

transformers for time series.

→ traditional deep learning methods for time series forecasting rely heavily on recurrent structures.

RNN → reads sequence step by step

LSTM → remembers longer forget noise

GRU → simplified memory unit

The above all share the same limitation.

they all process data sequentially.

issue-1 Long range dependency collapse

if something in the sequence from 100 steps ago impacts today.

→ An LSTM must carry that information through its recurrent updates.

→ Noise, randomness, small fluctuations slowly erode the memory.

→ Gradients shrink → older information disappears.

This makes LSTM short-sighted especially in financial markets.

issue-2 parallelism bottleneck

RNN force PyTorch to process time steps one by one which is slow and inefficient

transformers remove recurrence completely and replace it with

Self Attention.

Any time step can directly connect to any other time step
this makes model

⇒ Global

⇒ Parallel

⇒ Faster

⇒ Better at long seq

⇒ Better Learning patterns.

This is why transformers dominate NLP, audio, visual, TTS

Let take an example of 50 time steps.

instead of passing the sequence to one layer

the transformers say Let me look at every time step and see how important all the others are.

this comparison happens through three learned vectors.

Query (Q) → what this position wants to find

key (K) → what information this position contains

Value (V) → content to be shared

so here every time step becomes the weighted combination of all other time steps.

the sequence This gives deep global understanding of

Multi-head attention:

instead of single attention mechanism transformed
create multiple heads.

each head.

- => has its own Q , K , V projections.
- => looks at the sequence differently
- => capture different kind of relationships

one head might track.

- => short-term spikes
- > long-term cycles
- => volatility clusters
- > momentum shifts
- > cross-feature interactions

transformer encoder block.

Each encoder block has two major components.

A) Multi-head self-attention.

This part responsible for understanding relationships it takes the entire sequence and directly builds global dependencies.

B) feed forward network (FFN)

After the attention has mined information across time each time step is passed through a small neural network.

=> This expands the representation.

=> Applies nonlinearity.

=> compresses back to model dimension.

=> Why this is better for financial time series.

transformer encoders shine because they capture long range interactions.

- the impact of a federal reserve announcement 30 days ago on today's stock move.

=> they can handle multivariate inputs naturally.

=> they process the entire window in parallel

=> they are less prone to forgetting.

Because LSTM forgets distant events
Attention directly connects to them.

D-Linear Model

Despite of power of deep models like transformer
a simple linear model often performs better.

use few linear layers, do not overfit the noise.