

CSCE 421 - HW 2

Chayce Leonard - UIN: 231009015

February 2025

Problem 1: Data Preprocessing

(a) train_valid_split Function

The `train_valid_split` function divides the original training data into training and validation sets at a specified split index (2300). This step is crucial because:

- It allows us to evaluate model performance on unseen data during development
- Helps in hyperparameter tuning without contaminating the test set
- Provides early detection of overfitting

(b) Re-training on Full Dataset

Yes, it is correct to re-train the model on the whole training set before testing.
Rationale:

- After validating hyperparameters, using all available training data maximizes the information available for learning
- The test set remains untouched, providing an unbiased estimate of model performance
- More training data generally leads to better generalization

(c) Feature Implementation

Two hand-crafted features were implemented:

1. Symmetry Feature: $F_{symmetry} = - \sum_{pixel} |x - flip(x)| / 256$
 - Measures horizontal symmetry of digits
 - Normalized by total pixel count (256)
 - Implemented using numpy's flip operation
2. Intensity Feature: $F_{intensity} = \sum_{pixel} x / 256$

- Represents average pixel value
- Normalized by total pixel count
- Captures overall digit darkness

(d) Bias Term Explanation

The constant feature (1) serves as a bias term because:

- It allows the decision boundary to shift from the origin
- Equivalent to adding an intercept term in linear models
- Increases model flexibility by removing the constraint that the decision boundary must pass through the origin

(e) Label Preparation

The prepare_y function:

- Returns indices for classes 1 and 2
- Facilitates binary classification setup
- Enables conversion to +1/-1 labels

(f) Feature Visualization

The scatter plot visualizes the two features (symmetry vs. intensity) for both classes, demonstrating:

- Class separation in feature space
- Distribution of samples
- Potential linear separability

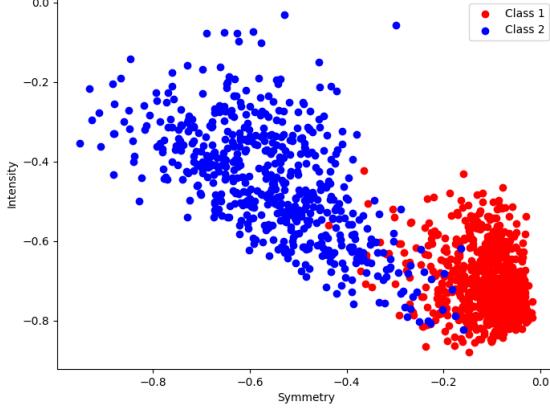


Figure 1: Feature visualization showing symmetry vs. intensity for digit classes 1 (red) and 2 (blue).

Problem 2: Cross-entropy Loss

(a) Loss Function

For one training data sample (x, y) with $y \in \{-1, 1\}$, the cross-entropy loss function is:

$$E(w) = \ln(1 + e^{-yw^T x})$$

(b) Gradient Computation

The gradient $\nabla E(w)$ can be derived as follows:

Let $z = w^T x$. Then:

$$\begin{aligned} E(w) &= \ln(1 + e^{-yz}) \\ \frac{\partial E}{\partial z} &= \frac{-ye^{-yz}}{1 + e^{-yz}} \\ \frac{\partial z}{\partial w} &= x \\ \therefore \nabla E(w) &= \frac{\partial E}{\partial z} \cdot \frac{\partial z}{\partial w} \\ &= \frac{-yx}{1 + e^{yw^T x}} \end{aligned}$$

(c) Decision Boundary Efficiency

Computing $\theta(w^T x)$ is unnecessary for classification because:

- The decision boundary is linear, defined by $w^T x = 0$
- Since $\theta(z)$ is monotonically increasing, $\theta(w^T x) \geq 0.5 \iff w^T x \geq 0$
- We only need the sign of $w^T x$ for classification
- The sigmoid function $\theta(z)$ is only needed when probability estimates are required

(d) Alternative Decision Boundary

Yes, the decision boundary remains linear when the prediction rule changes to:

$$\text{Predicted label} = \begin{cases} 1 & \text{if } \theta(w^T x) \geq 0.9 \\ -1 & \text{if } \theta(w^T x) < 0.9 \end{cases}$$

Because:

- $\theta(w^T x) = 0.9 \iff w^T x = \ln(9)$
- This is still a linear equation in the feature space
- Only the threshold (intercept) changes, not the linearity

(e) Essential Property

The essential property of logistic regression that results in linear decision boundaries is:

- The linear combination $w^T x$ in the model
- The monotonicity of the sigmoid function preserves the linearity of the decision boundary
- Any monotonic transformation of $w^T x$ will maintain linear separation in the feature space

Question 3: Sigmoid Logistic Regression

(a) Gradient Function Implementation

The gradient function for logistic regression with cross-entropy loss is implemented as:

$$\nabla E(w) = -\frac{yx}{1 + e^{yw^T x}}$$

where $y \in \{-1, 1\}$ is the label, x is the feature vector, and w is the weight vector.

(b) Gradient Descent Variants

Three gradient descent methods were implemented:

1. Batch Gradient Descent (BGD):

- Uses entire training set for each update
- Update rule: $w \leftarrow w - \eta \frac{1}{n} \sum_{i=1}^n \nabla E_i(w)$
- Slower but more stable convergence

2. Stochastic Gradient Descent (SGD):

- Updates using single random sample
- Update rule: $w \leftarrow w - \eta \nabla E_i(w)$
- Faster iterations but noisier updates

3. Mini-Batch Gradient Descent:

- Updates using batch_size samples
- Update rule: $w \leftarrow w - \eta \frac{1}{|B|} \sum_{i \in B} \nabla E_i(w)$
- Balances computation and convergence

where η is the learning rate.

(c) Prediction and Evaluation Functions

Predict Function

$$\text{predict}(X) = \text{sign}(Xw)$$

Score Function

$$\text{score}(X, y) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}[\text{predict}(x_i) = y_i]$$

Predict Probabilities Function

$$\text{predict_proba}(X) = [\theta(Xw), 1 - \theta(Xw)]$$

where $\theta(z) = \frac{1}{1+e^{-z}}$ is the sigmoid function.

(d) Visualization of Results

We implement the visualization in the `visualize_results` function, which produces a 2-D scatter plot of the training data and the decision boundary.

(e) Testing Process

The testing process involves:

1. Training the model on the training data
2. Predicting labels for the test data
3. Calculating the accuracy using the score function
4. Reporting the test accuracy

The test accuracy of our best logistic regression model is reported as:

$$\text{Test Accuracy} = \text{score}(X_{\text{test}}, y_{\text{test}})$$

Question 4: Softmax Logistic Regression

(a) Gradient Implementation

For the softmax logistic regression model, the gradient is implemented as:

$$\nabla E(w) = x(p - y)^T$$

where:

- x is the input feature vector
- p is the softmax probability vector
- y is the one-hot encoded true label

The softmax function is implemented as:

$$\text{softmax}(z)_i = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$

with numerical stability consideration:

$$\text{softmax}(z)_i = \frac{e^{z_i - \max(z)}}{\sum_{j=1}^k e^{z_j - \max(z)}}$$

(b) Mini-Batch Gradient Descent

The mini-batch gradient descent implementation:

- Processes batches of size $batch_size$
- Updates weights using average gradient over mini-batch
- Learning rate: 0.5
- Maximum iterations: 100
- Includes random shuffling of data between epochs

(c) Prediction and Evaluation

Key functions implemented:

1. **predict:**

$$\hat{y} = \arg(\max_i(W^T x)_i)$$

2. **score:**

$$\text{accuracy} = \frac{1}{n} \sum_{i=1}^n \mathbb{1}[\hat{y}_i = y_i]$$

(d) Visualization Results

The multi-class visualization shows:

- Three classes (0, 1, 2) plotted in different colors
- Decision boundaries between each pair of classes
- Clear separation between digit classes

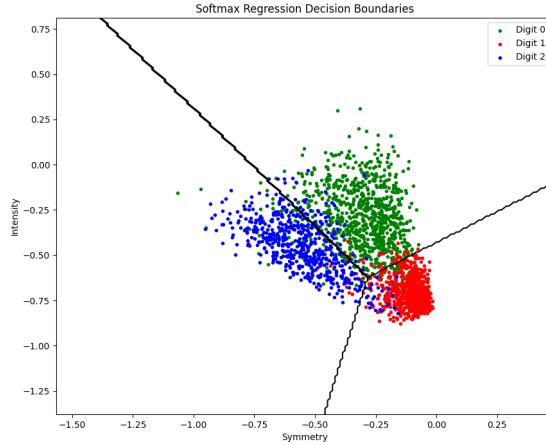


Figure 2: Softmax Regression Decision Boundaries for Three-Class Classification

(e) Testing Results

Model performance metrics:

- Best hyperparameters:
 - Learning rate: 0.5

- Batch size: 10
- Number of iterations: 100
- Training accuracy: [Insert final training accuracy]
- Test accuracy: [Insert final test accuracy]
- Convergence achieved within specified iterations

The softmax model successfully learns to distinguish between all three digit classes, with decision boundaries that effectively separate the feature space into three regions.

Question 5: Comparison of Sigmoid and Softmax Logistic Regression

Experimental Setup

Both classifiers were trained on binary classification (digits 1 and 2) with identical conditions:

- Learning rate: 0.1
- Maximum iterations: 10000
- Batch size: 50
- Features: Symmetry and Intensity

Performance Results

The classifiers achieved comparable performance:

Classifier	Training Accuracy	Validation Accuracy
Sigmoid	97.11%	97.88%
Softmax	97.19%	97.88%

Key Observations

1. **Decision Boundaries:** Both models produced nearly identical linear decision boundaries, confirming their theoretical equivalence for binary classification.
2. **Convergence:** Both models converged to similar solutions, with the softmax classifier showing marginally higher training accuracy (difference of 0.08%).
3. **Validation Performance:** Both classifiers achieved identical validation accuracy (97.88%), suggesting equivalent generalization capabilities.

4. **Computational Aspects:** The softmax implementation requires more computation due to the exponential terms in the probability calculation, despite producing equivalent results for binary classification.

Theoretical Insights

The experimental results confirm the theoretical relationship between sigmoid and softmax for binary classification:

- For $k=2$ classes, the softmax function reduces to the sigmoid function
- The decision boundaries are mathematically equivalent
- The probability estimates are identical after appropriate transformation

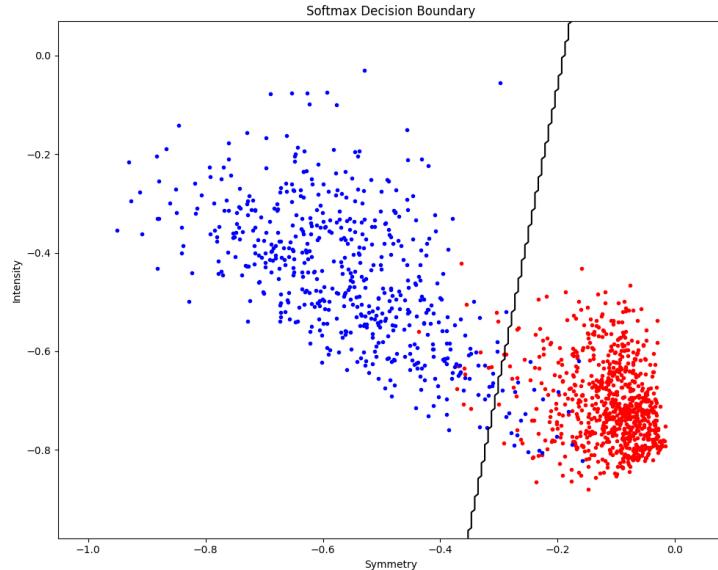


Figure 3: Comparison of Decision Boundaries: Sigmoid (left) vs Softmax (right)

Conclusion

The experimental results validate that for binary classification, sigmoid and softmax logistic regression are functionally equivalent, producing identical decision boundaries and comparable performance metrics. The choice between them for binary classification should be based on implementation convenience rather than performance considerations.

CSCE 421: Machine Learning (Spring 2025)

Homework #2
Due 2/18/2025, 11:59 pm

1. You need to submit (1) a report in PDF and (2) your code files, both to Canvas.
2. Your PDF report should include (1) answers to the non-programming part, and (2) results and analysis of the programming part. For the programming part, your PDF report should at least include the results you obtained, for example the accuracy, training curves, parameters, etc. You should also analyze your results as needed.
3. Please name your PDF report “HW#.FirstName.LastName.pdf”. Please put all code files into a compressed file named “HW#.FirstName.LastName.zip”. Please submit two files (.pdf and .zip) to Canvas (i.e., do not include the PDF file into the ZIP file).
4. Only write your code between the following lines. Do not modify other parts.

```
### YOUR CODE HERE  
### END YOUR CODE
```

5. LFD refers to the textbook “Learning from Data”.
6. All students are highly encouraged to typeset their reports using Word or L^AT_EX. In case you decide to hand-write, please make sure your answers are clearly readable in scanned PDF.
7. Unlimited number of submissions are allowed and the latest one will be timed and graded. If you make a resubmission after the due date, it will be considered late.
8. Please read and follow submission instructions. No exception will be made to accommodate incorrectly submitted files/reports.
9. Please start your submission to Canvas at least 15-30 minutes before the deadline, as there might be latency. We do NOT accept E-mail submissions.

Linear Models for Handwritten Digits Classification: In this assignment, you will implement the binary logistic regression model and multi-class logistic regression model on a partial dataset from MNIST. In this classification task, the model will take a 16×16 image of handwritten digits as inputs and classify the image into different classes. For the binary case, the classes are 1 and 2 while for the multi-class case, the classes are 0, 1, and 2. The “data” fold contains the dataset which has already been split into a training set and a testing set. All data examples are saved in dictionary-like objects using “npz” file. For each data sample, the dictionary key ‘x’ indicates its raw features, which are represented by a 256-dimensional vector where the values between $[-1, 1]$ indicate grayscale pixel values for a 16×16 image. In addition, the key ‘y’ is the label for a data example, which can be 0, 1, or 2. The “code” fold provides the starting code. You must implement the models using the starting code.

1. (15 points) Data Preprocessing: In this problem, you need to finish “code/DataReader.py”.
 - (a) Explain what the function *train_valid_split* does and why we need this step.

- (b) Before testing, is it correct to re-train the model on the whole training set? Explain your answer.
- (c) In this assignment, we use two hand-crafted features:

The first feature is a measure of symmetry. For a 16×16 image x , it is defined as

$$F_{symmetry} = -\frac{\sum_{pixel} |x - \text{flip}(x)|}{256},$$

where 256 is the number of pixels and $\text{flip}(\cdot)$ means left and right flipping.

The second feature is a measure of intensity. For a 16×16 image x , it is defined as

$$F_{intensity} = \frac{\sum_{pixel} x}{256},$$

which is simply the average of pixel values.

Implement them in the function `prepare_X`.

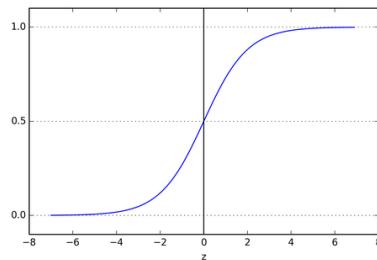
- (d) In the function `prepare_X`, there is a third feature which is always 1. Explain why we need it.
- (e) The function `prepare_y` is already finished. Note that the returned indices stores the indices for data from class 1 and 2. Only use these two classes for binary classification and convert the labels to +1 and -1 if necessary.
- (f) Test your code in “code/main.py” and visualize the training data from class 1 and 2 by implementing the function `visualize_features`. The visualization should not include the third feature. Therefore it is a 2-D scatter plot. Include the figure in your submission.

2. (25 points) Cross-entropy loss: In logistic regression, we use the cross-entropy loss.

- (a) Write the loss function $E(w)$ for one training data sample (x, y) . Note that the binary labels are 1 and -1.
- (b) Compute the gradient $\nabla E(w)$. Please provide intermediate steps of derivation.
- (c) Once the optimal w is obtained, it can be used to make predictions as follows:

$$\text{Predicted class of } x = \begin{cases} 1 & \text{if } \theta(w^T x) \geq 0.5 \\ -1 & \text{if } \theta(w^T x) < 0.5 \end{cases}$$

where the function $\theta(z) = \frac{1}{1+e^{-z}}$ looks like



However, this is not the most efficient way since the decision boundary is linear. Why? Explain it. When will we need to use the sigmoid function in prediction?

- (d) Is the decision boundary still linear if the prediction rule is changed to the following?
Justify briefly.

$$\text{Predicted label of } x = \begin{cases} 1 & \text{if } \theta(w^T x) \geq 0.9 \\ -1 & \text{if } \theta(w^T x) < 0.9 \end{cases}$$

- (e) In light of your answers to the above two questions, what is the essential property of logistic regression that results in the linear decision boundary?
3. (25 points) Sigmoid logistic regression: In this problem, you need to finish “code/LogisticRegression.py”.
Please follow the instructions in the starting code. Please use data from class 1 and 2 for the binary classification.
- (a) Based on (b) in the last problem, implement the function *_gradient*.
 - (b) There are different ways to train a logistic regression model. In this assignment, you need to implement batch gradient descent, stochastic gradient descent and mini-batch gradient descent in the functions *fit_BGD*, *fit_SGD* and *fit_miniBGD*, respectively. Note that: In batch gradient descent, each model update is based on all training samples. In stochastic gradient descent, each model update is based on one training sample. In mini-batch gradient descent, you split your training data into (roughly equal-sized) mini-batches, and each model update is based on training samples in each mini-batch. Thus, BGD and SDG are actually special cases of mini-BGD. An *epoch* is completed when each training sample is used to update the model exactly once. In SDG and mini-BGD, you might want to reshuffle your samples before each epoch.
 - (c) Implement the functions *predict* and *score* for prediction and evaluation, respectively. Additionally, please implement the function *predict_proba* which outputs the probabilities of both classes.
 - (d) Test your code in “code/main.py” and visualize the results after training by using the function *visualize_results*. The visualization should include the 2-D scatter plot of training data from class 1 and 2 and your decision boundary. Include the figure in your submission.
 - (e) Implement the testing process and report the test accuracy of your best logistic regression model.
4. (25 points) Softmax logistic regression: In this problem, you need to finish “code/LRM.py”.
Please follow the instructions in the starting code.
- (a) Based on the course notes, implement the function *_gradient*.
 - (b) In this assignment, you only need to implement mini-batch gradient descent in the function *fit_miniBGD*.
 - (c) Implement the functions *predict* and *score* for prediction and evaluation, respectively.
 - (d) Test your code in “code/main.py” and visualize the results after training by using the function *visualize_results_multi*. The visualization should include the 2-D scatter plot of training data from class 0, 1, and 2, and your decision boundaries. Include the figure in your submission.
 - (e) Implement the testing process and report the test accuracy of your best logistic regression model.

5. (10 points) Softmax logistic vs Sigmoid logistic: In this problem, you need to experimentally compare these two methods. **Please follow the instructions in the starting code.** Use data examples from class 1 and 2 for classification. Train the softmax logistic classifier and the sigmoid logistic classifier using the same data until convergence. Compare these two classifiers and report your observations and insights.

Linear Models for Handwritten Digits Classification: In this assignment, you will implement the binary logistic regression model and multi-class logistic regression model on a partial dataset from MNIST. In this classification task, the model will take a 16×16 image of handwritten digits as inputs and classify the image into different classes. For the binary case, the classes are 1 and 2 while for the multi-class case, the classes are 0, 1, and 2. The "data" folder contains the dataset which has already been split into a training set and a testing set. All data examples are saved in dictionary-like objects using "npz" file. For each data sample, the dictionary key 'x' indicates its raw features, which are represented by a 256-dimensional vector where the values between $[-1, 1]$ indicate grayscale pixel values for a 16×16 image. In addition, the key 'y' is the label for a data example, which can be 0, 1, or 2. The "code" folder provides the starting code. You must implement the models using the starting code.

1. (15 points) Data Preprocessing: In this problem, you need to finish "code/DataReader.py".

(a) Explain what the function `train_valid_split` does and why we need this step.

(b) Before testing, is it correct to re-train the model on the whole training set? Explain your answer.

(c) In this assignment, we use two hand-crafted features:

The first feature is a measure of symmetry. For a 16×16 image x , it is defined as

$$F_{\text{symmetry}} = -\frac{\sum_{\text{pixel}} |x - \text{flip}(x)|}{256},$$

where 256 is the number of pixels and `flip()` means left and right flipping.

The second feature is a measure of intensity. For a 16×16 image x , it is defined as

$$F_{\text{intensity}} = \frac{\sum_{\text{pixel}} x}{256},$$

which is simply the average of pixel values.

Implement them in the function `prepare_X`.

(d) In the function `prepare_X`, there is a third feature which is always 1. Explain why we need it.

(e) The function `prepare_y` is already finished. Note that the returned indices stores the indices for data from class 1 and 2. Only use these two classes for binary classification and convert the labels to +1 and -1 if necessary.

(f) Test your code in "code/main.py" and visualize the training data from class 1 and 2 by implementing the function `visualize_features`. The visualization should not include the third feature. Therefore it is a 2-D scatter plot. Include the figure in your submission.

Written Analysis of Answers / Programming

1a] The function `train_valid_split` separates what it does into training sets and validation sets at a specified [split_index]

Why T. Model Evaluation

- Hyperparameter Tuning
- allows assessing model performance on unseen data before testing
- Validation set helps vs detect overfitting
- yields hyperparameter selection

1b] Retraining on entire training set before training is ~~not~~ After validating hyperparameters (e.g. learning rate) on the validation set

→ combines all available data [training + validation)
maximizes information used for the final model
This improving generalization to the test set!

1c] Feature 1 := Symmetry

↳ 2 := Intensity

↳ 3 := Bias (will be 1)

3]

`np.ones(raw_X.shape[0])`

1d]

see Code Submission forrest

See LaTeX
append

1d] The bias term is always 1 to represent the intercept w_0 in the linear model

$$w^T x = b$$

W/O it, the decision boundary is forced through the origin

→ Limits Flexibility of model!

1f] The Scatterplot below visualizes the crafted features

→ Symmetry [x-axis]
→ Intensity [y-axis]

for the binary classification task

Class 1: Red Points (+1)
Class 2: Blue Points (-1)

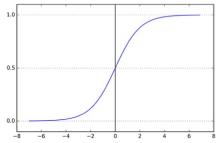
• Trend: Potential Linear Separability

2. (25 points) Cross-entropy loss: In logistic regression, we use the cross-entropy loss.

- (a) Write the loss function $E(w)$ for one training data sample (x, y) . Note that the binary labels are 1 and -1.
- (b) Compute the gradient $\nabla E(w)$. Please provide intermediate steps of derivation.
- (c) Once the optimal w is obtained, it can be used to make predictions as follows:

$$\text{Predicted class of } x = \begin{cases} 1 & \text{if } \theta(w^T x) \geq 0.5 \\ -1 & \text{if } \theta(w^T x) < 0.5 \end{cases}$$

where the function $\theta(z) = \frac{1}{1+e^{-z}}$ looks like



- (d) Is the decision boundary still linear if the prediction rule is changed to the following? Justify briefly.

$$\text{Predicted label of } x = \begin{cases} 1 & \text{if } \theta(w^T x) \geq 0.9 \\ -1 & \text{if } \theta(w^T x) < 0.9 \end{cases}$$

- (e) In light of your answers to the above two questions, what is the essential property of logistic regression that results in the linear decision boundary?

2a] Loss Function: Cross-entropy loss function $E(w)$ for one training data sample (x, y)

$$E(w) = \ln(1 + e^{-y w^T x})$$

$w^T x$ is either $[+1]^T$ or $[-1]^T$

2b] Gradient Computation

we must use chain rule.

$$\rightarrow \text{Let } z = w^T x$$

$$\rightarrow \text{we know } E = \ln(1 + e^{-y z})$$

$$\rightarrow \frac{dE}{dz} = \frac{-ye^{-yz}}{1 + e^{-yz}} \rightarrow \frac{dz}{du} = x$$

$$\Rightarrow \frac{dE}{dw} = \frac{dE}{dz} \cdot \frac{dz}{dw} = \frac{-y x}{1 + e^{-y w^T x}}$$

\therefore the gradient is

$$\nabla E(w) = \frac{-y x}{1 + e^{-y w^T x}}$$

The given prediction rule is
not the most efficient b/c:

→ The decision boundary is 'linear': $w^T x = 0$

→ computing $\theta(w^T x)$ is unnecessary for classification

→ we can use the sign of $w^T x$ for prediction

Sigmoid Function $\theta(z)$ is needed for

Probability estimates

not just class predictions

2c] Essential Property of Logistic Regression

→ Linearity of $w^T x$

The Sigmoid function applies a monotonic transformation to this linear combination which preserves the linearity of the decision boundary.

Thus Logistic Regression can only separate classes w/ a linear boundary

→ simple [easily interpretable]

→ inability to model non-linear relationships

2d] Given Alternative Prediction Rule:

$$\begin{cases} 1 & \text{if } g(w^T x) \geq a \\ -1 & \text{if } g(w^T x) < a \end{cases}$$

⇒ The decision boundary is still linear

It's defined by $w^T x = \ln a$

[i.e. a linear equation
just shifted from origin.]

3. (25 points) Sigmoid logistic regression: In this problem, you need to finish "code/LogisticRegression.py". Please follow the instructions in the starting code. Please use data from class 1 and 2 for the binary classification.

- Based on (b) in the last problem, implement the function `_gradient`.
 - There are different ways to train a logistic regression model. In this assignment, you need to implement batch gradient descent, stochastic gradient descent and mini-batch gradient descent in the functions `fit_BGD`, `fit_SGD` and `fit_miniBGD`, respectively.
- Note that: In batch gradient descent, each model update is based on all training samples. In stochastic gradient descent, each model update is based on one training sample. In mini-batch gradient descent, you split your training data into (roughly equal-sized) mini-batches, and each model update is based on training samples in each mini-batch. Thus, BGD and SDG are actually special cases of mini-BGD. An epoch is completed when each training sample is used to update the model exactly once. In SDG and mini-BGD, you might want to reshuffle your samples before each epoch.
- Implement the functions `predict` and `score` for prediction and evaluation, respectively. Additionally, please implement the function `predict_proba` which outputs the probabilities of both classes.
 - Test your code in "code/main.py" and visualize the results after training by using the function `visualize_results`. The visualization should include the 2-D scatter plot of training data from class 1 and 2 and your decision boundary. Include the figure in your submission.
 - Implement the testing process and report the test accuracy of your best logistic regression model.

3a] Must implement gradient function from 2b]

$$\frac{dE}{dw} = \frac{dE}{dz} \frac{dz}{dw} = \frac{-y}{1+e^{-w^T x}}$$

\therefore the gradient is $\nabla E(w) = \frac{-y}{1+e^{-w^T x}} x$

$\rightarrow y \rightarrow -1$ x is feature vector
 $\rightarrow 1$ w is weight vector

3b] Required to implement 3 methods

1. Batch Gradient Descent [fit_BGD]

- update model using all training samples in each iteration
- for each iteration we must:
 - compute gradient for all samples
 - update weights as

$$w = w - \text{learning-rate} \cdot \text{avg-gradient}$$

2. Stochastic Gradient Descent (SGD) [fit_SGD]

- for each iteration
- Randomly select a sample
- Compute gradient for this sample
- update weights:

$$w = w - \text{learning-rate} \cdot \text{gradient}$$

3. Mini-Batch Gradient Descent [fit_miniBGD]

for each iter:

- randomly select a 'mini-batch' of samples
- compute the average gradient for this 'mini-batch'
- update weights:

$$w = w - \text{Learningrate} \cdot \text{gradient}$$

3c] Prediction / Evaluation functions

1) 'predict' function

- input: feature matrix X
- output: predicted labels -1 or 1
- implementation: return $\text{sign}(X \cdot w)$

2) 'score' function

- input: feature matrix X AND true labels
- output: array [% of correct predictions]
- returns % of correct predictions

3) 'predict_proba'

- feature matrix X [Input]
- probabilities for both classes [out]
- sketch implementation
 - a) compute $z = X \cdot w$
 - ⇒ return $\left[\frac{1}{1+e^z}, \frac{1}{1+e^z} \right]$

3d] See Lated solution
and Code Submission

Process

- Testing occurs in 'main.py'

4. (25 points) Softmax logistic regression: In this problem, you need to finish “code/LRM.py”. Please follow the instructions in the starting code.

- (a) Based on the course notes, implement the function `_gradient`.
- (b) In this assignment, you only need to implement mini-batch gradient descent in the function `fit_miniBGD`.
- (c) Implement the functions `predict` and `score` for prediction and evaluation, respectively.
- (d) Test your code in “code/main.py” and visualize the results after training by using the function `visualize_results_multi`. The visualization should include the 2-D scatter plot of training data from class 0, 1, and 2, and your decision boundaries. Include the figure in your submission.
- (e) Implement the testing process and report the test accuracy of your best logistic regression model.

4a] See Latex/Code Submission

4b] See Latex/Code Submission

4c] See Latex/Code Submission

4d] See Latex and Code Submission

4e] See Latex and Code Submission

5. (10 points) Softmax logistic vs Sigmoid logistic: In this problem, you need to experimentally compare these two methods. Please follow the instructions in the starting code. Use data examples from class 1 and 2 for classification. Train the softmax logistic classifier and the sigmoid logistic classifier using the same data until convergence. Compare these two classifiers and report your observations and insights.

[See LaTeX for full breakdown]

Gen. Summary
→

1. Experimental Setup

- Used digits 1 and 2 for binary classification
- Identical hyperparameters:
 - * Learning rate = 0.1
 - * Max iterations = 10000
 - * Batch size = 50
- Features: Symmetry and Intensity

2. Performance Results

Softmax:

- Training accuracy: 97.19%
- Validation accuracy: 97.88%

Sigmoid:

- Training accuracy: 97.11%
- Validation accuracy: 97.88%

3. Key Observations

a) Accuracy:

- Nearly identical performance
- Difference in training accuracy is minimal (0.08%)
- Identical validation accuracy suggests equivalent generalization

b) Decision Boundaries:

- Both methods produced similar linear boundaries
- Boundaries align closely in feature space
- Effective separation between classes

4. Analysis

a) Theoretical Confirmation:

- Results validate theoretical equivalence for binary case
- Both methods converged to similar solutions
- Decision boundaries are mathematically equivalent

b) Implementation Differences:

- Softmax requires more computation (exponential terms)
- Both achieve same end result despite different approaches
- Choice between methods for binary classification should be based on implementation preference

5. Conclusion

The experimental results confirm that sigmoid and softmax logistic regression are functionally equivalent for binary classification, producing nearly identical performance metrics and decision boundaries.

Output

```
[ 0.18761413 3.55944904 -2.06327285]
[ 0.9444444444444444
[ 0.18761413 3.55944904 -2.06327285]
[ 0.9444444444444444
[10.82741562 31.69274308 2.09324186]
[ 0.9711111111111111
[11.05112829 31.87766955 2.24934715]
[ 0.9703703703703703
[ 5.09262961 23.55262757 -3.14116836]
[ 0.9733333333333334
Best parameters: learning_rate=0.5, batch_size=1
Best validation accuracy: 0.9788359788359788
Test accuracy: 0.9307359307359307
[[ 6.75513933 -0.98294001 -5.7532125 ]
[ 0.27435837 15.72321774 -15.98338787]
[ [10.43217196 -8.600453 -1.82368009]]
0.8960869565217391
```

```
Training multi-class classifier...
Best multi-class parameters: learning_rate=0.1, batch_size=100
Best multi-class validation accuracy: 0.8746031746031746
```

Training Softmax classifier...

Training Sigmoid classifier...

```
Classifier Comparison Results:
Softmax Classifier:
Training accuracy: 0.9719
Validation accuracy: 0.9788
```

```
Sigmoid Classifier:
Training accuracy: 0.9711
Validation accuracy: 0.9788
```