

Topik : 4.2. Integrasi TF Privacy ke Model

Objective : Implementasikan optimizer dengan DP-SGD

Task : Tambahkan dp_keras_optimizer ke model FL

Source : https://github.com/tensorflow/privacy/blob/master/tutorials/mnist_dpsgd_tutorial.py

```
PROBLEMS 8 OUTPUT DEBUG CONSOLE TERMINAL PORTS
PS C:\KP\MATERI\4.2 Integrasi TF Privacy ke Model> python -m venv venv
>>
PS C:\KP\MATERI\4.2 Integrasi TF Privacy ke Model> |
```

```
>>
PS C:\KP\MATERI\4.2 Integrasi TF Privacy ke Model> venv\Scripts\activate
(venv) PS C:\KP\MATERI\4.2 Integrasi TF Privacy ke Model> |
```

```
Successfully uninstalled tensorflow-estimator-2.15.0
(venv) PS C:\KP\MATERI\4.2 Integrasi TF Privacy ke Model> pip install "tensorflow==2.15.*" "tensorflow-estimator==2.15.*" "tensorflow-privacy==0.8.*" absl-py -U
>>
Collecting tensorflow==2.15.*
  Downloading tensorflow-2.15.1-cp310-cp310-win_amd64.whl.metadata (3.4 kB)
Collecting tensorflow-estimator==2.15.*
  Using cached tensorflow_estimator-2.15.0-py2.py3-none-any.whl.metadata (1.3 kB)
Collecting tensorflow-privacy==0.8.*
  Downloading tensorflow_privacy-0.8.12-py3-none-any.whl.metadata (962 bytes)
Requirement already satisfied: absl-py in c:\kp\materi\4.2 integrasi tf privacy ke model\venv\lib\site-packages (1.4.0)
Collecting absl-py
  Using cached absl_py-2.3.1-py3-none-any.whl.metadata (3.3 kB)
Collecting tensorflow-intel==2.15.1 (from tensorflow==2.15.*)
  Downloading tensorflow_intel-2.15.1-cp310-cp310-win_amd64.whl.metadata (4.9 kB)
```

1. Import dan Konfigurasi Awal

```
16 import time
17
18 from absl import app
19 from absl import flags
20 from absl import logging
21 import tensorflow as tf
22 tf.compat.v1.disable_eager_execution()
23 from tensorflow_estimator import estimator as tf_estimator
24
25 from tensorflow_privacy.privacy.analysis import compute_dp_sgd_privacy_lib
26 from tensorflow_privacy.privacy.optimizers import dp_optimizer
27 import mnist_dpsgd_tutorial_common as common
```

Mengaktifkan graph mode, menyiapkan estimator, DP-SGD, dan modul common yang berisi get_cnn_model() & make_input_fn()

2. Hyperparameter via absl.flags

```
29 flags.DEFINE_boolean(
30     'dpsgd', True, 'If True, train with DP-SGD. If False, '
31     'train with vanilla SGD.')
32 flags.DEFINE_float('learning_rate', .15, 'Learning rate for training')
33 flags.DEFINE_float('noise_multiplier', 1.1,
34     'Ratio of the standard deviation to the clipping norm')
35 flags.DEFINE_float('l2_norm_clip', 1.0, 'Clipping norm')
36 flags.DEFINE_integer('batch_size', 256, 'Batch size')
37 flags.DEFINE_integer('epochs', 30, 'Number of epochs')
38 flags.DEFINE_integer(
39     'microbatches', 256, 'Number of microbatches '
40     '(must evenly divide batch_size)')
41 flags.DEFINE_string('model_dir', None, 'Model directory')
42
43 FLAGS = flags.FLAGS
```

semua parameter pelatihan & privasi dikontrol dari sini.

3. Model function untuk melakukan estimator (cnn_model_fn)

```

46 def cnn_model_fn(features, labels, mode, params): # pylint: disable=unused-argument
47     """Model function for a CNN."""
48
49     # Define CNN architecture.
50     logits = common.get_cnn_model(features)
51
52     # Calculate loss as a vector (to support minibatches in DP-SGD).
53     vector_loss = tf.nn.sparse_softmax_cross_entropy_with_logits(
54         labels=labels, logits=logits)
55     # Define mean of loss across minibatch (for reporting through tf.Estimator).
56     scalar_loss = tf.reduce_mean(input_tensor=vector_loss)
57
58     # Configure the training op (for TRAIN mode).
59     if mode == tf.estimator.ModeKeys.TRAIN:
60         if FLAGS.dpsgd:
61             # Use DP version of GradientDescentOptimizer. Other optimizers are
62             # available in dp_optimizer. Most optimizers inheriting from
63             # tf.compat.v1.train.Optimizer should be wrappable in differentially
64             # private counterparts by calling dp_optimizer.optimizer_from_args().
65             optimizer = dp_optimizer.DPGradientDescentGaussianOptimizer(
66                 l2_norm_clip=FLAGS.l2_norm_clip,
67                 noise_multiplier=FLAGS.noise_multiplier,
68                 num_microbatches=FLAGS.microbatches,
69                 learning_rate=FLAGS.learning_rate)
70             opt_loss = vector_loss
61         else:
72             optimizer = tf.compat.v1.train.GradientDescentOptimizer(
73                 learning_rate=FLAGS.learning_rate)
74             opt_loss = scalar_loss
75
76         global_step = tf.compat.v1.train.get_global_step()
77         train_op = optimizer.minimize(loss=opt_loss, global_step=global_step)
78
79         # In the following, we pass the mean of the loss (scalar_loss) rather than
80         # the vector_loss because tf.estimator requires a scalar loss. This is only
81         # used for evaluation and debugging by tf.estimator. The actual loss being
82         # minimized is opt_loss defined above and passed to optimizer.minimize().
83         return tf.estimator.EstimatorSpec(
84             mode=mode, loss=scalar_loss, train_op=train_op)
85
86     # Add evaluation metrics (for EVAL mode).
87     elif mode == tf.estimator.ModeKeys.EVAL:
88         eval_metric_ops = {
89             'accuracy':
90                 tf.compat.v1.metrics.accuracy(
91                     labels=labels, predictions=tf.argmax(input=logits, axis=1))
92         }
93         return tf.estimator.EstimatorSpec(
94             mode=mode, loss=scalar_loss, eval_metric_ops=eval_metric_ops)
95

```

- **DP-SGD**: clipping L2 + Gaussian noise → wajib `opt_loss = vector_loss`.
- **Non-DP**: `opt_loss = scalar_loss`.
- Metrik yang dilaporkan saat evaluasi: **akurasi**.

4. Fungsi main: validasi, Estimator, loop train/eval, dan hitung ϵ

```

97 def main(unused_argv):
98     logging.set_verbosity(logging.INFO)
99     if FLAGS.dpsgd and FLAGS.batch_size % FLAGS.microbatches != 0:
100         raise ValueError('Number of microbatches should divide evenly batch_size')
101
102     # Instantiate the tf.Estimator.
103     mnist_classifier = tf.estimator.Estimator(
104         model_fn=cnn_model_fn, model_dir=FLAGS.model_dir)
105
106     # Training loop.
107     steps_per_epoch = 60000 // FLAGS.batch_size
108     for epoch in range(1, FLAGS.epochs + 1):
109         start_time = time.time()
110         # Train the model for one epoch.
111         mnist_classifier.train(
112             input_fn=common.make_input_fn('train', FLAGS.batch_size),
113             steps=steps_per_epoch)
114         end_time = time.time()
115         logging.info('Epoch %d time in seconds: %.2f', epoch, end_time - start_time)
116
117     # Evaluate the model and print results
118     eval_results = mnist_classifier.evaluate(
119         input_fn=common.make_input_fn('test', FLAGS.batch_size, 1))
120     test_accuracy = eval_results['accuracy']
121     print('Test accuracy after %d epochs is: %.3f' % (epoch, test_accuracy))
122
123     # Compute the privacy budget expended.
124     if FLAGS.dpsgd:
125         if FLAGS.noise_multiplier > 0.0:
126             eps, _ = compute_dp_sgd_privacy_lib.compute_dp_sgd_privacy(
127                 60000, FLAGS.batch_size, FLAGS.noise_multiplier, epoch, 1e-5)
128             print('For delta=1e-5, the current epsilon is: %.2f' % eps)
129         else:
130             print('Trained with DP-SGD but with zero noise.')
131     else:
132         print('Trained with vanilla non-private SGD optimizer')
133

```

- **Validasi**: `batch_size` harus habis dibagi `microbatches`.
- **Train/Eval per-epoch**.
- **Epsilon (ϵ)** dihitung kumulatif per-epoch → trade-off akurasi vs privasi.

5. Entry point

```
134
135     if __name__ == '__main__':
136         app.run(main)
137
```

menjalankan main() dan mengurai flag dari command line.

Ringkasan :

- **Estimator** memanggil `cnn_model_fn` untuk **TRAIN** (pilih DP-SGD vs SGD) dan **EVAL** (akurasi).
- **DP-SGD**: pakai **loss per-contoh** (`vector_loss`), **clipping L2** (`l2_norm_clip`), **noise Gaussian** (`noise_multiplier`), dan **microbatches**.
- Setiap epoch: **latih** → **uji** → **cetak akurasi**, dan kalau DP aktif **hitung ϵ** (semakin kecil = privasi lebih kuat).

Source : https://www.tensorflow.org/responsible_ai/privacy/tutorials/classification_privacy

Implement Differential Privacy with Tensorflow Privacy

Privacy differensial (DP) Adalah sebuah kerangka kerja untuk mengukur jaminan privasi yang diberikan oleh sebuah algoritma. Melalui lensa privasi differensial , kita dapat merancang algoritma machine learning yang dapat melatih model secara bertanggung jawab pada data pribadi. Pembelajaran dengan privasi diferensial memberikan jaminan privasi yang terukur, membantu mengurangi risiko terbukanya data sensitif dalam pembelajaran mesin. Secara intuitif, model yang dilatih dengan privasi diferensial tidak boleh dipengaruhi oleh satu contoh pelatihan tunggal, atau kumpulan kecil contoh pelatihan, dalam himpunan datanya. Hal ini membantu mengurangi risiko terbukanya data sensitif dalam ML.

Ide dasar dari pendekatan ini, yang disebut *differentially private stochastic gradient descent* (DP-SGD), Adalah melakukan modifikasi gradien yang akan digunakan dalam *stochastic gradient descent* (SGD), yang merupakan inti dari hampir semua algoritma deep learning. Model yang dilatih dengan DP-SGD memberikan jaminan privasi diferensial yang terbukti untuk data input mereka.

2 modifikasi yang dilakukan pada algoritma SGD standar :

1. Sensitivitas setiap gradien perlu dibatasi

Dengan kata lain, kita perlu membatasi sejauh mana setiap titik data pelatihan yang diambil dalam sebuah minibatch dapat memengaruhi perhitungan gradien dan pembaruan yang diterapkan pada parameter model. Hal ini dilakukan dengan *clipping* (memotong) setiap gradient yang dihitung pada setiap titik pelatihan.

2. Noise acak

Diambil sampelnya dan ditambahkan pada gradien yang sudah dipotong, untuk membuatnya secara statistic mustahil mengetahui apakah suatu titik data tertentu termasuk atau tidak dalam dataset pelatihan, dengan cara membandingkan pembaruan yang dilakukan oleh SGD Ketika beroperasi dengan atau tanpa titik data tersebut dalam dataset pelatihan.

Tutorial ini menggunakan tf.keras untuk melatih *convolutional neural network (CNN)* dalam mengenali digit tulisan tangan dengan optimizer DP-SGD yang disediakan oleh Pustaka Tensorflow Privacy. Tensorflow Privacy menyediakan kode untuk membungkus optimizer Tensorflow yang sudah ada untuk membuat varian yang mengimplementasikan DP-SGD.

1. SetUp

```
1 ##### SETUP #####
2 import tensorflow as tf
3
4
5
6 import numpy as np
7 tf.get_logger().setLevel('ERROR')
8
9 import tensorflow_privacy
10 from tensorflow_privacy.privacy.analysis import compute_dp_sgd_privacy
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

```
Successfully installed tensorflow-probability-0.19.0
(.venv) PS C:\KIP\WATERIV4.2 Integrasi TF Privacy ke Model> pip uninstall -y tensorflow tensorflow-intel tensorflow-estimator tensorflow-privacy tensorflow-probability
>> pip install tensorflow==2.15.0 tensorflow-estimator==2.15.0 tensorflow-privacy==0.9.0 tensorflow-probability==0.22.1
>>
Found existing installation: tensorflow 2.12.0
```

2. Load and pre-process the dataset

Muat dataset MNIST dan siapkan data untuk pelatihan.

```
13 ##### Load and pre-process the dataset #####
14
15 train, test = tf.keras.datasets.mnist.load_data()
16 train_data, train_label = train
17 test_data, test_label = test
18
19 train_data = np.array(train_data, dtype=np.float32) / 255
20 test_data = np.array(test_data, dtype=np.float32) / 255
21
22 train_data = train_data.reshape(train_data.shape[0], 28, 28, 1)
23 test_data = test_data.reshape(test_data.shape[0], 28, 28, 1)
24
25 train_label = np.array(train_label, dtype=np.int32)
26 test_label = np.array(test_label, dtype=np.int32)
27
28 train_label = tf.keras.utils.to_categorical(train_label, num_classes=10)
29 test_label = tf.keras.utils.to_categorical(test_label, num_classes=10)
30
31 assert train_data.min() == 0
32 assert train_data.max() == 1
33 assert test_data.min() == 0
34 assert test_data.max() == 1
```

```
(.venv) PS C:\KP\MATERI\4.2 Integrasi TF Privacy ke Model> python classification_privacy.py
2025-09-10 09:10:35.255575: I tensorflow/core/util/port.cc:113] oneDNN custom operations are on. You
_ONEDNN_OPTS=0*.
WARNING:tensorflow:From C:\KP\MATERI\4.2 Integrasi TF Privacy ke Model\venv\Lib\site-packages\keras
Downloading data from https://storage.googleapis.com/tensorflow/tf-keras-datasets/mnist.npz
11490434/11490434 [=====] - 4s 0us/step
(.venv) PS C:\KP\MATERI\4.2 Integrasi TF Privacy ke Model> █
```

3. Define the hyperparameters

Tetapkan nilai hiperparameter model pembelajaran.

```
37 #####
38 epochs = 3
39 batch_size = 250
```

DP-SGD memiliki tiga hyperparameter khusus privasi dan satu hyperparameter yang sudah ada yang harus anda atur :

1. l2_norm_clip (float)

Nilai maksimum norma Euclidean (L2) dari setiap gradient yang diterapkan untuk memperbarui parameter model. Hyperparameter ini digunakan untuk membatasi sensitivitas optimizer terhadap titik data pelatihan individu.

2. noise_multiplier (float)

Jumlah noise yang diambil sampelnya dan ditambahkan ke gradien selama pelatihan. Secara umum, semakin banyak noise maka semakin menghasilkan privasi yang lebih baik (meskipun tidak selalu dengan mengorbankan utilitas yang lebih rendah)

3. microbatches (int)

Setiap batch data dibagi ke dalam unit – unit lebih kecil yang disebut microbatch. Secara default, setiap microbatch sebaiknya hanya berisi satu contoh pelatihan. Hal ini memungkinkan kita untuk melakukan clipping gradien pada basis per-contoh, bukan setelah rata – rata di seluruh minibatch. Dengan demikian, efek (negatif) clipping pada sinyal yang ditemukan di gradien berkurang dan biasanya memaksimalkan utilitas. Namun, overhead komputasi dapat dikurangi dengan memperbesar ukuran microbatch sehingga mencakup lebih dari satu contoh pelatihan. Rata – rata gradien di seluruh contoh pelatihan ini kemudian di-clipping. Jumlah total contoh yang digunakan dalam satu batch (satu Langkah *gradient descent*) tetap sama. Jumlah microbatch harus membagi ukuran batch secara merata.

4. learning_rate (float)

Hyperparameter ini sudah ada dalam vanilla SGD. Semakin tinggi nilai learning rate, semakin besar pengaruh setiap pembaruan. Jika pembaruan sangat bising (misalnya saat noise aditif besar dibandingkan ambang clipping), learning rate yang rendah dapat membantu prosedur pelatihan untuk konvergen.

Rangkuman :

- `l2_norm_clip` → membatasi gradien biar tidak ada satu data mendominasi.
- `noise_multiplier` → mengatur seberapa besar noise ditambahkan demi privasi.
- `microbatches` → membagi batch jadi unit kecil supaya clipping lebih akurat.
- `learning_rate` → mengatur kecepatan update (harus seimbang dengan noise & clipping).

Gunakan nilai hiperparameter di bawah ini untuk mendapatkan model yang cukup akurat (akurasi pengujian 95%):

```
41 l2_norm_clip = 1.5
42 noise_multiplier = 1.3
43 num_microbatches = 250
44 learning_rate = 0.25
45
46 if batch_size % num_microbatches != 0 :
47     raise ValueError('Batch size should be an integer multiple of the number of microbatches')
48
```

4. Build the model

Tentukan jaringan saraf konvolusional sebagai model pembelajaran.

```
53 ##### Build the model #####
54 model = tf.keras.Sequential([
55     tf.keras.layers.Conv2D(16, 8, strides=2, padding='same', activation='relu', input_shape=(28,28,1)),
56     tf.keras.layers.MaxPool2D(2,1),
57     tf.keras.layers.Conv2D(32,4, strides=2, padding='valid', activation='relu'),
58     tf.keras.layers.MaxPool2D(2,1),
59     tf.keras.layers.Flatten(),
60     tf.keras.layers.Dense(32, activation='relu'),
61     tf.keras.layers.Dense(10)
62 ])
```

Tentukan pengoptimal dan fungsi kerugian untuk learning model. Hitung kerugian sebagai vector kerugian per-contoh, alih – alih sebagai rata – rata pada minibatch untuk mendukung manipulasi gradien pada setiap titik pelatihan.

```
63
64 optimizer = tensorflow_privacy.DPKerasSGDOptimizer(
65     l2_norm_clip=l2_norm_clip,
66     noise_multiplier=noise_multiplier,
67     num_microbatches=num_microbatches,
68     learning_rate=learning_rate
69 )
70 loss = tf.keras.losses.CategoricalCrossentropy(
71     from_logits = True, reduction=tf.losses.Reduction.NONE
72 )
```

5. Train the model

```
74
75
76 ##### Train the model #####
77 model.compile(optimizer=optimizer, loss=loss, metrics=['accuracy'])
78 model.fit(train_data,
79         train_label,
80         epochs=epochs,
81         validation_data=(test_data, test_label),
82         batch_size=batch_size)
```

```

ags.
Epoch 1/3
240/240 [=====] - 1115s 5s/step - loss: 0.9647 - accuracy: 0.7020 - val_loss: 0.4016 - val_accuracy: 0.8855
Epoch 2/3
240/240 [=====] - 633s 3s/step - loss: 0.3895 - accuracy: 0.8970 - val_loss: 0.3299 - val_accuracy: 0.9219
Epoch 3/3
240/240 [=====] - 610s 3s/step - loss: 0.3546 - accuracy: 0.9188 - val_loss: 0.3235 - val_accuracy: 0.9296
Traceback (most recent call last):

```

6. Measure the differential privacy guarantee

Melakukan analisis privasi untuk mengukur jaminan DP yang dicapai oleh algoritma pelatihan. Mengetahui Tingkat DP yang dicapai memungkinkan perbandingan objektif dari dua proses pelatihan untuk menentukan mana yang lebih menjaga privasi.

Analisis privasi mengukur seberapa jauh seorang lawan potensial dapat meningkatkan tebakan mereka tentang property titik data individu dengan mengamati hasil dari prosedur pelatihan (misalnya, pembaruan model dan parameter).

Jaminan ini kadang disebut sebagai anggaran privasi (privacy budget). Anggaran privasi yang lebih rendah membatasi lebih ketat kemampuan lawan untuk meningkatkan tebakan mereka. Hal ini memastikan jaminan privasi yang lebih kuat. Secara intuitif, ini karena lebih sulit bagi satu titik data pelatihan untuk memengaruhi hasil pembelajaran : misalnya, informasi dalam titik data pelatihan tidak bisa dihafal oleh algoritma ML, dan privasi individu yang menyumbangkan titik data tersebut ke dataset tetap terjaga.

Dalam tutorial ini, analisis privasi dilakukan dengan kerangka kerja *Renyi Differential Privacy* (RDP), yang merupakan relaksasi dari pure DP yang sangat cocok untuk DP-SGD.

Dua metrics yang digunakan untuk mengekspresikan jaminan DP dari sebuah algoritma ML :

1. Delta (δ)

Membatasi probabilitas jaminan privasi tidak berlaku. Aturannya Adalah nilainya lebih kecil dari kebalikan ukuran dataset pelatihan. Dalam tutorial ini, δ diset ke 10^{-5} karena dataset MNIST memiliki 60.000 titik data pelatihan.

2. Epsilon (ϵ)

Ini Adalah privacy budget. Mengukur kekuatan jaminan privasi dengan membatasi seberapa besar probabilitas keluaran model tertentu dapat bervariasi dengan menyertakan (atau menghapus) satu titik data pelatihan. Nilai ϵ yang lebih kecil berarti jaminan privasi lebih baik. Namun, nilai ϵ hanyalah batas atas; nilai yang besar masih bisa berarti privasi yang cukup baik dalam praktik.

Tensorflow privacy menyediakan sebuah alat, **compute_dp_sgd_privacy**, untuk menghitung nilai ϵ dengan nilai δ yang tetap dan hyperparameter berikut dari proses pelatihan:

1. Jumlah total titik data dalam data pelatihan, **n**
2. **Batch_size**
3. **noise_multiplier**
4. jumlah **epochs** pelatihan

```

87
88 ##### Measure the differential privacy guarantee #####
89 from tensorflow_privacy.privacy.analysis import compute_dp_sgd_privacy_lib
90
91 epsilon, best_alpha = compute_dp_sgd_privacy_lib.compute_dp_sgd_privacy(
92     n=train_data.shape[0],
93     batch_size=batch_size,
94     noise_multiplier=noise_multiplier,
95     epochs=epochs,
96     delta=1e-5
97 )
98
99 print(f"DP-SGD with  $\epsilon$  = {epsilon:.2f},  $\delta$  = 1e-5. For  $\alpha$  = {best_alpha}.")
100

```

```

(.venv) PS C:\KP\MATERI\4.2 Integrasi TF Privacy ke Model>
(.venv) PS C:\KP\MATERI\4.2 Integrasi TF Privacy ke Model> python classification_privacy.py
2025-09-10 11:33:43.345782: I tensorflow/core/util/port.cc:113] oneDNN custom operations are
e to floating-point round-off errors from different computation orders. To turn them off, se
WARNING:tensorflow:From C:\KP\MATERI\4.2 Integrasi TF Privacy ke Model\.venv\Lib\site-packag
oftmax_cross_entropy is deprecated. Please use tf.compat.v1.losses.sparse_softmax_cross_entr

2025-09-10 11:33:54.319346: I tensorflow/core/platform/cpu_feature_guard.cc:182] This Tensor
ns in performance-critical operations.
To enable the following instructions: SSE SSE2 SSE3 SSE4.1 SSE4.2 AVX2 FMA, in other operati
ags.
WARNING:absl:`compute_dp_sgd_privacy` is deprecated. It does not account for doubling of sen
ling, which is rarely used in practice. Please use `compute_dp_sgd_privacy_statement`, which
ute epsilon under different assumptions than those in `compute_dp_sgd_privacy_statement`, ca
DP-SGD with  $\epsilon$  = 0.56,  $\delta$  = 1e-5. For  $\alpha$  = 18.0.
(.venv) PS C:\KP\MATERI\4.2 Integrasi TF Privacy ke Model>

```

Alat tersebut melaporkan bahwa untuk hiperparameter yang dipilih di atas, model yang dilatih memiliki nilai sebesar 1,18.

Task : Tambahkan `dp_keras_optimizer` ke model FL

1. Import & hyperparameter DP klien

```

10 print("TensorFlow:", tf.__version__)
11 try:
12     from tensorflow_privacy.privacy.optimizers.dp_optimizer_keras import DPKerasSGDOptimizer
13     from tensorflow_privacy.privacy.analysis import compute_dp_sgd_privacy_lib
14     TF_PRIVACY_OK = True
15     print("TF-Privacy: OK (dp_keras_optimizer tersedia)")

```

```

TensorFlow: 2.14.1
TF-Privacy: OK (dp_keras_optimizer tersedia)

```

```

DP_CLIENT_LR           = 0.05
DP_L2_NORM_CLIP_CLIENT = 1.5
NOISE_WARMUP           = 0.001 # kecil supaya naik di awal
NOISE_MAIN              = 0.1

```


2. Dipasang ke model klien saat compile

```
# DP optimizer
opt = DPKerasSGDOptimizer(
    l2_norm_clip=DP_L2_NORM_CLIP_CLIENT,
    noise_multiplier=client_noise,
    num_microbatches=DP_NUM_MICROBATCHES, # = BATCH_SIZE
    learning_rate=DP_CLIENT_LR,
    momentum=0.9
)
loss = tf.keras.losses.BinaryCrossentropy(from_logits=False, reduction=tf.losses.Reduction.NONE)
local_model.compile(optimizer=opt, loss=loss, metrics=["accuracy"])
```

3. Dataset klien disiapkan agar kompatibel dengan DP

```
ds = ds.batch(BATCH_SIZE, drop_remainder=True).repeat(local_epochs).prefetch(tf.data.AUTOTUNE)
return ds
```

4. Setelah training lokal selesai, delta diklip & di-agregasi → inilah bagian FL-nya:

```
# Delta
w_local = local_model.get_weights()
delta = [wl - wg for wl, wg in zip(w_local, global_weights)]

# Clip user-level sebelum agregasi
delta_clipped, _ = clip_by_l2_norm(delta, SERVER_CLIP)
```

Hasil Training :

```
Federated Training (manual)
Round 01 | client_loss=0.4943 | global_loss=0.5345 | global_acc=0.7449 | ma_acc=0.7449 | client_noise=0.001
Round 02 | client_loss=0.4152 | global_loss=0.4404 | global_acc=0.7449 | ma_acc=0.7449 | client_noise=0.001
Round 03 | client_loss=0.3370 | global_loss=0.3723 | global_acc=0.7818 | ma_acc=0.7560 | client_noise=0.001
Round 04 | client_loss=0.2747 | global_loss=0.2796 | global_acc=0.8756 | ma_acc=0.7918 | client_noise=0.001
Round 05 | client_loss=0.2238 | global_loss=0.2243 | global_acc=0.9191 | ma_acc=0.8300 | client_noise=0.001
Round 06 | client_loss=0.1630 | global_loss=0.1881 | global_acc=0.9387 | ma_acc=0.8626 | client_noise=0.1
Round 07 | client_loss=0.1403 | global_loss=0.1756 | global_acc=0.9407 | ma_acc=0.8860 | client_noise=0.1
Round 08 | client_loss=0.1272 | global_loss=0.1682 | global_acc=0.9342 | ma_acc=0.9005 | client_noise=0.1
Round 09 | client_loss=0.1240 | global_loss=0.1645 | global_acc=0.9362 | ma_acc=0.9112 | client_noise=0.1
Round 10 | client_loss=0.1151 | global_loss=0.1629 | global_acc=0.9376 | ma_acc=0.9191 | client_noise=0.1
Round 11 | client_loss=0.1021 | global_loss=0.1589 | global_acc=0.9371 | ma_acc=0.9245 | client_noise=0.1
Round 12 | client_loss=0.0996 | global_loss=0.1506 | global_acc=0.9402 | ma_acc=0.9292 | client_noise=0.1
Round 13 | client_loss=0.0953 | global_loss=0.1472 | global_acc=0.9400 | ma_acc=0.9325 | client_noise=0.1
Round 14 | client_loss=0.0948 | global_loss=0.1369 | global_acc=0.9424 | ma_acc=0.9355 | client_noise=0.1
Round 15 | client_loss=0.0860 | global_loss=0.1282 | global_acc=0.9458 | ma_acc=0.9386 | client_noise=0.1

Training selesai (FedAvg manual + DP-SGD di klien).
WARNING:absl:compute_dp_sgd_privacy` is deprecated. It does not account for doubling of sensitivity with microbatching. To compute epsilon under different batch sizes, use `compute_dp_sgd_privacy_statement`, which provides appropriate context for the guarantee. To compute epsilon under different batch sizes, use `compute_dp_sgd_privacy_statement`.
Perkiraan DP-SGD klien (fase utama):  $\epsilon \approx 3361.26$ ,  $\delta = 1e-5$  ( $\alpha = 1.25$ )
(venv) ezranahumury@DESKTOP-80038IM:/mnt/c/KP/MATERI/4.2 Integrasi TF Privacy ke Model/task$
```

Measure the differential privacy

```
319 # =====
320 # 7) Estimasi epsilon untuk fase utama
321 # =====
322 try:
323     n = max(client_sizes)
324     eps, best_alpha = compute_dp_sgd_privacy_lib.compute_dp_sgd_privacy(
325         n=n,
326         batch_size=BATCH_SIZE,
327         noise_multiplier=NOISE_MAIN,
328         epochs=LOCAL_EPOCHS_MAIN,
329         delta=1e-5
330     )
331     print(f"Perkiraan DP-SGD klien (fase utama):  $\epsilon \approx \{eps:.2f\}$ ,  $\delta = 1e-5$  ( $\alpha = \{best\_alpha\}$ )")
332 except Exception as e:
333     print("Lewati perhitungan epsilon:", e)
334
```

raries directly.

Perkiraan DP-SGD klien (fase utama): $\epsilon \approx 3361.26$, $\delta = 1e-5$ ($\alpha = 1.25$)

(venv) ezranahumury@DESKTOP-8003BIM:/mnt/c/KP/MATERI/4.2 Integrasi TF Privacy ke Model/task\$

Estimasi ϵ , δ dengan `compute_dp_sgd_privacy_lib.compute_dp_sgd_privacy(...)` per klien.