

TỔNG LIÊN ĐOÀN LAO ĐỘNG VIỆT NAM  
TRƯỜNG ĐẠI HỌC TÔN ĐỨC THẮNG  
KHOA CÔNG NGHỆ THÔNG TIN



BÁO CÁO KẾT THÚC MÔN: KIẾN TẬP CÔNG NGHIỆP

# TÌM HIỂU VỀ TIME SERIES FORECASTING TRONG DEEP LEARNING

*Người hướng dẫn:* Nguyễn Thanh Hiên

*Người thực hiện:* Trần Quốc Lĩnh - 51703124

Lớp: 17050301

Khoá: 21

THÀNH PHỐ HỒ CHÍ MINH, NĂM 2020

## LỜI CẢM ƠN

Cảm ơn <bằng cấp> thầy Nguyễn Thanh Hiên đã hướng dẫn và tận tình chỉ bảo em trong suốt thời gian em thực tập ở công ty TNHH Tin học Đại Phát, cũng như trong khoảng thời gian em thực hiện bài báo cáo này.

## **BÀI BÁO CÁO ĐƯỢC HOÀN THÀNH TẠI TRƯỜNG ĐẠI HỌC TÔN ĐỨC THẮNG**

Em xin cam đoan đây là sản phẩm nghiên cứu của riêng em. Các nội dung nghiên cứu, kết quả trong đề tài này là trung thực và chưa công bố dưới bất kỳ hình thức nào trước đây. Những số liệu, hình ảnh được chính em thu thập từ các nguồn khác nhau có ghi rõ trong phần tài liệu tham khảo.

Ngoài ra, trong bài tiểu luận còn sử dụng một số nhận xét, đánh giá cũng như số liệu của các tác giả khác, cơ quan tổ chức khác đều có trích dẫn và chú thích nguồn gốc.

Nếu phát hiện có bất kỳ sự gian lận nào em xin hoàn toàn chịu trách nhiệm về nội dung bài tập lớn của mình. Trường đại học Tôn Đức Thắng không liên quan đến những vi phạm tác quyền, bản quyền do em gây ra trong quá trình thực hiện (nếu có).

*TP. Hồ Chí Minh, ngày 20 tháng 07 năm 2020*

*Tác giả*

*(Đã ký)*

*Trần Quốc Lĩnh*

# PHẦN XÁC NHẬN VÀ ĐÁNH GIÁ CỦA GIẢNG VIÊN

Phần xác nhận của GV hướng dẫn

---

---

---

---

---

---

Tp. Hồ Chí Minh, ngày 20 tháng 07 năm 2020  
(*kí và ghi họ tên*)

Phần đánh giá của GV giám sát

---

---

---

---

---

---

Tp. Hồ Chí Minh, ngày 20 tháng 07 năm 2020  
(*kí và ghi họ tên*)

## TÓM TẮT

Đây là một bài báo cáo, một bài tóm lược tổng quát và cơ bản nhất về time series forecasting trong deeplearning. Kèm theo đó bài báo cáo này cũng trình bày các phương pháp, các bài toán thực tế sử dụng time series forecasting.

## MỤC LỤC

|  |    |
|--|----|
| LỜI CẢM ƠN .....   | 2  |
| CAM KẾT .....  | 3  |
| ĐÁNH GIÁ CỦA GIÁO VIÊN .....   | 4  |
| TÓM TẮT .....  | 5  |
| MỤC LỤC .....  | 6  |
| DANH MỤC CHÚ THÍCH CÁC THUẬT NGỮ VÀ HÌNH ẢNH .....                           | 7  |
| CHƯƠNG I: BỐI CẢNH HIỆN TẠI VÀ TỔNG QUAN<br>VỀ TIME SERIES FORECASTING ..... | 8  |
| Time series .....  | 9  |
| Time series forecasting .....  | 10 |
| Những phương pháp được sử dụng trong<br>time series forecasting .....        | 10 |
| CHƯƠNG II: MỘT SỐ ỨNG DỤNG THỰC TẾ .....                                     | 13 |
| Dự đoán mức tăng trưởng của một loại cổ phiếu .....                          | 13 |
| Dự báo thời tiết .....   | 14 |
| Dự đoán doanh thu .....  | 15 |
| CHƯƠNG III: VÍ DỤ MINH HỌA DỰ BÁO THỜI TIẾT .....                            | 17 |
| CHƯƠNG IV: TỔNG KẾT .....  | 32 |
| TÀI LIỆU THAM KHẢO .....   | 33 |

# DANH MỤC CHÚ THÍCH

## Hình ảnh

Hình 1: Bảng thống kê sự tăng trưởng cổ phiếu (hình minh họa)

Hình 2: Biểu đồ ghi lại sự thay đổi nhiệt độ của hai thành phố Tokyo và Lodon qua các tháng.

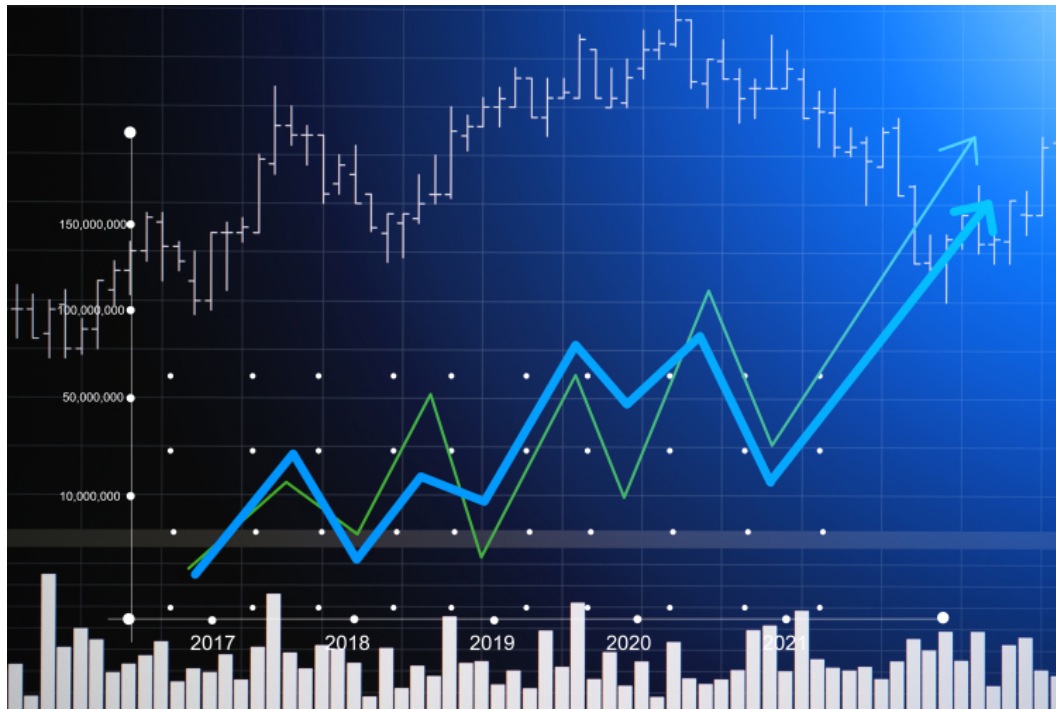
Hình 3: Bảng dự đoán mức tăng trưởng của cổ phiếu (hình minh họa)

Hình 4: Bảng minh họa kết quả dự báo thời tiết (hình minh họa)

Hình 5: Bảng dự đoán hạn mức doanh thu theo từng năm (hình minh họa)

# CHƯƠNG I: BỐI CẢNH HIỆN TẠI VÀ TỔNG QUAN VỀ TIME SERIES FORECASTING

Từ trước đến nay, việc xem xét, đánh giá những khả năng có thể xảy ra và tác động của nó đến các công ty, tổ chức, bộ máy là một điều hết sức quan trọng. Điều này là nhân tố quyết định cho sự tồn tại và phát triển của các cơ quan, đoàn thể. Dĩ nhiên nó đòi hỏi có một sự tỉ mỉ, một cái nhìn bao quát, một sự tính toán trong nhiều năm trời, và chỉ dành cho những chuyên gia thâm niên, những nhà nghiên cứu tài ba trong chính lĩnh vực đó. Một công việc không ít tốn kém cả về sức người lẫn sức của.



Hình 1: Bảng thống kê sự tăng trưởng cổ phiếu (hình minh họa)

Bài toán về việc sử dụng máy tính để tính toán, dự đoán tương lai, dự đoán kết quả có thể xảy ra dựa trên những dữ kiện đã ghi nhận từ quá khứ đến hiện tại, không phải là một điều gì mới mẻ. Nhưng trước đó do hiệu quả và khả năng dự đoán kém, nên bài toán này không được cộng đồng xã hội chú ý đến. Tuy nhiên, trong những năm trở lại đây, như một quả bom, nó đã bùng nổ và lan rộng khắp thế giới, đón nhận được sự quan tâm nhiệt tình của không chỉ cộng đồng khoa học

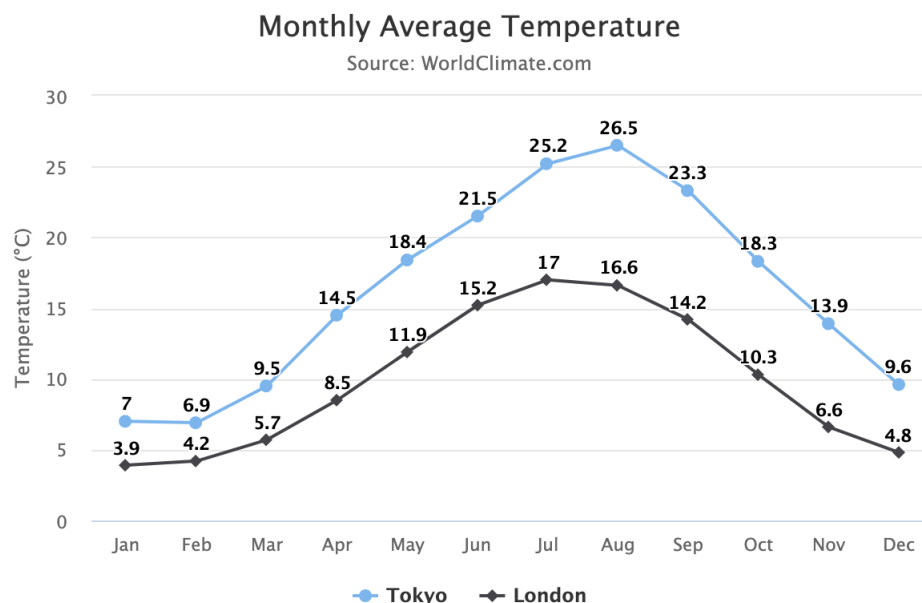


mà còn cả các ngành nghề, lĩnh vực khác (kinh tế, y học, khí tượng thủy văn,...). Ấy là nhờ vào sự xuất hiện của deep learning và việc ứng dụng deep learning vào việc tính toán, đánh giá, thay đổi và cải thiện tính chính xác của bài toán.

## Time series

Một time series là một chuỗi dữ liệu được đánh chỉ mục (một danh sách hoặc đồ thị) theo thứ tự thời gian. Thông thường các giá trị dữ liệu được xác định/đo đạc tại các mốc thời gian cách đều nhau. Chính vì thế, time series là một chuỗi dữ liệu rời rạc. Cụ thể là các thông tin về độ ẩm, hướng gió, lực gió, lượng mây,... được đo đạc mỗi giờ tại các trạm khí tượng thủy văn.

Chuỗi thời gian thường được biểu diễn dưới dạng biểu đồ đường. Chúng được sử dụng trong thống kê, xử lý tín hiệu, nhận dạng mẫu, kinh tế, tài chính, toán học, dự báo thời tiết, dự báo động đất, điện não đồ, kỹ thuật điều khiển, thiên văn học, kỹ thuật truyền thông và trong bất kỳ lĩnh vực khoa học ứng dụng - kỹ thuật nào liên quan đến các phép đo thời gian.



Hình 2: Biểu đồ ghi lại sự thay đổi nhiệt độ của hai thành phố Tokyo và London qua các tháng

Phân tích chuỗi thời gian nhằm mục đích phát hiện và tập hợp lại các yếu tố có ảnh hưởng đến giá trị của đối tượng quan sát theo thời gian.

Trong Time-series Data, có hai loại chính:

- Chuỗi thời gian thông thường (regular time series), loại thông thường được gọi là số liệu.
- Chuỗi thời gian bất thường (events) là những sự kiện.

## **Time series forecasting**

Như đã trình bày ở trên, dự đoán dựa trên chuỗi thời gian là việc dùng những dữ liệu đã có để dự đoán sự kiện hoặc một chuỗi sự kiện có thể xảy đến trong tương lai. Tuy nhiên, dự đoán thời gian là một việc khó khăn. Không như các bài toán phân lớp và hồi quy, bài toán chuỗi thời gian làm tăng độ phức tạp của thứ tự các quan sát hay sự phụ thuộc giữa các quan sát.

Người ta cũng đã dành không ít thời gian để xây dựng và điều chỉnh các mô hình, các phương pháp dự đoán thời gian sao cho sự kiện được dự đoán sai lệch thấp nhất có thể. Các phương pháp được trình bày bên dưới đây là những phương pháp phổ biến và có hiệu quả tốt nhất đã được kiểm chứng bằng thực nghiệm

## **Những phương pháp được sử dụng trong time series forecasting**

### **Phương pháp truyền thống**

#### **Autoregressive integrated moving average (ARIMA)**

Mô hình trung bình trượt, đồng liên kết, tự hồi quy ARIMA sử dụng chính các giá trị trong quá khứ của Y để giải thích cho Y hiện tại. Vì việc xác định mỗi một mô hình dựa trên sự phân tích dữ liệu cụ thể cho từng trường hợp mà không phụ thuộc hay bám lấy bất kỳ một lý thuyết tổng quát nào. Nên ARIMA có được tính chất linh hoạt và tiết kiệm chi phí.

### **Hạn chế, rào cản của phương pháp truyền thống**

- **Chỉ hiệu quả với dữ liệu hoàn chỉnh:** Một vài giá trị bị thiếu hụt có thể ảnh hưởng nặng nề đến mô hình.

- **Phụ thuộc vào mối quan hệ tuyến tính:** Các mô hình truyền thống không thích hợp với phân phối phức tạp.
- **Chỉ dùng được với dữ liệu đơn biến:** Trong khi thực tế, dữ liệu đầu vào hầu hết là đa biến.
- **Không chạy tốt với những dự đoán có chu kỳ dài:** Chúng thường chỉ tập trung vào việc dự đoán duy nhất một sự kiện (one-step).
- **Thời gian phụ thuộc phải là cố định:** Mối quan hệ giữa các quan sát thường luôn khác nhau tại các thời điểm khác nhau, do đó cần có sự cố định giữa các thời gian quan sát.

## Các phương pháp mới sử dụng deep learning

### 1. Multilayer Perceptrons (MLP)

MLP có những đặc trưng nổi trội như:

- Khả năng xấp xỉ các hàm phi tuyến tính. Mang lại một kết quả đầy hứa hẹn.
- Khả năng chịu lỗi. Học tốt kể cả khi bị thiếu sót dữ liệu.
- Có thể xử lý dữ liệu đầu vào là đa biến.
- Có khả năng dự đoán nhiều bước (multi-step).

### 2. Convolutional Neural Networks (CNN)

CNN thường được dùng trong việc phân lớp hình ảnh, bởi hiệu quả tuyệt vời của nó. Bên cạnh đó CNN cũng có khả năng trích xuất đặc trưng của dữ liệu đầu vào trong time series forecasting. CNN còn sở hữu đầy đủ những ưu điểm của MLP.

### 3. Recurrent Neural Networks(RNN)

RNN cung cấp khả năng xử lý trật tự của các quan sát và cho phép quan sát thể hiện sự phụ thuộc theo gian của nó. Khác với MLP, RNN sử dụng chính bộ nhớ của mình để xử lý chuỗi dữ liệu đầu vào. Nhờ đó mà RNN có thể xử lý các quan sát có độ trễ cố định về thời gian.

## So sánh, đánh giá các phương pháp

| Phương pháp  | Ưu điểm  | Hạn chế   |
|--------------|--|---|
| Truyền thống | <ul style="list-style-type: none"> <li>- Đơn giản</li> <li>- Dễ hiểu</li> </ul>  | <ul style="list-style-type: none"> <li>- Đòi hỏi dữ liệu hoàn chỉnh</li> <li>- Chỉ có mối quan hệ tuyến tính</li> <li>- Hiệu quả với đơn biến</li> <li>- Kém hiệu quả với chu kỳ dài</li> </ul>                 |
| MLP          | <ul style="list-style-type: none"> <li>- Xấp xỉ phi tuyến</li> <li>- Chịu được nhiễu (Noise)</li> <li>- Xử lý được đầu vào đa biến</li> <li>- Dự đoán được nhiều bước</li> </ul> | <ul style="list-style-type: none"> <li>- Require meaningful mapping between inputs and outputs</li> <li>- Rely on lag observations</li> <li>- Static mapping function and fixed outputs &amp; inputs</li> </ul> |
| CNN          | <ul style="list-style-type: none"> <li>- Tất cả ưu điểm của MLP</li> <li>- Kết xuất đặc trưng</li> </ul>   | <ul style="list-style-type: none"> <li>- Không thể học temporal dependence</li> </ul>   |
| RNN          | <ul style="list-style-type: none"> <li>- Tất cả ưu điểm của CNN</li> <li>- Học được temporal dependence</li> </ul>   | <ul style="list-style-type: none"> <li>- Còn đang tranh cãi</li> </ul>  |

## CHƯƠNG II: MỘT SỐ ỨNG DỤNG THỰC TẾ

### Dự đoán mức tăng trưởng của một loại cổ phiếu

Đây là một lĩnh vực mà không phải ai cũng am hiểu, không phải ai cũng thích hợp để tham gia vào. Nó đòi hỏi không chỉ tài năng nghiên cứu thị trường của người chơi cổ phiếu mà còn là những hành động quyết đoán dựa trên các tính toán của mình.

Để có thể dự đoán được giá trị của cổ phiếu trong khoảng thời gian nào đó sắp tới, người đầu tư chuyên nghiệp cần phải điều tra vô số các yếu tố về thị trường, chính trị, xã hội,... các yếu tố này sẽ ảnh hưởng trực tiếp đến giá của cổ phiếu. Tuy nhiên các yếu tố vừa kể không hề đơn giản để thu thập và tính toán.



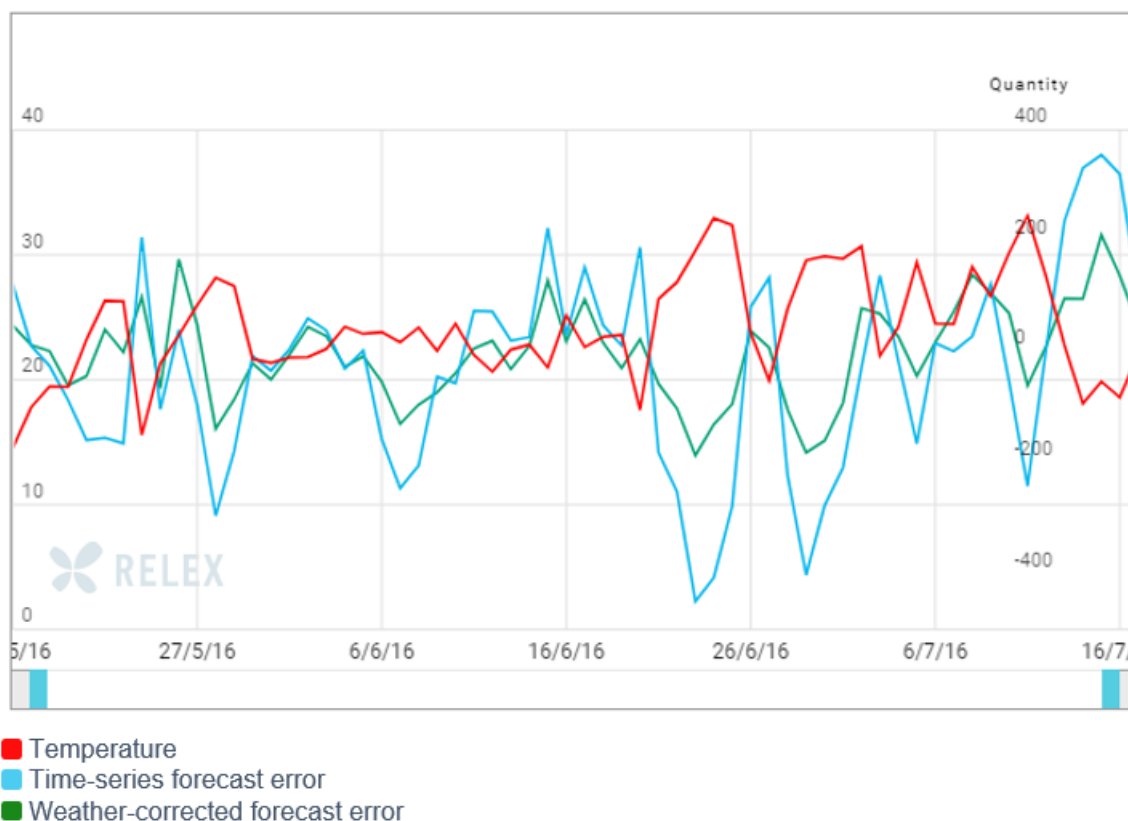
Hình 3: Bảng dự đoán mức tăng trưởng của cổ phiếu (Hình minh họa)

Nắm bắt được nhu cầu thực tế đó, người ta mới đưa, chỉnh sửa, phân tích các dữ liệu về thị trường, chính trị, xã hội sang dạng số liệu theo thời gian và áp dụng vào bài toán time series forecasting. Mặc dù vậy, có lẽ là do số liệu thu thập vẫn còn thiếu sót hoặc mô hình dự đoán vẫn chưa hoàn chỉnh nên kết quả dự đoán vẫn còn nhiều thiếu sót, không chiếm lấy được sự tin tưởng của người sử dụng. Tuy nhiên, các phương pháp mới vẫn được đang nghiên cứu sao cho kết quả ngày một chính xác hơn. Điều này mang đến một tương lai đầy hứa hẹn cho các nhà chơi cổ phiếu,

và cũng có thể là thu hút thêm nhiều người chơi cổ phiếu hơn nữa do tính đơn giản mà nó mang lại.

## Dự báo thời tiết

Dự báo thời tiết chính xác đem lại nhiều lợi ích cho không ít các ngành nghề, lĩnh vực trong xã hội. Theo phương pháp truyền thống, quá trình dự báo thời tiết thường trải qua các giai đoạn: Quan trắc, trao đổi số liệu, phân tích, dự báo và phân phát sản phẩm dự báo. Với dữ liệu đã có, người ta sẽ phải thực hiện hàng tỉ phép tính trước khi cho ra dữ liệu cuối cùng, do đó vô cùng tiêu tốn thời gian và chi phí. Đây cũng là quy trình cơ bản của các loại hình dịch vụ dự báo thời tiết.



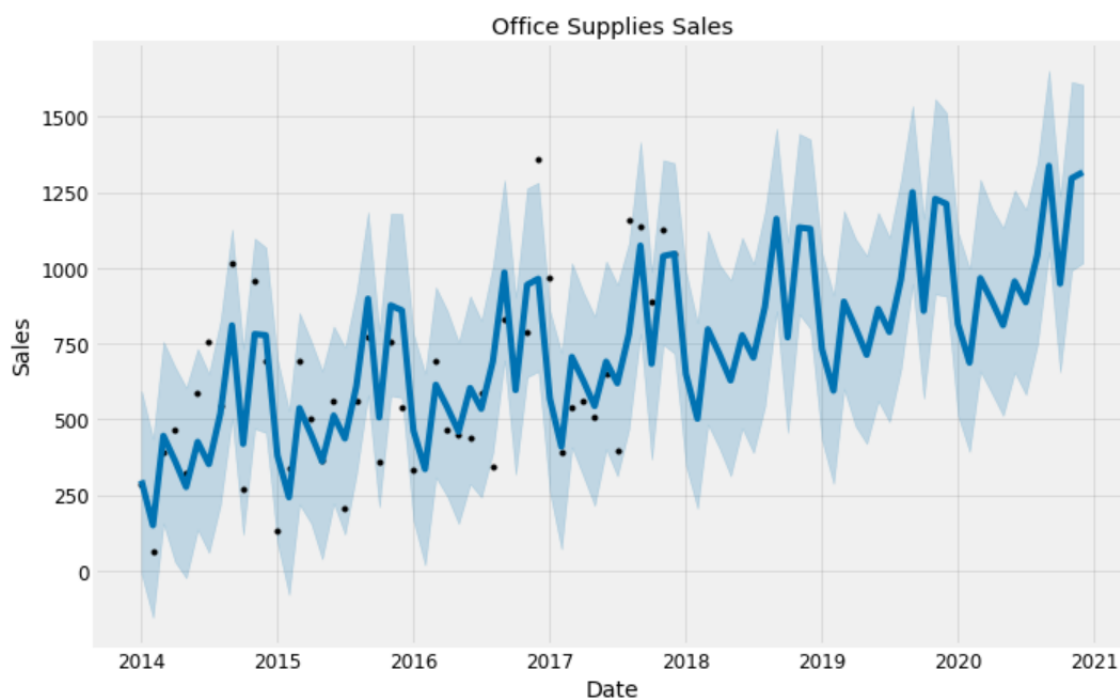
Hình 4: Bảng minh họa kết quả dự báo thời tiết (Hình minh họa)

Với việc thu thập, ghi lại các thông tin về trạng thái của khí quyển, nhiệt độ, độ ẩm, gió,... theo các mốc thời gian cụ thể. Đã đem lại tính tương thích cao cho bài toán series forecasting. Và đúng như mong đợi, hiệu quả mà time series forecasting

mang đến cho ngành dự báo thời tiết như một làn sóng mới về sự đơn giản, chính xác, tiết kiệm chi phí và nhanh chóng.

## Dự đoán doanh thu

Có thể nói việc tính toán doanh thu trong một sự kiện hoặc một giai đoạn thời gian sắp tới đối với các ngành nghề liên quan đến sản xuất, mua bán, dịch vụ,... là một trong những công việc đòi hỏi sự tỉ mỉ và khó khăn nhất. Bởi nó giúp nhà đầu tư biết rằng, liệu họ có nên đầu tư vào cái này thay vì cái kia không?, liệu họ có nên mở rộng cơ sở trong tháng sắp tới hay không? nếu mở thì nên mở thêm bao nhiêu cửa hàng? Tổng quát hơn, nó ảnh hưởng trực tiếp đến hành động và quyết định của nhà đầu tư, của người sở hữu cửa hàng, sở hữu doanh nghiệp,...



Hình 5: Bảng dự đoán hạn mức doanh thu theo từng năm (Hình minh họa)

Từ ngàn xưa, từ trước khi người kinh doanh bắt đầu sử dụng máy tính để dự đoán mức doanh thu của mình, để tính toán mức doanh thu họ thường xem xét các vấn đề loại mặt hàng nào đang yêu thích? đang bán chạy? Có bao nhiêu cửa hàng đang kinh doanh loại mặt hàng này trong khu vực? Họ có bán chạy không? Vô số câu hỏi như thế sẽ được đặt ra, và sau khi trả lời tất cả các câu hỏi đó, người ta

sẽ dự đoán (phần lớn các trường hợp là ước lượng) mức doanh thu mà họ thu được nếu người ta đầu tư vào loại mặt hàng, lĩnh vực đó. Nghe có vẻ đơn giản, nhưng để tính toán tốt, ít sai sót nhất thì số câu hỏi đặt ra phải lên tới hàng trăm, có khi hàng nghìn câu hỏi, và với mỗi câu hỏi người ta phải có câu trả lời dưới dạng số liệu (dĩ nhiên số liệu này phải gần với thực tế nhất có thể). Tuy nhiên phương pháp vừa nêu chỉ dành cho những doanh nghiệp khổng lồ, đa phần các doanh nghiệp vừa và nhỏ chỉ xem xét và đặt ra một vài câu hỏi trọng yếu, do đó mà các kết quả thu được cũng khác nhau.

Cho đến khi người ta bắt đầu nghĩ đến việc dùng máy tính để thay mình tính toán doanh thu, họ cũng dùng các phương pháp, các công cụ công nghệ cao thay họ thu thập dữ liệu, sau đó sắp xếp thành các chuỗi dữ liệu theo thời gian. Vâng, rất phù hợp với bài toán time series forecasting. Tương tự như trên, time series forecasting cũng cho ra các kết quả dự đoán có độ chính xác đến bất ngờ (ở một vài trường hợp hầu như trùng khớp). Nhờ vào điều này, không chỉ có các công ty, doanh nghiệp lớn mới có thể đưa ra các quyết định đúng đắn, tầm cỡ, mà còn mở rộng cơ hội cho các doanh nghiệp, công ty nhỏ lẻ khác. Góp phần cho sự phát triển chung của toàn xã hội.



# CHƯƠNG III: VÍ DỤ MINH HỌA DỰ BÁO THỜI TIẾT

## Chuẩn bị

### Môi trường

- Sử dụng python 3.8
- Trình soạn thảo Jupyter notebook
- Package: tensorflow, matplotlib, numpy, os, pandas

### Dữ liệu

- Tập dữ liệu thời tiết được cung cấp bởi [Max Planck Institute for Biogeochemistry](#).
- Tập dữ liệu này chứa 14 loại đặc điểm khác nhau của thời tiết như: nhiệt độ (temperature), áp suất không khí (atmospheric), độ ẩm (humidity),... Các giá trị này sẽ được ghi lại cứ mỗi 10 phút, bắt đầu từ năm 2009 đến 2016. Phần dữ liệu này đã được chỉnh sửa lại bởi François Chollet.

## Mục tiêu

Với dữ liệu đã có, sử dụng mô hình RNN để dự đoán nhiệt độ của 6 giờ kế tiếp qua hai trường hợp: đơn biến (univariate) và đa biến (multivariate)

## Code

### Khai báo các package cần thiết

```
1 import tensorflow as tf
2
3 import matplotlib as mpl
4 import matplotlib.pyplot as plt
5 import numpy as np
6 import os
7 import pandas as pd
8
9 mpl.rcParams['figure.figsize'] = (8, 6)
10 mpl.rcParams['axes.grid'] = False
```

## Tải dữ liệu

```
1 zip_path = tf.keras.utils.get_file(  
2 origin='https://storage.googleapis.com/tensorflow/tf-keras-  
   datasets/jena_climate_2009_2016.csv.zip',  
3 fname='jena_climate_2009_2016.csv.zip',  
4 extract=True)  
5 csv_path, _ = os.path.splitext(zip_path)  
6 df = pd.read_csv(csv_path)
```

## Xem thử bảng dữ liệu

```
1 df.head()
```

Nếu kết quả in ra có dạng như thế này, thì xem như việc tải dữ liệu thành công.

|   | Date Time              | p<br>(mbar) | T<br>(degC) | Tpot<br>(K) | Tdew<br>(degC) | rh<br>(%) | VPmax<br>(mbar) | VPact<br>(mbar) | VPdef<br>(mbar) | sh<br>(g/kg) | H2OC<br>(mmol/mol) | rho<br>(g/m**3) | wv<br>(m/s) | max<br>(m/s) |
|---|------------------------|-------------|-------------|-------------|----------------|-----------|-----------------|-----------------|-----------------|--------------|--------------------|-----------------|-------------|--------------|
| 0 | 01.01.2009<br>00:10:00 | 996.52      | -8.02       | 265.40      | -8.90          | 93.3      | 3.33            | 3.11            | 0.22            | 1.94         | 3.12               | 1307.75         | 1.03        | 1.75         |
| 1 | 01.01.2009<br>00:20:00 | 996.57      | -8.41       | 265.01      | -9.28          | 93.4      | 3.23            | 3.02            | 0.21            | 1.89         | 3.03               | 1309.80         | 0.72        | 1.50         |
| 2 | 01.01.2009<br>00:30:00 | 996.53      | -8.51       | 264.91      | -9.31          | 93.9      | 3.21            | 3.01            | 0.20            | 1.88         | 3.02               | 1310.24         | 0.19        | 0.63         |
| 3 | 01.01.2009<br>00:40:00 | 996.51      | -8.31       | 265.12      | -9.07          | 94.2      | 3.26            | 3.07            | 0.19            | 1.92         | 3.08               | 1309.19         | 0.34        | 0.50         |
| 4 | 01.01.2009<br>00:50:00 | 996.51      | -8.27       | 265.15      | -9.04          | 94.1      | 3.27            | 3.08            | 0.19            | 1.92         | 3.09               | 1309.00         | 0.32        | 0.63         |

Từ bảng trên ta thấy, cứ mỗi giờ thì có 6 quan sát được ghi lại. Vậy nên mỗi ngày sẽ có 144 ( $6 \times 24$ ) quan sát.

Với mục tiêu dự đoán nhiệt độ của 6 giờ kế tiếp. Ta thử chọn ra 5 ngày quan sát trước đó. Tức là ta sẽ lấy 720 ( $5 \times 144$ ) quan sát. Tương ứng ta sẽ có 2 tham số 'history\_size' biểu diễn cho số quan sát trước đó và 'target\_size' là số quan sát cần dự đoán.

Ta dùng 300,000 hàng đầu tiên của tập dữ liệu để làm tập dữ liệu huấn luyện, phần còn lại sẽ làm dữ liệu kiểm thử.

```
1 TRAIN_SPLIT = 300000  
2 tf.random.set_seed(13) #Setting seed to ensure reproducibility.
```

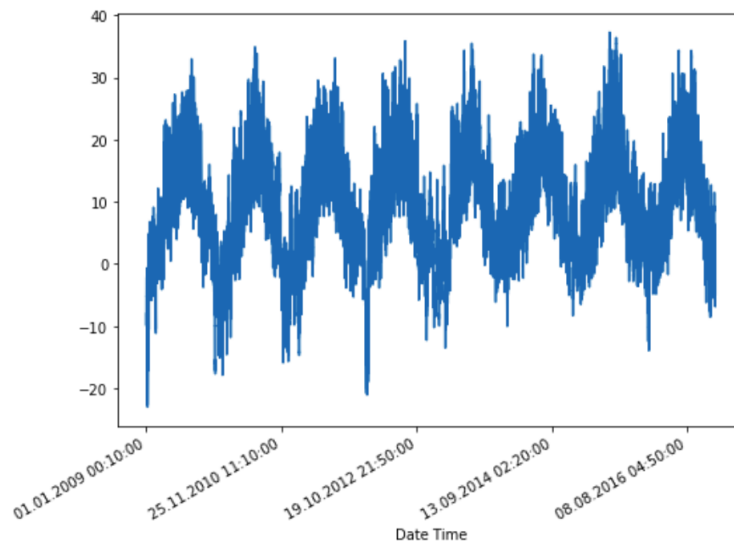
## Dự đoán với đơn biến

Với trường hợp này ta sẽ chỉ dùng một giá trị duy nhất là nhiệt độ, để dự đoán nhiệt độ trong tương lai.

```
1 uni_data = df['T (degC)']
2 uni_data.index = df['Date Time']
3 uni_data.head()
```

Xem thử dữ liệu

```
1 uni_data.plot(subplots=True)
2 uni_data = uni_data.values
```



Để cho mô hình hoạt động hiệu quả, chính xác hơn, ta cần chuẩn hóa dữ liệu trước khi cho đi huấn luyện.

```
1 uni_train_mean = uni_data[:TRAIN_SPLIT].mean()
2 uni_train_std = uni_data[:TRAIN_SPLIT].std()
3 uni_data = (uni_data - uni_train_mean) / uni_train_std
```

## Chế biến dữ liệu lại thành tập dữ liệu huấn luyện

```
1 def univariate_data(dataset, start_index, end_index, history_size,
2   target_size):
3     data = []
4     labels = []
5
6     start_index = start_index + history_size
7     if end_index is None:
8         end_index = len(dataset) - target_size
```

```

9     for i in range(start_index, end_index):
10         indices = range(i-history_size, i)
11         # Reshape data from (history_size,) to (history_size, 1)
12         data.append(np.reshape(dataset[indices], (history_size, 1)))
13         labels.append(dataset[i+target_size])
14     return np.array(data), np.array(labels)
15
16
17 univariate_past_history = 20
18 univariate_future_target = 0
19
20 x_train_uni, y_train_uni = univariate_data(univariate_data, 0, TRAIN_SPLIT
21     , univariate_past_history, univariate_future_target)
22 x_val_uni, y_val_uni = univariate_data(univariate_data, TRAIN_SPLIT, None,
23     univariate_past_history,
24     univariate_future_target)

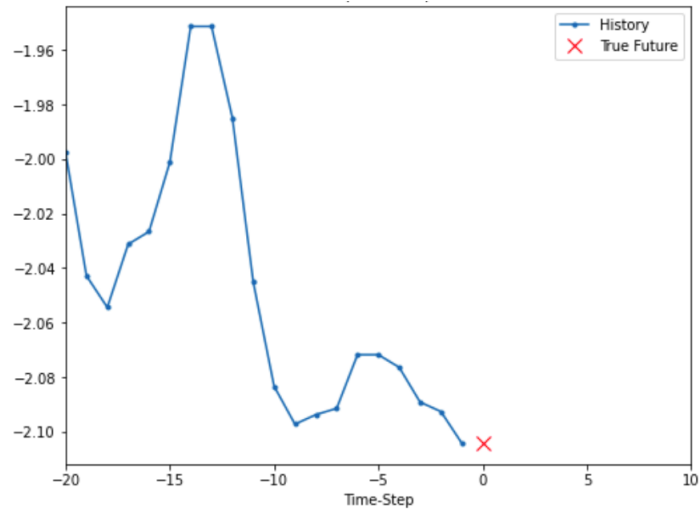
```

Vậy là ta đã chuẩn bị xong dữ liệu, hãy xem thử nó trông như thế nào.

```

1 def create_time_steps(length):
2     return list(range(-length, 0))
3
4 def show_plot(plot_data, delta, title):
5     labels = ['History', 'True Future', 'Model Prediction']
6     marker = ['.-', 'rx', 'go']
7     time_steps = create_time_steps(plot_data[0].shape[0])
8     if delta:
9         future = delta
10    else:
11        future = 0
12
13    plt.title(title)
14    for i, x in enumerate(plot_data):
15        if i:
16            plt.plot(future, plot_data[i], marker[i], markersize=10,
17                label=labels[i])
18        else:
19            plt.plot(time_steps, plot_data[i].flatten(), marker[i], label
20                =labels[i])
21    plt.legend()
22    plt.xlim([time_steps[0], (future+5)*2])
23    plt.xlabel('Time-Step')
24    return plt
25
26 show_plot([x_train_uni[0], y_train_uni[0]], 0, 'Sample Example')

```



## Tiến hành huấn luyện

```

1 BATCH_SIZE = 256
2 BUFFER_SIZE = 10000
3
4 train_univariate = tf.data.Dataset.from_tensor_slices((x_train_uni,
5 y_train_uni))
6 train_univariate = train_univariate.cache().shuffle(BUFFER_SIZE).
7 batch(BATCH_SIZE).repeat()
8
9 val_univariate = tf.data.Dataset.from_tensor_slices((x_val_uni,
10 y_val_uni))
11 val_univariate = val_univariate.batch(BATCH_SIZE).repeat()
12
13 simple_lstm_model = tf.keras.models.Sequential([
14 tf.keras.layers.LSTM(8, input_shape=x_train_uni.shape[-2:]),
15 tf.keras.layers.Dense(1)
16 ])
17
18 simple_lstm_model.compile(optimizer='adam', loss='mae')
19
20 simple_lstm_model = tf.keras.models.Sequential([
21 tf.keras.layers.LSTM(8, input_shape=x_train_uni.shape[-2:]),
22 tf.keras.layers.Dense(1)
23 ])
24
25 simple_lstm_model.compile(optimizer='adam', loss='mae')
26
27 EVALUATION_INTERVAL = 200
28 EPOCHS = 10

```

```

26
27 simple_lstm_model.fit(train_univariate, epochs=EPOCHS,
28 steps_per_epoch=EVALUATION_INTERVAL,
29 validation_data=val_univariate, validation_steps=50)

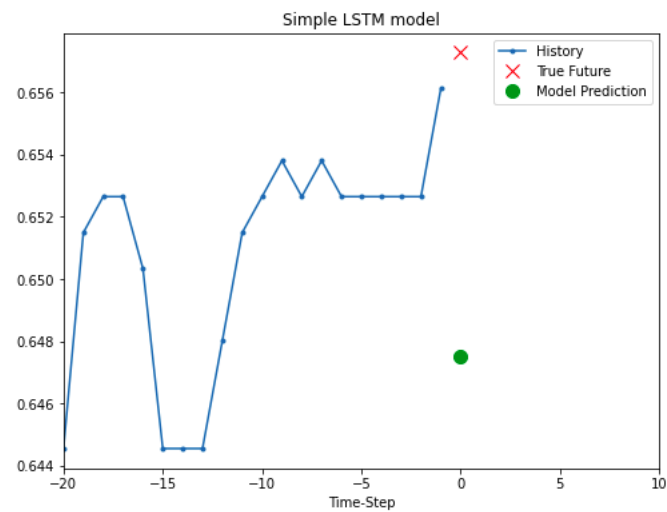
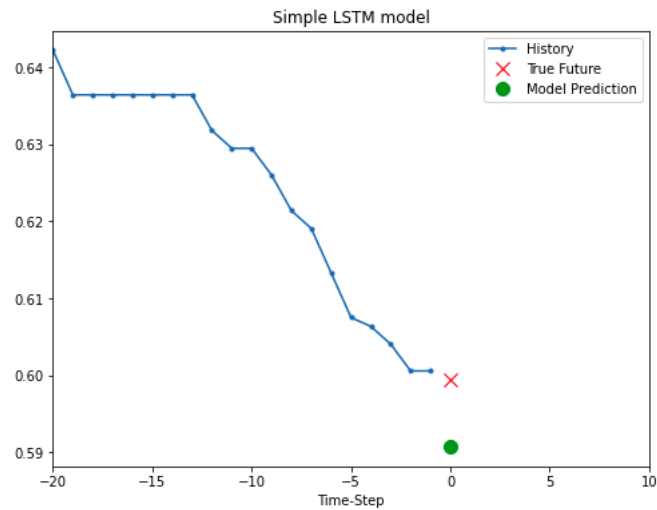
```

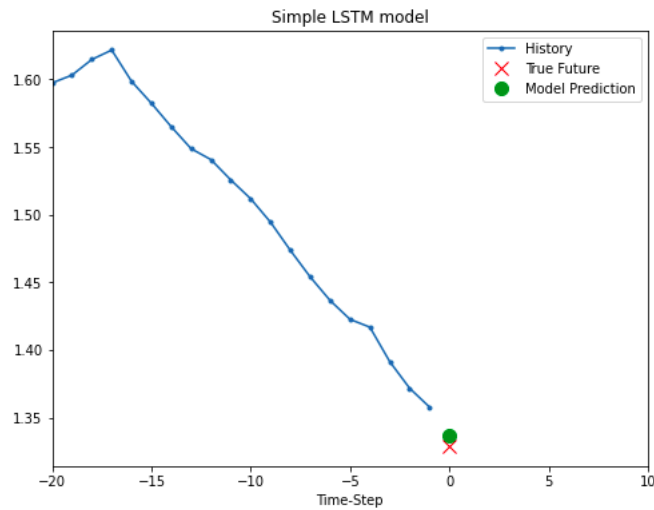
## Thử một vài dự đoán

```

1 for x, y in val_univariate.take(3):
2 plot = show_plot([x[0].numpy(), y[0].numpy()],
3 simple_lstm_model.predict(x)[0], 0, 'Simple LSTM model')
4 plot.show()

```





## Dự đoán với đa biến

### Dự đoán một bước (one-step)

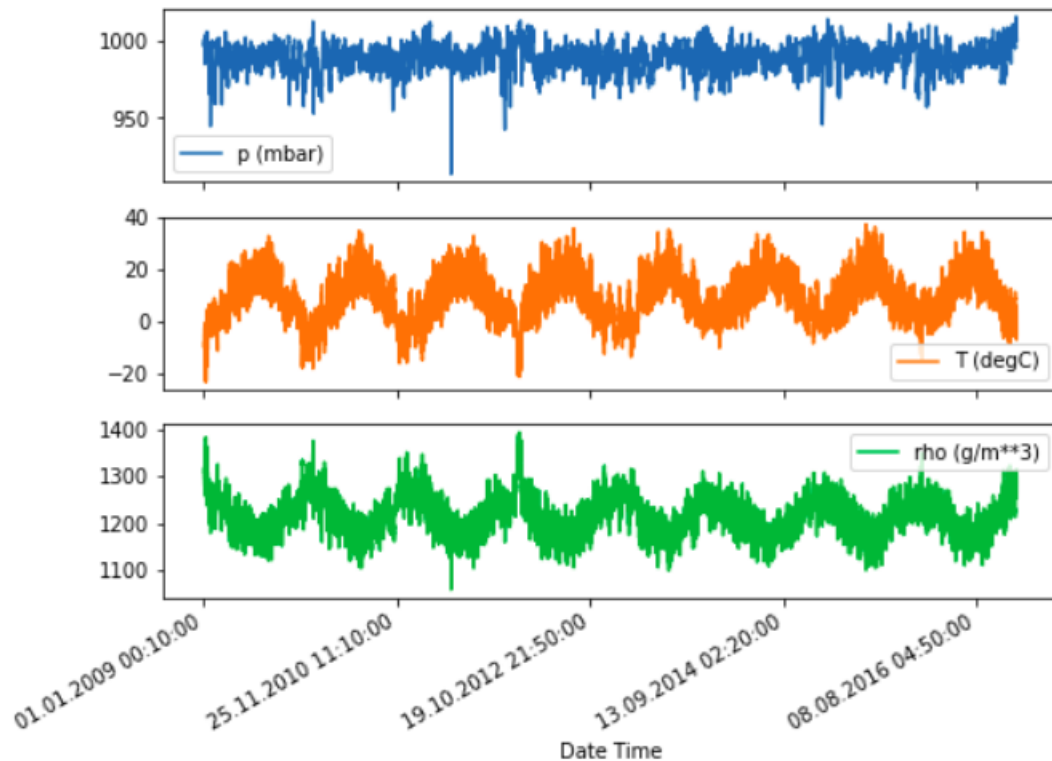
Với trường hợp này ta sẽ chỉ dùng 3 đặc trưng trong 14 đặc trưng trong tập dữ liệu, để dự đoán nhiệt độ trong tương lai.

```
1 features_considered = ['p (mbar)', 'T (degC)', 'rho (g/m**3)']
2
3 features = df[features_considered]
4 features.index = df['Date Time']
```

Xem thử dữ liệu

```
1 features.head()
2 features.plot(subplots=True)
```

|                     | p (mbar) | T (degC) | rho (g/m**3) |
|---------------------|----------|----------|--------------|
| Date Time           |          |          |              |
| 01.01.2009 00:10:00 | 996.52   | -8.02    | 1307.75      |
| 01.01.2009 00:20:00 | 996.57   | -8.41    | 1309.80      |
| 01.01.2009 00:30:00 | 996.53   | -8.51    | 1310.24      |
| 01.01.2009 00:40:00 | 996.51   | -8.31    | 1309.19      |
| 01.01.2009 00:50:00 | 996.51   | -8.27    | 1309.00      |



Tương tự như trên, ta cũng chuẩn hóa dữ liệu trước khi cho đi huấn luyện.

```

1 dataset = features.values
2 data_mean = dataset[:TRAIN_SPLIT].mean(axis=0)
3 data_std = dataset[:TRAIN_SPLIT].std(axis=0)
4
5 dataset = (dataset-data_mean)/data_std

```

**Chế biến dữ liệu lại thành tập dữ liệu huấn luyện**

```

1 def multivariate_data(dataset, target, start_index, end_index,
2     history_size,
3     target_size, step, single_step=False):
4     data = []
5     labels = []
6
7     start_index = start_index + history_size
8     if end_index is None:
9         end_index = len(dataset) - target_size
10
11     for i in range(start_index, end_index):
12         indices = range(i-history_size, i, step)
13         data.append(dataset[indices])

```



```

13
14     if single_step:
15         labels.append(target[i+target_size])
16     else:
17         labels.append(target[i:i+target_size])
18
19     return np.array(data), np.array(labels)
20
21 past_history = 720
22 future_target = 72
23 STEP = 6
24
25 x_train_single, y_train_single = multivariate_data(dataset, dataset
26    [:, 1], 0,
27     TRAIN_SPLIT, past_history,
28     future_target, STEP,
29     single_step=True)
30 x_val_single, y_val_single = multivariate_data(dataset, dataset[:,
31     1],
32     TRAIN_SPLIT, None, past_history,
33     future_target, STEP,
34     single_step=True)

```

## Tiến hành huấn luyện

```

1 train_data_single = tf.data.Dataset.from_tensor_slices((
2     x_train_single, y_train_single))
3
4 train_data_single = train_data_single.cache().shuffle(BUFFER_SIZE).
5     batch(BATCH_SIZE).repeat()
6
7 val_data_single = tf.data.Dataset.from_tensor_slices((x_val_single,
8     y_val_single))
9 val_data_single = val_data_single.batch(BATCH_SIZE).repeat()
10
11 single_step_model = tf.keras.models.Sequential()
12 single_step_model.add(tf.keras.layers.LSTM(32,
13     input_shape=x_train_single.shape[-2:]))
14 single_step_model.add(tf.keras.layers.Dense(1))
15
16 single_step_model.compile(optimizer=tf.keras.optimizers.RMSprop(),
17     loss='mae')
18
19 single_step_history = single_step_model.fit(train_data_single,
20     epochs=EPOCHS,
21     steps_per_epoch=EVALUATION_INTERVAL, validation_data=
22     val_data_single, validation_steps=50)

```

Xem quá biểu đồ biểu diễn độ lỗi

```
1 def plot_train_history(history, title):
2     loss = history.history['loss']
3     val_loss = history.history['val_loss']
4
5     epochs = range(len(loss))
6
7     plt.figure()
8
9     plt.plot(epochs, loss, 'b', label='Training loss')
10    plt.plot(epochs, val_loss, 'r', label='Validation loss')
11    plt.title(title)
12    plt.legend()
13
14    plt.show()
15
16 plot_train_history(single_step_history,
17 'Single Step Training and validation loss')
```



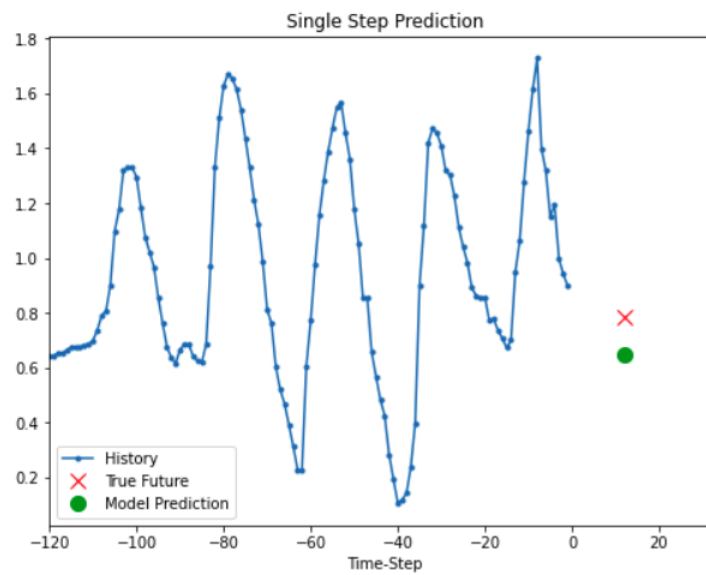
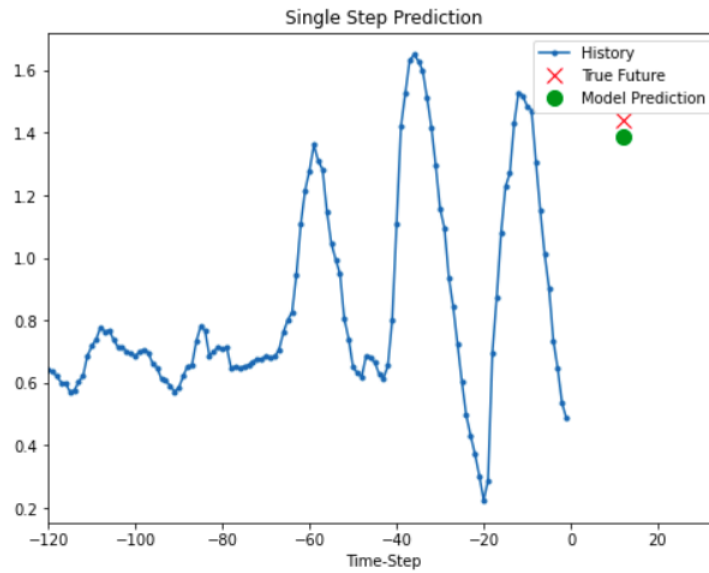
Thử một vài dự đoán

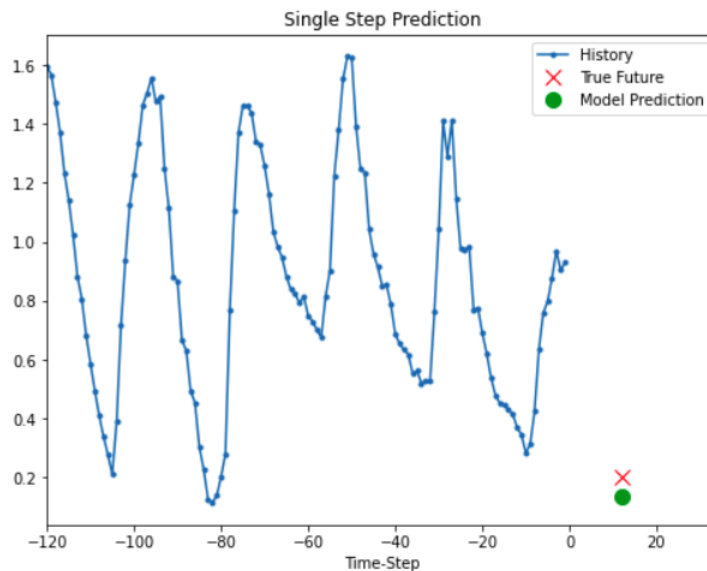
```
1 for x, y in val_data_single.take(3):
2     plot = show_plot([x[0][:, 1].numpy(), y[0].numpy()],
```

```

3 single_step_model.predict(x)[0]], 12,
4 'Single Step Prediction')
5 plot.show()

```





## Dự đoán nhiều bước (Multi-step)

Các bước bên dưới thực hiện tương tự như trên.

```

1 future_target = 72
2 x_train_multi, y_train_multi = multivariate_data(dataset,
3   dataset[:, 1], 0,
4   TRAIN_SPLIT, past_history,
5   future_target, STEP)
6 x_val_multi, y_val_multi = multivariate_data(dataset, dataset
7  [:, 1],
8   TRAIN_SPLIT, None, past_history,
9   future_target, STEP)
10
11 train_data_multi = tf.data.Dataset.from_tensor_slices((
12   x_train_multi, y_train_multi))
13 train_data_multi = train_data_multi.cache().shuffle(BUFFER_SIZE
14   ).batch(BATCH_SIZE).repeat()
15
16 val_data_multi = tf.data.Dataset.from_tensor_slices((
17   x_val_multi, y_val_multi))
18 val_data_multi = val_data_multi.batch(BATCH_SIZE).repeat()
19
20 multi_step_model = tf.keras.models.Sequential()
21 multi_step_model.add(tf.keras.layers.LSTM(32,
22   return_sequences=True,
23   input_shape=x_train_multi.shape[-2:]))

```

```

19 multi_step_model.add(tf.keras.layers.LSTM(16, activation='relu'
    ))
20 multi_step_model.add(tf.keras.layers.Dense(72))
21
22 multi_step_model.compile(optimizer=tf.keras.optimizers.RMSprop(
    clipvalue=1.0), loss='mae')
23
24 multi_step_history = multi_step_model.fit(train_data_multi,
    epochs=EPOCHS,
25 steps_per_epoch=EVALUATION_INTERVAL,
26 validation_data=val_data_multi,
27 validation_steps=50)

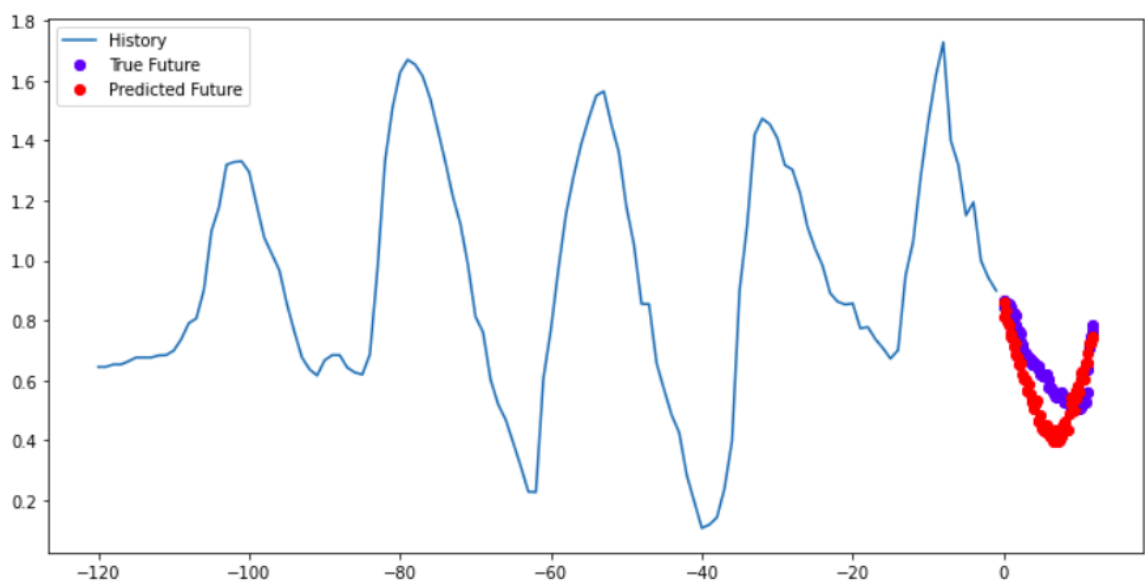
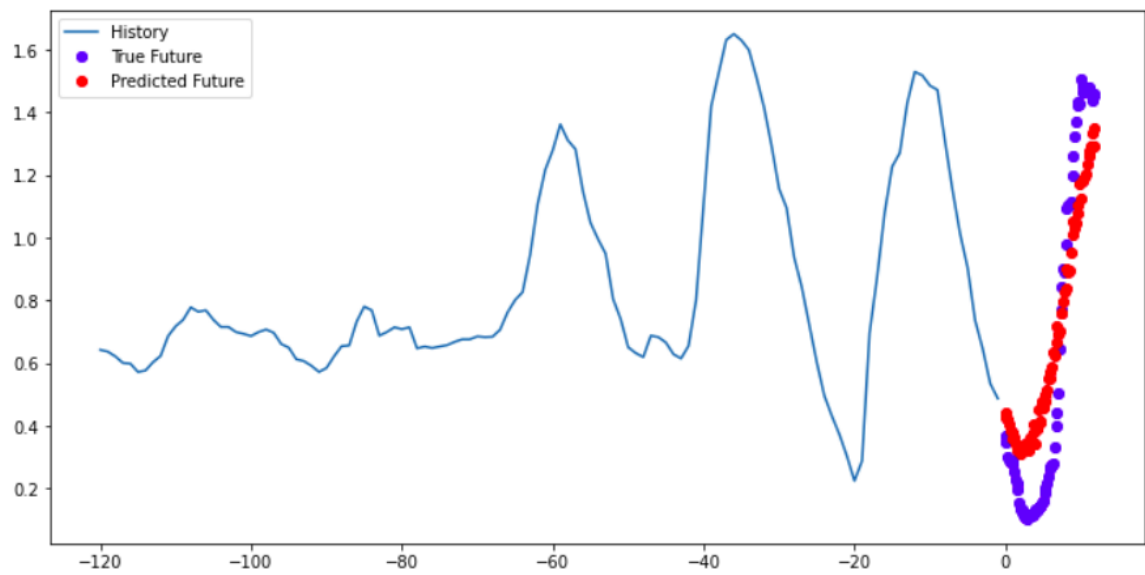
```

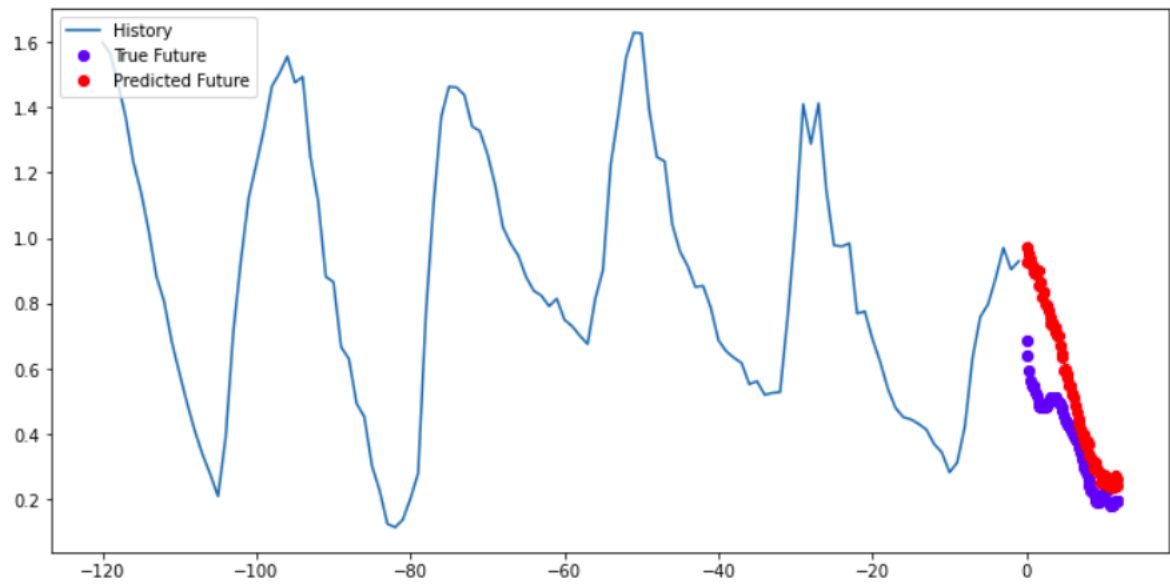
Xem vài kết quả dự đoán

```

1 def multi_step_plot(history, true_future, prediction):
2     plt.figure(figsize=(12, 6))
3     num_in = create_time_steps(len(history))
4     num_out = len(true_future)
5
6     plt.plot(num_in, np.array(history[:, 1]), label='History')
7     plt.plot(np.arange(num_out)/STEP, np.array(true_future), 'bo'
8         ,
9         label='True Future')
10    if prediction.any():
11        plt.plot(np.arange(num_out)/STEP, np.array(prediction), 'ro'
12            , label='Predicted Future')
13    plt.legend(loc='upper left')
14    plt.show()
15
16 for x, y in train_data_multi.take(1):
17     multi_step_plot(x[0], y[0], np.array([0]))

```





## Chương 4: TỔNG KẾT

### Kết quả.

### Kết quả đạt được.

Qua những nội dung được trình bày ở trên, ta có thể thấy bài báo cáo đã trình bày một cách cơ bản và đầy đủ nhất chủ đề tìm hiểu về time series forecasting. Qua đó ta có thể hiểu được thế nào là time series, ứng dụng và cách triển khai thuật toán. Nó mang nhiều giá trị hiện thực không chỉ trong nghiên cứu khoa học mà còn có thể áp dụng vào thực tiễn của cuộc sống.

### Hạn chế.

Bài báo cáo tuy đã trình bày tốt nhất có thể nhưng vẫn còn tồn tại một vài thiếu sót nhỏ. Các phần cần nêu theo mục đích của đề tài đã được hoàn thành nhưng chỉ ở mức là tổng quát, chung nhất, chưa đi vào chi tiết và trình bày cụ thể.

### Hướng phát triển.

Có nhiều hướng phát triển khác nhau cho các bài time series forecasting. Ngoài việc xây dựng thuật toán theo các phương pháp trên thì còn có nhiều phương pháp khác, được thiết kế lại để tối ưu hóa hiệu năng, độ chính xác,... cho những trường hợp và mục đích riêng biệt. Ngành công nghệ thông tin ngày càng phát triển, đặc biệt là các lĩnh vực trí tuệ nhân tạo, học máy, học sâu. Nên trong tương lai, chắc chắn sẽ có thêm nhiều cách tiếp cận mới, hiệu quả và chính xác hơn. Không chỉ hữu ích trong bài toán time series forecasting mà còn nhiều bài toán khác nữa.



## TÀI LIỆU THAM KHẢO

- [1]. <https://towardsdatascience.com/deep-learning-for-time-series-and-why-deep-learning-a6120b147d60>
- [2]. <https://www.slideshare.net/BeriDang/m-hnh-arima-103586847>
- [3]. <https://viblo.asia/p/time-series-data-gDVK2Qbv5Lj>
- [4]. [https://en.wikipedia.org/wiki/Time\\_series](https://en.wikipedia.org/wiki/Time_series)
- [5]. <https://arxiv.org/pdf/2004.13408.pdf>
- [6]. Deep Learning for Time Series Forecasting - Jason Brownlee
- [7]. <http://www.ijmlc.org/vol7/632-P17.pdf>
- [8]. [https://www.tensorflow.org/tutorials/structured\\_data/time\\_series](https://www.tensorflow.org/tutorials/structured_data/time_series)