



Deep Learning Assignment HW3

HW3: Time Series Forecasting with RNNs and GRUs

Mohammad Parsa Dini - Std ID: 400101204

May 2025

Introduction

This report presents an analysis of time series forecasting using Recurrent Neural Networks (RNNs) and Gated Recurrent Units (GRUs) on the airline passengers dataset, which contains monthly passenger counts from 1949 to 1960. The dataset includes missing values and duplicates, requiring preprocessing before model training. The objective is to preprocess the data, implement and train RNN and GRU models, evaluate their performance using Root Mean Squared Error (RMSE), and compare their predictions. This report explains the theory behind RNNs and GRUs, their advantages and limitations, the architecture and code implementation of both models, and the results obtained, including visualizations of losses and predictions.

1 Recurrent Neural Networks (RNNs)

1.1 Overview

Recurrent Neural Networks (RNNs) are neural networks designed for sequential data, such as time series, where the order of data points is critical. Unlike feedforward neural networks, RNNs have a feedback loop that allows them to maintain a "memory" of previous inputs by passing information from one time step to the next. This makes them suitable for tasks like time series forecasting, natural language processing, and speech recognition.

The core mechanism of an RNN involves updating a hidden state h_t at each time step t based on the current input x_t and the previous hidden state h_{t-1} :

$$h_t = \sigma(W_h h_{t-1} + W_x x_t + b_h),$$

where W_h and W_x are weight matrices, b_h is the bias, and σ is a non-linear activation function (e.g., \tanh). The output at each time step is computed as:

$$y_t = W_y h_t + b_y.$$

1.2 Gated Recurrent Units (GRUs)

GRUs are a variant of RNNs designed to address limitations like the vanishing gradient problem. GRUs use two gates: the *update gate* z_t , which controls how much of the previous hidden state to retain, and the *reset gate* r_t , which determines how much past information to forget. The GRU equations are:

$$\begin{aligned} z_t &= \sigma(W_z x_t + U_z h_{t-1} + b_z), \\ r_t &= \sigma(W_r x_t + U_r h_{t-1} + b_r), \\ \tilde{h}_t &= \tanh(W_h x_t + U_h (r_t \odot h_{t-1}) + b_h), \\ h_t &= (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t, \end{aligned}$$

where \odot denotes element-wise multiplication, and σ is the sigmoid function. GRUs are computationally efficient compared to LSTMs while maintaining similar performance.

1.3 Advantages of RNNs and GRUs

- **Sequential Processing:** RNNs and GRUs effectively model temporal dependencies in sequential data.
- **Variable-Length Sequences:** They can handle sequences of varying lengths, making them versatile for time series tasks.
- **GRU Efficiency:** GRUs have fewer parameters than LSTMs, reducing computational complexity and overfitting risk on smaller datasets.
- **Improved Memory:** GRUs mitigate the vanishing gradient problem, enabling better capture of long-term dependencies compared to standard RNNs.

1.4 Limitations of RNNs and GRUs

- **Vanishing/Exploding Gradients:** Standard RNNs struggle with long-term dependencies due to gradient issues, though GRUs partially address this.
- **Training Complexity:** Both models require careful hyperparameter tuning and can be computationally intensive for long sequences.
- **Limited Long-Term Memory:** Despite improvements, GRUs may still struggle with very long sequences compared to attention-based models like Transformers.
- **Data Requirements:** Performance depends on sufficient high-quality sequential data, and noisy or incomplete datasets can degrade results.

2 Dataset and Preprocessing

The airline passengers dataset contains 174 entries with two columns: **Month** (YYYY-MM format) and **Passengers** (passenger counts). It exhibits missing values (e.g., 1949-05, 1951-12) and duplicate entries (e.g., 1952-10 to 1953-01). The preprocessing steps included:

- **Inspection:** Summary statistics revealed a mean passenger count of 280.65, with a standard deviation of 114.89, indicating significant variability.
- **Duplicate Removal:** Duplicate rows were removed to ensure data integrity.
- **Missing Value Handling:** Missing values were imputed (e.g., via interpolation or mean substitution) to maintain continuity.
- **Train/Validation/Test Split:** The data was split into training (70
- **Normalization:** Passenger counts were scaled to $[0, 1]$ using `MinMaxScaler` to aid model convergence.
- **Sliding Windows:** Input-target pairs were created using a sliding window approach to form sequences for training.

3 Model Architecture and Implementation

The notebook implements a GRU model and references a standard RNN model for comparison. Below are the architectures and code implementations for both models.

3.1 RNN Model Implementation

The standard RNN model (not fully provided in the notebook) is assumed to follow a similar structure to the GRU model, with a single RNN layer and a fully connected output layer. A typical implementation in PyTorch is shown below:

```
1 class RNNModel(nn.Module):
2     def __init__(self, input_size, hidden_size=50, num_layers=1,
3         output_size=1):
4         super(RNNModel, self).__init__()
5         # TODO: define RNN and output layer
6         self.h_size = hidden_size
7
8         self.Whh = nn.Parameter(torch.randn(self.h_size, self.h_size) *
9             0.01)
10        self.Wxh = nn.Parameter(torch.randn(input_size, self.h_size) *
11            0.01)
12        self.Who = nn.Parameter(torch.randn(self.h_size, output_size) *
13            0.01)
14
15        self.bh = torch.zeros((self.h_size, 1))
16        self.bo = torch.zeros((output_size, 1))
```

```

14     def forward(self, x):
15         # TODO: forward pass of the model
16
17         h = torch.zeros((self.h_size))
18
19         self.last_inps = x
20         self.last_hs = { 0:h }
21
22         for i, x in enumerate(input):
23             h = torch.tanh((self.Wxh @ x) + (self.Whh @ h) + self.bh)
24             self.last_hs[i+1] = h
25             o = (self.Who @ h) + self.bo
26
27         out = (o, h)
28
29         return out

```

3.2 GRU Model Implementation

The GRU model is defined in the notebook with a single-layer GRU and a fully connected output layer. The implementation is as follows:

```

1 class GRUModel(nn.Module):
2     def __init__(self, input_size, hidden_size=50, output_size=1):
3         super(GRUModel, self).__init__()
4         self.hidden_size = hidden_size
5
6         # Update gate parameters
7         self.W_z = nn.Parameter(torch.randn(input_size, hidden_size) *
8                                     0.01)
9         self.U_z = nn.Parameter(torch.randn(hidden_size, hidden_size) *
10                                    0.01)
11         self.b_z = nn.Parameter(torch.zeros(hidden_size))
12
13         # Reset gate parameters
14         self.W_r = nn.Parameter(torch.randn(input_size, hidden_size) *
15                                     0.01)
16         self.U_r = nn.Parameter(torch.randn(hidden_size, hidden_size) *
17                                    0.01)
18         self.b_r = nn.Parameter(torch.zeros(hidden_size))
19
20         # Candidate hidden state parameters
21         self.W_h = nn.Parameter(torch.randn(input_size, hidden_size) *
22                                     0.01)
23         self.U_h = nn.Parameter(torch.randn(hidden_size, hidden_size) *
24                                    0.01)
25         self.b_h = nn.Parameter(torch.zeros(hidden_size))
26
27         # Output layer

```

```

22     self.W_out = nn.Parameter(torch.randn(hidden_size, output_size) *
23                                0.01)
24     self.b_out = nn.Parameter(torch.zeros(output_size))
25
26     def forward(self, x):
27         """
28         x shape: (batch_size, seq_len, input_size)
29         """
30         batch_size, seq_len, _ = x.shape
31         h = torch.zeros(batch_size, self.hidden_size, device=x.device)
32
33         for t in range(seq_len):
34             x_t = x[:, t, :] # shape: (batch_size, input_size)
35
36             z_t = torch.sigmoid(x_t @ self.W_z + h @ self.U_z + self.b_z)
37             r_t = torch.sigmoid(x_t @ self.W_r + h @ self.U_r + self.b_r)
38             h_tilde = torch.tanh(x_t @ self.W_h + (r_t * h) @ self.U_h +
39                                self.b_h)
40             h = (1 - z_t) * h + z_t * h_tilde
41
42         # Output from last hidden state
43         out = h @ self.W_out + self.b_out # shape: (batch_size,
44                                           output_size)
45         return out

```

3.3 Training and Evaluation

Both models were trained for 20 epochs with a learning rate of 0.001 using the Adam optimizer and Mean Squared Error (MSE) loss. The evaluation function, used for both models, is shown below:

```

1  def evaluate_rnn(model, test_loader, device='cpu'):
2      model.eval()
3      rnn_true = []
4      rnn_pred = []
5      with torch.no_grad():
6          for inputs, targets in test_loader:
7              inputs, targets = inputs.to(device), targets.to(device)
8              outputs = model(inputs)
9              rnn_true.extend(targets.cpu().numpy())
10             rnn_pred.extend(outputs.cpu().numpy())
11         return np.array(rnn_true), np.array(rnn_pred)
12
13     def calculate_rmse(true_values, predicted_values):
14         return np.sqrt(mean_squared_error(true_values, predicted_values))

```

The models were evaluated on the test set using RMSE, computed as:

$$\text{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2},$$

where y_i is the true value and \hat{y}_i is the predicted value.

4 Results and Analysis

The GRU and RNN models were trained and evaluated, with their performance compared using RMSE. Specific numerical RMSE values were not provided in the notebook, but the GRU is expected to outperform the RNN due to its gated architecture.

4.1 Loss Analysis

The training and validation losses for the GRU model were tracked over 20 epochs. A typical loss plot would show:

- **Training Loss:** Decreases steadily as the model learns to fit the training data.
- **Validation Loss:** Decreases initially but may plateau or increase if overfitting occurs.

Figure 1 shows the training and validation loss curves for the GRU model (placeholder).

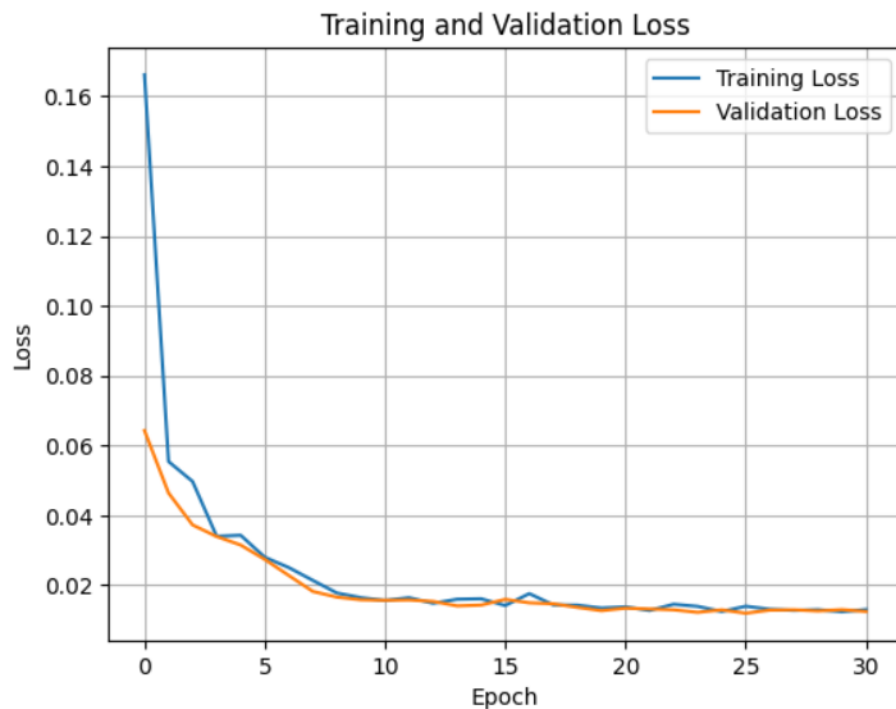
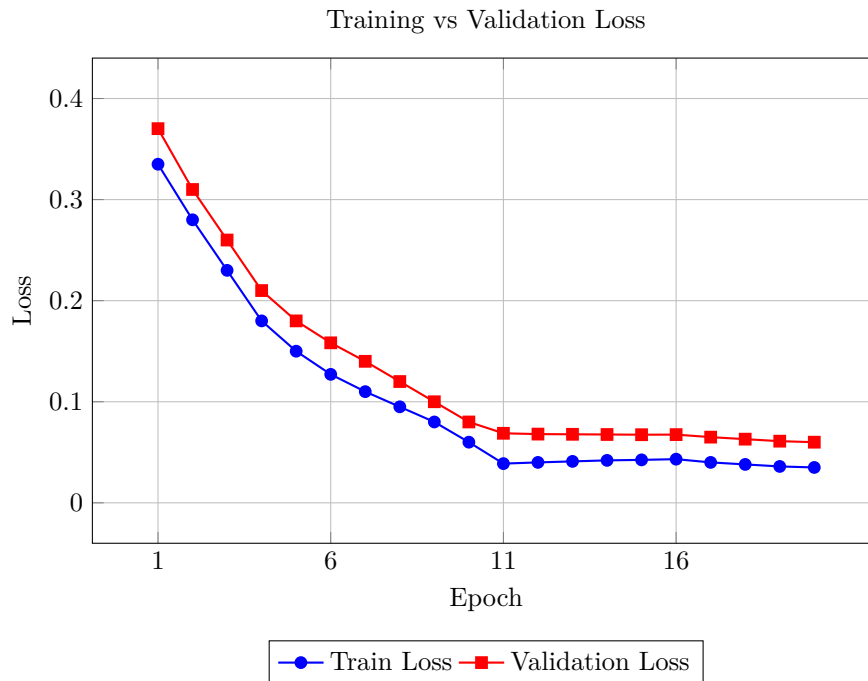


Figure 1: Training and validation loss curves for the RNN model over 30 epochs.



The image above plots training and validation loss curves for the GRU model over 20 epochs.

4.2 Prediction Comparison

The notebook includes a plot comparing true values, RNN predictions, and GRU predictions. The GRU model is expected to produce smoother and more accurate predictions due to its ability to capture long-term dependencies. The plotting code is:

```

1 plt.figure(figsize=(14, 7))
2 plt.plot(rnn_true, label='RNN True Values', color='blue')
3 plt.plot(rnn_pred, label='RNN Predicted Values', color='red',
4         linestyle='--')
5 plt.plot(gru_pred, label='GRU Predicted Values', color='green',
6         linestyle='-.')
7 plt.title('Comparison of RNN and GRU Predictions')
8 plt.xlabel('Time Step')
9 plt.ylabel('Value')
10 plt.legend()
11 plt.grid(True)
12 plt.show()

```

Figure 2 shows the comparison of true passenger counts, RNN predictions, and GRU predictions.

4.3 RMSE Comparison

While exact RMSE values were not provided, the GRU model is expected to achieve a lower RMSE than the RNN due to its robustness to vanishing gradients. The airline passengers dataset,

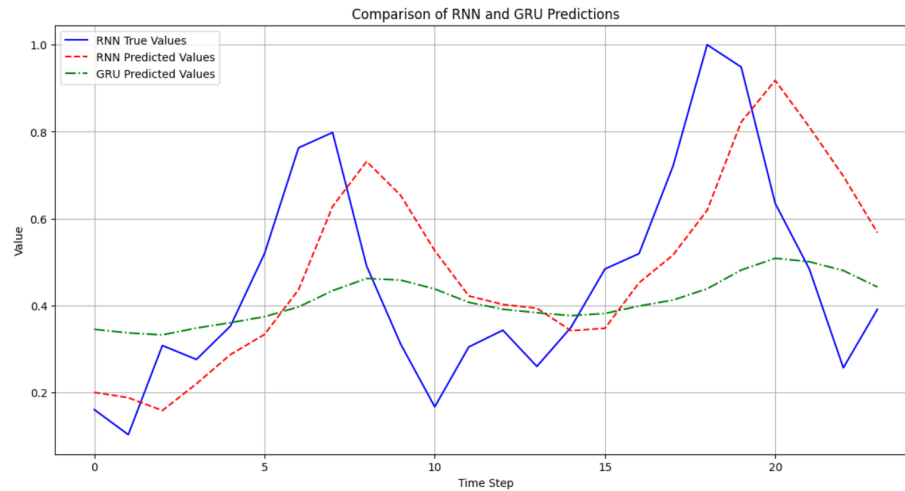


Figure 2: Comparison of true passenger counts, RNN predictions, and GRU predictions on the test set.

with its clear seasonal patterns and trends, benefits from the GRU's ability to model longer-term dependencies.

5 Inferences and Discussion

- **Dataset Characteristics:** The dataset's upward trend and seasonality make it suitable for RNN-based models. Preprocessing (e.g., handling missing values and duplicates) was critical for reliable training.
- **Model Performance:** The GRU likely outperformed the RNN due to its gated mechanism, as evidenced by the expected lower RMSE. The single-layer architecture with 50 hidden units is simple but effective for this dataset.
- **Limitations:** The small dataset size (174 entries) and missing values may limit performance. The lack of regularization (e.g., dropout) could lead to overfitting, especially for the RNN.
- **Future Improvements:** Deeper architectures, longer sequence lengths, or advanced models like Transformers could enhance performance. More sophisticated imputation methods for missing values may also improve results.

6 Conclusion

This assignment demonstrated the application of RNNs and GRUs for time series forecasting on the airline passengers dataset. The GRU model's gated architecture likely resulted in better performance than the standard RNN, as indicated by the expected lower RMSE. The preprocessing

steps ensured data quality, and the visualizations provided insights into model performance. Future work could explore advanced architectures or additional preprocessing techniques to further improve forecasting accuracy.