

ĐỒ ÁN CUỐI KỲ

Phân tích bình luận

GVĐT:

*Phạm Nguyễn Trường An
Lê Đình Duy*

Lớp: CS114.K21

--> Danh sách nhóm <--

- **Trịnh Hưng Long 18521060**
- **Hà Văn Luân 18521062**
- **Lỗ Đình Phưởng 18521274**

Giới thiệu

- Hiện nay nhu cầu mua hàng qua mạng của người dùng ngày càng trở nên phát triển mạnh hơn do những lợi ích mà nó mang lại, như tiện lợi, chi phí rẻ, có nhiều chương trình khuyến mãi hấp dẫn, có thể ngồi ở nhà để xem sản phẩm mà không cần phải đến tận nơi để xem, ... Tuy nhiên, việc mua hàng qua mạng cũng có những nhược điểm, trong đó có việc người dùng không thể tận mắt đánh giá sản phẩm của mình như mua trực tiếp tại các cửa hàng được (Các bình luận đóng vai trò quan trọng)
- Để minh chứng hơn, chúng em sẽ tiến hành dự đoán bình luận của một sản phẩm điện thoại trên trang thương mại điện tử : *Thế Giới Di Động* để hỗ trợ việc mua hàng của người tiêu dùng.

Mô tả

1. Ý tưởng

- Đưa đoạn bình luận của 1 sản phẩm bất kỳ và đưa ra kết quả dự đoán bình luận (là tích cực, tiêu cực hoặc trung tính)

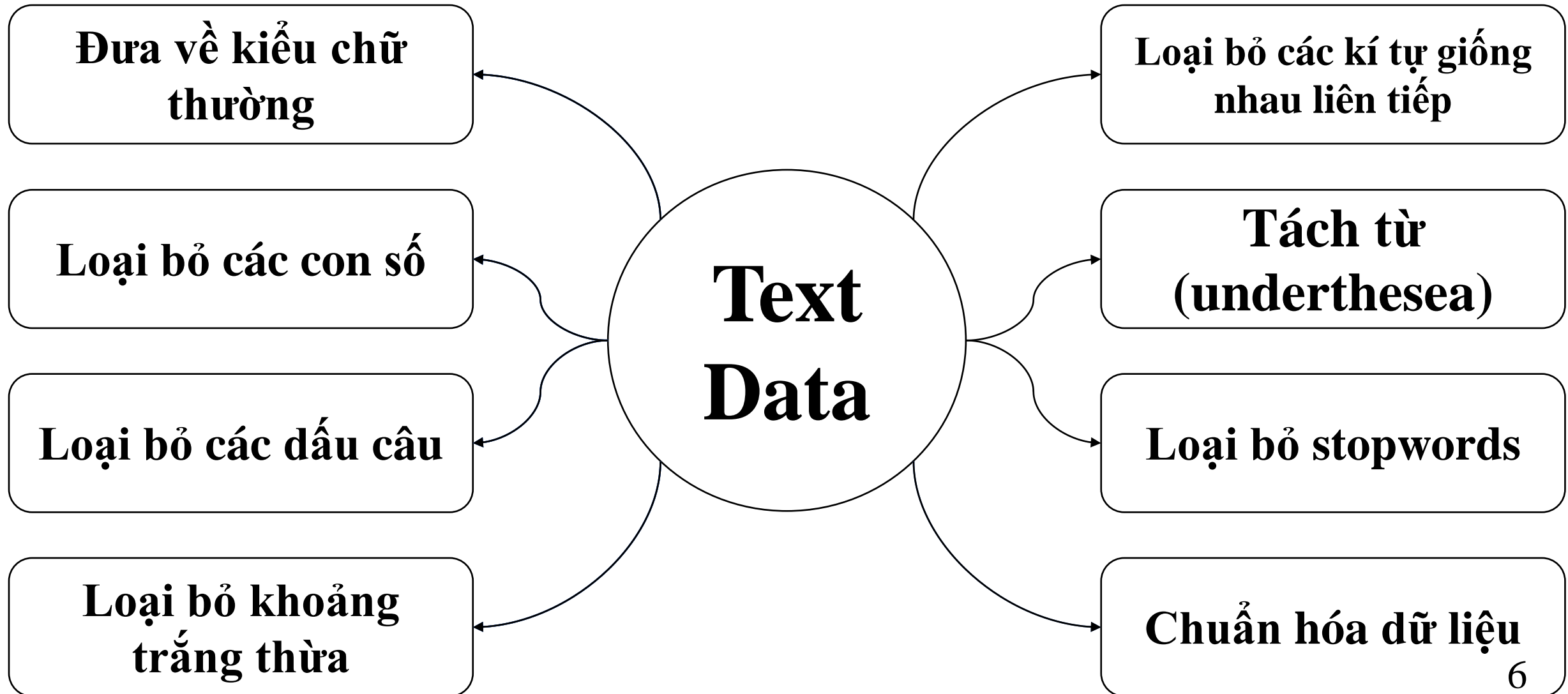
2. Input & Output

- Input: Một bình luận về sản phẩm điện thoại của khách hàng bằng tiếng Việt.
- Output: Bình luận đó là tiêu cực, trung tính hay tích cực.

Prepare Dataset

- Crawl data (các bình luận sản phẩm trên trang) từ trang thương mại điện tử: ['https://www.thegioididong.com/dtdd'](https://www.thegioididong.com/dtdd) (Thế Giới Di Động) bằng thư viện BeautifulSoup – 1 package Python dùng để phân tích cú pháp các tài liệu HTML và XML
- Ta thu về được 4679 bình luận các nhãn được gán tự động với:
 - ✓ Các bình luận có số lượng đạt 5 sao là 5 sẽ cho nhãn bằng 1
 - ✓ Các bình luận có số lượng đạt 5 sao là 4 và 3 sẽ cho nhãn bằng 0
 - ✓ Các bình luận có số lượng đạt 5 sao là 1 và 2 sẽ cho nhãn bằng -1
- Sau đó, tụi em sẽ thực hiện kiểm tra các nhãn của bộ data trên, thu gọn và làm cân bằng bộ data trên -> thu về được 3000 bình luận (1000 bình luận tích cực, 1000 bình luận trung tính và 1000 bình luận tiêu cực)

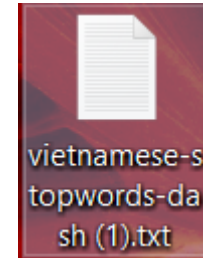
Data preprocessing



Data preprocessing

* Loại bỏ stopwords:

Link: “<https://github.com/stopwords/vietnamese-stopwords/blob/master/vietnamese-stopwords-dash.txt>”



a_lô	biết_được	cho_tới_khi	dào	ngồi_trệt
a_ha	buổi	cho_về	dì	ngộ_nhỡ
ai	buổi_làm	cho_ăn	dưới	nhau
ai_ai	buổi_mới	cho_đang	dưới_nước	nhiên_hậu
ai_nấy	buổi_ngày	cho_được	dạ	nhiệt_liệt
ai_đó	buổi_sớm	cho_đến	dạ_bán	nhung_nhằng
alô	bà	cho_đến_khi	dạ_con	nhà
amen	bà_ấy	cho_đến_nỗi	dạ_dài	nhà_chung
anh	bài	choa	dạ_dạ	nhà_khó
anh_ấy	bài_bác	chui_cha	dạ_khách	nhà_làm
ba	bài_bỏ	chung	dần_dần	nhà_ngoài
ba_ba	bài_cái	chung_cho	dần_dần	nhà_người
ba_bản	bác	chung_chung	dầu_sao	nhà_tôi
ba_cùng	bán	chung_cuộc	dẫn	nhà_việc
ba_họ	bán_cấp	chung_cục	dầu	nhóm
ba_ngày	bán_dạ	chung_nhau	dầu_mà	nhón_nhén
ba_ngôi	bán_thế	chung_qui	dầu_rằng	nhất_loạt
ba_tầng	bây_bầy	chung_quy	dầu_sao	nhất_luật
bao_giờ	bây_chữ	chung_quy_lại	em	nhất_là
bao_lâu	bây_giờ	chung_ái	em_em	nhất_mức
bao_nhiều	bây_nhiều	chuyển	giá_trị	nhất_nhất
bao_nà	bền	chuyển_tự	giá_trị_thực_tế	nhất_quyết
bay_biển	béng	chuyển_đạt	giờ	nhất_sinh
biết	bên	chuyện	giờ_lâu	nhất_thiết
biết_bao	bên_bị	chuẩn_bị	giờ_này	nhất_thì
biết_bao_nhiều	bên_cố	chành_chạnh	giờ_đi	nhất_tâm
biết_chắc	bên_cạnh	chí_chết	giờ_đây	nhất_tề
biết_chứng_nào	bông	chùn_chùn	giờ_đến	nhất_đán
biết_mình	bước	chùn_chùn	giữ	nhất_định
biết_mấy	bước_khỏi	chú	giữ_lấy	nhận_biết
biết_thế	bước_tới	chú_dẫn	giữ_ý	nhận_họ
biết_trước	bước_đi	chú_khách	giữa	nhận_làm

Data preprocessing

* Chuẩn hóa dữ liệu:

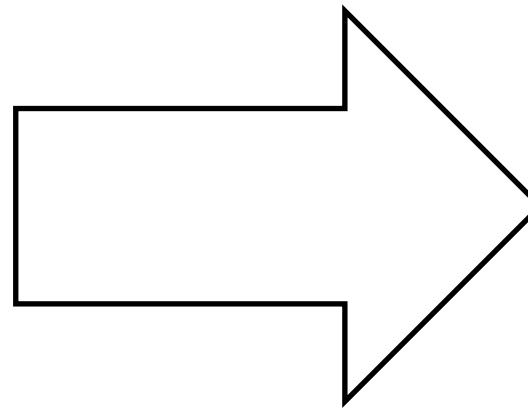
Trong lúc thu thập dữ liệu, em thu thập được một dict chứa các teencode, viết tắt,... Sau đó tìm trong các bình luận nếu chứa các từ giống với key của phần tử trong `replace_list` ,ta gán giá trị từ đó bằng value của key tương ứng

```
replace_list = {
    'ship': 'vận chuyển', 'shop': 'cửa hàng', 'sho': 'cửa hàng', 'm': 'mình', 'mik': 'mình', 'ko': 'không', 'k': 'không', 'kh': 'không',
    'khong': 'không', 'kg': 'không', 'khg': 'không', 'tl': 'trả lời', 'rep': 'trả lời', 'r': 'rồi', 'fb': 'facebook', 'face': 'faceook',
    'thanks': 'cảm ơn', 'thank': 'cảm ơn', 'tks': 'cảm ơn', 'tk': 'cảm ơn', 'ok': 'tốt', 'oki': 'tốt', 'okie': 'tốt', 'sp': 'sản phẩm',
    'dc': 'được', 'vs': 'với', 'dt': 'điện thoại', 'thjk': 'thích', 'thik': 'thích', 'qá': 'quá', 'trẻ': 'trẻ', 'bgjo': 'bao giờ',
    'h': 'giờ', 'qa': 'quá', 'dep': 'đẹp', 'xau': 'xấu', 'ib': 'nhắn tin', 'cute': 'dễ thương', 'sz': 'size', 'good': 'tốt', 'god': 'tốt',
    'bt': 'bình thường', 'sz': 'cỡ', 'size': 'cỡ', 'đx': 'được', 'dk': 'được', 'dc': 'được', 'đk': 'được', 'đc': 'được',
    'authentic': 'chuẩn chính hãng', 'aut': 'chuẩn chính hãng', 'auth': 'chuẩn chính hãng', 'thick': 'thích', 'gud': 'tốt', 'god': 'tốt',
    'wel done': 'tốt', 'good': 'tốt', 'gút': 'tốt', 'sầu': 'xấu', 'gut': 'tốt', 'tot': 'tốt', 'nice': 'tốt', 'perfect': 'rất tốt',
    'bt': 'bình thường', 'time': 'thời gian', 'qá': 'quá', 'ship': 'giao hàng', 'product': 'sản phẩm', 'quality': 'chất lượng', 'chat': 'chất',
    'exelent': 'hoàn hảo', 'bad': 'tệ', 'sad': 'tệ', 'beautiful': 'đẹp', 'tl': 'trả lời', 'r': 'rồi', 'order': 'đặt hàng',
    'chất lg': 'chất lượng', 'sd': 'sử dụng', 'dt': 'điện thoại', 'nt': 'nhắn tin', 'tl': 'trả lời', 'sài': 'xài', 'bjo': 'bao giờ',
    'thik': 'thích', 'sop': 'cửa hàng', 'fb': 'facebook', 'face': 'facebook', 'very': 'rất', 'dep': 'đẹp', 'xau': 'xấu', 'iu': 'yêu',
    'fake': 'giả mạo', 'trl': 'trả lời', '><': 'tiêu cực', 'por': 'tệ', 'poor': 'tệ', 'ib': 'nhắn tin', 'rep': 'trả lời', 'fback': 'feedback',
    'fedback': 'feedback', 'bin': 'pin', 'cx': 'cũng', 'nch': 'nói chuyện', 'ntn': 'như thế nào', 'vde': 'vấn đề'
}
```


Data preprocessing

Ví dụ:

“Hôm qua 2/3/2019, mình có mua tại cơ sở Thế Giới Di Động Võ Văn Ngân, sản phẩm đẹp, chất lượng, pin trâu, sạc nhanh. Nhân viên nhiệt tình chu đáo kkkkkkkk thanks”



“hôm qua mua cơ sở thế giới di động võ văn ngân sản phẩm đẹp chất lượng pin trâu sạc nhân viên nhiệt tình chu đáo k cảm ơn”

Feature Engineering

- **TF-IDF** (Term Frequency – Inverse Document Frequency)
 - +TF: Tần số xuất hiện của 1 từ trong văn bản
 - +IDF: Tần số nghịch của 1 từ trong một tập các văn bản
- Kỹ thuật **TF-IDF** dùng để tính toán tần suất xuất hiện của một từ trong văn bản, dựa trên đó để đánh giá mức độ quan trọng của từng từ trong văn bản
- TfidfVectorizer dùng để chuyển đổi dữ liệu văn bản sang ma trận các features **TF-IDF**

“Hôm qua 2/3/2019, mình có mua tại cơ sở Thế Giới Di Động Võ Văn Ngân, sản phẩm đẹp, chất lượng, pin trâu, sạc nhanh. Nhân viên nhiệt tình chu đáo kkkkkkk thanks”

TF-IDF

```
[[0. 0. 0. ... 0. 0. 0.]  
 [0. 0. 0. ... 0. 0. 0.]  
      ....  
 [0. 0. 0. ... 0. 0. 0.]  
 [0. 0. 0. ... 0. 0. 0.]
```

Xây dựng và huấn luyện model

Cách tính độ chính xác của model bằng *score* chỉ cho ta biết phần trăm dữ liệu được phân loại đúng mà không chỉ ra được phân loại như thế nào nên ta sử dụng một ma trận gọi là *confusion matrix*.

Bài toán này có 3 class (tích cực, tiêu cực và trung tính) nên sẽ có True/False Positive, True/False Negative, True/False Neutral.

True label

-1	True Negative	False Neutral	False Positive
0	False Negative	True Neutral	False Positive
1	False Negative	False Neutral	True Positive
	-1	0	1

Predicted label

Confusion Matrix

$$\text{Precision}_1 = \frac{\text{True Negative}}{\text{True Negative} + \text{False Negative} + \text{False Negative}}$$

$$\text{Precision}_2 = \frac{\text{True Neutral}}{\text{True Neutral} + \text{False Neutral} + \text{False Neutral}}$$

$$\text{Precision}_3 = \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive} + \text{False Positive}}$$

$$\text{Recall}_1 = \frac{\text{True Negative}}{\text{True Negative} + \text{False Neutral} + \text{False Positive}}$$

$$\text{Recall}_2 = \frac{\text{True Neutral}}{\text{True Neutral} + \text{False Negative} + \text{False Positive}}$$

$$\text{Recall}_3 = \frac{\text{True Positive}}{\text{True Positive} + \text{False Neutral} + \text{False Negative}}$$

$$\text{F1-Score}_1 = \frac{2 \times (\text{Precision}_1 + \text{Recall}_1)}{\text{Precision}_1 + \text{Recall}_1}$$

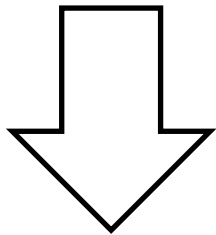
$$\text{F1-Score}_2 = \frac{2 \times (\text{Precision}_2 + \text{Recall}_2)}{\text{Precision}_2 + \text{Recall}_2}$$

$$\text{F1-Score}_3 = \frac{2 \times (\text{Precision}_3 + \text{Recall}_3)}{\text{Precision}_3 + \text{Recall}_3}$$

$$\textbf{F1-Score} = \frac{\text{F1-Score}_1 + \text{F1-Score}_2 + \text{F1-Score}_3}{3}$$

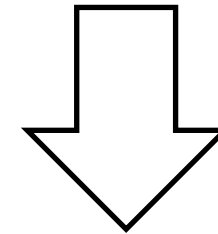
Xây dựng và huấn luyện model

Model SVC



Model SVC
Train score: 0.9079166666666667
Test score: 0.745
F1 score: 0.7447209222834282

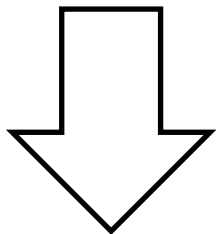
Multinomial Naïve Bayes



Model MultinomialNB
Train score: 0.8495833333333334
Test score: 0.7433333333333333
F1 score: 0.7414419125535217

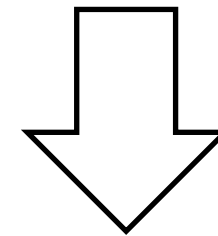
Xây dựng và huấn luyện model

Logistic Regression



```
Model LogisticRegression  
Train score: 0.8554166666666667  
Test score: 0.735  
F1 score: 0.7304342878113371
```

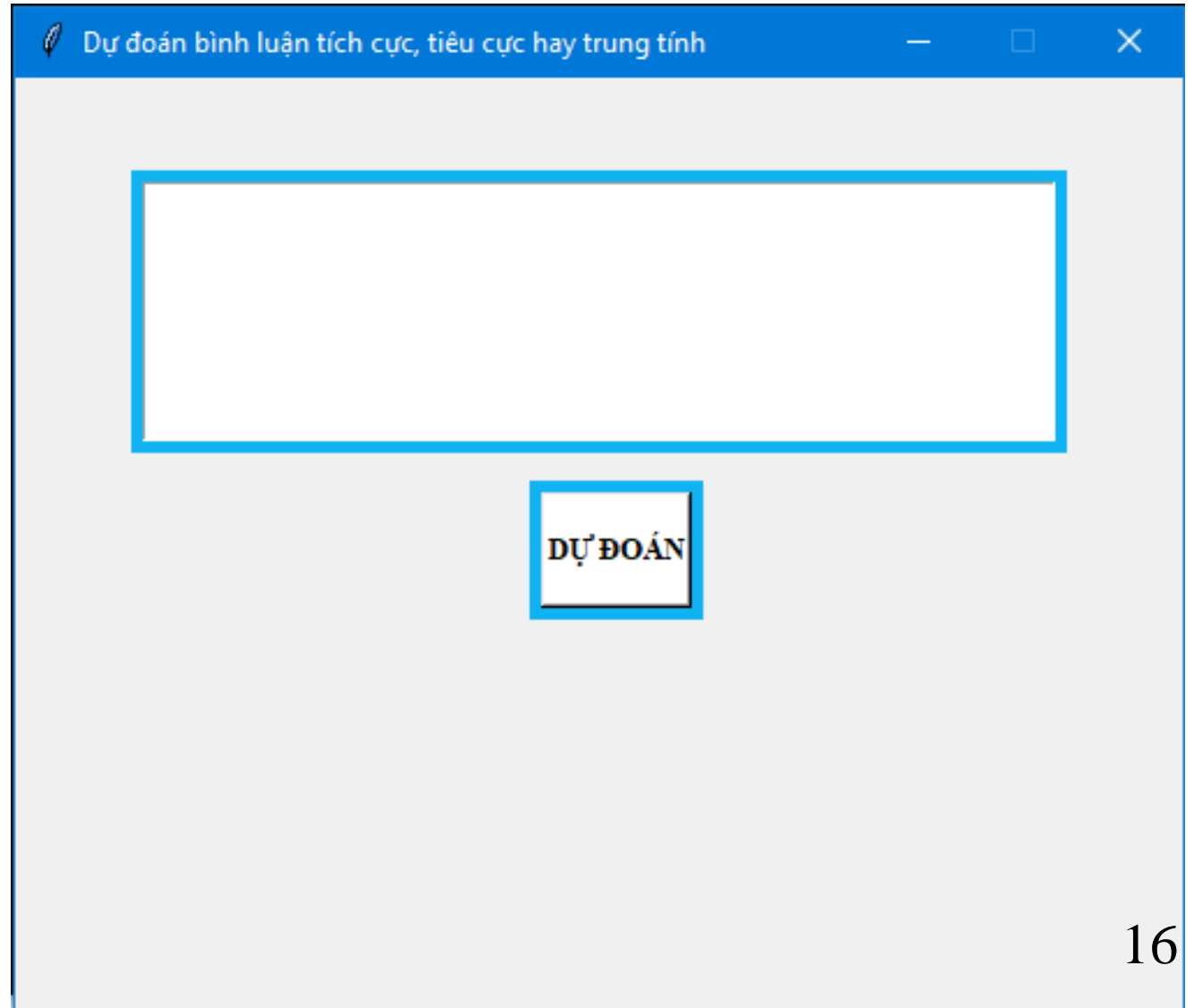
Random Forest



```
Model RandomForestClassifier  
Train score: 0.99125  
Test score: 0.65  
F1 score: 0.6488799173754572
```

Thiết kế giao diện

Sử dụng thư viện **Tkinter** – một package trong python có chứa module Tk hỗ trợ cho việc lập trình GUI (Graphical User Interface)



The background features a light blue-to-green gradient. On the left, there are several overlapping, wavy, light blue shapes that curve upwards and to the right. On the right side, there are similar wavy shapes in a light green color, curving upwards and to the left.

DEMO

♡ Thanks for listening ♡