

# Lec 1

## Introduction

### # NLP

What is NLP? Make computers understand & generate human languages

For...

- Human-computer
- Human-human } interaction
- Understand languages (syntax, etc.)

Applications

- Answer questions
- Translate
- Aid scientific research
  - comp social sci corpus analysis paper
  - agency / power analysis of characters in films

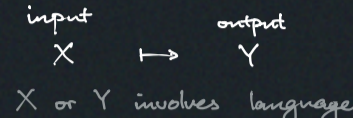
Models can make mistake in all of these

Class

- Build NLP systems
- When do these systems fail
- How to improve

### # NLP System Overview

General framework:



Ex.

- lang1  $\mapsto$  lang2
- question + choices  $\mapsto$  answer
- query + documents  $\mapsto$  relevant documents
- text  $\mapsto$  label
- image  $\mapsto$  caption

Methods

- Rule based
- Prompting
  - prompt a lang model, no training
- Fine-tuning

Data

- Rules/prompt  $\rightarrow$  maybe no data use intuition
  - spot-check to fix rules/prompt
- Rules/prompt with evaluation  $\rightarrow$  dev set + test set
  - 200-2000 examples
- Fine-tuning  $\rightarrow$  the more the better in general:
  - performance  $\uparrow$  linear
  - when data  $\uparrow$  quadratic

### # Make rule-based sentiment analysis \* Bad idea

Task: product review  $\mapsto$  { positive, neutral, negative }

1. Feature extraction

$$\vec{h} = f(x)$$

2. Score calculations

$$s = w \cdot \vec{h}$$

$\vec{s} = W\vec{h}$   
 $\uparrow$  weight matrix

3. Decision function

$$\hat{y} = \text{decide}(s)$$

good words = [ love, good, ... ]  
 bad words = [ hate, sad, ... ] } - designing this is complicated

bias = 1

$$\text{score} = 1.0 * \text{count good words} + 1.0 * \text{count bad words} + 0.5 * \text{bias}$$

not that simple

$$\hat{y} = \begin{cases} \text{positive} & \text{if score} > 0 \\ \text{neutral} & \text{if score} = 0 \\ \text{negative} & \text{if score} < 0 \end{cases}$$

To improve: comprehensive analysis

Issues:

- low-frequency words ...
  - ↳ manually add them
  - ↳ sentiment dictionary
  - ↳ use root form of words i.e. morphological analysis
- negation "not bad"  $\mapsto$  neutral / positive?
  - ↳ syntactic analysis to see what's negated
- Metaphor
- Someone sending reviews in another language
  - ↳ learn Japanese (?)

### # Machine learning approach - BOW

→ Learn feature extractor or weight function

Fixed feature extractor

I hate this movie

$$\begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} =: f(x) = \vec{h}$$

$$w \cdot \vec{h} =: \vec{s}$$

learn this

Structured perception training algorithm

for  $(x, y)$  in train set

if  $y = \text{neutral}$ : skip

$$\vec{h} = f(x)$$

$$\hat{y} = \text{predict}(\vec{h})$$

if  $\hat{y} \neq y$ :

upweight / downweight weight matrix

BOW problems

- Conjugation
- Word similarity
- Combined feature
  - love
  - don't hate
  - don't love
  - hate
- Sentence structure "but"

### # Neural network

Theoretically NN can model & solve any problem

### # Assignments

1. Build LLaMa
2. NLP task from scratch for specified task
  - ↳ collect data
  - ↳ modelling
  - ↳ evaluation
3. Survey + re-implementation
  - ↳ lit review
  - ↳ reimplement an NLP paper
4. New research
  - ↳ improve performance
  - ↳ apply technique to new task