# Detailed Explanation of LSTM + DQN (and PPO) Trading Workflow

Your Name

February 24, 2025

## Contents

# 1 Introduction

This document provides a step-by-step explanation of the Python program that trains:

1. An LSTM neural network (for price prediction) using Optuna for hyperparameter tuning.

2. A Deep Q-Network (DQN) agent (from Stable-Baselines3) in a custom Stock Trading Gym environment for reinforcement learning.

Additionally, it will cover:

- The benefits of potentially substituting the DQN algorithm with Proximal Policy Optimization (PPO).

- A detailed guide on how to modify the existing code to implement a PPO agent instead of a DQN agent.

# 2 Full Source Code

In this section, we display the full source code. Each component will be described in detail in subsequent sections.

Listing 1: Full program source code.

```python
import os
import sys
import argparse
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import logging
from tabulate import tabulate

from sklearn.preprocessing import MinMaxScaler
from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score

import tensorflow as tf
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import LSTM, Dense, Dropout, Bidirectional
from tensorflow.keras.optimizers import Adam, Nadam
from tensorflow.keras.callbacks import EarlyStopping, ReduceLROnPlateau
from tensorflow.keras.losses import Huber
from tensorflow.keras.regularizers import l2

import xgboost as xgb
import optuna
from optuna.integration import KerasPruningCallback

# For Reinforcement Learning
import gym
from gym import spaces
from stable_baselines3 import DQN
from stable_baselines3.common.vec_env import DummyVecEnv

# Suppress TensorFlow warnings
os.environ['TF_CPP_MIN_LOG_LEVEL'] = '2'  # Suppress INFO/WARNING

# Configure logging
logging.basicConfig(level=logging.INFO, format='%(asctime)s - %(levelname)s - %(message)s')


##############################
# 1. Data Loading / Indicators
##############################
def load_data(file_path):
    logging.info(f"Loading data from: {file_path}")
    try:
        df = pd.read_csv(file_path, parse_dates=['time'])
    except FileNotFoundError:
```

```python
            logging.error(f"File not found: {file_path}")
            sys.exit(1)
        except pd.errors.ParserError as e:
            logging.error(f"Error parsing CSV file: {e}")
            sys.exit(1)
        except Exception as e:
            logging.error(f"Unexpected error: {e}")
            sys.exit(1)

    rename_mapping = {
        'time': 'Date',
        'open': 'Open',
        'high': 'High',
        'low': 'Low',
        'close': 'Close'
    }
    df.rename(columns=rename_mapping, inplace=True)

    logging.info(f"Data columns after renaming: {df.columns.tolist()}")

    df.sort_values('Date', inplace=True)
    df.reset_index(drop=True, inplace=True)
    logging.info("Data loaded and sorted successfully.")
    return df


def compute_rsi(series, window=14):
    delta = series.diff()
    gain = delta.where(delta > 0, 0).rolling(window=window).mean()
    loss = -delta.where(delta < 0, 0).rolling(window=window).mean()
    RS = gain / (loss + 1e-9)
    return 100 - (100 / (1 + RS))


def compute_macd(series, span_short=12, span_long=26, span_signal=9):
    ema_short = series.ewm(span=span_short, adjust=False).mean()
    ema_long = series.ewm(span=span_long, adjust=False).mean()
    macd_line = ema_short - ema_long
    signal_line = macd_line.ewm(span=span_signal, adjust=False).mean()
    return macd_line - signal_line  # histogram


def compute_obv(df):
    signed_volume = (np.sign(df['Close'].diff()) * df['Volume']).fillna(0)
    return signed_volume.cumsum()


def compute_adx(df, window=14):
    """Pseudo-ADX approach using rolling True Range / Close."""
    df['H-L'] = df['High'] - df['Low']
    df['H-Cp'] = (df['High'] - df['Close'].shift(1)).abs()
    df['L-Cp'] = (df['Low'] - df['Close'].shift(1)).abs()
    tr = df[['H-L','H-Cp','L-Cp']].max(axis=1)
    tr_rolling = tr.rolling(window=window).mean()

    adx_placeholder = tr_rolling / (df['Close'] + 1e-9)
    df.drop(['H-L','H-Cp','L-Cp'], axis=1, inplace=True)
    return adx_placeholder


def compute_bollinger_bands(series, window=20, num_std=2):
    sma = series.rolling(window=window).mean()
    std = series.rolling(window=window).std()
    upper = sma + num_std * std
    lower = sma - num_std * std
    bandwidth = (upper - lower) / (sma + 1e-9)
    return upper, lower, bandwidth


def compute_mfi(df, window=14):
    typical_price = (df['High'] + df['Low'] + df['Close']) / 3
    money_flow = typical_price * df['Volume']
    prev_tp = typical_price.shift(1)
```

3

```python
120         flow_pos = money_flow.where(typical_price > prev_tp, 0)
121         flow_neg = money_flow.where(typical_price < prev_tp, 0)
122         pos_sum = flow_pos.rolling(window=window).sum()
123         neg_sum = flow_neg.rolling(window=window).sum()
124         mfi = 100 - (100 / (1 + pos_sum/(neg_sum+1e-9)))
125         return mfi


128     def calculate_technical_indicators(df):
129         logging.info("Calculating technical indicators...")
130
131         df['RSI'] = compute_rsi(df['Close'], window=14)
132         df['MACD'] = compute_macd(df['Close'])
133         df['OBV'] = compute_obv(df)
134         df['ADX'] = compute_adx(df)
135
136         upper_bb, lower_bb, bb_width = compute_bollinger_bands(df['Close'], window=20, num_std=2)
137         df['BB_Upper'] = upper_bb
138         df['BB_Lower'] = lower_bb
139         df['BB_Width'] = bb_width
140
141         df['MFI'] = compute_mfi(df, window=14)
142
143         df['SMA_5'] = df['Close'].rolling(window=5).mean()
144         df['SMA_10'] = df['Close'].rolling(window=10).mean()
145         df['EMA_5'] = df['Close'].ewm(span=5, adjust=False).mean()
146         df['EMA_10'] = df['Close'].ewm(span=10, adjust=False).mean()
147
148         df['STDDEV_5'] = df['Close'].rolling(window=5).std()
149         df.dropna(inplace=True)
150         logging.info("Technical indicators calculated successfully.")
151         return df


154     ##############################
155     # 2. ARGUMENT PARSING
156     ##############################
157     def parse_arguments():
158         parser = argparse.ArgumentParser(description='Train LSTM and DQN models for stock trading.
                  ')
159         parser.add_argument('csv_path', type=str, help='Path to the CSV data file (with columns
                  time,open,high,low,close,volume).')
160         return parser.parse_args()


163     ##############################
164     # 3. MAIN
165     ##############################
166     def main():
167         # 1) Parse args
168         args = parse_arguments()
169         csv_path = args.csv_path
170
171         # 2) Load data & advanced indicators
172         data = load_data(csv_path)
173         data = calculate_technical_indicators(data)
174
175         # EXCLUDE 'Close' from feature inputs
176         feature_columns = [
177             'SMA_5', 'SMA_10', 'EMA_5', 'EMA_10', 'STDDEV_5',
178             'RSI', 'MACD', 'ADX', 'OBV', 'Volume', 'Open', 'High', 'Low',
179             'BB_Upper','BB_Lower','BB_Width','MFI'
180         ]
181         target_column = 'Close'
182         data = data[['Date'] + feature_columns + [target_column]]
183         data.dropna(inplace=True)
184
185         # 3) Scale
186         scaler_features = MinMaxScaler()
187         scaler_target = MinMaxScaler()
188
189         X_all = data[feature_columns].values
190         y_all = data[[target_column]].values
```

```python
191
192     X_scaled = scaler_features.fit_transform(X_all)
193     y_scaled = scaler_target.fit_transform(y_all).flatten()
194
195     # 4) Create LSTM Sequences
196     def create_sequences(features, target, window_size=15):
197         X_seq, y_seq = [], []
198         for i in range(len(features) - window_size):
199             X_seq.append(features[i:i+window_size])
200             y_seq.append(target[i+window_size])
201         return np.array(X_seq), np.array(y_seq)
202
203     window_size = 15
204     X, y = create_sequences(X_scaled, y_scaled, window_size)
205
206     # 5) Train/Val/Test Split
207     train_size = int(len(X)*0.7)
208     val_size = int(len(X)*0.15)
209     test_size = len(X) - train_size - val_size
210
211     X_train = X[:train_size]
212     y_train = y[:train_size]
213     X_val   = X[train_size:train_size+val_size]
214     y_val   = y[train_size:train_size+val_size]
215     X_test  = X[train_size+val_size:]
216     y_test  = y[train_size+val_size:]
217
218     logging.info(f"Scaled training features shape: {X_train.shape}")
219     logging.info(f"Scaled validation features shape: {X_val.shape}")
220     logging.info(f"Scaled testing features shape: {X_test.shape}")
221     logging.info(f"Scaled training target shape: {y_train.shape}")
222     logging.info(f"Scaled validation target shape: {y_val.shape}")
223     logging.info(f"Scaled testing target shape: {y_test.shape}")
224
225     # 6) GPU/CPU Config
226     def configure_device():
227         gpus = tf.config.list_physical_devices('GPU')
228         if gpus:
229             try:
230                 for gpu in gpus:
231                     tf.config.experimental.set_memory_growth(gpu, True)
232                 logging.info(f"{len(gpus)} GPU(s) detected and configured.")
233             except RuntimeError as e:
234                 logging.error(e)
235         else:
236             logging.info("No GPU detected, using CPU.")
237
238     configure_device()
239
240     # 7) Build LSTM
241     def build_advanced_lstm(input_shape, hyperparams):
242         model = Sequential()
243         for i in range(hyperparams['num_lstm_layers']):
244             return_seqs = (i < hyperparams['num_lstm_layers'] - 1)
245             model.add(Bidirectional(
246                 LSTM(hyperparams['lstm_units'],
247                     return_sequences=return_seqs,
248                     kernel_regularizer=tf.keras.regularizers.l2(0.001)
249                 ), input_shape=input_shape if i==0 else None))
250             model.add(Dropout(hyperparams['dropout_rate']))
251
252         model.add(Dense(1, activation='linear'))
253
254         # Optimizer
255         if hyperparams['optimizer'] == 'Adam':
256             opt = Adam(learning_rate=hyperparams['learning_rate'], decay=hyperparams['decay'])
257         elif hyperparams['optimizer'] == 'Nadam':
258             opt = Nadam(learning_rate=hyperparams['learning_rate'])
259         else:
260             opt = Adam(learning_rate=hyperparams['learning_rate'])
261
262         model.compile(optimizer=opt, loss=Huber(), metrics=['mae'])
263         return model
```

```python
264
265      # 8) Optuna Tuning
266      def objective(trial):
267          num_lstm_layers = trial.suggest_int('num_lstm_layers', 1, 3)
268          lstm_units = trial.suggest_categorical('lstm_units', [32, 64, 96, 128])
269          dropout_rate = trial.suggest_float('dropout_rate', 0.1, 0.5)
270          learning_rate = trial.suggest_loguniform('learning_rate', 1e-5, 1e-2)
271          optimizer_name = trial.suggest_categorical('optimizer', ['Adam', 'Nadam'])
272          decay = trial.suggest_float('decay', 0.0, 1e-4)
273
274          hyperparams = {
275              'num_lstm_layers': num_lstm_layers,
276              'lstm_units': lstm_units,
277              'dropout_rate': dropout_rate,
278              'learning_rate': learning_rate,
279              'optimizer': optimizer_name,
280              'decay': decay
281          }
282
283          model_ = build_advanced_lstm((X_train.shape[1], X_train.shape[2]), hyperparams)
284
285          early_stop = EarlyStopping(monitor='val_loss', patience=10, restore_best_weights=True)
286          lr_reduce  = ReduceLROnPlateau(monitor='val_loss', factor=0.5, patience=5, min_lr=1e
                 -6)
287
288          cb_prune = KerasPruningCallback(trial, 'val_loss')
289
290          history = model_.fit(
291              X_train, y_train,
292              epochs=100,
293              batch_size=16,
294              validation_data=(X_val, y_val),
295              callbacks=[early_stop, lr_reduce, cb_prune],
296              verbose=0
297          )
298          val_mae = min(history.history['val_mae'])
299          return val_mae
300
301      logging.info("Starting hyperparameter optimization with Optuna...")
302      study = optuna.create_study(direction='minimize')
303      study.optimize(objective, n_trials=50)  # might take a long time
304
305      best_params = study.best_params
306      logging.info(f"Best Hyperparameters from Optuna: {best_params}")
307
308      # 9) Train Best LSTM
309      best_model = build_advanced_lstm((X_train.shape[1], X_train.shape[2]), best_params)
310      early_stop2 = EarlyStopping(monitor='val_loss', patience=20, restore_best_weights=True)
311      lr_reduce2  = ReduceLROnPlateau(monitor='val_loss', factor=0.5, patience=5, min_lr=1e-6)
312
313      logging.info("Training the best LSTM model with optimized hyperparameters...")
314      history = best_model.fit(
315          X_train, y_train,
316          epochs=300,
317          batch_size=16,
318          validation_data=(X_val, y_val),
319          callbacks=[early_stop2, lr_reduce2],
320          verbose=1
321      )
322
323      # 10) Evaluate
324      def evaluate_model(model, X_test, y_test):
325          logging.info("Evaluating model...")
326          # Predict scaled
327          y_pred_scaled = model.predict(X_test).flatten()
328          y_pred_scaled = np.clip(y_pred_scaled, 0, 1)  # clamp if needed
329          # Inverse
330          y_pred = scaler_target.inverse_transform(y_pred_scaled.reshape(-1,1)).flatten()
331          y_test_actual = scaler_target.inverse_transform(y_test.reshape(-1,1)).flatten()
332
333          mse = mean_squared_error(y_test_actual, y_pred)
334          rmse = np.sqrt(mse)
335          mae = mean_absolute_error(y_test_actual, y_pred)
```

```python
336             r2   = r2_score(y_test_actual, y_pred)
337
338             # Directional accuracy
339             direction_actual = np.sign(np.diff(y_test_actual))
340             direction_pred   = np.sign(np.diff(y_pred))
341             directional_accuracy = np.mean(direction_actual == direction_pred)
342
343             logging.info(f"Test MSE: {mse}")
344             logging.info(f"Test RMSE: {rmse}")
345             logging.info(f"Test MAE: {mae}")
346             logging.info(f"Test R2 Score: {r2}")
347             logging.info(f"Directional Accuracy: {directional_accuracy}")
348
349             # Plot
350             plt.figure(figsize=(14, 7))
351             plt.plot(y_test_actual, label='Actual Price')
352             plt.plot(y_pred, label='Predicted Price')
353             plt.title('Actual vs Predicted Prices')
354             plt.xlabel('Time Step')
355             plt.ylabel('Price')
356             plt.legend()
357             plt.grid(True)
358             plt.savefig('actual_vs_predicted.png')
359             plt.close()
360             logging.info("Actual vs Predicted plot saved as 'actual_vs_predicted.png'")
361
362             # Tabulate first 40 predictions
363             table_data = []
364             for i in range(min(40, len(y_test_actual))):
365                 table_data.append([i, round(y_test_actual[i],2), round(y_pred[i],2)])
366             headers = ["Index", "Actual Price", "Predicted Price"]
367             print(tabulate(table_data, headers=headers, tablefmt="pretty"))
368
369             return mse, rmse, mae, r2, directional_accuracy
370
371     mse, rmse, mae, r2, directional_accuracy = evaluate_model(best_model, X_test, y_test)
372
373     # 11) Save
374     best_model.save('optimized_lstm_model.h5')
375     import joblib
376     joblib.dump(scaler_features, 'scaler_features.save')
377     joblib.dump(scaler_target, 'scaler_target.save')
378     logging.info("Model and scalers saved as 'optimized_lstm_model.h5', 'scaler_features.save
                ', and 'scaler_target.save'.")
379
380     #######################################
381     # 12) Reinforcement Learning Environment
382     #######################################
383     class StockTradingEnv(gym.Env):
384         """
385         A simple stock trading environment for OpenAI Gym
386         """
387         metadata = {'render.modes': ['human']}
388
389         def __init__(self, df, initial_balance=10000):
390             super().__init__()
391             self.df = df.reset_index()
392             self.initial_balance = initial_balance
393             self.balance = initial_balance
394             self.net_worth = initial_balance
395             self.max_steps = len(df)
396             self.current_step = 0
397             self.shares_held = 0
398             self.cost_basis = 0
399
400             # We re-use feature_columns from above
401             # (Excluding 'Close' from the observation)
402             # Actions: 0=Sell, 1=Hold, 2=Buy
403             self.action_space = spaces.Discrete(3)
404
405             # Observations => advanced feature columns + 3 additional (balance, shares,
                    cost_basis)
406             self.observation_space = spaces.Box(
```

```python
                low=0,
                high=1,
                shape=(len(feature_columns) + 3,),
                dtype=np.float32
            )

    def reset(self):
        self.balance = self.initial_balance
        self.net_worth = self.initial_balance
        self.current_step = 0
        self.shares_held = 0
        self.cost_basis = 0
        return self._next_observation()

    def _next_observation(self):
        obs = self.df.loc[self.current_step, feature_columns].values
        # Simple normalization by max to keep it [0,1]
        obs = obs / np.max(obs) if np.max(obs)!=0 else obs

        additional = np.array([
            self.balance / self.initial_balance,
            self.shares_held / 100.0,
            self.cost_basis / self.initial_balance
        ])
        return np.concatenate([obs, additional])

    def step(self, action):
        current_price = self.df.loc[self.current_step, 'Close']

        if action == 2:  # Buy
            total_possible = self.balance // current_price
            shares_bought = total_possible
            if shares_bought > 0:
                self.balance -= shares_bought * current_price
                self.shares_held += shares_bought
                self.cost_basis = (
                    (self.cost_basis * (self.shares_held - shares_bought)) +
                    (shares_bought * current_price)
                ) / self.shares_held

        elif action == 0:  # Sell
            if self.shares_held > 0:
                self.balance += self.shares_held * current_price
                self.shares_held = 0
                self.cost_basis = 0

        self.net_worth = self.balance + self.shares_held * current_price
        self.current_step += 1

        done = (self.current_step >= self.max_steps - 1)
        reward = self.net_worth - self.initial_balance

        obs = self._next_observation()
        return obs, reward, done, {}

    def render(self, mode='human'):
        profit = self.net_worth - self.initial_balance
        print(f"Step: {self.current_step}")
        print(f"Balance: {self.balance}")
        print(f"Shares held: {self.shares_held} (Cost Basis: {self.cost_basis})")
        print(f"Net worth: {self.net_worth}")
        print(f"Profit: {profit}")

def train_dqn_agent(env):
    logging.info("Training DQN Agent...")
    try:
        model = DQN(
            'MlpPolicy',
            env,
            verbose=1,
            learning_rate=1e-3,
            buffer_size=10000,
            learning_starts=1000,
```

```
480            batch_size=64,
481            tau=1.0,
482            gamma=0.99,
483            train_freq=4,
484            target_update_interval=1000,
485            exploration_fraction=0.1,
486            exploration_final_eps=0.02,
487            tensorboard_log="./dqn_stock_tensorboard/"
488        )
489        model.learn(total_timesteps=100000)
490        model.save("dqn_stock_trading")
491        logging.info("DQN Agent trained and saved as 'dqn_stock_trading.zip'.")
492        return model
493    except Exception as e:
494        logging.error(f"Error training DQN Agent: {e}")
495        sys.exit(1)
496
497    # Initialize RL environment
498    logging.info("Initializing and training DQN environment...")
499    trading_env = StockTradingEnv(data)
500    trading_env = DummyVecEnv([lambda: trading_env])
501
502    # Train
503    dqn_model = train_dqn_agent(trading_env)
504
505    logging.info("All tasks complete. Exiting.")
506
507
508 if __name__ == "__main__":
509    main()
```

# 3 Step-by-Step Explanation

## 3.1 Imports and Configuration

- `numpy`, `pandas`, `sklearn` are used for data manipulation and scaling.

- `matplotlib` and `seaborn` are used for plotting.

- `tensorflow.keras` is used for building the LSTM model.

- `optuna` is used for hyperparameter tuning.

- `gym` and `stable_baselines3` are used for the reinforcement learning (RL) portion of the code.

- The environment variable `TF_CPP_MIN_LOG_LEVEL` suppresses TensorFlow messages to keep logs clean.

## 3.2 Data Loading and Preprocessing (Functions in Step 1)

1. **load_data(file_path)**:

   - Reads a CSV file, expecting columns `time`, `open`, `high`, `low`, `close`, `volume`.
   - Renames these columns to `Date`, `Open`, `High`, `Low`, `Close`.
   - Sorts the DataFrame by date, resets the index, and returns the processed DataFrame.

2. **Technical Indicators Functions**: The code defines several functions to compute widely used technical indicators:

   - `compute_rsi(series, window=14)`: Calculates the Relative Strength Index over a given window.
   - `compute_macd(series, ...)`: Computes the Moving Average Convergence Divergence (MACD) histogram.
   - `compute_obv(df)`: Computes On-Balance Volume.
   - `compute_adx(df, window=14)`: Returns a placeholder ADX-like measure.

- `compute_bollinger_bands(series, ...)`: Returns Bollinger Bands and the bandwidth.
- `compute_mfi(df, window=14)`: Computes the Money Flow Index.

3. **`calculate_technical_indicators(df)`**:

   - Applies all of the above functions on the `df`.
   - Creates short/medium moving averages like SMA, EMA.
   - Removes any resulting NaN values from the DataFrame.
   - Returns the DataFrame with all computed indicators.

## 3.3 Argument Parsing (Step 2)

- Uses `argparse` to read a single argument, `csv_path`, which specifies the path to the CSV data file.

## 3.4 Main Function (Step 3)

The **main** function orchestrates the entire workflow:

1. **Parse arguments**: Calls `parse_arguments()` to get `csv_path`.

2. **Load and process data**:

   - Calls `load_data(csv_path)` to load the CSV.
   - Calls `calculate_technical_indicators(data)` to augment it with technical features.

3. **Setup feature and target columns**: The code excludes `Close` from the features (since it is the prediction target) and identifies the target as the `Close` column.

4. **Scaling**:

   - Instantiates two `MinMaxScalers`: one for features and one for the target.
   - Fits and transforms the data to a 0-1 range.

5. **Sequence Creation** (`create_sequences`):

   - Converts a time series of scaled features/targets into 3D sequences (samples × window_size × number_of_features).
   - Correspondingly, the targets become a vector of length (`original_length - window_size`).

6. **Train-Validation-Test Split**: Splits the resulting 3D array (and 1D targets) as 70% train, 15% val, and 15% test.

7. **GPU/CPU Configuration**: The function `configure_device()` attempts to set memory growth on any detected GPUs.

## 3.5 Building and Tuning the LSTM (Steps 7–9)

1. **`build_advanced_lstm`**:

   - Takes `input_shape` (the shape of the input sequence) and a `hyperparams` dictionary.
   - Builds a stacked or repeated LSTM architecture, each wrapped in a Bidirectional layer, with dropout in between.
   - A final Dense layer of size 1 is added to output the price prediction.
   - Compilation uses `Huber` loss and one of the possible optimizers (`Adam` or `Nadam`) with various learning rate/decay settings.

2. **Optuna Tuning** (`objective` function):

- Suggests a search space for the number of LSTM layers, units, dropout rate, learning rate, etc.
- Builds and trains the model with early stopping, learning rate reduction, and a pruning callback.
- Returns the best validation `MAE` found.
- This tuning is done in `study.optimize(..., n_trials=50)`.

3. **Train the best model**:

- Retrieves the best hyperparameters from the study, builds a new model, and trains it again for up to 300 epochs (with patience).

## 3.6 Evaluation and Saving the Model (Steps 10–11)

1. **Evaluation**:

- Predicts on test data.
- Inverse-transforms these predictions using `scaler_target`.
- Computes MSE, RMSE, MAE, $R\hat{2}$, and a directional accuracy metric.
- Saves a plot of Actual vs. Predicted.
- Prints a table of the first 40 predictions for inspection.

2. **Saving**:

- Saves the trained model as `optimized_lstm_model.h5`.
- Saves the fitted `MinMaxScalers` via `joblib`.

## 3.7 Reinforcement Learning Environment and DQN Training (Step 12)

- **StockTradingEnv**:
  - Inherits from `gym.Env`.
  - Observations are a combination of scaled technical features and some additional info (balance ratio, shares held, cost basis).
  - Action space has three discrete actions: `Sell (0)`, `Hold (1)`, and `Buy (2)`.
  - `step(action)` adjusts balance, shares, and net worth accordingly and returns a reward of `net_worth - initial_balance`.

- **train_dqn_agent(env)**:
  - Creates a DQN model with MLP policy from `stable_baselines3`.
  - Trains for 100,000 timesteps.
  - Saves the resulting model to disk as `dqn_stock_trading`.

# 4 Replacing DQN with PPO

In many cases, **Proximal Policy Optimization (PPO)** can be more robust or sample-efficient than DQN for continuous or more complex action spaces (and sometimes even for discrete spaces). It uses a different approach to update the policy in a clipped fashion, preserving stability. Here is a step-by-step guide on how to replace the current DQN approach with PPO in your code:

## 4.1 Potential Benefits of PPO over DQN

1. **Policy Gradient Method**: PPO is an on-policy algorithm that directly optimizes a policy. This can be more stable in some environments.

2. **Clipping Mechanism**: PPO uses a clipping in the objective function that prevents large updates from one iteration to the next, making training more stable.

3. **Sample Efficiency**: In certain setups, PPO can reuse on-policy samples effectively, especially if you carefully manage the rollout steps.

## 4.2    Modifications in the Code

To make the switch, you only need to modify a few lines where DQN is used:

1. **Import PPO**:

```python
# From:
from stable_baselines3 import DQN

# To:
from stable_baselines3 import PPO
```

2. **Create the PPO Model Instead of DQN**:

```python
def train_ppo_agent(env):
    logging.info("Training PPO Agent...")
    try:
        model = PPO(
            'MlpPolicy',
            env,
            verbose=1,
            learning_rate=3e-4,    # typical default for PPO
            n_steps=2048,          # number of steps to run for each environment
            batch_size=64,
            gamma=0.99,
            gae_lambda=0.95,
            clip_range=0.2,
            ent_coef=0.0,
            tensorboard_log="./ppo_stock_tensorboard/"
        )
        model.learn(total_timesteps=100000)
        model.save("ppo_stock_trading")
        logging.info("PPO Agent trained and saved as 'ppo_stock_trading.zip'.")
        return model
    except Exception as e:
        logging.error(f"Error training PPO Agent: {e}")
        sys.exit(1)
```

3. **Instantiate the PPO Training**:

   - Instead of calling `train_dqn_agent`, call `train_ppo_agent`:

```python
logging.info("Initializing and training PPO environment...")
trading_env = StockTradingEnv(data)
trading_env = DummyVecEnv([lambda: trading_env])

# Train
ppo_model = train_ppo_agent(trading_env)
```

Everything else regarding the environment remains the same. The main difference is in how the policy is updated and how the rollout buffer is handled internally.

# 5    Conclusion

By following the explanations in this document, you should now:

- Understand the structure and purpose of each code block in the LSTM + DQN trading workflow.

- Know how to easily switch from DQN to PPO by changing the import statements and the agent creation and training code.

- Appreciate some of the benefits (and potential pitfalls) of using PPO instead of DQN for reinforcement learning tasks in stock trading environments.

**Disclaimer**: Reinforcement learning models in algorithmic trading should be carefully tested under realistic constraints, including transaction costs, slippage, and risk management rules. Past performance does not guarantee future results.