# A Gentle Introduction to Natural Language Processing
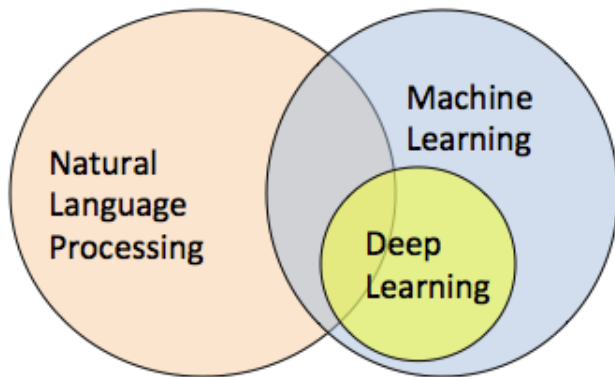
Natalia Klyueva

The Hong Kong Polytechnic University
Chinese and Bilingual Studies

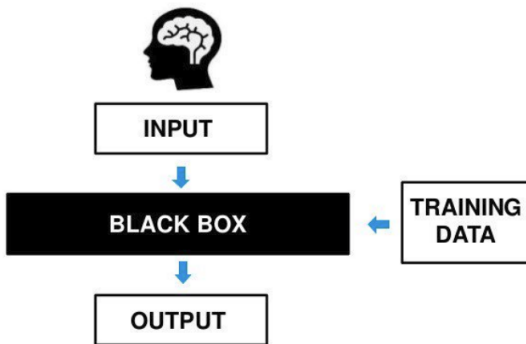natalia.klyueva@polyu.edu.hk

June 25, 2018

picture taken from https://rutumulkar.com/blog/2016/NLP-ML

## MACHINE LEARNING



picture taken from https://twitter.com/KirkDBorne

# Canonical Examples of ML data

| Size of House | Lot Size (acre) | # of Bedrooms | # of Bathrooms | Price of House |
|---|---|---|---|---|
| 950 | 2.5 | 2 | 1 | $127,325 |
| 1,535 | 1.5 | 2 | 2 | $156,570 |
| 1,605 | 2.25 | 3 | 1.5 | $158,895 |
| 1,905 | 2.5 | 2 | 1.5 | $200,025 |
| 2,057 | 2.25 | 3 | 2 | $230,384 |
| 2,227 | 2.75 | 3 | 2 | $233,835 |
| 3,150 | 1 | 4 | 2 | $261,420 |
| 3,620 | 3 | 4 | 3 | $433,500 |

| | A | B | C | D | E | F | G | |
|---|---|---|---|---|---|---|---|---|
| 1 | age | job | marital | education | default | housing | loan | con |
| 2 | 56 | housemaid | married | basic.4y | no | no | no | tel |
| 3 | 57 | services | married | high.scho | unknown | no | no | tel |
| 4 | 37 | services | married | high.scho | no | yes | no | tel |
| 5 | 40 | admin. | married | basic.6y | no | no | no | tel |
| 6 | 56 | services | married | high.scho | no | no | yes | tel |
| 7 | 45 | services | married | basic.9y | unknown | no | no | tel |
| 8 | 59 | admin. | married | professio | no | no | no | tel |
| 9 | 41 | blue-collar | married | unknown | unknown | no | no | tel |
| 10 | 24 | technician | single | professio | no | yes | no | tel |
| 11 | 25 | services | single | high.scho | no | yes | no | tel |
| 12 | 41 | blue-collar | married | unknown | no | no | tel |

# Natural Language Processing

- The same as in Machine Learning tasks above, but...
- features are: words, sentences, paragraphs, documents etc.
- Canonical NLP tasks:
    - search engines
    - Machine Translation
    - face recognition systems
    - chatbots
    - natural language generation
    - opinion mining (sentiment analysis)

How can machine understand meaning? Let's start with the
basic notion - 'sense'

- How can machine understand meaning of a word?
- image recognition task: how to relate from an image (real
  world object) to a word/sentence
  - "Is this a cat on the picture?"
  - FB: "Do you want to tag FRIENDS_NAME
    FRIENDS_SURNAME?"
- given a picture, give a description of an object

# Image recognition: near-NLP task

What object do we(humans) and they(machines) refer to?



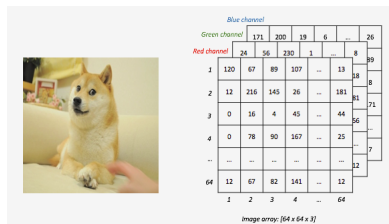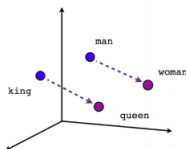| Machine-generated (but turker prefered) | a group of motorcycles parked next to a motorcycle |
|---|---|
| Human-annotated (but turker not prefered) | two girls wearing are wearing short skirts and one of them sits on a motorcycle while the other stands nearby |

- We identify object by its features (is furry, has 4 paws, tail, special form of ears etc...)? Or by observations?
- Representation for a machine: all images are 64x64 pixels, and each pixel is a certain value of RGB
- Input: array of numbers



- Output: 97% it is a dog
- That is how to represent a picture.
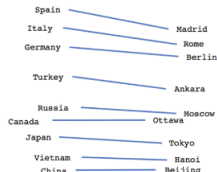- And how do we represent words?

# Word representation

- Obviously, we need numbers to make calculations easier and faster
- We need to assign some numerical representation for each word
- I->id1; see->id2; a->id3 cat->id4 ... Is it a good one????
- Small revolution in NLP: word2vec
  https://www.tensorflow.org/tutorials/word2vec
- Input: words,sentences,text,documents. Output: vector



Male-Female          Verb tense          Country-Capital

- '2+2' of NLP: king−man+women=queen
- https://projector.tensorflow.org/

- Predict the word:



- For doc2vec - the same principle, with document information included
  https://cs.stanford.edu/~quocle/paragraph_vector.pdf

- Likes/dislikes on facebook, wechat; swipes right/left; comments on products/services on internet
- User generated content. Comes for free and in large quantity
- The questions are: "Do you like this or dislike this?", "What are your emotions"

- Sentiment analysis of the text
- On the user generated content, you can perform various machine learning experiments and model "sentiment"
- Input: review -> MODEL -> Output:sentiment/emotion

# Our data: openrice

- Sentiment analysis of the text
- On the user generated content, you can perform various machine learning experiments and model "sentiment"
- Input: review -> MODEL -> Output:sentiment/emotion
- Openrice scrapped
- About 360,000 comments; 5,000 restaurant IDs

# Data formatting

## Our data: openrice

- The simple (baseline) task: given a comment, predict the 'value' (1-5)
- Input: document embeddings, classifier from sklearn
- Demo: https://github.com/natalink/openrice_ annotations/blob/master/clf_openrice.ipynb
- For the small sample (12506, on github) that I put online, the accuracy of the model was about 0.4
- If we select only very positive(rank 5) or negative (rank 1) comments, the accuracy was 0.7.
- For the whole data with rank value [1, 5] (26549 reviews) accuracy was around 0.850

```python
model = Doc2Vec.load("ALL_picc_doc2vec.vec")
print("\n","Most similar to 'good' ", model.most_similar('好')) # just to test the model
print ("most similar to  'Hong Kong' ", model.most_similar('香港') )
print ("most similar to 'kowloon' ", model.most_similar('九龍') )
print("Most similar to 'dimsum': ", model.most_similar('點心')) # just to test the model
print ("china + hk - england =???", model.most_similar(positive=['中國', '香港'], negative=['英國']) )
print ("spicy + Sechuan - Hong Kong =???", model.most_similar(positive=['四川', '辣'], negative=['香港']))
```

```
 Most similar to 'good'  [('幾', 0.6377090215682983), ('好好', 0.6115480661392212), ('幾好', 0.582214593887329
1), ('呀', 0.5795080661773682), ('超級', 0.5684177279472351), ('夠', 0.5599690675735474), ('非常', 0.552945137
0239258), ('仲好', 0.5522783398628235), ('食落夠', 0.545662522315979), ('勁', 0.5417361259460449)]
most similar to  'Hong Kong'  [('台灣', 0.7651838064193726), ('小店', 0.7570662498474121), ('中菜', 0.75396239
75753784), ('一間', 0.74861741065979), ('天堂', 0.737130880355835), ('連鎖', 0.7272865772247314), ('老店', 0.72
64119982719421), ('近年', 0.7263711094856262), ('館', 0.7198368310928345), ('食店', 0.7153961658477783)]
most similar to 'kowloon'  [('海港', 0.9411194324493408), ('灣仔', 0.9407057166099548), ('上環', 0.94050866365
43274), ('大圍', 0.939861536026001), ('北角', 0.9396387338638306), ('邨', 0.9392281770706177), ('萬', 0.931288
7191772461), ('兆', 0.9311912655830383), ('中環', 0.9305603540810908), ('天后', 0.9292474985122681)]
Most similar to 'dimsum':  [('包類', 0.7647587060928345), ('食品', 0.7600960731506348), ('小菜', 0.74586737155
91431), ('午市', 0.7368004322052002), ('其式', 0.7218303680419922), ('粥品', 0.7156221866607666), ('熟食', 0.71
18027806282043), ('老式', 0.6852641701698303), ('供應', 0.6790927648544312), ('廣東', 0.67188446559524536)]
china + hk - england =???  [('小店', 0.721519947052002), ('懷舊', 0.7202714681625366), ('舊式', 0.7020691037178
04), ('高檔', 0.6987070441246033), ('街坊', 0.6961661577224731), ('格局', 0.6949720978736877), ('裝修', 0.68691
36095046997), ('風格', 0.6842846274375916), ('老店', 0.683087944984436), ('酒吧', 0.6787234544754028)]
spicy + Sechuan - Hong Kong =???  [('酸菜', 0.8346952795982361), ('酸辣', 0.8040406107902527), ('麻辣', 0.78325
20008087158), ('小辣', 0.7822674512863159), ('胡椒', 0.7609094977378845), ('香辣', 0.7529984712600708), ('白肉'
, 0.739170491695404), ('紅油', 0.7313030958175659), ('鳳爪', 0.7268850803375244), ('少辣', 0.7263926267623901)]
```