



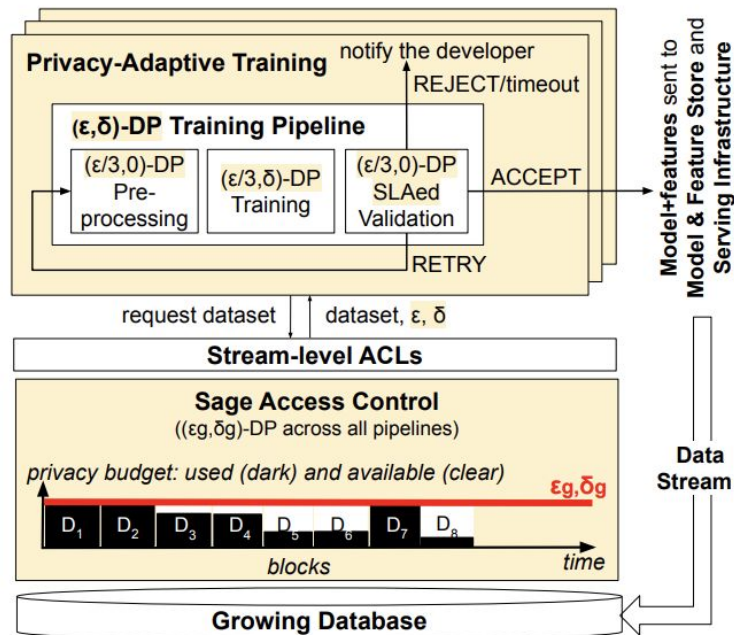
Privacy Accounting and Quality Control in Sage

Whys is DP needed with ML?

- ML datasets could leak specifics about individual entries in their training sets.
- Prevent featurization of dataset
 - Membership inference
 - Reconstruction attacks

Q: Why can't you just train a ML model using PINQ?

Sage Access Control & privacy adaptive training



Leverages the idea that the growing database is not static but growing, keeps training models endlessly on sensitive data stream

Challenges

Privacy Utility trade-off:

- Less accurate results that fail to meet the quality targets more often than w/o DP.
- low -quality models whose validations succeed by chance.

Splitting the data

- User-Level: based on user ID
 - Use incrementing userID's, max stored
 - New blocks are only created when new users join
- Event-level: splitting on time
 - days , months, etc.

Taxi Example

```
1 def preprocessing_fn(inputs, epsilon):
2     dist_01 = tft.scale_to_0_1(inputs["distance"], 0, 100)
3     speed_01 = tft.scale_to_0_1(inputs["speed"], 0, 100)
4     hour_of_day_speed = group_by_mean
5     sage.dp_group_by_mean(
6         inputs["hour_of_day"], speed_01, 24, epsilon, 1.0)
7     return {"dist_scaled": dist_01,
8            "hour_of_day": inputs["hour_of_day"],
9            "hour_of_day_speed": hour_of_day_speed,
10           "duration": inputs["duration"]}
11
12 def trainer_fn(hparams, schema, epsilon, delta): [...]
13     feature_columns = [numeric_column("dist_scaled"),
14                       numeric_column("hour_of_day_speed"),
15                       categorical_column("hour_of_day", num_buckets=24)]
16     estimator = \
17         tf.estimator.DNNRegressor(sage.DPDNNRegressor(
18             config=run_config,
19             feature_columns=feature_columns,
20             dnn_hidden_units=hparams.hidden_units,
21             privacy_budget=(epsilon, delta))
22     return tfx.executors.TrainingSpec(estimator, ...)
23
24 def validator_fn(epsilon):
25     model_validator = \
26         tfx.components.ModelValidator(sage.DPModelValidator(
27             examples=examples_gen.outputs.output,
28             model=trainer.outputs.output,
29             metric_fn = _MSE_FN, target = _MSE_TARGET,
30             epsilon=epsilon, confidence=0.95, B=1)
31     return model_validator
32
33 def dp_group_by_mean(key_tensor, value_tensor, nkeys,
34                     epsilon, value_range):
35     key_tensor = tf.dtypes.cast(key_tensor, tf.int64)
36     ones = tf.fill(tf.shape(key_tensor), 1.0)
37     dp_counts = group_by_sum(key_tensor, ones, nkeys) \
38         + laplace(0.0, 2/epsilon, nkeys)
39     dp_sums = group_by_sum(
40         key_tensor, value_tensor, nkeys) \
41         + laplace(0.0, value_range * 2/epsilon, nkeys)
42     return tf.gather(dp_sums/dp_counts, key_tensor)
```

- Preprocessing_fn: makes aggregate features i.e distance of ride, hour of day
 - Dp_group_by_mean:
 - Number of times key appears
 - Sum of values associated w/ key
 - Each data point has one key

Sage Access Control : requirements for composition theory

- R1: Multiple training pipelines w/ differing amounts of data needed for performance
- R2: Adaptivity in choice of queries, DP parameters and data subsets
- R3: Some models are ran periodically w/ new data and others are retired

Failed Methods: which rules do these violate?

1. Query across the entire stream:
 - $\epsilon_d = \epsilon_1 + \epsilon_2 + \epsilon_3$
2. Queries split in to subqueries and each run DP on individual blocks, results aggregated
3. A new data point is allocated to one of the waiting queries, which consumes entire privacy budget.

Block Composition Theory cont.

- Splits data into disjoint blocks adaptively chosen($R1, R2$)
- Privacy loss of three queries will be max of $\epsilon1 + \epsilon2$, and $\epsilon2 + \epsilon3$
- New blocks $D5$ arrive w/ privacy loss of zero($R3$)

System can run endlessly by training new models on new data!

Q: What does it mean for
DP parameters to be
chosen Adaptively?

Adaptive Parameters

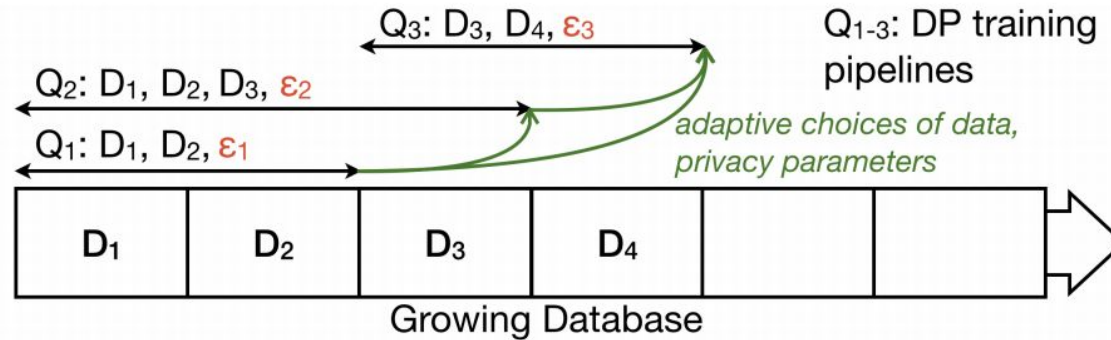


Fig. 3. Characteristics of Data Interaction in ML.

Privacy-Adaptive Training

- To improve DP quality:
 - Increase privacy budget(ϵ , δ) or increase dataset size
- Accept: prediction target reached
- Retry: more data needed for assessment
- Reject: model will never reach target w/ sample size/privacy requirements

Discuss:

Q: What assumptions are made about the data? In what cases could Sage potentially not perform well?