# Intent Classifier with Facebook fastText

Facebook Developer Circle, Malang

22 February 2017

Bayu Aldi Yansyah

Data Scientist at Sale Stock

Developer Circles
from facebook

UB FILKOM
FAKULTAS ILMU KOMPUTER

SALE STOCK
ENGINEERING

# OUR GOALS

1.  Understand how Machine Learning is applied in Messenger bot development.
2.  Understand what is fastText and why it is important.
3.  Able to create intent classifier model with fastText.

Aska is our Messenger bot
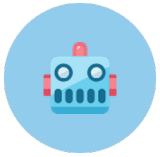
# HUMAN

Understand the context and extract the available information easily.

Intent                  : Book Flight Ticket

Departure Location      : Surabaya

Departure Time          : Minggu pagi

Arrival Location        : Jakarta

Quantity                : 2

Aku pesan tiket surabaya-jakarta buat minggu pagi ya

Untuk 2 org

SALE STOCK
ENGINEERING

# OUR BOT

He need to know the intent of the chat first and extract the available information.

**1.** What is the intent?    **2.** Extract information

Aku pesan tiket surabaya-jakarta buat minggu pagi ya

Untuk 2 org

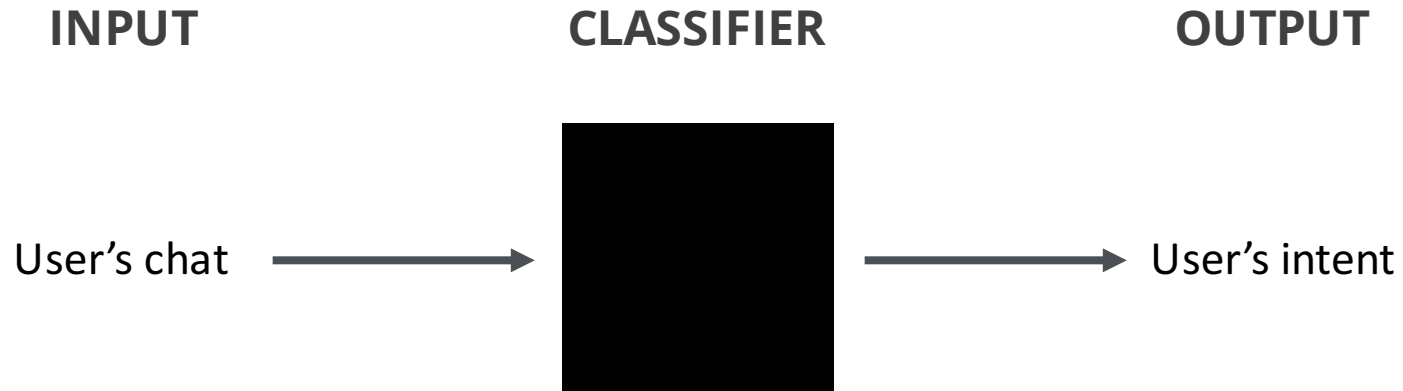| Book Flight Ticket |
| Book Hotel |
| Ask a Question |
| Buy a Product |
| ... |

DepLoc      : Surabaya

DepTime    : Minggu pagi

ArrLoc      : Jakarta

Quantity    : 2

# INTENT CLASSIFIER

**INPUT**　　　　　　　**CLASSIFIER**　　　　　　**OUTPUT**

User's chat　→　[ ]　→　User's intent

Example of user's intent: Book a flight ticket, book a hotel, buy a product and asking a specific questions.

# INTENT CLASSIFIER

- Intent Classifier can be seen as Document Classification/Text Classification task

- Text categorization (a.k.a. text classification) is the task of assigning predefined categories to free-text documents.

- We can use Machine Learning approach to solve this task.

- Popular method is Supervised learning.

# SUPERVISED LEARNING

| DATA | | ALGORITHM | MODEL | |
|---|---|---|---|---|
| INPUT | OUTPUT | TRAINING | EVALUATE | DEPLOY |

- Example data: chat message as an input and the label of intent as the output.

- Popular algorithms for text classification:
    - Linear classifiers: Support Vector Machine (Joachims, 1998)
    - Convolutional Neural Networks based: Kim 2014 & Zhang and LeCun 2015

- Example of the metrics to evaluate: Precision & Recall

SALE STOCK
ENGINEERING

# SUPERVISED LEARNING

- Neural Network based models achieve very good performance in practice, but it is **time consuming to train and test**. It need large resources for large datasets.

- Linear classifiers are simple and often obtain state-of-the-art performance if the right features is used. It doesn't need large resources resource. However, **linear classifiers do not share parameters among features and classes**.

- fastText is trying to solve these problems.

# FASTTEXT

- fastText is a C++ library for efficient learning of word representations and sentence classification.
- Word Representation learning:
    - Continuous Bag-Of-Words (CBOW)
    - Continuous Skip-gram
- Text Classification: A simple linear model where using averaged word representation as text representation.

# WORD REPRESENTATION LEARNING

- Word representation learning or word embedding is a technique to represent a word as a vector.
- The result of word embedding frequently referred as "word vector" or "distributed representation of words".
- There are 3 main approaches to word embedding:
    1. Neural Networks model based
    2. Dimensionality reduction based
    3. Probabilistic model based
- The idea of these approaches are to learn vector representations of words in an unsupervised manner.
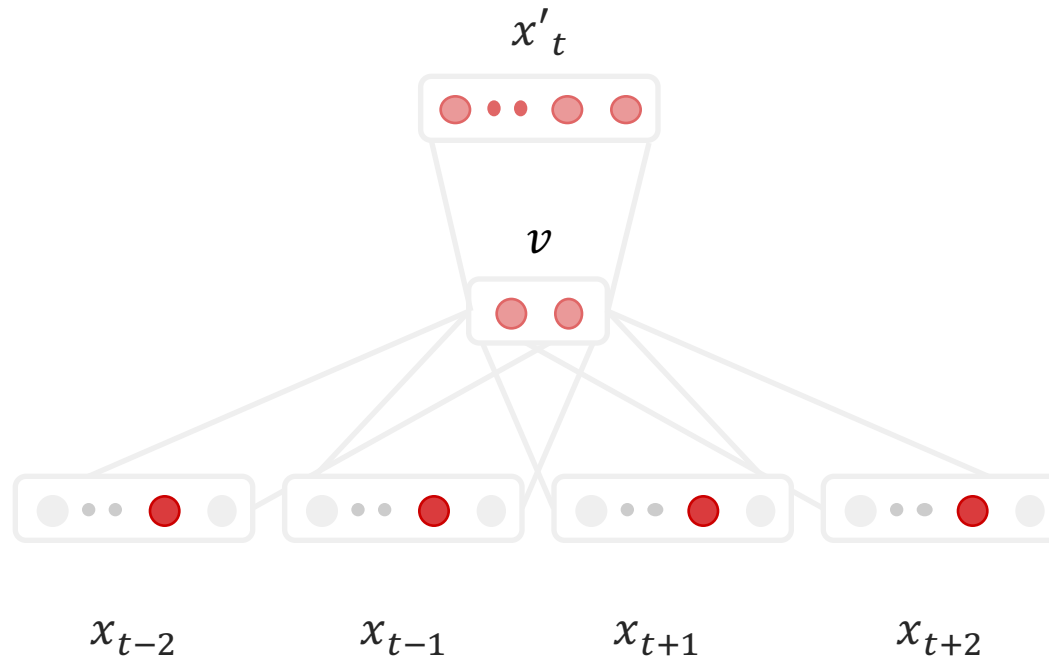
# CONTINOUS BAG-OF-WORDS

| ... | Keren | Sale | | bisa | bayar | dirumah | ... |
|-----|-------|------|-|------|-------|---------|-----|
| | $w_{t-2}$ | $w_{t-1}$ | $w_t$ | $w_{t+1}$ | $w_{t+2}$ | | |

- A simple neural networks model.
- The training data is a sequence of words $w_1, w_2, ..., w_T$ for $w_t \in V$
- The model is trying predict the word $w_t$ based on the surrounding context ($n$ words from left: $w_{t-1}, w_{t-2}$ and $n$ words from the right: $w_{t-1}, w_{t-2}$).
- There are no hidden layer in this model.
- Projection layer is averaged across input words.

# CONTINOUS BAG-OF-WORDS



Visualization of CBOW model

# CONTINOUS SKIP-GRAM

| ... | Keren | | | bisa | | | ... |

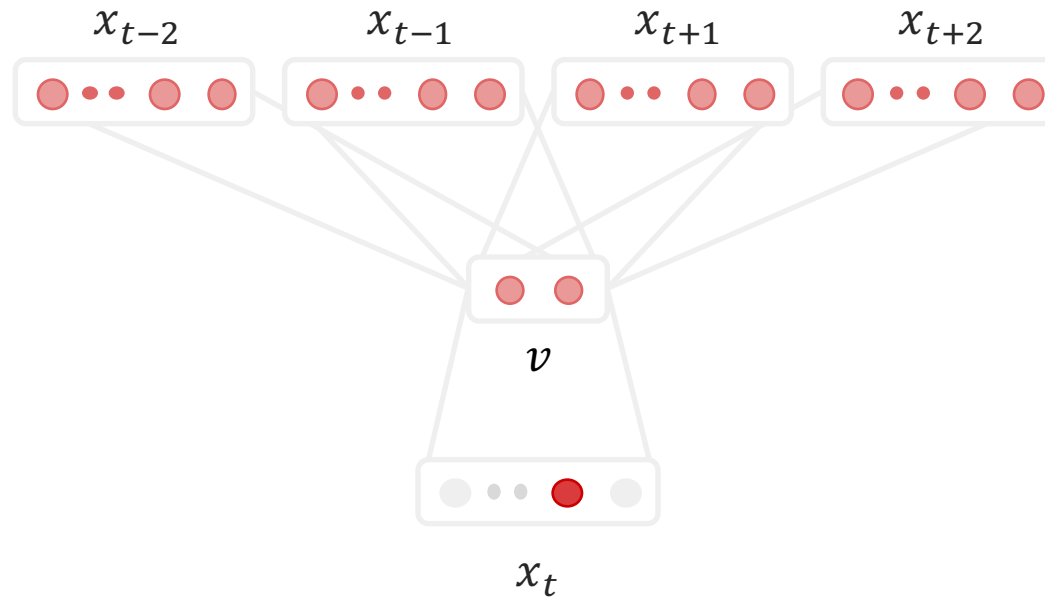$$w_{t-2} \quad w_{t-1} \quad w_t \quad w_{t+1} \quad w_{t+2}$$

- The training data is a sequence of words $w_1, w_2, \ldots, w_T$ for $w_t \in V$
- The model is trying predict the surrounding context ($n$ words from left: $w_{t-1}, w_{t-2}$ and $n$ words from the right: $w_{t-1}, w_{t-2}$) based on the word $w_t$.

# CONTINOUS SKIP-GRAM

Visualization of Skip-gram model

# WORD REPRESENTATION LEARNING

More detailed information about word embedding available at

"Clustering semantically similar words"

https://github.com/pyk/talks

# LINEAR CLASSIFIER

- The architecture is similar with Continuous Bag-of-Words model, where the middle of words is replaced by the label.
- The word representations are averaged into text representation, which is in turn fed to a linear classifier.

# FASTTEXT

Step by step to create intent classifier model with fastText.

It's super easy.

Install the `fasttext(1)` first

https://github.com/facebookresearch/fasttext

(Only support Linux and OSX)

# FASTTEXT

Prepare the training data:

1.  Convert the intent as numeric.

    For example: Booking flight = 0, Buy product = 1

2.  Preprocess the data: Tokenization, remove symbols etc

3.  Create the training data with the following format:

    `prefix_labelnum token1 token2 token3 …`

# FASTTEXT

Create the model:

1. Run:

```
./fasttext supervised -input train.dat \
        -output intent_classifier -dim 100 -lr 0.1 \
        -wordNgrams 2 -minCount 1 -bucket 2000000 \
        -epoch 5 -thread 4
```

- `train.dat` is the training data
- output model: `intent_classifier.bin`

# DEPLOY

If you are Python programmer, you can use fastText.py

https://github.com/salestock/fastText.py to deploy the model.

- Create a API service with Flask

- Load & serve the model

- Create new endpoint for the inference step

# REFERENCES

- Enriching Word Vectors with Subword Information by Bojanowski et al (2016) https://arxiv.org/abs/1607.04606
- Bag of Tricks for Efficient Text Classification by Joulin et al (2016) https://arxiv.org/abs/1607.01759

# THANKS

Questions? bay@artificialintelligence.id

Slide: https://github.com/pyk/talks

Sale Stock: https://careers.salestock.io

SALE STOCK
ENGINEERING