

Risk Evaluation of LDP scheme LoPub against Variational Autoencoder

Andres Hernandez-Matamoros* and Hiroaki Kikuchi

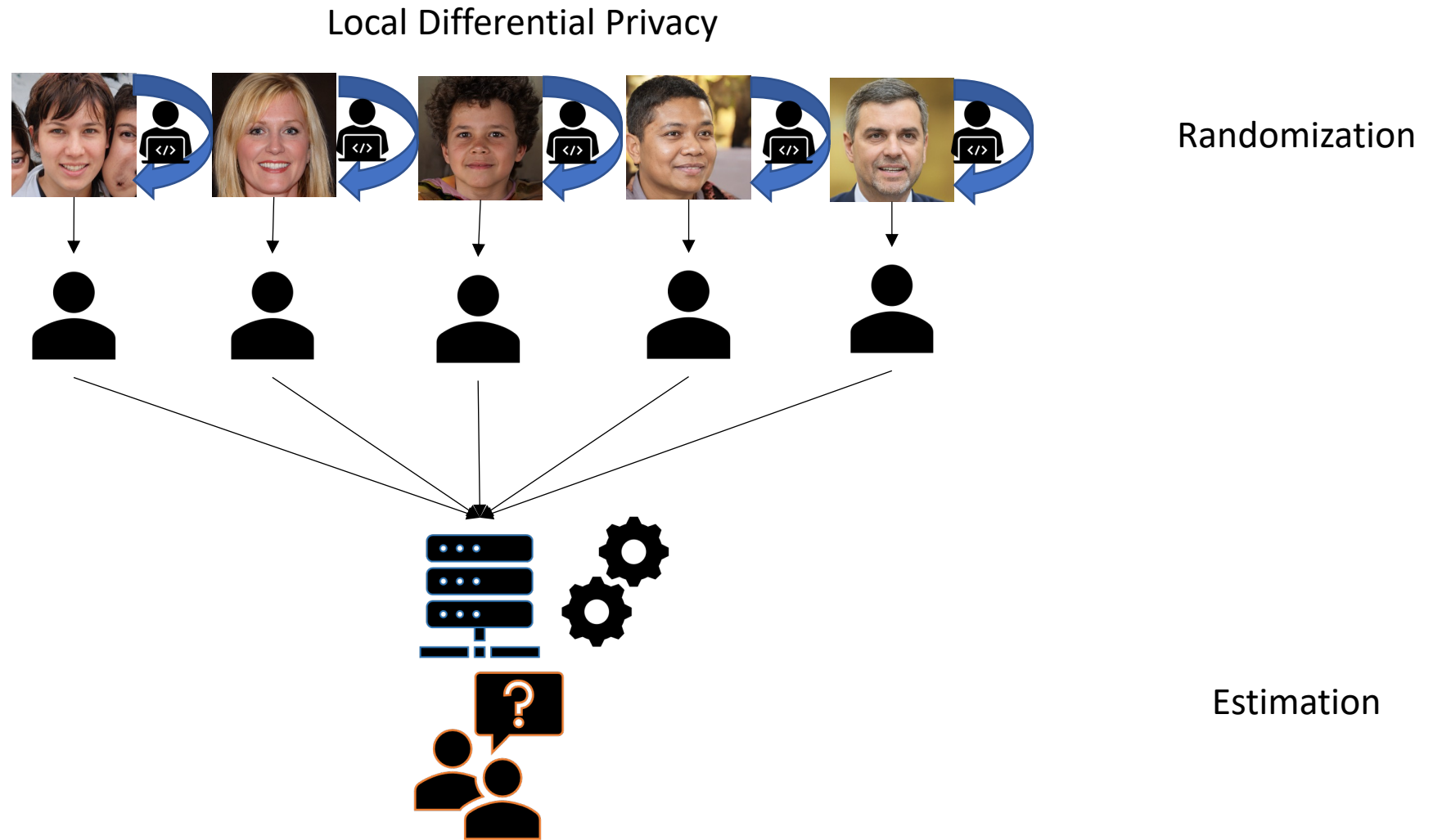
Meiji University

Computer Security Symposium 2022

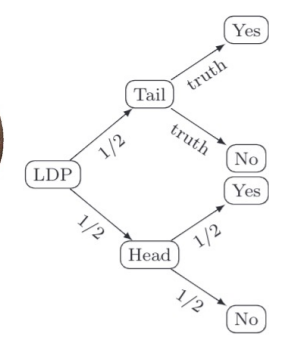
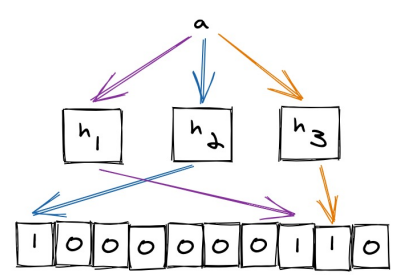
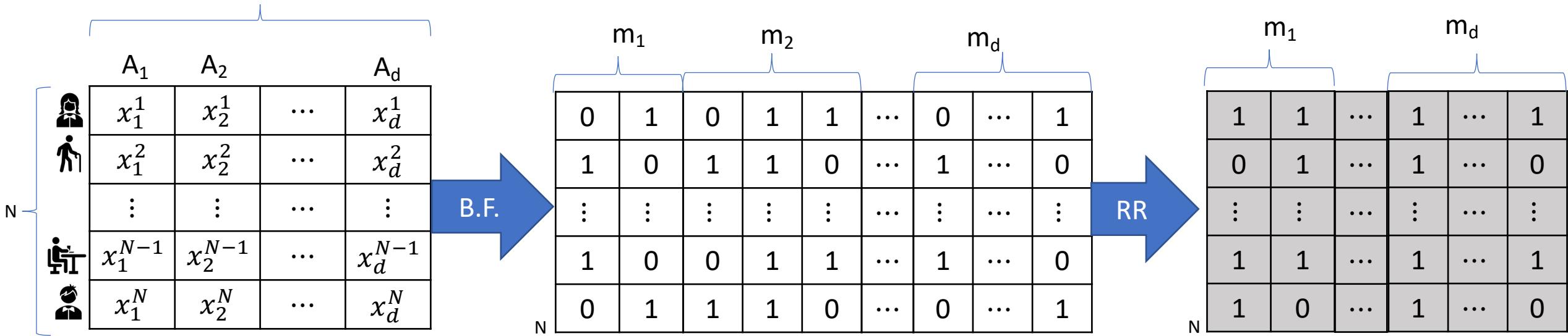
2022/10/24

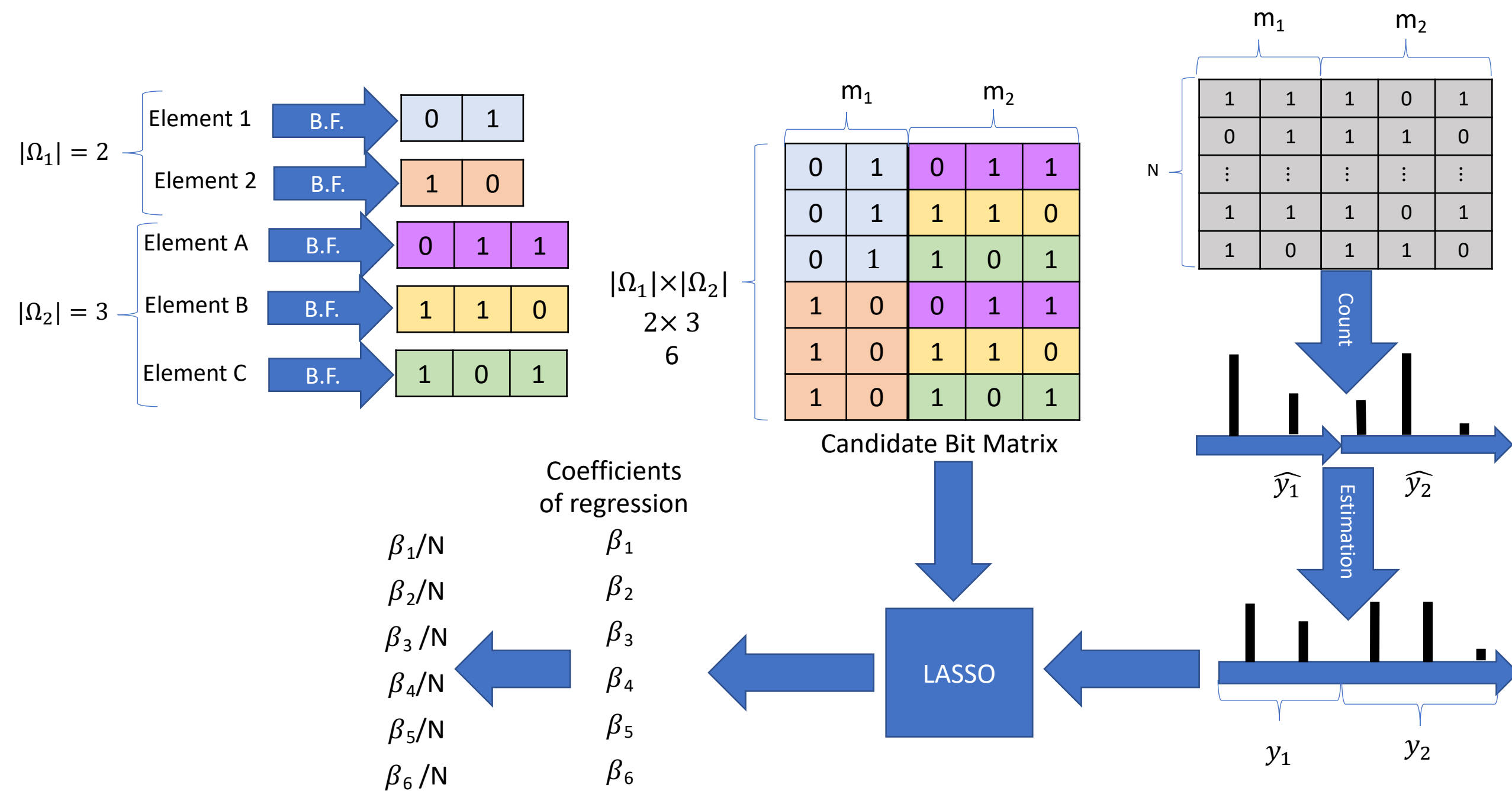
*matamoros@meiji.ac.jp

What is LDP?



d(Attributes)





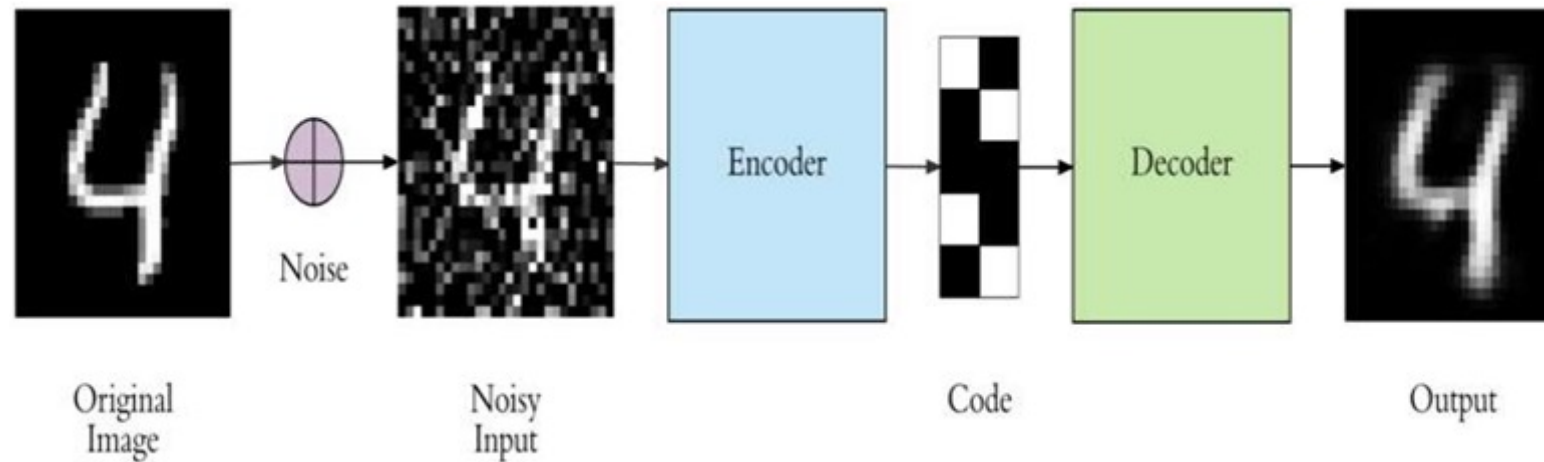
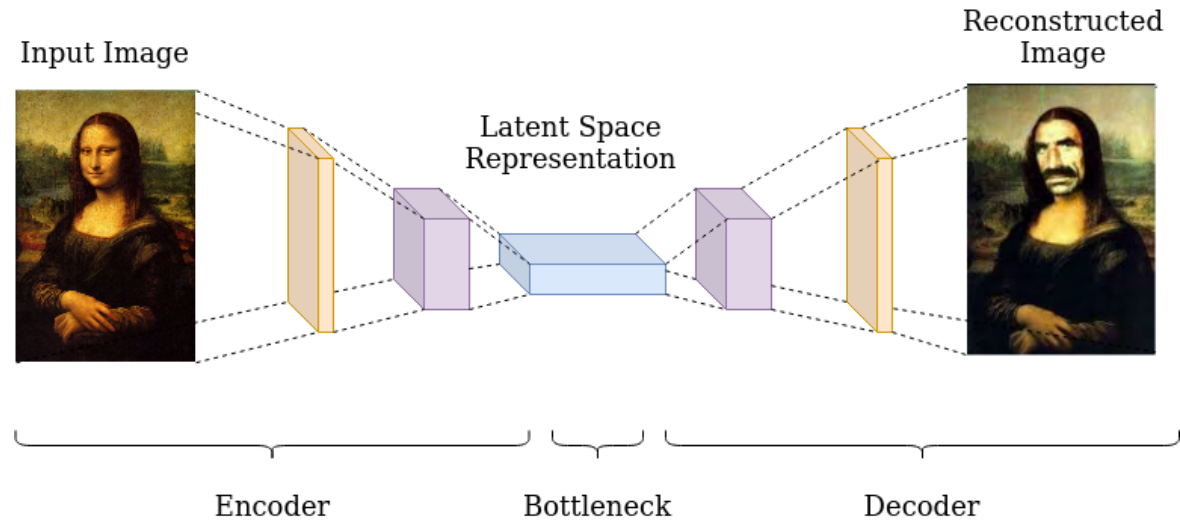
Pro vs Cons

- ✓ One/two-dimensional probability distributions can be efficiently estimated through the Lasso regression-based algorithm.
- ✗ The k -dimensional distribution estimations in LoPub still suffer from the low data utility when k is large.

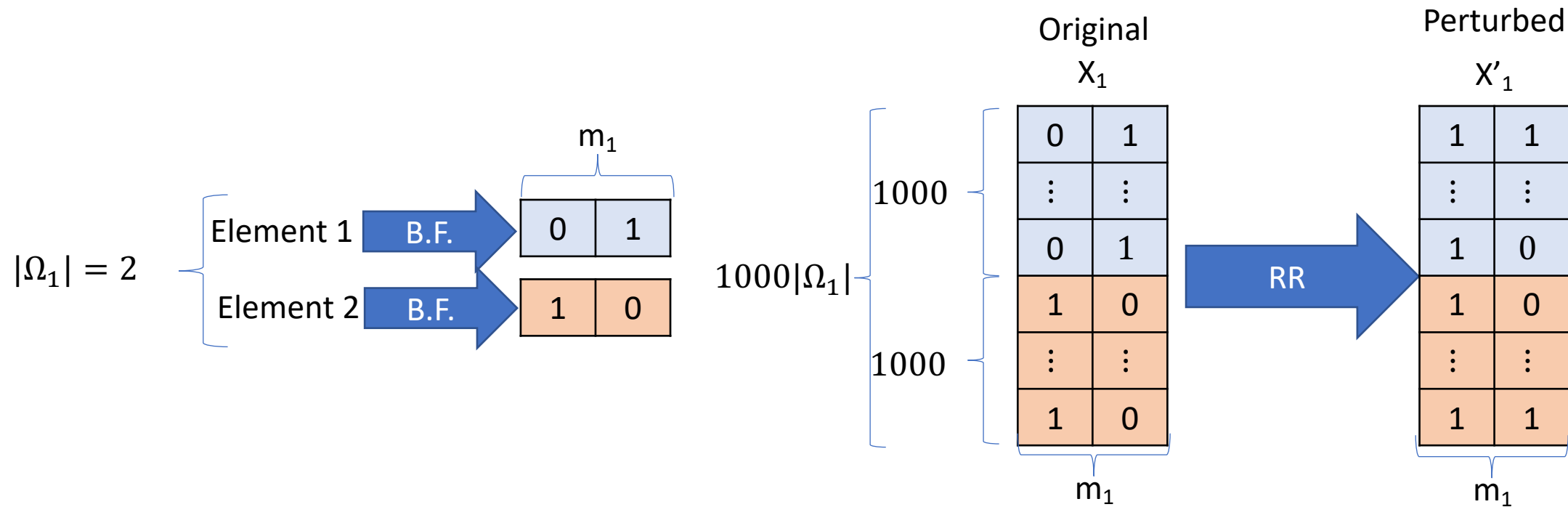
Proposal

	LoPub	Ours
Randomization	BF+RR	
Estimation	LASSO	VAE

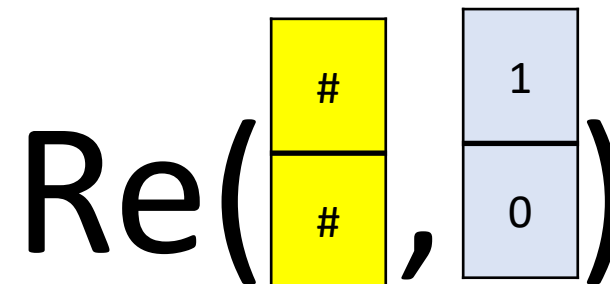
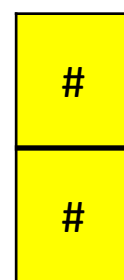
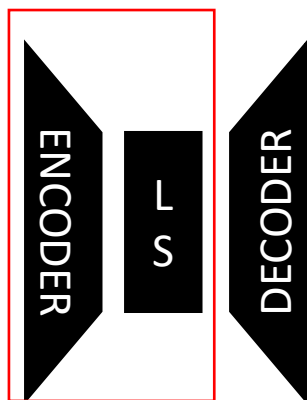
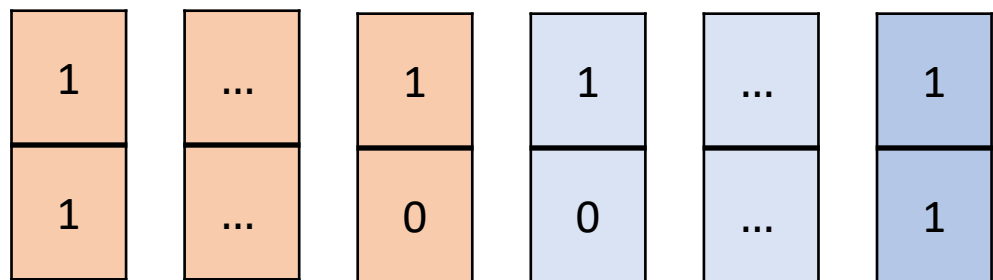
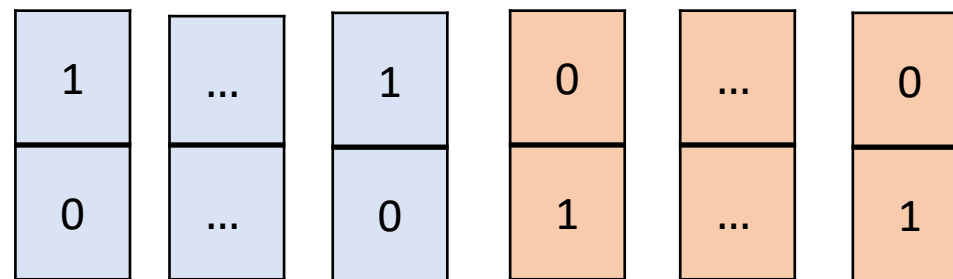
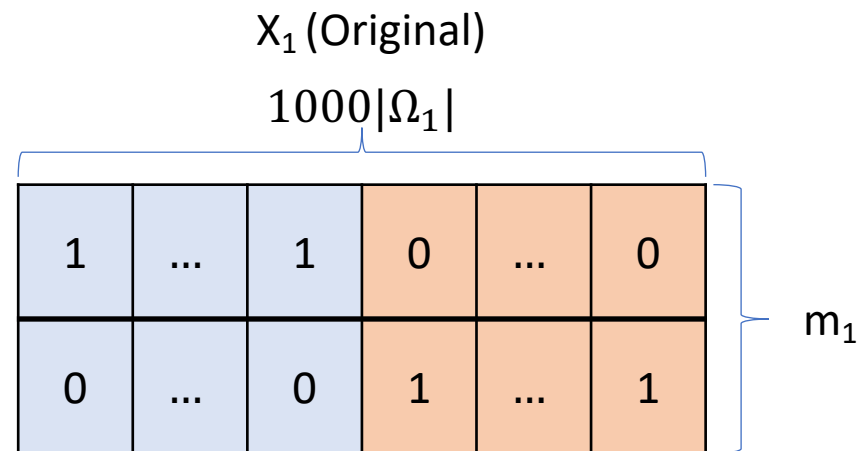
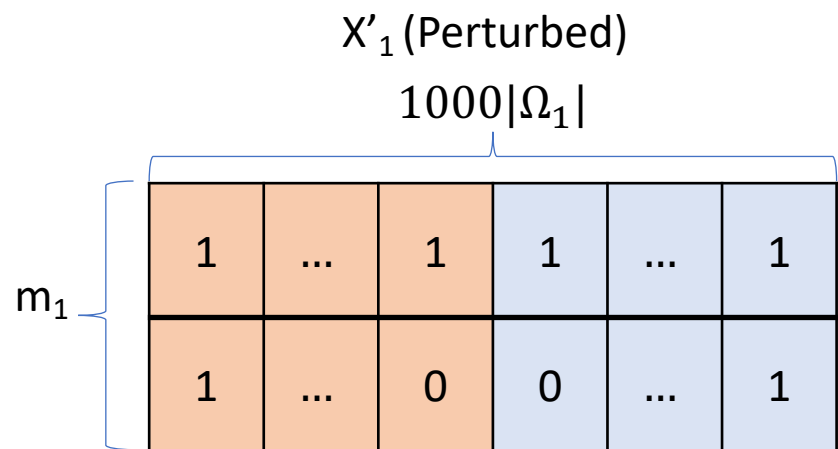
VAE-Applications



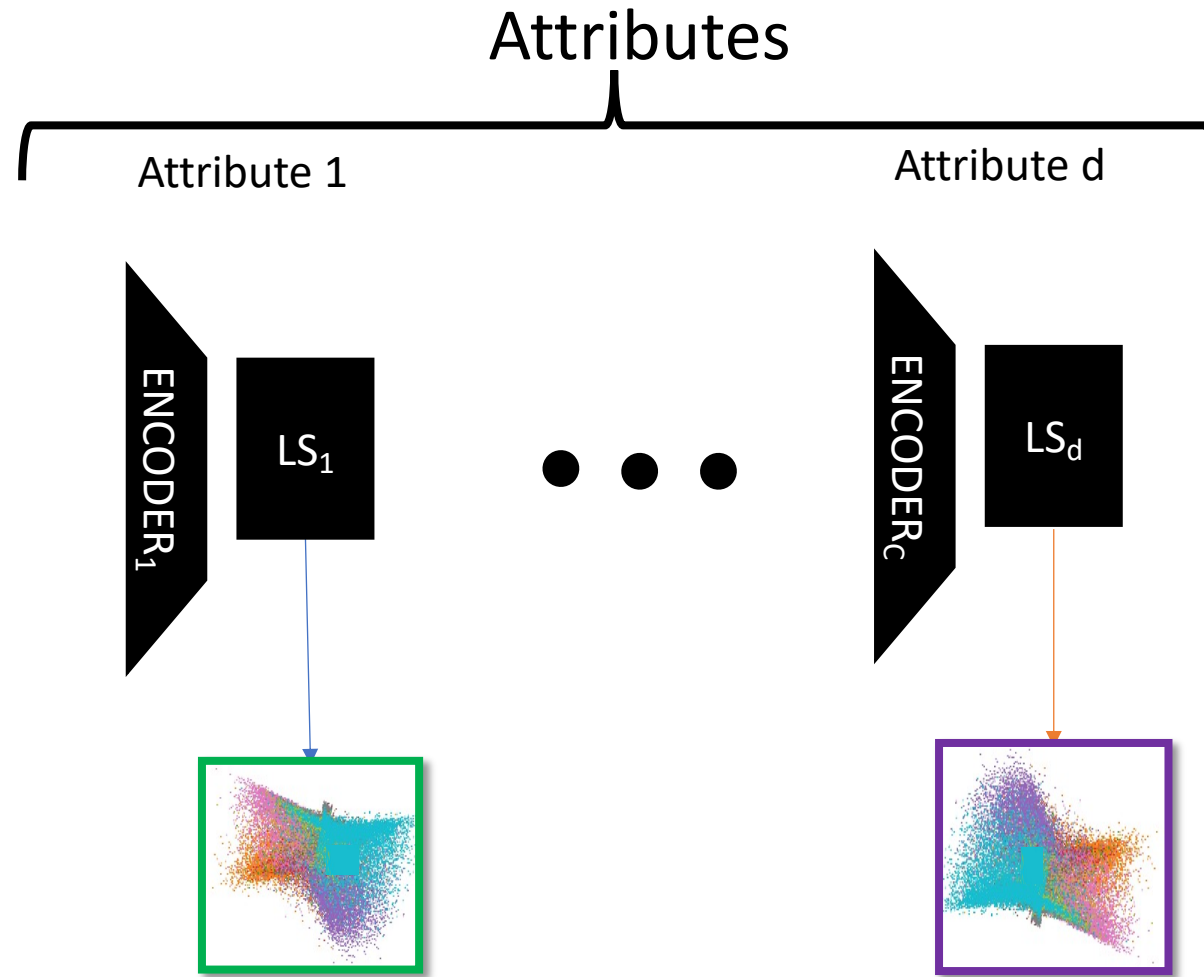
VAE Training-Dataset



VAE Training

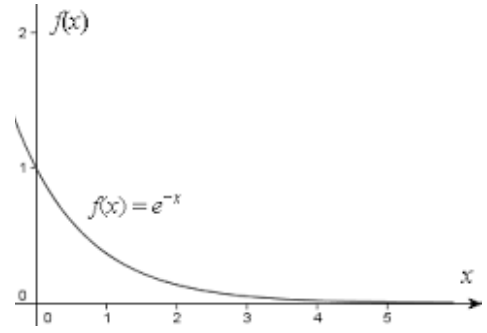
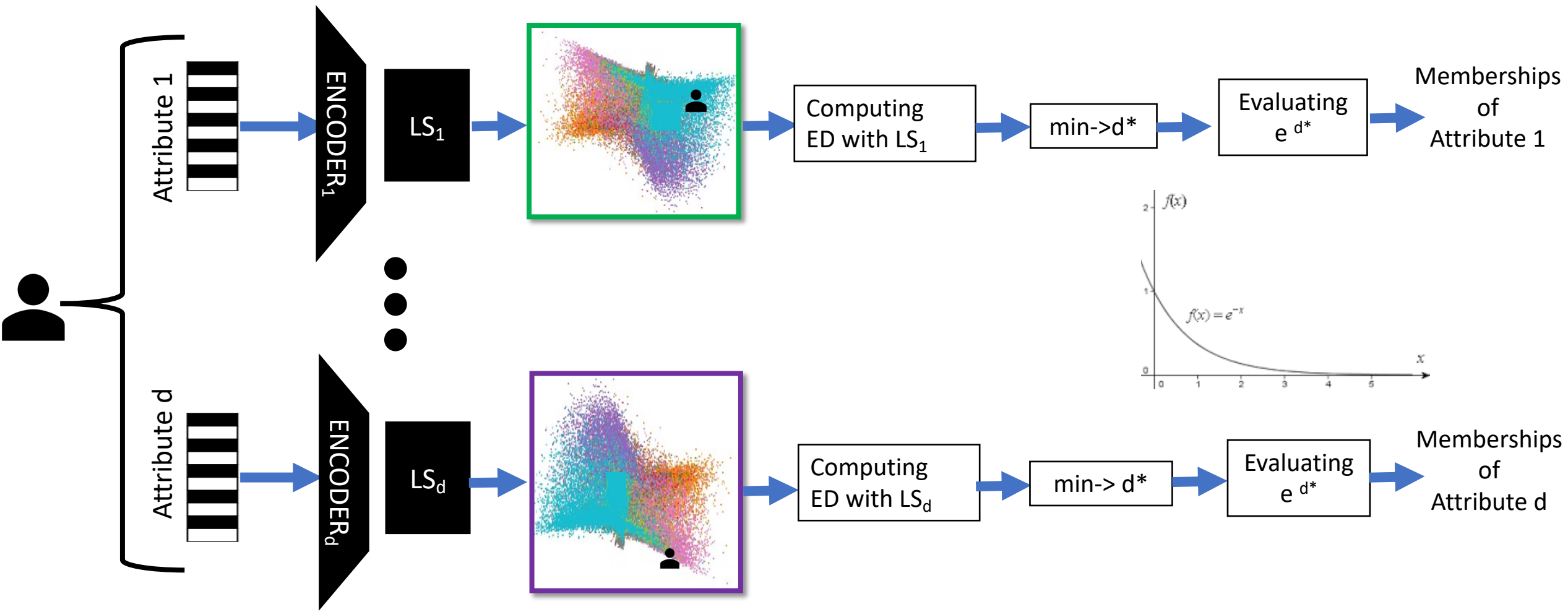


Latent Space's Attributes



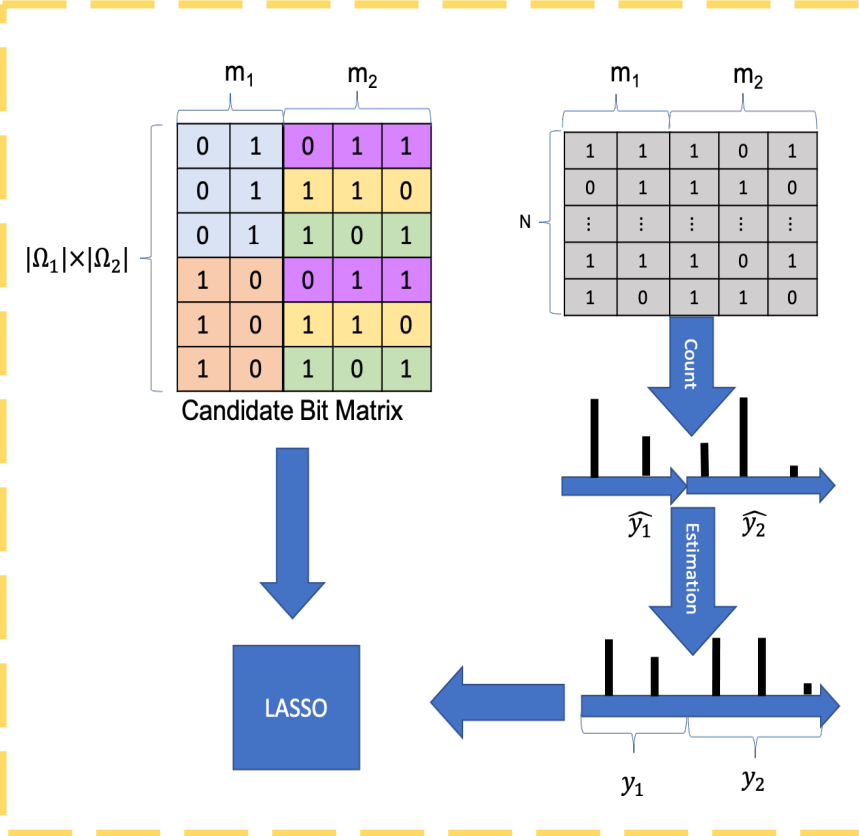
Latent Space Examples

Evaluation



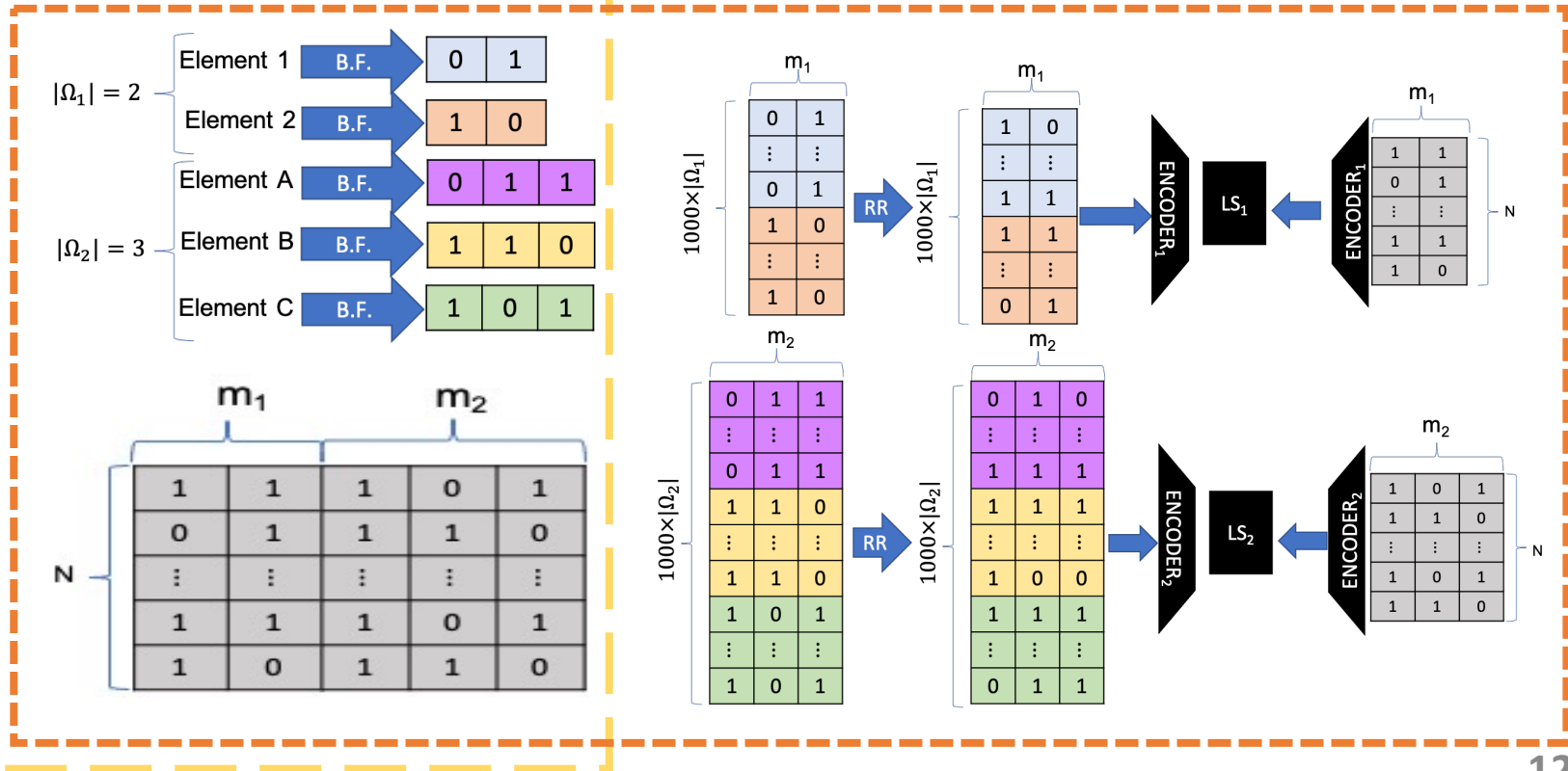
LoPub vs Proposal

LoPub



User
 Encode
 Perturb

Proposal

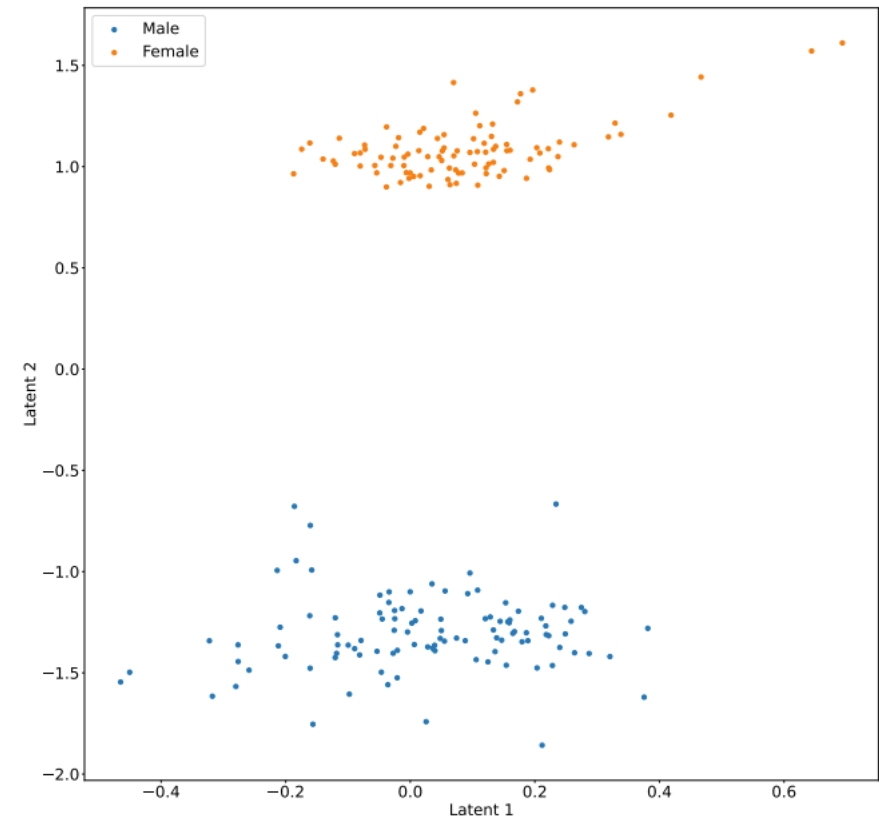
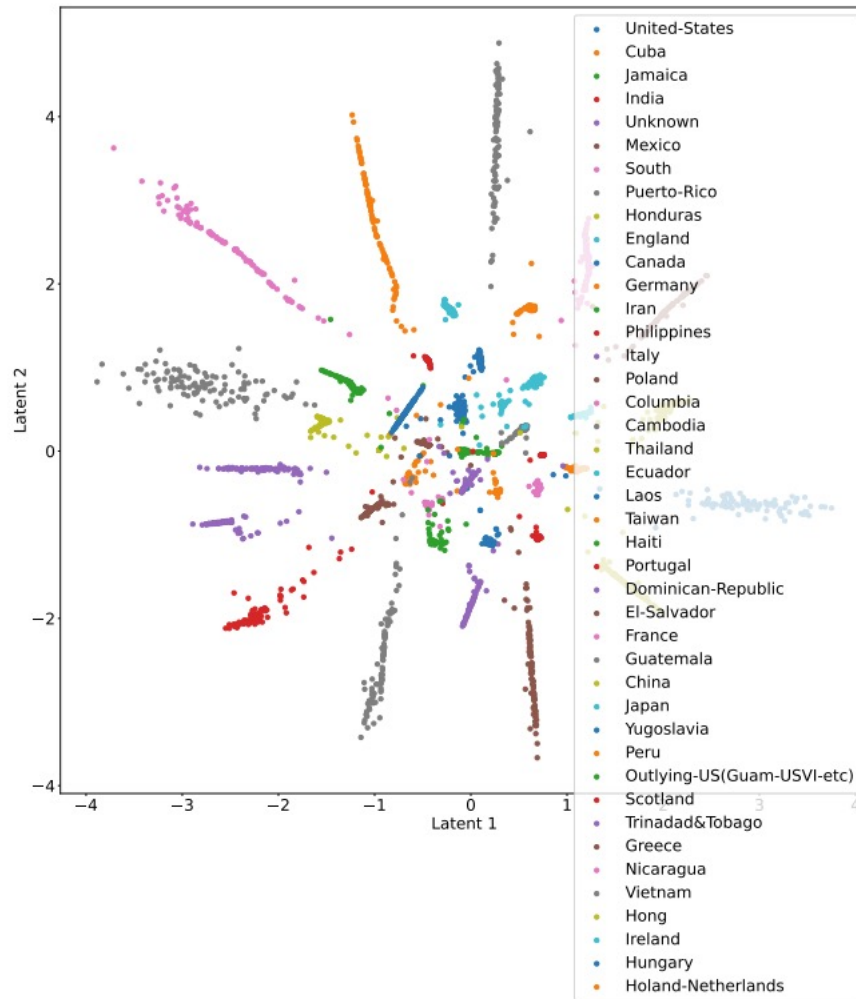


Datasets

Dataset	Users	Attributes
NHANES	4189	$ \omega_{Genre} = 2$
		$ \omega_{Race} = 5$
		$ \omega_{Education} = 5$
		$ \omega_{Marital} = 6$
		$ \omega_{Q_m} = 4$
Adult	32561	$ \omega_{Workclass} = 9$
		$ \omega_{Education} = 16$
		$ \omega_{MaritalStatus} = 7$
		$ \omega_{Occupation} = 15$
		$ \omega_{Relationship} = 6$
		$ \omega_{Race} = 5$
		$ \omega_{Sex} = 2$
		$ \omega_{Country} = 42$
$ \omega_{Target} = 2$		
Bank	45211	$ \omega_{Job} = 12$
		$ \omega_{Marital} = 3$
		$ \omega_{Education} = 4$
		$ \omega_{Default} = 2$
		$ \omega_{Housing} = 2$
		$ \omega_{Loan} = 2$
		$ \omega_{Contact} = 3$
		$ \omega_{Month} = 12$
		$ \omega_{Outcome} = 4$
$ \omega_Y = 2$		

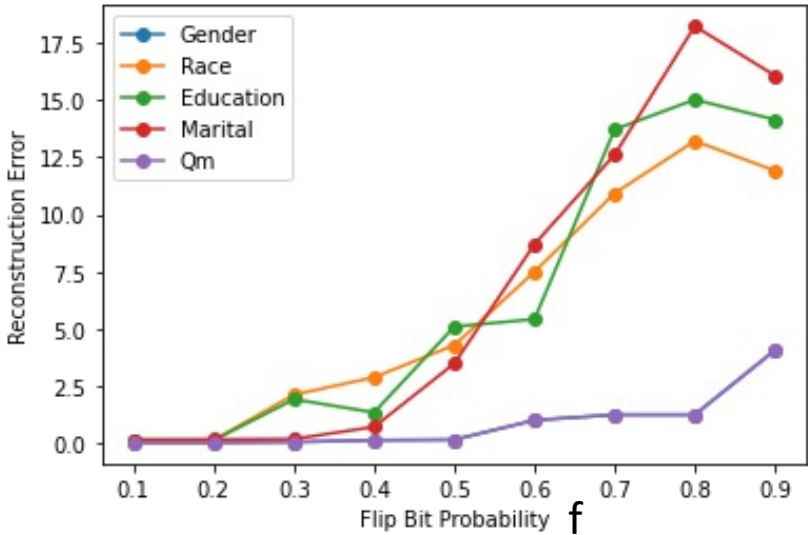
Latent Space

Examples of Latent Spaces in 2D

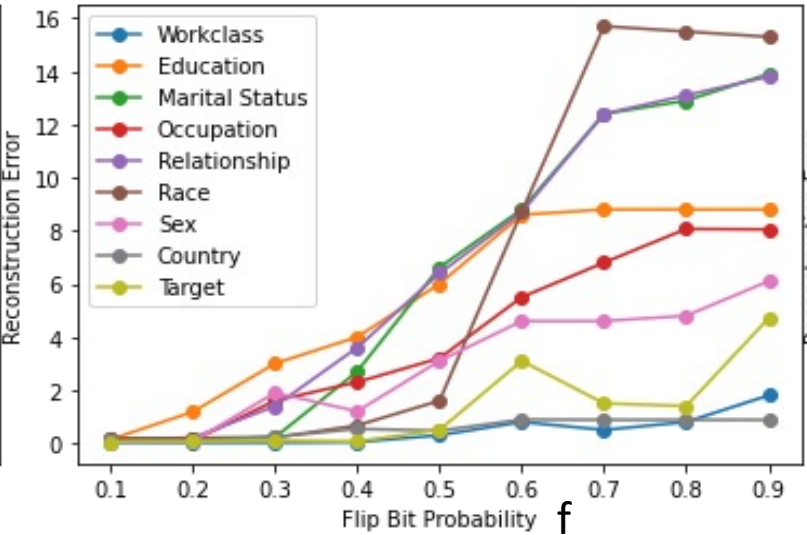


Country, Adults Dataset

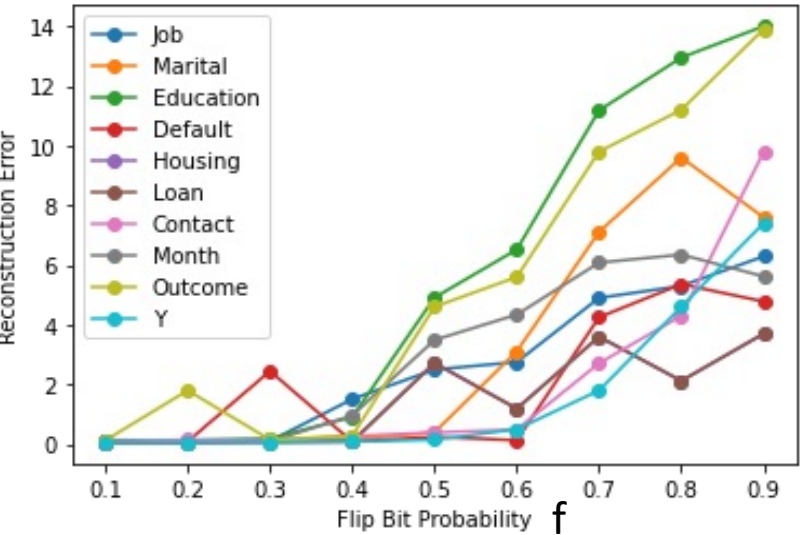
Reconstruction Error



NHANES Dataset



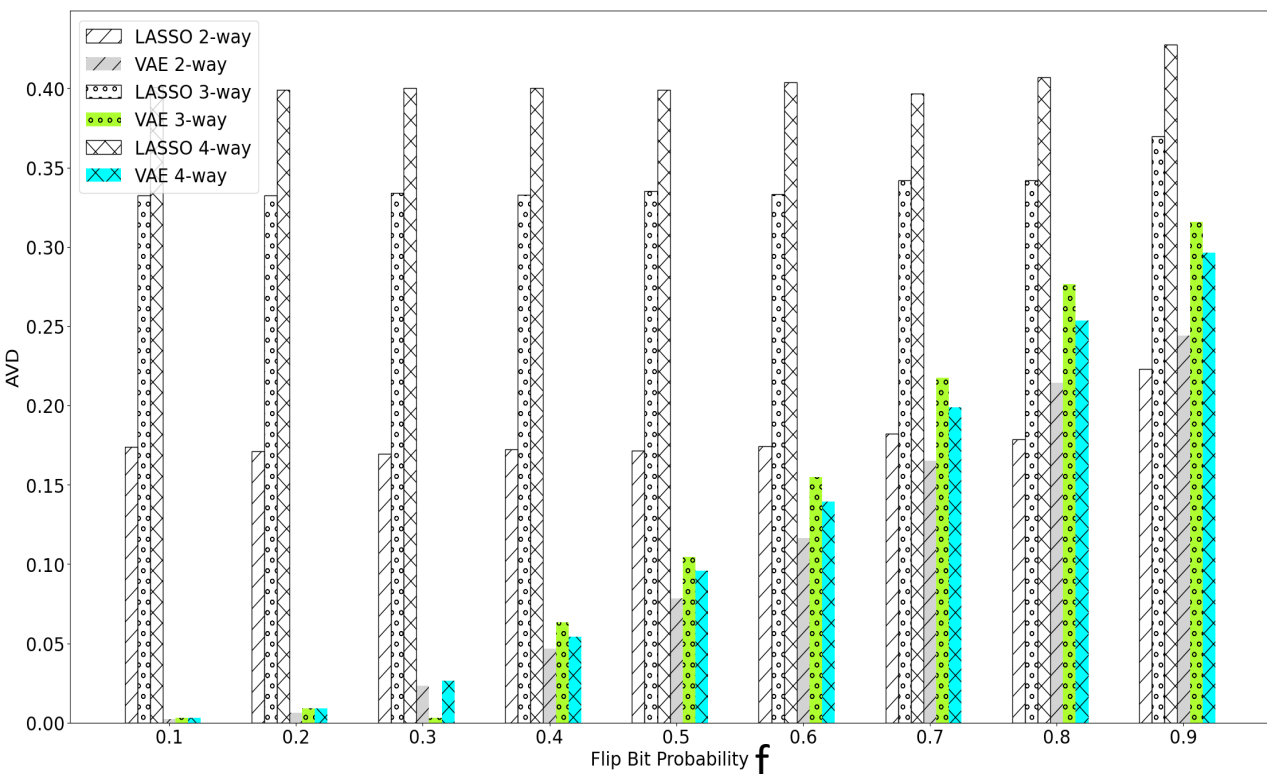
Adult Dataset



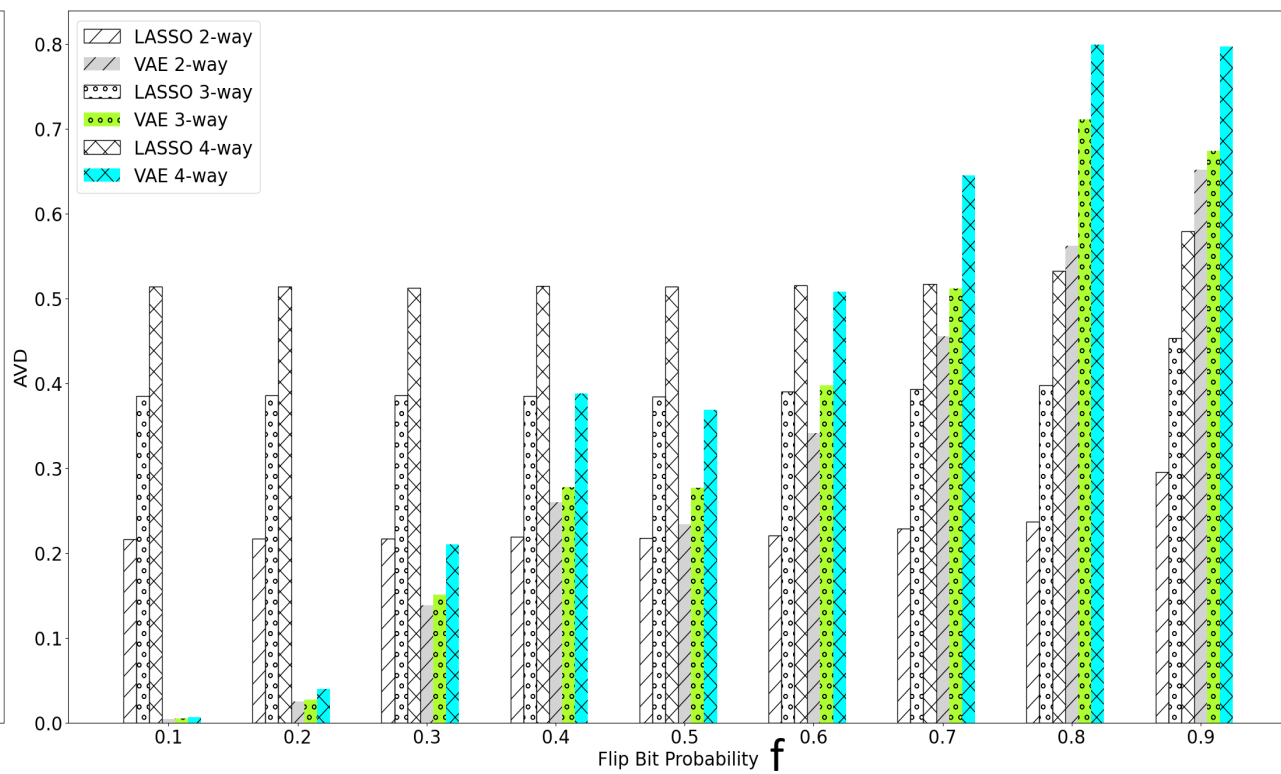
Bank Dataset

Accuracy K-way

NHANES Dataset



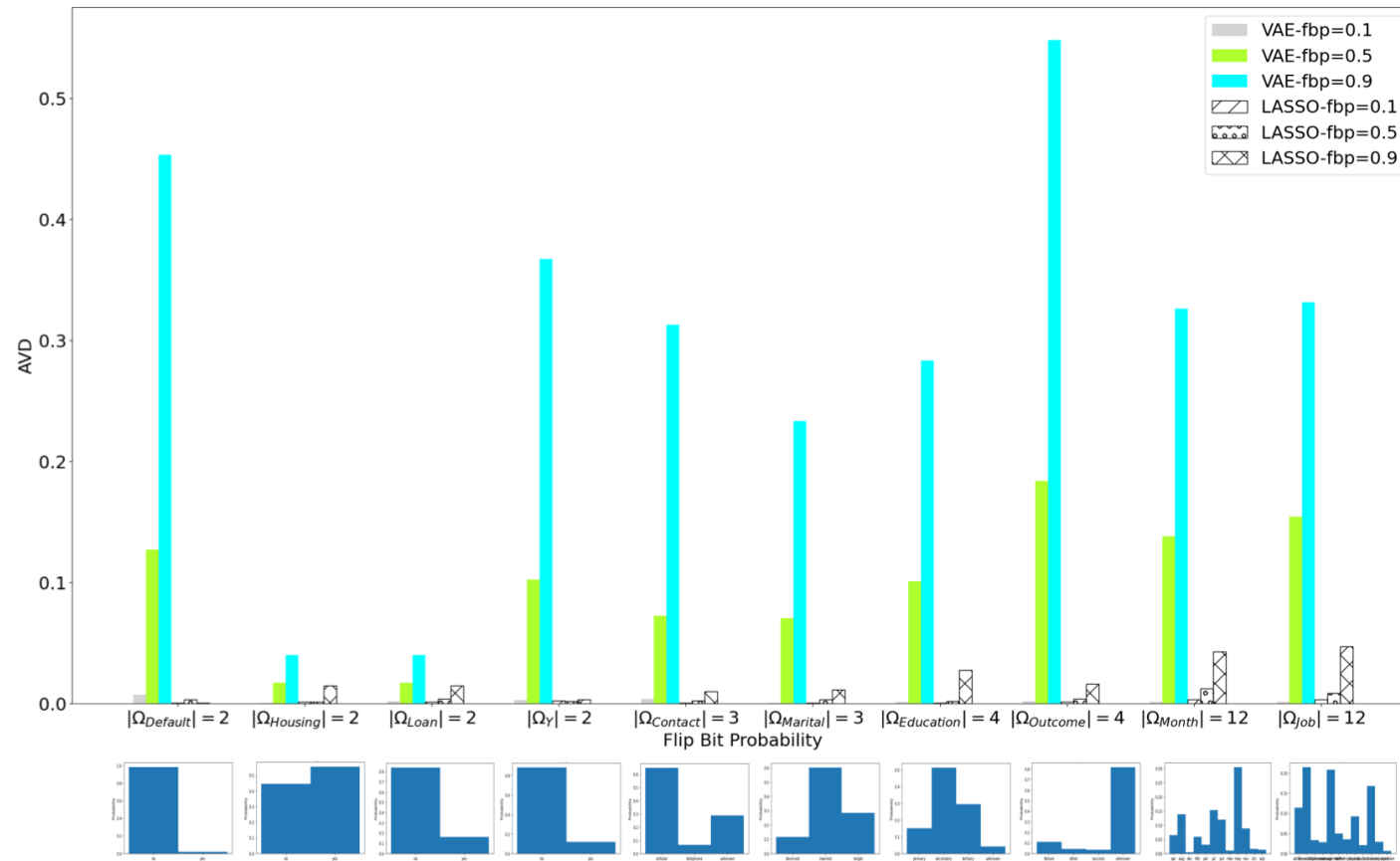
Adult Dataset



$$Dist_{AVD}(P, Q) = \frac{1}{2} \sum_{\omega \in \Omega} |P(\omega) - Q(\omega)|$$

Bank cardinality

Accuracy vs Cardinality



Cardinality

Distribution of elements

Bank Dataset

Conclusions

- VAE performs better in small dataset. (less than 5k)
- Cardinality and the distribution of the attribute impact the performance.
- Future study varying quantities of users, distributions, and cardinalities to quantify the performance of an LDP schemes.
- Try specific VAE for attributes with different cardinality.