

Deep Learning - Musterlösung Übung 2

Multi-Layer Perceptrons und Backpropagation

Fachhochschule Südwestfalen

23. Oktober 2025

Hinweise zur Musterlösung

Diese Musterlösung bietet detaillierte mathematische Herleitungen und vollständige Implementierungen. Alternative Lösungsansätze sind oft ebenfalls korrekt.

1 Backpropagation-Algorithmus - Lösungen

1.1 Aufgabe 1.1: Mathematische Herleitung

Gegeben ist ein 3-Schicht-MLP mit Binary Cross-Entropy Loss:

$$L = -[y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})] \quad (1)$$

Teil (a): Gradient der Ausgabe

$$\frac{\partial L}{\partial z^{(3)}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z^{(3)}} \quad (2)$$

$$\frac{\partial L}{\partial \hat{y}} = -\frac{y}{\hat{y}} + \frac{1-y}{1-\hat{y}} \quad (3)$$

$$\frac{\partial \hat{y}}{\partial z^{(3)}} = \sigma(z^{(3)})(1 - \sigma(z^{(3)})) = \hat{y}(1 - \hat{y}) \quad (4)$$

Einsetzen:

$$\frac{\partial L}{\partial z^{(3)}} = \left(-\frac{y}{\hat{y}} + \frac{1-y}{1-\hat{y}} \right) \hat{y}(1 - \hat{y}) \quad (5)$$

$$= -y(1 - \hat{y}) + (1 - y)\hat{y} \quad (6)$$

$$= -y + y\hat{y} + \hat{y} - y\hat{y} \quad (7)$$

$$= \hat{y} - y \quad (8)$$

$$\frac{\partial L}{\partial z^{(3)}} = \hat{y} - y$$

Teil (b): Gradienten der Ausgabeschicht

$$\frac{\partial L}{\partial \mathbf{w}^{(3)}} = \frac{\partial L}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial \mathbf{w}^{(3)}} = (\hat{y} - y) \mathbf{a}^{(2)} \quad (9)$$

$$\frac{\partial L}{\partial b^{(3)}} = \frac{\partial L}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial b^{(3)}} = \hat{y} - y \quad (10)$$

Teil (c): Gradient für vorherige Schicht

$$\frac{\partial L}{\partial \mathbf{a}^{(2)}} = \frac{\partial L}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial \mathbf{a}^{(2)}} = (\hat{y} - y) \mathbf{w}^{(3)} \quad (11)$$

Teil (d): Gradienten für Hidden Layer 2

$$\frac{\partial L}{\partial \mathbf{z}^{(2)}} = \frac{\partial L}{\partial \mathbf{a}^{(2)}} \odot \frac{\partial \mathbf{a}^{(2)}}{\partial \mathbf{z}^{(2)}} \quad (12)$$

$$= (\hat{y} - y) \mathbf{w}^{(3)} \odot \mathbf{a}^{(2)} \odot (1 - \mathbf{a}^{(2)}) \quad (13)$$

Definiere $\delta^{(2)} = \frac{\partial L}{\partial \mathbf{z}^{(2)}}$:

$$\frac{\partial L}{\partial \mathbf{W}^{(2)}} = \delta^{(2)} (\mathbf{a}^{(1)})^T \quad (14)$$

$$\frac{\partial L}{\partial \mathbf{b}^{(2)}} = \delta^{(2)} \quad (15)$$

Teil (e): Gradienten für Hidden Layer 1

$$\frac{\partial L}{\partial \mathbf{a}^{(1)}} = (\mathbf{W}^{(2)})^T \delta^{(2)} \quad (16)$$

$$\delta^{(1)} = \frac{\partial L}{\partial \mathbf{a}^{(1)}} \odot \mathbf{a}^{(1)} \odot (1 - \mathbf{a}^{(1)}) \quad (17)$$

$$\frac{\partial L}{\partial \mathbf{W}^{(1)}} = \delta^{(1)} \mathbf{x}^T \quad (18)$$

$$\frac{\partial L}{\partial \mathbf{b}^{(1)}} = \delta^{(1)} \quad (19)$$

1.2 Aufgabe 1.2: Numerisches Beispiel

Gegeben:

$$\mathbf{x} = \begin{pmatrix} 0.5 \\ 0.8 \end{pmatrix}, \quad y = 1 \quad (20)$$

$$\mathbf{W}^{(1)} = \begin{pmatrix} 0.2 & 0.1 \\ -0.3 & 0.4 \\ 0.6 & -0.2 \end{pmatrix}, \quad \mathbf{b}^{(1)} = \begin{pmatrix} 0.1 \\ -0.2 \\ 0.3 \end{pmatrix} \quad (21)$$

Forward Pass:

Layer 1:

$$\mathbf{z}^{(1)} = \mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)} \quad (22)$$

$$= \begin{pmatrix} 0.2 & 0.1 \\ -0.3 & 0.4 \\ 0.6 & -0.2 \end{pmatrix} \begin{pmatrix} 0.5 \\ 0.8 \end{pmatrix} + \begin{pmatrix} 0.1 \\ -0.2 \\ 0.3 \end{pmatrix} \quad (23)$$

$$= \begin{pmatrix} 0.18 \\ 0.12 \\ 0.14 \end{pmatrix} \quad (24)$$

$$\mathbf{a}^{(1)} = \sigma(\mathbf{z}^{(1)}) = \begin{pmatrix} 0.545 \\ 0.530 \\ 0.535 \end{pmatrix} \quad (25)$$

Layer 2:

$$\mathbf{z}^{(2)} = \mathbf{W}^{(2)} \mathbf{a}^{(1)} + \mathbf{b}^{(2)} \quad (26)$$

$$= \begin{pmatrix} 0.4 & -0.1 & 0.3 \\ 0.2 & 0.5 & -0.4 \end{pmatrix} \begin{pmatrix} 0.545 \\ 0.530 \\ 0.535 \end{pmatrix} + \begin{pmatrix} 0.1 \\ -0.1 \end{pmatrix} \quad (27)$$

$$= \begin{pmatrix} 0.419 \\ 0.195 \end{pmatrix} \quad (28)$$

$$\mathbf{a}^{(2)} = \sigma(\mathbf{z}^{(2)}) = \begin{pmatrix} 0.603 \\ 0.549 \end{pmatrix} \quad (29)$$

Output Layer:

$$z^{(3)} = \mathbf{w}^{(3)T} \mathbf{a}^{(2)} + b^{(3)} \quad (30)$$

$$= 0.6 \cdot 0.603 + (-0.3) \cdot 0.549 + 0.2 \quad (31)$$

$$= 0.397 \quad (32)$$

$$\hat{y} = \sigma(0.397) = 0.598 \quad (33)$$

Loss:

$$L = -[1 \cdot \log(0.598) + 0 \cdot \log(0.402)] = -\log(0.598) = 0.515 \quad (34)$$

Backward Pass:

Output Layer:

$$\frac{\partial L}{\partial z^{(3)}} = 0.598 - 1 = -0.402 \quad (35)$$

$$\frac{\partial L}{\partial \mathbf{w}^{(3)}} = -0.402 \begin{pmatrix} 0.603 \\ 0.549 \end{pmatrix} = \begin{pmatrix} -0.242 \\ -0.221 \end{pmatrix} \quad (36)$$

$$\frac{\partial L}{\partial b^{(3)}} = -0.402 \quad (37)$$

Hidden Layer 2:

$$\frac{\partial L}{\partial \mathbf{a}^{(2)}} = -0.402 \begin{pmatrix} 0.6 \\ -0.3 \end{pmatrix} = \begin{pmatrix} -0.241 \\ 0.121 \end{pmatrix} \quad (38)$$

$$\delta^{(2)} = \begin{pmatrix} -0.241 \\ 0.121 \end{pmatrix} \odot \begin{pmatrix} 0.603 \cdot 0.397 \\ 0.549 \cdot 0.451 \end{pmatrix} = \begin{pmatrix} -0.058 \\ 0.030 \end{pmatrix} \quad (39)$$

Hidden Layer 1:

$$\frac{\partial L}{\partial \mathbf{a}^{(1)}} = \begin{pmatrix} 0.4 & 0.2 \\ -0.1 & 0.5 \\ 0.3 & -0.4 \end{pmatrix} \begin{pmatrix} -0.058 \\ 0.030 \end{pmatrix} = \begin{pmatrix} -0.017 \\ 0.021 \\ -0.029 \end{pmatrix} \quad (40)$$

2 MLP-Implementierung - Musterlösung

2.1 Aufgabe 2.1: MLP-Klasse

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 class MLP:
5     def __init__(self, layer_sizes, activation='sigmoid',
6                 learning_rate=0.01):
7         self.layer_sizes = layer_sizes
8         self.activation_name = activation
9         self.learning_rate = learning_rate
10        self.num_layers = len(layer_sizes)
11
12        # Initialize weights and biases
13        self.weights = {}
14        self.biases = {}
15        self._initialize_weights()
16
17        # Store for backpropagation
18        self.cache = {}
19        self.losses = []
20
21    def _initialize_weights(self):
22        """Xavier/He initialization"""
23        for i in range(1, self.num_layers):
24            # Xavier initialization for sigmoid, He for ReLU
25            if self.activation_name == 'relu':
26                # He initialization
27                std = np.sqrt(2.0 / self.layer_sizes[i-1])
28            else:
29                # Xavier initialization
30                std = np.sqrt(1.0 / self.layer_sizes[i-1])
31
32            self.weights[i] = np.random.normal(0, std,
33                (self.layer_sizes[i], self.layer_sizes[i-1]))
34            self.biases[i] = np.zeros((self.layer_sizes[i], 1))

```

```
34
35 def _activation(self, z):
36     """Activation function"""
37     if self.activation_name == 'sigmoid':
38         return self._sigmoid(z)
39     elif self.activation_name == 'relu':
40         return self._relu(z)
41     elif self.activation_name == 'tanh':
42         return np.tanh(z)
43     else:
44         raise ValueError(f"Unknown activation: {self.
45                             activation_name}")
46
47 def _activation_derivative(self, z):
48     """Derivative of activation function"""
49     if self.activation_name == 'sigmoid':
50         s = self._sigmoid(z)
51         return s * (1 - s)
52     elif self.activation_name == 'relu':
53         return (z > 0).astype(float)
54     elif self.activation_name == 'tanh':
55         return 1 - np.tanh(z)**2
56     else:
57         raise ValueError(f"Unknown activation: {self.
58                             activation_name}")
59
60 def _sigmoid(self, z):
61     """Numerically stable sigmoid"""
62     return np.where(z >= 0,
63                     1 / (1 + np.exp(-z)),
64                     np.exp(z) / (1 + np.exp(z)))
65
66 def _relu(self, z):
67     return np.maximum(0, z)
68
69 def _forward(self, X):
70     """Forward pass with caching for backpropagation"""
71     self.cache = {}
72     A = X.T # Shape: (features, samples)
73     self.cache[0] = A
74
75     for i in range(1, self.num_layers):
76         Z = self.weights[i] @ A + self.biases[i]
77         if i == self.num_layers - 1: # Output layer
78             A = self._sigmoid(Z) # Always sigmoid for binary
79                                     classification
80         else: # Hidden layers
81             A = self._activation(Z)
82
83     self.cache[i] = {'Z': Z, 'A': A}
84     A = self.cache[i]['A']
```

```

82
83     return A
84
85     def _backward(self, X, y):
86         """Backward pass - compute gradients"""
87         m = X.shape[0] # Number of samples
88         y = y.reshape(1, -1) # Shape: (1, samples)
89
90         # Get final output
91         A_final = self.cache[self.num_layers - 1]['A']
92
93         # Gradients storage
94         gradients = {}
95
96         # Output layer gradient (assuming binary cross-entropy)
97         dZ = A_final - y # Shape: (1, samples)
98
99         for i in range(self.num_layers - 1, 0, -1):
100             # Get previous layer activation
101             A_prev = self.cache[i-1] if i == 1 else self.cache[i-1]['A']
102
103             # Compute gradients
104             gradients[f'dW{i}'] = (1/m) * dZ @ A_prev.T
105             gradients[f'db{i}'] = (1/m) * np.sum(dZ, axis=1, keepdims=True)
106
107             if i > 1: # Not input layer
108                 # Compute dZ for previous layer
109                 dA_prev = self.weights[i].T @ dZ
110                 Z_prev = self.cache[i-1]['Z']
111                 dZ = dA_prev * self._activation_derivative(Z_prev)
112
113         return gradients
114
115     def _update_parameters(self, gradients):
116         """Update weights and biases using gradients"""
117         for i in range(1, self.num_layers):
118             self.weights[i] -= self.learning_rate * gradients[f'dW{i}']
119             self.biases[i] -= self.learning_rate * gradients[f'db{i}']
120
121     def _compute_loss(self, y_true, y_pred):
122         """Binary cross-entropy loss"""
123         m = y_true.shape[0]
124         # Clip predictions to prevent log(0)
125         y_pred = np.clip(y_pred, 1e-15, 1 - 1e-15)
126         loss = -(1/m) * np.sum(y_true * np.log(y_pred) +
127                                (1 - y_true) * np.log(1 - y_pred))
128         return loss

```

```
129
130 def train(self, X, y, epochs=1000, batch_size=None, verbose=False
131 ):
132     """Training with mini-batch gradient descent"""
133     if batch_size is None:
134         batch_size = X.shape[0] # Full batch
135
136     for epoch in range(epochs):
137         # Shuffle data
138         indices = np.random.permutation(X.shape[0])
139         X_shuffled = X[indices]
140         y_shuffled = y[indices]
141
142         epoch_loss = 0
143         num_batches = 0
144
145         # Mini-batch training
146         for i in range(0, X.shape[0], batch_size):
147             X_batch = X_shuffled[i:i+batch_size]
148             y_batch = y_shuffled[i:i+batch_size]
149
150             # Forward pass
151             y_pred = self._forward(X_batch)
152
153             # Compute loss
154             loss = self._compute_loss(y_batch, y_pred.T)
155             epoch_loss += loss
156             num_batches += 1
157
158             # Backward pass
159             gradients = self._backward(X_batch, y_batch)
160
161             # Update parameters
162             self._update_parameters(gradients)
163
164             # Average loss for epoch
165             avg_loss = epoch_loss / num_batches
166             self.losses.append(avg_loss)
167
168             if verbose and (epoch + 1) % 100 == 0:
169                 print(f"Epoch {epoch+1}/{epochs}, Loss: {avg_loss:.4f
170                       }")
171
172 def predict(self, X):
173     """Make predictions"""
174     y_pred = self._forward(X)
175     return (y_pred.T > 0.5).astype(int).flatten()
176
177 def predict_proba(self, X):
178     """Predict probabilities"""
179     return self._forward(X).T.flatten()
```

```
178
179     def score(self, X, y):
180         """Compute accuracy"""
181         predictions = self.predict(X)
182         return np.mean(predictions == y)
183
184     def plot_loss(self):
185         """Plot training loss"""
186         plt.figure(figsize=(10, 6))
187         plt.plot(self.losses)
188         plt.title('Training Loss')
189         plt.xlabel('Epoch')
190         plt.ylabel('Loss')
191         plt.grid(True)
192         plt.show()
```

2.2 Aufgabe 2.2: Experimentelle Evaluierung

XOR-Problem lösen:

```
1 # XOR Dataset
2 X_xor = np.array([[0, 0], [0, 1], [1, 0], [1, 1]])
3 y_xor = np.array([0, 1, 1, 0])
4
5 # Create and train MLP
6 mlp_xor = MLP([2, 4, 1], activation='sigmoid', learning_rate=0.1)
7 mlp_xor.train(X_xor, y_xor, epochs=5000, verbose=True)
8
9 # Test predictions
10 print("XOR Results:")
11 for i in range(len(X_xor)):
12     pred = mlp_xor.predict_proba(X_xor[i:i+1])[0]
13     print(f"Input: {X_xor[i]}, Target: {y_xor[i]}, Prediction: {pred:.3f}")
14
15 # Plot decision boundary
16 def plot_decision_boundary(model, X, y, title="Decision Boundary"):
17     plt.figure(figsize=(10, 8))
18
19     # Create mesh
20     h = 0.01
21     x_min, x_max = X[:, 0].min() - 0.1, X[:, 0].max() + 0.1
22     y_min, y_max = X[:, 1].min() - 0.1, X[:, 1].max() + 0.1
23     xx, yy = np.meshgrid(np.arange(x_min, x_max, h),
24                           np.arange(y_min, y_max, h))
25
26     # Predict on mesh
27     mesh_points = np.c_[xx.ravel(), yy.ravel()]
28     Z = model.predict_proba(mesh_points)
29     Z = Z.reshape(xx.shape)
30
```



```

31     # Plot
32     plt.contourf(xx, yy, Z, levels=50, alpha=0.8, cmap=plt.cm.RdYlBu)
33     plt.colorbar(label='Prediction Probability')
34
35     # Plot data points
36     colors = ['blue', 'red']
37     for i in range(len(X)):
38         plt.scatter(X[i, 0], X[i, 1], c=colors[y[i]], s=100,
39                     edgecolors='black')
40
41     plt.title(title)
42     plt.xlabel('x1')
43     plt.ylabel('x2')
44     plt.grid(True)
45     plt.show()
46 plot_decision_boundary(mlp_xor, X_xor, y_xor, "XOR Decision Boundary"
47 )
48 mlp_xor.plot_loss()

```

Spiralen-Datensatz:

```

1 def make_spirals(n_samples=200, noise=0.1):
2     """Create spiral dataset"""
3     t = np.linspace(0, 4*np.pi, n_samples//2)
4     x1 = t * np.cos(t) + noise * np.random.randn(n_samples//2)
5     y1 = t * np.sin(t) + noise * np.random.randn(n_samples//2)
6     x2 = -t * np.cos(t) + noise * np.random.randn(n_samples//2)
7     y2 = -t * np.sin(t) + noise * np.random.randn(n_samples//2)
8
9     X = np.vstack([np.column_stack([x1, y1]), np.column_stack([x2, y2
10 ])])
11     y = np.hstack([np.zeros(n_samples//2), np.ones(n_samples//2)])
12     return X, y
13
14 # Generate spiral data
15 X_spiral, y_spiral = make_spirals(n_samples=400, noise=0.3)
16
17 # Normalize data
18 X_spiral = (X_spiral - X_spiral.mean(axis=0)) / X_spiral.std(axis=0)
19
20 # Train MLP
21 mlp_spiral = MLP([2, 16, 8, 1], activation='relu', learning_rate
22                 =0.01)
23 mlp_spiral.train(X_spiral, y_spiral, epochs=2000, batch_size=32,
24                 verbose=True)
25
26 print(f"Spiral Accuracy: {mlp_spiral.score(X_spiral, y_spiral):.3f}")
27 plot_decision_boundary(mlp_spiral, X_spiral, y_spiral, "Spiral
28 Decision Boundary")

```

3 Vertiefende Fragen - Lösungen

3.1 Aufgabe 3.1: Vanishing Gradient Problem

Sigmoid-Problem: Die Sigmoid-Ableitung hat Maximum bei $x = 0$:

$$\sigma'(0) = \sigma(0)(1 - \sigma(0)) = 0.5 \cdot 0.5 = 0.25 \quad (41)$$

In tiefen Netzwerken werden Gradienten durch Multiplikation mit ≤ 0.25 exponentiell kleiner:

$$\frac{\partial L}{\partial W^{(1)}} \propto \prod_{i=2}^L W^{(i)} \sigma'(z^{(i)}) \leq (0.25)^{L-1} \quad (42)$$

ReLU-Lösung:

$$\text{ReLU}'(x) = \begin{cases} 1 & x > 0 \\ 0 & x \leq 0 \end{cases} \quad (43)$$

Für positive Aktivierungen ist der Gradient konstant 1, verhindert Vanishing.

Initialisierung:

- **Xavier:** $W \sim \mathcal{N}(0, \frac{1}{n_{in}})$ für Sigmoid/Tanh
- **He:** $W \sim \mathcal{N}(0, \frac{2}{n_{in}})$ für ReLU
- Verhindert zu große/kleine Aktivierungen

3.2 Aufgabe 3.2: Praktische Probleme

Overfitting-Demo:

```

1 # Small dataset
2 X_small = X_xor
3 y_small = y_xor
4
5 # Overparameterized network
6 mlp_over = MLP([2, 50, 50, 1], learning_rate=0.1)
7
8 # Train with validation split
9 val_losses = []
10 train_losses = []
11
12 for epoch in range(1000):
13     # Train
14     mlp_over.train(X_small, y_small, epochs=1, verbose=False)
15
16     # Track losses
17     if epoch % 10 == 0:
18         train_pred = mlp_over.predict_proba(X_small)
19         train_loss = mlp_over._compute_loss(y_small, train_pred)
20         train_losses.append(train_loss)
21
22     # Validation on same data (for demo)
23     val_losses.append(train_loss)

```

```
24 |
25 | # Plot overfitting
26 | plt.figure(figsize=(10, 6))
27 | plt.plot(range(0, 1000, 10), train_losses, label='Training Loss')
28 | plt.plot(range(0, 1000, 10), val_losses, label='Validation Loss')
29 | plt.title('Overfitting Demo')
30 | plt.xlabel('Epoch')
31 | plt.ylabel('Loss')
32 | plt.legend()
33 | plt.grid(True)
34 | plt.show()
```

Zusätzliche Implementierungshinweise

Numerische Stabilität

- Gradient clipping: $\nabla W = \text{clip}(\nabla W, -\theta, \theta)$
- Batch normalization vor Aktivierungen
- Learning rate scheduling: $\eta_t = \eta_0 / (1 + \alpha t)$

Debugging-Techniken

- Gradient checking: Numerische vs. analytische Gradienten
- Aktivierung-Monitoring: Histogramme der Aktivierungen
- Weight-Monitoring: Norm der Gewichtsmatrizen