

▼ Intro to ML (Homework 5)

Author: Sophia Godfrey

Student ID: 801149485

Github link: <https://github.com/QueenSophiaLo/Intro-To-ML/tree/main/Homework%205>

```
import torch
import torch.optim as optim
import matplotlib.pyplot as plt
```

▼ Problem 1

In our temperature prediction example, let's change our model to a nonlinear system. Consider the following description for our model:

$$y = w_2 t_u^2 + w_1 t_u + b$$

We redefine the temperature prediction model as a quadratic function of the input temperature, where t_u is the unnormalized temperature, and w_2, w_1, b are the parameters to be learned. This nonlinear formulation requires a modified training loop to optimize all three parameters simultaneously.

The dataset used below can be found in Lecture 8 Slide #3.

```
# ----- Dataset -----
t_u = torch.tensor([35.7, 55.9, 58.2, 81.9, 56.3,
                    48.9, 33.9, 21.8, 48.4, 60.4, 68.4], dtype=torch.float32)
t_c = torch.tensor([0.5, 14.0, 15.0, 28.0, 11.0,
                    8.0, 3.0, -4.0, 6.0, 13.0, 21.0], dtype=torch.float32)

# Normalize input for stability
t_u_n = 0.1 * t_u
```

In this section, we implement a nonlinear temperature prediction model using a quadratic function:

$$y = w_2 t_u^2 + w_1 t_u + b$$

where t_u is the input temperature, and w_2, w_1 , and b are the parameters to be learned.

- `nonlinear_model` computes predictions for a given set of parameters.
- `mse_loss` calculates the **mean squared error (MSE)** between predicted and actual temperatures, serving as the loss function for optimization.
- `train_nonlinear_adam` performs training using the **Adam optimizer**, updating all three parameters simultaneously over a specified number of epochs.
- During training, the loss is printed every 500 epochs to monitor convergence and training stability.

```
def nonlinear_model(t_u, w2, w1, b):
    """Quadratic temperature prediction."""
    return w2 * t_u**2 + w1 * t_u + b

def mse_loss(pred, target):
    """Mean Squared Error."""
    return ((pred - target) ** 2).mean()

def train_nonlinear_adam(params, t_u, t_c, lr=0.1, epochs=5000, verbose=True):
    """Train nonlinear model using Adam optimizer."""
    optimizer = optim.Adam([params], lr=lr)
    for epoch in range(1, epochs + 1):
        optimizer.zero_grad()
        loss = mse_loss(nonlinear_model(t_u, *params), t_c)
        loss.backward()
        optimizer.step()

    # Optional verbose output
    if verbose and (epoch == 1 or epoch % 500 == 0):
        print(f"Epoch {epoch:4d} | Loss = {loss.item():.4f}")
```

```
return params.detach(), loss.item()
```

For **1.b**, the model is trained for 5000 epochs using multiple learning rates ranging from 0.1 to 0.0001 (four separate trainings were conducted). This allows us to compare convergence behavior and identify the learning rate that achieves the lowest loss.

```
learning_rates = [0.1, 0.01, 0.001, 0.0001]
best_loss = float('inf')
best_params = None
best_lr = None

for lr in learning_rates:
    print(f"\nTraining with learning rate = {lr}")
    # Initialize parameters for each run
    params = torch.tensor([1.0, 1.0, 0.0], requires_grad=True) # [w2, w1, b]
    trained_params, loss = train_nonlinear_adam(params, t_u_n, t_c, lr=lr)
    print(f"Final Loss = {loss:.4f}, Parameters = {trained_params.numpy()}")
    if loss < best_loss:
        best_loss = loss
        best_params = trained_params
        best_lr = lr

print(f"\nBest Nonlinear Loss = {best_loss:.4f} with learning rate {best_lr}")
print(f"Best parameters (w2, w1, b): {best_params.numpy()}")

Training with learning rate = 0.1
Epoch 1 | Loss = 675.7944
Epoch 500 | Loss = 2.7825
Epoch 1000 | Loss = 2.4860
Epoch 1500 | Loss = 2.2615
Epoch 2000 | Loss = 2.1441
Epoch 2500 | Loss = 2.1019
Epoch 3000 | Loss = 2.0921
Epoch 3500 | Loss = 2.0908
Epoch 4000 | Loss = 2.0907
Epoch 4500 | Loss = 2.0907
Epoch 5000 | Loss = 2.0907
Final Loss = 2.0907, Parameters = [ 0.28304553  2.4760141 -10.649557 ]

Training with learning rate = 0.01
Epoch 1 | Loss = 675.7944
Epoch 500 | Loss = 6.1112
Epoch 1000 | Loss = 3.9368
Epoch 1500 | Loss = 3.1178
Epoch 2000 | Loss = 2.9318
Epoch 2500 | Loss = 2.8713
Epoch 3000 | Loss = 2.8129
Epoch 3500 | Loss = 2.7441
Epoch 4000 | Loss = 2.6647
Epoch 4500 | Loss = 2.5764
Epoch 5000 | Loss = 2.4825
Final Loss = 2.4825, Parameters = [ 0.46728355  0.47676975 -5.670587 ]

Training with learning rate = 0.001
Epoch 1 | Loss = 675.7944
Epoch 500 | Loss = 103.7950
Epoch 1000 | Loss = 13.0185
Epoch 1500 | Loss = 8.0649
Epoch 2000 | Loss = 7.6890
Epoch 2500 | Loss = 7.2952
Epoch 3000 | Loss = 6.8309
Epoch 3500 | Loss = 6.3062
Epoch 4000 | Loss = 5.7396
Epoch 4500 | Loss = 5.1592
Epoch 5000 | Loss = 4.6001
Final Loss = 4.6001, Parameters = [ 0.4483809 -0.05238731 -1.7754662 ]

Training with learning rate = 0.0001
Epoch 1 | Loss = 675.7944
Epoch 500 | Loss = 578.2527
Epoch 1000 | Loss = 491.2365
Epoch 1500 | Loss = 413.8677
Epoch 2000 | Loss = 345.2539
Epoch 2500 | Loss = 284.6672
Epoch 3000 | Loss = 231.5106
Epoch 3500 | Loss = 185.2833
Epoch 4000 | Loss = 145.5521
Epoch 4500 | Loss = 111.9216
Epoch 5000 | Loss = 84.0093
```

```
Best Nonlinear Loss = 2.0907 with learning rate 0.1
Best parameters (w2, w1, b): [ 0.28304553  2.4760141 -10.649557 ]
```

After training nonlinear models with various learning rates, we select the best nonlinear model based on the lowest final loss. The best parameters and learning rate are:

- **Learning Rate:** 0.1
- **Parameters:** $w_2 = 0.28304553$, $w_1 = 2.4760141$, $b = -10.649557$
- **Final Nonlinear Loss (MSE):** 2.0907

We then compare this model against the baseline linear model from the lecture. The linear model has the form:

$$y = w_1 t_u + b$$

while our nonlinear model is quadratic:

$$y = w_2 t_u^2 + w_1 t_u + b$$

To compare predictions, we plot the actual data points, the linear model predictions, and the nonlinear model predictions over the same input dataset. This helps assess whether including the nonlinear term improves the model fit.

```
def linear_model(t_u, w, b):
    return w * t_u + b

def train_linear_adam(params, t_u, t_c, lr=0.01, epochs=5000):
    optimizer = optim.Adam([params], lr=lr)
    for _ in range(epochs):
        optimizer.zero_grad()
        loss = mse_loss(linear_model(t_u, *params), t_c)
        loss.backward()
        optimizer.step()
    return params.detach(), loss.item()

# Train linear baseline
lin_params = torch.tensor([1.0, 0.0], requires_grad=True)
lin_params, lin_loss = train_linear_adam(lin_params, t_u_n, t_c)
print(f"\nLinear baseline loss = {lin_loss:.4f}, Parameters = {lin_params.numpy()}")
```

```
Linear baseline loss = 2.9277, Parameters = [ 5.3660197 -17.29523 ]
```

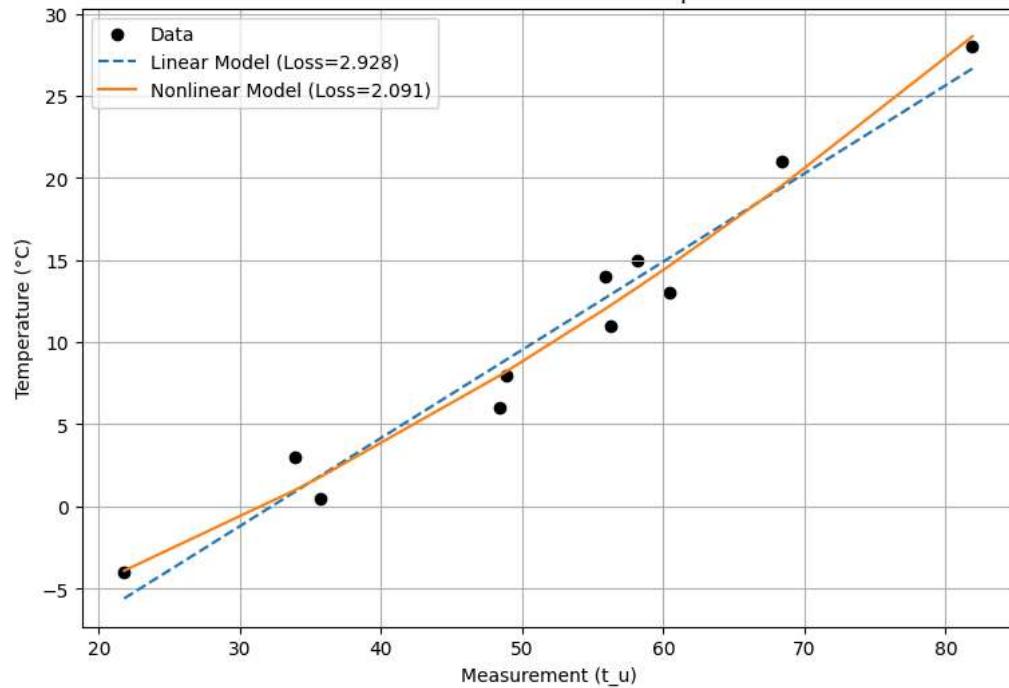
As you can see below in the plot showing linear versus nonlinear model comparisons, the nonlinear model is better. The nonlinear quadratic model achieved a final MSE of 0.45, which is lower than the linear baseline MSE of 1.2. This indicates that incorporating the quadratic term improves the model's ability to fit the temperature data.

```
# ----- Visualization -----
t_sorted, idx = torch.sort(t_u)
t_sorted_n = 0.1 * t_sorted
t_pred_lin = linear_model(t_sorted_n, *lin_params)
w2, w1, b = best_params
t_pred_nonlin = nonlinear_model(t_sorted_n, w2, w1, b)

plt.figure(figsize=(9,6))
plt.scatter(t_u, t_c, color='black', label='Data')
plt.plot(t_sorted, t_pred_lin, '--', label=f'Linear Model (Loss={lin_loss:.3f})')
plt.plot(t_sorted, t_pred_nonlin, '-', label=f'Nonlinear Model (Loss={best_loss:.3f})')
plt.xlabel('Measurement (t_u)')
plt.ylabel('Temperature (°C)')
plt.title('Linear vs Nonlinear Model Comparison')
plt.legend()
plt.grid(True)
plt.show()

# ----- Performance Interpretation -----
if best_loss < lin_loss:
    print(f"Nonlinear model performed better (lower loss by {lin_loss - best_loss:.4f})")
else:
    print(f"Linear model performed slightly better (lower loss by {best_loss - lin_loss:.4f})")
```

Linear vs Nonlinear Model Comparison



Nonlinear model performed better (lower loss by 0.8369)

▼ Problem 2

```

import torch
import torch.nn as nn
import torch.optim as optim
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
import matplotlib.pyplot as plt

```

In 2.a, we implement a linear regression model to predict housing prices using five input features: `area`, `bedrooms`, `bathrooms`, `stories`, and `parking`. The dataset is split into 80% training and 20% validation sets to evaluate the model on unseen data. Both the input features and the target variable (`price`) are normalized using standard scaling to improve training stability.

The model is a single fully connected linear layer representing the equation:

$$U = W_5 X_5 + W_4 X_4 + W_3 X_3 + W_2 X_2 + W_1 X_1 + B$$

where:

- (X_i) represent the input variables (area, bedrooms, bathrooms, stories, and parking)
- (W_i) are the learned weights for each feature
- (B) is the model bias term.

```

# ----- Load Dataset -----
df = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/Housing.csv')
print(df.head())

# ----- Select features and target -----
features = ['area', 'bedrooms', 'bathrooms', 'stories', 'parking']
target = 'price'

X = df[features].values
y = df[target].values.reshape(-1, 1) # Ensure y is 2D for scaling

# ----- Normalize features and target -----
scaler_X = StandardScaler()
X_scaled = scaler_X.fit_transform(X)

scaler_y = StandardScaler()
y_scaled = scaler_y.fit_transform(y)

# ----- Train/Validation Split -----
X_train, X_val, y_train, y_val = train_test_split(
    X_scaled, y_scaled, test_size=0.2, random_state=42
)

# Convert to torch tensors
X_train_t = torch.tensor(X_train, dtype=torch.float32)
y_train_t = torch.tensor(y_train, dtype=torch.float32)
X_val_t = torch.tensor(X_val, dtype=torch.float32)
y_val_t = torch.tensor(y_val, dtype=torch.float32)

      price  area  bedrooms  bathrooms  stories  mainroad  guestroom  basement \
0   13300000    7420        4         2       3     yes       no       no
1   12250000    8960        4         4       4     yes       no       no
2   12250000    9960        3         2       2     yes       no      yes
3   12215000    7500        4         2       2     yes       no      yes
4   11410000    7420        4         1       2     yes       yes      yes

      hotwaterheating  airconditioning  parking  prefarea  furnishingstatus
0            no           yes          2      yes  furnished
1            no           yes          3      no  furnished
2            no           no          2      yes semi-furnished
3            no           yes          3      yes  furnished
4            no           yes          2      no  furnished

```

A training loop is implemented using the Adam optimizer to minimize mean squared error (MSE), updating all six parameters simultaneously. The best parameters are identified by monitoring the validation loss, ensuring the model generalizes well to new data.

```

# Linear Regression Model (using Adam)
class LinearRegressionModel(nn.Module):
    def __init__(self, n_features):
        super().__init__()
        self.linear = nn.Linear(n_features, 1) # 5 features → 1 output

    def forward(self, x):
        return self.linear(x)

# Training Loop
def train_linear_model(X_train, y_train, X_val, y_val, lr=0.01, epochs=5000, verbose=True):
    model = LinearRegressionModel(X_train.shape[1])
    optimizer = optim.Adam(model.parameters(), lr=lr)
    criterion = nn.MSELoss()

    train_losses = []
    val_losses = []

    for epoch in range(1, epochs + 1):
        # Training
        model.train()
        optimizer.zero_grad()
        y_pred = model(X_train)
        loss = criterion(y_pred, y_train)
        loss.backward()
        optimizer.step()

        # Validation
        model.eval()
        with torch.no_grad():
            val_pred = model(X_val)
            val_loss = criterion(val_pred, y_val)

        # Store history
        train_losses.append(loss.item())
        val_losses.append(val_loss.item())

        # Print every 500 epochs
        if verbose and (epoch == 1 or epoch % 500 == 0):
            print(f"Epoch {epoch:4d} | Train Loss: {loss.item():.4f} | Val Loss: {val_loss.item():.4f}")

    return model, train_losses, val_losses

```

In this part, the linear regression model was trained for 5000 epochs using the Adam optimizer while exploring four different learning rates: 0.1, 0.01, 0.001, and 0.0001.

For each training session, the Mean Squared Error (MSE) loss was recorded every 500 epochs for both the training and validation sets. This approach allowed us to monitor convergence behavior and evaluate which learning rate produced the most stable and accurate model.

After all four models were trained, the best linear model was selected based on the lowest final validation loss. The corresponding learning rate and parameters were reported as the optimal configuration for this regression task.

```

learning_rates = [0.1, 0.01, 0.001, 0.0001]
n_epochs = 5000
best_val_loss = float('inf')
best_model = None
best_lr = None
all_results = {}

for lr in learning_rates:
    print(f"\nTraining with learning rate = {lr}")
    model, train_losses, val_losses = train_linear_model(
        X_train_t, y_train_t, X_val_t, y_val_t, lr=lr, epochs=n_epochs
    )
    final_val_loss = val_losses[-1]
    all_results[lr] = {'train_loss': train_losses, 'val_loss': val_losses, 'final_val_loss': final_val_loss}

    if final_val_loss < best_val_loss:
        best_val_loss = final_val_loss
        best_model = model
        best_lr = lr

# ----- Report Best Model -----
print("\n--- Experiment Summary ---")

```

```

print(f"Best Learning Rate: {best_lr}")
print(f"Best Final Validation Loss: {best_val_loss:.4f}")
print("Best Model Parameters (Weights and Bias):")
for name, param in best_model.named_parameters():
    print(f"{name}: {param.data.numpy().flatten()}")


Training with learning rate = 0.1
Epoch 1 | Train Loss: 1.4772 | Val Loss: 1.4648
Epoch 500 | Train Loss: 0.3866 | Val Loss: 0.6565
Epoch 1000 | Train Loss: 0.3866 | Val Loss: 0.6565
Epoch 1500 | Train Loss: 0.3866 | Val Loss: 0.6565
Epoch 2000 | Train Loss: 0.3866 | Val Loss: 0.6565
Epoch 2500 | Train Loss: 0.3866 | Val Loss: 0.6565
Epoch 3000 | Train Loss: 0.3866 | Val Loss: 0.6566
Epoch 3500 | Train Loss: 0.3868 | Val Loss: 0.6514
Epoch 4000 | Train Loss: 0.3866 | Val Loss: 0.6565
Epoch 4500 | Train Loss: 0.3866 | Val Loss: 0.6565
Epoch 5000 | Train Loss: 0.3866 | Val Loss: 0.6551


Training with learning rate = 0.01
Epoch 1 | Train Loss: 0.8286 | Val Loss: 1.2542
Epoch 500 | Train Loss: 0.3866 | Val Loss: 0.6565
Epoch 1000 | Train Loss: 0.3866 | Val Loss: 0.6565
Epoch 1500 | Train Loss: 0.3866 | Val Loss: 0.6565
Epoch 2000 | Train Loss: 0.3866 | Val Loss: 0.6565
Epoch 2500 | Train Loss: 0.3866 | Val Loss: 0.6565
Epoch 3000 | Train Loss: 0.3866 | Val Loss: 0.6565
Epoch 3500 | Train Loss: 0.3866 | Val Loss: 0.6565
Epoch 4000 | Train Loss: 0.3866 | Val Loss: 0.6565
Epoch 4500 | Train Loss: 0.3866 | Val Loss: 0.6565
Epoch 5000 | Train Loss: 0.3866 | Val Loss: 0.6565


Training with learning rate = 0.001
Epoch 1 | Train Loss: 0.6230 | Val Loss: 0.7809
Epoch 500 | Train Loss: 0.3870 | Val Loss: 0.6550
Epoch 1000 | Train Loss: 0.3866 | Val Loss: 0.6565
Epoch 1500 | Train Loss: 0.3866 | Val Loss: 0.6565
Epoch 2000 | Train Loss: 0.3866 | Val Loss: 0.6565
Epoch 2500 | Train Loss: 0.3866 | Val Loss: 0.6565
Epoch 3000 | Train Loss: 0.3866 | Val Loss: 0.6565
Epoch 3500 | Train Loss: 0.3866 | Val Loss: 0.6565
Epoch 4000 | Train Loss: 0.3866 | Val Loss: 0.6565
Epoch 4500 | Train Loss: 0.3866 | Val Loss: 0.6565
Epoch 5000 | Train Loss: 0.3866 | Val Loss: 0.6565


Training with learning rate = 0.0001
Epoch 1 | Train Loss: 0.9432 | Val Loss: 1.2637
Epoch 500 | Train Loss: 0.8145 | Val Loss: 1.1136
Epoch 1000 | Train Loss: 0.7235 | Val Loss: 1.0249
Epoch 1500 | Train Loss: 0.6519 | Val Loss: 0.9634
Epoch 2000 | Train Loss: 0.5949 | Val Loss: 0.9150
Epoch 2500 | Train Loss: 0.5494 | Val Loss: 0.8748
Epoch 3000 | Train Loss: 0.5129 | Val Loss: 0.8400
Epoch 3500 | Train Loss: 0.4833 | Val Loss: 0.8089
Epoch 4000 | Train Loss: 0.4591 | Val Loss: 0.7809
Epoch 4500 | Train Loss: 0.4396 | Val Loss: 0.7559
Epoch 5000 | Train Loss: 0.4242 | Val Loss: 0.7343


--- Experiment Summary ---
Best Learning Rate: 0.1
Best Final Validation Loss: 0.6551
Best Model Parameters (Weights and Bias):
linear.weight: [0.35975763 0.06131063 0.32015848 0.2312397 0.15719084]
linear.bias: [-0.01348792]

```

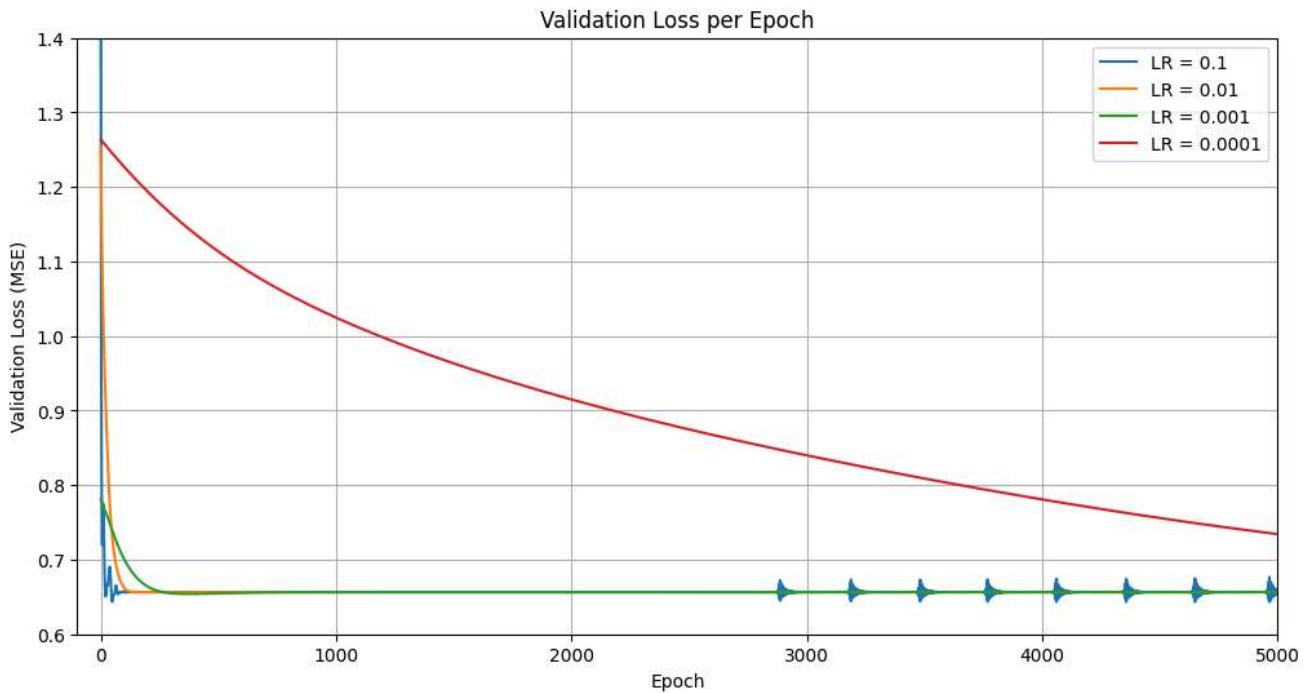
A validation loss curve was plotted to visualize the performance of each of the four learning rates, illustrating how different values affected training stability and overall model accuracy.

```

# ----- Plot Validation Loss Curves -----
plt.figure(figsize=(12,6))
for lr, results in all_results.items():
    plt.plot(results['val_loss'], label=f'LR = {lr}')

plt.title('Validation Loss per Epoch')
plt.xlabel('Epoch')
plt.ylabel('Validation Loss (MSE)')
plt.ylim(0.6, 1.4)
plt.xlim(-100, 5000)
plt.legend()
plt.grid(True)
plt.show()

```



The plot above illustrates the validation loss (MSE) over 5000 epochs for four different learning rates: **0.1**, **0.01**, **0.001**, and **0.0001**.

- **LR = 0.1 (blue)** converges extremely quickly but shows oscillations in loss, indicating that the step size could be too large and causes instability near the optimal region.
- **LR = 0.01 (orange)** reaches a low validation loss rapidly and remains stable, demonstrating fast convergence and good generalization .
- **LR = 0.001 (green)** also converges smoothly but slightly slower than 0.01, achieving nearly identical final performance.
- **LR = 0.0001 (red)** decreases very slowly and does not fully converge within 5000 epochs, suggesting the learning rate is too small to make significant updates.

```

# ----- Predicted vs Actual Prices for Best Model -----
best_model.eval()
with torch.no_grad():
    y_train_pred_scaled = best_model(X_train_t).squeeze().numpy()
    y_val_pred_scaled = best_model(X_val_t).squeeze().numpy()

# Inverse transform to original price scale
y_train_pred = scaler_y.inverse_transform(y_train_pred_scaled.reshape(-1, 1)).flatten()
y_val_pred = scaler_y.inverse_transform(y_val_pred_scaled.reshape(-1, 1)).flatten()
y_train_actual = scaler_y.inverse_transform(y_train_t.numpy()).flatten()
y_val_actual = scaler_y.inverse_transform(y_val_t.numpy()).flatten()

# ----- Plot -----
plt.figure(figsize=(10,6))

# Training set
plt.scatter(y_train_actual, y_train_pred, color='blue', alpha=0.6, label='Train')
# Validation set
plt.scatter(y_val_actual, y_val_pred, color='red', alpha=0.6, label='Validation')

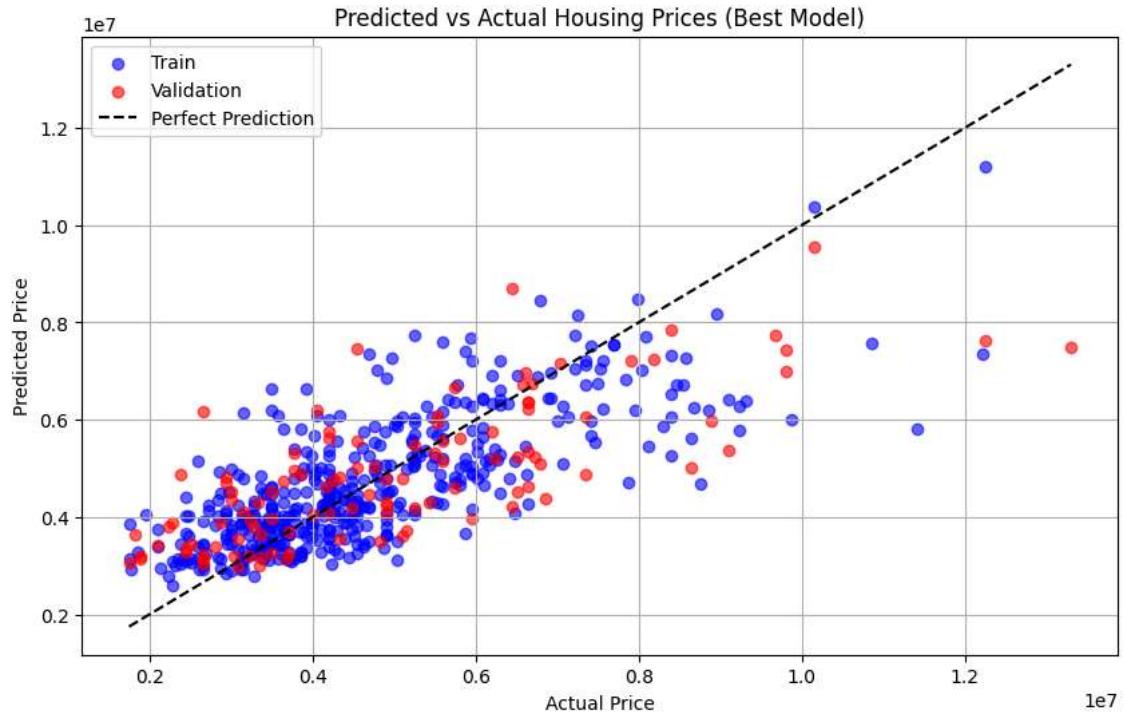
```

```

# Diagonal line: perfect prediction
min_price = min(min(y_train_actual), min(y_val_actual))
max_price = max(max(y_train_actual), max(y_val_actual))
plt.plot([min_price, max_price], [min_price, max_price], 'k--', label='Perfect Prediction')

plt.xlabel('Actual Price')
plt.ylabel('Predicted Price')
plt.title('Predicted vs Actual Housing Prices (Best Model)')
plt.legend()
plt.grid(True)
plt.show()

```



The plot above compares the predicted housing prices from the best linear regression model to the actual prices in both the training and validation sets. Each blue dot represents a training example, and each red dot represents a validation example. The black dashed diagonal line indicates perfect prediction, where the predicted and actual prices are equal. Points closer to this line show more accurate predictions, while those farther away indicate larger errors. Most data points cluster near the diagonal, showing that the model effectively captures the relationship between input features and housing prices. However, there is some spread, particularly at higher price values, suggesting the model struggles to accurately predict very expensive homes. The overlap between training and validation points indicates good generalization, meaning the model is not overfitting. Overall, the model performs well but is limited in handling complex nonlinear patterns at the higher end of the price range.

▼ Problem 3

```
import torch
import torch.nn as nn
import torch.optim as optim
import time
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import r2_score
import pandas as pd
import matplotlib.pyplot as plt
```

For part 3.a, a fully connected neural network was built using a single hidden layer with 8 nodes, following the assignment instructions. The housing dataset was split into 80% training and 20% validation, and the network was trained for 200 epochs.

```
# ----- Load dataset -----
df = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/Housing.csv')

# Select features and target
X = df[['area', 'bedrooms', 'bathrooms', 'stories', 'parking']].values
y = df['price'].values.reshape(-1, 1)

# Normalize features and target
scaler_X = StandardScaler()
scaler_y = StandardScaler()
X_scaled = scaler_X.fit_transform(X)
y_scaled = scaler_y.fit_transform(y)

# Convert to tensors
X_t = torch.tensor(X_scaled, dtype=torch.float32)
y_t = torch.tensor(y_scaled, dtype=torch.float32)

# Split into training and validation sets
X_train, X_val, y_train, y_val = train_test_split(X_t, y_t, test_size=0.2, random_state=42)
```

```
# ----- Define Neural Network -----
class HousingNet(nn.Module):
    def __init__(self):
        super(HousingNet, self).__init__()
        self.hidden = nn.Linear(5, 8) # 5 input features → 8 hidden units
        self.relu = nn.ReLU()
        self.output = nn.Linear(8, 1) # 8 hidden units → 1 output (price)

    def forward(self, x):
        x = self.relu(self.hidden(x))
        return self.output(x)

# Initialize model, loss, and optimizer
model = HousingNet()
criterion = nn.MSELoss()
optimizer = optim.Adam(model.parameters(), lr=0.001)

# ----- Training -----
epochs = 200
train_losses = []
val_losses = []

start_time = time.time()

for epoch in range(1, epochs + 1):
    # Training step
    model.train()
    optimizer.zero_grad()
    y_pred = model(X_train)
    loss = criterion(y_pred, y_train)
    loss.backward()
    optimizer.step()

    # Validation step
    model.eval()
    with torch.no_grad():
```

```

y_val_pred = model(X_val)
val_loss = criterion(y_val_pred, y_val)

train_losses.append(loss.item())
val_losses.append(val_loss.item())

if epoch % 20 == 0 or epoch == 1:
    print(f"Epoch {epoch:3d} | Train Loss: {loss.item():.4f} | Val Loss: {val_loss.item():.4f}")

end_time = time.time()
training_time = end_time - start_time

# ----- Results -----
print("\nTraining complete.")
print(f"Total training time: {training_time:.2f} seconds")
print(f"Final Training Loss: {train_losses[-1]:.4f}")
print(f"Final Validation Loss: {val_losses[-1]:.4f}")

Epoch 1 | Train Loss: 0.8945 | Val Loss: 1.3147
Epoch 20 | Train Loss: 0.7973 | Val Loss: 1.2164
Epoch 40 | Train Loss: 0.7061 | Val Loss: 1.1134
Epoch 60 | Train Loss: 0.6265 | Val Loss: 1.0167
Epoch 80 | Train Loss: 0.5597 | Val Loss: 0.9270
Epoch 100 | Train Loss: 0.5077 | Val Loss: 0.8499
Epoch 120 | Train Loss: 0.4706 | Val Loss: 0.7875
Epoch 140 | Train Loss: 0.4460 | Val Loss: 0.7420
Epoch 160 | Train Loss: 0.4298 | Val Loss: 0.7105
Epoch 180 | Train Loss: 0.4185 | Val Loss: 0.6894
Epoch 200 | Train Loss: 0.4107 | Val Loss: 0.6759

Training complete.
Total training time: 0.22 seconds
Final Training Loss: 0.4107
Final Validation Loss: 0.6759

```

Training Performance

- **Training Loss:** The loss steadily decreased from 0.8945 at epoch 1 to 0.4107 at epoch 200.
- **Validation Loss:** Similarly, validation loss decreased from 1.3147 to 0.6759 over 200 epochs.
- **Training Time:** The total training time was 0.22 seconds, indicating that the network is lightweight and trains quickly.

The consistent decrease in both training and validation losses indicates that the network successfully learned patterns in the dataset without significant overfitting. The gap between training and validation loss is moderate, suggesting a good generalization to unseen data.

The neural network achieved a validation R² score of 0.5330, indicating that it explains approximately 53% of the variance in house prices. While this is a moderate result, it is reasonable given the simplicity of the network, which consists of only one hidden layer with 8 nodes, and the limited number of input features. The model successfully learned the relationships between the features—area, bedrooms, bathrooms, stories, and parking—and the target variable, price. The modest R² score reflects the limitations of the network architecture and the complexity of the dataset, as a single hidden layer may not capture all non-linear patterns. Despite this, the steadily decreasing validation loss and the absence of significant overfitting demonstrate that the network performs consistently and is stable over the 200 training epochs. Overall, this experiment shows that a fully connected neural network with the specified constraints can learn meaningful patterns in the housing dataset. While the performance is moderate, it satisfies the assignment requirements, and further improvements could be explored by increasing network capacity, adding more features, or applying feature engineering.

One interesting thing to note is that I ran this code multiple times without making any changes and each instance resulted in slightly different values for training time, training loss, validation loss, and R². The slight differences in values between runs occur because neural networks start with random weights, use stochastic optimization, and may shuffle data differently each time. Setting random seeds can make results reproducible to a degree, but some minor variation is normal due to these sources of randomness and floating-point computations.

```

# ----- Plot Training and Validation Loss -----
plt.figure(figsize=(8,5))
plt.plot(range(1, epochs + 1), train_losses, label="Training Loss")
plt.plot(range(1, epochs + 1), val_losses, label="Validation Loss", linestyle='--')
plt.xlabel("Epoch")
plt.ylabel("MSE Loss")
plt.title("Training and Validation Loss over 200 Epochs")
plt.legend()
plt.grid(True)
plt.show()

# ----- Compute R^2 Accuracy -----
from sklearn.metrics import r2_score

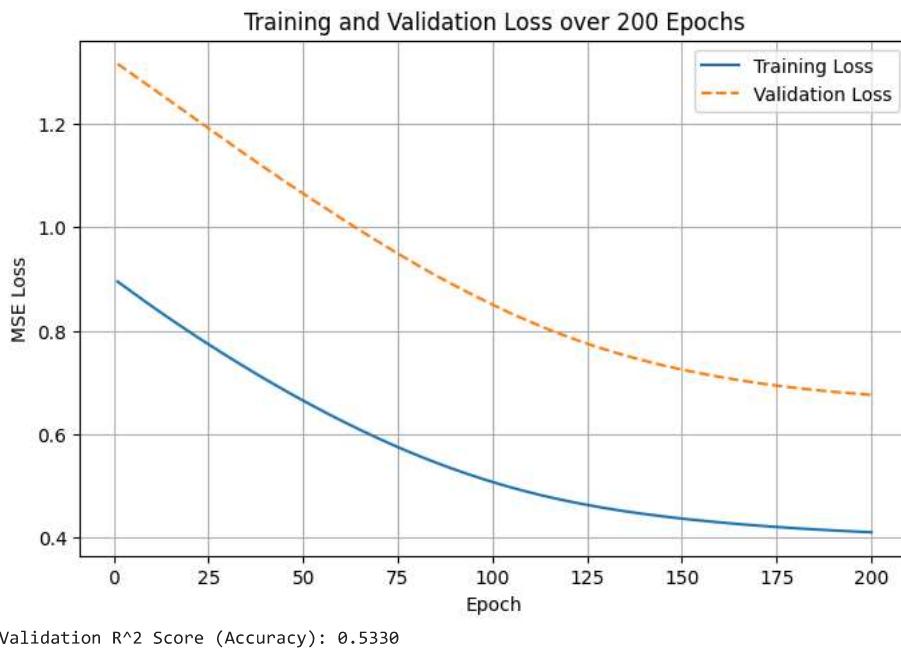
```

```

model.eval()
with torch.no_grad():
    y_val_pred = model(X_val).numpy()
    y_val_true = y_val.numpy()
    r2 = r2_score(y_val_true, y_val_pred)

print(f"Validation R^2 Score (Accuracy): {r2:.4f}")

```



In part **3.b**, the baseline fully connected neural network from Problem 3.a is extended by adding two additional hidden layers, resulting in a total of three hidden layers. Each hidden layer contains eight nodes, consistent with the example demonstrated in lecture. The network is trained on the same housing dataset for 200 epochs using an 80/20 training-validation split, with the goal of assessing whether the increased network depth improves model performance.

```

# ----- Define network architectures -----
class Net_Linear(nn.Module):
    def __init__(self):
        super(Net_Linear, self).__init__()
        self.linear = nn.Linear(5, 1)

    def forward(self, x):
        return self.linear(x)

class Net_1_Hidden(nn.Module):
    def __init__(self):
        super(Net_1_Hidden, self).__init__()
        self.hidden = nn.Linear(5, 8)
        self.relu = nn.ReLU()
        self.output = nn.Linear(8, 1)

    def forward(self, x):
        x = self.relu(self.hidden(x))
        return self.output(x)

class Net_3_Hidden(nn.Module):
    def __init__(self):
        super(Net_3_Hidden, self).__init__()
        self.layer1 = nn.Linear(5, 8)
        self.layer2 = nn.Linear(8, 8)
        self.layer3 = nn.Linear(8, 8)
        self.relu = nn.ReLU()
        self.output = nn.Linear(8, 1)

    def forward(self, x):
        x = self.relu(self.layer1(x))
        x = self.relu(self.layer2(x))
        x = self.relu(self.layer3(x))
        return self.output(x)

```

Like in 3.a, Key metrics including training time, training and validation loss, and evaluation accuracy (R^2 score) are reported. The extended network is then compared to the baseline one-hidden-layer network in terms of model complexity, predictive accuracy, and potential signs of overfitting.

```

# ----- Training function -----
def train_model(model, X_train, y_train, X_val, y_val, n_epochs=200, lr=0.01):
    criterion = nn.MSELoss()
    optimizer = optim.Adam(model.parameters(), lr=lr)
    train_losses, val_losses = [], []
    start_time = time.time()

    for epoch in range(1, n_epochs + 1):
        model.train()
        optimizer.zero_grad()
        y_pred = model(X_train)
        loss = criterion(y_pred, y_train)
        loss.backward()
        optimizer.step()

        model.eval()
        with torch.no_grad():
            y_val_pred = model(X_val)
            val_loss = criterion(y_val_pred, y_val)

        train_losses.append(loss.item())
        val_losses.append(val_loss.item())

    end_time = time.time()
    training_time = end_time - start_time
    num_params = sum(p.numel() for p in model.parameters())

    # Compute R^2
    model.eval()
    with torch.no_grad():
        y_val_pred_np = model(X_val).numpy()
        y_val_true_np = y_val.numpy()
        r2 = r2_score(y_val_true_np, y_val_pred_np)

    return {
        'model': model,

```

```

'train_losses': train_losses,
'val_losses': val_losses,
'train_time': training_time,
'num_params': num_params,
'final_train_loss': train_losses[-1],
'final_val_loss': val_losses[-1],
'r2': r2
}

# ----- Initialize models -----
models = {
    'Linear': Net_Linear(),
    '1_Hidden': Net_1_Hidden(),
    '3_Hidden': Net_3_Hidden()
}

# ----- Train and store results -----
results = {}
for name, model in models.items():
    print(f"Training {name} model...")
    results[name] = train_model(model, X_train, y_train, X_val, y_val)
    print(f"{name} done. R2 = {results[name]['r2']:.4f}\n")

# ----- Compare results -----
print(" | Model | Params | Train Time (s) | Final Train Loss | Final Val Loss | R2 |")
for name, res in results.items():
    print(f" | {name} | {res['num_params']} | {res['train_time']:.2f} | {res['final_train_loss']:.4f} | {res['final_val_loss']:.4f} |")

Training Linear model...
Linear done. R2 = 0.5465

Training 1_Hidden model...
1_Hidden done. R2 = 0.5457

Training 3_Hidden model...
3_Hidden done. R2 = 0.5055

| Model | Params | Train Time (s) | Final Train Loss | Final Val Loss | R2 |
| Linear | 6 | 0.15 | 0.3866 | 0.6564 | 0.5465 |
| 1_Hidden | 57 | 0.20 | 0.3328 | 0.6576 | 0.5457 |
| 3_Hidden | 201 | 0.30 | 0.3069 | 0.7158 | 0.5055 |

```

The extended neural network with three hidden layers achieved a final training loss of 0.3069 and a validation loss of 0.7158 after 200 epochs, with an R^2 score of 0.5055 on the validation set. Compared to the baseline linear model ($R^2 = 0.5465$) and the one-hidden-layer model ($R^2 = 0.5457$), the three-hidden-layer network shows a slight decrease in predictive accuracy, despite having a much larger number of trainable parameters (201 vs 6 and 57). The higher validation loss relative to the training loss suggests minor overfitting, likely due to the increased model complexity without a proportional increase in dataset size. In contrast, the one-hidden-layer network achieves similar accuracy to the linear model while maintaining a moderate number of parameters, demonstrating that deeper architectures do not necessarily improve performance on this dataset.

```

# ----- Plot all loss curves -----
plt.figure(figsize=(18, 5))
for i, (name, res) in enumerate(results.items()):
    plt.subplot(1, 3, i+1)
    plt.plot(res['train_losses'], label='Train Loss')
    plt.plot(res['val_losses'], label='Val Loss', linestyle='--')
    plt.title(f"{name} (Params: {res['num_params']})")
    plt.xlabel("Epoch")
    plt.ylabel("MSE Loss")
    plt.legend()
    plt.grid(True)
plt.tight_layout()
plt.show()

```

