FastShap++

Differentially Private FastSHAP for Federated Learning Model Explainability

Background

Federated Learning eXplainability FastSHAP

Federated Learning

Federated Learning - FedAvg



Federated Learning - FedAvg



Flower Framework



Powerful Abstractions

High Flexibility

Fast Prototyping

Fast Deployment

explainable Artificial Intelligence

eXplainable Artificial Intelligence

Why explanations are produced?

Red Team - Validation Centered

- Research on Data
- Explore Models
- Debug

Blue Team - Human Centered

- responsi**B**le
- Legal issues
- tr**U**stfulness in predictions
- Ethical issues

Explanation Taxonomy



Shapley Values as Feature Importance Explanations



Feature Importance Explanation



Privacy Enhancing Techniques

Differential Privacy

The higher the epsilon the less noise



FastShap Architecture



Shapley values as FI Explanations

$$\begin{split} \phi_{i}(v) &= \frac{1}{d} \sum_{s_{i} \neq 1} {\binom{d-1}{\mathbf{1}^{\top} s}}^{-1} \left(v(s+e_{i}) - v(s) \right)^{\text{Shapley value}} \\ \hline \phi(v_{x,y}) &= \operatorname*{arg\,min}_{\phi_{x,y}} \mathbb{E} \left[\left(v_{x,y}(\mathbf{s}) - v_{x,y}(\mathbf{0}) - \mathbf{s}^{\top} \phi_{x,y} \right)^{2} \right] \\ &\qquad \text{s.t.} \quad \mathbf{1}^{\top} \phi_{x,y} = v_{x,y}(\mathbf{1}) - v_{x,y}(\mathbf{0}), \\ \hline \phi_{\text{fast}}(\mathbf{x}, \mathbf{y}; \theta^{*}) &= \phi(v_{\mathbf{x}, \mathbf{y}}) \text{ almost surely in } p(\mathbf{x}, \mathbf{y}) \\ \hline \end{split}$$

Imitating Black Box with masks (Surrogate)

- D_{KL}:Kullback Leibler Divergence
- gamma: Black Box
- theta: model's parameters
- y: bb prediction
- m: mask function
- x: instance

$\mathcal{L}(\beta) = \mathbb{E}_{p(x)} \mathbb{E}_{p(b)} \left[D_{KL}(\gamma(x;\theta) || \hat{\gamma}(y | m(x,b);\beta) \right]$

- b: mask
- beta: surrogate parameters

Predicting SHAP Values (Explainer)

• Unif(y): uniform distribution over the classes

$$\mathcal{L}(\eta) = \mathbb{E}_{p(x) \operatorname{Unif}(y)} \mathbb{E}_{p(b)} \left[\left(v_{x,y}(b) - v_{x,y}(0) - b^T \phi_{fast}(x,y;\eta) \right)^2 \right]$$

Predicted: What the explainer thinks will happen when using certain features Actual: What actually happens when those feature combinations are tested Challenges And Current approaches

XAI in Federated Learning Context

Current Solutions and Challenges



The majority of Explanation methods (specifically SHAP based) needs some data

• Kernel Shap Explainers need a reference dataset

shap.KernelExplainer

class shap.KernelExplainer(model data, feature_names=None, link='identity', **kwargs)

- Data cannot be shared in FL
- Every client has its data, therefore local Shap Explanations differs
- Having a fixed, shared reference dataset makes the computation easier
 - Yet does not protect the local data against privacy attacks

Current Approaches

Specifically to SHAP

- Federated Fuzzy C Means as Background Data Points
 - Shares Centroids
- Aggregation of explanations or Explanations as a Service
 - Share Explanations

Both have strong points and weak points

Our solution approaches the challenge differently



Differentially Private FastSHAP for Federated Learning Model Explainability Quick explanations Privacy Enhancing Technology mechanism Share only and solely the model's weights



Privacy Protected - Differential Privacy









Experiments

Dataset

- Dutch: 70 Clients
- Income: 51 Clients
- Employment: 51 Clients

Modalities

- Vanilla Pipeline benchmark
- Semi-private (only S)
- Full-Private
 - Multiple levels of privacy (only E)

Measures

- BlackBox
 - Accuracy
- Surrogate
 - Fidelity
- Explainer
 - ShapGaps
 - L2-distance
 - Cosine Similarity
 - Feature Agreement
 - Sign Agreement
 - Ranking Correlation
 - Delta Faithfulness



DP impact in Centralized

	Accuracy BB	Fidelity Surrogate
Dutch	0.83 ± 0.01	0.97 ± 0.01
Dutch (DP)	0.82 ± 0.01	0.91 ± 0.01
ACSI	0.77 ± 0.01	0.97 ± 0.01
ACSI (DP)	0.77 ± 0.02	0.95 ± 0.01
ACSE	0.80 ± 0.01	0.95 ± 0.01
ACSE (DP)	0.74 ± 0.03	0.90 ± 0.01

average and standard deviation on three runs

Degradation in performance for the Black Box (using DP as PET)

Small yet visible Impact of DP on Surrogate model

Those results say that epsilon-privacy guarantee come at the cost of accuracy or fidelity

Centralized vs Federated

Measuring the differences when explaining the Black Box Federated using

• Centralized FastShap

Compared with

• Federated FastShap

	ℓ_2 Dist. (\downarrow)	Cosine Sim. (†)	Feat. Agr. (†)	Sign Agr. (†)	Rank Corr. (†)	Δ Faith (\downarrow)
Dutch	0.03 ± 0.02	0.99 ± 0.01	0.87 ± 0.13	0.87 ± 0.13	0.84 ± 0.11	0.02 ± 0.03
ACSI	0.09 ± 0.05 0.10 ± 0.05	0.81 ± 0.22 0.80 ± 0.19	0.77 ± 0.14 0.66 ± 0.16	0.70 ± 0.17 0.64 ± 0.17	0.70 ± 0.17 0.65 ± 0.16	0.13 ± 0.10 0.21 ± 0.18

Explanation Measures - 1



Explanation Measures - 2



Future Work

Next experiments with

- For image data
- Introduction of Fairness
 - Private, Fair and Federated
 - Fair explanations metrics
 - Fairness as transferable, stable portable property
 - Measured with surrogate
 - Fidelity
 - And explainer
 - Specific measures

