

LSTM's and Corn Price Elasticity derived from Supply Factors within the USA

Kenji Clair

Agenda

- Thesis
- Application in Quant Finance
- Long Term Goal
- Data
- Model Assumptions
- RNN capabilities
- What is an LSTM?
- Code
- Key Findings
- Model Issues

Weather patterns within the top 10 corn producing states in the US can be used to derive corn spot price fluctuations by using a Recurrent Neural Network Architecture to detect temporal relationships varying through time.

Thought Process

- The US produces 32% of the Worlds Corn.
- The top 10 corn producing states (primarily from the Midwest) are responsible for all of this corn production.
- Linear/polynomial regression would oversimplify relationship missing key details.

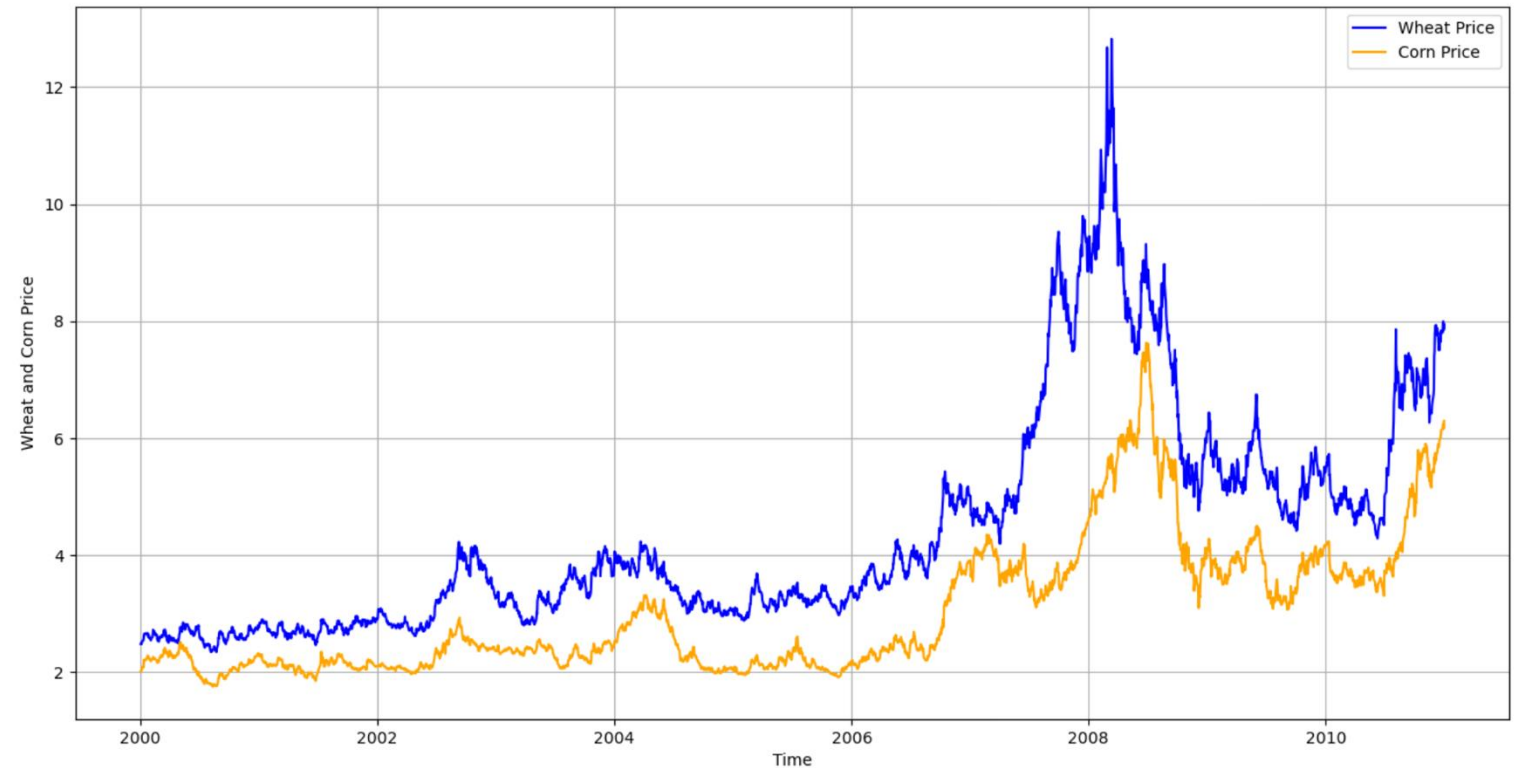


Pairs Trading

- Find detect the divergence of corn and wheat.
- When corn prices diverge from wheat, short the overvalued asset and go long on the undervalued asset.
- Assets prices will converge in the future.

Agriculture Indexes

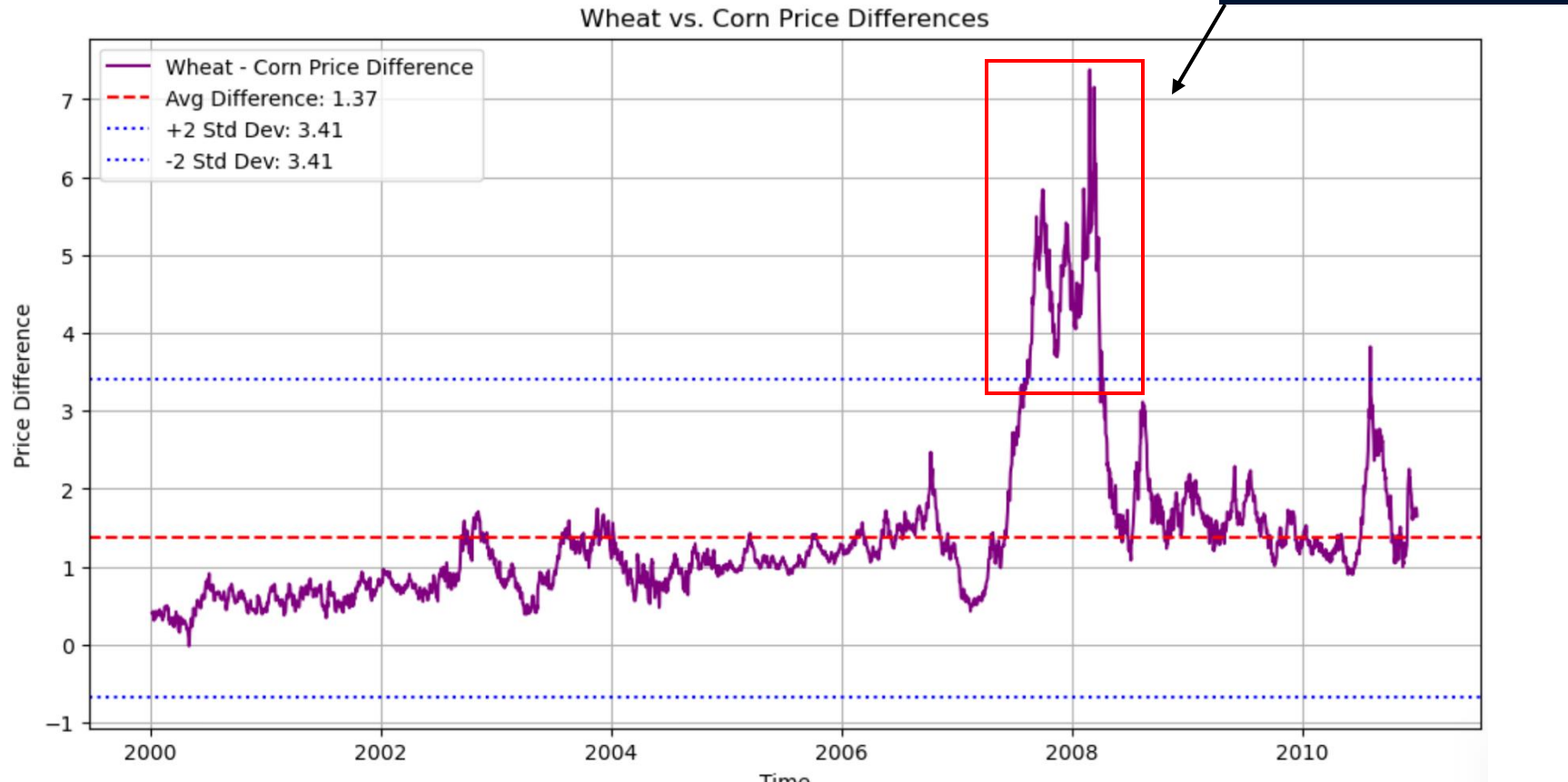
- S&P GSCI Agriculture Index
- S&P/ASX Agribusiness Index



Ultimate Long-Term Goal

Detect price divergence of wheat and corn

When the difference between wheat and corn passes confidence interval threshold ($\pm 2SD$), it is counted as an anomalous or an outlier in its prevailing pattern.



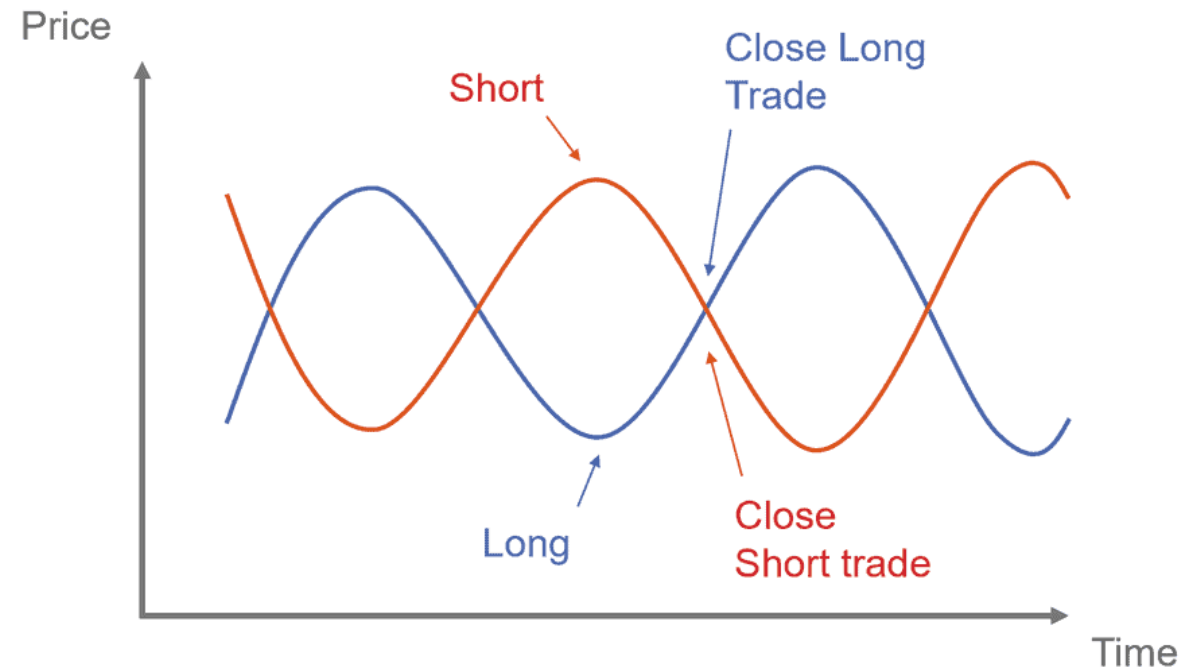
| | time | wheat_price | corn_price | lat | lon | max_day_temp | effective_degree_days | ice_days | heavy_rain_days |
|--------|------------|-------------|------------|-------|---------|--------------|-----------------------|----------|-----------------|
| 0 | 2000-01-05 | 2.4975 | 2.0300 | 36.75 | -103.75 | 293.35178 | 6.794754 | 1.0 | 0.0 |
| 1 | 2000-01-05 | 2.4975 | 2.0300 | 36.75 | -103.25 | 295.86730 | 10.660324 | 1.0 | 0.0 |
| 2 | 2000-01-05 | 2.4975 | 2.0300 | 36.75 | -102.75 | 297.86926 | 12.471822 | 1.0 | 0.0 |
| 3 | 2000-01-05 | 2.4975 | 2.0300 | 36.75 | -102.25 | 298.21940 | 13.270727 | 1.0 | 0.0 |
| 4 | 2000-01-05 | 2.4975 | 2.0300 | 36.75 | -101.75 | 298.43073 | 14.187780 | 2.0 | 0.0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 311370 | 2010-12-15 | 7.6475 | 5.8425 | 48.75 | -82.75 | 274.35205 | 0.000000 | 9.0 | 0.0 |
| 311371 | 2010-12-15 | 7.6475 | 5.8425 | 48.75 | -82.25 | 274.23660 | 0.000000 | 9.0 | 0.0 |
| 311372 | 2010-12-15 | 7.6475 | 5.8425 | 48.75 | -81.75 | 274.08870 | 0.000000 | 9.0 | 0.0 |
| 311373 | 2010-12-15 | 7.6475 | 5.8425 | 48.75 | -81.25 | 274.41687 | 0.000000 | 9.0 | 0.0 |
| 311374 | 2010-12-15 | 7.6475 | 5.8425 | 48.75 | -80.75 | 274.77618 | 0.000000 | 9.0 | 0.0 |

311375 rows x 9 columns

Equation Representation: $corn_price = wheat_price + lat + lon + max_day_temp + effective_degree_days + ice_days + heavy_rain_days + bias$

Model Assumptions

- Supply factors influence pricing of commodities
- Certain relationships between securities tend to revert to their mean or exhibit predictable patterns over time
- The top 10 US corn producing states have a substantial outcome on corn prices
- Severe weather events influence corn supply
- Agriculture is an efficient market



Recurrent Neural Network Capabilities

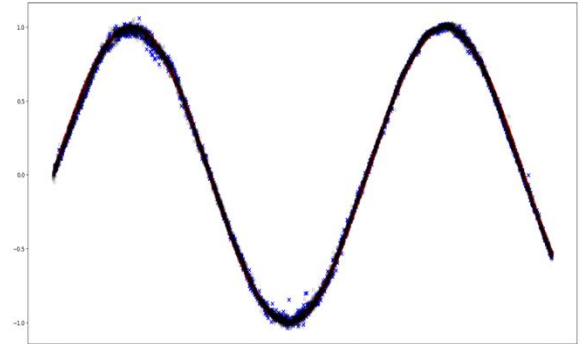
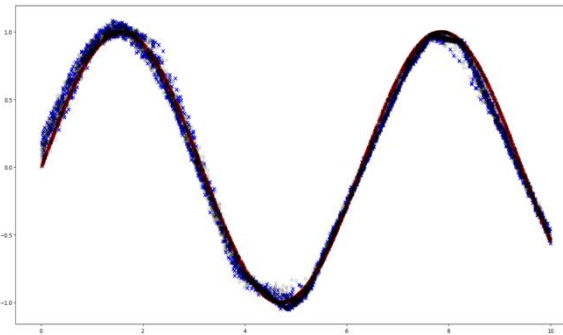
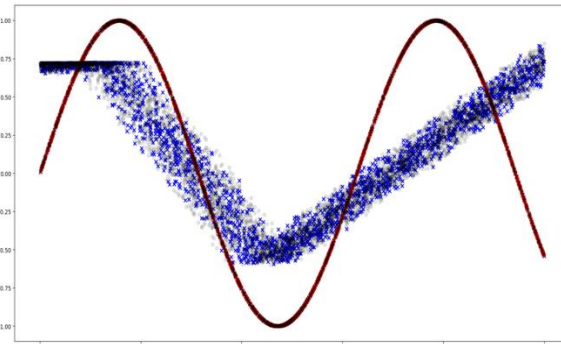
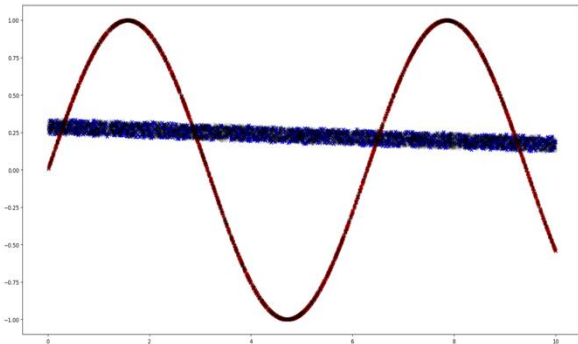
- Universal Function Approximator
- Through many training iterations and a large dataset, it can very accurately detect nonlinear patterns

Training iterations with nonlinear activation layers

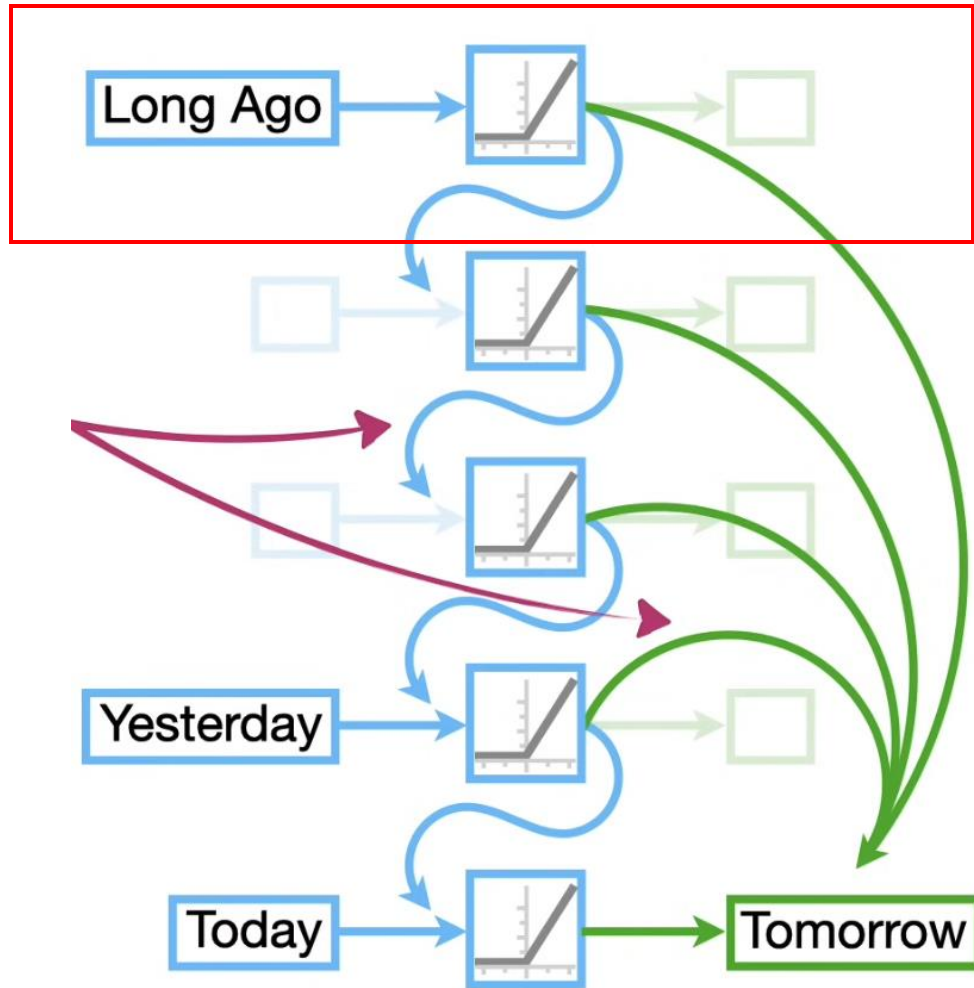
Start



Result



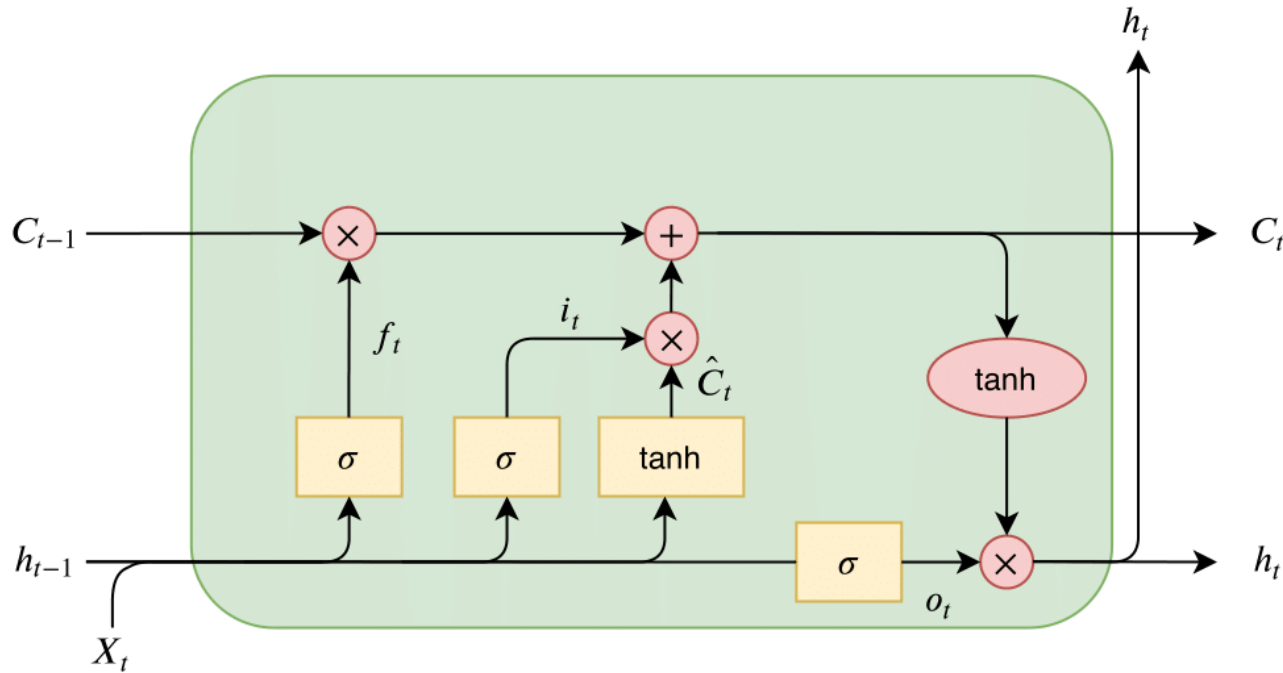
How does an LSTM work?



Type of Recurrent Neural Network

Utilizes a gated architecture to determine relevance of long term and short-term data

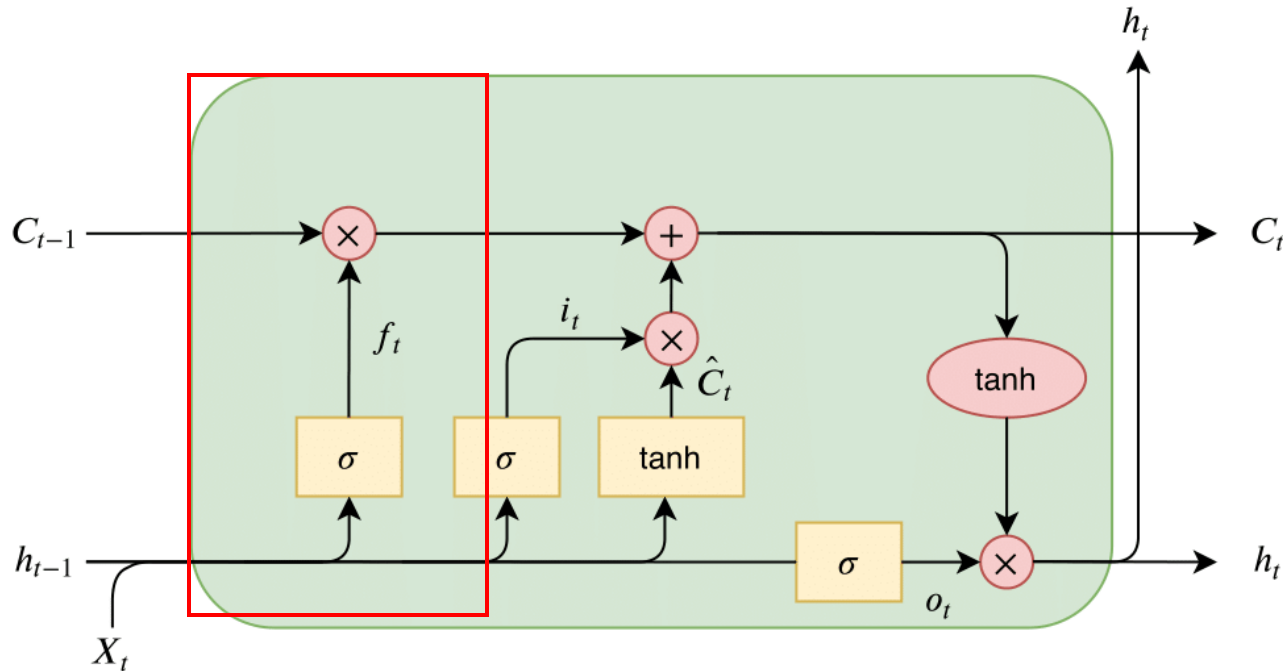
Output today can depend on data from long ago or data from yesterday



Tanh

$$f(x) = \frac{(e^x - e^{-x})}{(e^x + e^{-x})} = \text{range}[-1,1]$$

$$\sigma(x) = \frac{1}{1 + e^{-x}} = \text{range}[0,1]$$

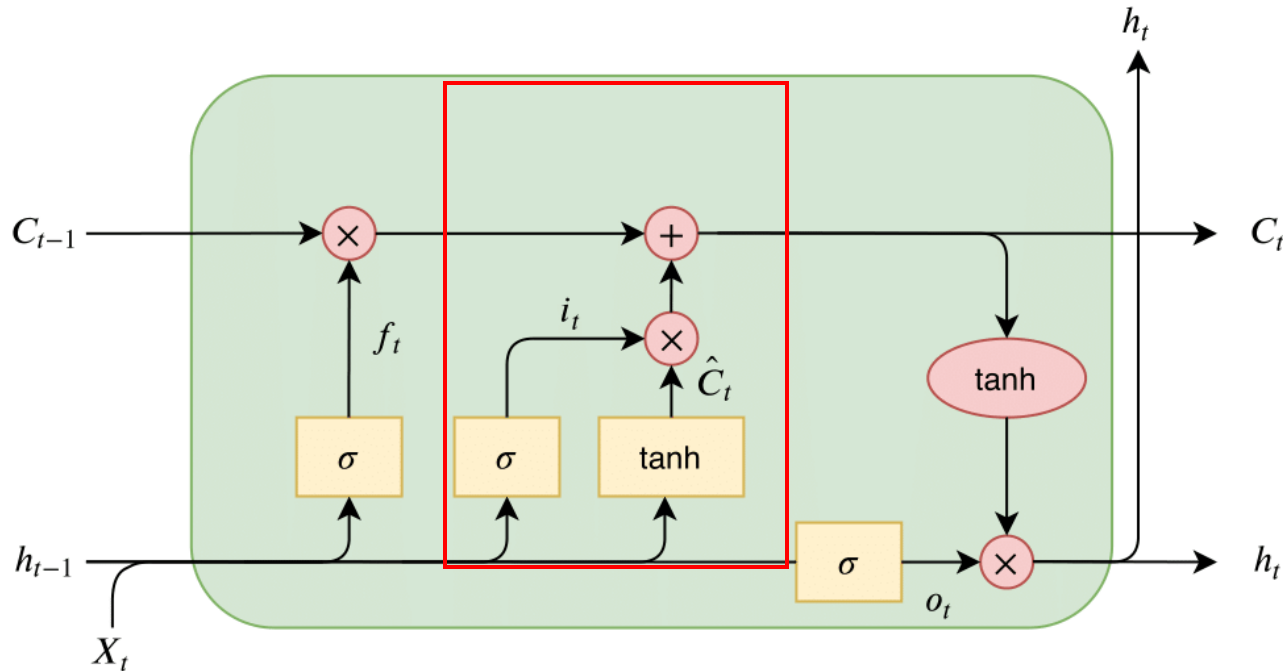


Forget gate – Short term memory determines percentage of long term memory relevant
To make new prediction

Tanh

$$f(x) = \frac{(e^x - e^{-x})}{(e^x + e^{-x})} = \text{range}[-1,1]$$

$$\sigma(x) = \frac{1}{1 + e^{-x}} = \text{range}[0,1]$$

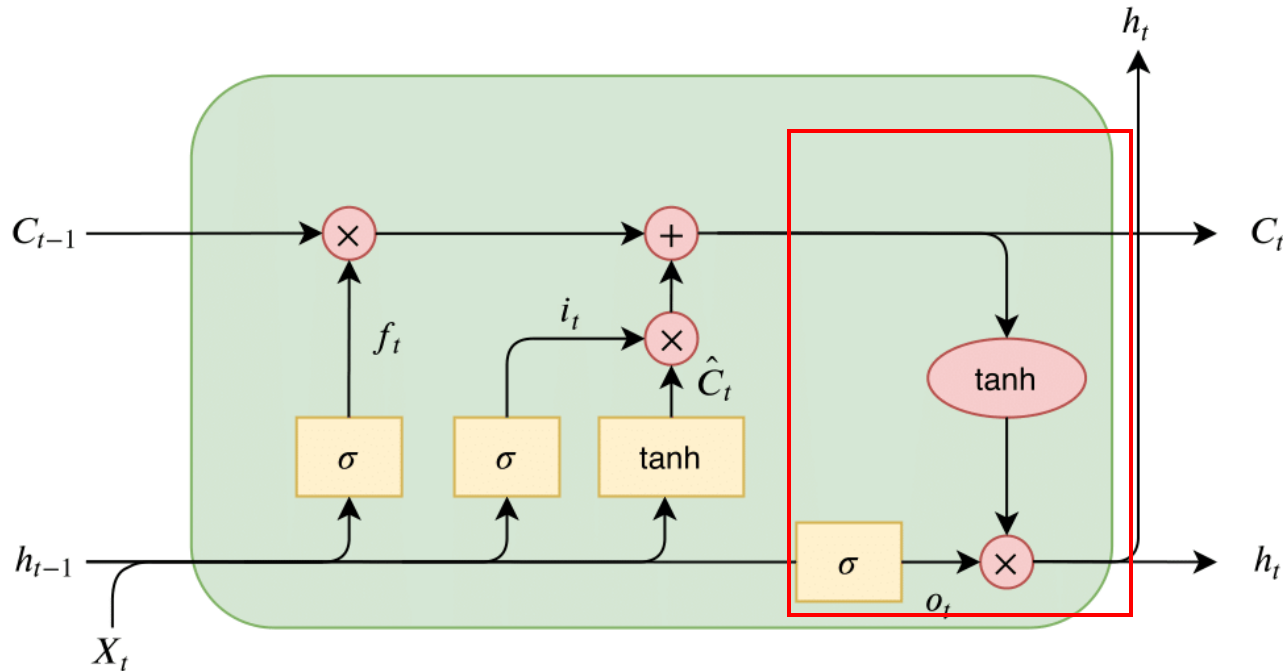


Input gate – Short term memory determines newly updated long-term memory

Tanh

$$f(x) = \frac{(e^x - e^{-x})}{(e^x + e^{-x})} = \text{range}[-1,1]$$

$$\sigma(x) = \frac{1}{1 + e^{-x}} = \text{range}[0,1]$$



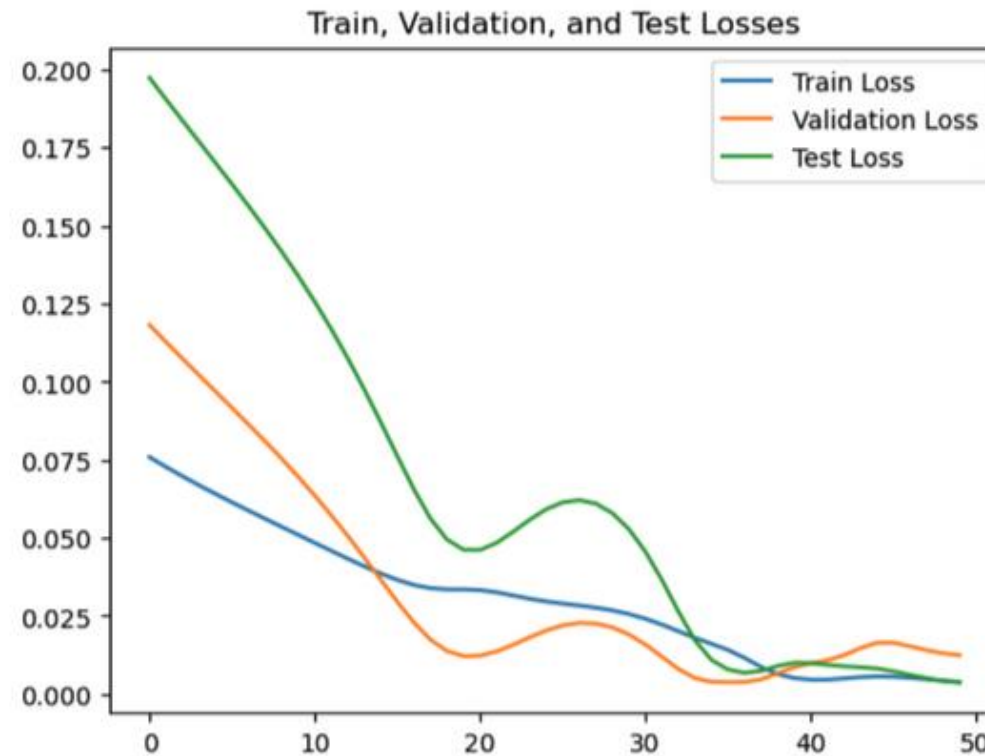
Output gate – Long term memory determines newly updated short-term memory (which is the output)

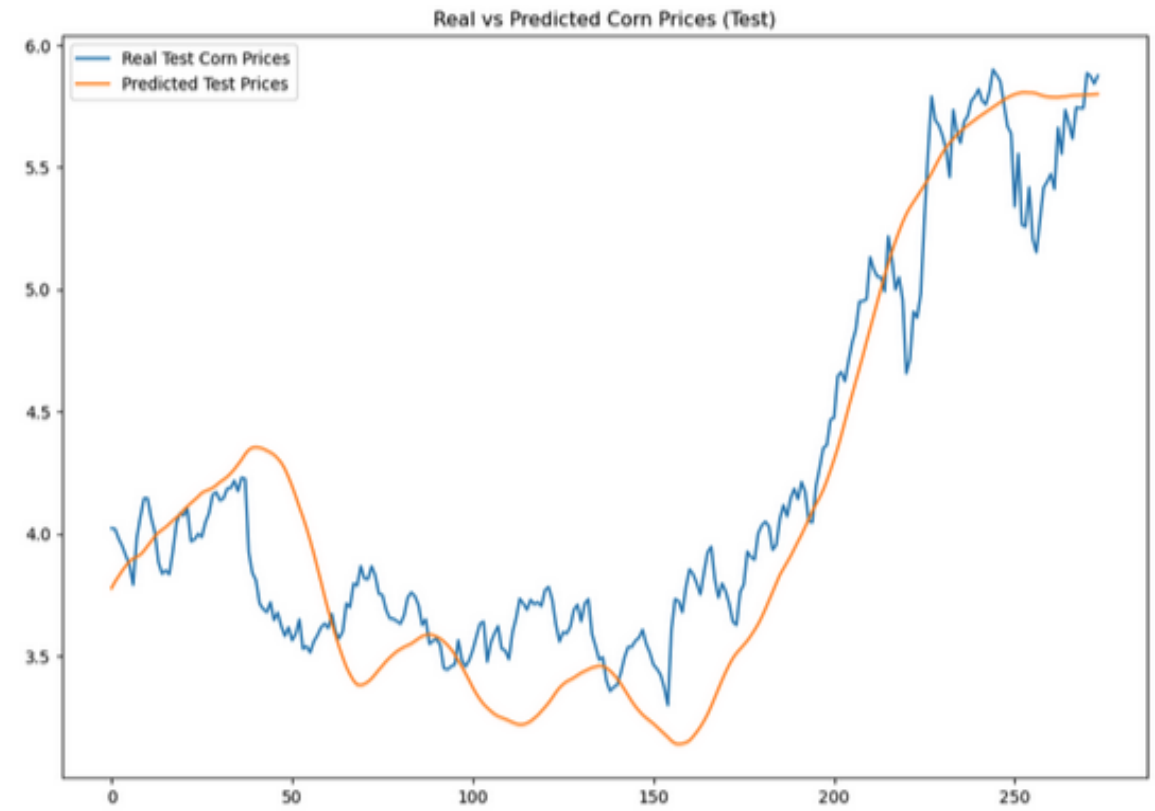
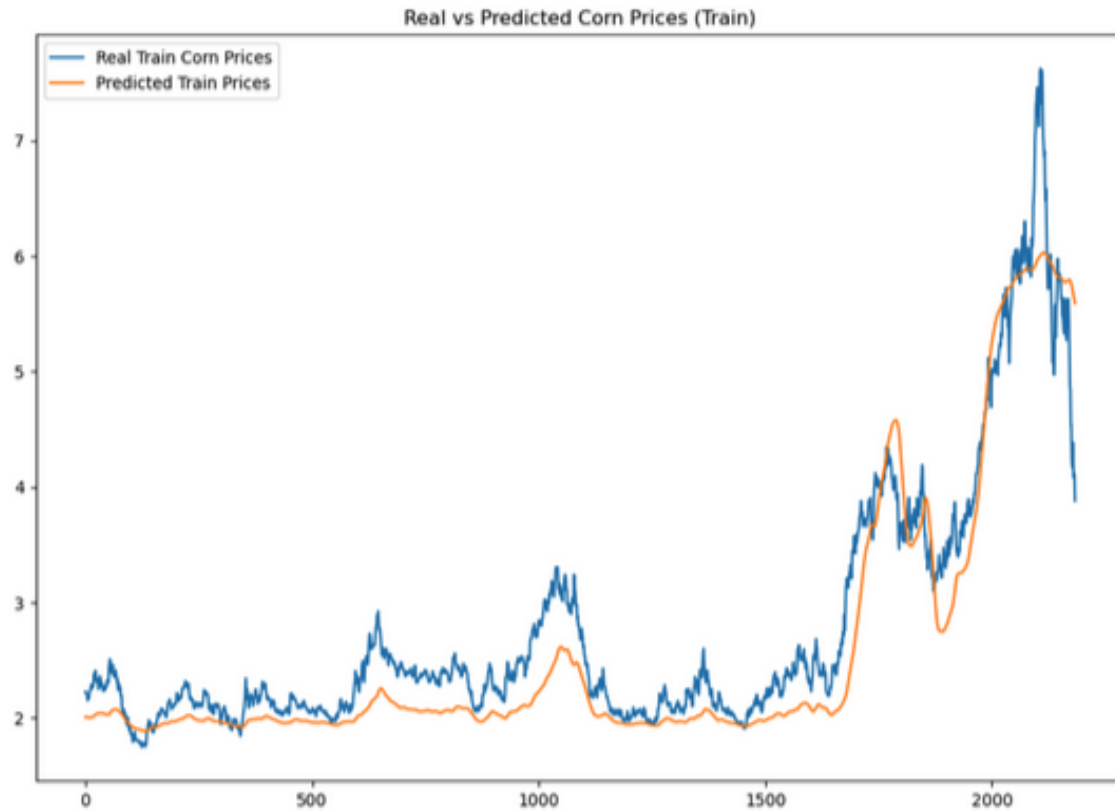
Tanh

$$f(x) = \frac{(e^x - e^{-x})}{(e^x + e^{-x})} = \text{range}[-1,1]$$

$$\sigma(x) = \frac{1}{1 + e^{-x}} = \text{range}[0,1]$$

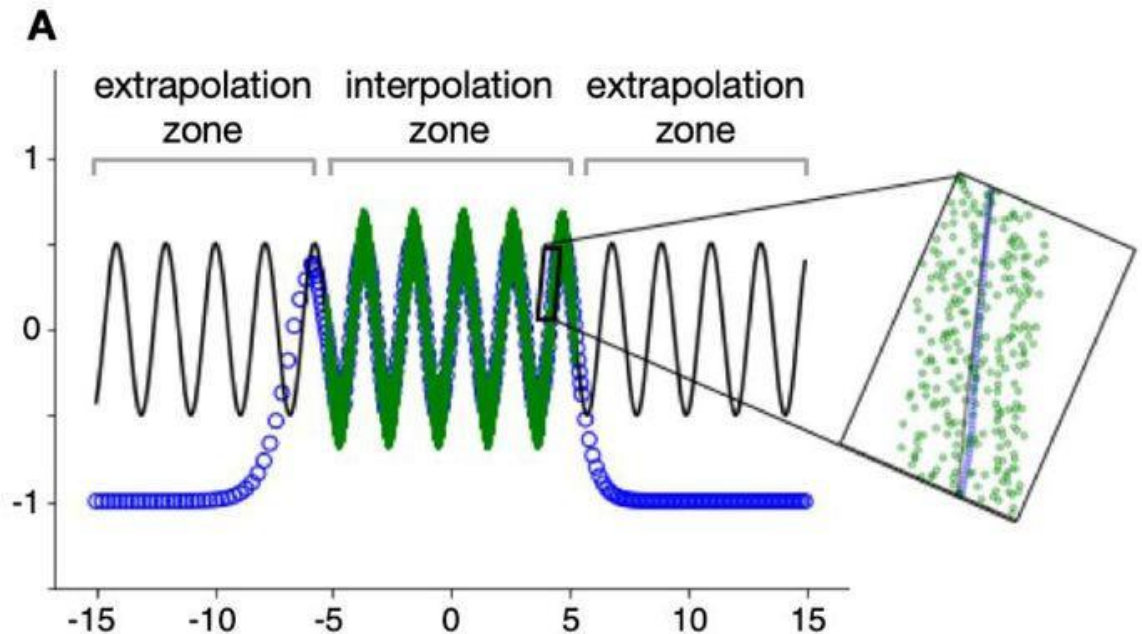
- Convergence of train, test and validation RMSE through 50 training epochs
- Low model error (accurate predictions)
- Not too overfit but predictions are slightly lagged



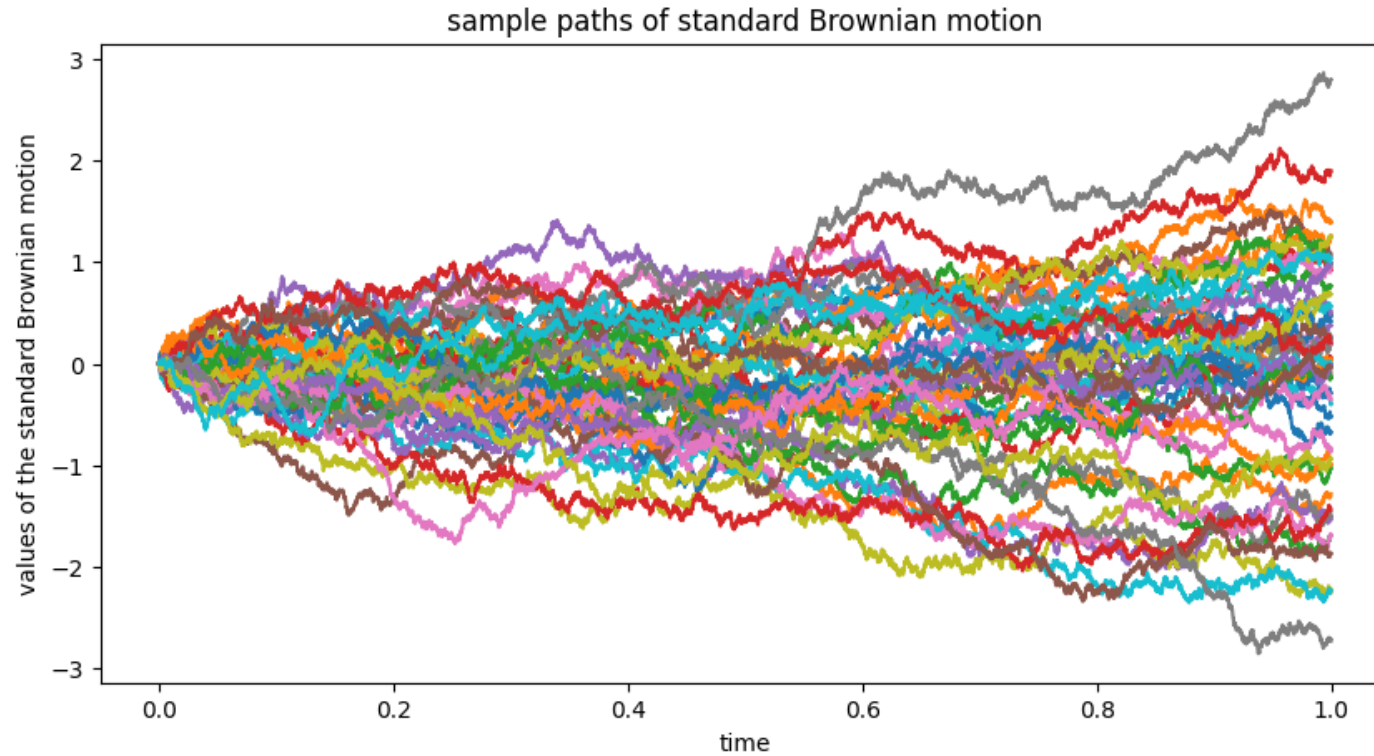


Key Findings

- Model can very accurately predict corn prices between 2000-2010
- If overfit, not a problem (can be tested for)
- **Evaluation Metric - Root Mean Square Error (RMSE)**
- Low test, validation and training RMSE (implies model accuracy)
- Convergence of train, test, validation RMSE (not overfit)



Violation of Brownian Motion



AB Testing to compare models: Select corn prices through Wiener process sampling within price range -> compared to LSTM output to see which model is more accurate

To prevent overfitting:

- Further fine tune model
- Incorporate more data that accounts for other supply factors other than weather (business cycles, market conditions)
- Incorporate wheat price prediction model to predict divergence more accurately
- Need to account for lagged predictionsa

Thank You

