# Project Title: Clothing Item Classification Using CNN and Spark

By: Feven Tefera
Mihret Kemal

# Goal

- Classify clothing items into 10 categories using Fashion MNIST

- Address challenges:

  - Overfitting - with drop out and early stopping

  - Computational efficiency - with Spark

# Real World Application

- Online clothing stores like **Amazon**, and  **Zara**  can use the model to classify product images into categories

- Visual search capabilities where users can find visually similar items based on images **(Google, Pinterest)**

- **Amazon** or **eBay** could use automated captions for product descriptions based on images
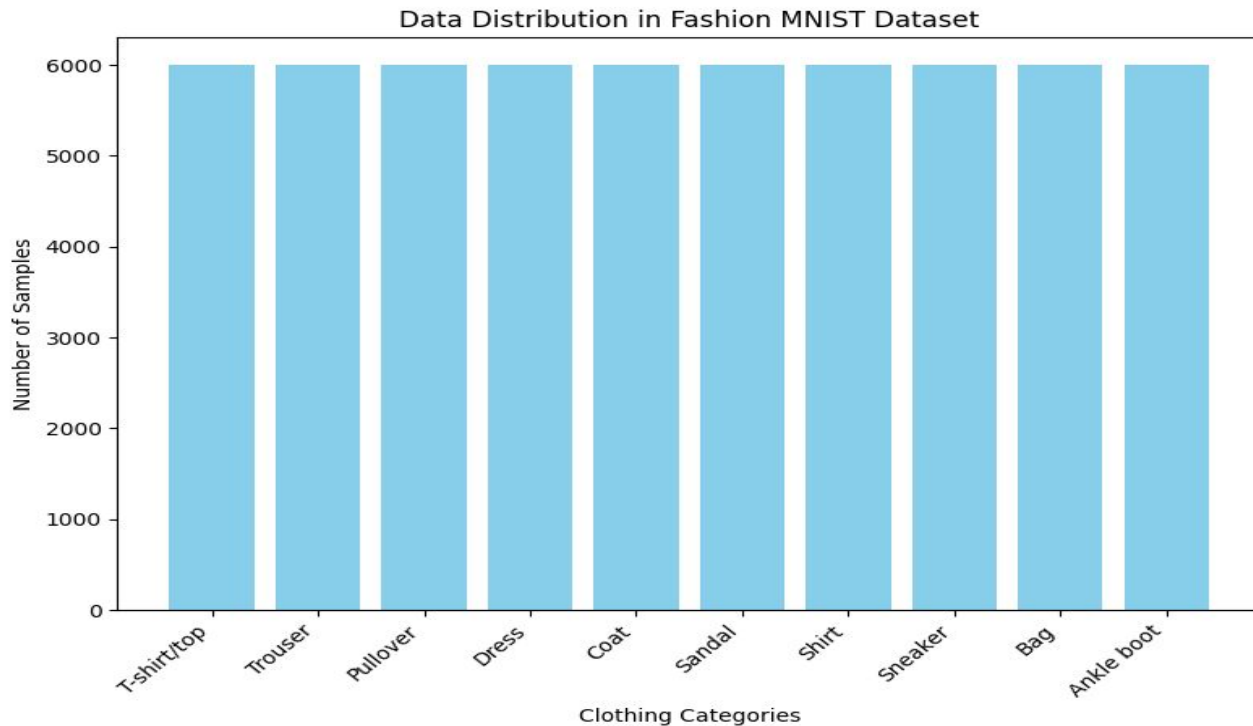
# Dataset Overview

- **Dataset**: [Fashion MNIST](#)

- **Categories**: T-shirt/top, Trouser, Pullover, Dress, Coat, Sandal, Shirt, Sneaker, Bag, Ankle Boot

- **Sample Size**:

  - 60,000 training images

  - 10,000 test images

**Data split:** 90% - training set, and 10% validation set

# Data Distribution

- Equal distribution across 10 categories.



Data Distribution in Fashion MNIST Dataset

# Implementation Tools

- **Libraries**: PySpark, TensorFlow, Keras, NumPy, Sklearn, MatplotLib

- **Processing**:

  - Spark for distributed data handling

  - Keras for CNN design and training

# CNN Model Architecture

```
Model Architecture:
Model: "sequential_25"
```

| Layer (type) | Output Shape | Param # |
|---|---|---|
| conv2d_50 (Conv2D) | (None, 26, 26, 32) | 320 |
| max_pooling2d_50 (MaxPooling2D) | (None, 13, 13, 32) | 0 |
| conv2d_51 (Conv2D) | (None, 11, 11, 64) | 18,496 |
| max_pooling2d_51 (MaxPooling2D) | (None, 5, 5, 64) | 0 |
| dropout_50 (Dropout) | (None, 5, 5, 64) | 0 |
| flatten_25 (Flatten) | (None, 1600) | 0 |
| dense_50 (Dense) | (None, 128) | 204,928 |
| dropout_51 (Dropout) | (None, 128) | 0 |
| dense_51 (Dense) | (None, 10) | 1,290 |

```
Total params: 225,034 (879.04 KB)
Trainable params: 225,034 (879.04 KB)
Non-trainable params: 0 (0.00 B)
```

# Tuning Hyperparameters

- Learning rate: {0.01, 0.001}

- Batch size: {64,32}

- Dropout: {0.3,0.5}

- We used **Grid search** with 3 fold cv for hyperparameter tuning
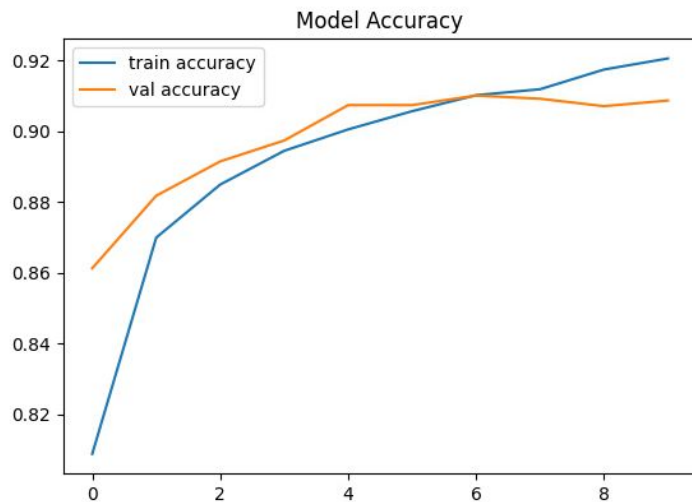
# Best Hyperparameters

- Learning rate: 0.001

- Batch size: 32

- Dropout: 0.3

```
Best parameters: {'batch_size': 32, 'model__dropout_rate': 0.3, 'model__learning_rate': 0.001}
Best accuracy: 0.9096000000000001

Epoch 1/10
1875/1875 ━━━━━━━━━━━━━━ 156s 82ms/step - accuracy: 0.7414 - loss: 0.7043 - val_accuracy: 0.8612 - val_loss: 0.3818
Epoch 2/10
1875/1875 ━━━━━━━━━━━━━━ 202s 82ms/step - accuracy: 0.8646 - loss: 0.3694 - val_accuracy: 0.8817 - val_loss: 0.3280
Epoch 3/10
1875/1875 ━━━━━━━━━━━━━━ 207s 85ms/step - accuracy: 0.8862 - loss: 0.3123 - val_accuracy: 0.8914 - val_loss: 0.2983
Epoch 4/10
1875/1875 ━━━━━━━━━━━━━━ 171s 68ms/step - accuracy: 0.8938 - loss: 0.2844 - val_accuracy: 0.8973 - val_loss: 0.2786
Epoch 5/10
1875/1875 ━━━━━━━━━━━━━━ 144s 70ms/step - accuracy: 0.9005 - loss: 0.2667 - val_accuracy: 0.9073 - val_loss: 0.2575
Epoch 6/10
1875/1875 ━━━━━━━━━━━━━━ 141s 69ms/step - accuracy: 0.9059 - loss: 0.2501 - val_accuracy: 0.9073 - val_loss: 0.2519
Epoch 7/10
1875/1875 ━━━━━━━━━━━━━━ 142s 69ms/step - accuracy: 0.9084 - loss: 0.2409 - val_accuracy: 0.9100 - val_loss: 0.2425
Epoch 8/10
1875/1875 ━━━━━━━━━━━━━━ 141s 69ms/step - accuracy: 0.9122 - loss: 0.2280 - val_accuracy: 0.9091 - val_loss: 0.2462
Epoch 9/10
1875/1875 ━━━━━━━━━━━━━━ 141s 68ms/step - accuracy: 0.9185 - loss: 0.2196 - val_accuracy: 0.9070 - val_loss: 0.2547
Epoch 10/10
1875/1875 ━━━━━━━━━━━━━━ 142s 68ms/step - accuracy: 0.9214 - loss: 0.2065 - val_accuracy: 0.9086 - val_loss: 0.2506
```

# Performance Metrics

- **Test Accuracy**: 0.908599

- **Test Loss**: 0.2506

# Model Performance on Two VMs

Time: 3.5 hours

# Combined Performance of Four VMs:

Time: 2.2 hours

# Challenges and Solutions

**Challenges**:

- Preventing overfitting
- Computational resource constraints

**Solutions**:

- Dropout layers and early stopping(patience = 3)
- Efficient training using Spark

# Thank you!

# Any Questions?