

Scaling:

Normalization (min max scaler)

Maps features to [0, 1]

keeps shape but compresses outliers. use when you need bounded inputs · when absolute range matters.

Standardization (standard scaler)

→ transforms features to zero mean and unit variance  
→ often preferred for neural nets because it centers and scales inputs so gradient descent behaves better.

walk forward:

fit only on training data to avoid data leakage.

when doing walk forward you must fit the scaler using only the current training window before transforming training inputs and the single forecast input. if you fit it on the entire series, future information leaks into training.

inverse transform for evaluation:

Always invert predictions back to the original scale before computing the RMSE, MAE, MAPE are avg economic metric. for returns you may want to reapply transformations to get price forecasts.

Multivariate features:

As we have the multiple features then we need to standardize each feature individually.

## walk forward Validation

two main patterns:

### 1. Expanding window (growing training set)

→ Start with train  $t=1 \dots T_0$  Predict  $T_{0+1}$

Next fold train till  $T_{0+1}$  to predict  $T_{0+2}$

⇒ uses all available prior information each step

→ as the train set  $\rightarrow$  can become expensive.

### 2. Rolling window (fixed window size)

→ train on window  $t=t_0 \dots t_{t_0+w-1}$  to predict  $t_{t_0+w}$   
then shift window forward.

⇒ keeps training size fixed, adapts to recent regime change

⇒ loses older information, needs careful window size tuning.

\* Scaler leakage fitting Scaler on the whole dataset

\* Lookahead features don't use columns built using future info.

\* insufficient training examples relative to lookback

\* Model reinitialization vs warm start. reinitializing each step is safer but expensive. warm-starting keep weights can be efficient but can leak information indirectly if you tune on future validation carefully.

## Choosing initial training size.

→ Should be large enough to allow the model to learn stable patterns. daily data 200+ may be minimal for high frequency instruments to adjust accordingly.

## Sequence building and lookback.

→ Lookback controls how many past time steps the model sees. typical values 10 - 60 daily data depending on pattern length (test multiple values).

## Batching:

→ use shuffle = false for training time series if you want each batch to preserve contiguous sequences.

→ Alternatively use randomized mini batches that respect sequence boundaries

## Hidden states & statefulness:

→ PyTorch, defaults to stateless LSTM. hidden state is not preserved across forward calls unless you manually keep it. for walk-forward training it is usual to treat each sequence independently.

→ If you use stateful LSTM, then must carefully reset hidden states between epochs and between training vs validation

## Regularization & optimization

- we can use dropout in case of overfitting
- monitor training loss and validation loss

warm-start vs full retrain:

- full retrain: each fold safer and avoids cumulative bias but expensive to run
- fine tune but to ensure hyperparameters and validation are not tuned on future data.

## Evaluation Metrics.

RMSE / MAE : measure magnitude of error

MAPE : problematic when actuals near zero.

feature engineering for time series:

→ lagged returns

→ rolling statistics, rolling mean, rolling std, skewness

using a coarser re-training frequency then still make daily predictions using the latest model.

→ using a rolling window with warm start

limit model complexity for frequent retraining  
parallelize folds where possible.