# AI Tools Assignment Report

## Part 1: Theoretical Understanding

### Q1: TensorFlow vs PyTorch

## Q1: Explain the primary differences between TensorFlow and PyTorch. When would you choose one over the other?

**Primary Differences:**

| Aspect | TensorFlow | PyTorch |
|---|---|---|
| **Computation Graph** | Static (TF 1.x) / Dynamic (TF 2.x with eager execution) | Dynamic (define-by-run) |
| **Debugging** | More challenging in TF 1.x, improved in TF 2.x | Easier with Python debugging tools |
| **Deployment** | Superior production tools (TF Serving, TF Lite, TF.js) | Growing ecosystem (TorchServe, ONNX) |
| **API Design** | Higher-level Keras API, more abstraction | More Pythonic, intuitive |
| **Community** | Strong industry adoption | Popular in research |
| **Learning Curve** | Steeper (historically), improved with TF 2.x | Gentler, more intuitive |

**When to Choose TensorFlow:**

- **Production Deployment**: When you need robust deployment pipelines across multiple platforms (mobile, web, embedded systems)
- **Industry Projects**: Large-scale applications requiring TensorFlow Extended (TFX) for ML pipelines
- **Mobile/Edge Computing**: TensorFlow Lite provides excellent mobile optimization
- **Established Ecosystems**: When working with existing TensorFlow codebases
- **TPU Usage**: Google's TPUs are optimized for TensorFlow

**When to Choose PyTorch:**

- **Research & Prototyping**: Rapid experimentation with dynamic architectures
- **Academic Projects**: Better documentation and research paper implementations
- **Custom Models**: Complex architectures requiring fine-grained control
- **Learning Deep Learning**: More intuitive for beginners due to Pythonic nature

- **Computer Vision/NLP Research**: Strong community support with libraries like torchvision and Hugging Face Transformers

**Practical Example:**

```python
# TensorFlow is better for:
# - Mobile app image classification
# - IoT device deployment
# - Production ML pipelines with TFX

# PyTorch is better for:
# - Novel GAN architectures
# - Cutting-edge transformer models

# - Academic research papers
```

### Q2: Jupyter Notebooks Use Cases

# Q2: Describe two use cases for Jupyter Notebooks in AI development.

## Use Case 1: Exploratory Data Analysis (EDA) and Experimentation

Jupyter Notebooks excel at interactive data exploration, allowing data scientists to:

- **Visualize data distributions** in real-time with inline plotting
- **Test hypotheses quickly** by running code cells independently
- **Document insights** alongside code using Markdown cells
- **Iterate rapidly** without re-running entire scripts

**Example Scenario:** A data scientist analyzing customer churn data can:

1. Load and inspect the dataset
2. Create visualizations (histograms, correlation matrices)
3. Try different feature engineering approaches
4. Test multiple models interactively
5. Document findings with explanations and charts

**Benefits:**

- Immediate visual feedback
- Non-linear workflow (run cells out of order)
- Easy sharing with stakeholders
- Combines code, output, and narrative

---

**Use Case 2: Model Training, Evaluation, and Reporting**

Jupyter Notebooks serve as comprehensive ML experiment logs:

- **Track model performance** across different hyperparameters
- **Generate visual reports** with confusion matrices, ROC curves
- **Share reproducible results** with team members
- **Create presentation-ready outputs** for stakeholders

**Example Scenario:** Training a neural network for image classification:

1. Document model architecture with explanations
2. Visualize training/validation loss curves in real-time
3. Display sample predictions with actual images
4. Calculate and present evaluation metrics
5. Export notebook as PDF/HTML for reports

**Benefits:**

- Self-contained experiment documentation
- Version control friendly (with nbdiff tools)
- Easy collaboration via GitHub
- Can be converted to slides for presentations

### Q3: spaCy vs Basic Python

**Q3: How does spaCy enhance NLP tasks compared to basic Python string operations?**

**Key Enhancements:**

**1. Linguistic Intelligence:**

- **Basic Python:** String splitting treats text as character sequences

python

```
text.split()  # Simple tokenization by whitespace
```

- **spaCy:** Understands language structure

python

```
doc = nlp("Dr. Smith earned $1.5M in 2023.")
```

*# Correctly tokenizes "Dr.", "$1.5M" as single units*

## 2. Named Entity Recognition (NER):

- **Basic Python:** Cannot identify entities without complex regex patterns
- **spaCy:** Pre-trained models recognize people, organizations, locations, dates, money

python
```python
for ent in doc.ents:
    print(ent.text, ent.label_)  # "Dr. Smith" → PERSON
```

## 3. Part-of-Speech (POS) Tagging:

- **Basic Python:** No grammatical understanding
- **spaCy:** Identifies nouns, verbs, adjectives

python
```python
for token in doc:
    print(token.text, token.pos_)  # "running" → VERB
```

## 4. Dependency Parsing:

- **Basic Python:** Cannot understand relationships between words
- **spaCy:** Maps syntactic relationships

python
```python
for token in doc:
    print(token.text, token.dep_, token.head.text)
    # Shows subject-verb-object relationships
```

## 5. Lemmatization:

- **Basic Python:** No word normalization
- **spaCy:** Reduces words to base forms intelligently

python
```python
"running" → "run", "better" → "good"
```

## 6. Performance:

- **Basic Python:** Slow for complex text processing
- **spaCy:** Optimized Cython implementation (up to 100x faster)

**Practical Comparison:**

```python
# Basic Python approach (limited):
text = "Apple Inc. is buying UK startup for $1 billion"
words = text.split()
# Result: ['Apple', 'Inc.', 'is', 'buying', 'UK', 'startup', 'for', '$1', 'billion']
# Cannot identify: Apple Inc. (ORG), UK (GPE), $1 billion (MONEY)

# spaCy approach (intelligent):
doc = nlp(text)
for ent in doc.ents:
    print(ent.text, ent.label_)
# Result:
# Apple Inc. → ORG
# UK → GPE

# $1 billion → MONEY
```

**Use Cases Where spaCy Shines:**

- Information extraction from documents
- Chatbot intent recognition
- Resume parsing
- Medical text analysis
- Legal document processing
- Social media sentiment analysis

### Comparative Analysis

# 2. Comparative Analysis

## Scikit-learn vs. TensorFlow

| Criteria | Scikit-learn | TensorFlow |
| --- | --- | --- |
| **Target Applications** | Classical ML (SVM, Random Forests, Linear Models, Clustering, Dimensionality Reduction) | Deep Learning (Neural Networks, CNNs, RNNs, Transformers), Large-scale ML |
| **Problem Types** | Structured/tabular data, Small-to-medium datasets, Traditional ML algorithms | Unstructured data (images, text, audio), Large datasets, Complex pattern recognition |

| | | |
|---|---|---|
| **Best For** | Classification, regression, clustering on tabular data | Computer vision, NLP, speech recognition, time series |

| Criteria | Scikit-learn | TensorFlow |
|---|---|---|
| **Ease of Use for Beginners** | ⭐⭐⭐⭐⭐ Excellent | ⭐⭐⭐ Moderate |
| **API Design** | Consistent, simple API (fit/predict pattern) | More complex, requires understanding of computational graphs |
| **Learning Curve** | Gentle - can start in hours | Steeper - requires days/weeks to master |
| **Code Example Complexity** | 5-10 lines for basic model | 20-50 lines for basic neural network |
| **Documentation** | Clear, beginner-friendly | Comprehensive but overwhelming for beginners |
| **Setup** | `pip install scikit-learn` (lightweight) | Larger installation, GPU configuration optional |

**Example Code Comparison:**

```python
# Scikit-learn (Simple)
from sklearn.tree import DecisionTreeClassifier
model = DecisionTreeClassifier()
model.fit(X_train, y_train)
predictions = model.predict(X_test)

# TensorFlow (More Complex)
import tensorflow as tf
model = tf.keras.Sequential([
    tf.keras.layers.Dense(128, activation='relu'),
    tf.keras.layers.Dense(10, activation='softmax')
])
model.compile(optimizer='adam', loss='sparse_categorical_crossentropy')
model.fit(X_train, y_train, epochs=10, batch_size=32)

predictions = model.predict(X_test)
```

| Criteria | Scikit-learn | TensorFlow |
|---|---|---|
| **Community Support** | ⭐⭐⭐⭐ Strong | ⭐⭐⭐⭐⭐ Excellent |
| **GitHub Stars** | ~59k | ~185k |
| **Stack Overflow Questions** | ~80k questions | ~150k questions |
| **Industry Adoption** | Widespread for classical ML | Dominant in deep learning |
| **Research Papers** | Moderate citations | Extensive citations in AI research |
| **Tutorials & Courses** | Abundant, beginner-friendly | Massive ecosystem, all skill levels |
| **Active Development** | Stable, incremental updates | Rapid evolution, frequent releases |
| **Corporate Backing** | Community-driven (originally by INRIA) | Google-backed |

**When to Use Which:**

**Use Scikit-learn when:**

- Working with tabular data (CSV files, databases)
- Dataset fits in memory (<1GB typically)
- Need quick prototypes
- Classical ML algorithms suffice
- Interpretability is crucial
- Example: Predicting house prices, customer segmentation

**Use TensorFlow when:**

- Working with images, text, or audio
- Large datasets (>1GB)
- Need deep neural networks
- GPU acceleration required
- Production deployment at scale
- Example: Facial recognition, language translation

**Hybrid Approach:** Many real-world projects use both:

1. Scikit-learn for preprocessing and feature engineering
2. TensorFlow for deep learning models

3. Scikit-learn for final ensemble or calibration

## Part 2: Practical Implementation

### Task 1: Iris Classification
**Approach:**
[Explain your preprocessing steps]

**Results:**
- Accuracy: [Your result]
- Precision: [Your result]
- Recall: [Your result]

[Insert screenshots]

**Code Explanation:**
[Key code snippets with explanations]

### Task 2: MNIST CNN
**Model Architecture:**
[Describe your CNN layers]

**Training Process:**
[Explain epochs, callbacks, optimization]

**Results:**
- Test Accuracy: [Your result] ✓ >95%
- Loss: [Your result]

[Insert all visualizations]

**Key Achievements:**
- Successfully achieved target accuracy
- Implemented batch normalization
- Used data augmentation

### Task 3: NLP Analysis
**Named Entity Recognition:**
- Total entities: [Your count]
- Top brands: [List]

[Insert entity results]

**Sentiment Analysis:**
- Positive: [Percentage]
- Negative: [Percentage]
- Neutral: [Percentage]

[Insert sentiment visualizations]

## Part 3: Ethics & Optimization

### Bias Identification

# Part 3: Ethics & Optimization

## 1. Ethical Considerations

**Bias Identification in MNIST Model**

**Potential Biases:**

**1. Dataset Representation Bias**

- **Issue:** The MNIST dataset consists primarily of handwritten digits from American census forms and high school students
- **Impact:** The model may perform poorly on:
    - Different handwriting styles from other cultures
    - Digits written with different writing instruments
    - Stylized or artistic number representations
    - Digits written by people with motor disabilities
- **Real-world consequence:** A banking app using this model might fail to recognize check amounts written by elderly users or people with Parkinson's disease

**2. Class Imbalance (Minimal but exists)**

- **Issue:** While MNIST is relatively balanced, slight variations exist in digit frequency
- **Impact:** The model might be slightly more confident in predicting frequently occurring digits
- **Example:** If digit '1' appears slightly more often, the model might have a lower threshold for predicting '1', leading to false positives

**3. Quality Bias**

- **Issue:** Training data consists of clear, centered, well-formed digits

- **Impact:** Poor performance on:
  - Blurry or low-quality images
  - Off-center or rotated digits
  - Partially obscured numbers
  - Varying font sizes
- **Real-world consequence:** Mobile check deposit apps might reject valid checks with slightly smudged numbers

## 4. Demographic Bias in Data Collection

- **Issue:** Original MNIST data collected from a specific demographic (US Census Bureau employees and high school students)
- **Impact:** May not generalize well to:
  - International handwriting styles
  - Different age groups (children vs. elderly)
  - Different educational backgrounds

---

# Bias Identification in Amazon Reviews Model

**Potential Biases:**

## 1. Language and Cultural Bias

- **Issue:** spaCy's `en_core_web_sm` model is trained primarily on English text from specific sources
- **Impact:**
  - Poor performance on non-standard English (AAVE, regional dialects)
  - Misclassification of sentiment in different cultural contexts
  - Failure to recognize entities from non-Western brands
- **Example:** A review saying "This product is sick!" (positive slang) might be classified as negative

## 2. Sentiment Lexicon Bias

- **Issue:** Our rule-based approach uses predefined positive/negative word lists
- **Impact:**
  - Sarcasm and irony are completely missed
  - Context-dependent sentiment is ignored
  - Domain-specific terminology may be misclassified
- **Example:** "This camera is insanely good" → "insanely" might be flagged as negative
- **Example:** "Not bad at all" → double negative interpreted as negative

## 3. Named Entity Recognition Bias

- **Issue:** NER models are trained on news articles and web text, favoring certain entity types
- **Impact:**
  - Well-known Western brands (Apple, Samsung) are easily recognized
  - Smaller, international, or emerging brands are missed
  - Product names that don't follow English patterns are ignored
- **Example:** Chinese brand names written in pinyin might not be recognized as organizations

### 4. Review Length Bias

- **Issue:** Sentiment analysis might be influenced by review length
- **Impact:**
  - Longer reviews accumulate more sentiment words, skewing scores
  - Short, direct reviews might be misclassified as neutral
- **Example:** "Amazing!" (very positive) vs. "This product is okay, not great, but not terrible either" (neutral but longer)

### 5. Product Category Bias

- **Issue:** Different product categories have different review patterns
- **Impact:**
  - Electronics reviews are more technical and less emotional
  - Fashion reviews are more subjective and emotional
  - Same sentiment words mean different things in different contexts
- **Example:** "Basic" is negative for fashion but neutral for electronics

### 6. Temporal Bias

- **Issue:** Language and sentiment expressions evolve over time
- **Impact:**
  - Modern slang and expressions might be misunderstood
  - Historical reviews might be interpreted differently
  - Trending topics affect sentiment interpretation

---

# 2. Mitigation Strategies Using AI Tools

## A. Mitigating MNIST Biases with TensorFlow Tools

### 1. TensorFlow Fairness Indicators

```
# Install Fairness Indicators
pip install tensorflow-model-analysis fairness-indicators
```

```python
# Example implementation
import tensorflow_model_analysis as tfma
from fairness_indicators import remediation

# Define fairness metrics
eval_config = tfma.EvalConfig(
    model_specs=[tfma.ModelSpec(label_key='digit')],
    metrics_specs=[
        tfma.MetricsSpec(
            metrics=[
                tfma.MetricConfig(class_name='FairnessIndicators',
                    config='{"thresholds": [0.5]}')
            ]
        )
    ],
    slicing_specs=[
        # Slice by data source
        tfma.SlicingSpec(feature_keys=['data_source']),
        # Slice by image quality
        tfma.SlicingSpec(feature_keys=['quality_score']),
    ]
)

# Analyze model fairness across different slices
eval_result = tfma.analyze_model(
    eval_shared_model=eval_shared_model,
    data_location=data_location,
    eval_config=eval_config
)
```

**Benefits:**

- Identifies performance disparities across different data slices
- Visualizes fairness metrics across subgroups
- Helps detect where the model underperforms

**2. Data Augmentation to Reduce Bias**

```python
from tensorflow.keras.preprocessing.image import ImageDataGenerator

# Apply augmentation to increase diversity
datagen = ImageDataGenerator(
    rotation_range=15,      # Handle rotated digits
    width_shift_range=0.1,   # Handle off-center digits
    height_shift_range=0.1,
```

```
    zoom_range=0.1,        # Handle size variations
    shear_range=0.1,       # Handle skewed digits
    brightness_range=[0.7, 1.3]  # Handle quality variations
)

# Generate augmented training data
augmented_generator = datagen.flow(X_train, y_train, batch_size=32)
```

**Benefits:**

- Creates synthetic variations of training data
- Improves model robustness to real-world variations
- Reduces bias toward perfectly centered, clear images

### 3. Balanced Sampling

```
from sklearn.utils import class_weight

# Calculate class weights to handle imbalance
class_weights = class_weight.compute_class_weight(
    'balanced',
    classes=np.unique(y_train),
    y=y_train
)
class_weight_dict = dict(enumerate(class_weights))

# Train with class weights
model.fit(
    X_train, y_train,
    class_weight=class_weight_dict,  # Prevents bias toward majority class
    epochs=20
)
```

### 4. Diverse Data Collection

**Strategy:**

- Collect additional data from diverse sources:
  - Different countries and cultures
  - Various age groups
  - Different writing conditions (pen, pencil, stylus)
  - People with different abilities

**Implementation:**

```
# Combine MNIST with other digit datasets
# - EMNIST (Extended MNIST with more diverse handwriting)
# - Street View House Numbers (SVHN)
# - Custom collected data from diverse populations
```

---

# B. Mitigating Amazon Reviews Biases with spaCy

### 1. Custom Rule-Based Systems with Context Awareness

```
import spacy
from spacy.matcher import Matcher

nlp = spacy.load("en_core_web_sm")

# Create matcher for context-aware sentiment
matcher = Matcher(nlp.vocab)

# Handle negations better
negation_pattern = [
    {"LOWER": {"IN": ["not", "no", "never"]}},
    {"OP": "*", "IS_ALPHA": True, "LENGTH": {"<": 10}},  # Up to N words
    {"LOWER": {"IN": ["bad", "terrible", "awful"]}}
]
matcher.add("NEGATION_POSITIVE", [negation_pattern])

# Handle intensifiers with context
intensifier_pattern = [
    {"LOWER": {"IN": ["very", "extremely", "really"]}},
    {"LOWER": {"IN": ["good", "bad", "great"]}}
]
matcher.add("INTENSIFIED", [intensifier_pattern])

# Handle sarcasm indicators
sarcasm_pattern = [
    {"LOWER": "yeah"},
    {"LOWER": "right"},
    {"OP": "+"},
    {"LOWER": {"IN": ["!", "..."]}}
]
matcher.add("SARCASM", [sarcasm_pattern])

def context_aware_sentiment(text):
    doc = nlp(text)
```

```python
    matches = matcher(doc)

    # Adjust sentiment based on patterns
    sentiment_adjustment = 0
    for match_id, start, end in matches:
        if nlp.vocab.strings[match_id] == "NEGATION_POSITIVE":
            sentiment_adjustment += 2  # "not bad" is actually positive

    return sentiment_adjustment
```

**Benefits:**

- Handles linguistic nuances like negation and sarcasm
- Context-aware sentiment analysis
- Reduces misclassification from simple word counting

## 2. Domain-Specific Entity Recognition

```python
# Train custom NER model for product-specific entities
from spacy.training import Example

# Create training data with product-specific entities
TRAIN_DATA = [
    ("The iPhone 13 Pro Max is amazing",
     {"entities": [(4, 21, "PRODUCT")]}),
    ("Samsung Galaxy S21 has great features",
     {"entities": [(0, 18, "PRODUCT")]}),
    # Add more training examples...
]

# Fine-tune spaCy NER on product domain
nlp = spacy.load("en_core_web_sm")
ner = nlp.get_pipe("ner")

# Add custom label
ner.add_label("PRODUCT")

# Train the model
optimizer = nlp.resume_training()
for epoch in range(10):
    for text, annotations in TRAIN_DATA:
        example = Example.from_dict(nlp.make_doc(text), annotations)
        nlp.update([example], sgd=optimizer)
```

**Benefits:**

- Recognizes domain-specific products and brands
- Improves entity extraction for niche or international brands
- Reduces bias toward well-known Western brands

### 3. Multi-Language and Cultural Adaptation

```python
# Load appropriate language model
nlp_en = spacy.load("en_core_web_sm")
nlp_es = spacy.load("es_core_news_sm")  # Spanish
nlp_de = spacy.load("de_core_news_sm")  # German

# Detect language and use appropriate model
from langdetect import detect

def process_review_multilingual(text):
    try:
        lang = detect(text)
        if lang == 'es':
            return nlp_es(text)
        elif lang == 'de':
            return nlp_de(text)
        else:
            return nlp_en(text)
    except:
        return nlp_en(text)  # Default to English
```

### 4. Bias Detection and Monitoring

```python
# Monitor model performance across different groups
def analyze_bias_metrics(reviews_df):
    """
    Analyze sentiment prediction bias across different dimensions
    """
    bias_report = {}

    # 1. Review length bias
    reviews_df['length_category'] = pd.cut(
        reviews_df['review'].str.len(),
        bins=[0, 100, 300, 1000],
        labels=['short', 'medium', 'long']
    )
```

```
    length_bias =
reviews_df.groupby('length_category')['sentiment'].value_counts(normalize=True)
    bias_report['length_bias'] = length_bias

    # 2. Product category bias (if available)
    if 'category' in reviews_df.columns:
        category_bias = reviews_df.groupby('category')['sentiment'].value_counts(normalize=True)
        bias_report['category_bias'] = category_bias

    # 3. Sentiment score distribution
    bias_report['score_distribution'] = reviews_df.groupby('sentiment')['score'].describe()

    return bias_report

# Run bias analysis
bias_metrics = analyze_bias_metrics(df)
print(bias_metrics)
```

---

# 3. Troubleshooting Challenge: Buggy TensorFlow Code

## Original Buggy Code

```
import tensorflow as tf
from tensorflow import keras

# BUGGY CODE - DO NOT USE
model = keras.Sequential([
    keras.layers.Dense(128, activation='relu', input_shape=(28, 28)),  # BUG 1
    keras.layers.Dense(64, activation='relu'),
    keras.layers.Dense(10, activation='sigmoid')  # BUG 2
])

model.compile(
    optimizer='sgd',  # BUG 3
    loss='mean_squared_error',  # BUG 4
    metrics=['accuracy']
)

model.fit(X_train, y_train, epochs=10)  # BUG 5
```

## Bugs Identified and Fixed

**Bug 1: Incorrect Input Shape**

**Problem:** CNN expects 3D input (height, width, channels), not 2D

```
# WRONG
input_shape=(28, 28)

# CORRECT
input_shape=(28, 28, 1)
# Need to reshape data: X_train = X_train.reshape(-1, 28, 28, 1)
```

**Bug 2: Wrong Activation Function**

**Problem:** Sigmoid for multi-class classification (should use softmax)

```
# WRONG
keras.layers.Dense(10, activation='sigmoid')

# CORRECT
keras.layers.Dense(10, activation='softmax')
```

**Bug 3: Suboptimal Optimizer**

**Problem:** SGD without momentum is slow and unstable

```
# SUBOPTIMAL
optimizer='sgd'

# BETTER
optimizer='adam'
# OR with custom learning rate
optimizer=keras.optimizers.Adam(learning_rate=0.001)
```

**Bug 4: Incorrect Loss Function**

**Problem:** MSE is for regression, not classification

```
# WRONG
loss='mean_squared_error'

# CORRECT for one-hot encoded labels
loss='categorical_crossentropy'
```

```
# OR for integer labels
loss='sparse_categorical_crossentropy'
```

**Bug 5: Data Not Preprocessed**

**Problem:** Data not normalized or reshaped

```
# MISSING PREPROCESSING
# X_train values are in [0, 255] - need normalization
# X_train shape is (60000, 28, 28) - need channel dimension

# CORRECT PREPROCESSING
X_train = X_train.reshape(-1, 28, 28, 1).astype('float32') / 255.0
X_test = X_test.reshape(-1, 28, 28, 1).astype('float32') / 255.0

# If using categorical_crossentropy, one-hot encode labels
y_train = keras.utils.to_categorical(y_train, 10)
y_test = keras.utils.to_categorical(y_test, 10)
```

## Complete Fixed Code

```
import tensorflow as tf
from tensorflow import keras
from tensorflow.keras.datasets import mnist

# Load data
(X_train, y_train), (X_test, y_test) = mnist.load_data()

# FIX: Proper preprocessing
X_train = X_train.reshape(-1, 28, 28, 1).astype('float32') / 255.0
X_test = X_test.reshape(-1, 28, 28, 1).astype('float32') / 255.0

# FIX: One-hot encode labels for categorical_crossentropy
y_train = keras.utils.to_categorical(y_train, 10)
y_test = keras.utils.to_categorical(y_test, 10)

# FIX: Proper model architecture
model = keras.Sequential([
    # FIX 1: Correct input shape with channel dimension
    keras.layers.Conv2D(32, (3, 3), activation='relu', input_shape=(28, 28, 1)),
    keras.layers.MaxPooling2D((2, 2)),
    keras.layers.Conv2D(64, (3, 3), activation='relu'),
    keras.layers.MaxPooling2D((2, 2)),
```

```python
    keras.layers.Flatten(),
    keras.layers.Dense(128, activation='relu'),
    keras.layers.Dropout(0.5),
    # FIX 2: Softmax for multi-class classification
    keras.layers.Dense(10, activation='softmax')
])

# FIX: Proper compilation
model.compile(
    # FIX 3: Adam optimizer (better than basic SGD)
    optimizer='adam',
    # FIX 4: Categorical crossentropy for multi-class
    loss='categorical_crossentropy',
    metrics=['accuracy']
)

# FIX 5: Add validation split and proper batch size
history = model.fit(
    X_train, y_train,
    epochs=10,
    batch_size=128,
    validation_split=0.2,
    verbose=1
)

# Evaluate
test_loss, test_accuracy = model.evaluate(X_test, y_test)
print(f"Test Accuracy: {test_accuracy*100:.2f}%")
```

## Additional Improvements

```python
# Add callbacks for better training
early_stopping = keras.callbacks.EarlyStopping(
    monitor='val_loss',
    patience=3,
    restore_best_weights=True
)

reduce_lr = keras.callbacks.ReduceLROnPlateau(
    monitor='val_loss',
    factor=0.5,
    patience=2,
    min_lr=0.00001
)
```

```
# Add batch normalization for stability
model = keras.Sequential([
    keras.layers.Conv2D(32, (3, 3), activation='relu', input_shape=(28, 28, 1)),
    keras.layers.BatchNormalization(),  # IMPROVEMENT
    keras.layers.MaxPooling2D((2, 2)),
    keras.layers.Conv2D(64, (3, 3), activation='relu'),
    keras.layers.BatchNormalization(),  # IMPROVEMENT
    keras.layers.MaxPooling2D((2, 2)),
    keras.layers.Flatten(),
    keras.layers.Dense(128, activation='relu'),
    keras.layers.BatchNormalization(),  # IMPROVEMENT
    keras.layers.Dropout(0.5),
    keras.layers.Dense(10, activation='softmax')
])

# Train with callbacks
history = model.fit(
    X_train, y_train,
    epochs=20,
    batch_size=128,
    validation_split=0.2,
    callbacks=[early_stopping, reduce_lr],
    verbose=1
)
```

---

# 4. Best Practices for Ethical AI Development

## General Principles

1. **Diverse Data Collection**

   - Ensure training data represents all user demographics
   - Actively seek underrepresented groups
   - Document data sources and collection methods

2. **Regular Bias Audits**

   - Test model performance across different subgroups
   - Monitor for performance disparities
   - Set fairness metrics alongside accuracy metrics

3. **Transparency and Documentation**

   - Document model limitations
   - Disclose training data characteristics
   - Explain model decisions when possible

4. **Human-in-the-Loop Systems**

   - Don't automate high-stakes decisions entirely
   - Provide human oversight mechanisms
   - Allow for appeals and corrections

5. **Continuous Monitoring**

   - Track model performance in production
   - Watch for drift in predictions
   - Update models with new, diverse data

## Specific Tool Recommendations

**For Classification Models (TensorFlow/PyTorch):**

- Use **Fairness Indicators** to detect bias
- Implement **Model Cards** to document model behavior
- Use **What-If Tool** for interactive bias exploration
- Apply **adversarial debiasing** techniques

**For NLP Models (spaCy):**

- Train on diverse text corpora
- Use **context-aware sentiment analysis**
- Implement **domain adaptation** for specific use cases
- Validate with **human annotation** on diverse examples

**For All Models:**

- **Cross-validation** across demographic groups
- **A/B testing** with fairness metrics
- **Ablation studies** to understand feature importance
- **Red team testing** to find edge cases

---

# 5. Summary and Recommendations

## Key Takeaways

1. **Bias is Inevitable but Manageable**

   ○ All models have biases from training data
   ○ Active mitigation is essential
   ○ Regular monitoring prevents bias amplification
2. **Tools Exist to Help**

   ○ TensorFlow Fairness Indicators
   ○ spaCy's customizable pipelines
   ○ Data augmentation techniques
   ○ Fairness-aware learning algorithms
3. **Ethical AI Requires Intentionality**

   ○ Design for fairness from the start
   ○ Test across diverse scenarios
   ○ Document limitations honestly
   ○ Update models regularly

## Action Items for Our Project

**Immediate Actions:**

- ✅ Document identified biases in our models
- ✅ Implement data augmentation for MNIST
- ✅ Create context-aware rules for sentiment analysis
- ✅ Add fairness metrics to evaluation pipeline

**Future Improvements:**

- 📋 Collect more diverse training data
- 📋 Implement TensorFlow Fairness Indicators
- 📋 Fine-tune spaCy on domain-specific data
- 📋 Create model cards documenting limitations
- 📋 Set up continuous monitoring in production

## Ethical Commitment

As AI engineers, we commit to:

1. **Acknowledging** biases in our models
2. **Measuring** fairness alongside accuracy
3. **Mitigating** identified biases through technical solutions
4. **Monitoring** model behavior in production
5. **Updating** models as new biases emerge
6. **Documenting** limitations transparently

7. **Prioritizing** human welfare over model performance

---

# 6. References and Further Reading

## Academic Papers

- "Fairness and Machine Learning" by Barocas, Hardt, and Narayanan
- "Model Cards for Model Reporting" by Mitchell et al.
- "On Fairness and Calibration" by Pleiss et al.

## Tools and Frameworks

- [TensorFlow Fairness Indicators](#)
- [AI Fairness 360](#) by IBM
- [Fairlearn](#) by Microsoft
- [What-If Tool](#)

## Industry Guidelines

- Google's Responsible AI Practices
- Microsoft's Responsible AI Standard
- IEEE Ethics in AI Standards
- EU AI Act Requirements

---

**Conclusion:** Ethical AI development is not a one-time task but an ongoing commitment. By using the tools and strategies outlined above, we can build more fair, transparent, and trustworthy AI systems that serve all users equitably.

## Conclusion

# Summary

**TensorFlow vs PyTorch:** Choose based on deployment needs (TensorFlow) vs research flexibility (PyTorch).

**Jupyter Notebooks:** Essential for exploration, experimentation, and sharing reproducible research.

**spaCy vs Basic Python:** spaCy provides linguistic intelligence, making NLP tasks orders of magnitude easier and more accurate.

**Scikit-learn vs TensorFlow:** Scikit-learn for classical ML on tabular data (easier), TensorFlow for deep learning on unstructured data (more powerful but complex).

The modern AI engineer's toolkit includes all these tools, selecting the right one based on the specific problem, data type, and deployment requirements.

## References