

## 1 Honor Code

- Collaborator: None
- *I certify that all solutions are entirely in my own words and that I have not looked at another student's solutions. I have given credit to all external sources I consulted.*

Signature: \_\_\_\_\_

## 2 Logistic Regression with Newton's Method

1. Derive  $\frac{ds(\gamma)}{d\gamma}$

$$\begin{aligned}\frac{ds(\gamma)}{d\gamma} &= \frac{d}{dx}\left(\frac{1}{x}\right) \frac{d}{d\gamma}(1 + e^{-\gamma}) \quad \text{where } x = 1 + e^{-\gamma} \\ &= -\frac{1}{x^2}(-e^{-\gamma}) \\ &= \frac{e^{-\gamma}}{(1 + e^{-\gamma})^2} \\ &= s(\gamma)(1 - s(\gamma))\end{aligned}$$

Derive  $\nabla_w J(w)$

$$\begin{aligned}\nabla_w J(w) &= -\sum_{i=1}^n \left( \frac{y_i}{s_i} \frac{ds_i}{dw} - \frac{1 - y_i}{1 - s_i} \frac{ds_i}{dw} \right) \\ &= -\sum_{i=1}^n \left( \frac{y_i}{s_i} - \frac{1 - y_i}{1 - s_i} \right) \frac{ds_i}{dw} \\ &= -\sum_{i=1}^n \frac{y_i(1 - s_i) - s_i(1 - y_i)}{s_i(1 - s_i)} \frac{ds_i}{d\gamma_i} \frac{d\gamma_i}{dw} \\ &= -\sum_{i=1}^n \frac{y_i(1 - s_i) - s_i(1 - y_i)}{s_i(1 - s_i)} s_i(1 - s_i) \mathbf{x}_i \\ &= -\sum_{i=1}^n (y_i - s_i) \mathbf{x}_i \\ &= -X^T(\mathbf{y} - \mathbf{s}) \quad \text{where } \mathbf{s} = \begin{bmatrix} s_1 & s_2 & \dots & s_n \end{bmatrix}^T\end{aligned}$$

2. Derive  $\nabla_w^2 J(w)$

$$\begin{aligned}\nabla_w^2 J(w) &= \nabla_w [\nabla_w J(w)]^T \\ &= \nabla_w [-X^T(\mathbf{y} - \mathbf{s})]^T \\ &= \nabla_w [-(\mathbf{y}^T - \mathbf{s}^T)X] \\ &= (\mathbf{s}^T(\mathbf{1} - \mathbf{s})X)^T X \\ &= X^T(1 - \mathbf{s})^T \mathbf{s} X \\ &= X^T \Omega X \quad \text{where } \Omega = \begin{bmatrix} s_1(1 - s_1) & 0 & \dots & 0 \\ 0 & s_2(1 - s_2) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & s_n(1 - s_n) \end{bmatrix}\end{aligned}$$

3. By the update rule of Newton's method derived in the lecture,

$$\begin{aligned}\mathbf{w}^{(n+1)} &= \mathbf{w}^{(n)} - (\nabla_{\mathbf{w}^{(n)}}^2 J(\mathbf{w}^{(n)}))^{-1} \nabla_{\mathbf{w}^{(n)}} J(\mathbf{w}^{(n)}) \\ &= \mathbf{w}^{(n)} + (X^T \Omega^{(n)} X)^{-1} (X^T(\mathbf{y} - \mathbf{s}^{(n)}))\end{aligned}$$

4.

$$\mathbf{s}^{(0)} = \begin{bmatrix} 0.9478 \\ 0.8808 \\ 0.8022 \\ 0.5250 \end{bmatrix}$$

$$\mathbf{w}^{(1)} = \begin{bmatrix} 1.3247 \\ 3.0499 \\ -6.8291 \end{bmatrix}$$

$$\mathbf{s}^{(1)} = \begin{bmatrix} 0.9474 \\ 0.9746 \\ 0.0315 \\ 0.1044 \end{bmatrix}$$

$$\mathbf{w}^{(2)} = \begin{bmatrix} 1.3660 \\ 4.1585 \\ -9.1996 \end{bmatrix}$$

### 3 Wine Classification with Logistic Regression

1.

$$\mathbf{w}^{(n+1)} = \mathbf{w}^{(n)} + \epsilon[2\lambda\mathbf{w}^{(n)} - X^T(y - \mathbf{s}^{(n)})]$$

2. The "train\_bgd\_decay\_False" (blue line) in the Figure 1 shows cost function versus iterations of Batch Gradient Descent(BGD) with the following parameters.

$$\epsilon = 10^{-7}$$

$$\lambda = 1$$

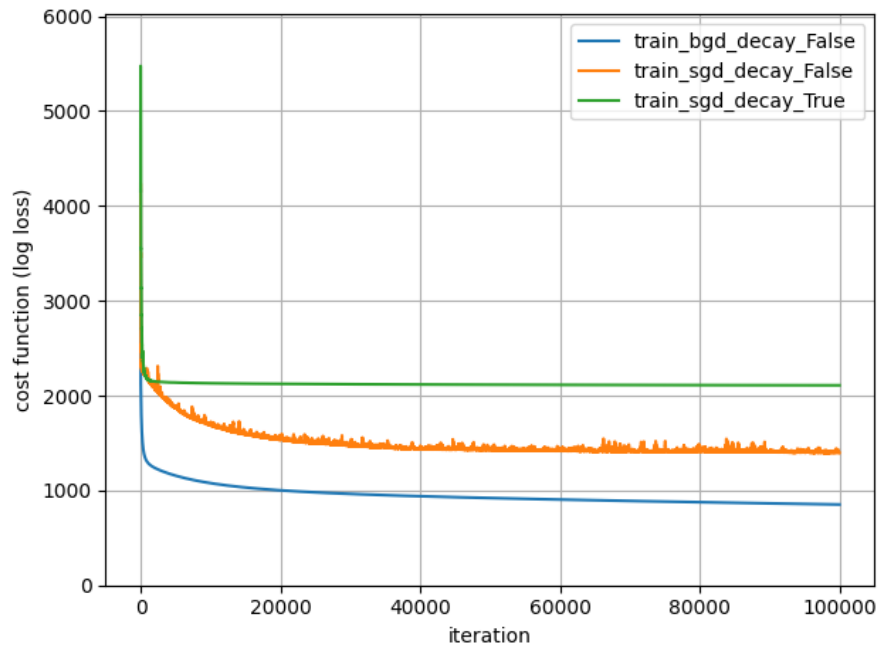


Figure 1: Update iterations and cost function

(Code is in the appendix. Corresponding pages are selected at Grade Scope submission)

3.

$$\mathbf{w}^{(n+1)} = \mathbf{w}^{(n)} + \epsilon[2\lambda\mathbf{w}^{(n)} - \mathbf{x}_i^T(\mathbf{y}_i - \mathbf{s}_i^{(n)})]$$

4. The "train\_sgd\_decay\_False" (orange line) in the Figure 1 shows cost function versus iterations of Stochastic Gradient Descent(SGD) with the following parameters.

$$\epsilon = 10^{-5}$$

$$\lambda = 1$$

BGD converges more quickly than SGD because BGD uses all input for one update while SGD uses only one randomly selected input. If the input size is much larger, SGD may perform better because SGD takes shorter computational time for one iteration than BGD.

(Code is in the appendix. Corresponding pages are selected at Grade Scope submission)

5. The "train\_sgd\_decay\_True" (green line) in the Figure 1 shows cost function versus iterations of SGD with the following parameters and shrinking learning rate.

$$\epsilon = 10^{-3}$$

$$\lambda = 1$$

SGD with shrinking learning rate needs more iterations to converge than SGD with static learning rate. Shrinking learning rate has possibility to reach the better optimum point by searching optimum from roughly to finely. However, in this problem, the shrinking policy ( $\epsilon_t = \frac{\delta}{t}$ ) updates the learning rate too drastically.

(Code is in the appendix. Corresponding pages are selected at Grade Scope submission)

- 6.
- Kaggle username: Takuma Kinoshita
  - Best score: 0.99396
  - Preprocess: Standardization ( $x' = \frac{x-\mu}{\sigma}$ ) to all features to remove the scale difference
  - Model: Logistic Regression with  $l_2$  regularization ( $C = 1$ ), fixed learning rate ( $10^{-5}$ ) and  $10^6$  times update iterations

(Code is in the appendix. Corresponding pages are selected at Grade Scope submission)

## 4 Convergence of Batch Gradient Descent in Logistic Regression

1. Derive  $\nabla_w J(w)$ .

$$\begin{aligned}
 J(\mathbf{w}) &= \frac{\lambda \|\mathbf{w}\|^2}{2} - \mathbf{y}^T \ln(s(X\mathbf{w})) - (1 - \mathbf{y})^T \ln(1 - s(X\mathbf{w})) \\
 &= \frac{\lambda \|\mathbf{w}\|^2}{2} - \sum_{i=1}^n y_i \ln(s(\mathbf{x}_i^T \mathbf{w})) - (1 - y_i) \ln(1 - s(\mathbf{x}_i^T \mathbf{w})) \\
 \nabla_{\mathbf{w}} J(\mathbf{w}) &= \lambda \mathbf{w} - \sum_{i=1}^n (y_i - s(\mathbf{x}_i^T \mathbf{w})) \mathbf{x}_i \quad (\text{from Question 2.1}) \\
 &= \lambda \mathbf{w} - X^T (\mathbf{y} - s(X\mathbf{w}))
 \end{aligned}$$

The update rule of batch gradient descent is:

$$\begin{aligned}
 \mathbf{w}^{(t+1)} &= \mathbf{w}^{(t)} - \epsilon \nabla_{\mathbf{w}} J(\mathbf{w}) \\
 &= \mathbf{w}^{(t)} - \epsilon [\lambda \mathbf{w}^{(t)} - X^T (\mathbf{y} - s(X\mathbf{w}^{(t)}))] \\
 &= g(\mathbf{w}^{(t)})
 \end{aligned}$$

2. Derive  $\nabla_w^2 J(w)$ .

$$\begin{aligned}
 \nabla_{\mathbf{w}}^2 J(\mathbf{w}) &= \nabla_{\mathbf{w}} [\nabla_{\mathbf{w}} J(\mathbf{w})]^T \\
 &= \nabla_{\mathbf{w}} [\lambda \mathbf{w}^T - (\mathbf{y}^T - s(X\mathbf{w})^T) X] \\
 &= \lambda I + X^T \Omega X \quad (\text{from Question 2.2})
 \end{aligned}$$

where  $\Omega = \begin{bmatrix} s_1(1-s_1) & 0 & \dots & 0 \\ 0 & s_2(1-s_2) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & s_n(1-s_n) \end{bmatrix}$

First,  $\lambda I$  is positive definite because  $\lambda > 0$ . Second,  $\Omega$  is positive definite because all diagonal entities of diagonal matrix  $\Omega$  is positive. When  $\Omega$  is positive definite,  $X^T \Omega X$  is positive semidefinite. Therefore,  $\lambda I + X^T \Omega X$  is positive definite.

3. By Mean Value Theorem, there exist  $w_\rho$  satisfying:

$$\begin{aligned}
 \frac{g(\mathbf{w}_1)_i - g(\mathbf{w}_2)_i}{\mathbf{w}_{1i} - \mathbf{w}_{2i}} &= \nabla g(\mathbf{w}_\rho)_i \quad (i = 1, \dots, d) \\
 g(\mathbf{w}_1)_i - g(\mathbf{w}_2)_i &= \nabla g(\mathbf{w}_\rho)_i (\mathbf{w}_{1i} - \mathbf{w}_{2i}) \\
 \|g(\mathbf{w}_1) - g(\mathbf{w}_2)\|^2 &= \|\nabla g(\mathbf{w}_\rho)^T (\mathbf{w}_1 - \mathbf{w}_2)\|^2
 \end{aligned}$$

By Cauchy-Schwarz inequality,

$$\|\nabla g(\mathbf{w}_\rho)(\mathbf{w}_1 - \mathbf{w}_2)\|^2 \leq \|\nabla g(\mathbf{w}_\rho)\|^2 \|\mathbf{w}_1 - \mathbf{w}_2\|^2$$

By adjusting  $\epsilon$ , there exist  $\|\nabla g(\mathbf{w}_\rho)\|^2 \in (0, 1)$ . Therefore,  $g(\cdot)$  is a contraction.

$$\|g(\mathbf{w}_1) - g(\mathbf{w}_2)\|^2 \leq \rho \|\mathbf{w}_1 - \mathbf{w}_2\|^2$$

4.

$$\begin{aligned}
\|\mathbf{w}^{(t)} - \mathbf{w}^{(t+1)}\| &\leq \rho \|\mathbf{w}^{(t-1)} - \mathbf{w}^{(t)}\| \\
&\leq \rho^2 \|\mathbf{w}^{(t-2)} - \mathbf{w}^{(t-1)}\| \\
&\vdots \\
&\leq \rho^{t-1} \|\mathbf{w}^{(0)} - \mathbf{w}^{(1)}\|
\end{aligned}$$

Because  $\rho \in (0, 1)$ ,

$$\lim_{t \rightarrow \infty} \rho^{t-1} \|\mathbf{w}^{(0)} - \mathbf{w}^{(1)}\| = 0$$

Thus,

$$\lim_{t \rightarrow \infty} \|\mathbf{w}^{(t)} - \mathbf{w}^{(t+1)}\| = 0$$

Therefore,  $\mathbf{w}^{(t)}$  converges to the unique minimizer.

5. Find  $\epsilon$  which minimizes  $\rho = \|\nabla g(\mathbf{w}_\rho)\|^2$

6. From Question 4-2, when  $\lambda = 0$ ,

$$\nabla_{\mathbf{w}^2} J(\mathbf{w}) = X^T \Omega X$$

Because  $X^T \Omega X$  is positive semidefinite, there are more than one optimal points. Therefore, the logistic regression problem does not converge.

## 5 A Bayesian Interpretation of Lasso

1. By Bayes' Theorem,

$$\begin{aligned}
 f(w | (\mathbf{x}_i, y_i)_{i \in [n]}) &= \frac{f(\mathbf{x}_{i \in [n]}, y_{i \in [n]} | \mathbf{w}) f(\mathbf{w})}{f(\mathbf{x}_{i \in [n]}, y_{i \in [n]})} \\
 &= \frac{f(y_{i \in [n]} | \mathbf{x}_{i \in [n]}, \mathbf{w}) f(\mathbf{x}_{i \in [n]}) f(\mathbf{w})}{f(\mathbf{x}_{i \in [n]}, y_{i \in [n]})} \\
 &\propto f(y_{i \in [n]} | \mathbf{x}_{i \in [n]}, \mathbf{w}) f(\mathbf{w})
 \end{aligned}$$

2. From Question 5.1

$$\begin{aligned}
 l(w) &= \ln f(\mathbf{w} | \mathbf{x}_{i \in [n]}, y_{i \in [n]}) \\
 &\propto f(y_{i \in [n]} | \mathbf{x}_{i \in [n]}, \mathbf{w}) f(\mathbf{w}) \\
 &= \ln \prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(y_i - \mathbf{w} \cdot \mathbf{x}_i)^2}{2\sigma^2}} \prod_{j=1}^d \frac{e^{-\frac{|w_j|}{b}}}{2b} \\
 &\propto - \sum_{i=1}^n \frac{(y_i - \mathbf{w} \cdot \mathbf{x}_i)^2}{2\sigma^2} + \sum_{j=1}^d \frac{|w_j|}{b} \\
 &\propto - \sum_{i=1}^n (y_i - \mathbf{w} \cdot \mathbf{x}_i)^2 + \frac{2\sigma^2}{b} \sum_{j=1}^d |w_j|
 \end{aligned}$$

Thus, when  $\frac{2\sigma^2}{b} = \lambda$ , maximizing  $l(\mathbf{w})$  is equivalent of minimizing  $\sum_{i=1}^n (y_i - \mathbf{w} \cdot \mathbf{x}_i)^2 + \lambda \|\mathbf{w}\|_1$



## 6 $l_1$ -regularization, $l_2$ -regularization, and Sparsity

1.

$$\begin{aligned}
 J_1(\mathbf{w}) &= \|X\mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|_1 \\
 &= \mathbf{y}^T \mathbf{y} - (X\mathbf{w})^T \mathbf{y} - \mathbf{y}^T (X\mathbf{w}) + (X\mathbf{w})^T X X \mathbf{w} + \lambda \|\mathbf{w}\|_1 \\
 &= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T X \mathbf{w} + \mathbf{w}^T X^T X \mathbf{w} + \lambda \|\mathbf{w}\|_1 \\
 &= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T X \mathbf{w} + n\mathbf{w}^T \mathbf{w} + \lambda \|\mathbf{w}\|_1 \\
 &= \|\mathbf{y}\|^2 + \sum_{i=1}^d -2\mathbf{y}^T \mathbf{x}_{*i} w_i + n w_i^2 + \lambda |w_i| \\
 &= \|\mathbf{y}\|^2 + \sum_{i=1}^d f(\mathbf{x}_{*i}, w_i)
 \end{aligned}$$

2. When  $w_i^* > 0$ ,

$$\begin{aligned}
 \nabla_{\mathbf{w}} J_1(\mathbf{w})_i &= -2\mathbf{y}^T \mathbf{x}_{*i} + 2n w_i + \lambda = 0 \\
 \Rightarrow w_{*i} &= \frac{2\mathbf{y}^T \mathbf{x}_{*i} + \lambda}{2n}
 \end{aligned}$$

When  $w_i^* < 0$ ,

$$\begin{aligned}
 \nabla_{\mathbf{w}} J_1(\mathbf{w})_i &= -2\mathbf{y}^T \mathbf{x}_{*i} + 2n w_i - \lambda = 0 \\
 \Rightarrow w_{*i} &= \frac{2\mathbf{y}^T \mathbf{x}_{*i} - \lambda}{2n}
 \end{aligned}$$

When  $w_i^* = 0$ , by subgradient,

$$\begin{aligned}
 \nabla_{\mathbf{w}} J_1(\mathbf{w})_i &= [-2\mathbf{y}^T \mathbf{x}_{*i} + 2n w_i - \lambda, -2\mathbf{y}^T \mathbf{x}_{*i} + 2n w_i + \lambda] = 0 \\
 \Rightarrow &\begin{cases} -2\mathbf{y}^T \mathbf{x}_{*i} + 2n w_i - \lambda \leq 0 \\ -2\mathbf{y}^T \mathbf{x}_{*i} + 2n w_i + \lambda \geq 0 \end{cases} \\
 \Rightarrow &-\lambda \leq 2\mathbf{y}^T \mathbf{x}_{*i} + 2n w_i \leq \lambda \\
 \Rightarrow &-\lambda \leq 2\mathbf{y}^T \mathbf{x}_{*i} \leq \lambda \quad (w_i^* = 0)
 \end{aligned}$$

3. Similary,

$$\begin{aligned}
 J_2(\mathbf{w}) &= \|X\mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|_2 \\
 &= \|\mathbf{y}\|^2 + \sum_{i=1}^d -2\mathbf{y}^T \mathbf{x}_{*i} w_i + n w_i^2 + \lambda w_i^2 \\
 \nabla_{\mathbf{w}} J_2(\mathbf{w})_i &= -2\mathbf{y}^T \mathbf{x}_{*i} + 2n w_i - 2\lambda w_i = 0 \\
 \Rightarrow w_i^* &= \frac{\mathbf{y}^T \mathbf{x}_{*i}}{n + \lambda}
 \end{aligned}$$

When  $w_i^\# = 0$ ,

$$w_i^* = \frac{\mathbf{y}^T \mathbf{x}_{*i}}{n + \lambda} = 0 \\ \Rightarrow \mathbf{y}^T \mathbf{x}_{*i} = 0$$

4.  $w^*$  ( $l_1$  regularization) is more likely to be sparse than  $w^\#$  ( $l_2$  regularization). In  $l_1$  regularization, to be  $w_i = 0$ ,  $\mathbf{y}^T \mathbf{x}_{*i}$  need to in the range between  $-\lambda$  and  $\lambda$ . However, in  $l_2$  regularization,  $\mathbf{y}^T \mathbf{x}_{*i}$  need to be exactly 0.

## 7 Appendix: Source Codes

### 1. utils.py (Question 3.2, 3.4, 3.5, 3.6)

```

1 from collections import defaultdict
2
3 import numpy as np
4 import pandas as pd
5 import scipy.io as sio
6 from scipy.special import expit
7 import matplotlib.pyplot as plt
8 from sklearn.metrics import accuracy_score
9
10 from tqdm import tqdm
11
12
13 def load_data(path):
14     data = sio.loadmat(path)
15     return data['X'], data['X_test'], data['y']
16
17
18 def my_log_loss(y, s):
19     t1 = y.T @ np.log(s)
20     t2 = (~y.astype(bool)).astype(int).T @ np.log(1 - s)
21     return -float(t1 + t2)
22
23
24 def plot_costs(costs_df, fig_path):
25     fig, ax = plt.subplots()
26     costs_df.plot(ax=ax)
27     ax.set_xlabel('iteration')
28     ax.set_ylabel('cost function (log loss)')
29     ax.set_ylim(0, float(costs_df.max().max()) * 1.1)
30     ax.grid()
31     fig.tight_layout()
32     fig.savefig(fig_path)
33
34
35 class LogisticRegression():
36
37     def __init__(self, C, lr, lr_decay=False, max_iter=1e2, solver='sgd'):
38         self.C = C
39         self.lr = lr
40         self.lr_decay = lr_decay
41         self.solver = solver
42         self.max_iter = max_iter
43         self.solver = solver
44         if lr_decay:
45             self.lr_init = lr
46
47     def fit(self, x_tr, y_tr, verbose=False):
48         # initialize
49         ones = np.ones((x_tr.shape[0], 1))
50         x_tr = np.hstack((x_tr, ones))
51         self.w = np.zeros((x_tr.shape[1], 1))
52
53         costs = defaultdict(list)

```

```

54         for i in tqdm(range(self.max_iter)):
55             # output prob
56             s = expit(x_tr @ self.w)
57
58             # cost function
59             if verbose:
60                 cost_train = my_log_loss(y_tr, s)
61                 costs[f'train_{self.solver}_decay_{self.lr_decay}'].append(
62                     cost_train)
63
64             # learning rate decay
65             if self.lr_decay:
66                 self.lr = self.lr_init / (i + 1)
67
68             # update weight
69             if self.solver == 'bgd':
70                 self.w = self.w - self.lr * (2 * self.C * self.w - x_tr.T @ (y_tr -
71                     s))
72             if self.solver == 'sgd':
73                 rand_i = np.random.randint(0, x_tr.shape[0], 1)
74                 self.w = self.w - self.lr * (2 * self.C * self.w - (y_tr[rand_i] -
75                     s[rand_i]) * x_tr[rand_i].T)
76
77             if verbose:
78                 return pd.DataFrame(costs)
79
80         def predict_proba(self, x):
81             ones = np.ones((x.shape[0], 1))
82             x = np.hstack((x, ones))
83             return expit(x @ self.w).ravel()
84
85         def predict(self, x):
86             prob = self.predict_proba(x)
87             return np.round(prob).ravel()
88
89         def score(self, x, y):
90             pred = self.predict(x)
91             return accuracy_score(y, pred)

```

## 2. gd\_logreg.py (Question 3.2, 3.4, 3.5)

```

1  import numpy as np
2  import pandas as pd
3  from pathlib import Path
4
5  from utils import load_data, LogisticRegression, plot_costs
6
7
8  if __name__ == "__main__":
9      # set seed
10     np.random.seed = 289
11
12     # load data
13     data_path = Path('data/data.mat')
14     x_train, x_test, y_train = load_data(data_path)
15

```

```

16 # parameters
17 bgd_params = {
18     'C': 1,
19     'lr': 1e-7,
20     'max_iter': int(1e5),
21     'solver': 'bgd'
22 }
23 sgd_params = {
24     'C': 1,
25     'lr': 1e-5,
26     'max_iter': int(1e5),
27     'solver': 'sgd'
28 }
29 decay_params = {
30     'C': 1,
31     'lr': 1e-3,
32     'lr_decay': True,
33     'max_iter': int(1e5),
34     'solver': 'sgd'
35 }
36
37 # models
38 model = LogisticRegression(**bgd_params)
39 bgd_costs_df = model.fit(x_train, y_train, verbose=True)
40
41 model = LogisticRegression(**sgd_params)
42 sgd_costs_df = model.fit(x_train, y_train, verbose=True)
43
44 model = LogisticRegression(**decay_params)
45 decay_costs_df = model.fit(x_train, y_train, verbose=True)
46
47 # plot
48 costs_df = pd.concat([bgd_costs_df, sgd_costs_df, decay_costs_df], axis=0)
49 fig_path = Path('figures/iterations_vs_costs.png')
50 plot_costs(costs_df, fig_path)

```

### 3. kaggle.py (Question 3.6)

```

1 from pathlib import Path
2 from collections import defaultdict
3
4 import numpy as np
5 import pandas as pd
6 from sklearn.model_selection import StratifiedKFold
7 from sklearn.metrics import accuracy_score
8 from sklearn.preprocessing import StandardScaler
9
10 from utils import load_data, plot_costs, LogisticRegression
11 from save_csv import results_to_csv
12
13
14 def cross_validation(x_train, y_train, params, n_splits, random_state):
15     skf = StratifiedKFold(n_splits=n_splits, shuffle=True, random_state=seed)
16     scores = defaultdict(list)
17
18     for tr_idx, va_idx in skf.split(x_train, y_train):
19         x_tr, x_va = x_train[tr_idx], x_train[va_idx]

```

```
20     y_tr, y_va = y_train[tr_idx], y_train[va_idx]
21
22     model = LogisticRegression(**params)
23     model.fit(x_tr, y_tr)
24
25     scores['acc_tr'].append(model.score(x_tr, y_tr))
26     scores['acc_va'].append(model.score(x_va, y_va))
27
28     return pd.DataFrame(scores)
29
30
31 if __name__ == "__main__":
32
33     # set seed
34     seed = 289
35     np.random.seed = seed
36
37     # load data
38     data_path = Path('data/data.mat')
39     x_train, x_test, y_train = load_data(data_path)
40
41     # parameters
42     params = {
43         'C': 1,
44         'lr': 1e-5,
45         'lr_decay': False,
46         'max_iter': int(1e6),
47         'solver': 'bgd'
48     }
49
50     # preprocessing
51     scaler = StandardScaler()
52     scaler.fit(np.vstack((x_train, x_test)))
53     x_train = scaler.transform(x_train)
54     x_test = scaler.transform(x_test)
55
56     # evaluate
57     scores_df = cross_validation(x_train, y_train, params=params, n_splits=5,
58     random_state=seed)
59     print(scores_df)
60     print(scores_df.mean(axis=0))
61
62     # predict
63     model = LogisticRegression(**params)
64     model.fit(x_train, y_train)
65     y_pred = model.predict(x_test)
66
67     # save
68     results_to_csv(y_pred)
```