE	$0.685 {\scriptstyle \pm .003}$	$0.434 {\pm .006}$	$0.638 {\scriptstyle \pm .003}$	$0.493 {\pm} .006$	$0.784 {\pm .010}$	$0.824 {\pm} .003$	$0.912 {\scriptstyle \pm .001}$	$0.501 {\scriptstyle \pm .012}$
CTABGAN	_	$0.392 {\pm .006}$	$0.575 {\pm .004}$	-	-	—	$0.889 {\scriptstyle \pm .002}$	$0.906 {\scriptstyle \pm .019}$
CTABGAN+	$0.509 {\scriptstyle \pm .011}$	$0.406 {\scriptstyle \pm .009}$	$0.664 {\scriptstyle \pm .002}$	$0.504 {\pm} .005$	$0.797 {\scriptstyle \pm .005}$	$0.444 {\scriptstyle \pm.014}$	$0.892 {\scriptstyle \pm .002}$	$0.798 {\scriptstyle \pm .021}$
SMOTE	$0.803 {\scriptstyle \pm .002}$	$0.658 {\scriptstyle \pm .007}$	$0.722 {\scriptstyle \pm .001}$	$0.662 {\scriptstyle \pm .004}$	$0.812 {\scriptstyle \pm .002}$	$0.842 {\scriptstyle \pm .004}$	$0.932 {\scriptstyle \pm .001}$	$0.913 {\scriptstyle \pm .007}$
TabDDPM	$0.713 {\scriptstyle \pm .002}$	$0.597 {\pm} .006$	$0.722 {\scriptstyle \pm .001}$	$0.677 {\scriptstyle \pm .010}$	$0.809 {\scriptstyle \pm .002}$	$0.833 {\scriptstyle \pm.014}$	$0.936 {\scriptstyle \pm.001}$	$0.904 {\scriptstyle \pm .009}$
Real	$0.837 {\pm .001}$	$0.636 \pm .007$	$0.724 \pm .001$	$0.662 \pm .003$	$0.814 \pm .001$	$0.907 {\pm} .002$	$0.934 \pm .000$	$0.898 \pm .006$

Privacy



Q1 problem statement

- Training set: dataset that contains labeled real data used to train TabDPPM and synthetic data
- Test set: real data used to train TabDPPM and synthetic data.
- Task: Binary classification -> discriminate between synthetic data generated by TabDPPM and real data used to train TabDPPM.
- Intuition: If unable to discriminate between synthetic and real data then we expect TabDPPM to have a high machine learning efficiency.
- Can also be used to measure privacy as the distinguishability between real and synthetic data
- Goal: Interpret the classifier model to gain insights on potential methods that could lead to improvements in machine learning efficiency for TabDPPM.

Classification model performance

Classification Models	Raw Categorical Data	Encoded Categorical Data	Encoded Categorical Data + Scaled Numerical Data
Logistic Regression	N/A	N/A	57%
XGBoost	86%	85%	75%
Random Forest	N/A	75%	N/A

Feature importance





Difficult feature to recreate



Difficult features to recreate (cont)

Features	AUC
All features	86%
["euribor_3m", "jour_semaine", "indice_consommation", "nbr_employes"]	85%
['age', 'statut_marital', 'education', 'en_defaut', 'proprietaire', ''pret_consommation', 'contact', 'mois', 'duree', 'nbr_actuel', 'njours', 'nbr_precedent', 'resultat_precedent', 'souscription']	55%

What makes a feature hard to learn?

Important numerical features





Unimportant numerical features





NaN Values



Potential future improvement

- Different preprocessing step for numerical features -> particularly for multimodal distributions
- Different method to encode NaNs in numerical features
- Tuning process guided by machine learning efficiency + lambda (membership inference attack)

Q2 problem statement

Research question : Is it possible to reconstruct a real-world observation starting from synthetic data?

The problem consisted on predicting labels for a real -incomplete- test set (%25 of data missing). The training set was complete but synthetic.

Two ideas were considered:

- An imputation method for imputing the missing values.
- Training models with training data injected with missing values.

The baseline performance

We explored several classifiers (SVM, Logistic Regression, XGBoost, KNN, lightGBM).

We evaluate the performance of several models on a validation partition:

Model	Accuracy - validation	Acc. Complete Data		
Majority Label	0.60082			
KNN	0.73297	0.75926		
SVM	0.73354	0.76235		
XGBoost	0.71914	0.75720		
lightGBM	0.71296	0.77058		

We observe around a 5% reduction in performance due to the missing data.

Injecting missing data in the training set

- We analyzed the distribution of missing values in the test set and inject missing values in the training set following the same distribution.
- Intuition: the model will learn to classify in scenarios with incomplete data.
- We can inject the same noise in a subset of the train set to get a **validation partition**.





- Issue: The injection of missing values is random and some configurations may be better than others.
- This can be addressed through a ensemble of models trained using different configurations of incomplete train data. We do this by varying the seed for the null injection.



MICE imputation

- 1. Impute missing values in each feature with temporary data derived from the non-missing values available for that variable (example, mean).
- 2. Remove the temporary data for a given feature and regress it using the other features that we do have (or that we have "imputed")
- 3. Use the fitted regression model to predict the missing values in that feature.
- 4. Repeat 2-3 iteratively for each of the variables that still have missing values.
- 5. Perform several 1-4 cycles.

W Issue with MICE: there is some variation depending on the order the features are imputed in the iterations. (the random seed)

We propose a Multi-MICE approach that uses several MICE-imputed versions of the test set to be predicted by a given model, the final label predictions are generated through a soft-voting mechanism over the predictions on the test .



Combining both approaches

• These two approaches can be easily combined by simply replacing the classifiers in the Null Injection ensemble with multi-MICE blocks





Experimental results on validation set

		Accuracy -	%
Model	Method	validation	improvement
SVM	none	0.73297	0.000
SVM+XGBoost+GaussianNB+lightGBM	None	0.74332	1.412
SVM	null injection	0.7737	5.552
SVM	multi-MICE	0.7634	4.148
SVM	null injection +	0 7757	5 832
		0.7757	0.002

Experimental results on test set

		Accuracy -	%
Model	Method	validation	improvement
SVM	none	0.7063	0.000
SVM+XGBoost+GaussianNB+lightGBM	none	0.7222	2.248
SVM	null injection	0.7460	5.619
SVM	multi-MICE	0.7698	8.990
	null injection +	0 7040	7.000
SVM	multi-MICE	0.7619	7.866
	null injection +		
SVM+XGBoost+KNN	multi-MICE	0.7937	12.361

Conclusion on Q2 problem

- We observe that it is possible to reconstruct to some extent some features of a real-world observation with only synthetic data as training.
- Through MICE and Null-Injection in the training dataset we dampen the decrease originated to the missing features in the test set.