



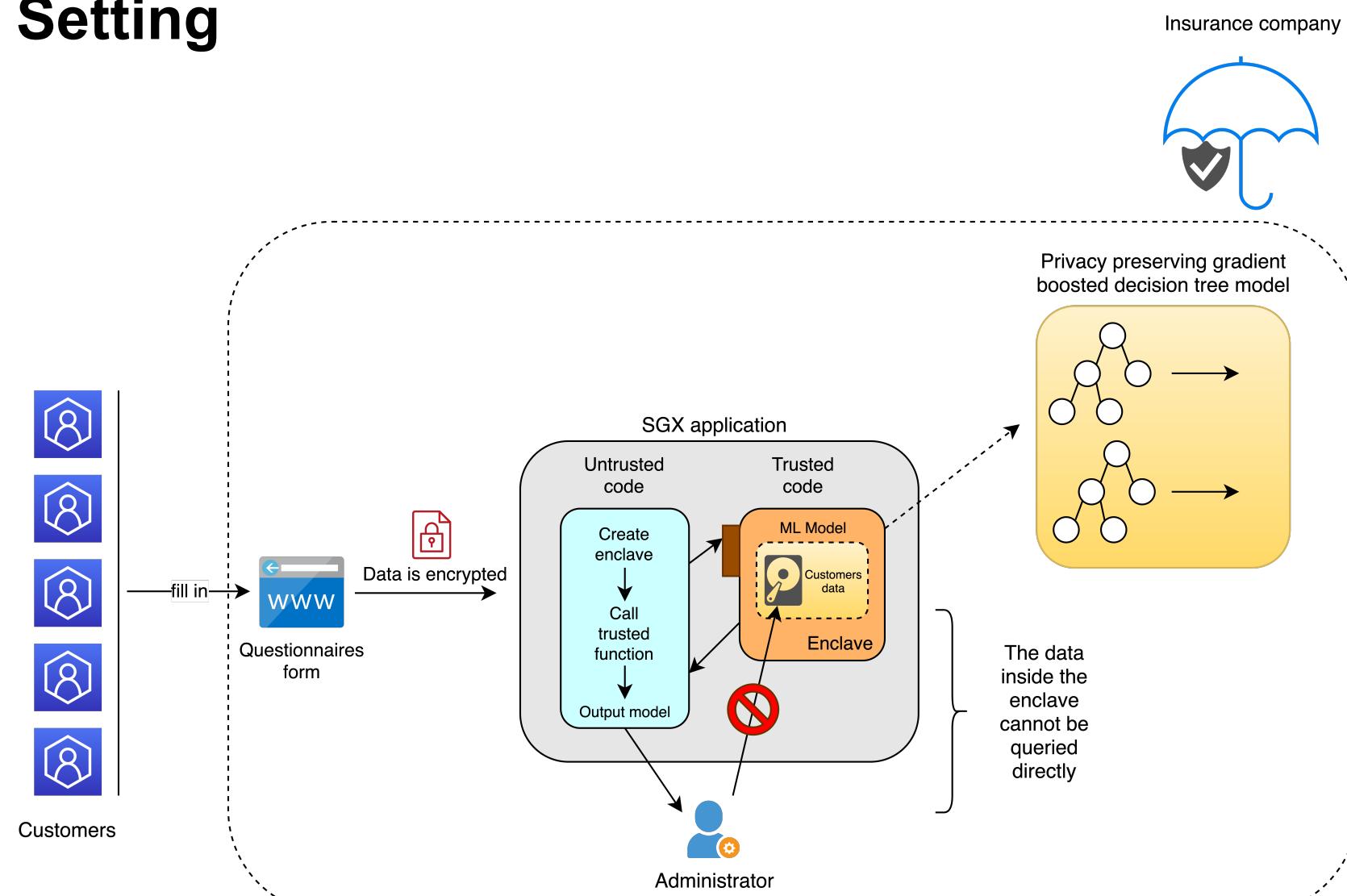
Privacy-Preserving Machine Learning for Cyber Insurance

Theo Giovanna (ETH), Prof. Dr. Esfandiar Mohammadi (Lübeck), Dr. Kari Kostiainen (ETH)

Motivation

- Insurance companies have come up with cyber insurance products
 - Evaluating cyber risk is hard
 - Use of questionnaires
 - Customers are not so willing to provide full details about their security practices
 - They need privacy guarantees
 - Lack of any historical data
 - About 10 real samples

Example Setting



Chosen Approach

- Differentially Private Gradient Boosted Decision Trees (DP-GBDT)
 - Great performance
 - Appropriate for low memory environments such as within an SGX's enclave
 - Easily explainable due to the decision trees' structure
 - Can work with little data (unlike e.g. DNN)
 - Provable privacy guarantees through differential privacy
 - Existing state of the art for DP-GBDT
 - AAAI - Qinbin et al. - Privacy-Preserving Gradient Boosting Decision Trees
 - <https://arxiv.org/pdf/1911.04209.pdf>

Thesis Focus

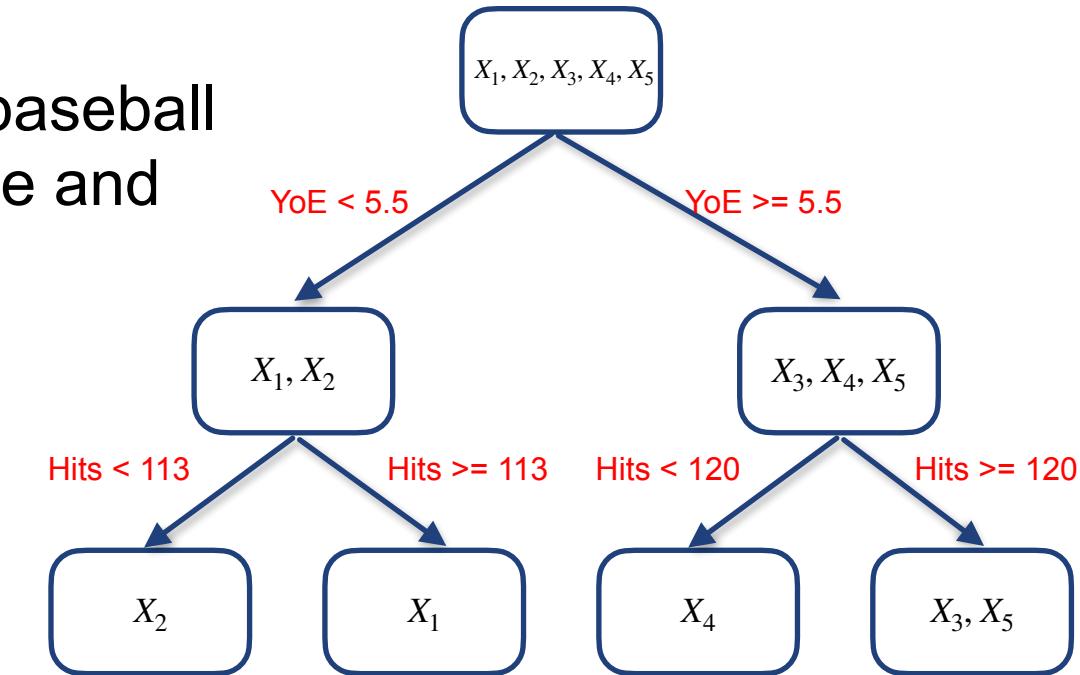
- Shed some light on the cyber insurance problem
 - **Contribution 1.1:** Implemented state of the art DP-GBDT from the literature, and tuned it specifically to cyber insurance
 - **Contribution 1.2:** Proposed a new decision tree induction algorithm to better performances over small datasets
- Since historical data is lacking, generate synthetic data that mimics cyber insurance questionnaire's answers
 - **Contribution 2:** Proposed a synthetic data generation algorithm based on Bayesian networks
- Results:
 - **Contribution 3.1:** Accuracy evaluation
 - **Contribution 3.2:** Privacy evaluation

Background

Background - Decision Trees

- Decision Tree
 - Goal: given a set of inputs X_1, \dots, X_n , predict a response or class Y for an unseen input X_i
 - Example (adapted from [RCS46](#)): predict a baseball player's salary based on years of experience and number of hits

	YoE	Hits	Salary
X_1	4	113	5.11
X_2	2	90	4.5
X_3	7	162	7.9
X_4	5.5	105	6.1
X_5	6	120	6.5



Background - Decision Trees

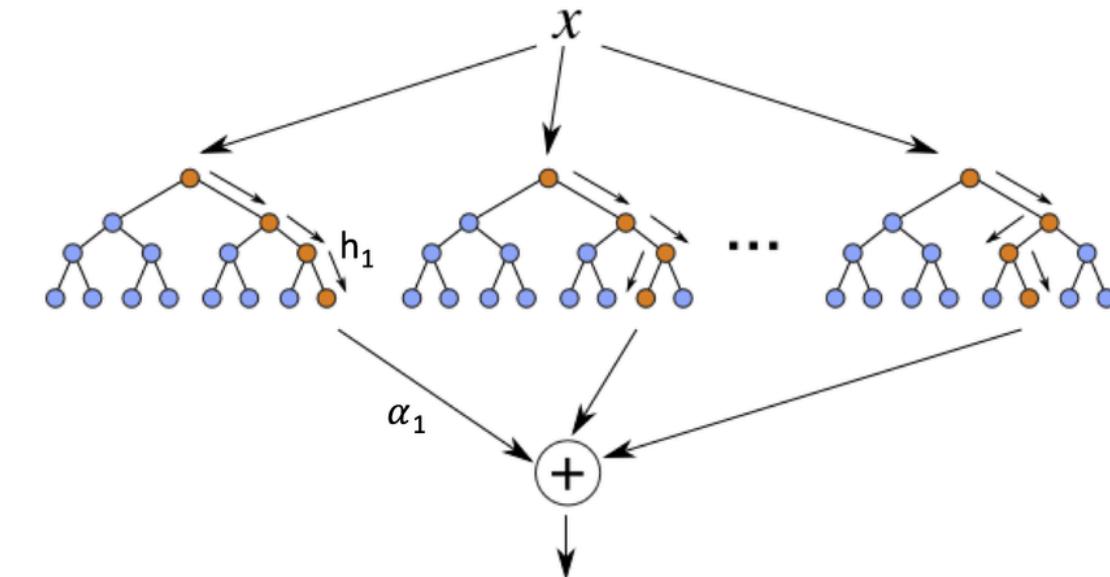
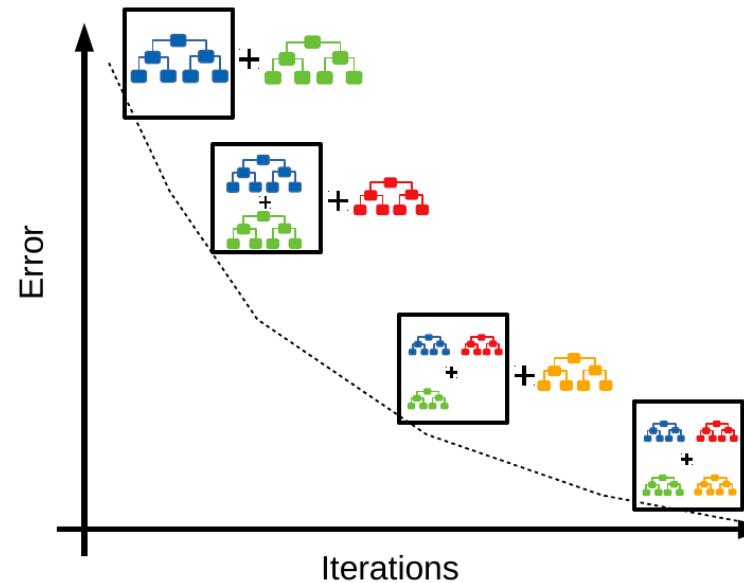
- Problems: decision trees are
 - Prone to overfitting
 - Can get complex when there's many class labels or when the tree depth is very large
 - A change in the training set might result in a complete different decision tree

Background - GBDT

- Solution: Gradient Boosted Decision Trees
 - *Ensemble* method: combine many weak learners (here, many simple decision trees) to build the final model - more trees allow for error correction
 - Each tree is fit on the residual of the previous tree
- Trees stop growing when termination criteria is reached
 - Depth
 - Gain
 - Number of samples in node

Background - GBDT

- Each new tree helps in reducing the overall error
- Use all trees to make final prediction



Background - Differential Privacy

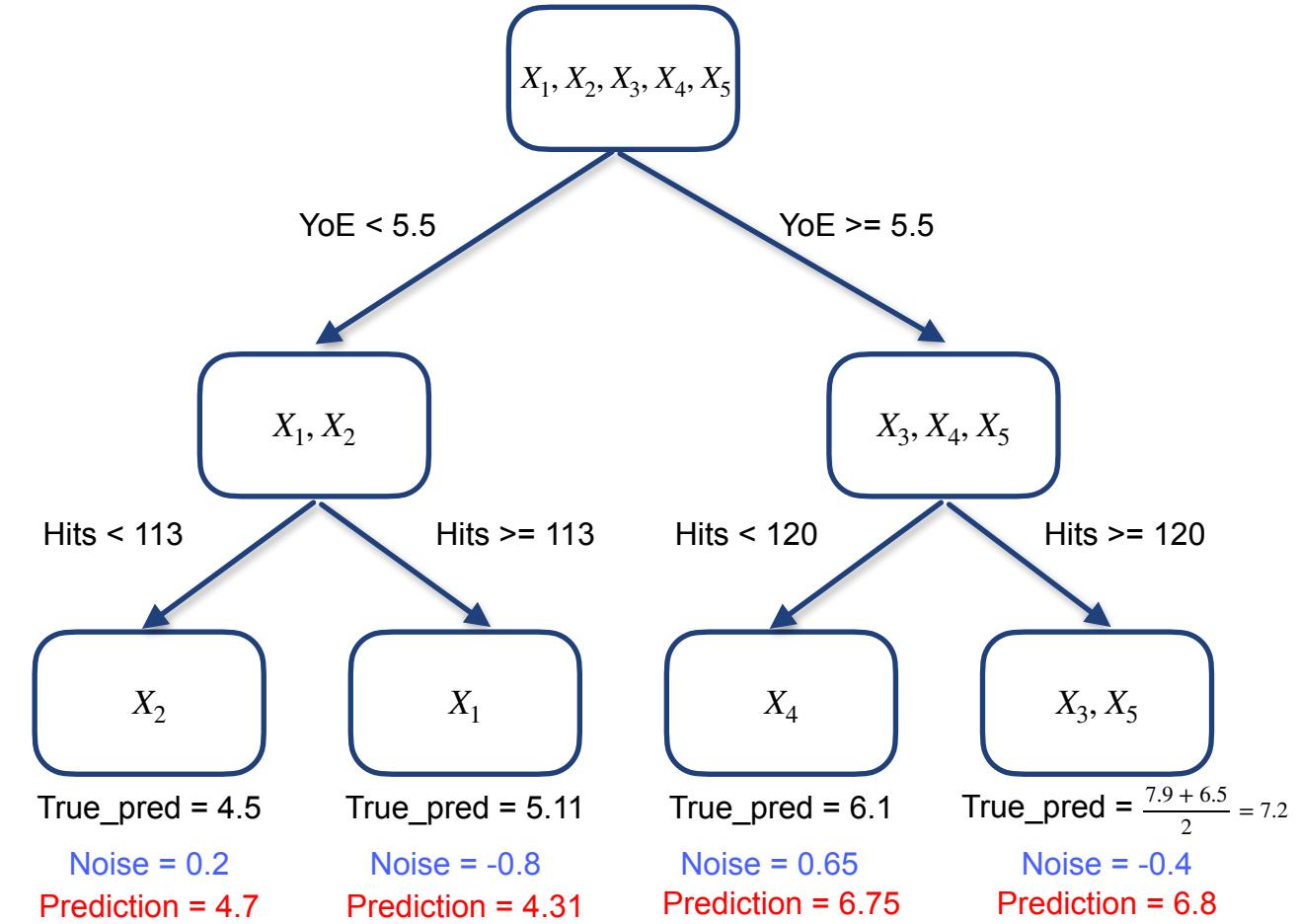
Definition 1. (ε -Differential Privacy) Let ε be a positive real number and f be a randomized function. The function f is said to provide ε -differential privacy if, for any two datasets D and D' that differ in a single record and any output O of function f ,

$$\Pr[f(D) \in O] \leq e^\varepsilon \cdot \Pr[f(D') \in O] \quad (5)$$

- The ratio between the output probability distributions for D and D' is strictly bounded by e^ε
- More intuitively, we aim to learn about a dataset, yet don't learn anything about a particular datapoint
- The influence of every single data point on the learned model is *deniable*
- Privacy 'level' controlled by *privacy budget* ϵ

Background - Differential Privacy

- Applying this to GBDTs: select the attribute and attribute's value to split on with the exponential mechanism:
 - Essentially, the exponential mechanism is a randomised `arg_max` which outputs high gain splits with high probability
- Once the termination criteria is reached within a branch, draw some random noise from a Laplacian distribution, and add it to the leaf node's prediction



Background - Differential Privacy

- Challenges:
 - How much noise is too little? Too much?
 - How to design models that query the data as less as possible while learning sufficient information
 - Each query consumes some of the privacy budget
 - The noise should not outweigh the data itself, which happens when the privacy budget is too low or when there's too little data

Contribution 1: Implementing and Improving DP-GBDT

Implementing DP-GBDT (Qinbin et al., 2020)

Algorithm 2: Differentially private GBDTs training process [25]

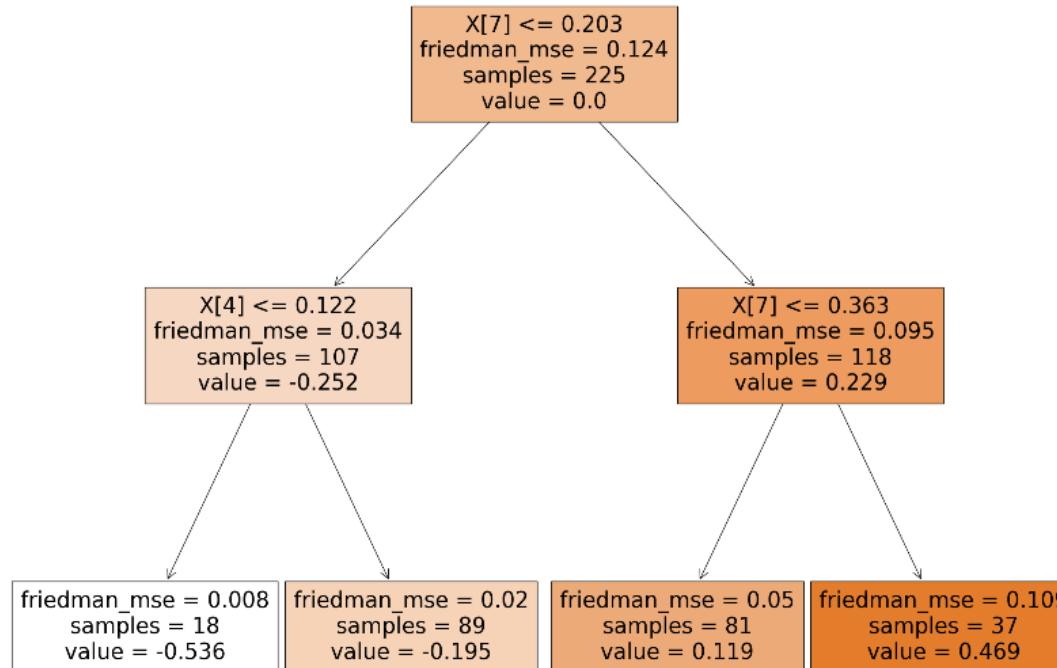
Input: $X = X_1, \dots, X_n$: instances, $y = y_1, \dots, y_n$: labels
Input: λ : regularisation parameter, d_{max} : maximum depth, η : learning rate
Input: T : total number of trees, l : loss function, ε : privacy budget
Output: An ensemble of trained differentially private decision trees.

- 1 $\varepsilon_t = \varepsilon$ \triangleright *Each tree is trained on a disjoint subset of the dataset, so we can apply Theorem 2.8*
- 2 **for** $t = 1$ **to** T **do**
- 3 Update gradients of all training instances on loss l
- 4 $\varepsilon_{leaf} = \frac{\varepsilon_t}{2}$, $\varepsilon_{node} = \frac{\varepsilon_t}{2d_{max}}$
- 5 **for** $d = 1$ **to** d_{max} **do**
- 6 **for each node in current depth do**
- 7 **for each split value i do**
- 8 $G_i \leftarrow \frac{(\sum_{i \in I_L} g_i)^2}{|I_L| + \lambda} + \frac{(\sum_{i \in I_R} g_i)^2}{|I_R| + \lambda}$ \triangleright *Equation 2.3*
- 9 $P_i \leftarrow \exp\left(\frac{\varepsilon_{node} \cdot G_i}{2\Delta G}\right)$ \triangleright *Theorem 2.6*
- 10 Split node on split value i , where i is chosen with probability $P_i / \sum_j P_j$
- 11 **for each leaf node i do**
- 12 $V_i \leftarrow \eta \left(-\frac{\sum_{i \in I} g_i}{|I| + \lambda} + Lap(0, \Delta V / \varepsilon_{leaf}) \right)$ \triangleright *Equation 2.4 and Theorem 2.5*

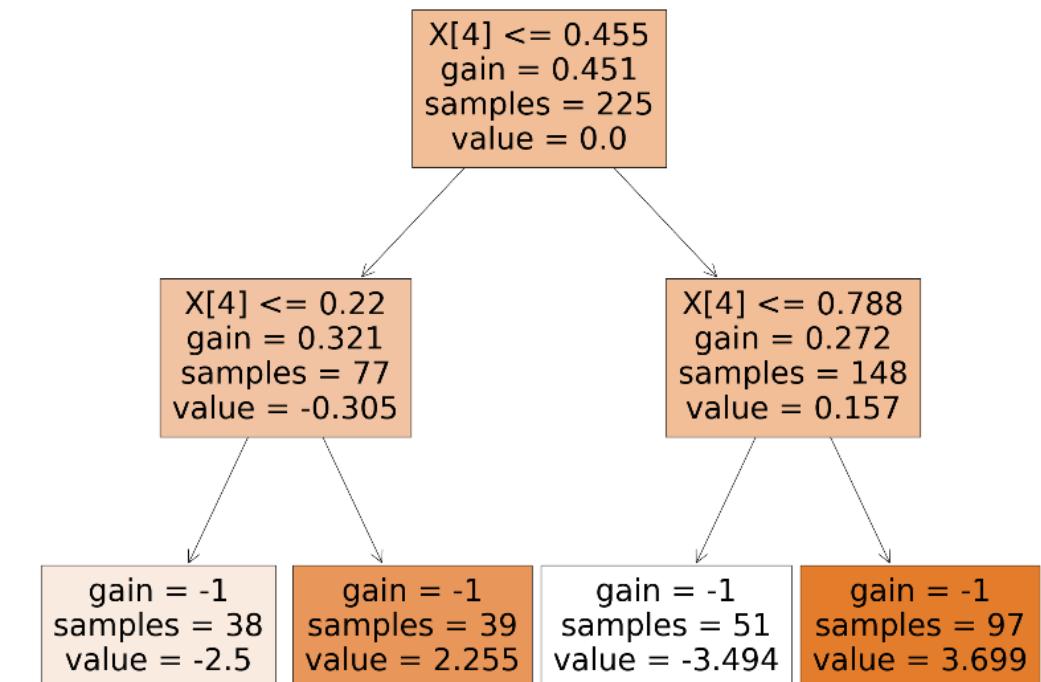
- Additionally:
 - Gradient Based Data Filtering
 - Geometric Leaf Clipping

Impact of DP

Dataset: abalone, samples = 300



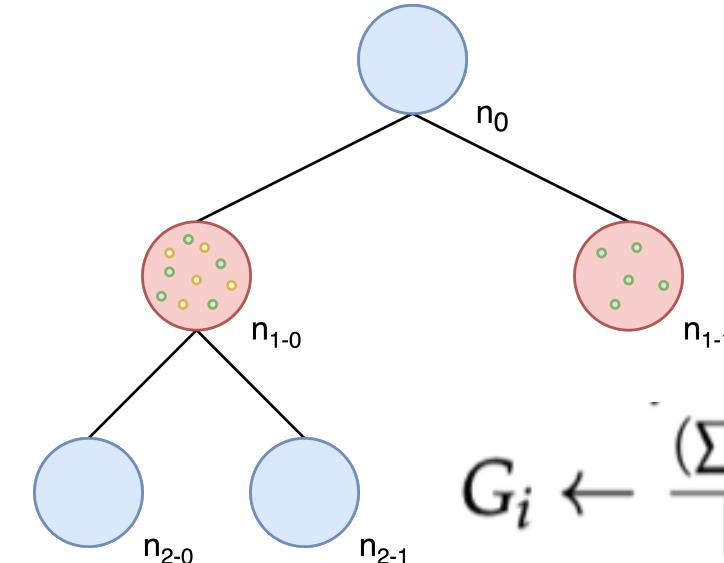
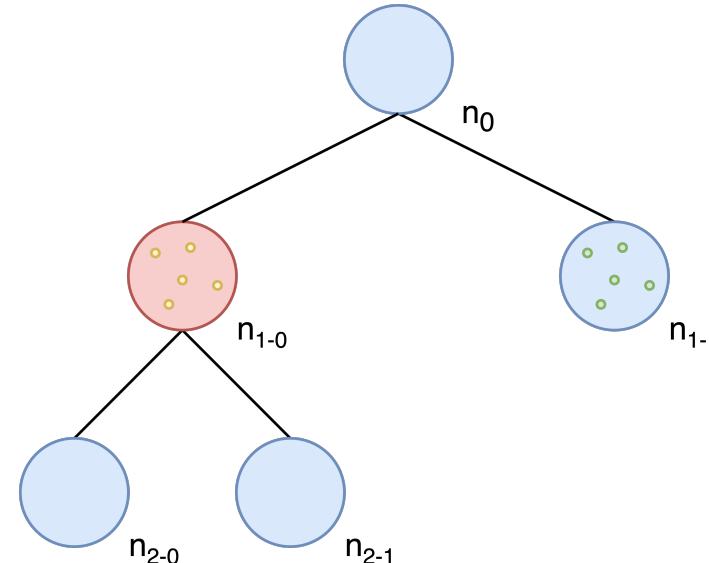
Vanilla



Differentially private

Improving DP-GBDT

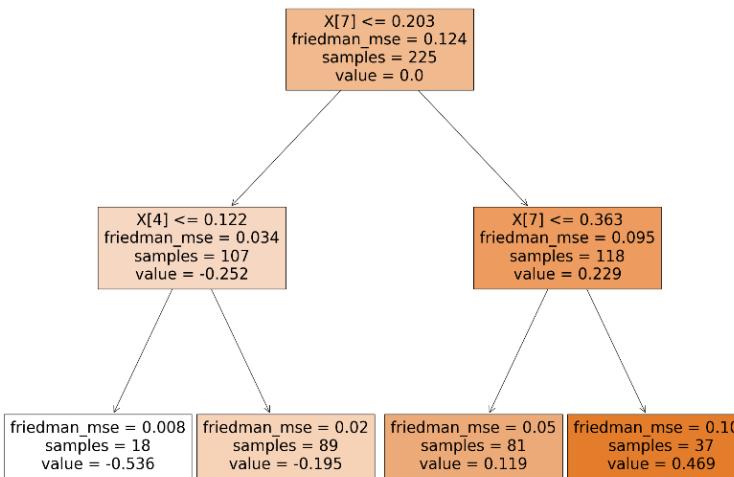
- How can we enhance the quality of the splits when there isn't much data available?
 - Look for more data!
 - ‘2-nodes’ algorithm: instead of splitting a node based on its samples only, split it based on its samples + the ones from its sibling



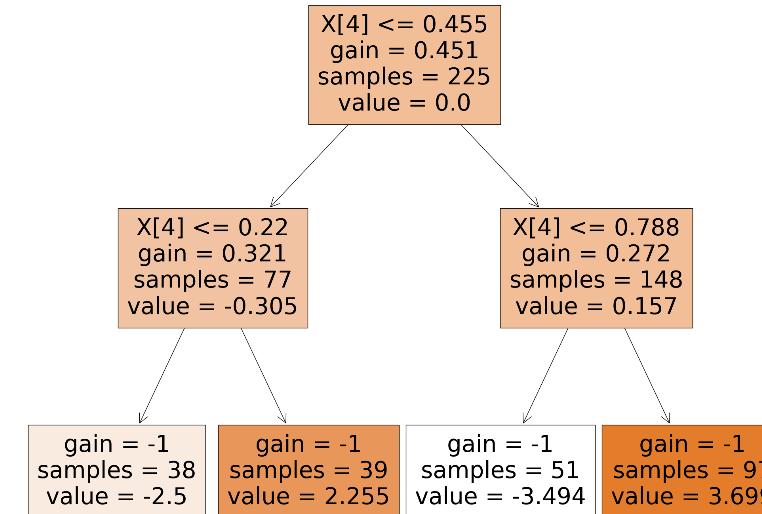
$$G_i \leftarrow \frac{(\sum_{i \in I_L} g_i)^2}{|I_L| + \lambda} + \frac{(\sum_{i \in I_R} g_i)^2}{|I_R| + \lambda}$$

Improving DP-GBDT

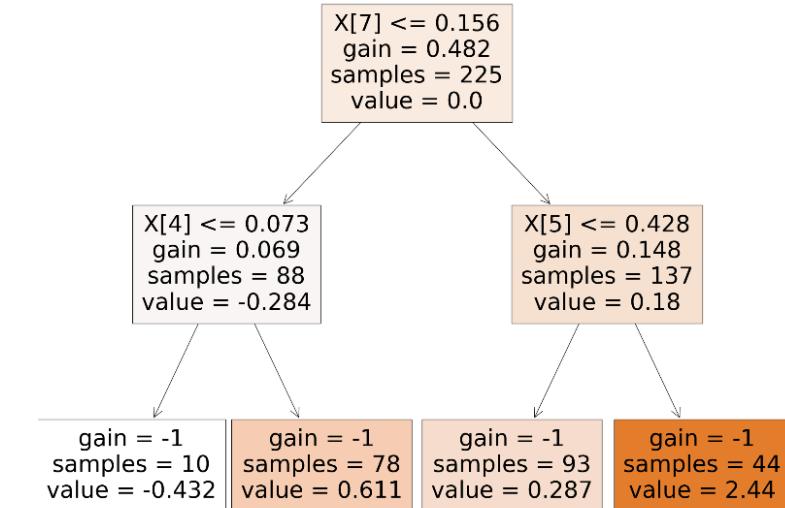
Dataset: abalone, samples = 300



Vanilla



Differentially private (e=1)



Differentially private - 2-nodes (e=1)

- What to look for:
 - Attribute selected for the split
 - Leaf value

Improving DP-GBDT

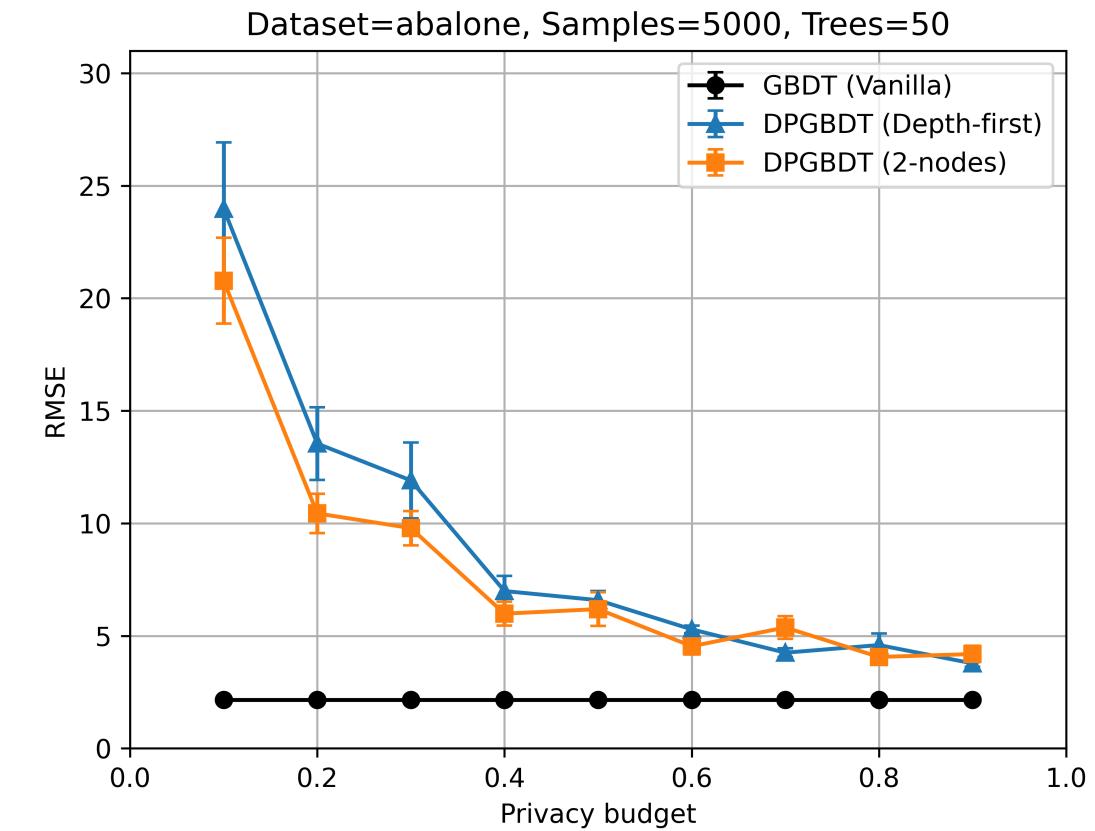
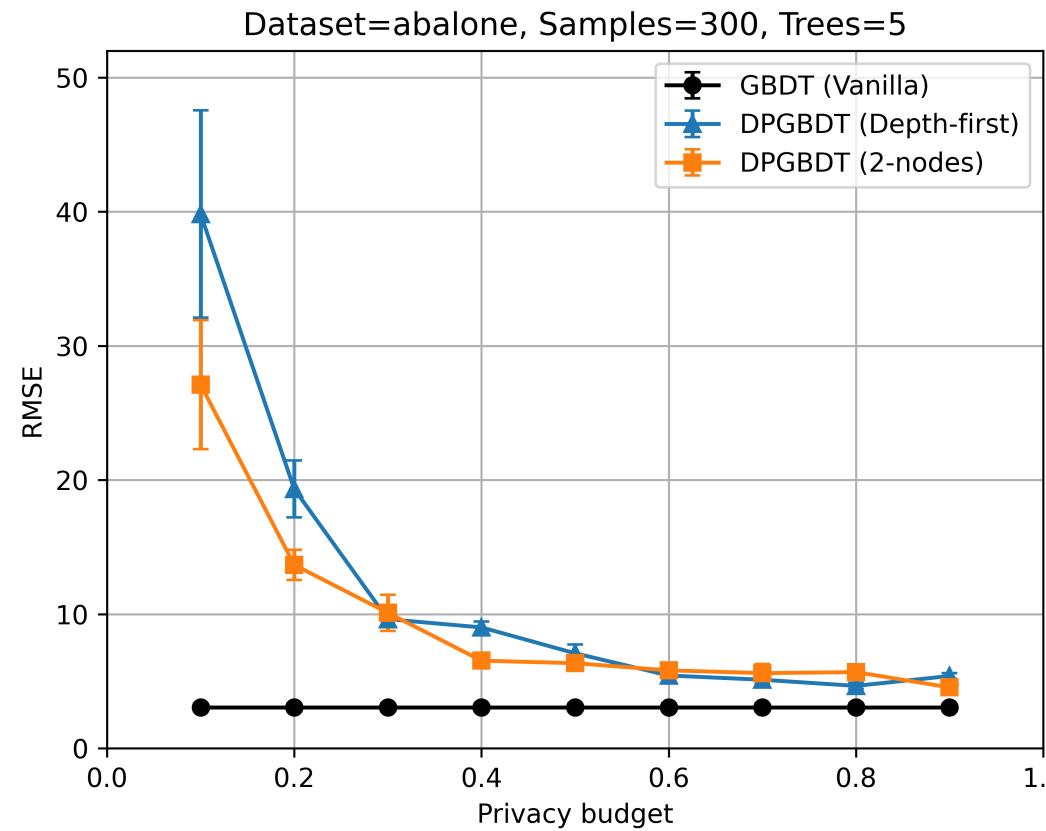
- Leaf values:



Improving DP-GBDT

- Why not look into more nodes?
 - More data querying, privacy budget would be consumed too rapidly
 - Nodes in other sub-trees may share complete different attributes compared to the direct sibling
 - Until that point in the tree, the node and its sibling share identical features
- Is it still differentially private?
 - We query the data twice as much compared to the original implementation, hence we only need to adapt our privacy budget consumption accordingly

Improving DP-GBDT



Contribution 2: Synthetic Data Generation

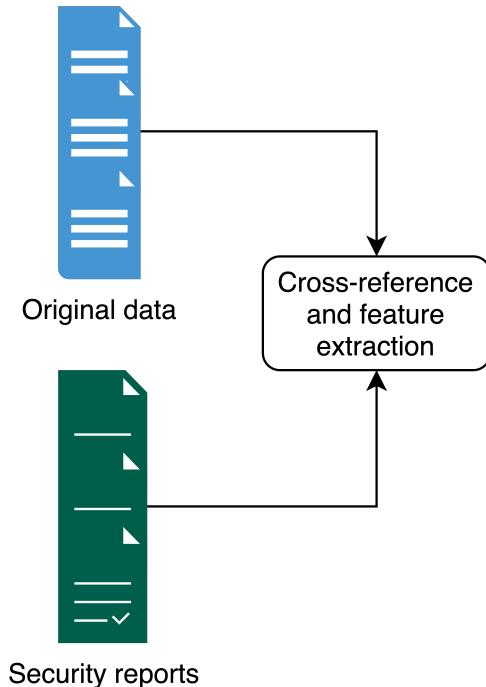
Synthetic Data Generation

- Real data only available in small amount, so we have to extrapolate them
 - 10 samples provided
 - Several thousands needed to evaluate the model

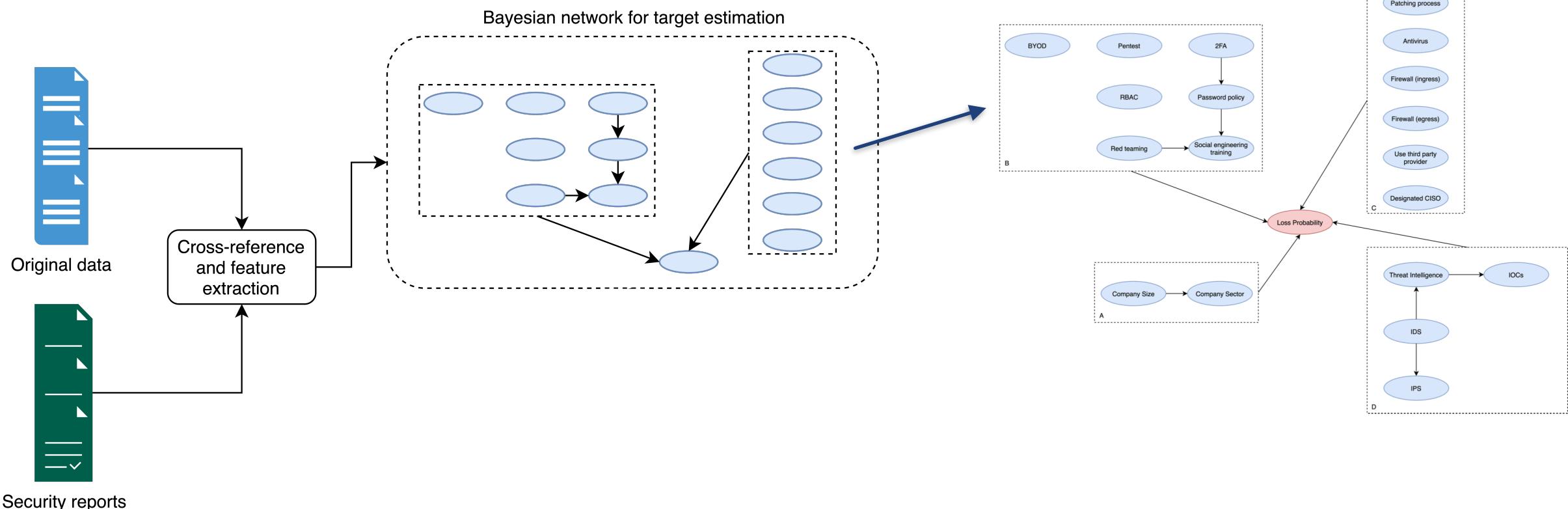
Synthetic Data Generation

- Real data only available in small amount, so we have to extrapolate them
 - 10 samples provided
 - Several thousands needed to evaluate the model
- Goal: build a dataset which:
 - features mimic questionnaire answers
 - targets mimic
 - 1) the probability for a customer of suffering a loss
 - 2) the potential cost involved if a customer has a security incident
 - is sufficiently complex for evaluation

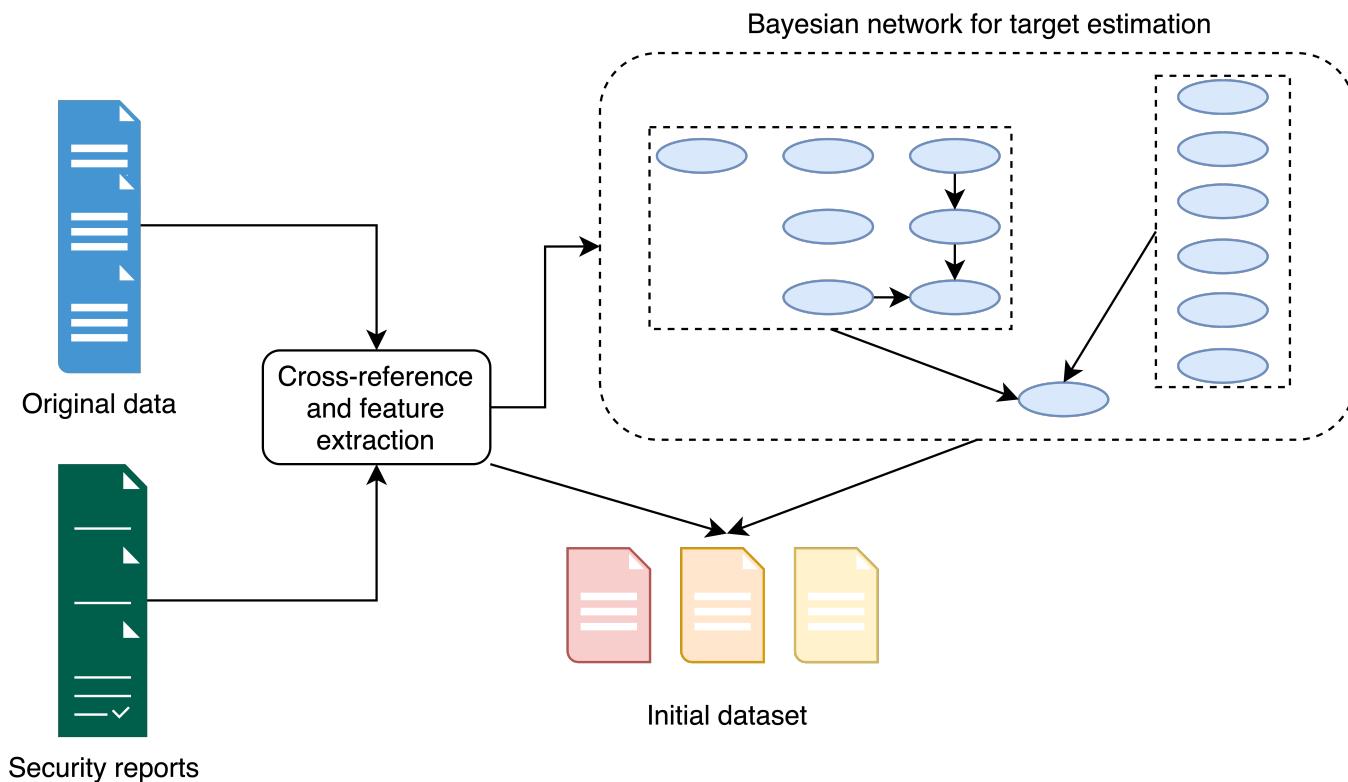
Synthetic Data Generation



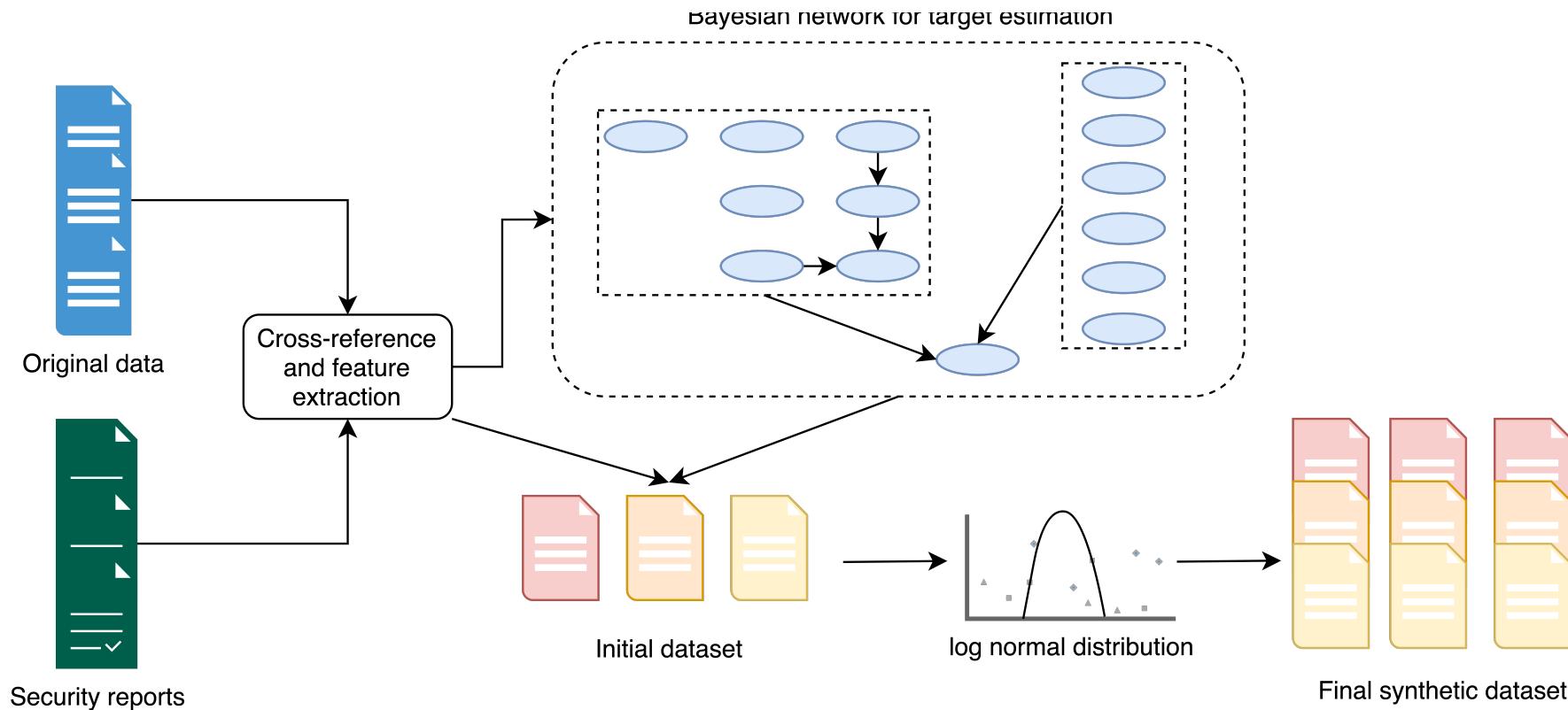
Synthetic Data Generation



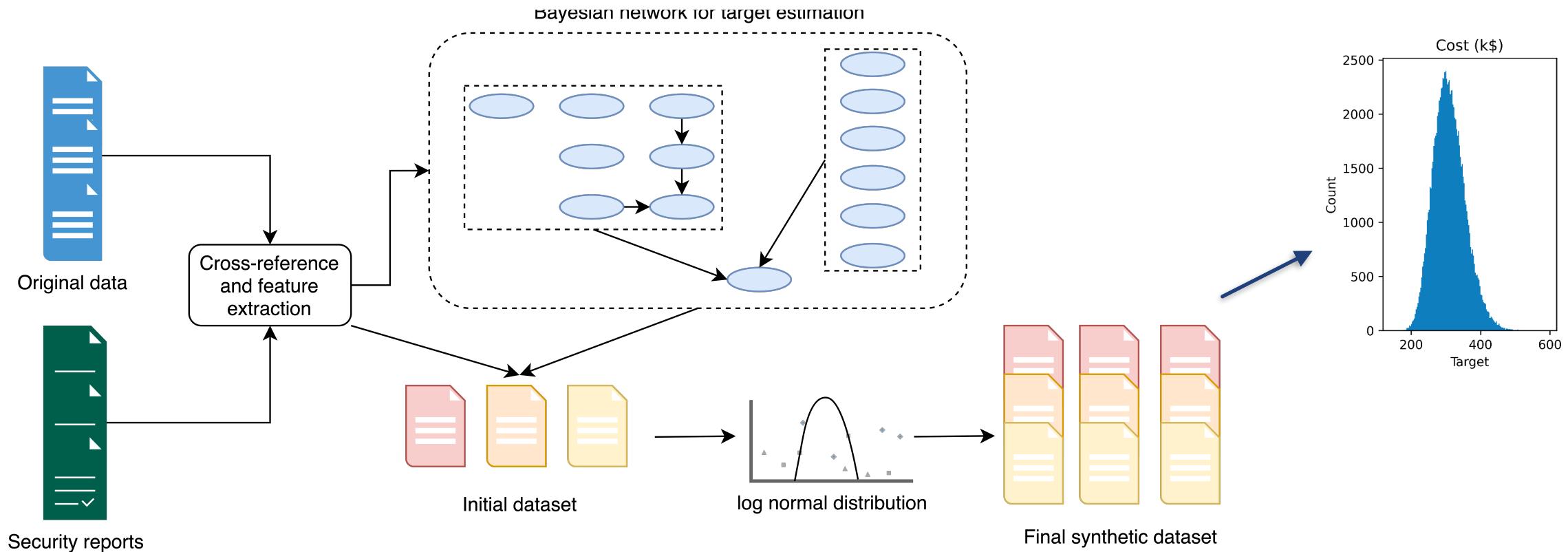
Synthetic Data Generation



Synthetic Data Generation



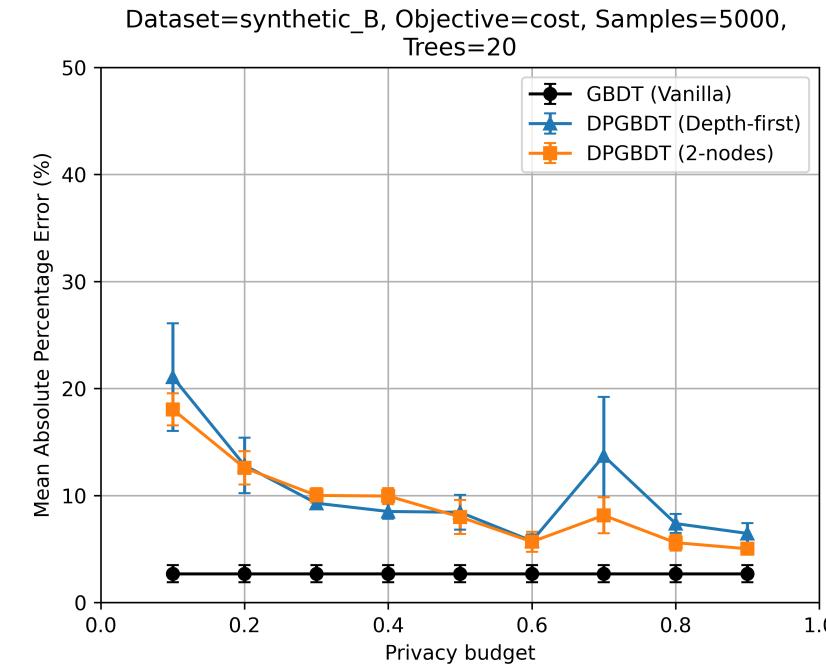
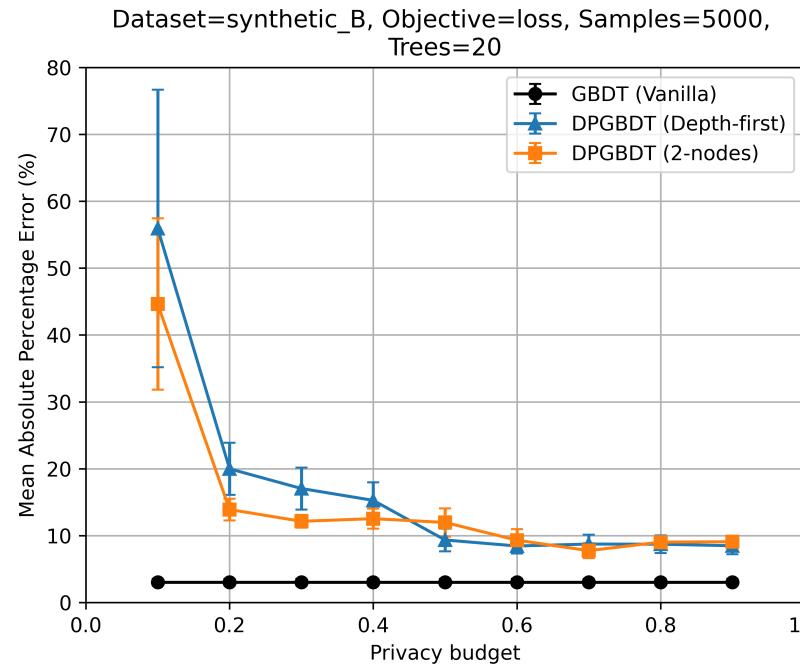
Synthetic Data Generation



Contribution 3: Accuracy and Privacy Evaluation

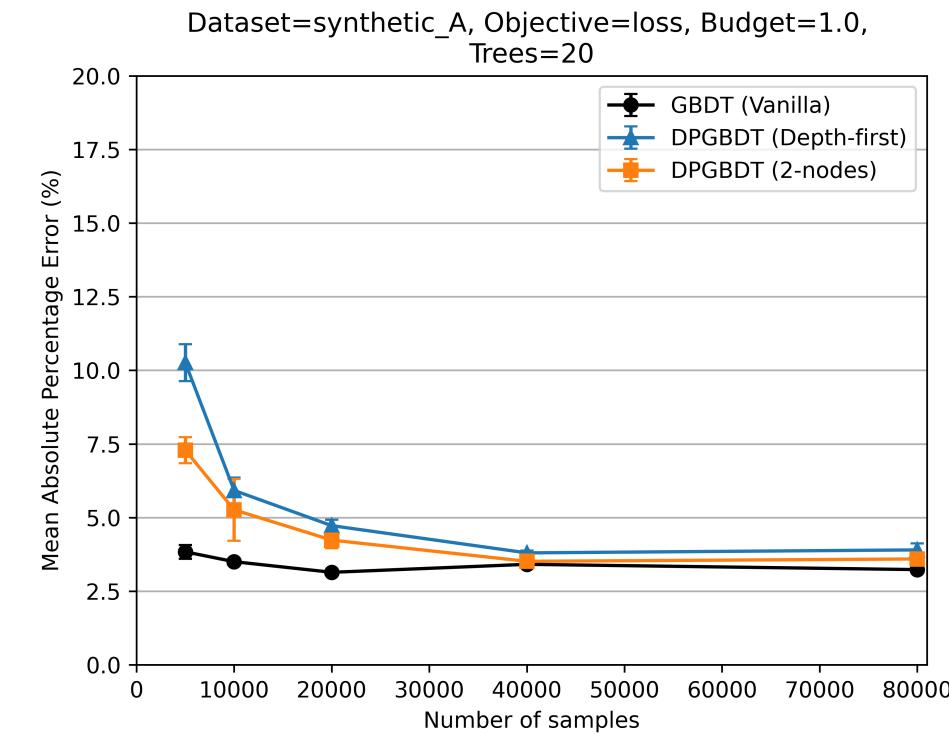
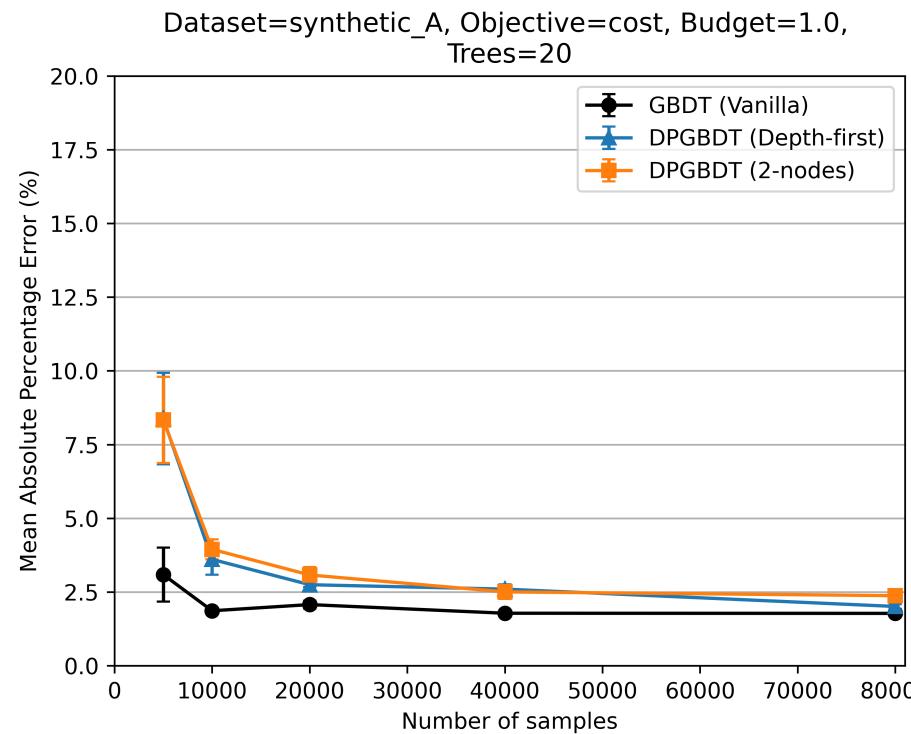
Evaluation

- For the evaluation, we report the Mean Average Percentage Error (MAPE) for each target.
- E.g.: if the true value is a 100, and the predicted value is 110, then the MAPE score would be 10%.



Evaluation

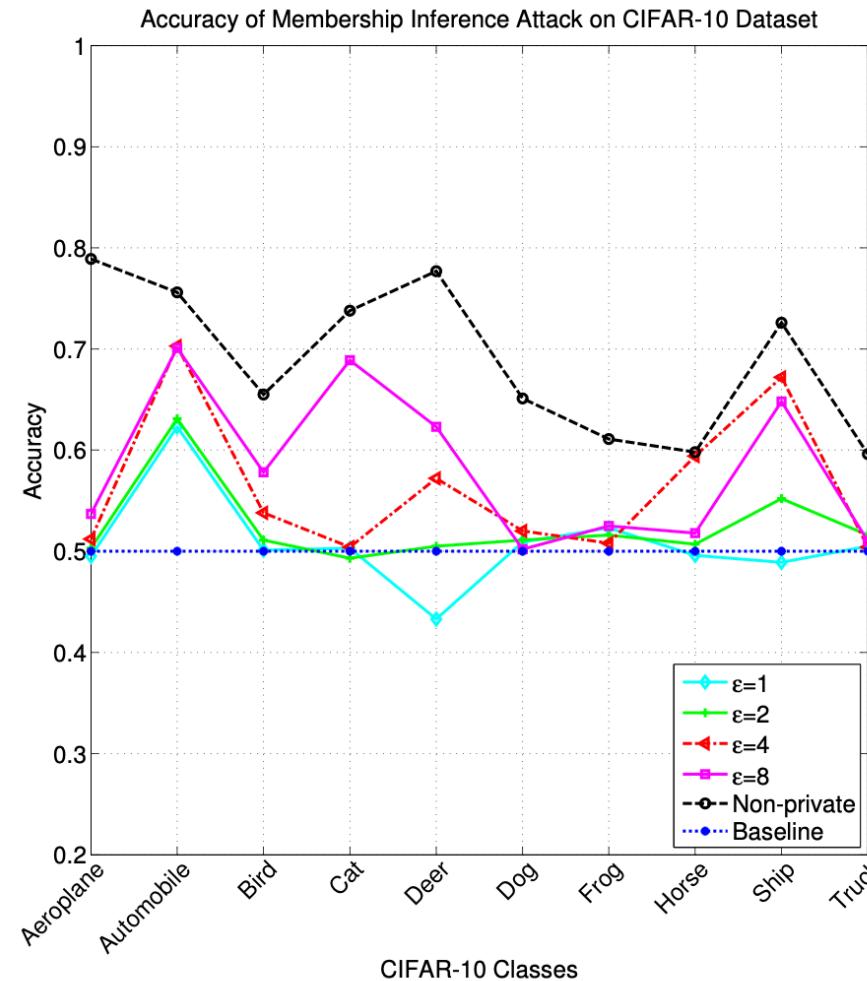
- The more samples, the better the model can approach non-DP results



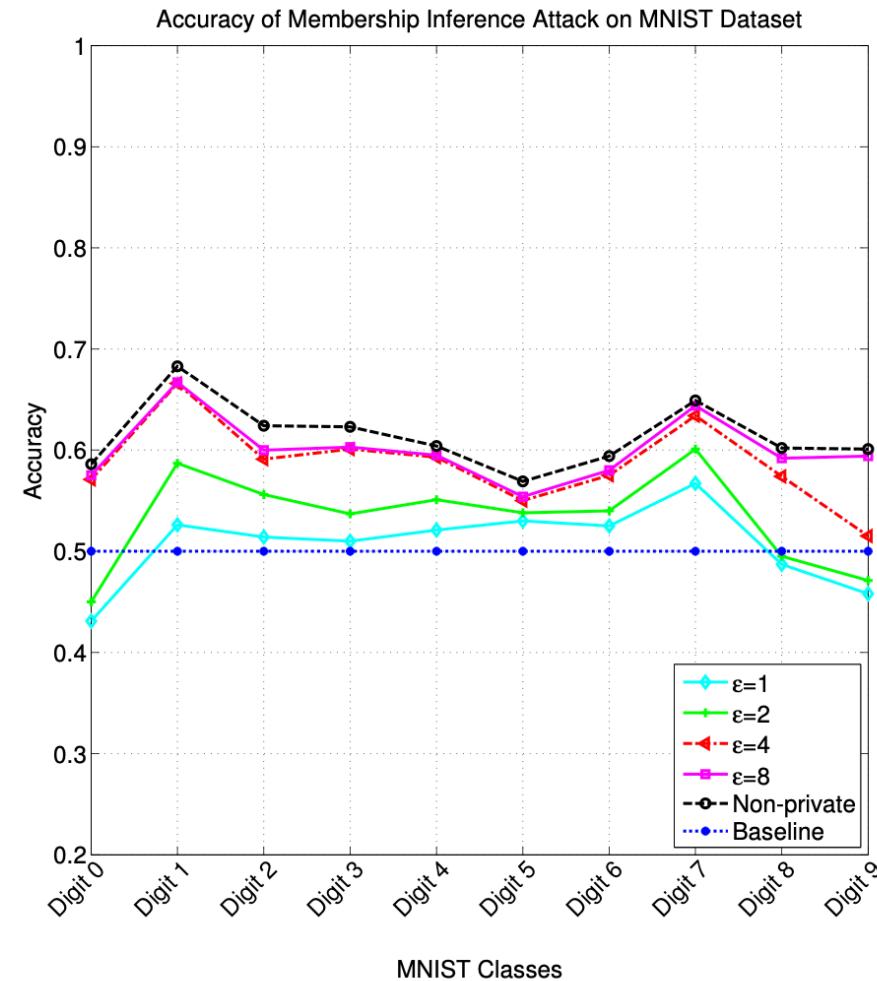
Privacy Attacks

- Example attacks:
 - Model inversion attack: attacker tries to extract data points from the model
 - Very hard
 - **Membership Inference Attack (MIA)**: attacker tries to identify if a given instance was part of the training set or not
 - Easier from the attacker's perspective, harder to defend against from the defender's perspective
 - We want to achieve 50% attack accuracy (i.e. random guess, no effectiveness)

Privacy Attacks



(a)



(b)

Conclusion

- We showed that state of the art DP can be applied to a specific business problem and provide good results
 - Accuracy enhanced for small datasets
- Trade-off accuracy vs. size of the dataset and privacy constraints
- There's a clear need for proper privacy attacks that target regression models

Future Work & Research Direction

- How to deploy such system
 - Retraining the model will ultimately break privacy, as querying the data again is a budget-consuming task
 - Enclave needs to be data-aware to block re-training
 - SGX side-channel resilience
- How to properly evaluate the model's privacy preserving features
 - Develop new attacks. Membership inference might be too strong for Zurich, but perhaps other attacks such as attribute inference attack could be developed
- GDF raises questions
 - Bug in Qinbin et al.'s proof found post-submission
 - Preliminary research indicates that this could be fixed with a cost of $\epsilon * 1.5$



Q&A

Appendix - Feature extraction

Feature (F)	$P(F)$	$P(F loss)$
<i>company_size</i>	small: 0.68 large: 0.32	
<i>company_sector</i>	See Table 5.1	
<i>company Allows BYOD</i>	0.44	0.14
<i>company Does pentest</i>	0.36	0.23
<i>company Uses RBAC</i>	0.85	0.67
<i>company Does red teaming</i>	0.33	
<i>company Has 2FA</i>	0.78	
<i>company Has password policy</i>	0.88	
<i>company Trains employees against social engineering</i>	0.70	
<i>company Has antivirus</i>	0.95	0.81
<i>company Has patching process</i>	0.60	0.828
<i>company Has ingress firewall</i>	0.83	0.84
<i>company Has egress firewall</i>	0.61	0.76
<i>company Uses 3rd party</i>	0.87	0.79
<i>company Has ciso</i>	0.67	0.68
<i>company Has threat intelligence team</i>	0.38	
<i>company Monitors IOCs</i>	0.44	
<i>company Has IDS</i>	0.94	
<i>company Has IPS</i>	0.81	
<i>company Separates systems</i>	0.79	
<i>company Does daily backup</i>	0.91	
<i>company Has recovery plan</i>	0.33	
<i>company Has incident response team</i>	0.46	

Feature (F)	$Pr(F)$	$Pr(F loss)$	$Pr(F SME)$	$Pr(F large)$
<i>company_is_in_hospitality</i>			0.03	0.04
<i>company_is_a_public_entity</i>			0.03	
<i>company_is_in_other</i>			0.11	0.09
<i>company_is_a_nonprofit</i>			0.05	
<i>company_is_in_education</i>			0.05	0.08
<i>company_is_in_technology</i>			0.06	0.03
<i>company_is_in_manufacturing</i>			0.08	0.03
<i>company_is_in_financial_services</i>			0.09	0.15
<i>company_is_in_energy</i>				0.04
<i>company_is_in_retail</i>			0.09	0.24
<i>company_is_in_healthcare</i>			0.19	0.26

Sector (S)	$Pr(PII S)$	$Pr(PCI S)$	$Pr(PHI S)$	$Pr(credentials S)$	$Pr(other S)$
<i>Hospitality</i>	0.44			0.34	0.23
<i>Public Entity</i>	0.51			0.33	0.34
<i>Other</i>	0.81			0.36	0.42
<i>Education</i>	0.75			0.30	0.23
<i>Technology</i>	0.69			0.41	0.34
<i>Manufacturing</i>	0.49	0.20		0.55	0.25
<i>Financial Services</i>	0.77	0.32		0.35	0.35
<i>Energy</i>	0.41	0.68		0.41	0.35
<i>Retail</i>	0.49	0.47		0.27	0.25
<i>Healthcare</i>	0.77		0.67	0.18	0.18
<i>Professional Services</i>	0.75			0.45	0.32

Appendix - Privacy Attacks

- Problems:
 - Current research work only focuses on classification models
 - Doesn't apply to our task
 - Doesn't apply to our model
 - Adapting the attacks to regression models is non-trivial
 - They rely on the class-probability vector outputted by classification models, which is lacking in regression decision trees
 - Some assumptions about the attacker might not hold
 - We can only look at the impact of DP in classification models from the literature