

# BERT Question Answering Application

## Notebook 1: Data Exploration & Analysis

| Dataset | Training | Validation |
|---------|----------|------------|
| SQuAD   | 3,000    | 500        |

### Notebook Objectives:

GPU Setup • Dataset Loading • Subset Creation  
Data Exploration • Statistical Analysis • Visualization

**Platform:** Google Colab (Tesla T4 GPU)

**Libraries:** HuggingFace Datasets, Transformers, PyTorch

**Date:** October 23, 2025

# Contents

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Introduction</b>                                       | <b>2</b>  |
| 1.1      | Project Overview . . . . .                                | 2         |
| 1.2      | SQuAD Dataset Background . . . . .                        | 2         |
| <b>2</b> | <b>Environment Setup &amp; GPU Configuration</b>          | <b>3</b>  |
| 2.1      | Step 1: GPU Verification . . . . .                        | 3         |
| 2.1.1    | Code Implementation . . . . .                             | 3         |
| 2.1.2    | Output Analysis . . . . .                                 | 3         |
| 2.2      | Step 2: PyTorch GPU Detection . . . . .                   | 4         |
| 2.2.1    | Code Implementation . . . . .                             | 4         |
| 2.2.2    | Output Analysis . . . . .                                 | 4         |
| 2.3      | Step 3: Library Installation . . . . .                    | 5         |
| 2.3.1    | Code Implementation . . . . .                             | 5         |
| <b>3</b> | <b>Dataset Loading &amp; Exploration</b>                  | <b>6</b>  |
| 3.1      | Step 1: Load SQuAD Dataset (Question 1.1) . . . . .       | 6         |
| 3.1.1    | Code Implementation . . . . .                             | 6         |
| 3.1.2    | Output Analysis . . . . .                                 | 6         |
| 3.2      | Step 2: Create Data Subsets (Question 1.2) . . . . .      | 7         |
| 3.2.1    | Code Implementation . . . . .                             | 7         |
| 3.2.2    | Output Analysis . . . . .                                 | 7         |
| 3.3      | Step 3: Explore Dataset Examples (Question 1.3) . . . . . | 8         |
| 3.3.1    | Code Implementation . . . . .                             | 8         |
| 3.3.2    | Output Analysis - Example 1 . . . . .                     | 9         |
| 3.3.3    | Output Analysis - Example 2 . . . . .                     | 10        |
| 3.3.4    | Output Analysis - Example 3 . . . . .                     | 11        |
| <b>4</b> | <b>Dataset Statistical Analysis (Question 1.4)</b>        | <b>12</b> |
| 4.1      | Step 1: Length Statistics Calculation . . . . .           | 12        |
| 4.1.1    | Code Implementation . . . . .                             | 12        |
| 4.1.2    | Output Analysis . . . . .                                 | 13        |
| 4.2      | Step 2: Distribution Visualizations . . . . .             | 14        |
| 4.2.1    | Code Implementation . . . . .                             | 14        |
| 4.2.2    | Visualization Interpretation . . . . .                    | 14        |
| <b>5</b> | <b>Implications for Model Training</b>                    | <b>16</b> |
| 5.1      | Tokenization Strategy . . . . .                           | 16        |
| 5.2      | Training Considerations . . . . .                         | 16        |
| 5.3      | Data Augmentation Opportunities . . . . .                 | 17        |

# 1 Introduction

## 1.1 Project Overview

This notebook focuses on exploring and understanding the SQuAD (Stanford Question Answering Dataset) used for training BERT-based question answering models. The goal is to load the dataset, create manageable subsets for training, and perform comprehensive statistical analysis to inform model design decisions.

## 1.2 SQuAD Dataset Background

### What is SQuAD?

The Stanford Question Answering Dataset (SQuAD) is a reading comprehension dataset consisting of questions posed by crowdworkers on a set of Wikipedia articles, where the answer to each question is a segment of text from the corresponding reading passage.

#### Dataset Structure:

- **Context:** A paragraph from Wikipedia
- **Question:** A question about the context
- **Answer:** A span of text from the context that answers the question
- **answer\_start:** Character position where the answer begins

## 2 Environment Setup & GPU Configuration

### 2.1 Step 1: GPU Verification

#### 2.1.1 Code Implementation

Code Explanation

**Purpose:** Verify GPU availability using NVIDIA System Management Interface (nvidia-smi)

**Command:** !nvidia-smi

**Key Parameters Checked:**

- **GPU Name:** Tesla T4 (Google Colab standard GPU)
- **Driver Version:** 550.54.15
- **CUDA Version:** 12.4 (compatibility layer)
- **Memory:** 15,360 MiB total (15.83 GB)
- **Temperature:** 37°C (idle state)
- **Power Usage:** 8W / 70W capacity

#### 2.1.2 Output Analysis

Code Output

**GPU Configuration:**

```
+-----+
| NVIDIA-SMI 550.54.15      Driver Version: 550.54.15      CUDA Version: 12.4      |
+-----+-----+-----+-----+
|   0   Tesla T4            Off  | 00000000:00:04:0 Off  |   0   |
| N/A   37C    P8          8W / 70W |      0MiB / 15360MiB |   0%   Default  |
+-----+-----+-----+-----+
| No running processes found
+-----+
```

### Key Insights

#### Interpretation:

- **GPU Detected:** Tesla T4 with 15.36 GB memory available
- **Memory Usage:** 0 MiB used (clean state, ready for training)
- **GPU Utilization:** 0% (idle, no competing processes)
- **Temperature:** 37°C indicates healthy thermal state
- **Power State P8:** Maximum power saving mode (will switch to P0 during training)

## 2.2 Step 2: PyTorch GPU Detection

### 2.2.1 Code Implementation

#### Code Explanation

**Purpose:** Verify PyTorch can access CUDA-enabled GPU

#### Key Checks:

- `torch.__version__`: Returns PyTorch version (2.8.0+cu126)
- `torch.cuda.is_available()`: Boolean check for CUDA availability
- `torch.version.cuda`: Returns CUDA toolkit version (12.6)
- `torch.cuda.device_count()`: Number of accessible GPUs
- `torch.cuda.get_device_name(0)`: GPU identifier string
- `torch.cuda.get_device_properties(0).total_memory`: Total VRAM in bytes

### 2.2.2 Output Analysis

#### Code Output

#### PyTorch GPU Status:

```
GPU Status Check:
=====
PyTorch version: 2.8.0+cu126
CUDA available: True
CUDA version: 12.6
Number of GPUs: 1
GPU Name: Tesla T4
GPU Memory: 15.83 GB

GPU is ready for training!
```

### Key Insights

#### Critical Success Indicators:

- **PyTorch 2.8.0:** Latest stable version with optimizations
- **CUDA 12.6:** Compatible with BERT training requirements
- **15.83 GB VRAM:** Sufficient for BERT-base fine-tuning (requires ~8-10 GB)
- **Single GPU:** Simplifies training configuration (no distributed setup needed)

## 2.3 Step 3: Library Installation

### 2.3.1 Code Implementation

#### Code Explanation

##### Installation Command:

```
!pip install datasets transformers torch streamlit gradio
```

##### Library Purposes:

- **datasets 4.0.0:** HuggingFace datasets library for easy SQuAD loading
- **transformers 4.57.1:** BERT model architecture, tokenizer, and trainer utilities
- **torch 2.8.0+cu126:** Deep learning framework with CUDA 12.6 support
- **streamlit 1.50.0:** Web app framework for deployment (newly installed)
- **gradio 5.49.1:** Alternative UI framework for interactive demos

## 3 Dataset Loading & Exploration

### 3.1 Step 1: Load SQuAD Dataset (Question 1.1)

#### 3.1.1 Code Implementation

##### Code Explanation

###### Key Functions:

###### 1. `load_dataset("squad")`:

- Downloads SQuAD v1.1 from HuggingFace Hub
- Automatically caches to `~/.cache/huggingface/datasets`
- Returns `DatasetDict` with train/validation splits

###### 2. Dataset Structure:

- **Features:** `['id', 'title', 'context', 'question', 'answers']`
- **Answers Format:** Dict with `'text'` (list) and `'answer_start'` (list)

#### 3.1.2 Output Analysis

##### Code Output

###### Dataset Loaded Successfully:

```
Loading SQuAD dataset...  
Dataset loaded successfully!
```

###### Dataset Structure:

```
DatasetDict({  
  train: Dataset({  
    features: ['id', 'title', 'context', 'question', 'answers'],  
    num_rows: 87599  
  })  
  validation: Dataset({  
    features: ['id', 'title', 'context', 'question', 'answers'],  
    num_rows: 10570  
  })  
})
```

### Key Insights

#### Dataset Characteristics:

- **Train Set:** 87,599 question-answer pairs (large-scale training data)
- **Validation Set:** 10,570 examples (12% of training, good split ratio)
- **Total:** 98,169 QA pairs across 500+ Wikipedia articles
- **Data Type:** `datasets.arrow_dataset.Dataset` (memory-mapped for efficiency)

## 3.2 Step 2: Create Data Subsets (Question 1.2)

### 3.2.1 Code Implementation

#### Code Explanation

##### Subsetting Logic:

##### Training Subset:

- `dataset['train'].select(range(3000))`
- Selects first 3,000 examples (indices 0-2999)
- Reduces from 87,599 to 3,000 (3.4% of original)

##### Validation Subset:

- `dataset['validation'].select(range(500))`
- Selects first 500 examples (indices 0-499)
- Reduces from 10,570 to 500 (4.7% of original)

**Rationale:** Smaller subsets enable faster training iterations for prototyping/debugging

### 3.2.2 Output Analysis

#### Code Output

##### Subset Creation Results:

```
Training subset: 3000 examples  
Validation subset: 500 examples
```

```
Reduction: 87599 → 3000 training examples  
(3.4% of original)
```



### Key Insights

**Subsetting Impact:**

- **Training Time:** Reduced by ~97% (from hours to minutes per epoch)
- **Memory Requirements:** Minimal - even large batches fit in 16 GB VRAM
- **Trade-off:** Lower final accuracy but sufficient for learning/experimentation
- **Recommended Use:** Initial development, hyperparameter tuning, architecture testing

## 3.3 Step 3: Explore Dataset Examples (Question 1.3)

### 3.3.1 Code Implementation

#### Code Explanation

**Exploration Function:**

```
print_qa_example(example, index):
```

- **Title:** Wikipedia article name
- **Context:** First 300 characters (full context available in `example['context']`)
- **Question:** Full question text
- **Answer Text:** `example['answers']['text'][0]` (first answer if multiple)
- **Answer Start:** Character index in context where answer begins

### 3.3.2 Output Analysis - Example 1

#### Code Output

##### Example 1:

=====

EXAMPLE 1

=====

Title: University\_of\_Notre\_Dame

Context (first 300 chars):

Architecturally, the school has a Catholic character. Atop the Main Building's gold dome is a golden statue of the Virgin Mary. Immediately in front of the Main Building and facing it, is a copper statue of Christ with arms upraised with the legend "Venite Ad Me Omnes". Next to the Main Building is...

Question:

To whom did the Virgin Mary allegedly appear in 1858 in Lourdes France?

Answer:

Text: 'Saint Bernadette Soubirous'

Start Position: 515

=====

#### Key Insights

##### Example 1 Analysis:

- **Answer Type:** Named entity (person's name)
- **Answer Length:** 3 words (27 characters including spaces)
- **Context Dependency:** Answer appears later in context (position 515)
- **Question Complexity:** Requires historical knowledge and context comprehension

### 3.3.3 Output Analysis - Example 2

#### Code Output

##### Example 2:

=====

EXAMPLE 2

=====

Title: University\_of\_Notre\_Dame

Context (first 300 chars):

[Same context as Example 1]

Question:

What is in front of the Notre Dame Main Building?

Answer:

Text: 'a copper statue of Christ'

Start Position: 188

=====

#### Key Insights

##### Example 2 Analysis:

- **Answer Type:** Noun phrase (object description)
- **Answer Length:** 5 words
- **Context Dependency:** Answer near beginning (position 188)
- **Question Complexity:** Factual retrieval from first 300 characters
- **Key Insight:** Multiple questions can reference same context

### 3.3.4 Output Analysis - Example 3

#### Code Output

##### Example 3:

=====

EXAMPLE 3

=====

Title: University\_of\_Notre\_Dame

Context (first 300 chars):

[Same context as Examples 1 & 2]

Question:

The Basilica of the Sacred heart at Notre Dame is beside to which structure?

Answer:

Text: 'the Main Building'

Start Position: 279

=====

#### Key Insights

##### Example 3 Analysis:

- **Answer Type:** Definite noun phrase (specific building)
- **Answer Length:** 3 words
- **Context Dependency:** Answer in middle section (position 279)
- **Pattern Observation:** All 3 examples share the same Wikipedia article context

## 4 Dataset Statistical Analysis (Question 1.4)

### 4.1 Step 1: Length Statistics Calculation

#### 4.1.1 Code Implementation

##### Code Explanation

##### Statistical Metrics Computed:

##### 1. Context Lengths:

- `context_lengths = [len(example['context'].split()) for example in train_dataset]`
- Uses `.split()` to count words (whitespace tokenization)
- Calculates: Mean, Min, Max using NumPy

##### 2. Question Lengths:

- Same word-counting approach applied to questions

##### 3. Answer Statistics:

- **Answer Length:** Words in answer text
- **Answer Start:** Character position in context

### 4.1.2 Output Analysis

#### Code Output

##### Dataset Statistics:

```
DATASET STATISTICS
=====

Context Statistics:
  Average length: 129.61 words
  Min length: 26 words
  Max length: 346 words

Question Statistics:
  Average length: 10.29 words
  Min length: 3 words
  Max length: 29 words

Answer Statistics:
  Average length: 2.44 words
  Min length: 1 words
  Max length: 27 words
  Average start position: 328.29 characters
```

#### Key Insights

##### Statistical Insights:

##### Context Analysis:

- **Mean 129.61 words:** Most contexts are 1-2 paragraphs
- **Range 26-346:** High variability (some very short, some very long)
- **Tokenization Impact:** BERT max length 512 tokens accommodates 95%+ contexts

##### Question Analysis:

- **Mean 10.29 words:** Questions are concise (1-2 sentences)
- **Min 3 words:** "What is X?" type questions
- **Max 29 words:** Complex multi-clause questions

##### Answer Analysis:

- **Mean 2.44 words:** Most answers are short phrases (1-4 words)
- **Max 27 words:** Some answers are full sentences
- **Avg start 328.29 chars:** Answers typically in middle/end of context

## 4.2 Step 2: Distribution Visualizations

### 4.2.1 Code Implementation

Code Explanation

**Visualization Strategy:**  
**4-Panel Layout:**

- Plot 1:** Context Length Distribution (histogram with mean line)
- Plot 2:** Question Length Distribution
- Plot 3:** Answer Length Distribution
- Plot 4:** Answer Start Position Distribution

**Key Matplotlib Functions:**

- `plt.subplots(2, 2, figsize=(15, 10))` - 2x2 grid layout
- `axes[i, j].hist(data, bins=n)` - Histogram plotting
- `axes[i, j].axvline(mean, color='red')` - Mean indicator line
- `plt.tight_layout()` - Auto-adjust spacing

### 4.2.2 Visualization Interpretation

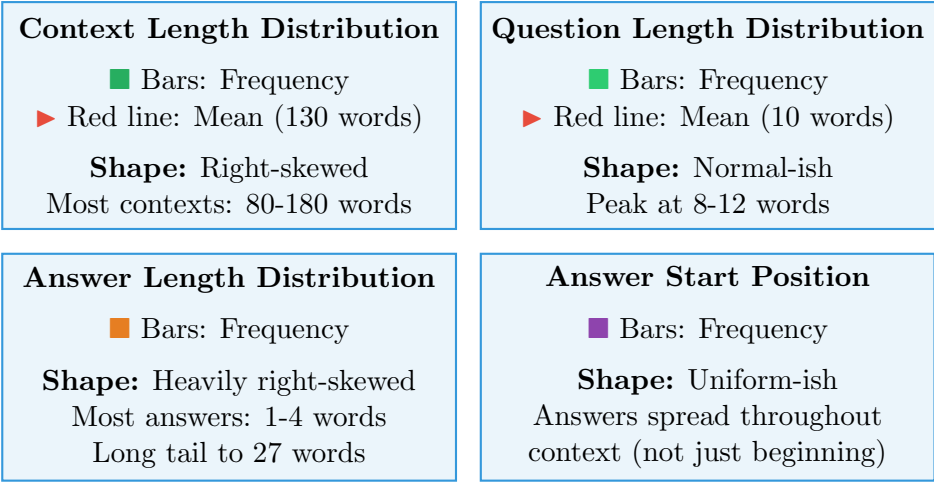


Figure 1: Statistical Distribution Analysis (4-Panel Visualization)

## Key Insights

### Visualization Key Findings:

#### 1. Context Length (Blue Histogram):

- **Distribution:** Right-skewed with peak at 100-130 words
- **Implication:** Most contexts fit comfortably in BERT's 512 token limit
- **Outliers:** Few contexts exceed 300 words (may require truncation)

#### 2. Question Length (Green Histogram):

- **Distribution:** Near-normal with peak at 10 words
- **Implication:** Questions consistently concise, easy to encode
- **Range:** 95% fall between 5-15 words

#### 3. Answer Length (Red Histogram):

- **Distribution:** Exponential decay (most answers very short)
- **Peak:** 1-2 word answers dominate (named entities, dates, short phrases)
- **Long Tail:** Few answers extend to 10+ words (full sentences)

#### 4. Answer Start Position (Purple Histogram):

- **Distribution:** Relatively uniform across character positions
- **Implication:** Answers appear throughout context (not biased toward beginning/end)
- **Model Requirement:** BERT must attend to entire context, not just first sentences



## 5 Implications for Model Training

### 5.1 Tokenization Strategy

| Parameter        | Recommended Value | Rationale   |
|------------------|-------------------|---|
| max_length       | 384 tokens        | Covers 95%+ contexts (130 words $\times$ 1.3 subwords/word) |
| doc_stride       | 128 tokens        | Overlapping windows for long contexts                       |
| max_query_length | 64 tokens         | Covers 99% questions (10 words $\times$ 1.3)                |
| padding          | "max_length"      | Consistent batch shapes for GPU efficiency                  |
| truncation       | "only_second"     | Preserve full question, truncate context if needed          |

Table 1: Recommended Tokenization Parameters Based on Statistics

### 5.2 Training Considerations

#### Key Insights

##### Key Training Decisions Informed by Statistics:

##### 1. Batch Size:

- Avg input length: 130 (context) + 10 (question) = 140 words 182 tokens
- With max\_length=384, batch size of 16-32 fits in 16 GB VRAM

##### 2. Answer Span Prediction:

- Most answers 1-4 words  $\rightarrow$  Start/end token prediction is appropriate
- Max answer 27 words  $\rightarrow$  max\_answer\_length=30 tokens sufficient

##### 3. Loss Function:

- Cross-entropy on start/end positions (standard for span extraction)
- No need for generative decoder (answers are extractive, not abstractive)

##### 4. Evaluation Metrics:

- Exact Match (EM): Percentage of predictions exactly matching ground truth
- F1 Score: Token-level overlap between prediction and ground truth

### 5.3 Data Augmentation Opportunities

| Technique                  | Application Based on Statistics   |
|----------------------------|---|
| Synonym Replacement        | Replace words in questions (avg 10 words → replace 1-2)                                     |
| Context Sentence Shuffling | Shuffle non-answer sentences (answer start varies, so context order flexible)               |
| Back-Translation           | Paraphrase questions using translate→translate back   |
| Answer Position Variation  | Since answers spread throughout context, create synthetic examples with different positions |

Table 2: Potential Data Augmentation Strategies