
Federated Learning and Differential Privacy with Generative Encoding

David Banh^{1,*}

1 AskExplain

* david.b@askexplain.com

Abstract

A new method proposes to use a vector projection method to simultaneously encode the samples and features of multiple datasets into an aligned subspace that generates a private key from any given public key using the matrix decomposition. Importantly, the public keys contain vital information about all samples - with one clear concept, the feature space is encoded and made private. Given a dataset now with encoded features, this aligned and dimensionally reduced dataset can be secured based on the principles of cryptography and differential privacy for federated learning in a wide range of fields. Provided both the feature space and the existing data remain private, the integrity of the samples should be secure.

GitHub at: <https://github.com/AskExplain/gcproc>

Model

Using a Generative Encoding method called Generalised Canonical Procrustes (gcproc - see GitHub link in Abstract), multiple datasets can be aligned into a projected subspace with learned parameters. The model for this is given as:

$$X = L^T Z_0 u^T$$

$$Y = K^T Z_0 v^T$$

It is equivalent to Singular Value Decomposition (SVD) when the parameters L , K , u , and v are orthogonal - for example, $LL^T = I$.

Dimension reduction using the properties of SVD is expressed as:

$$Xu = L^T Z_0$$

$$Yv = K^T Z_0$$

If the left-hand side of the expressions are made private, then the original data X and the feature parameters u are non-identifiable, even if right-hand side $K^T Z_0$ and $L^T Z_0$ are made public.

Given the useful property of SVD, Yv and Xu are the dimensionally reduced forms of X and Y , useful information for federated learning secured by the principles of differential privacy.

1 Federated Learning

In federated learning the datasets X and Y are usually stored securely and locally, with computations done elsewhere - similarly for Generative Encoding via gepoch both the datasets and the feature encoding parameter (u for X , and v for Y) are made private. This leaves the parameters Z_0 , L and K to be made public.

The question then becomes - how can both X and Y be aligned to the same subspace for further analysis, that is, useful analysis for the dimensionally reduced data $L^T Z_0$, and $K^T Z_0$?

The solution to this is to generate Z_0 publicly, even randomly generated if so wished. The parameter hereby termed the "code", Z_0 , can be passed from public servers to the local secure servers such that L and K are computed locally, then to be assigned for public use by data science practitioners for federated learning. Of most vital importance, the data X and Y as well as the feature encoding parameters u and v are stored privately, thus ensuring all information in the original space of X and Y are held securely with integrity.

This system provides a public key for each dataset: K and L , and a private key for each dataset: v and u . The code Z_0 is a special key in that it can be generated privately or publicly based on the intent of use. However, the code must be known and used with the public key K and the private key v to reconstruct the original dataset $Y = K^T Z_0 v^T$, or $X = L^T Z_0 u^T$.

2 Differential Privacy

The public data passed on for federated learning are: $K^T Z_0$ and $L^T Z_0$, with the original datasets Y and X held securely on a private server along with the feature encoding parameters v and u . Provided the feature parameters are secure and non-identifiable from the public data $K^T Z_0$, and $L^T Z_0$ then the dimension reductions ideally ensure differentially privacy.

Even if Z_0 is known and thus K and L become identifiable through matrix inversion of Z_0 , the original datasets Y and X are not compromised provided v and u are not released publicly.

The use of the code Z_0 , and the presence of a public key with the need of a private key to reconstruct the original data ensures security through knowledge based on the foundations of cryptography.