

**DA5401 (Jul - Nov 2025)**

## **Data Challenge Project Report**

Saggurthi Dinesh

BE21B032

### **Table of Contents**

**1. Introduction**

**2. Data Engineering**

**3. Exploratory Data Analysis (EDA)**

**4. Sampling Strategy**

**5. Model Architecture**

**6. Loss Functions**

**7. Sampling Strategy & Weighting**

**8. Encoder Fine-tuning**

**9. Distribution Bias Fixes**

**10. Training Improvements**

**11. Inference Pipeline**

**12. Results & Performance**

## Introduction:

### Problem Statement:

The challenge requires predicting a continuous fitness score (ranging from 0 to 10) that measures how well a given evaluation metric matches with a prompt-response pair. This is fundamentally a **regression problem** evaluated using **Root Mean Squared Error (RMSE)** as the primary metric.

### Key Challenges:

1. **Asymmetric Input Representation:** The problem combines pre-computed metric embeddings (768-dimensional) with raw text data (prompts and responses).
2. **Multilingual Data:** Text data spans multiple languages including Tamil, Hindi, and English.
3. **Severe Class Imbalance:** Over 91% of training data consists of scores 9-10, with rare occurrences of scores 0-8.
4. **Discretized Target:** While scores are integers 0-10, they represent a continuous scale for RMSE evaluation.

### Solution Approach:

My approach integrates several sophisticated deep learning techniques, including multilingual sentence transformer encoding for text processing, diverse model architectures (Standard, Ordinal, and Heteroscedastic regression), innovative data augmentation strategies (synthetic negatives), and advanced sampling and weighting schemes. Furthermore, I employ encoder fine-tuning through SimCSE-style contrastive learning and comprehensively address dataset bias.

### Evaluation Metric:

The primary evaluation metric is **Root Mean Squared Error (RMSE)** calculated on a hidden test set:

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

RMSE = Where ( $y_i$ ) is the true score and ( $\hat{y}_i$ ) is the predicted score for sample (i).

## 1. Data Engineering

The data engineering process was focused on transforming the raw text and metric data into a feature-rich format suitable for a deep learning mode.

### 1.1 Text Feature Engineering

A vital step involved merging diverse text fields. This necessitated a flexible method, utilizing unified prompt/response concatenation for multilingual encoding, due to varying field names across different records.

This combined text was then encoded using the **“paraphrase-multilingual-mpnet-base-v2 SentenceTransformer”**, converting each sample into a **768-dimensional vector** (text\_embedding).

This multilingual model was essential for handling text in Tamil, Hindi, and English.

### 1.2 Metric Feature Engineering

The metric\_name for each sample was used as a key to look up its corresponding 768-dimensional vector from the pre-computed `metric_name_embeddings.npy` file and this lookup avoids redundant encoding. This approach is computationally efficient and ensures that identical metric names always map to the same embedding vector.

### 1.3 Data Augmentation (Synthetic Negatives)

To help the model learn what a bad metric/text pairing looks like, I have employed a novel data augmentation strategy. This technique creates “negative examples” by deliberately misaligning metrics with prompt-response pairs, teaching the model to distinguish between good and poor matches.. This process involves duplicating existing training samples and replacing the original metric with a randomly selected different one from a predetermined list. A low score (1-6) is then assigned to this “misaligned” sample, generating a synthetic negative example where the metric and prompt do not match, indicating a poor pairing. These synthetic negatives are subsequently incorporated into the original training data.

## 2. Exploratory Data Analysis (EDA)

The exploratory data analysis revealed critical insights that drove nearly all modeling decisions. The primary focus was understanding the distribution of target scores and identifying potential biases.

## 2.1 Score Distribution Analysis

### Key Finding: Severe Score Imbalance

#### Typical Training Distribution:

Score	Count	Percentage
0	~50	1.0%
1-8	~200-400	4-8% each
9	~2,800	56.0%
10	~1,450	29.0%
Total	5,000	100%

#### Critical Observations:

1. **Over 85% of data is scores 9-10:** This extreme imbalance means a naive model would simply predict 9 for everything and achieve low training RMSE
2. **Rare scores 0-8:** These represent less than 15% of training data but are critical for generalization
3. **Distribution mismatch:** The training distribution heavily skews toward high scores, but the test set may have different characteristics

## 2.2 Impact of EDA on Modeling Decisions

This single finding drove almost every key decision in our pipeline:

#### 1. Sampling Strategy:

##### StratifiedKFold:

This method ensures validation sets maintain representative score distributions. And also prevents biased validation. As, without stratification, validation might contain only high scores.

##### Implementation:

```
from sklearn.model_selection import StratifiedKFold

# Stratified split to maintain score distribution
targets = [int(float(rec['score'])) for rec in train_data]
skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
all_folds = list(skf.split(range(len(train_data)), targets))

train_idx, val_idx = all_folds[0]
print(f"Train: {len(train_idx)}, Val: {len(val_idx)}")
```

## **2. Loss Function Design:**

**Sample weighting:** Up-weights rare scores.

**Heteroscedastic modeling:** Expresses uncertainty on ambiguous/imbalanced data.

**Focal loss consideration:** Designed for imbalanced classification.

## **3. Evaluation Strategy:**

**Per-bin RMSE analysis:** Tracks performance on each score (0-10), prevents mode collapse (model predicting only mean/median), and reveals weak score ranges.

## **4. Data Augmentation:**

- **Synthetic negatives:** Directly addresses scarcity of low-score examples.
- **Targeted generation:** Creates diverse negative examples to balance training.

## **5. Model Architecture:**

- **Heteroscedastic regression:** Learns to express uncertainty on rare/ambiguous samples.
- **Bilinear interactions:** Captures complex metric-text relationships despite limited low-score data.

### *2.3 Statistical Insights*

#### **Training Data Characteristics:**

- **Mean score:** ~9.1
- **Median score:** 9
- **Standard deviation:** ~0.8
- **Skewness:** Heavily right-skewed (most scores clustered at high end)

#### **Implications of Data Characteristics:**

- **Violation of Standard Regression Assumptions:** The data does not conform to typical regression assumptions such as normality and balanced classes.

- **Extreme Class Imbalance:** Models must be specifically developed to address the significant imbalance between classes.
- **Rigorous Validation Requirements:** Validation procedures will need to incorporate careful stratification to reflect validation results if the test data distribution differs from the validation set.

### 3. Sampling Strategy

Building on EDA insights, the sampling strategy was designed to mitigate the severe score imbalance identified.

#### **Stratified Train/Validation Split:**

- Used StratifiedKFold from sklearn.model\_selection
- The data was split into n\_cv\_folds (e.g., 5) folds
- Stratification was performed on the integer score target to ensure each validation fold is representative
- For single-model training, used the first fold (all\_folds[0]) to create a single train/validation split

#### **Weighted Sampling (Loss Weighting):**

Instead of traditional over/undersampling, implemented a sample weighting at the loss function level:

- Calculated the frequency of each score (0-10) in the training set
- Assigned weights using a Square-Root Inverse Frequency formula:  

$$\text{weight} = (\max\_count / \text{count})^{** 0.5}$$
- This method aggressively up-weights rare scores (e.g., 0, 1, 10) and down-weights common scores (e.g., 5, 6, 7)
- These score\_weights were passed to the training loop and used to multiply the loss for each sample and a per-bin analysis approach.

### 4. Model Architecture

The problem sits at the intersection of classification (11 discrete classes), ordinal regression (scores have a natural order), and continuous regression (predicting values on 0-10 scale). I implemented and tested three distinct architectural approaches.

## 4.1 Shared Architecture Components

All models share a common feature extraction backbone:

```
import torch
import torch.nn as nn
import torch.nn.functional as F
class BaseArchitecture:
    def __init__(self, embedding_dim=768, hidden_dims=[512, 256, 128],
dropout=0.4):
        self.embedding_dim = embedding_dim

        # Bilinear layer for second-order metric-text interactions
        self.bilinear = nn.Bilinear(embedding_dim, embedding_dim,
embedding_dim)

        # MLP with LayerNorm
        self.fc_layers = nn.ModuleList()
        self.layer_norms = nn.ModuleList()
        input_dim = embedding_dim * 3 # metric + text + bilinear

        for hidden_dim in hidden_dims:
            self.fc_layers.append(nn.Linear(input_dim, hidden_dim))
            self.layer_norms.append(nn.LayerNorm(hidden_dim))
            input_dim = hidden_dim

        self.dropout = nn.Dropout(dropout)

    def shared_forward(self, metric_emb, text_emb):
        # Compute bilinear similarity features
        bilinear_out = self.bilinear(metric_emb, text_emb)

        # Concatenate all features
        combined = torch.cat([metric_emb, text_emb, bilinear_out], dim=1)

        # Pass through MLP with LayerNorm
        x = combined
        for fc, ln in zip(self.fc_layers[:-1], self.layer_norms[:-1]):
            x = fc(x)
            x = ln(x) # LayerNorm before activation
            x = F.relu(x)
            x = self.dropout(x)

        # Final layer
        x = self.fc_layers[-1](x)
        x = self.layer_norms[-1](x)
        x = F.relu(x)
        x = self.dropout(x)

    return x
```

## Key Design Choices:

- **Bilinear Layer**: This layer is used to capture second-order interactions between the metric and text embeddings.
- **Layer Normalization (LayerNorm)**: Applied after each Fully Connected (FC) layer to ensure training stability.
- **MLP Structure**: The Multi-Layer Perceptron (MLP) architecture reduces dimensionality through the following sequence: [2304 → 512 → 256 → 128].
- **Dropout**: A dropout rate of 0.4 is implemented for regularization.
- **ReLU Activation**: The standard non-linear activation function used is ReLU.

### 4.2 StandardMetricMatchingModel

Treats the problem as 11-class classification, then computes expected score.

```
class StandardMetricMatchingModel(nn.Module):
    """Standard classification model with expected value prediction."""

    def __init__(self, embedding_dim=768, hidden_dims=[512, 256, 128],
                 dropout=0.4):
        super().__init__()
        self.embedding_dim = embedding_dim
        self.bilinear = nn.Bilinear(embedding_dim, embedding_dim,
                                    embedding_dim)

        self.fc_layers = nn.ModuleList()
        self.layer_norms = nn.ModuleList()
        input_dim = embedding_dim * 3

        for hidden_dim in hidden_dims:
            self.fc_layers.append(nn.Linear(input_dim, hidden_dim))
            self.layer_norms.append(nn.LayerNorm(hidden_dim))
            input_dim = hidden_dim

        # Output 11 logits for classification
        self.output = nn.Linear(input_dim, 11)
        self.dropout = nn.Dropout(dropout)

    def forward(self, metric_emb, text_emb):
        bilinear_out = self.bilinear(metric_emb, text_emb)
        combined = torch.cat([metric_emb, text_emb, bilinear_out], dim=1)

        x = combined
```

```

        for fc, ln in zip(self.fc_layers[:-1], self.layer_norms[:-1]):
            x = fc(x)
            x = ln(x)
            x = F.relu(x)
            x = self.dropout(x)

        x = self.fc_layers[-1](x)
        x = self.layer_norms[-1](x)
        x = F.relu(x)
        x = self.dropout(x)

    return self.output(x)

def predict_score(self, metric_emb, text_emb):
    """Convert logits to expected continuous score."""
    logits = self.forward(metric_emb, text_emb)
    probs = F.softmax(logits, dim=1)

    # Expected value: sum(i * p_i) for i in 0..10
    expected_score = torch.sum(
        probs * torch.arange(11, device=metric_emb.device).float(),
        dim=1
    )
    return expected_score

```

**Loss:** CrossEntropyLoss or HybridLoss

**Strengths:** Simple, interpretable, stable training

**Weaknesses:** Treats scores as independent classes (ignores ordinality)

### 4.3 OrdinalMetricMatchingModel

Models cumulative probabilities  $P(\text{score} \leq k)$  to respect ordinal structure.

```

class OrdinalMetricMatchingModel(nn.Module):
    """Ordinal regression model with cumulative logits."""

    def __init__(self, embedding_dim=768, hidden_dims=[512, 256, 128],
                 dropout=0.4):
        super().__init__()
        self.embedding_dim = embedding_dim
        self.bilinear = nn.Bilinear(embedding_dim, embedding_dim,
                                    embedding_dim)

        self.fc_layers = nn.ModuleList()
        self.layer_norms = nn.ModuleList()
        input_dim = embedding_dim * 3

        for hidden_dim in hidden_dims:
            self.fc_layers.append(nn.Linear(input_dim, hidden_dim))
            self.layer_norms.append(nn.LayerNorm(hidden_dim))

```

```

    input_dim = hidden_dim

    # Output 11 cumulative logits  $P(score \leq k)$ 
    self.output = nn.Linear(input_dim, 11)
    self.dropout = nn.Dropout(dropout)

def forward(self, metric_emb, text_emb):
    bilinear_out = self.bilinear(metric_emb, text_emb)
    combined = torch.cat([metric_emb, text_emb, bilinear_out], dim=1)

    x = combined
    for fc, ln in zip(self.fc_layers[:-1], self.layer_norms[:-1]):
        x = fc(x)
        x = ln(x)
        x = F.relu(x)
        x = self.dropout(x)

    x = self.fc_layers[-1](x)
    x = self.layer_norms[-1](x)
    x = F.relu(x)
    x = self.dropout(x)

    return self.output(x)

def predict_score(self, metric_emb, text_emb):
    """Convert cumulative logits to expected score."""
    logits = self.forward(metric_emb, text_emb)
    cum_probs = torch.sigmoid(logits)  #  $P(score \leq k)$ 

    # Convert to per-class probabilities
    probs = torch.zeros_like(cum_probs)
    probs[:, 0] = cum_probs[:, 0]  #  $P(score = 0)$ 

    for k in range(1, 11):
        probs[:, k] = cum_probs[:, k] - cum_probs[:, k-1]  #  $P(score = k)$ 

    probs = torch.clamp(probs, min=0)
    probs = probs / (probs.sum(dim=1, keepdim=True) + 1e-8)

    # Expected score
    expected_score = torch.sum(
        probs * torch.arange(11, device=metric_emb.device).float(),
        dim=1
    )
    return torch.clamp(expected_score, 0.0, 10.0)

```

**Loss:** OrdinalRegressionLoss (cumulative BCE)

**Strengths:** Theoretically sound for ordered classes, respects score ordering

**Weaknesses:** More complex, harder to train

#### 4.4 HeteroscedasticMatchingModel (Selected)

Predicts both mean and log-variance, allowing uncertainty modeling.

```
class HeteroscedasticMatchingModel(nn.Module):

    def __init__(self, embedding_dim=768, hidden_dims=[512, 256, 128],
dropout=0.4):
        super().__init__()
        self.embedding_dim = embedding_dim
        self.bilinear = nn.Bilinear(embedding_dim, embedding_dim,
embedding_dim)

        self.fc_layers = nn.ModuleList()
        self.layer_norms = nn.ModuleList()
        input_dim = embedding_dim * 3

        for hidden_dim in hidden_dims:
            self.fc_layers.append(nn.Linear(input_dim, hidden_dim))
            self.layer_norms.append(nn.LayerNorm(hidden_dim))
            input_dim = hidden_dim

        self.mean_head = nn.Linear(input_dim, 1)
        self.logvar_head = nn.Linear(input_dim, 1)
        self.dropout = nn.Dropout(dropout)

    def forward(self, metric_emb, text_emb):
        bilinear_out = self.bilinear(metric_emb, text_emb)
        combined = torch.cat([metric_emb, text_emb, bilinear_out], dim=1)

        x = combined
        for fc, ln in zip(self.fc_layers[:-1], self.layer_norms[:-1]):
            x = fc(x)
            x = ln(x)
            x = F.relu(x)
            x = self.dropout(x)

        x = self.fc_layers[-1](x)
        x = self.layer_norms[-1](x)
        x = F.relu(x)
        x = self.dropout(x)

        mean = self.mean_head(x)
        logvar = self.logvar_head(x)
        return mean.squeeze(-1), torch.clamp(logvar, -10.0, 10.0).squeeze(-1)

    def predict_score(self, metric_emb, text_emb):
        mean, _ = self.forward(metric_emb, text_emb)
        return torch.clamp(mean, 0.0, 10.0)
```

The HeteroscedasticLoss function, defined as  $\text{precision} \times \text{squared\_error} + \text{logvar}$ , was chosen for its ability to model uncertainty effectively, especially with noisy or imbalanced data. It also aligns well with RMSE and demonstrated superior performance on the leaderboard.

## 5. Loss Functions

I have implemented and tested multiple loss functions to handle the imbalanced regression problem.

### 5.1 CrossEntropyLoss

Standard multi-class classification loss.

```
# Simple CrossEntropyLoss with expected value
criterion = nn.CrossEntropyLoss()
logits = model(metric_emb, text_emb)
loss = criterion(logits, target)

# Prediction
probs = F.softmax(logits, dim=1)
expected_score = torch.sum(probs * torch.arange(11, device=device).float(),
dim=1)
```

### 5.2 OrdinalRegressionLoss

Cumulative link loss for ordinal data.

```
class OrdinalRegressionLoss(nn.Module):
    """Cumulative link loss for ordered classes."""

    def __init__(self):
        super().__init__()
        self.bce_loss = nn.BCEWithLogitsLoss()

    def forward(self, logits, target, sample_weights=None,
               return_per_sample=False):
        """
        Args:
            logits: (B, 11) - Logits for  $P(\text{score} \leq k)$ 
            target: (B,) - integer targets 0-10
        """
        batch_size = logits.shape[0]
        targets_cum = torch.zeros_like(logits)

        # Create cumulative targets:  $P(\text{score} \leq k) = 1$  for  $k \geq \text{target}$ 
        for i in range(batch_size):
```

```

k = target[i].item()
targets_cum[i, :k+1] = 1.0
targets_cum[i, k+1:] = 0.0

# BCE Loss per sample
bce_per_sample = F.binary_cross_entropy_with_logits(
    logits, targets_cum, reduction='none'
)
loss_per_sample = bce_per_sample.mean(dim=1)

# Apply sample weights
if sample_weights is not None:
    loss_per_sample = loss_per_sample * sample_weights

if return_per_sample:
    return loss_per_sample

return loss_per_sample.mean()

```

### 5.3 HybridLoss

Combines classification and regression objectives.

```

class HybridLoss(nn.Module):
    """Combines CE and MSE for classification + regression."""

    def __init__(self, ce_weight=1.0, mse_weight=0.1):
        super().__init__()
        self.ce_loss = nn.CrossEntropyLoss(reduction='none')
        self.mse_weight = mse_weight
        self.ce_weight = ce_weight

    def forward(self, logits, target, sample_weights=None,
               return_per_sample=False):
        # CrossEntropyLoss
        ce = self.ce_loss(logits, target)

        # MSE on expected score
        probs = F.softmax(logits, dim=1)
        expected_scores = torch.sum(
            probs * torch.arange(11, device=logits.device).float(),
            dim=1
        )
        mse = (expected_scores - target.float()) ** 2

        # Combined loss per sample
        loss_per_sample = self.ce_weight * ce + self.mse_weight * mse

        # Apply sample weights
        if sample_weights is not None:

```

```

        loss_per_sample = loss_per_sample * sample_weights

    if return_per_sample:
        return loss_per_sample

    return loss_per_sample.mean()

```

## 5.4 FocalLoss

Focuses on hard examples for imbalanced data.

```

class FocalLoss(nn.Module):
    """Focal loss for imbalanced data."""

    def __init__(self, alpha=1.0, gamma=2.0):
        super().__init__()
        self.alpha = alpha
        self.gamma = gamma

    def forward(self, logits, target):
        ce = F.cross_entropy(logits, target, reduction='none')
        pt = torch.exp(-ce)
        return (self.alpha * (1 - pt) ** self.gamma * ce).mean()

```

## 5.5 HeteroscedasticLoss (Selected)

Models prediction uncertainty.

```

class HeteroscedasticLoss(nn.Module):

    def __init__(self):
        super().__init__()

    def forward(self, mean, logvar, target, sample_weights=None,
               return_per_sample=False):

        target_float = target.float()
        variance = torch.exp(logvar)
        precision = 1.0 / (variance + 1e-6)
        squared_error = (mean - target_float) ** 2
        loss_per_sample = precision * squared_error + logvar

        if sample_weights is not None:
            loss_per_sample = loss_per_sample * sample_weights

        if return_per_sample:
            return loss_per_sample

        return loss_per_sample.mean()

```

## **Mathematical Formulation:**

The loss function, denoted as (L), is defined as:

$$[L = (y - \mu)^2 + \sigma^2]$$

Where ( $\mu$ ) represents the predicted mean, ( $\sigma^2$ ) represents the predicted variance (derived from log-variance), and (y) signifies the target value.

## **6. Sampling Strategy & Weighting**

To address severe class imbalance, I use stratified sampling and sqrt inverse frequency weighting.

### **6.1 StratifiedKFold**

Used 5-fold StratifiedKFold maintaining score distribution.

### **6.2 Sample Weighting Implementation**

#### **Configuration:**

- **max\_weight:** The maximum allowed weight, capped at 8.0.
- **min\_weight:** The minimum allowed weight, set at 0.5.
- **boost\_high\_scores:** A boolean value, set to True, which provides an additional boost for scores ranging from 9 to 10.
- **sample\_weight\_temp:** The temperature used for calculating weights, set at 0.5.

#### **Example Weights:**

Score	Count	Weight	Notes
0	50	~8.0	Maximum up-weighting
5	300	~2.6	Moderate
9	2800	~1.0→1.5	Most common, boosted
10	1450	~1.2→1.6	Second most, boosted

Square root weighting offers several advantages: It provides increased stability compared to inverse frequency weighting, reduces loss variability by mitigating extreme weights, and effectively preserves the signal for rare classes.

## 7. Encoder Fine-tuning

To adapt the multilingual sentence transformer to the competition domain, I employed SimCSE-style contrastive fine-tuning. This technique improves text embeddings before main model training.

### 7.1 SimCSE-Style Contrastive Learning

SimCSE creates positive pairs by treating a sentence as similar to itself (with dropout), learning better representations through contrastive learning.

#### Implementation:

```
def fine_tune_encoder_simcse(text_encoder, train_data, device, epochs=3,
batch_size=64, lr=2e-5):

    texts = [combine_text_fields(rec) for rec in train_data]
    examples = [InputExample(texts=[t, t], label=1.0) for t in texts]

    train_dataloader = TorchDataLoader(
        examples,
        shuffle=True,
        batch_size=batch_size,
        drop_last=False
    )

    num_steps = len(train_dataloader) * epochs
    warmup_steps = max(100, num_steps // 10)

    train_loss = losses.MultipleNegativesRankingLoss(model=text_encoder)
    # Fine-tune using SentenceTransformer.fit() API
    output_dir_ft = CONFIG['encoder_dir']
    os.makedirs(output_dir_ft, exist_ok=True)

    text_encoder.fit(
        train_objectives=[(train_dataloader, train_loss)],
        epochs=epochs,
        warmup_steps=warmup_steps,
        optimizer_params={'lr': lr},
        show_progress_bar=True,
        output_path=output_dir_ft,
        use_amp=True # Mixed precision for faster training
    )

    return text_encoder
```

## 7.2 Loss Function Selection

### MultipleNegativesRankingLoss (Preferred):

- More effective for large batches (64+)
- Treats other samples in the batch as negatives
- Faster convergence than ContrastiveLoss

### ContrastiveLoss (Fallback):

- Used if MultipleNegativesRankingLoss fails
- Requires explicit positive/negative pairs
- Still effective but slower

## 7.3 Configuration

```
CONFIG = {
    'fine_tune_encoder': True,
    'encoder_ft_epochs': 3, # Small number of epochs for adaptation
    'encoder_ft_lr': 2e-5, # Much smaller than head LR (1e-4)
    'encoder_ft_batch_size': 64, # Larger batch for better contrastive
    learning}
```

## 7.4 Impact

- **Domain adaptation:** Better captures competition-specific text patterns
- **Multilingual improvement:** Enhanced handling of Tamil, Hindi, and English
- **Performance boost:** Typically 0.1-0.2 RMSE improvement on validation set

## 8. Distribution Bias Fixes

A significant train/validation mismatch emerged: training RMSE was 0.93 but leaderboard RMSE was ~4.6.

### 8.1 Problem Analysis

**Root cause:** - Training: 91% scores 9–10 - Predictions: 99% scores 9–10 - Test likely includes more 0–8 labels - Result: collapse to score 9 with large errors on low scores

#### Evidence:

Training distribution:	56% score 9, 29% score 10
Model predictions:	85% score 9, 14% score 10 (too conservative)
Leaderboard RMSE:	~4.6 (model fails on unseen low scores)

## 8.2 Solution 1: Exponential Reweighting

Up-weight rare scores via sqrt inverse frequency.

**Configuration:**

```
CONFIG = {
    'use_weights': True,
    'max_weight': 8.0, # Increased from 3.0 for stronger upweighting
    'boost_high_scores': True, # 1.5x boost for score 9, 1.3x for score 10
}
```

## 8.3 Solution 2: Synthetic Negatives

Generate low-score examples by misaligning metrics with prompt-response pairs.

```
CONFIG = {
    'use_synthetic_negatives': True,
    'synthetic_ratio': 0.10, # 10% of training data
    'synthetic_score_range': (1, 6), # Low scores for mismatched pairs
}
```

## 8.4 Solution 3: Heteroscedastic Model

Model both mean and variance to handle uncertainty.

**Heteroscedastic loss:**

**Before fixes:** - Validation RMSE: 0.93 - Leaderboard RMSE: 4.6 - Prediction spread: 85% score 9, 14% score 10

**After fixes:** - Validation RMSE: 1.0–1.1 - Leaderboard RMSE: **3.534** - Prediction spread: 60–70% score 9, better 8–10 coverage

**Trade-off:** Validation RMSE rises slightly due to more diverse predictions with **23% improvement in leaderboard RMSE** confirms alignment with test distribution.

# 9. Training Improvements

Several improvements were applied to stabilize and speed training.

## 9.1 Architecture Enhancements

**Regularization:**

- **Layer Normalization:** Applied after fully connected layers and before ReLU for stable training and improved convergence.
- **Dropout:** Increased from 0.3 to 0.4 to prevent overfitting.

- **Weight Decay (L2 regularization):** A  $1e-5$  weight decay with Adam optimizer encourages smaller weights and reduces overfitting.

### Optimization and Training Control:

- **Early Stopping:** Training halts if validation RMSE doesn't improve for 20 epochs, saving the best model.
- **Learning Rate Scheduling:** ReduceLROnPlateau halves the learning rate (down to  $1e-6$ ) if validation RMSE plateaus for 5 epochs.
- **Mixed Precision Training (AMP):** `torch.cuda.amp` accelerates training and reduces memory using lower precision.
- **Gradient Clipping:** With a value of 1.0, it prevents exploding gradients, applied after the backward pass (with `scaler.unscale_()` for AMP).

## 9.2 Configuration Summary

```
CONFIG = {
    # Architecture
    'dropout': 0.4,  # Increased from 0.3
    'weight_decay': 1e-5,  # L2 regularization

    # Training
    'batch_size': 64,
    'learning_rate': 1e-4,
    'epochs': 100,
    'patience': 20,  # Early stopping
    'use_amp': True,  # Mixed precision

    # Optimization
    'lr_patience': 5,  # LR scheduler patience
    'min_lr': 1e-6,  # Minimum Learning rate
    'gradient_clip': 1.0,  # Gradient clipping

    # Checkpointing
    'checkpoint_dir': '/kaggle/working/checkpoints',
    'save_interval': 10,
    'auto_save': True,
}
```

**Impact:** Faster training, less overfitting, better generalisation.

## 10. Inference Pipeline

The inference pipeline applies the trained model to the test dataset to generate test predictions for submission while keeping memory, corruption handling, and domain adaptation in mind.

### 10.1 Safe Data Handling

All JSON inputs (metric\_names.json, metric\_name\_embeddings.npy, test\_data.json) are loaded with a `safe_load_json()` helper that handles `JSONDecodeError` and `UnicodeDecodeError`. This prevents crashes on partially corrupted files and attempts alternative encodings or partial recovery before failing explicitly.

### 10.2 Text combination & encoder handling

A small helper `combine_text_fields(record)` normalises variable record schemas (system\_prompt, user\_prompt/prompt, response/expected\_response) into a single text string per sample. The pipeline supports two encoder modes:

- **Auto-finetune:** load base `paraphrase-multilingual-mpnet-base-v2`, fine-tune using SimCSE-style contrastive learning on `train + test` data (if available) with `MultipleNegativesRankingLoss` (fallback to `ContrastiveLoss`) through the SentenceTransformer `fit()` API.
- **Load finetuned encoder:** if a cached encoder exists, it's loaded instead.

Finetuning uses a safe DataLoader wrapper and computes reasonable `warmup_steps` ( $\approx 10\%$  of total steps). After encoding, the encoder is explicitly deleted and `torch.cuda.empty_cache()` is called to free GPU memory.

### 10.3 Encoding test texts

`text_encoder.encode(...)` runs in batches (configurable). The encoder returns 768-d text embeddings (numpy), which are paired with pre-computed metric embeddings (lookup via `metric_names_map`) to create a `MetricMatchingDataset`

### 10.4 Model Loading with Flexible Checkpoints

The heteroscedastic model architecture (bilinear + MLP  $\rightarrow$  mean & logvar heads) is instantiated and loaded from either:

- Full checkpoint dict with `model_state_dict`, or

- Raw state dict saved by `torch.save(model.state_dict())`. Errors are caught and surfaced with clear messages.

## *10.5 Post-processing & submission*

Predicted means are clamped to [0,10], rounded to integers (to match evaluation), and re-ordered to match original sample IDs. The final CSV contains `ID` (1..N) and `score`. The pipeline prints distribution statistics for quick sanity checks.

# **11. Results & Performance**

## *11.1 Training Performance*

The heteroscedastic model with bias mitigation converged quickly and achieved strong results.

### **Best Model Performance:**

**Best validation RMSE:** 1.6959 (Epoch 47)

**Final RMSE:** 1.7323 (Epoch 67)

**Early stopping:** After epoch 67 (patience=20)

**Training accuracy:** ~81%

### **Per-Bin RMSE Analysis:**

```
Score  0: RMSE= 2.10, Count=    50
Score  1: RMSE= 1.85, Count=    89
Score  2: RMSE= 1.50, Count=   156
Score  8: RMSE= 1.15, Count=   450
Score  9: RMSE= 0.90, Count= 2800
Score 10: RMSE= 0.95, Count= 1450
Overall: RMSE= 1.02, Total= 5000
```

## *11.2 Leaderboard Performance*

**Final Leaderboard RMSE:** 3.534

This represents a **23% improvement** from the initial 4.6 baseline through systematic bias mitigation and uncertainty modeling.

