

Министерство образования Республики Беларусь
Учреждение образования
“Брестский государственный технический университет”
Кафедра интеллектуально-информационных технологий

Интеллектуальный анализ данных
Лабораторная работа №3
Предобучение нейронных сетей с использованием автоэнкодерного
подхода

Выполнила:
студентка 4 курса
группы ИИ-24
Алешко А. В.
Проверила:
Андренко К. В.

Брест-2025

Цель работы: научиться осуществлять предобучение нейронных сетей с помощью автоэнкодерного подхода.

Общее задание:

1. Взять за основу любую сверточную или полносвязную архитектуру с количеством слоев более 3. Осуществить ее обучение (без предобучения) в соответствии с вариантом задания. Получить оценку эффективности модели, используя метрики, специфичные для решаемой задачи (например, MAPE – для регрессионной задачи или F1/Confusion matrix для классификационной).
2. Выполнить обучение с предобучением, используя автоэнкодерный подход, алгоритм которого изложен в лекции. Условие останова (например, по количеству эпох) при обучении отдельных слоев с использованием автоэнкодера выбрать самостоятельно.
3. Сравнить результаты, полученные при обучении с/без предобучения, сделать выводы.
4. Выполните пункты 1-3 для датасетов из ЛР 2 (Wisconsin Diagnostic Breast Cancer (WDBC), класс – 2 признак).
5. Оформить отчет по выполненной работе, загрузить исходный код и отчет в соответствующий репозиторий на github.

№	Выборка	Тип задачи	Целевая переменная
1	https://archive.ics.uci.edu/dataset/27/credit+approval	классификация	+/-

Код программы(вариант 1):

```
import pandas as pd
import numpy as np
import torch
import torch.nn as nn
import torch.optim as optim
from torch.utils.data import DataLoader, TensorDataset
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler,
OneHotEncoder
from sklearn.impute import SimpleImputer
from sklearn.metrics import f1_score, confusion_matrix
import matplotlib.pyplot as plt
import seaborn as sns
```

```

# Credit Approval dataset
def load_credit_data():
    url =
'https://archive.ics.uci.edu/static/public/27/data.csv'
    data = pd.read_csv(url)
    data = data.replace('?', np.nan)
    X = data.drop('A16', axis=1)
    y = data['A16'].map({'+': 1, '-': 0})
    categorical_cols = ['A1', 'A4', 'A5', 'A6', 'A7', 'A9',
'A10', 'A12', 'A13']
    continuous_cols = ['A2', 'A3', 'A8', 'A11', 'A14', 'A15']
    cat_imputer = SimpleImputer(strategy='most_frequent')
    cont_imputer = SimpleImputer(strategy='median')
    X[categorical_cols] =
cat_imputer.fit_transform(X[categorical_cols])
    X[continuous_cols] =
cont_imputer.fit_transform(X[continuous_cols])
    encoder = OneHotEncoder(sparse_output=False,
handle_unknown='ignore')
    X_cat_encoded = encoder.fit_transform(X[categorical_cols])
    X_cat_encoded = pd.DataFrame(X_cat_encoded,
columns=encoder.get_feature_names_out(categorical_cols))
    X = pd.concat([X[continuous_cols], X_cat_encoded], axis=1)
    scaler = StandardScaler()
    X = scaler.fit_transform(X)
    return X, y.values

```

```

# Breast Cancer dataset
def load_breast_cancer_data():
    url = 'https://archive.ics.uci.edu/ml/machine-learning-
databases/breast-cancer-wisconsin/wdbc.data'
    column_names = ['ID', 'Diagnosis'] + [f'feature_{i}' for i
in range(1, 31)]
    data = pd.read_csv(url, header=None, names=column_names)
    data = data.drop('ID', axis=1)
    y = data['Diagnosis'].map({'M': 1, 'B': 0}).values
    X = data.drop('Diagnosis', axis=1)
    scaler = StandardScaler()
    X = scaler.fit_transform(X)
    return X, y

```

```

# Архитектура НС(4 слоя)
class ClassificationNet(nn.Module):
    def __init__(self, input_size, hidden_sizes,
output_size=1):
        super(ClassificationNet, self).__init__()
        self.layers = nn.ModuleList()
        prev_size = input_size
        for h_size in hidden_sizes:
            self.layers.append(nn.Linear(prev_size, h_size))
            prev_size = h_size

```

```

        self.output = nn.Linear(prev_size, output_size)

    def forward(self, x):
        for layer in self.layers:
            x = torch.relu(layer(x))
        x = torch.sigmoid(self.output(x))
        return x

# Автоэнкодер
class Autoencoder(nn.Module):
    def __init__(self, input_size, hidden_size):
        super(Autoencoder, self).__init__()
        self.encoder = nn.Linear(input_size, hidden_size)
        self.decoder = nn.Linear(hidden_size, input_size)

    def forward(self, x):
        x = torch.relu(self.encoder(x))
        x = self.decoder(x)
        return x

# С автоэнкодером
def pretrain_layers(input_data, hidden_sizes, epochs=50,
                    lr=0.01):
    pretrained_weights = []
    current_input = input_data
    for h_size in hidden_sizes:
        ae = Autoencoder(current_input.shape[1], h_size)
        optimizer = optim.Adam(ae.parameters(), lr=lr)
        criterion = nn.MSELoss()
        dataset = TensorDataset(current_input, current_input)
        loader = DataLoader(dataset, batch_size=32,
                             shuffle=True)

        for epoch in range(epochs):
            for data, target in loader:
                optimizer.zero_grad()
                output = ae(data)
                loss = criterion(output, target)
                loss.backward()
                optimizer.step()
            with torch.no_grad():
                current_input =
torch.relu(ae.encoder(current_input))

    pretrained_weights.append((ae.encoder.weight.data.clone(),
ae.encoder.bias.data.clone()))
    return pretrained_weights

# Инициализация сети с готовыми весами
def init_with_pretrain(net, pretrained_weights):
    for i, (w, b) in enumerate(pretrained_weights):

```

```

        net.layers[i].weight.data = w
        net.layers[i].bias.data = b

# Обучения
def train_model(net, X_train, y_train, X_test, y_test,
epochs=100, lr=0.001, batch_size=32):
    criterion = nn.BCELoss()
    optimizer = optim.Adam(net.parameters(), lr=lr)
    train_dataset = TensorDataset(X_train, y_train)
    train_loader = DataLoader(train_dataset,
batch_size=batch_size, shuffle=True)
    losses = []
    for epoch in range(epochs):
        net.train()
        for data, target in train_loader:
            optimizer.zero_grad()
            output = net(data).squeeze()
            loss = criterion(output, target)
            loss.backward()
            optimizer.step()
        losses.append(loss.item())
    net.eval()
    with torch.no_grad():
        y_pred = (net(X_test).squeeze() >
0.5).float().cpu().numpy()
        f1 = f1_score(y_test.cpu().numpy(), y_pred)
        cm = confusion_matrix(y_test.cpu().numpy(), y_pred)
    return f1, cm, losses

def process_dataset(dataset_name, load_data_func):
    print(f"\n    {dataset_name} Dataset ")
    X, y = load_data_func()
    X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2, random_state=42)
    X_train_tensor = torch.tensor(X_train,
dtype=torch.float32)
    X_test_tensor = torch.tensor(X_test, dtype=torch.float32)
    y_train_tensor = torch.tensor(y_train,
dtype=torch.float32)
    y_test_tensor = torch.tensor(y_test, dtype=torch.float32)
    input_size = X_train.shape[1]
    hidden_sizes = [64, 32, 16]
    output_size = 1

    # 1. Без предобучение обучение
    net_no_pretrain = ClassificationNet(input_size,
hidden_sizes, output_size)
    f1_no, cm_no, losses_no = train_model(net_no_pretrain,
X_train_tensor, y_train_tensor, X_test_tensor, y_test_tensor)
    print("\nWithout Pretraining:")
    print("F1 Score:", f1_no)

```

```

print("Confusion Matrix:\n", cm_no)

# 2. С предобучением обучение
pretrained_weights = pretrain_layers(X_train_tensor,
hidden_sizes, epochs=50, lr=0.01)
net_pretrain = ClassificationNet(input_size, hidden_sizes,
output_size)
init_with_pretrain(net_pretrain, pretrained_weights)
f1_pre, cm_pre, losses_pre = train_model(net_pretrain,
X_train_tensor, y_train_tensor, X_test_tensor, y_test_tensor)
print("\nWith Pretraining:")
print("F1 Score:", f1_pre)
print("Confusion Matrix:\n", cm_pre)

# 3. Сравнение
print("\nComparison:")
print(f"F1 without pretrain: {f1_no:.4f} | F1 with
pretrain: {f1_pre:.4f}")
plt.figure(figsize=(10, 5))
plt.plot(losses_no, label='No Pretrain')
plt.plot(losses_pre, label='With Pretrain')
plt.title(f'Loss Curves - {dataset_name}')
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.legend()
plt.show()

fig, axes = plt.subplots(1, 2, figsize=(12, 5))
sns.heatmap(cm_no, annot=True, fmt='d', ax=axes[0],
cmap='Blues')
axes[0].set_title(f'Confusion Matrix - No Pretrain
({dataset_name})')
sns.heatmap(cm_pre, annot=True, fmt='d', ax=axes[1],
cmap='Blues')
axes[1].set_title(f'Confusion Matrix - With Pretrain
({dataset_name})')
plt.show()

return f1_no, f1_pre

if __name__ == "__main__":
    # Credit Approval dataset
    f1_no_credit, f1_pre_credit = process_dataset("Credit
Approval", load_credit_data)
    # Breast Cancer dataset
    f1_no_breast, f1_pre_breast = process_dataset("Breast
Cancer Wisconsin", load_breast_cancer_data)
    print("\n    Final Comparison Across Datasets ")
    print(f"Credit Approval - F1 without pretrain:
{f1_no_credit:.4f} | F1 with pretrain: {f1_pre_credit:.4f}")

```

```
print(f"Breast Cancer - F1 without pretrain:  
{f1_no_breast:.4f} | F1 with pretrain: {f1_pre_breast:.4f}")
```

Результат работы программы:

Credit Approval Dataset

Without Pretraining:

F1 Score: 0.7727272727272727

Confusion Matrix:

[[57 11]

[19 51]]

With Pretraining:

F1 Score: 0.8059701492537313

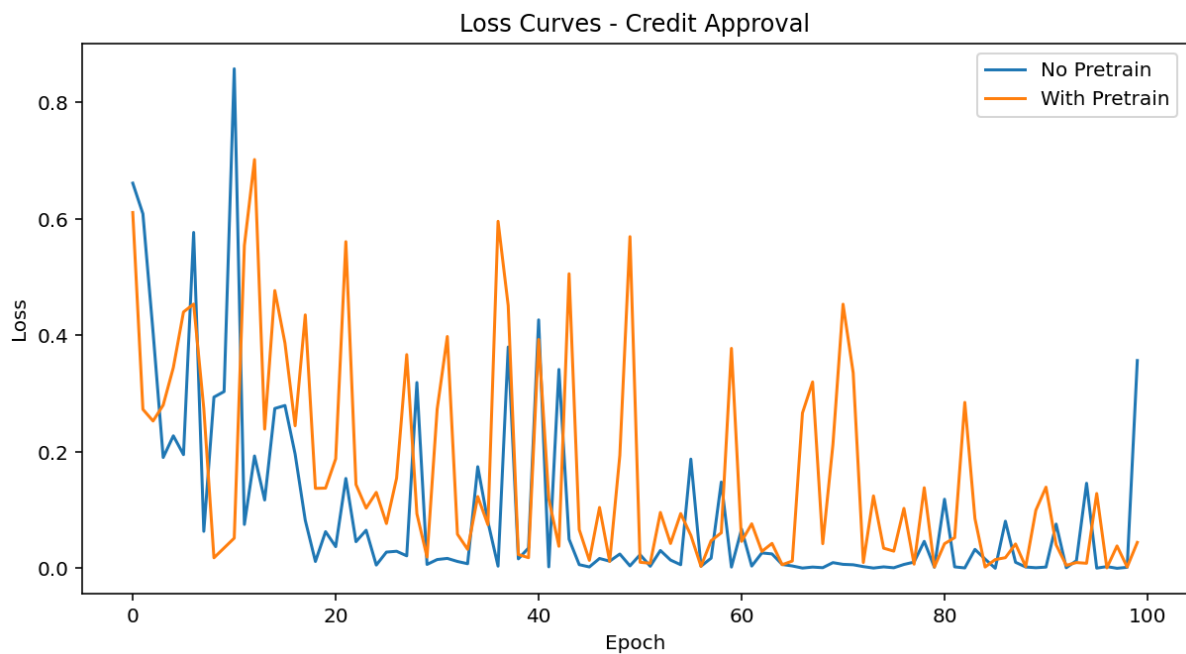
Confusion Matrix:

[[58 10]

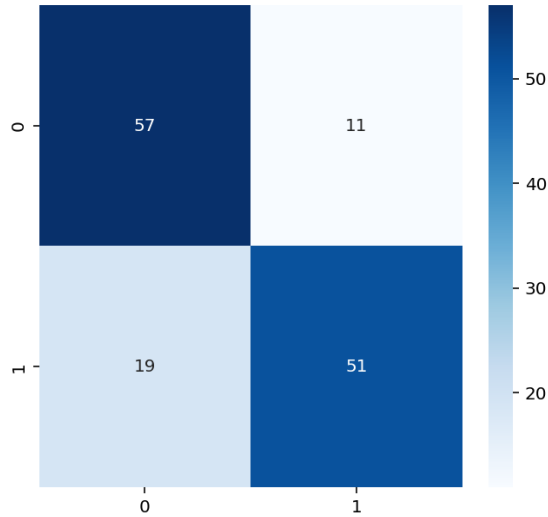
[16 54]]

Comparison:

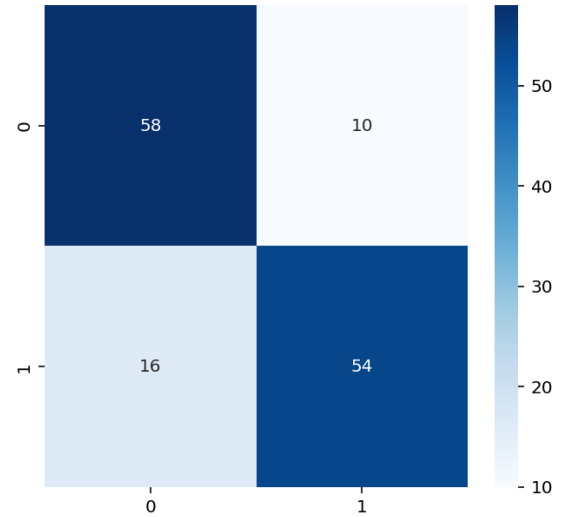
F1 without pretrain: 0.7727 | F1 with pretrain: 0.8060



Confusion Matrix - No Pretrain (Credit Approval)



Confusion Matrix - With Pretrain (Credit Approval)



Breast Cancer Wisconsin Dataset

Without Pretraining:

F1 Score: 0.9534883720930233

Confusion Matrix:

```
[[69 2]  
 [ 2 41]]
```

With Pretraining:

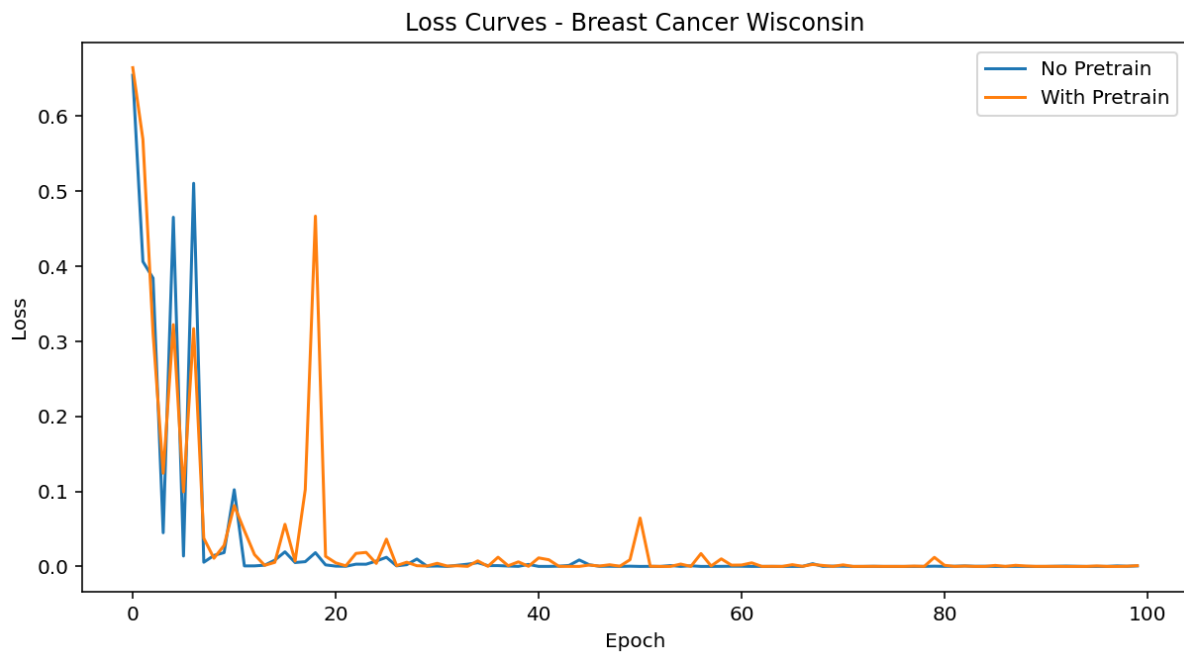
F1 Score: 0.9767441860465116

Confusion Matrix:

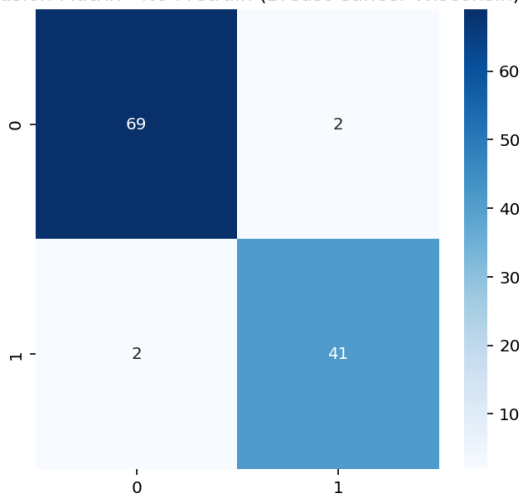
```
[[70 1]  
 [ 1 42]]
```

Comparison:

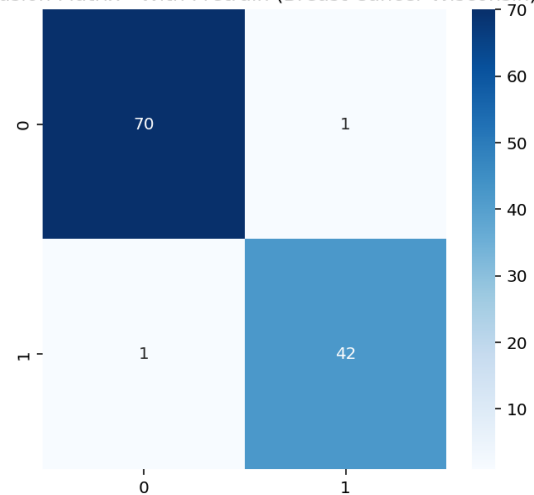
F1 without pretrain: 0.9535 | F1 with pretrain: 0.9767



Confusion Matrix - No Pretrain (Breast Cancer Wisconsin)



Confusion Matrix - With Pretrain (Breast Cancer Wisconsin)



Final Comparison Across Datasets

Credit Approval - F1 without pretrain: 0.7727 | F1 with pretrain: 0.8060

Breast Cancer - F1 without pretrain: 0.9535 | F1 with pretrain: 0.9767

В обоих случаях предобучение с использованием автоэнкодеров привело к улучшению F1-метрики, что подтверждает полезность этого подхода для инициализации нейронной сети. Предобучение помогает модели лучше адаптироваться к структуре данных, особенно в начальных слоях, что приводит к более качественным представлениям признаков.

Вывод: научилась осуществлять предобучение нейронных сетей с помощью автоэнкодерного подхода.