

Ethics & Optimization

1. **Ethical Considerations** Identify potential biases in your MNIST or Amazon Reviews model. How could tools like TensorFlow Fairness Indicators or spaCy's rule-based systems mitigate these biases?

Ethical Considerations: Bias in MNIST and Amazon Reviews Models

Potential Biases:

MNIST Model:

- **Dataset Bias:** MNIST is generally considered balanced, but variants like Biased-MNIST introduce artificial correlations (e.g., associating certain digits with specific background colors).
- **Impact:** Models may learn to rely on background color or other spurious features instead of digit shape, leading to poor generalization and unfair performance on unbiased or real-world data.
- **Minority Representation:** If certain digits are underrepresented, the model may perform worse on those classes.

Amazon Reviews Model:

- **Demographic Bias:** Reviews may reflect the language, preferences, or cultural context of majority groups, leading to sentiment models that misinterpret or undervalue minority voices.
- **Labelling Bias:** Sentiment labels or NER annotations may be inconsistent or reflect annotator subjectivity, affecting model fairness and accuracy.

Mitigating Bias with Tools:

TensorFlow Fairness Indicators:

- These tools evaluate model performance across different subgroups (e.g., gender, age, language) to detect disparities¹.
- By visualizing metrics like accuracy, precision, and recall for each subgroup, developers can identify and address unfair outcomes.
- Models can then be retrained with balanced datasets, reweighted loss functions, or adversarial debiasing to reduce bias.

spaCy's Rule-Based Systems:

- Rule-based approaches in spaCy allow customization of entity recognition and sentiment analysis pipelines to correct systematic errors or biases (e.g., ensuring product names from minority-owned brands are recognized)¹.
- Custom patterns and post-processing rules can be used to supplement or override model predictions where bias is detected.

2. 2. Troubleshooting Challenge Buggy Code: A provided TensorFlow script has errors (e.g., dimension mismatches, incorrect loss functions). Debug and fix the code.

Step 1: Identify Dimension Mismatches

- Check Tensor Shapes:
Print the shapes of your tensors before operations (e.g., before passing data to layers or loss functions). Use `tensor.shape` or `print(tensor.shape)` to inspect dimensions.
- Common Issues:
 - Inputs to layers (e.g., `Dense`, `Conv2D`) do not match the expected shape.
 - Labels and predictions have incompatible shapes for the chosen loss function.
 - Matrix multiplication or concatenation with mismatched dimensions.

Example Fix:

```
import tensorflow as tf

# Example: Dense layer expects input of shape (batch_size, features)

model = tf.keras.Sequential([

    tf.keras.layers.Dense(10, input_shape=(5,)), # expects input with 5 features

])

input_data = tf.random.uniform([1, 5]) # Correct shape

output = model(input_data) # No error

# If you use input_data = tf.random.uniform([1, 6]), you'll get a ValueError[4].
```

Step 2: Use Assertions and Debugging Tools

Assertions:

Use `tf.debugging.assert_shapes` to enforce expected tensor shapes during runtime.

python

```
tf.debugging.assert_shapes([
    (input_data, ('batch', 5)), # batch size can be any, features must be 5
])
```

- **Print Statements:**
Add print statements before key operations to confirm shapes.

Step 3: Check and Correct Loss Functions

- **Classification:**
For multi-class classification, use `SparseCategoricalCrossentropy` if your labels are integers, or `CategoricalCrossentropy` if your labels are one-hot encoded.
- **Regression:**
Use `MeanSquaredError` or similar for regression tasks.

Example Fix:

python

```
# For integer labels (e.g., [0, 1, 2])
```

```
loss_fn = tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True)
```

```
# For one-hot labels (e.g., [1, 0, 0])
```

```
# loss_fn = tf.keras.losses.CategoricalCrossentropy(from_logits=True)
```

Step 4: General Debugging Strategy

- **Check Data Loading and Preprocessing:**
Ensure data is properly shaped and pre-processed before feeding into the model.
- **Layer Output Shapes:**
Print or inspect output shapes after each layer in your model.

References

- Use print statements and `tensor.shape` for shape inspection.

- Use `tf.debugging.assert_shapes` for runtime shape checks.
- Match input data shapes to layer expectations and loss function requirements.