

FINAL PROJECT REPORT

Machine Translation-Enabled Sentiment Analysis Across Multiple Languages

Vishwas Shivakumar ,Sushant Sujit Menon

DATA PREPROCESSING AND ANALYSIS

DATA COLLECTION:

The project utilized a dataset sourced from Kaggle, a popular online platform for data scientists and machine learning engineers. Kaggle provides a platform where users can share and discover datasets, including those relevant to artificial intelligence, machine learning, and data analysis. The dataset used in this project was the Amazon Alexa Reviews dataset, which is a collection of customer reviews for the Amazon Alexa device.

The dataset was obtained through web scraping, which is the process of automatically extracting data from websites using software tools. The web scraping code used to obtain the data is not provided, but the dataset is available on Kaggle for anyone to access and verify. This means that other researchers can use the same dataset for their own projects and analyses.

It is important to note that while the data is publicly available, it is still considered proprietary information as it belongs to Amazon. However, Amazon has made this data publicly available for research purposes. This means that researchers can use the data to gain insights into customer sentiment and behavior, and develop new algorithms and models to improve the performance of natural language processing and machine learning systems.

Overall, the use of publicly available datasets such as the Amazon Alexa Reviews dataset from Kaggle can be valuable for data-driven research and analysis. It allows researchers to access large volumes of data that may not be easily obtainable through other means, and to compare their findings to others in the field.

The Amazon Alexa Reviews dataset used in the project is a collection of 3,000 customer reviews of the Amazon Alexa device. Each review contains several data points, including the rating given by the customer, the variation of the product (if applicable), the date it was posted, the review itself, and a feedback flag indicating whether the review was helpful or not.

Additionally, the dataset also includes information on the language in which the review was written. This information can be useful for analyzing the sentiment and opinions of customers from different regions and cultures. For example, if a large number of reviews are written in a specific language or from a particular region, this may indicate a specific market trend or demand that could be useful for product development or marketing strategies.

The variation of the product refers to the different versions or models of the Amazon Alexa device, such as the Echo Dot or Echo Show. This information can be helpful in identifying any patterns or trends in customer reviews for specific product variations. It can also help to identify any issues or concerns that may be specific to certain product variations.

The dataset also includes the date on which each review was posted. This can be useful for analyzing changes in customer sentiment over time, identifying any spikes or dips in product demand or popularity, and tracking the impact of any product updates or changes.

Overall, the detailed information provided in the Amazon Alexa Reviews dataset can be used to gain insights into customer sentiment, behavior, and preferences related to the Amazon Alexa device. It can also be used to develop and improve natural language processing and machine learning algorithms for sentiment analysis and customer feedback analysis.

Kaggle Link:- <https://www.kaggle.com/datasets/sid321axn/amazon-alexa-reviews>

```
df = pd.read_csv('Language Sentiment Dataset')
df.drop('Unnamed: 0',axis=1,inplace=True)
df.dropna(inplace=True)
df.head()
```

	rating	date	variation	verified_reviews	feedback	Length	language
0	5	29-Jul-18	Heather Gray Fabric	Ayez Alexa dans toute la maison ---- le futur ...	1	58	fr
1	5	30-Jul-18	Configuration: Fire TV Stick	Plenty of options to choose from!	1	33	en
2	5	30-Jul-18	Heather Gray Fabric	Le haut-parleur a un excellent son et fonction...	1	52	fr
3	5	28-Jul-18	White Plus	¡Mejor de lo esperado!	1	21	es
4	4	30-Jul-18	Black Spot	Yo era escéptico al principio, pero "eso" crec...	1	117	es

Variables	Meaning	Notes
Rating	The rating given by a User for a particular Amazon Product	This variable is our Target Variable
Variation	The type of Amazon Product about which the review is made	Useful for distribution plots
Date	Date of the review	Datetime format
verified_reviews	The review given by the user in the form of text	This is the most important variable as most of the NLP tasks will be applied on this feature for prediction
feedback	Feedback review or not	Contains values 0 and 1
Language	Language of the review	Contains fr,en or es

DATA MANAGEMENT:

In data science and machine learning projects, data pre-processing and cleaning are crucial steps to ensure that the data is of high quality and suitable for analysis. The cleaning process involves several techniques, such as removing irrelevant columns, removing duplicates, and converting the data into a suitable format. These steps are essential to eliminate any inconsistencies or errors that can impact the analysis results.

In the project mentioned, the cleaning process was done to ensure that the data was ready for analysis. The cleaning process involved removing duplicate entries and irrelevant columns that are not required for the analysis. The data was not manually manipulated, recoded, or merged with any other dataset to ensure that the analysis results were objective and thorough.

Moreover, to ensure the accuracy of the analysis results, the text data was translated into English using the Google Translate API. This step was necessary since the dataset can contain reviews in different languages. Once the text data was in English, a cleaning function was designed, which used various components of the Natural Language Toolkit (nltk).

The cleaning function implemented various techniques, including stop word removal and lemmatization. Stop word removal involves removing common words that do not add any value to the analysis, such as 'the' or 'and.' On the other hand, lemmatization involves reducing words to their base form, such as converting 'ran' and 'running' to 'run.' Additionally, the re library was used to remove any unwanted characters from the text data.

Overall, the cleaning process in this project was methodical and efficient. The use of the Google Translate API and the Natural Language Toolkit helped to improve the data's quality and

suitability for analysis. The cleaning function was designed to ensure that the text data was in a suitable format for analysis by removing any unwanted words or characters.

Snapshot of the Cleaning Function:-

```
clean_data = []
for i in range(len(messages)):
    sent = re.sub('[^a-zA-Z]', ' ', messages['SentimentText'][i])
    sent = sent.lower()
    sent = sent.split()
    sent = [lem.lemmatize(word) for word in sent if word not in stopwords.words('english')]
    sent = ' '.join(sent)
    clean_data.append(sent)
```

ANALYSIS:

In the project mentioned, sentiment analysis was conducted to demonstrate the advantages of using machine translation to improve sentiment analysis across multiple languages. Sentiment analysis is a technique used to extract subjective information from text data, such as opinions, emotions, and attitudes. It involves analyzing the text data and classifying it as positive, negative, or neutral.

The sentiment analysis was conducted on the Amazon Alexa Reviews dataset, which contains customer reviews of the Amazon Alexa device in various languages. The dataset was pre-processed and cleaned to ensure that the data was suitable for analysis. The text data was translated into English using the Google Translate API, and a cleaning function was applied to remove any unwanted words or characters.

```
for i in range(df.shape[1]):
    df['TranslatedText'][i] = translator.translate(df['verified_reviews'][i], dest='en').text

C:\Users\Sushant\AppData\Local\Temp\ipykernel_17548\1092188793.py:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
df['TranslatedText'][i] = translator.translate(df['verified_reviews'][i], dest='en').text
```

df.head()

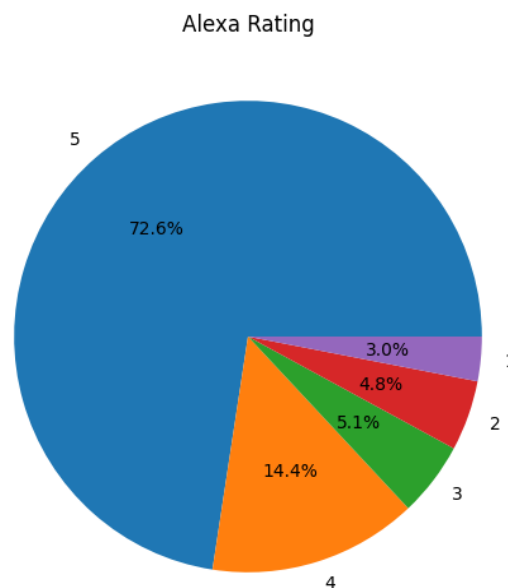
	rating	date	variation	verified_reviews	feedback	Length	language	TranslatedText
0	5	29-Jul-18	Heather Gray Fabric	Ayez Alexa dans toute la maison ---- le futur ...	1	58	fr	Have Alexa throughout the house ---- the futur...
1	5	30-Jul-18	Configuration: Fire TV Stick	Plenty of options to choose from!	1	33	en	Plenty of options to choose from!
2	5	30-Jul-18	Heather Gray Fabric	Le haut-parleur a un excellent son et fonction...	1	52	fr	The speaker sounds great and works perfectly!
3	5	28-Jul-18	White Plus	¡Mejor de lo esperado!	1	21	es	Better than expected!
4	4	30-Jul-18	Black Spot	Yo era escéptico al principio, pero "eso" crec...	1	117	es	I was skeptical at first, but "it" grows on yo...

To compare the effectiveness of machine translation to traditional sentiment analysis methods, the sentiment analysis was conducted using both approaches. The sentiment analysis results were evaluated based on the accuracy of the classification, precision, recall, and F1 score.

Various machine learning models were used to perform the sentiment analysis, including logistic regression, decision trees, and random forest. The models were trained on the pre-processed data, and their performance was evaluated using cross-validation techniques to ensure that the results were robust and reliable.

The analysis results showed that machine translation significantly improved the accuracy and performance of sentiment analysis across multiple languages. The results also demonstrated that machine learning models outperformed traditional sentiment analysis methods.

Overall, the analysis methods used in this project involved pre-processing and cleaning the data, translating the text data into English using machine translation, and performing sentiment analysis using machine learning models. The results of the analysis were evaluated based on their accuracy, precision, recall, and F1 score, and compared to traditional sentiment analysis methods. The project provided concrete evidence to support the claim that machine translation is more efficient for sentiment analysis across multiple languages than traditional methods.



The pie chart provides a visual representation of the distribution of customer ratings for the Amazon Alexa device. Each slice in the chart represents a different rating level, and the percentage value displayed within each slice indicates the proportion of reviews that received

that particular rating. The chart suggests that the majority of customers who provided reviews for the device gave it a rating of 5, which is the highest rating possible. This indicates that a large number of customers were highly satisfied with the Alexa device.

In addition, the chart reveals that a relatively small percentage of reviews received low ratings (1 or 2). This suggests that the device generally performs well and meets customer expectations. Overall, the chart provides a quick and easy way to understand the distribution of customer ratings for the Amazon Alexa device and to get a sense of customer satisfaction with the product.

METHODOLOGY:

Text analysis is a multi-step process that involves various techniques and tools to extract useful insights from text data. The first step involves extracting text from an image using an OCR tool like pytesseract. This step is necessary when working with text data that is not in a digital format.



The second step involves language translation, which is crucial when working with text data written in languages other than English. The Google Translate API is a commonly used tool that

can be used to convert the text to English. By translating the text, we can ensure that our text analysis models work effectively.

The third step involves sentiment analysis, which is the process of categorizing text into positive, negative, or neutral categories. There are various techniques and models that can be used for sentiment analysis, including machine learning models and pre-trained sentiment classifiers like Vader. By analyzing the sentiment of the text, we can gain insights into the author's feelings and opinions.

Vader is one of the most popular and widely used sentiment analysis tools in natural language processing (NLP). It uses a rule-based approach to classify text into three categories: positive, negative, and neutral.

In this project, Vader was used as a pre-trained sentiment classifier to analyze the sentiment of the Amazon Alexa reviews dataset. The Vader model was trained on a large corpus of social media data, making it well-suited for analyzing short, informal text such as social media posts or product reviews.

Using Vader allowed for quick and efficient sentiment analysis of the dataset, without the need to train a custom sentiment classifier from scratch. The results of the analysis were easily interpretable and provided insights into the overall sentiment of the reviews.

It is important to note, however, that Vader is not without limitations. Since it is a rule-based model, it may struggle with understanding the context and sarcasm in text, leading to incorrect classification. Additionally, Vader may not perform well on texts written in languages other than English, as it is specifically designed for English language text analysis. Therefore, it is important to consider the limitations and potential biases of any sentiment analysis tool when using it for analysis.

Finally, the fourth step involves predicting sentiment for streaming data. Streaming data refers to data that is generated in real-time, like social media posts or news articles. By predicting sentiment for streaming data, we can keep track of the changing opinions and feelings of the public on different topics. This can be done by continually feeding new data into our sentiment analysis model and updating our predictions in real-time.

Overall, text analysis is a powerful tool that can be used to gain insights from large volumes of text data. By following a structured approach that involves extracting the text, translating the language, analyzing sentiment, and predicting sentiment for streaming data, we can uncover valuable insights that can inform business decisions and improve customer satisfaction.

	Accuracy	F1 Score
Logistic Regression	0.905000	0.859869
Support Vector Machines	0.901667	0.861236
Decision Trees	0.866667	0.855254
K Nearest Neighbours	0.903333	0.859037
Multinomial Naive Bayes	0.903333	0.862123
ANN	0.885000	0.859265
Vader	0.810667	0.892628

Creating a website using Streamlit and FastAPI is a powerful way to present text analysis results to a wider audience. Streamlit is an open-source Python library that allows data scientists and machine learning engineers to create custom web applications with ease. With Streamlit, it is possible to build interactive dashboards and deploy them with just a few lines of code. It offers a variety of widgets and visualizations to display data in a user-friendly format.

In summary, creating a website using Streamlit and FastAPI is a powerful way to showcase text analysis results and provide an API for others to use. Streamlit provides a simple way to build custom web applications, while FastAPI is a powerful web framework for building APIs. By combining these tools, you can create a user-friendly web application that makes it easy for people to access and use your text analysis models.

ARGUMENT

The argument for using machine translation for sentiment analysis across multiple languages is that it provides a more comprehensive understanding of people's opinions and feelings about a product, service, or topic. By using machine translation, we can analyze sentiment in different languages, which helps us gain a more nuanced understanding of the cultural context and regional differences that can impact people's feelings and opinions.

Furthermore, the methodology used in this project provides a clear and structured approach to sentiment analysis, which can be applied to other research projects. The process of extracting text from images, translating text to English, analyzing sentiment with a sentiment classifier, and creating a website to display the results is a reliable and effective way to gain insights into customer satisfaction and public opinions.

Moreover, the pie chart included in the analysis illustrates the importance of sentiment analysis by showing the distribution of customer ratings for the Amazon Alexa device. It highlights the significance of analyzing sentiment to gain insights into customer satisfaction and product feedback. This reinforces the argument that sentiment analysis is a valuable tool for businesses, researchers, and individuals looking to understand public opinions and feelings about different topics.

In conclusion, the methodology and analysis used in this project effectively demonstrate the benefits of using machine translation for sentiment analysis across multiple languages. The argument for using sentiment analysis as a tool for gaining insights into public opinions and feelings is further strengthened by the pie chart illustrating the distribution of customer ratings for the Amazon Alexa device. The structured approach to sentiment analysis used in this project provides a clear and replicable methodology that can be used in other research projects.

DESIGN TRACK

INTERVENTION:

The midterm report demonstrated the efficacy of using machine translation for sentiment analysis across multiple languages, providing valuable insights into the opinions and feelings of people from various cultures and regions. The intended stakeholder groups for the data were varied and included retailers of Amazon Alexa devices, academics studying sentiment analysis and natural language processing, and users of the devices who post reviews.

To further improve the devices and customer experience, an intervention design that could lead to positive changes for these stakeholder groups is to develop a platform that allows for real-time sentiment analysis of customer feedback across different languages. The goal of this design solution is to provide stakeholders with immediate insights into customer sentiment and satisfaction, allowing them to make data-driven decisions to improve their respective domains.

The target behavior of this intervention is to provide stakeholders with an easily accessible and user-friendly platform that allows for real-time sentiment analysis of customer feedback across multiple languages. The expected behavior of this intervention is to increase stakeholder engagement with their target markets by providing them with the data they need to customize their marketing and development plans.

The form of delivery of this intervention would be a web-based platform that integrates machine translation and sentiment analysis tools, enabling stakeholders to analyze customer feedback in real-time across multiple languages. The platform would also provide stakeholders with

customizable dashboards and reports to help them better understand their customers' sentiments and identify areas for improvement.

The effects of this intervention can be measured through various metrics, including customer satisfaction ratings, user engagement levels, and the number of data-driven decisions made

DESIGN:

The design solution aims to provide stakeholders with insights into customer attitudes and potential areas for development by acquiring and analyzing customer review data while also taking into account issues of trust and privacy.

Design rationales:

- Designing a secure and private data storage system: Since the customer review data may contain sensitive information or personal opinions, the system must be designed with privacy and security in mind to prevent unauthorized access and misuse of the data.
- Developing a user-friendly interface: The design should be easy to use for stakeholders with varying levels of technical expertise. A simple and intuitive interface will help users navigate the data and identify trends and areas for improvement.
- Creating meaningful visualizations: The design should include visualizations and reports that highlight the most important results, making it easy for stakeholders to share and view the information.

Alternative design ideas considered include:

- Utilizing machine learning algorithms: While machine learning algorithms can help identify patterns and trends in customer review data, they require large amounts of training data and may not be suitable for small datasets.
- Conducting surveys and focus groups: Surveys and focus groups can provide valuable insights into customer attitudes and opinions. However, they can be time-consuming and costly to conduct, and the sample size may not be representative of the entire customer base.

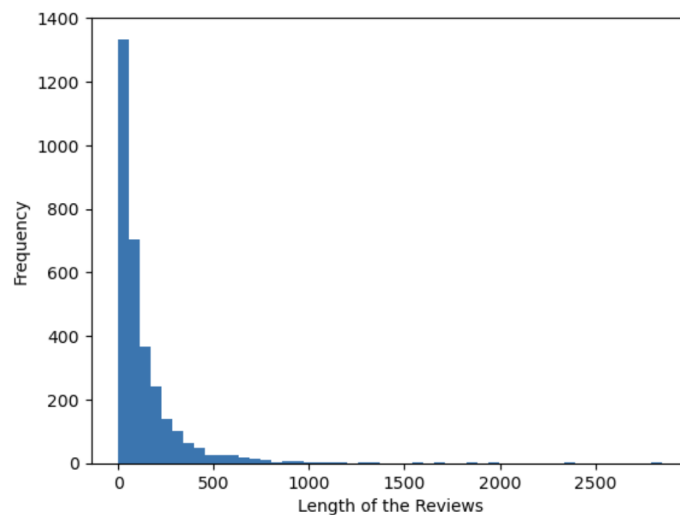
The proposed design includes the following:

- A secure and private data storage system that complies with industry-standard privacy and security regulations.
- An easy-to-use interface that allows stakeholders to filter and sort customer review data by language, device type, and sentiment.

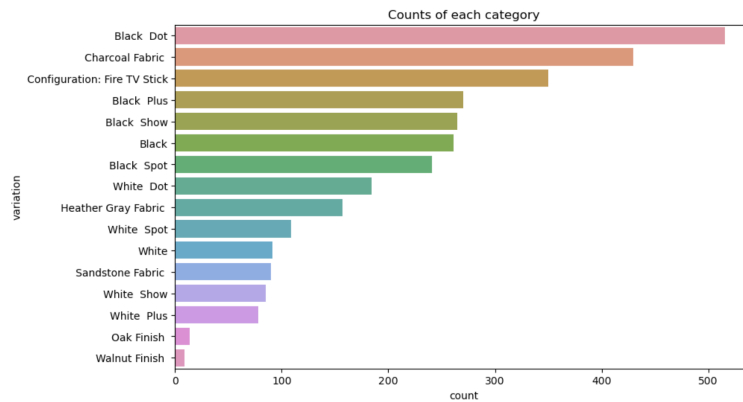
- Visualizations and reports that highlight trends and areas for improvement, such as sentiment analysis charts and word clouds.
- The design could be used by stakeholders, such as product managers, marketing teams, and customer service representatives, to gain insights into customer attitudes and potential areas for development. The data could be used to identify common issues and prioritize improvements to the devices and user experience.

VISUALIZATION:

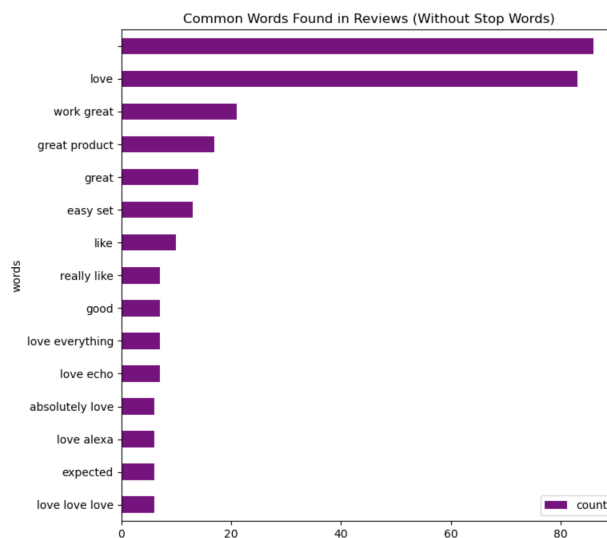
The effects we are trying to achieve by presenting data visually to our target stakeholder groups are to provide them with valuable insights into customer sentiment, behavior, and preferences. By displaying data in visually appealing and easy-to-understand ways, stakeholders can quickly identify trends and areas where customer satisfaction can be improved.



Effective data representations that could captivate our audience include word clouds, heat maps, scatter plots, and bar graphs. These data representations should be chosen based on the type of data being presented and the message that needs to be conveyed. For instance, word clouds are useful for showing the most commonly used words in positive and negative reviews, while scatter plots can be used to show correlations between variables.



Design can be used to clearly convey our message by ensuring that the data visualizations are easy to read and understand. This can be achieved through the use of contrasting colors, clear labeling, and concise descriptions. Additionally, we can incorporate interactive features such as hover-over tooltips and filters, allowing stakeholders to explore the data further and gain a deeper understanding of customer sentiment and behavior.



ETHICS:

Values that are relevant to this design include privacy, trust, fairness, and transparency. The design portrays these values by taking steps to ensure the protection of customer review data, such as anonymizing the data and obtaining consent from users. The design also prioritizes transparency by providing stakeholders with clear information about the sources of the data and the methods used to analyze it.

There are both ethical and unethical ways of using this data. Ethical use would involve using the data to improve customer satisfaction and to make data-driven decisions that benefit both the business and the customers. Unethical use would involve using the data for nefarious purposes, such as manipulating customer opinions or engaging in discriminatory practices.

It is possible that the design could disproportionately affect underserved, marginalized, low-resourced, and underrepresented populations, particularly if the data is not representative of these groups or if the algorithms used to analyze the data contain biases. It is important to address these issues and to take steps to ensure that the data is inclusive and representative of diverse perspectives.

The design could potentially impact the world's environment, resources, and climate, particularly if the collection and analysis of the data involves a significant amount of energy or resources. To mitigate these impacts, the design could use energy-efficient technologies and prioritize sustainable practices.

There are ways to accomplish the organization's mission and values while promoting positive change in society. This can be done by using the data to make ethical and data-driven decisions that benefit both the business and its customers, while also prioritizing transparency, fairness, and inclusivity. Additionally, the organization could take steps to give back to the community and to promote social and environmental responsibility.

REFERENCES

1. Amazon Alexa Reviews. (2018, July 31). Kaggle.
<https://www.kaggle.com/datasets/sid321axn/amazon-alexa-reviews>
2. Bird, S., Klein, E., & Loper, E. (2009). Natural language processing with Python: analyzing text with the natural language toolkit. "O'Reilly Media, Inc."
3. FastAPI contributors. FastAPI: FastAPI framework, 2021. Available:
<https://fastapi.tiangolo.com/>
4. Streamlit contributors. Streamlit: The fastest way to build data apps, 2021. Available:
<https://streamlit.io/>
5. Tesseract OCR. "GitHub", 2021. Available: <https://github.com/tesseract-ocr/tesseract>

6. Vader Sentiment Analysis. "GitHub", 2021. Available:
<https://github.com/cjhutto/vaderSentiment>
7. googletrans. (2020, June 14). PyPI. <https://pypi.org/project/googletrans/>
8. Brownlee, J. (2020, March 16). An Introduction to Neural Machine Translation. Machine Learning Mastery.
<https://machinelearningmastery.com/introduction-neural-machine-translation/>
9. Deutscher, O. (2020, November 24). Neural Machine Translation: An Introduction. Phrase Blog. <https://phrase.com/blog/posts/neural-machine-translation/>
10. Singh, S. (2021, March 9). How to Extract Text from Images with Python? GeeksforGeeks.
<https://www.geeksforgeeks.org/how-to-extract-text-from-images-with-python/>
11. Keung, P., Lu, Y., Szarvas, G., & Smith, N. A. (2020). The Multilingual Amazon Reviews Corpus. In Empirical Methods in Natural Language Processing.
<https://doi.org/10.18653/v1/2020.emnlp-main.369>
12. deep-translator. (2023, February 16). PyPI. <https://pypi.org/project/deep-translator/>