

# Q&A Session Summary

Here is a complete log of the questions you've asked and the answers I've provided during our conversation.

## Question 1: Data Splitting with Heavy Right-Tail

**Q:** "Given the nature of the heavy right-tail data, how would you do a split?"

**A:** Given the heavy right-tail in the Number of items\_scanned data, the primary concern for splitting the data is to ensure that the rare but important transactions (those with a very high number of items) are properly represented in both the training and testing sets. A simple random split could accidentally place most of these outliers in one set, leading to poor model performance.

Here are the recommended approaches to splitting the data.

### Method 1: Stratified Split (Recommended)

This method directly addresses the problem of the heavy-tailed distribution. **Stratified splitting** ensures that the distribution of Number of items\_scanned is preserved across the training and test sets.

#### How to do it:

1. **Bin the Data:** Since Number of items\_scanned is a continuous variable, you first need to create discrete categories or "bins". A good way to do this is by using quantiles (e.g., deciles, which would create 10 bins with an equal number of transactions in each). This ensures that transactions from the long tail are grouped into their own bin.
2. **Perform the Split:** Use these new bins as the basis for stratification. When you perform the train-test split, you instruct the function to maintain the same percentage of samples for each bin in both the train and test datasets.

This guarantees that both your training and test sets will have a representative sample of transactions, from the most common (few items) to the most extreme (many items). You could also stratify by Loyalty\_tier at the same time for even better representation.

### Method 2: Time-Based Split (A Critical Consideration)

The presence of a Seq\_no column implies that the data is ordered by time. If the goal is to build a model that predicts future transactions (a common business objective), you **must** use a time-based split to simulate a real-world scenario and avoid data leakage.

#### How to do it:

1. **Sort the Data:** Ensure the entire dataset is sorted by Seq\_no in ascending order.
2. **Split by Time:** Designate the first portion of the data for training and the latter portion

for testing. For example, with 1.7 million transactions and an 80/20 split:

- **Training Set:** The first 80% of transactions (approx. 1.36 million).
- **Test Set:** The last 20% of transactions (approx. 0.34 million).

**Caution:** While this is the correct method for forecasting, you lose the guarantee of a balanced distribution for Number of items\_scanned. After performing a time-based split, you should still analyze the distributions in your train and test sets to ensure the test set isn't fundamentally different in a way that would make evaluation unreliable (a phenomenon known as concept drift).

## Conclusion

- For a **general-purpose model** where time is not a predictive factor, use a **stratified split** based on binned Number of items\_scanned. This is the most direct solution to the heavy-tail problem.
- For a **forecasting model** that predicts future behavior, you must use a **time-based split** based on Seq\_no.

## Question 2: Heavy Right Tail Explained

**Q:** "Whats heavy right taik"

**A:** A **heavy right tail** describes a specific shape in a set of data when you graph it. It means that while most of your data points are clustered on the left side (representing lower values), there are a significant number of extreme high values that stretch out far to the right. 

Think of it like this:

- **Tail:** The "tails" are the far ends of a distribution graph.
- **Right Tail:** This means the tail stretches out to the **right**, in the direction of higher values. This is also called a "right-skewed" distribution.
- **Heavy:** This is the key part. "Heavy" means there are **more data points in the tail than you'd expect** in a standard distribution (like a bell curve). These aren't just one or two flukes; they are a meaningful number of extreme outliers.

## A Simple Example: Income 💰

Imagine you look at the income of 1,000 people.

- Most people will earn a similar, average amount. This creates a large "lump" on the left side of the graph.
- However, a few individuals might be CEOs or investors who earn extremely high incomes. These few, super-high values create the long, heavy tail stretching out to the right.

Because of this heavy tail, the **average income** would be much higher than what most people actually earn, making the **median** (the middle value) a much better representation of the typical person's income.

In your previous problem about grocery transactions, the "Number of items\_scanned" had a heavy right tail because most shoppers buy a small number of items, but a few (like businesses or people hosting large events) buy a massive number of items, creating those extreme values in the tail.

### Question 3: Why Not Random Split?

**Q: "Whh not random split here"**

**A:** A random split is a poor choice here because it ignores two critical features of your dataset, creating a high risk of building an unreliable model.

#### The Risk with Heavy Tails: Unrepresentative Samples

The main problem is the **heavy right-tail** in the Number of items\_scanned. This means you have a small but important number of "extreme" transactions (e.g., a business buying 500 items).

When you split randomly, you're relying on chance. It's very likely that these rare, extreme transactions could end up disproportionately in one set:

- **Scenario 1 (Worse): Outliers end up only in the Test Set.** Your model has never been trained on what a massive transaction looks like. When it encounters them during testing, it will perform poorly, and you won't know why. The model is essentially blind to a crucial segment of your data.
- **Scenario 2: Outliers end up only in the Training Set.** Your model learns from the extreme cases, but your test set is now "too easy" because it contains no large transactions. You can't validate how well the model handles them, leading to an overly optimistic performance score.

A **stratified split** solves this by forcing the train and test sets to have the same percentage of extreme transactions as the original dataset, guaranteeing both are representative.

#### The Risk with Time Data: Data Leakage

The Seq\_no column implies the data is a **time series**—transactions are listed in chronological order. Randomly shuffling this data is a major mistake in data science called **lookahead bias**.

By mixing past and future data, your model can "cheat." It might learn from information in the training set that would not have been available at the time of a transaction in the test set.

**Example:** Imagine a new loyalty program starts on transaction #1,000,000. A random split might put a transaction from *after* the program started into the training set, and a transaction from *before* it started into the test set. The model could use its "future knowledge" of the program to make an impossible prediction about the past.

This results in a model that looks fantastic during testing but fails in the real world. A

**time-based split** (e.g., train on the first 80%, test on the most recent 20%) prevents this by simulating how the model would actually be used in production.

## Question 4: Cross-Entropy vs. MSE/MAE in Classification

**Q:** "In classification tasks, we often use Binary Cross-Entropy (BCE) or Cross-Entropy (CE) as the loss function. Why don't we use distance-based metrics like Mean Squared Error (MSE) or Mean Absolute Error (MAE) instead? What are the benefits of using BCE/CE compared to MSE/MAE?"

**A:** In classification, we don't use distance-based loss functions like MSE or MAE because they are **not designed to measure the error of probabilistic predictions**, which can lead to major issues during model training. Cross-Entropy (BCE/CE) is used instead because it is specifically designed to evaluate how well a model's predicted probabilities match the true class labels.

### Why MSE and MAE Fail for Classification

Using a distance-based metric like Mean Squared Error (MSE) is a poor fit for classification for two main reasons:

#### 1. It Penalizes Confidently Wrong Predictions Too Little

In classification, models output a probability (e.g., "80% chance this is a cat"). If the answer is actually "cat" (label = 1) and the model confidently predicts "dog" (probability of cat = 0.01), this is a huge mistake that needs a strong correction.

However, MSE's penalty is small. The squared error would be  $(1 - 0.01)^2 = 0.98$ . If it was slightly less confident (e.g., probability = 0.1), the error would be  $(1 - 0.1)^2 = 0.81$ . The difference in penalty is small, even though the first mistake is much worse.

This leads to the **vanishing gradient problem**. The model receives a very weak signal (a small gradient) to correct major errors, causing it to learn extremely slowly or get stuck.

#### 2. It Creates a Difficult Learning Landscape

For the activation functions used in classification (like Sigmoid or Softmax), using MSE creates a non-convex loss landscape. This means the landscape is filled with many "local minima" (bad solutions) and flat plateaus. An optimization algorithm can easily get stuck in one of these bad solutions and fail to find the best possible answer.

### The Benefits of Cross-Entropy (BCE/CE)

Cross-Entropy works better because it directly addresses the failures of MSE.

#### 1. It Severely Penalizes Confidently Wrong Predictions

Cross-Entropy uses a logarithmic ( $\log$ ) function, which gives it a powerful property. As the

model's predicted probability for the **correct** class gets closer to 0, its loss value skyrockets towards infinity.

- **Correct Class Label:** 1 ("cat")
- **Model Prediction (Bad):** Predicts 0.01 for "cat". CE Loss =  $-\log(0.01) \approx 4.6\$$
- **Model Prediction (Terrible):** Predicts 0.0001 for "cat". CE Loss =  $-\log(0.0001) \approx 9.2\$$

This massive loss creates a large gradient, giving the model a strong "push" to correct its confident mistakes and learn much more effectively.

## 2. It Creates an Easy Learning Landscape

When paired with a Sigmoid or Softmax activation function, Cross-Entropy loss creates a **convex** loss landscape. This landscape looks like a simple bowl with only one global minimum at the bottom. This makes it very easy for the optimization algorithm to navigate directly to the best possible solution without getting stuck.

### Simple Analogy

Think of it this way:

- **MSE/MAE** is like using a **ruler**. It's great for measuring the physical distance between two points (a regression problem).
- **Cross-Entropy** is like a **specialized exam**. It doesn't just measure if you're right or wrong; it measures your confidence and heavily penalizes you for being confidently incorrect.

You wouldn't use a ruler to score an exam, and you shouldn't use a distance metric for a probability-based classification task.

## Question 5: Dog vs. Cat Classification Example

**Q:** "Compare them for a dog cat"

**A:** Of course. Let's compare them using a simple dog vs. cat classification scenario.

Imagine our model's job is to look at an image and output a single probability: "How likely is this a cat?"

- **Label for Cat:** 1
- **Label for Dog:** 0

### Scenario: A Confident but Incorrect Prediction

Let's say we show our model an image that is **clearly a cat** . The correct label is 1.

However, our untrained model makes a terrible mistake and is **99% confident it's a dog**. This means its output probability for "cat" is just 0.01.

Let's see how the two loss functions react to this major error.

### Mean Squared Error (MSE) Reaction 📈

MSE calculates the simple squared distance.

- **Formula:**  $(\text{True Label} - \text{Predicted Probability})^2$
- **Calculation:**  $(1 - 0.01)^2 = (0.99)^2 \approx \mathbf{0.98}$

The loss value is 0.98. This doesn't seem very high for such a confident and completely wrong prediction.

### Cross-Entropy (BCE) Reaction 🎓

Cross-Entropy calculates the "surprise" or error in the probability.

- **Formula:**  $-\log(\text{Predicted Probability for Correct Class})$
- **Calculation:**  $-\log(0.01) \approx \mathbf{4.61}$

The loss value is 4.61, which is **nearly 5 times higher** than the MSE loss!

This shows that Cross-Entropy heavily penalizes the model for being both confident and wrong, creating a much stronger signal that a major correction is needed.

## Why This Difference Matters: The Learning Signal

The loss value is used to create a "learning signal" (the gradient) that tells the model how to adjust itself.

- **With MSE:** The learning signal for this mistake is **tiny**. The model barely gets any feedback that it made a huge error, so it learns very slowly or gets stuck. This is the **vanishing gradient problem**.
- **With Cross-Entropy:** The huge loss value creates a **strong learning signal**. The model gets a powerful push in the right direction, forcing it to learn quickly from its confident mistakes.

In short, for classifying cats and dogs, **MSE is like a teacher who barely reacts to a wrong answer, while Cross-Entropy is like a teacher who provides firm, clear correction, leading to much faster and more effective learning.**

## Question 6: Benefits of Cross-Entropy (Summary)

Q: "Benefits"

A: Certainly. Here are the key benefits of using Cross-Entropy (BCE/CE) for classification tasks, summarized from our discussion:

- **Faster & More Reliable Training:** It provides a strong learning signal (gradient) when the model makes confident mistakes, directly preventing the "vanishing gradient" problem that often stalls training when using MSE.

- **Better Error Penalties:** It penalizes a confident wrong prediction much more heavily than a hesitant one. This focuses the model's learning on correcting its most significant errors first.
- **Easier Optimization:** When used with standard classification outputs (like a softmax or sigmoid function), it creates a convex loss landscape. This means it's like a simple bowl with one bottom, making it easy for the model to find the best possible solution without getting stuck.
- **Mathematically Suited for the Task:** Cross-Entropy is fundamentally designed in information theory to measure the difference between two probability distributions—exactly what you're doing in classification when comparing the model's predicted probabilities to the true labels.

## Question 7: Backpropagation Calculation

**Q:** "Please perform backpropagation step by step to calculate the gradients with respect to the following variables: [w0, x0, w1, x1, w2]."

**A:** Of course. Here is a step-by-step calculation of the gradients using backpropagation.

Backpropagation works by starting from the final output and moving backward through the graph, calculating the gradient at each step using the chain rule. The gradient of the final output with respect to itself is always 1.

For clarity, let's label the intermediate results from right to left, using the green values provided in the graph.

- $L = 0.27$  (Final Output)
- $a = 3.7$
- $b = 2.7$
- $c = 1$
- $d = -1$
- $e = 2$
- $g = -2$
- $h = 4$

We will calculate the gradient of the final output  $L$  with respect to each variable.

### Step 1: Gradient at the $1/x$ node

The last operation is  $L = 1/a$ . We need to find the gradient of  $L$  with respect to  $a$ .

- **Upstream Gradient:**  $\frac{\partial L}{\partial L} = 1$
- **Local Gradient:** The derivative of  $1/a$  with respect to  $a$  is  $-1/a^2$ .
- **Gradient:**  $\frac{\partial L}{\partial a} = \frac{\partial L}{\partial L} \cdot \frac{\partial L}{\partial a} = 1 \cdot -1/(3.7)^2 = -1/13.69 \approx -0.073$

### Step 2: Gradient at the $+1$ node

The operation is  $a = b + 1$ . We pass the gradient from a back to b.

- **Upstream Gradient:**  $\frac{\partial L}{\partial a} = -0.073$
- **Local Gradient:** The derivative of  $b+1$  with respect to b is  $1$ .
- **Gradient:**  $\frac{\partial L}{\partial b} = \frac{\partial L}{\partial a} \cdot 1 = -0.073 \cdot 1 = \mathbf{-0.073}$

### Step 3: Gradient at the exp node

The operation is  $b = \exp(c)$ . We pass the gradient from b back to c.

- **Upstream Gradient:**  $\frac{\partial L}{\partial b} = -0.073$
- **Local Gradient:** The derivative of  $\exp(c)$  is  $\exp(c)$ . From the graph, we know  $\exp(c) = b = 2.7$ .
- **Gradient:**  $\frac{\partial L}{\partial c} = \frac{\partial L}{\partial b} \cdot \exp(c) = -0.073 \cdot 2.7 \approx \mathbf{-0.197}$

### Step 4: Gradient at the \*-1 node

The operation is  $c = d \cdot (-1)$ . We pass the gradient from c back to d.

- **Upstream Gradient:**  $\frac{\partial L}{\partial c} = -0.197$
- **Local Gradient:** The derivative of  $d \cdot (-1)$  with respect to d is  $-1$ .
- **Gradient:**  $\frac{\partial L}{\partial d} = \frac{\partial L}{\partial c} \cdot (-1) = -0.197 \cdot (-1) = \mathbf{0.197}$

### Step 5: Gradient at the $d = e + w_2$ node

The operation is addition. The gradient from d is passed back equally to e and  $w_2$  because the local gradient for both inputs of an addition gate is 1.

- **Upstream Gradient:**  $\frac{\partial L}{\partial d} = 0.197$
- **Gradient for e:**  $\frac{\partial L}{\partial e} = \frac{\partial L}{\partial d} \cdot 1 = \mathbf{0.197}$
- **Gradient for  $w_2$ :**  $\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial d} \cdot 1 = \mathbf{0.197}$

### Step 6: Gradient at the $e = g + h$ node

This operation is also addition. The gradient from e is passed back equally to g and h.

- **Upstream Gradient:**  $\frac{\partial L}{\partial e} = 0.197$
- **Gradient for g:**  $\frac{\partial L}{\partial g} = \frac{\partial L}{\partial e} \cdot 1 = \mathbf{0.197}$
- **Gradient for h:**  $\frac{\partial L}{\partial h} = \frac{\partial L}{\partial e} \cdot 1 = \mathbf{0.197}$

### Step 7: Gradients for the top \* node ( $g = w_0 * x_0$ )

The operation is multiplication. The local gradient for one input is the value of the other input.

- **Upstream Gradient:**  $\frac{\partial L}{\partial g} = 0.197$

- **Gradient for w0:**  $\frac{\partial L}{\partial w_0} = \frac{\partial L}{\partial g} \cdot x_0 = 0.197 \cdot (-1) = \mathbf{-0.197}$
- **Gradient for x0:**  $\frac{\partial L}{\partial x_0} = \frac{\partial L}{\partial g} \cdot w_0 = 0.197 \cdot 2 = \mathbf{0.394}$

### Step 8: Gradients for the bottom \* node ( $h = w_1 * x_1$ )

The operation is multiplication, just like the one above.

- **Upstream Gradient:**  $\frac{\partial L}{\partial h} = 0.197$
- **Gradient for w1:**  $\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial h} \cdot x_1 = 0.197 \cdot (-2) = \mathbf{-0.394}$
- **Gradient for x1:**  $\frac{\partial L}{\partial x_1} = \frac{\partial L}{\partial h} \cdot w_1 = 0.197 \cdot (-2) = \mathbf{-0.394}$

### Final Calculated Gradients

Here are the final gradients for the requested variables:

- $\frac{\partial L}{\partial w_0} = \mathbf{-0.197}$
- $\frac{\partial L}{\partial x_0} = \mathbf{0.394}$
- $\frac{\partial L}{\partial w_1} = \mathbf{-0.394}$
- $\frac{\partial L}{\partial x_1} = \mathbf{-0.394}$
- $\frac{\partial L}{\partial w_2} = \mathbf{0.197}$

### Question 8: Backpropagation Verification

**Q:** "Are w0 and x0 correct?"

**A:** Yes, the calculated gradients for w0 and x0 are correct.

Here's a quick check of the logic for that step.

For the multiplication node  $g = w_0 * x_0$ , the chain rule states that the gradient flowing back from  $g$  is multiplied by the local gradient. For a multiplication gate, the local gradient for one input is simply the value of the **other input**.

The upstream gradient flowing into this node was  $\frac{\partial L}{\partial g} \approx 0.197$ .

- For w0:  
The gradient is the upstream gradient multiplied by the value of x0.  
$$\frac{\partial L}{\partial w_0} = (\text{gradient from } g) \cdot x_0 = 0.197 \cdot (-1) = \mathbf{-0.197}$$
- For x0:  
The gradient is the upstream gradient multiplied by the value of w0.  
$$\frac{\partial L}{\partial x_0} = (\text{gradient from } g) \cdot w_0 = 0.197 \cdot 2 = \mathbf{0.394}$$

Both calculations hold up.

## Question 9: Dice Rolling Probability Problem

**Q:** "Three players each roll a fair six-sided die (faces numbered 1-6) in a round. After each round, the smallest number among the three rolls is recorded. After 10,000 rounds, approximate the total of all recorded smallest values? Explain your reasoning."

**A:** To solve this, we can approximate the total by finding the **expected value** of the smallest number in a single round and then multiplying it by the number of rounds (10,000).

### Reasoning and Calculation

Let  $Y$  be the random variable representing the smallest of the three dice rolls. We need to find the expected value of  $Y$ , or  $E[Y]$ . A straightforward way to calculate the expected value for a variable with non-negative integer outcomes is using the formula:

$$\text{\$\$} E[Y] = \sum_{k=1}^{\infty} P(Y \geq k) \text{\$\$}$$

#### 1. Find the Probability $P(Y \geq k)$

The smallest roll  $Y$  is greater than or equal to a number  $k$  if and only if **all three dice** show a number greater than or equal to  $k$ . Since the dice rolls are independent, we can multiply their probabilities.

The probability that a single fair die roll is greater than or equal to  $k$ ,  $P(\text{die} \geq k)$ , is:

- $P(\text{die} \geq 1) = 6/6$
- $P(\text{die} \geq 2) = 5/6$
- $P(\text{die} \geq 3) = 4/6$
- $P(\text{die} \geq 4) = 3/6$
- $P(\text{die} \geq 5) = 2/6$
- $P(\text{die} \geq 6) = 1/6$

Now, we can find  $P(Y \geq k)$  by cubing these probabilities:

- $P(Y \geq 1) = (6/6)^3 = 1$
- $P(Y \geq 2) = (5/6)^3 = 125/216$
- $P(Y \geq 3) = (4/6)^3 = 64/216$
- $P(Y \geq 4) = (3/6)^3 = 27/216$
- $P(Y \geq 5) = (2/6)^3 = 8/216$
- $P(Y \geq 6) = (1/6)^3 = 1/216$

#### 2. Calculate the Expected Value $E[Y]$

Now, we sum these probabilities to get the expected value of the smallest roll in one round:

$$\text{\$\$} E[Y] = P(Y \geq 1) + P(Y \geq 2) + P(Y \geq 3) + P(Y \geq 4) + P(Y \geq 5) + P(Y \geq 6) \text{\$\$} E[Y] =$$

$$\frac{1}{216} + \frac{1}{125} + \frac{1}{64} + \frac{1}{27} + \frac{1}{8} + \frac{1}{1} = \frac{1}{216} + \frac{1}{125} + \frac{1}{64} + \frac{1}{27} + \frac{1}{8} + \frac{1}{1} = \frac{441}{216}$$

Simplifying the fraction (dividing by 9) gives:

$$E[Y] = \frac{49}{24} \approx 2.0417$$

So, on average, the smallest number recorded in each round is about 2.0417.

### 3. Approximate the Total for 10,000 Rounds

Finally, we multiply the expected value per round by the total number of rounds:

$$Total\ Sum \approx E[Y] \times 10,000 = \frac{49}{24} \times 10,000 = \frac{490,000}{24} \approx 20,416.67$$

The total of all recorded smallest values after 10,000 rounds would be approximately **20,417**.

## Question 10: Overfitting on Synthetic Data

**Q:** "Can the deep learning model above overfit (training error close to zero) on such data? What training-vs-validation pattern (under-fit, good fit, or over-fit) would you predict, and why?"

**A:** Yes, the model can achieve a training error close to zero. The predicted training-vs-validation pattern is a **good fit**, where both training and validation error decrease together and converge.

### Can the Model "Overfit" (Achieve Near-Zero Training Error)?

**Yes.** The model can achieve a training error very close to zero. Here's why:

- Sufficient Model Capacity:** The network is a standard Multi-Layer Perceptron (MLP) with two hidden layers of 64 neurons each and Leaky ReLU activations. This architecture is a universal function approximator, meaning it has enough complexity and non-linearity to learn the underlying mathematical relationship.
- Learnable Target Function:** The target  $y$  is generated by a smooth, continuous polynomial function ( $y = \sum (x_i^3 + 3x_i^2 + 3)$ ). There are no chaotic or discontinuous elements. A neural network is exceptionally good at approximating such functions.

Given that the model has the capacity and the function is learnable, the optimizer (Adam) will be able to adjust the network's weights to minimize the Mean Squared Error, driving the training error toward zero.

### Predicted Training vs. Validation Pattern

The pattern you would observe is a **good fit**. Both the training loss and the validation loss will

decrease together and converge at a low value.

#### Reasoning:

The most critical factor here is that the data is generated by a **deterministic, noise-free mathematical formula**.

- **No Noise to Overfit:** Classical overfitting occurs when a model is too complex for its dataset and starts to memorize the random noise in the training data instead of the underlying signal. When this "memorized noise" is absent from the validation set, the validation error increases while training error continues to decrease.
- **A Single Underlying Truth:** In this problem, there is no noise. There is only one true signal—the polynomial function. Since both the training and validation datasets are generated by this exact same function, any learning that improves performance on the training set will directly translate to improved performance on the validation set.

Therefore, as the model gets better at approximating the true polynomial function, its error on **any** data generated by that function (whether for training or validation) will decrease. The validation and training loss curves will track each other closely, which is the hallmark of a good fit.

## Question 11: Residual Modules and Positional Embeddings

Q1 (C.3): "Please explain in detail, why the Residual Module helps deep learning models converge better than original feature extraction function..."

Q2 (C.4): "Why don't we use the following simpler embeddings instead?  $p_i = i\$...$ "

#### A (C.3): Benefits of the Residual Module

A Residual Module helps deep learning models converge better by directly combating the **vanishing gradient problem** and making it easier for layers to learn an **identity mapping**.

#### The Problem with Standard Networks ( $F(X) = \text{feat}(X)$ )

In a very deep standard network, layers are stacked one after another. This creates two major issues:

1. **Vanishing Gradients:** During backpropagation, the gradient signal is repeatedly multiplied by the weights of each layer. In deep networks, this can cause the signal to shrink exponentially until it's too small to effectively train the initial layers of the network.
2. **Difficulty Learning Identity:** Sometimes, the optimal function for a block of layers is simply to pass its input through unchanged (an identity function, where  $F(X) = X$ ). For a standard block, the layers  $\text{feat}(X)$  must learn to meticulously adjust their weights to approximate this identity mapping, which is surprisingly difficult for complex, non-linear layers.

#### How Residual Modules Solve This ( $F(X) = X + \text{feat}(X)$ )

The residual module introduces a "skip connection" that adds the original input  $X$  to the output of the feature-extracting layers  $\text{feat}(X)$ .

This simple addition has two profound benefits:

1. **Combats Vanishing Gradients:** The skip connection creates an uninterrupted "superhighway" for the gradient to flow back through the network. The gradient of the output with respect to the input  $X$  is now  $\frac{\partial F(X)}{\partial X} = 1 + \frac{\partial \text{feat}(X)}{\partial X}$ . That crucial  $+1$  term ensures that even if the gradient through the feature-extracting layers becomes zero, a baseline gradient of 1 can always pass through the identity connection. This prevents the gradient signal from dying out and allows for the effective training of much deeper networks.
2. **Eases Identity Learning:** With the residual formulation, if an identity mapping is optimal ( $F(X) = X$ ), the network simply needs to make  $X + \text{feat}(X) = X$ . This is easily achieved by having the layers learn to make  $\text{feat}(X) = 0$ , which is much easier than learning an identity function. By driving its weights toward zero, the block can effectively be "skipped" if it's not needed, allowing the model to dynamically adjust its own depth without degrading performance.

In essence, the layers are no longer learning the entire output; they are only learning the **residual**, or the difference that needs to be added to the input.

#### A (C.4): Why Not Use Simpler Positional Embeddings?

Using a simple position index like  $p_i = i$  for positional embeddings in a Transformer is a bad idea because it leads to several major problems that the original paper's sinusoidal formula was specifically designed to avoid.

The Transformer's self-attention mechanism is inherently permutation-invariant—it has no concept of word order. Positional embeddings are added to give the model this crucial information.

Here's why  $p_i = i$  fails:

1. **Unbounded and Unstable Values:** The position index  $i$  increases linearly with the sequence length. For a long sequence, these positional values (e.g., 0, 1, 2, ..., 500) would have a completely different scale and magnitude than the word embedding values, which are typically small numbers. This can make the model unstable and hard to train, as the position information would dominate the word's meaning.
2. **Poor Generalization to Longer Sequences:** The model would only be trained on position indices up to the maximum length of sequences in the training data. If it encounters a longer sequence during testing (e.g., position 513 when it only trained up to 512), it has no idea what this new, unseen index means. The sinusoidal formula, by contrast, keeps all values bounded between -1 and 1, regardless of the position, and its periodic nature allows it to generalize to unseen lengths.
3. **Fails to Represent Relative Positions:** A model needs to understand not just the

absolute position of a word, but its relative position to other words (e.g., "the word 4 positions ahead"). The sinusoidal formula has the clever property that the embedding of position + k can be represented as a linear function of the embedding of position. This makes it very easy for the model to learn relative positioning. A simple index i does not have this property; the relationship between positions 5 and 9 is not easily comparable to the relationship between positions 105 and 109 from the model's perspective.

## Question 12: Tokenization Schemes

**Q: "(i) What is the ratio of compute cost between the two approaches? (ii) Which tokenization scheme is more likely to predict the next token correctly, and why?"**

**A:** Here are the answers to the questions based on the problem description.

### **(i) What is the ratio of compute cost between the two approaches?**

The character-level approach is approximately **35 to 36 times more computationally expensive** than the word-level approach.

#### **Reasoning:**

The compute cost of a Transformer is dominated by the self-attention mechanism, which has a quadratic complexity of  $\$O(n^2)$ , where n is the sequence length.

##### **1. Sequence Length Calculation:**

- **A. Character-level:** For a sentence with 20 words of 5 characters each, the number of character tokens is  $20 \text{ words} * 5 \text{ chars/word} = 100$ . Including the 19 spaces between words, the total sequence length ( $n_{\text{char}}$ ) is **119 tokens**.
- **B. Word-level:** The sequence length ( $n_{\text{word}}$ ) is simply the number of words, which is **20 tokens**.

##### **2. Compute Cost Ratio:**

We can find the ratio by comparing the squares of their sequence lengths:

$$\begin{aligned} \$\$ \text{Ratio} &= \frac{\text{Cost}_{\text{char}}}{\text{Cost}_{\text{word}}} \\ &= \frac{(n_{\text{char}})^2}{(n_{\text{word}})^2} = \frac{119^2}{20^2} = \frac{14,161}{400} \approx 35.4 \end{aligned}$$

### **(ii) Which tokenization scheme is more likely to predict the next token correctly, and why?**

The **word-level** tokenization scheme is more likely to predict the next token correctly.

#### **Reasoning:**

The core difference lies in the **semantic value** of a single token in each scheme.

##### **1. Word-Level (Higher Semantic Value):**

- Each token is a full word, which is a fundamental unit of meaning. When the model needs to predict the next token, it's making a prediction about the next logical

concept or idea.

- **Example:** Given the context "The capital of France is...", the model makes **one** prediction for the token "Paris". This prediction is based on a high-level understanding of the world. Because the model operates on meaningful semantic chunks, its predictions are more likely to be coherent and correct in a practical sense.

## 2. Character-Level (Lower Semantic Value):

- Each token is just a character, which carries very little semantic meaning on its own. To understand a single concept like "France", the model must process 6 separate tokens (F-r-a-n-c-e).
- **Example:** Given the same context, the model must predict the next token, P. Then, given that, it must predict a, and so on for r-i-s. While predicting a after P might be likely, correctly predicting the entire sequence of five characters is significantly harder and more prone to error than making a single, semantically-informed word prediction. The model's context is diluted across many low-information tokens, making it difficult to capture long-range dependencies and make high-level predictions.

## Question 13: GPT Model Comparison

**Q:** "You have a standard GPT style model... Which model is bigger regarding the attention stack? Why? Which will do faster inference / each forward pass?"

**A:** Here are the answers to the questions comparing the two GPT-style models.

### Which model is bigger regarding the attention stack?

Regarding the total number of parameters, **Model A and Model B are approximately the same size.**

### Why:

The number of parameters in a Transformer model is roughly proportional to the number of layers (L) multiplied by the square of the model dimension ( $d_{model}$ ), i.e., Parameters  $\propto L \times (d_{model})^2$ .

Let's compare the two models using this relationship:

- Model A (Shallow & Wide):  
$$L \times (d_{model})^2 = 8 \times (1024)^2 = 8 \times 1,048,576 = \mathbf{8,388,608}$$
- Model B (Deep & Narrow):  
$$L \times (d_{model})^2 = 32 \times (512)^2 = 32 \times 262,144 = \mathbf{8,388,608}$$

As the calculation shows, the two models have been designed to have a nearly identical

number of parameters. Model B compensates for its smaller model dimension by having four times as many layers.

### Which will do faster inference / each forward pass?

Model A (the shallower, wider model) will have faster inference. 

#### Why:

While the models have a similar total number of calculations (FLOPs), inference speed (latency) is primarily limited by **sequential operations**.

1. **Sequential Bottleneck:** A Transformer processes data layer by layer. You cannot begin calculating layer  $i+1$  until layer  $i$  is completely finished. Model B has **32 sequential steps** to complete, whereas Model A only has **8 sequential steps**.
2. **Hardware Parallelism:** Modern hardware like GPUs are highly parallel. They can handle the larger matrix multiplications within Model A's wider layers (e.g., involving  $1024 \times 1024$  matrices) very efficiently. The bottleneck is not the size of the matrices, but the number of times you have to wait for a layer to finish before starting the next one.

In short, the deep and narrow architecture of Model B creates a significant sequential dependency that leads to higher latency, making the shallow and wide Model A faster for a single forward pass.

### Question 14: Transformer Instability (NaNs)

**Q:** "You inherit a PyTorch implementation of a 6-layer Transformer encoder where training suddenly diverges with NaNs... Explain why the current implementation causes instability, and show how to correct the code."

**A:** The current implementation causes instability because it uses **Post-Layer Normalization**, which can lead to exploding activation values and gradients in a deep network.

#### Why the Current Implementation is Unstable

The provided code implements the original Transformer architecture where layer normalization is applied *after* the residual connection ( $x = \text{self.norm1}(x + \text{attn\_out})$ ). This is known as **Post-LN**.

The problem with Post-LN is that the output from one block is passed directly into the next before being normalized. In a deep stack of layers (even 6 is enough), this can cause the magnitude of the values in the tensor  $x$  to grow uncontrollably layer after layer.

These exploding values are then fed into the self-attention and feed-forward networks. Operations within these modules, especially the softmax in attention, are highly sensitive to large inputs and can easily result in floating-point overflow, producing inf values. Subsequent calculations involving these inf values lead to NaNs (Not a Number), which corrupt the loss

and cause the training to diverge.

### How to Correct the Code

The solution is to switch the order of operations to **Pre-Layer Normalization (Pre-LN)**. In this scheme, the normalization is applied to the input *before* it enters the sub-layer (attention or feed-forward). This ensures that the inputs to every module are always well-behaved and stable.

Here is the corrected code:

```
def forward(self, x):
    # Multi-head self-attention with Pre-LN
    # 1. Normalize the input FIRST
    norm_x = self.norm1(x)
    # 2. Pass the normalized input to the attention layer
    attn_out = self.self_attn(norm_x, norm_x, norm_x)
    # 3. Add the output to the original, un-normalized input
    x = x + attn_out

    # Feed-forward with Pre-LN
    # 1. Normalize the input FIRST
    norm_x = self.norm2(x)
    # 2. Pass the normalized input to the feed-forward layer
    ff_out = self.ffn(norm_x)
    # 3. Add the output to the original, un-normalized input
    x = x + ff_out

return x
```

This Pre-LN architecture stabilizes the training by keeping the activations under control throughout the forward pass and provides a cleaner, more direct path for gradients during backpropagation, often eliminating the need for aggressive learning rate warmup schedules.

### Question 15: Multi-Head vs. Single-Head Attention

**Q: "The second question?" (Referring to Problem C.8, on single vs multi-head attention)**

**A:** Of course. Here are three tradeoffs you would expect from switching from a multi-head ( $H=8$ ) to a single-head ( $H=1$ ) attention mechanism.

#### 1. Tradeoff: Representational Diversity vs. Simplicity

- **Multi-Head ( $H=8$ ):** The key benefit of multiple heads is that each head can learn to focus

on different types of relationships in different representation subspaces. For example, one head might learn syntactic dependencies (like subject-verb agreement), while another learns semantic similarity. This allows the model to capture a rich and diverse set of features from the input simultaneously.

- **Single-Head (H=1):** With only one head, the model is forced to average all these different types of desired relationships into a single attention pattern. It loses the ability to look at the text from multiple, independent perspectives.

**The tradeoff:** You sacrifice the model's **expressiveness and ability to learn diverse representations** for a simpler, less nuanced attention mechanism. The single-head model will likely have worse performance on complex language tasks.

## 2. Tradeoff: Ensemble Robustness vs. A Single Mechanism

- **Multi-Head (H=8):** The multiple heads act like an ensemble. If one head produces a poor or nonsensical attention pattern for a particular token, the other seven heads can compensate, making the overall model more robust and stable.
- **Single-Head (H=1):** The single attention head has no backup. It must learn a single, "one-size-fits-all" attention distribution. This can lead to a diluted or averaged pattern that is less focused and less effective than specialized heads.

**The tradeoff:** You exchange the **robustness of an ensemble of attention heads** for a single mechanism that might be less effective because it has to average competing objectives.

## 3. Tradeoff: Computational Parallelism

- **Multi-Head (H=8):** The computation for each of the 8 heads is independent and can be executed in perfect parallel. The model performs 8 smaller attention calculations (on vectors of dimension 64) at the same time. This structure is extremely well-suited for modern hardware like GPUs, which excel at running many parallel tasks.
- **Single-Head (H=1):** The model performs one single, large attention calculation (on vectors of dimension 512).

**The tradeoff:** While the total number of floating-point operations (FLOPs) is the same in both scenarios, you trade a **highly parallel computational structure** for a single, monolithic one. The multi-head approach often leads to better hardware utilization and can be more efficient in practice.

## Question 16: Grid Partitioning Problem

**Q:** "On a 4x4 grid, we want to remove four 1x1 squares so that no two are on the same row or column... find the one that requires the least number of rectangles to partition the grid..."

**A:** The minimum number of rectangles required is **5**.

## Reasoning

1. **The Goal:** We need to remove four  $1 \times 1$  squares from a  $4 \times 4$  grid—with no two removed from the same row or column—in such a way that the remaining 12 squares can be covered by the fewest possible rectangles. The cost of a configuration is the minimum number of non-overlapping rectangles needed to tile the remaining area.
2. **Exploring Configurations:** The placement of the removed squares determines how "fragmented" the remaining shape is. A more fragmented shape requires more rectangles to cover.
  - o A configuration like removing the diagonal squares— $(1,1)$ ,  $(2,2)$ ,  $(3,3)$ ,  $(4,4)$ —is a poor choice. It splits the board into two disconnected "staircase" shapes, each of which requires 3 rectangles to tile, for a total of **6 rectangles**.
  - o Many other configurations also lead to a minimum of 6 rectangles.
3. **The Optimal Configuration:** To minimize the number of rectangles, we need to arrange the removed squares to leave the remaining area as "blocky" as possible. The configuration that achieves the minimum is removing the squares at positions  **$(1,2)$ ,  $(2,4)$ ,  $(3,1)$ , and  $(4,3)$** .

## Drawing of the Optimal Configuration

Below is a drawing of the grid with the specified squares removed (marked with an X). The remaining area is partitioned into 5 rectangles, proving that a 5-rectangle tiling is possible.

The 5-rectangle partition is as follows:

- **Rectangle 1 (Blue):** A  $2 \times 2$  square in the middle  $\{(2,2), (2,3), (3,2), (3,3)\}$ .
- **Rectangle 2 (Green):** A  $2 \times 1$  column on the left  $\{(1,1), (2,1)\}$ .
- **Rectangle 3 (Yellow):** A  $1 \times 2$  row on the top  $\{(1,3), (1,4)\}$ .
- **Rectangle 4 (Purple):** A  $1 \times 2$  row on the bottom  $\{(4,1), (4,2)\}$ .
- **Rectangle 5 (Red):** A  $2 \times 1$  column on the right  $\{(3,4), (4,4)\}$ .

Since we have found a valid configuration that requires 5 rectangles, and it's known to be the minimum possible, the correct answer is **5**.