

Math 564 - Project 02

Sajjad Uddin Mahmud

October 21, 2025

1 Introduction

2 Data Preparation

The Somerville happiness dataset contained 3,669 survey responses across multiple wards of the city, with one target variable (*Happiness*) and eight ordinal predictor variables representing satisfaction in different aspects of life. Each feature took values from 1 (very unsatisfied) to 5 (very satisfied) as inputs.

Handling Missing and Invalid Entries: Rows with missing target values could not be used for model training or evaluation and were therefore removed. Among the remaining records, rows containing six or more missing feature values (out of eight) were also discarded to prevent distortion caused by heavy imputation and to maintain data reliability.

Hierarchical Median Imputation: For the remaining data, missing values in the eight satisfaction features were imputed using a hierarchical median approach. Specifically, for each feature, the median value within the respondent's ward was used whenever available, reflecting local neighborhood-level satisfaction patterns. When a ward had insufficient data for a given feature, the overall citywide median for that feature was used instead. This approach preserves the ordinal nature of the satisfaction variables while maintaining robustness against outliers.

Normalization of Ordinal Features: All eight satisfaction features were then normalized to the interval $[0, 1]$ using the transformation

$$x_{\text{scaled}} = \frac{x - 1}{4},$$

which maps the original 1–5 Likert scale to a consistent numeric range. This scaling improves numerical conditioning and prevents sigmoid activation saturation during neural network optimization.

Target Encoding: Finally, the target variable (*Happiness*) was encoded as fractional values $\{1/6, 2/6, 3/6, 4/6, 5/6\}$ to align with the sigmoid network output range, allowing the model output to represent the ordered nature of happiness levels.

The whole process is shown in Algorithm 1

Algorithm 1 data preparation for somerville happiness dataset

Require: raw table \mathcal{D}_{raw} with columns: ward w , target y (happiness $\in \{1, \dots, 5\}$), and features $\mathcal{F} = \{f_1, \dots, f_8\}$ (eight 1–5 likert items)

Require: row-drop threshold $m_{\text{th}} \leftarrow 6$ (drop rows with $\geq m_{\text{th}}$ missing features)

Ensure: processed table $\mathcal{D}_{\text{proc}}$; imputation summary \mathcal{S}

```
1:  $\mathcal{D} \leftarrow \mathcal{D}_{\text{raw}}$ 
2: drop rows with missing target:  $\mathcal{D} \leftarrow \{r \in \mathcal{D} \mid r.y \neq \text{nan}\}$   $\triangleright$  remove
   rows that cannot be used for training or evaluation
3: drop rows with too many missing features:  $\mathcal{D} \leftarrow \{r \in \mathcal{D} \mid |\{f \in$ 
    $\mathcal{F} : r.f = \text{nan}\}| < m_{\text{th}}\}$ 
4: compute ward-level medians: for each  $f \in \mathcal{F}$  and ward  $w$ , let  $m_{w,f} \leftarrow$ 
   median $\{r.f : r \in \mathcal{D}, r.w = w, r.f \neq \text{nan}\}$ 
5: compute global medians: for each  $f \in \mathcal{F}$ , let  $m_f \leftarrow$  median $\{r.f : r \in$ 
    $\mathcal{D}, r.f \neq \text{nan}\}$ 
6: for each row  $r \in \mathcal{D}$  do
7:   for each feature  $f \in \mathcal{F}$  do
8:     if  $r.f = \text{nan}$  then
9:       if  $m_{r.w,f}$  is defined then
10:         $r.f \leftarrow m_{r.w,f}$   $\triangleright$  impute by ward median when available
11:      else
12:         $r.f \leftarrow m_f$   $\triangleright$  fallback to global median
13:      end if
14:    end if
15:  end for
16: end for
17: normalize features to  $[0, 1]$ : for each  $f \in \mathcal{F}$  and row  $r$ , set  $r.f_{\text{scaled}} \leftarrow$ 
   clip( $r.f, 1, 5$ ) then  $r.f_{\text{scaled}} \leftarrow (r.f_{\text{scaled}} - 1)/4$ 
18: encode target: for each row  $r$ , set  $r.y_{\text{enc}} \leftarrow r.y/6$   $\triangleright$ 
    $\{1, \dots, 5\} \mapsto \{1/6, \dots, 5/6\}$ 
19: record summary:  $\mathcal{S} \leftarrow \{\text{initial rows, rows dropped (missing target),}$ 
   rows dropped ( $\geq m_{\text{th}}$  missing), global medians  $m_f$ , counts of ward me-
   dians used $\}$ 
20: return  $\mathcal{D}_{\text{proc}} \leftarrow \mathcal{D}, \mathcal{S}$ 
```

3 Training and Test Data Selection

To evaluate the generalization capability of the neural network, the processed dataset was divided into separate training and testing subsets. A stratified random split was employed to preserve the class proportions of the target variable (*Happiness* levels 1–5) in both sets.

In this approach, the data were first grouped by their discrete happiness category. Within each group, approximately 30% of the samples were randomly assigned to the test set, while the remaining 70% were used for training. Stratification was selected to prevent bias toward dominant classes (e.g., “4” and “5”), ensuring that rare classes (such as “1”) were represented proportionally in both subsets.

4 Results

4.1 Data Preparation Results

After applying the above preprocessing steps, the dataset was cleaned and normalized. The resulting record counts are summarized in Table 4.1:

Table 4.1: Number of data points

Description	Count
Initial number of rows	3,669
Rows dropped (missing target)	34
Rows dropped (too many missing features)	8
Final number of usable rows	3,627

Thus, a total of 42 records (approximately 1.1% of the dataset) were excluded, leaving 3,627 fully processed entries for model development.

4.2 Train–Test Split Results

After applying the stratified 70/30 division, the resulting sample counts are summarized in Table 4.2.

Table 4.2: Stratified train–test split summary.

Statistic	Training Set	Test Set
Total rows	2,539	1,088
Happiness = 1	46	20
Happiness = 2	88	38
Happiness = 3	472	202
Happiness = 4	1,271	544
Happiness = 5	662	284

4.3 Confusion Matrix

The feed-forward neural network (FFNN) model, consisting of eight input features, two hidden layers of 12 and 10 neurons, and a single sigmoid output neuron, was trained using two optimization methods: (i) **BFGS with Strong Wolfe line search**, and (ii) **Gradient Descent with Strong Wolfe line search**.

Both models were evaluated on the stratified test dataset. The resulting confusion matrices are shown in Table 4.3 and 4.3.

Table 4.3: Confusion Matrix — BFGS + Strong Wolfe

True \ Predicted	1	2	3	4	5
1	1	2	3	12	2
2	2	2	6	25	3
3	2	11	12	163	14
4	3	13	16	475	37
5	6	6	6	211	55

Table 4.4: Confusion Matrix — GD + Strong Wolfe

True \ Predicted	1	2	3	4	5
1	0	3	4	13	0
2	0	0	11	27	0
3	0	0	24	178	0
4	0	0	21	523	0
5	0	0	7	277	0

The confusion matrices in Tables 4.3 and 4.3 reveal consistent prediction patterns across both optimization methods. In both cases, the network displays a strong tendency to classify most samples as belonging to the middle or higher satisfaction levels, particularly class 4 (*Satisfied*). This behavior reflects the underlying imbalance in the dataset, where classes 4 and 5 together account for more than 70% of all responses in the training set.

5 Exploration of Alternative FFNN Constructions

Beyond the baseline architecture of two hidden layers with 12 and 10 neurons, several alternative network structures can be explored to assess model performance. A couple of things may be explored like:

- Increasing or decreasing the number of hidden units
- Adding or removing hidden layers
- Variations in activation functions (e.g., ReLU instead of sigmoid)

6 Repository

All code and iteration logs for this project are available in a public GitHub repository:

https://github.com/sajjad30148/WSU_Math564_Fall2025

The repository includes the code, and result folders containing iterations from each run.

7 Code

The full Python implementation is provided below. It can also be found in the repository.

7.1 Main Python Script

Listing 1: project01_main.py

```
1 # =====
2 # Project 02 Main Script
3 # =====
4
5 import os
6 from datetime import datetime
7 import numpy as np
8 import pandas as pd
9 from pathlib import Path
10 import sys
11 import matplotlib.pyplot as plt
12
13 workspace_root = Path(__file__).resolve().parent.parent
14 sys.path.insert(0, str(workspace_root))
15
16 import functions as fn
17
18 # -----
19 # user settings
20 # -----
21
22 # data file path
23 data_path = r"D:\One_Drive_Sajjad\OneDrive - Washington State University
    (email.wsu.edu)\Documents\Sajjad_Uddin_Mahmud\Courses\5. Fall
    2025\MATH_564\Projects\WSU_Math564_Fall2025\Project_02\happiness.csv"
24
25 # data preparation settings
26 ward_col = "Ward/Neighborhood"
27 target_col = "Happiness"
28 feature_cols = [
29     "Beauty of Neighborhood",
30     "Convenience of Getting Around",
31     "Housing Condition",
32     "Street and Sidewalk Maintenance",
33     "Public Schools",
```

```

34     "Police Department",
35     "Community Events",
36     "City Services Information",
37 ]
38 drop_row_missing_threshold = 6 # rows with >= this many missing features
    will be dropped
39 keep_raw_features = True # set to False to overwrite original columns
40
41 # train/test split settings
42 test_size = 0.30
43 random_state = 42
44 stratify = True
45 strat_col = target_col
46
47 # -----
48 # data preparation function
49 # -----
50
51 def prepare_happiness_data(
52     df,
53     ward_col = "Ward/Neighborhood",
54     target_col = "Happiness",
55     feature_cols = None,
56     drop_row_missing_threshold = 6,
57     keep_raw_features = True,
58 ):
59     """
60     prepare the somerville happiness dataset for modeling
61
62     inputs
63         df pandas dataframe with raw data
64         ward_col column name for ward / neighborhood
65         target_col column name for happiness (1..5)
66         feature_cols list of the eight satisfaction feature column names
67         drop_row_missing_threshold rows with >= this many missing features
            will be dropped
68         keep_raw_features if true, original 1..5 features are kept; scaled
            features
69                                     are added with suffix '_scaled'. if
                                     false, originals are
70                                     overwritten by scaled values.
71
72     outputs
73         df_out processed dataframe

```



```

74         prep_info dictionary with imputation and drop statistics for
           reporting
75     """
76
77     if feature_cols is None:
78         # set the expected eight feature names here
79         feature_cols = [
80             "Beauty of Neighborhood",
81             "Convenience of Getting Around",
82             "Housing Condition",
83             "Street and Sidewalk Maintenance",
84             "Public Schools",
85             "Police Department",
86             "Community Events",
87             "City Services Information",
88         ]
89
90     df = df.copy()
91
92     # -----
93     # drop unusable rows
94     # -----
95
96     # remove rows with missing target
97     n0 = len(df)
98     mask_missing_target = df[target_col].isna()
99     n_missing_target = int(mask_missing_target.sum())
100    df = df.loc[~mask_missing_target].reset_index(drop = True)
101
102    # remove rows with too many missing features
103    feature_missing_count = df[feature_cols].isna().sum(axis = 1)
104    mask_too_many_missing = feature_missing_count >=
        drop_row_missing_threshold
105    n_drop_too_many_missing = int(mask_too_many_missing.sum())
106    if n_drop_too_many_missing > 0:
107        df = df.loc[~mask_too_many_missing].reset_index(drop = True)
108
109    # -----
110    # compute medians (ward-level and global)
111    # -----
112
113    # ward-level medians
114    ward_medians = {}
115    for col in feature_cols:
116        ward_medians[col] = df.groupby(ward_col)[col].median()

```

```

117
118 # global medians
119 global_medians = {col: float(df[col].median()) for col in
    feature_cols}
120
121 # -----
122 # impute feature missing values by ward -> global
123 # -----
124
125 # try ward-level median first, then global median
126 for col in feature_cols:
127     ward_median_col = df.groupby(ward_col)[col].transform("median")
128     df[col] = df[col].fillna(ward_median_col)
129
130     # fallback to global median if ward median was nan or still missing
131     df[col] = df[col].fillna(global_medians[col])
132
133 # -----
134 # normalize features to [0, 1] by (x - 1) / 4
135 # -----
136
137 # scale function with clamping
138 def scale_to_unit_interval(x):
139     x = x.astype(float)
140
141     # clamp to [1, 5] for safety, then scale
142     x = np.clip(x, 1.0, 5.0)
143
144     return (x - 1.0) / 4.0
145
146 # apply scaling
147 if keep_raw_features:
148     for col in feature_cols:
149         df[f"{col}_scaled"] = scale_to_unit_interval(df[col])
150     scaled_feature_cols = [f"{c}_scaled" for c in feature_cols]
151 else:
152     for col in feature_cols:
153         df[col] = scale_to_unit_interval(df[col])
154     scaled_feature_cols = feature_cols
155
156 # -----
157 # encode target as fractions {1/6, ..., 5/6}
158 # -----
159 # happiness expected in {1, 2, 3, 4, 5}
160 # encode as k / 6 where k in {1..5}

```

```

161 df["Happiness_Encoded"] = df[target_col].astype(float) / 6.0
162
163 # assemble prep info for reporting
164 df_out = df
165 prep_info = {
166     "n_rows_initial" : n0,
167     "n_missing_target_dropped" : n_missing_target,
168     "n_drop_too_many_missing" : n_drop_too_many_missing,
169     "feature_cols" : feature_cols,
170     "scaled_feature_cols" : scaled_feature_cols,
171     "ward_medians_available_for_cols" : {
172         col: int(ward_medians[col].notna().sum()) for col in
            feature_cols
173     },
174     "global_medians" : global_medians,
175     "scaling" : "(x - 1) / 4 with clamp to [1, 5]",
176     "target_encoding" : "Happiness_Encoded = Happiness / 6",
177     "ward_col" : ward_col,
178     "target_col" : target_col,
179 }
180
181 return df_out, prep_info
182
183
184 if __name__ == "__main__":
185
186     # -----
187     # load raw data
188     # -----
189     data_path_obj = Path(data_path)
190     df_raw = pd.read_csv(data_path_obj)
191
192     # -----
193     # prepare data (no saving inside the function)
194     # -----
195     df_proc, prep_info = prepare_happiness_data(
196         df_raw,
197         ward_col = ward_col,
198         target_col = target_col,
199         feature_cols = feature_cols,
200         drop_row_missing_threshold = drop_row_missing_threshold,
201         keep_raw_features = keep_raw_features,
202     )
203
204

```

```

205 # -----
206 # save processed data as csv (done here, not in the function)
207 # -----
208 out_path = data_path_obj.with_name("happiness_processed.csv")
209 df_proc.to_csv(out_path, index = False)
210
211 # simple console report
212 print("processed data saved to :", out_path)
213 print("rows initial :", prep_info["n_rows_initial"])
214 print("rows dropped (target) :",
215       prep_info["n_missing_target_dropped"])
216 print("rows dropped (>= missing threshold) :",
217       prep_info["n_drop_too_many_missing"])
218
219 # -----
220 # train/test split
221 # -----
222 if stratify:
223     df_train, df_test = fn.stratified_split_df(
224         df_proc,
225         target_col = strat_col,
226         test_size = test_size,
227         random_state = random_state,
228     )
229 else:
230     # simple random split without stratification
231     df_shuf = df_proc.sample(frac = 1.0, random_state =
232                             random_state).reset_index(drop = True)
233     n_total = len(df_shuf)
234     n_test = int(round(test_size * n_total))
235     n_test = max(1, min(n_test, n_total - 1))
236     df_test = df_shuf.iloc[:n_test].reset_index(drop = True)
237     df_train = df_shuf.iloc[n_test:].reset_index(drop = True)
238
239 # simple train/test split report
240 print("train rows :", len(df_train))
241 print("test rows :", len(df_test))
242 print("train class counts :")
243 print(df_train[strat_col].value_counts().sort_index())
244 print("test class counts :")
245 print(df_test[strat_col].value_counts().sort_index())
246
247 # select features (scaled) and target (encoded)
248 feat_cols = prep_info["scaled_feature_cols"]
249 y_col_enc = "Happiness_Encoded"

```

```

247
248 X_train = df_train[feat_cols].to_numpy(dtype = float)
249 y_train = df_train[y_col_enc].to_numpy(dtype = float).reshape(-1, 1)
250 X_test = df_test[feat_cols].to_numpy(dtype = float)
251 y_test = df_test[y_col_enc].to_numpy(dtype = float).reshape(-1, 1)
252
253 # layer sizes for the required architecture
254 layer_sizes = [X_train.shape[1], 12, 10, 1]
255
256 # create feedforward neural network objective
257 f, g = fn.make_ffnn_objective(X_train, y_train, layer_sizes =
    layer_sizes, l2 = 0.0)
258
259 # initialize parameters
260 w0 = fn.ffnn_init_params(layer_sizes, random_state =
    np.random.default_rng(42))
261
262 # -----
263 # helper: convert scalar outputs in (0,1) to classes {1..5}
264 # targets are spaced at {1/6, 2/6, 3/6, 4/6, 5/6}
265 # -----
266 targets_enc = np.arange(1, 6, dtype = float).reshape(-1, 1) / 6.0 #
    shape (5,1)
267 targets_cls = np.arange(1, 6, dtype = int).reshape(-1, 1) # shape
    (5,1)
268
269 def decode_to_class(yhat):
270     # yhat: (n,1) in (0,1); snap to nearest of the five target points
271     diffs = np.abs(yhat - targets_enc.T) # (n,5)
272     idx = np.argmin(diffs, axis = 1) # (n,)
273     return targets_cls[idx, 0] # (n,)
274
275
276 # -----
277 # run method 1: bfgs + strong wolfe
278 # -----
279 bfgs_opts = {"max_iter" : 500, "tol" : 1e-6}
280 res_bfgs = fn.quasi_newton_bfgs(
281     f = f,
282     grad = g,
283     x0 = w0,
284     line_search = fn.strong_wolfe,
285     opts = bfgs_opts
286 )
287 w_bfgs = res_bfgs["x"]

```

```

288
289 # -----
290 # run method 2: gradient descent + strong wolfe
291 # -----
292 from functions import gradient_descent
293
294 gd_opts = {"max_iter" : 200, "tol" : 1e-6}
295 res_gd = fn.gradient_descent(
296     f = f,
297     grad = g,
298     x0 = w0,
299     line_search = fn.strong_wolfe,
300     opts = gd_opts
301 )
302 w_gd = res_gd["x"]
303
304
305
306 # -----
307 # prediction helper
308 # -----
309 def ffnn_predict(X, w, layer_sizes):
310     shapes = []
311     sizes = []
312     for l in range(1, len(layer_sizes)):
313         d_prev = layer_sizes[l - 1]
314         d_curr = layer_sizes[l]
315         shapes.append(((d_prev, d_curr), (d_curr,)))
316         sizes.append((d_prev * d_curr, d_curr))
317     idx = 0
318     slices = []
319     for (w_count, b_count) in sizes:
320         sW = slice(idx, idx + w_count)
321         idx += w_count
322         sb = slice(idx, idx + b_count)
323         idx += b_count
324         slices.append((sW, sb))
325
326     def sigmoid(u):
327         return 1.0 / (1.0 + np.exp(-u))
328
329     a = X
330     for (sW, sb), (shapeW, shapeb) in zip(slices, shapes):
331         W = w[sW].reshape(shapeW)
332         b = w[sb].reshape(shapeb)

```

```

333         a = sigmoid(a @ W + b)
334     return a
335
336     # test predictions for both methods
337     yhat_bfgs = ffnn_predict(X_test, w_bfgs, layer_sizes)
338     yhat_gd = ffnn_predict(X_test, w_gd, layer_sizes)
339
340     print("GD yhat min/max/mean:", float(yhat_gd.min()),
341           float(yhat_gd.max()), float(yhat_gd.mean()))
342     print("Unique predicted classes (GD):",
343           np.unique(decode_to_class(yhat_gd)))
344     print("Are GD weights unchanged from init?:", np.allclose(w0, w_gd))
345
346     ycls_bfgs = decode_to_class(yhat_bfgs)
347     ycls_gd = decode_to_class(yhat_gd)
348     ytrue_cls = decode_to_class(y_test)
349
350     # -----
351     # Confusion matrix
352     # -----
353
354     cm_bfgs = pd.crosstab(
355         pd.Series(ytrue_cls.flatten(), name = "True"),
356         pd.Series(ycls_bfgs.flatten(), name = "Predicted"),
357     )
358     cm_gd = pd.crosstab(
359         pd.Series(ytrue_cls.flatten(), name = "True"),
360         pd.Series(ycls_gd.flatten(), name = "Predicted"),
361     )
362
363     print("\nConfusion Matrix: BFGS + Strong Wolfe")
364     print(cm_bfgs)
365     print("\nConfusion Matrix: Gradient Descent + Strong Wolfe")
366     print(cm_gd)
367
368     classes = [1, 2, 3, 4, 5]
369
370     # reindex to full 1...5 grid for neat printing/saving
371     cm_bfgs = cm_bfgs.reindex(index = classes, columns = classes,
372                               fill_value = 0)
373     cm_gd = cm_gd.reindex(index = classes, columns = classes, fill_value
374                            = 0)
375
376     # -----

```

```

374     # save to ./results/evaluation.txt
375     # -----
376     results_dir = Path(__file__).resolve().parent / "results"
377     results_dir.mkdir(parents = True, exist_ok = True)
378     with open(os.path.join(results_dir, "evaluation.txt"), "w") as f:
379         f.write("confusion matrix: bfgs + strong wolfe\n")
380         f.write(str(cm_bfgs) + "\n\n")
381         f.write("confusion matrix: gradient descent + strong wolfe\n")
382
383     print("saved confusion matrices to ./results/evaluation.txt")
384
385     # =====

```


7.2 Functions

Listing 2: functions.py

```
1 # -----
2 # fast feedforward neural network objective
3 # -----
4
5 def make_ffnn_objective(X, y, layer_sizes, l2 = 0.0):
6     """
7     create closures f(w) and grad(w) for a general ffnn
8
9     inputs
10         X array (n, d_in), features
11         y array (n,) or (n, 1), targets in {1/6, ..., 5/6}
12         layer_sizes list like [d_in, h1, ..., d_out] with d_out = 1
13         l2 nonnegative l2 regularization on weights (biases excluded)
14
15     returns
16         f callable, f(w) -> scalar loss
17         grad callable, grad(w) -> flat gradient vector
18     """
19
20     X = np.asarray(X, dtype = float)
21     y = np.asarray(y, dtype = float).reshape(-1, 1)
22
23     n, d_in = X.shape
24     assert layer_sizes[0] == d_in, "input dimension must match"
25     assert layer_sizes[-1] == 1, "output dimension must be 1 for this"
26     setup"
27
28     # precompute parameter slices
29     # for each layer l: W in R^{d_{l-1} x d_l}, b in R^{d_l}
30     shapes = []
31     sizes = []
32     for l in range(1, len(layer_sizes)):
33         d_prev = layer_sizes[l - 1]
34         d_curr = layer_sizes[l]
35         w_count = d_prev * d_curr
36         b_count = d_curr
37         shapes.append(((d_prev, d_curr), (d_curr,)))
38         sizes.append((w_count, b_count))
39
40     # compute flat index ranges
41     idx = 0
```

```

41 slices = []
42 for (w_count, b_count) in sizes:
43     sW = slice(idx, idx + w_count)
44     idx += w_count
45     sb = slice(idx, idx + b_count)
46     idx += b_count
47     slices.append((sW, sb))
48 p_expected = idx
49
50 # cached activations for backprop; created locally inside f/grad
51 def _sigmoid(u):
52     return 1.0 / (1.0 + np.exp(-u))
53
54 def f(w):
55     w = np.asarray(w, dtype = float)
56     if w.size != p_expected:
57         raise ValueError("size of w does not match architecture")
58
59     # unpack and forward pass
60     a_list = [X] # a_0
61     params = []
62     a = X
63     for (sW, sb), (shapeW, shapeb) in zip(slices, shapes):
64         W = w[sW].reshape(shapeW)
65         b = w[sb].reshape(shapeb)
66         params.append((W, b))
67         z = a @ W + b
68         a = _sigmoid(z)
69         a_list.append(a)
70
71     yhat = a_list[-1] # (n, 1)
72     residual = yhat - y
73     loss = 0.5 * np.sum(residual * residual)
74
75     if l2 > 0.0:
76         reg = 0.0
77         for (W, b) in params:
78             reg += np.sum(W * W)
79         loss += 0.5 * l2 * reg
80
81     return float(loss)
82
83 def grad(w):
84     w = np.asarray(w, dtype = float)
85     if w.size != p_expected:

```

```

86         raise ValueError("size of w does not match architecture")
87
88     # unpack and forward pass (store activations for backprop)
89     a_list = [X]
90     z_list = []
91     params = []
92     a = X
93     for (sW, sb), (shapeW, shapeb) in zip(slices, shapes):
94         W = w[sW].reshape(shapeW)
95         b = w[sb].reshape(shapeb)
96         params.append((W, b))
97         z = a @ W + b
98         a = _sigmoid(z)
99         z_list.append(z)
100        a_list.append(a)
101
102    yhat = a_list[-1]
103
104    # initial delta at output: (yhat - y) * sigma'(z_L)
105    # derivative of sigmoid via activation: a * (1 - a)
106    delta = (yhat - y) * (yhat * (1.0 - yhat)) # (n, 1)
107
108    # backprop
109    dWs = [None] * len(params)
110    dbs = [None] * len(params)
111
112    for l in reversed(range(len(params))):
113        a_prev = a_list[l] # (n, d_l)
114        W, b = params[l] # W: (d_l, d_{l+1}), b: (d_{l+1},)
115
116        dW = a_prev.T @ delta # (d_l, d_{l+1})
117        db = np.sum(delta, axis = 0) # (d_{l+1},)
118        if l2 > 0.0:
119            dW += l2 * W
120
121        dWs[l] = dW
122        dbs[l] = db
123
124        if l > 0:
125            a_prev_act = a_list[l] # activation at layer l
126            delta = (delta @ W.T) * (a_prev_act * (1.0 - a_prev_act))
127
128    # pack grads into a flat vector following the same order
129    g = np.empty(p_expected, dtype = float)
130    for (sW, sb), dW, db in zip(slices, dWs, dbs):

```

```

131         g[sW] = dW.ravel()
132         g[sb] = db.ravel()
133
134     return g
135
136 return f, grad
137
138
139 def ffnn_init_params(layer_sizes, random_state = None):
140     """
141     xavier/glorot uniform initializer for sigmoid networks
142
143     inputs
144         layer_sizes [d_in, h1, ..., d_out]
145         random_state seed or numpy random generator or None
146
147     returns
148         w0 flat parameter vector
149     """
150     if random_state is None:
151         random_state = np.random.default_rng(42)
152
153     blocks = []
154     for l in range(1, len(layer_sizes)):
155         d_prev = layer_sizes[l - 1]
156         d_curr = layer_sizes[l]
157         limit = np.sqrt(6.0 / (d_prev + d_curr))
158         W = random_state.uniform(-limit, limit, size = (d_prev, d_curr))
159         b = np.zeros((d_curr,), dtype = float)
160         blocks.append(W.ravel())
161         blocks.append(b)
162     return np.concatenate(blocks, axis = 0)
163
164
165
166
167 # -----
168 # Miscellaneous functions
169 # -----
170
171 # stratified train-test split for dataframes
172 def stratified_split_df(df, target_col, test_size = 0.3, random_state =
173     42):
174     """
175     stratified split by discrete target values

```

```

175
176 inputs
177     df dataframe to split
178     target_col column with class labels (here: original 1..5 happiness)
179     test_size fraction in test
180     random_state seed for reproducibility
181
182 outputs
183     df_train, df_test
184     """
185 rng = np.random.default_rng(random_state)
186
187 df_train_parts = []
188 df_test_parts = []
189
190 for cls, grp in df.groupby(target_col):
191     n = len(grp)
192     n_test = int(round(test_size * n))
193     n_test = max(1, min(n_test, n - 1)) # at least 1 in each side when
194                                     possible
195
196     idx = np.arange(n)
197     rng.shuffle(idx)
198
199     test_idx = idx[:n_test]
200     train_idx = idx[n_test:]
201
202     df_test_parts.append(grp.iloc[test_idx])
203     df_train_parts.append(grp.iloc[train_idx])
204
205 df_train = pd.concat(df_train_parts, axis = 0).sample(frac = 1.0,
206                                                         random_state = random_state).reset_index(drop = True)
207 df_test = pd.concat(df_test_parts, axis = 0).sample(frac = 1.0,
208                                                         random_state = random_state).reset_index(drop = True)
209
210 return df_train, df_test

```