

Topik : 4.4. Finalisasi Prototipe

Objective : Siapkan pipeline end-to-end: dataset → FL → DP → hasil

Task : Uji ulang semua modul, pastikan reproducible



1. Dataset (Load Data)

```
17
18 # =====
19 # 0) File input
20 # =====
21 DINSOS_CSV = "dinsos.csv"
22 DUKCAPIL_CSV = "dukcapil.csv"
23 KEMENKES_CSV = "kemenkes.csv"
24
25 # =====
26 # 1) Load CSV
27 # =====
28 def load_or_fail(path):
29     p = Path(path)
30     if not p.exists():
31         raise FileNotFoundError(f"File tidak ditemukan: {path}")
32     return pd.read_csv(p)
33
34 dinsos = load_or_fail(DINSOS_CSV)
35 dukcapil = load_or_fail(DUKCAPIL_CSV)
36 kemenkes = load_or_fail(KEMENKES_CSV)
37
```

2. Labeling

```
38 # =====
39 # 2) Labeling rules
40 # =====
41 def label_dinsos(row):
42     jt = row.get("jumlah_tanggungan", np.nan)
43     ph = row.get("penghasilan", np.nan)
44     kr = row.get("kondisi_rumah", "")
45     if pd.isna(jt): jt = 0
46     if pd.isna(ph): ph = 0
47     if ph < 2_000_000: return 1
48     if (2_000_000 <= ph < 3_500_000) and jt >= 4: return 1
49     if (kr in ("tidak layak", "semi permanen", "sangat sederhana")) and (ph < 5_000_000): return 1
50     if (jt >= 5) and (ph < 6_000_000): return 1
51     if ph >= 8_000_000: return 0
52     return 0
53
54 def label_dukcapil(row):
55     u = row.get("umur", np.nan)
56     sp = row.get("status_pekerjaan", "")
57     st = row.get("status_pernikahan", "")
58     if pd.isna(u): u = 0
59     if u > 65: return 1
60     if sp in ("pengangguran", "buruh", "pekerja informal"): return 1
61     if (sp=="wirausaha") and (u>55): return 1
62     if (st in ("cerai", "janda", "duda")) and (spi=="pegawai tetap"): return 1
63     if (sp=="pegawai tetap") and (25<u<55) and (st=="menikah"): return 0
64     if (u<25) and (sp=="wirausaha"): return 0
65     return 0
66
67 def label_kemenkes(row):
68     rp = row.get("riwayat_penyakit", "")
69     sg = row.get("status_gizi", "")
70     t = row.get("tinggi_cm", np.nan)
71     b = row.get("berat_kg", np.nan)
72     bmi = None
73     if pd.notna(t) and pd.notna(b) and t > 0:
74         bmi = b / ((t/100.0)**2)
75     if rp in ("kronis", "jantung", "asma", "diabetes", "disabilitas"): return 1
76     if sg in ("kurang", "stunting", "gizi buruk"): return 1
77     if (bmi is not None) and (bmi < 18.5): return 1
78     if (bmi is not None) and ((bmi < 17) or (bmi > 35)): return 1
79     return 0
80
81 # Apply label
82 dinsos["layak_subsid"] = dinsos.apply(label_dinsos, axis=1)
83 dukcapil["layak_subsid"] = dukcapil.apply(label_dukcapil, axis=1)
84 kemenkes["layak_subsid"] = kemenkes.apply(label_kemenkes, axis=1)
85
```

3. Preprocessing

```

187 # =====
188 # 3) Preprocess (one-hot + min-max scaling global)
189 # =====
190 cat_cols_all = ["kondisi_rumah", "status_pekerjaan", "status_pernikahan", "riwayat_penyakit", "status_gizi"]
191 num_cols_all = ["jumlah_tanggungan", "penghasilan", "umur", "tinggi_cm", "berat_kg"]
192
193 for df in (dinsos, dukcapil, kemenkes):
194     if ("tinggi_cm", "berat_kg").issubset(df.columns):
195         df["BMI"] = df["berat_kg"] / ((df["tinggi_cm"] / 100.0) ** 2)
196     else:
197         df["BMI"] = np.nan
198     if "BMI" not in num_cols_all:
199         num_cols_all.append("BMI")
200
201 def union_categories(series_list):
202     cats = set()
203     for s in series_list:
204         if s is not None:
205             for v in s.dropna().unique().tolist():
206                 if isinstance(v, str): cats.add(v)
207     return sorted(list(cats))
208
209 global_vocab = {
210     "kondisi_rumah": union_categories([dinsos.get("kondisi_rumah")]),
211     "status_pekerjaan": union_categories([dukcapil.get("status_pekerjaan")]),
212     "status_pernikahan": union_categories([dukcapil.get("status_pernikahan")]),
213     "riwayat_penyakit": union_categories([kemenkes.get("riwayat_penyakit")]),
214     "status_gizi": union_categories([kemenkes.get("status_gizi")]),
215 }
216
217 def encode_and_scale(df):
218     for c in num_cols_all:
219         if c not in df: df[c] = np.nan
220     for c in cat_cols_all:
221         if c not in df: df[c] = np.nan
222     oh_parts = []
223     for col, vocab in global_vocab.items():
224         for v in vocab:
225             oh_parts.append(("{}_{}".format(col, v), (df[col] == v).astype(float)))
226     oh_df = pd.DataFrame([na.values for na, a in oh_parts], index=df.index) if oh_parts else pd.DataFrame(index=df.index)
227     num_df = df[num_cols_all].copy()
228     for c in num_cols_all:
229         if not np.isfinite(num_df[c]).any(): num_df[c] = 0.0
230         else: num_df[c] = num_df[c].fillna(num_df[c].median())
231     X_raw = pd.concat([num_df, oh_df], axis=1)
232     y = df["layak_subsidir"].astype(int).values
233     return X_raw, y
234
235 X_dinsos_raw, y_dinsos = encode_and_scale(dinsos)
236 X_dukcapil_raw, y_dukcapil = encode_and_scale(dukcapil)
237 X_kemenkes_raw, y_kemenkes = encode_and_scale(kemenkes)
238
239 all_X = pd.concat([X_dinsos_raw, X_dukcapil_raw, X_kemenkes_raw], axis=0)
240 mins, maxs = all_X.min(axis=0), all_X.max(axis=0)
241 rng = (maxs - mins).replace(0, 1.0)
242 def scale_like_global(X): return (X - mins) / rng
243
244 X_dinsos = scale_like_global(X_dinsos_raw).fillna(0.0)
245 X_dukcapil = scale_like_global(X_dukcapil_raw).fillna(0.0)
246 X_kemenkes = scale_like_global(X_kemenkes_raw).fillna(0.0)
247
248 FEATURE_COLS = list(X_dinsos.columns)

```

4. Dataset per klien

```

249 # =====
250 # 4) Dataset per klien
251 # =====
252 LOCAL_EPOCHS = 1
253 BATCH_SIZE = 64
254 DP_NUM_MICROBATCHES = BATCH_SIZE
255
256 def to_ds(X, y, local_epochs, shuffle=True):
257     X = X.values.astype("float32")
258     y = y.astype("float32").reshape(-1,1)
259     ds = tf.data.Dataset.from_tensor_slices((X,y))
260     if shuffle:
261         ds = ds.shuffle(buffer_size=min(len(y), 2848), reshuffle_each_iteration=True)
262     ds = ds.batch(BATCH_SIZE, drop_remainder=True).repeat(local_epochs).prefetch(tf.data.AUTOTUNE)
263     return ds
264
265 client_data = [(X_dinsos, y_dinsos), (X_dukcapil, y_dukcapil), (X_kemenkes, y_kemenkes)]
266 client_sizes = [len(y_dinsos), len(y_dukcapil), len(y_kemenkes)]
267
268 def make_client_datasets(local_epochs, shuffle=True):
269     return [to_ds(X, y, local_epochs, shuffle) for (X, y) in client_data]
270

```

5. Model

```

271 # =====
272 # 5) Model
273 # =====
274 def build_model(input_dim):
275     i = tf.keras.Input(shape=(input_dim,))
276     x = tf.keras.layers.Dense(64, activation="relu")(i)
277     x = tf.keras.layers.Dense(32, activation="relu")(x)
278     o = tf.keras.layers.Dense(1, activation="sigmoid")(x)
279     return tf.keras.Model(i, o)
280
281 def get_weights(model): return [w.copy() for w in model.get_weights()]
282 def set_weights(model, weights): model.set_weights([w.copy() for w in weights])
283
284 def evaluate_global(model):
285     X_all = np.vstack([X_dinsos.values, X_dukcapil.values, X_kemenkes.values]).astype("float32")
286     y_all = np.concatenate([y_dinsos, y_dukcapil, y_kemenkes]).astype("float32").reshape(-1,1)
287     loss = tf.keras.losses.BinaryCrossentropy(from_logits=False)
288     y_pred = model.predict(X_all, batch_size=256, verbose=0)
289     l = float(loss(y_all, y_pred).numpy())
290     acc = float(np.mean((y_pred >= 0.5) == y_all))
291     return l, acc
292

```

6. Federated learning

```
193 # -----
194 # 6) Federated Training Loop (funct1)
195 # -----
196 def clip_by_l2_norm(tensors, clip):
197     s = 0.0
198     for t in tensors: s += np.sum(np.square(t))
199     norm = float(np.sqrt(s)) + 1e-12
200     if norm <= clip: return tensors, 1.0
201     factor = clip / norm
202     return [t * factor for t in tensors], factor
203
204 def federated_train(noise_multiplier, rounds=5):
205     DP_CLIENT_LR = 0.05
206     DP_L2_NORM_CLIP_CLIENT = 1.5
207     SERVER_CLIP = 5.0
208     SERVER_LR = 1.0
209
210     tf.keras.utils.set_random_seed(42)
211     np.random.seed(42)
212
213     global_model = build_model(len(FEATURE_COLS))
214     global_model.compile(optimizer="sgd", loss="binary_crossentropy", metrics=["accuracy"])
215     global_weights = get_weights(global_model)
216
217     acc_log, loss_log = [], []
218
219     for round_idx in range(1, rounds+1):
220         client_ds = make_client_datasets(local_epochs=LOCAL_EPOCHS, shuffle=True)
221         clipped_weighted_sums = None
222         total_weight = 0.0
223
224         for k, ds in enumerate(client_ds):
225             local_model = build_model(len(FEATURE_COLS))
226             set_weights(local_model, global_weights)
227             opt = DPKerasSGDOptimizer(
228                 l2_norm_clip=DP_L2_NORM_CLIP_CLIENT,
229                 noise_multiplier=noise_multiplier,
230                 num_microbatches=BATCH_SIZE,
231                 learning_rate=DP_CLIENT_LR,
232                 momentum=0.9
233             )
234             loss = tf.keras.losses.BinaryCrossentropy(from_logits=False, reduction=tf.keras.losses.Reduction.NONE)
235             local_model.compile(optimizer=opt, loss=loss, metrics=["accuracy"])
236             local_model.fit(ds, epochs=1, verbose=0)
237
238             w_local = local_model.get_weights()
239             delta = [w_l - w_g for w_l, w_g in zip(w_local, global_weights)]
240             delta_clipped, _ = clip_by_l2_norm(delta, SERVER_CLIP)
241
242             weight_k = float(client_sizes[k])
243             if clipped_weighted_sums is None:
244                 clipped_weighted_sums = [d * weight_k for d in delta_clipped]
245             else:
246                 for i in range(len(delta_clipped)):
247                     clipped_weighted_sums[i] += delta_clipped[i] * weight_k
248             total_weight += weight_k
249
250         avg_delta = [cw / (total_weight+1e-12) for cw in clipped_weighted_sums]
251         avg_delta = [SERVER_LR * d for d in avg_delta]
252         global_weights = [w_g + d for w_g, d in zip(global_weights, avg_delta)]
253         set_weights(global_model, global_weights)
254
255         gl_loss, gl_acc = evaluate_global(global_model)
256         acc_log.append(gl_acc)
257         loss_log.append(gl_loss)
258         print(f"[noise_multiplier:.3f] Round {round_idx} | acc={gl_acc:.4f} | loss={gl_loss:.4f}")
259
260     try:
261         n = max(client_sizes)
262         eps, _ = compute_dp_sgd_privacy_lib.compute_dp_sgd_privacy(
263             n=n,
264             batch_size=BATCH_SIZE,
265             noise_multiplier=noise_multiplier,
266             epochs=LOCAL_EPOCHS,
267             delta=1e-5
268         )
269     except Exception:
270         eps = np.nan
271
272     return np.mean(acc_log), np.mean(loss_log), eps, acc_log, loss_log
273
```

7. Differential Privacy

```
local_model = build_model(len(FEATURE_COLS))
set_weights(local_model, global_weights)
opt = DPKerasSGDOptimizer(
    l2_norm_clip=DP_L2_NORM_CLIP_CLIENT,
    noise_multiplier=noise_multiplier,
    num_microbatches=BATCH_SIZE,
    learning_rate=DP_CLIENT_LR,
    momentum=0.9
)

try:
    n = max(client_sizes)
    eps, _ = compute_dp_sgd_privacy_lib.compute_dp_sgd_privacy(
        n=n,
        batch_size=BATCH_SIZE,
        noise_multiplier=noise_multiplier,
        epochs=LOCAL_EPOCHS,
        delta=1e-5
    )
```

8. Evaluasi

```
274 # -----
275 # 7) Evaluasi Trade-off
276 # -----
277 noise_list = [0.01, 0.05, 0.1, 0.5, 1.0]
278 results = []
279 rounds = 5
280
281 for nm in noise_list:
282     acc, loss, eps, acc_log, loss_log = federated_train(noise_multiplier=nm, rounds=rounds)
283     results.append({"NoiseMultiplier": nm, "Accuracy": acc, "Loss": loss, "Epsilon": eps})
284
285 df_results = pd.DataFrame(results)
286 print("\nHasil Evaluasi Trade-off:")
287 print(df_results)
```

9. Plot Trade-off

```
288 # -----
289 # 8) Plot Trade-off (Accuracy vs Epsilon)
290 # -----
291 # -----
292 plt.figure(figsize=(7,5))
293 sns.lineplot(data=df_results, x="Epsilon", y="Accuracy", marker="o")
294 for i, row in df_results.iterrows():
295     plt.text(row["Epsilon"], row["Accuracy"], f"nm={row['NoiseMultiplier']}", fontsize=9)
296
297 plt.xlabel("Epsilon (ε) → Privasi ↓")
298 plt.ylabel("Accuracy → Utility ↑")
299 plt.title("Trade-off Federated Learning (DP-SGD)")
300 plt.grid(True)
301 plt.tight_layout()
302 plt.savefig("tradeoff_fl_accuracy_epsilon.png", dpi=150)
303 print("Grafik trade-off disimpan sebagai tradeoff_fl_accuracy_epsilon.png")
304
```

Pipeline end-to-end :

1. Dataset (load data)

Tahap pertama Adalah mengambil data mentah dari tiga client :

- Dinsos → berisi informasi jumlah_tanggungan, penghasilan, kondisi rumah
- Dukcapil → berisi informasi umur, status_pekerjaan, status_pernikahan
- Kemenkes → berisi informasi Riwayat_penyakit, status_gizi, tinggi, berat

Data ini di-load dari file CSV menggunakan `pandas.read_csv()`. Jika file tidak ada, program berhenti agar reproducibility terjaga.

2. Labeling

Data mentah tidak memiliki target langsung. Maka dibuat aturan untuk menentukan label layak_subsidi (1 = layak, 0 = tidak layak).

- Output tahap ini : kolom baru layak_subsidi pada tiap dataset

3. Pre-processing

Agar model dapat dilatih :

- Numerik : nilai yang hilang diisi mendian, lalu di-*scale* pakai min-max normalisasi global
- Kategori : diubah menjadi one-hot encoding berdasarkan vocab gabungan (supaya konsisten antar klien)
- Fitur Tambahan : BMI dihitung dari tinggi & berat badan

Output : matriks fitur numerik dan vector label yang sudah bersih

4. Dataset per Klien

Setiap sumber data dianggap sebagai client dalam Federated Learning.

- Data tiap klien dikonversi ke tf.data.Dataset kemudian dilakukan batching, shuffling, repeat sesuai dengan local epochs
- Tiga dataset klien siap dipakai dalam loop training

Hal ini mensimulasikan kondisi nyata Dimana data tidak dikumpulllkan disatu tempat, melainkan dilatih secara terdistribusi

5. Model

Menggunakan arsitektur :

- Input layer sesuai jumlah fitur
- Hidden layer : 64 neuron (ReLU) → 32 neuron (ReLU)
- Output layer : 1 neuron sigmoid (binary classification)

Penerapan model yang cukup ringan sehingga cocok untuk melakukan federated learning.

6. Federated Learning

Implementasi manual Federated Averaging (FedAvg) :

- Tiap klien dimulai dari bobot global
- Klien melatih model local pada datasetnya dengan optimizer DP-SGD
- Delta bobot hasil training di-clip dengan L2 norm agar tidak terlalu besar

- Server melakukan agregasi delta bobot dari semua klien (weighted average berdasarkan ukuran dataset)
- Update bobot global → looping untuk beberapa round

7. Differential privacy

Untuk menjamin privasi data tiap klien, maka digunakan :

- DP-SGD Optimizer → menambahkan *noise Gaussian* ke gradien, plus clipping gradien
- Menggunakan parameter penting
 - o L2_norm_clip → batas clipping gradien
 - o Noise_multiplier → besar noise yang ditambahkan
- Estimasi privasi dihitung dengan menggunakan `compute_dp_sgd_privacy` → menghasilkan nilai ϵ (epsilon), semakin kecil artinya privasi lebih kuat

8. Evaluasi

Setelah setiap federated round :

- Model global dievaluasi pada gabungan semua data klien
- Metrics : accuracy dan loss
- Disimpan log per round agar bisa dianalisis konsistensinya
- Dicatat pula nilai rata – rata akurasi/loss untuk tiap noise multiplier

9. Plot Trade-off

Hasil evaluasi divisualisasikan :

- Grafik Accuracy vs Epsilon (ϵ) → menunjukkan trade-off antara utility(akurasi) dan privacy (ϵ).
- Tiap titik diberi label nm = noise_multiplier
- Grafik disimpan dalam file PNG agar reproducible