

**TRƯỜNG ĐẠI HỌC KHOA HỌC TỰ NHIÊN TP HỒ CHÍ MINH**

**MÔN HỌC: KHOA HỌC DỮ LIỆU**

# **BÁO CÁO ĐỒ ÁN CUỐI KỲ**



**Nhóm 21**

1612272 - Trần Nhật Huy

1512222 - Nguyễn Duy Hưng

## **I.     Họp nhóm và phân công công việc**

- **Lần 1**

Ngày: 24/11/2019

Địa điểm: Facebook

<b>Công việc</b>	<b>Ngày thực hiện</b>	<b>Người làm</b>
Crawl dữ liệu	26/11/2019	Nguyễn Duy Hưng
Viết báo cáo	30/11/2019	Trần Nhật Huy

- **Lần 2**

Ngày: 02/01/2020

Địa điểm: Facebook

<b>Công việc</b>	<b>Ngày thực hiện</b>	<b>Người làm</b>
Thống nhất kế hoạch làm	02/01/2020	Nguyễn Duy Hưng Trần Nhật Huy
Viết kế hoạch và chỉnh sửa	03/01/2020	Nguyễn Duy Hưng Trần Nhật Huy
Thí nghiệm 1	05/01/2020	Trần Nhật Huy
Thí nghiệm 2	05/01/2020	Nguyễn Duy Hưng
Thí nghiệm 3	05/01/2020	Trần Nhật Huy
Thí nghiệm 4, 5, 6	06/01/2020	Nguyễn Duy Hưng
Thí nghiệm 07	06/01/2020	Trần Nhật Huy
Viết báo cáo	06/01/2020	Nguyễn Duy Hưng
Slide	06/01/2020	Trần Nhật Huy

## **II.    Bài toán đặt ra và hướng giải quyết**

### **1. Bài toán**

- Nhóm muốn dự đoán độ ẩm của ngày hôm sau dựa vào dữ liệu của ngày hôm trước
- Input: Dữ liệu thời tiết của ngày hôm trước (temperature, humidity, ...)
- Output: Humidity của ngày hôm sau
- Lợi ích: Chúng ta có thể dự đoán trước được humidity của một ngày nào đó. Phục vụ cho nông nghiệp hoặc dự báo thời tiết
- Nguồn gốc: Nhóm tự đặt ra câu hỏi này

## **2. Thu thập dữ liệu**

- Dữ liệu của nhóm được lấy từ API Dark Sky từ 01-01-2010 đến 31-12-2011 ở thành phố Hồ Chí Minh, Việt Nam

## **3. Tiền xử lý dữ liệu**

- Điền giá trị thiếu, nhóm sử dụng SimpleImputer với các strategy lần lượt là mean, median, most\_frequent
- Lựa chọn đặc trưng nhóm sử dụng PCA, SelectKBest, Pearson, Feature Importance + với tự tạo thêm đặc trưng ở một số thí nghiệm
- Chuẩn hóa dữ liệu nhóm sử dụng StandardScalar

## **4. Các thí nghiệm đã thực hiện**

- Thí nghiệm 1

Nhóm tách thuộc tính time thành day và month. Bỏ thuộc tính visibility, apparentTemperature. Dùng PCA với n\_components từ [1:số lượng cột].

Độ lỗi MAE trên tập validation là 0.054. Không tốt vì std của cột Humidity là 0.08.

- Thí nghiệm 2

Nhóm giữ nguyên các thuộc tính ban đầu, không thêm hay xóa bất kỳ cột nào. Độ lỗi MAE trên tập validation là 0.055

- Thí nghiệm 3

Nhóm lựa chọn đặc trưng bằng Pearson's Correlation

Độ lỗi MAE trên tập validation là 0.056

- Thí nghiệm 4

Nhóm lựa chọn đặc trưng bằng SelectKBest

Độ lỗi MAE trên tập validation là 0.056

- Thí nghiệm 5

Nhóm lựa chọn đặc trưng bằng Feature Importance

Độ lỗi MAE trên tập validation là 0.057

- Thí nghiệm 6

Nhóm sử dụng mô hình Neural Network

Độ lỗi MAE trên tập validation là 0.07 với  $\alpha=0.1$  và hidden layer là 1

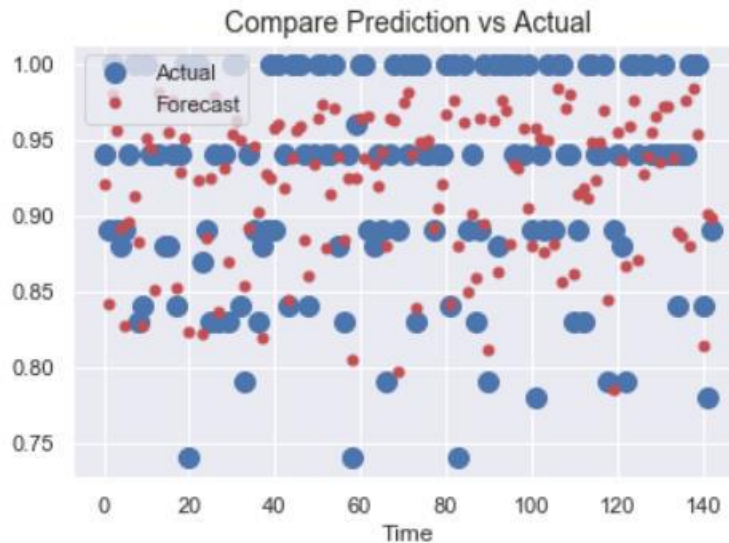
Độ lỗi của mô hình này còn lớn hơn trước. Có thể bị overfitting

- Thí nghiệm 7

Nhóm dự đoán có thể sử dụng dữ liệu của một ngày chưa đủ nên thử thêm độ trễ của k ngày trước đó. Nhóm thử nghiệm với  $k = 7$

Độ lỗi MAE trên tập validation là 0.045

### **III. Tổng kết và đánh giá đóng góp của các thành viên trong nhóm**



- Độ lỗi là 0.045
- Các thí nghiệm mà nhóm thực hiện có kết quả không được tốt
- Nhóm nghĩ có thể nguyên nhân là:
  - Dữ liệu có nhiều nhưng nhóm xử lý chưa tốt
  - Mô hình không phù hợp, bị underfitting
- Link github: [https://github.com/hungnd08/DS\\_Final\\_Project](https://github.com/hungnd08/DS_Final_Project)
- Đánh giá đóng góp của các thành viên trong nhóm

Ở mục contributors account “hungnd08” bị thiếu mất 6 commit vì em kết nối với project bằng SSH key và trong file .gitconfig trên máy cấu hình name em để là “hungnd” và email của em khác với email đăng ký account trên github nên khi commit github nó không tính vào mục contributors, phía dưới em có chụp thêm phần lịch sử commit để thầy xem ạ.

## Learn Git and GitHub without any code!

Using the Hello World guide, you'll start a branch, write comments, and open a pull request.

[Read the guide](#)

hungnd08 / DS\_Final\_Project

Unwatch 2 Star 0 Fork 0

[Code](#) [Issues 1](#) [Pull requests 0](#) [Actions](#) [Projects 0](#) [Wiki](#) [Security](#) [Insights](#) [Settings](#)

Pulse

Contributors

Community

Traffic

Commits

Code frequency

Dependency graph

Network

Forks

Nov 10, 2019 – Jan 8, 2020

Contributions: Commits

Contributions to master, excluding merge commits



Search or jump to...Pull requestsIssuesMarketplaceExplore

Learn Git and GitHub without any code!

Using the Hello World guide, you'll start a branch, write comments, and open a pull request.

Read the guide

hungnd08 / DS\_Final\_Project

Unwatch2★Star0🔒Fork0

CodeIssuesPull requestsActionsProjectsWikiSecurityInsightsSettings

Branch: master

Commits on Jan 7, 2020

Fix methods

hungnd08 committed yesterday

8795c91

add .gitignore

thathuy13388 committed 2 days ago

285c43b

- Update folder

thathuy13388 committed 2 days ago

a809fca

Commits on Jan 6, 2020

Fix report file

hungnd08 committed 2 days ago

Verifiedc758075

Add report file

hungnd08 committed 2 days ago

Verifieda6c88ea

- Fix Neural Network

thathuy13388 committed 2 days ago

d982af9

Merge branch 'try'

thathuy13388 committed 2 days ago

7968711

- Delete file

thathuy13388 committed 2 days ago

ae95abf

- Finish slide v.1

thathuy13388 committed 2 days ago

3a530a5

- Update slide

thathuy13388 committed 2 days ago

287e75b

- Merge try

thathuy13388 committed 2 days ago

c7a2bc9

- Add slide

thathuy13388 committed 2 days ago

8e683d1

Fix selectBest method and add comments

hungnd committed 2 days ago

4e23999

Try another methods

hungnd committed 2 days ago

7058b17

Try another methods

hungnd committed 2 days ago

c28cc05

- Add method 2

thathuy13388 committed 2 days ago

623803b

- Finish method 1

thathuy13388 committed 2 days ago

a8663d2

Commits on Jan 5, 2020

- Finish first version

thathuy13388 committed 3 days ago

8b18310

- Add preprocessing transformer

thathuy13388 committed 3 days ago

e045546

Commits on Jan 4, 2020

Handling data

hungnd committed 4 days ago

a5788fa

Commits on Jan 3, 2020

- Update

thathuy13388 committed 6 days ago

06c8493

- Update plan

thathuy13388 committed 6 days ago

53214db

Commits on Dec 8, 2019

- Sua lai cau hoi ma nhom dat ra

thathuy13388 committed on Dec 8, 2019

2b94309

Commits on Dec 2, 2019

- Sua lai file bao cao va lam ra cau hoi

thathuy13388 committed on Dec 2, 2019

87ad28a

Commits on Nov 30, 2019

- Sua lai dia diem thu thap du lieu trong file bao cao

thathuy13388 committed on Nov 28, 2019

b37803c

- Them file bao cao lan 1

thathuy13388 committed on Nov 20, 2019

6ac756c

Commits on Nov 26, 2019

Merge branch 'master' of github.com:hungnd08/DS\_Final\_Project

hungnd committed on Nov 26, 2019

6a3ef1f

Collect weather data

hungnd committed on Nov 26, 2019

8ae922f

Commits on Nov 14, 2019

Initial commit

hungnd08 committed on Nov 14, 2019

Verified8754c78

NeverOlder

© 2020 GitHub, Inc. TermsPrivacySecurityStatusHelp

ContactGitHubPricingAPITrainingBlogAbout

