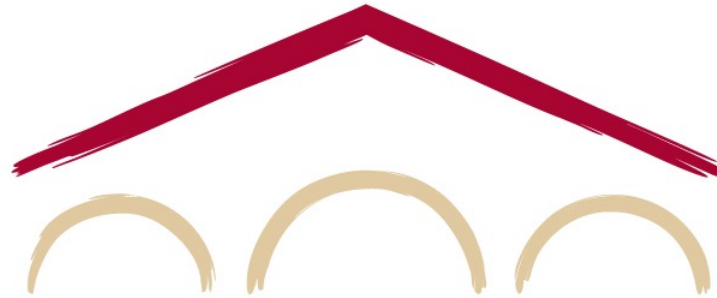# Natural Language Processing with Deep Learning
# CS224N/Ling284

Diyi Yang / Tatsunori Hashimoto

**Lecture 1: Introduction and Word Vectors**

# Lecture Plan

Lecture 1: Introduction and Word Vectors

1. The course (10 mins)
2. Human language and word meaning (15 mins)
3. Word2vec introduction (15 mins)
4. Word2vec objective function gradients (25 mins)
5. Optimization basics (5 mins)
6. Looking at word vectors (10 mins or less)

Key learning today: The (astounding!) result that word meaning can be represented rather well by a (high-dimensional) vector of real numbers

# Course logistics in brief

- Instructor: Diyi Yang, Tatsunori Hashimoto

- Head TA: Nelson Liu

- Course Manager: John Cho

- TAs: Many wonderful people! See website

- Time: Tu/Th 4:30–5:50 Pacific time, Nvidia Aud. ($\rightarrow$ video)

- Email list: cs224n-win2324-staff@lists.stanford.edu

- We've put a lot of other important information on the class webpage. Please read it!
  - http://cs224n.stanford.edu/
    a.k.a., http://www.stanford.edu/class/cs224n/
  - TAs, syllabus, help sessions/office hours, Ed (for all course questions/discussion)
    - Office hours start **Wednesday!**
    - Python/numpy and then PyTorch tutorials: First two Fridays. First is 4:30-5:20, Skilling Auditorium.
  - Slide PDFs uploaded before each lecture

## Instructors

Diyi Yang

Tatsunori Hashimoto

## Course Manager

John Cho

## Teaching Assistants

Nelson Liu (Head TA)

Anirudh Sriram

Annabelle Wang

Arvind Mahankali

Caleb Ziems

Cheng Chang

David Lim

Hamza El Boudali

Heidi Zhang

Kaylee Burns

Olivia Lee

Soumya Chatterjee

Tathagat Verma

Tianyi Zhang

Tim Dai

Tony Wang

Yuan Gao

Yuhui Zhang

# What do we hope to teach?  (A.k.a. "learning goals")

1.  The foundations of the effective modern methods for deep learning applied to NLP
    *   Basics first, then key methods used in NLP in 2023: Word vectors, feed-forward networks, recurrent networks, attention, encoder-decoder models, transformers, pretraining, post-training (RLHF, SFT), efficient adaptation, benchmarking and evaluation, human centered NLP, etc.

2.  A big picture understanding of human languages and the difficulties in understanding and producing them via computers

3.  An understanding of and **ability to build systems** (in PyTorch) for some of the major problems in NLP:
    *   Word meaning, dependency parsing, machine translation, question answering

# Course work and grading policy

- 5 x 1-week Assignments: 6% + 4 x 12%: 54%
  - **HW1 is released today! Due next Tuesday! At 4:30 p.m.**
  - Submitted to Gradescope in Canvas (i.e., using @stanford.edu email for your Gradescope account)
- Final Default or Custom Course Project (1–3 people): 43%
  - Project proposal: 5%, milestone: 5%, poster or web summary: 3%, report: 30%
- Participation: 3%
  - Guest lecture reactions, Ed, course evals, karma – see website!
- Late day policy
  - 6 free late days; afterwards, 1% off course grade per day late
  - Assignments not accepted more than 3 days late per assignment unless given permission in advance

# Course work and grading policy

- Collaboration policy:
  - Please read the website and the Honor Code! Understand allowed collaboration and how to document it: Don't take code off the web; acknowledge working with other students; write your own assignment solutions
- AI tools policy
  - Must independently submit their solutions to CS224N homework
  - Collaboration with AI tools is allowed; however, the direct solicitation is strictly prohibited
  - Employing AI tools to substantially complete assignments will be considered a violation of the Honor Code (see Generative AI Policy Guidance here for more details)

# High-Level Plan for Assignments (to be completed individually!)

- Hw1 is hopefully an easy on ramp – a Jupyter/IPython Notebook

- Hw2 is pure Python (numpy) but expects you to do (multivariate) calculus, so you really understand the basics

- Hw3 introduces PyTorch, building a feed-forward network for dependency parsing

- Hw4 and Hw5 use PyTorch on a GPU (Google Cloud)

  - Libraries like PyTorch, Tensorflow, and Jax are now the standard tools of DL

- For Final Project, more details presented later, but you either:

  - Do the default project, which is a question answering system

    - Open-ended but an easier start; a good choice for many

  - Propose a custom final project, which we approve

    - You will receive feedback from a **mentor** (TA/prof/postdoc/PhD)

  - Can work in teams of 1–3; can use any language/packages

# Lecture Plan

1. The course (10 mins)
2. Human language and word meaning (15 mins)
3. Word2vec introduction (15 mins)
4. Word2vec objective function gradients (25 mins)
5. Optimization basics (5 mins)
6. Looking at word vectors (10 mins or less)

# Trained on text data, neural machine translation is quite good!



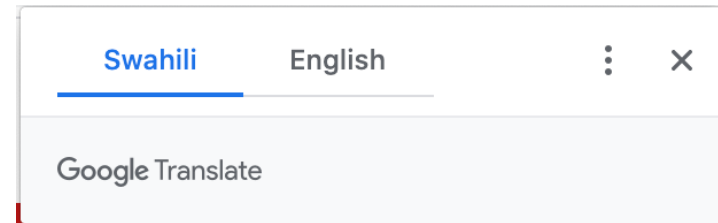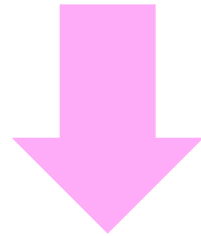https://kiswahili.tuko.co.ke/

**Malawi yawapoteza mawaziri 2 kutokana na maafa ya COVID-19**

TUKO.co.ke imefahamishwa kuwa waziri wa serikali ya mitaa Lingson Belekanyama na mwenzake wa uchukuzi Sidik Mia walifariki dunia ndani ya saa mbili tofauti.

Swahili English

Google Translate

**Malawi loses 2 ministers due to COVID-19 disaster**

TUKO.co.ke has been informed that local government minister Lingson Belekanyama and his transport counterpart Sidik Mia died within two separate hours.

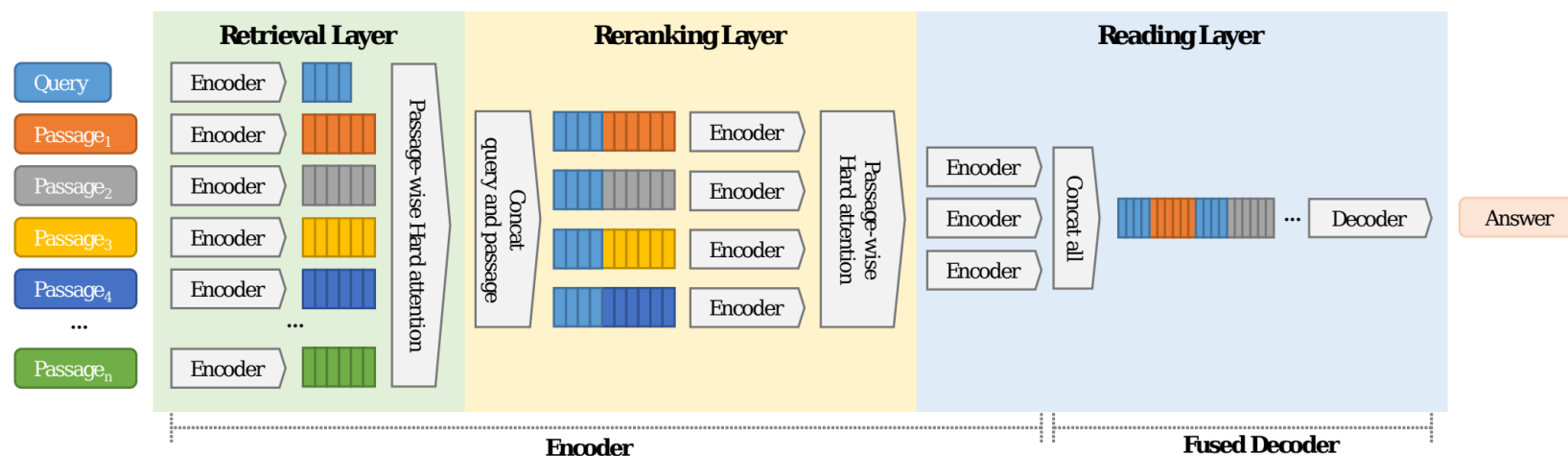# Free-text question answering: Next gen search

when did Kendrick lamar's
    first album come out?

July 2, 2011

These are my leftover songs you all can have them. I'm going to put my best out. My best effort. I'm trying to look for an album in 2012."[44] In June 2011, Lamar released "Ronald Reagan Era (His Evils)", a cut from *Section.80*, featuring Wu-Tang Clan leader RZA.[45] On July 2, 2011, Lamar released *Section.80*, his first independent album. The album features guest appearances from GLC, Colin Munroe, Schoolboy Q, and Ab-Soul, while the production was handled by Top Dawg in-house

E.g., YONO (Lee et al. 2021, https://arxiv.org/pdf/2112.07381.pdf)
    uses a T5-Large model fine-tuned for QA 3 times to run entire QA pipeline

# GPT-3: A first step on the path to foundation models    OpenAI

The SEC said, "Musk, your tweets are a
blight.
They really could cost you your job,
if you don't stop all this tweeting at night."

Then Musk cried, "Why?
The tweets I wrote are not mean,
I don't use all-caps
and I'm sure that my tweets are clean."

"But your tweets can move markets
and that's why we're sore.
You may be a genius and a billionaire,
but it doesn't give you the right to
be  a  bore!"

**S: I broke the window.**

**Q: What did I break?**

**S: I gracefully saved the day.**

**Q: What did I gracefully save?**

S: I gave John flowers.

Q: Who did I give flowers to?

S: I gave her a rose and a guitar.

Q: Who did I give a rose and a guitar to?

*How many users have signed up since the start of 2020?*

SELECT count(id) FROM users
WHERE created_at > '2020-01-01'

*What is the average number of influencers each user is subscribed to?*

SELECT avg(count) FROM ( SELECT user_id, count(*)
FROM subscribers GROUP BY user_id )
AS avg_subscriptions_per_user

A illustration from a graphic novel. A bustling city street under the shine of a full moon. The sidewalks bustling with pedestrians enjoying the nightlife. At the corner stall, a young woman with fiery red hair, dressed in a signature velvet cloak, is haggling with the grumpy old vendor. the grumpy vendor, a tall, sophisticated man is wearing a sharp suit, sports a noteworthy moustache is animatedly conversing on his steampunk telephone.

# ChatGPT, GPT-4, and more

**OpenAI**



A — Hey please draft a polite mail to explain my boss Jeremy that I would not be able to come to office for next 2 days because my 9 year song Peter is angry with me that I am not giving him much time. So I have plans to take him out for a trip. I will return and complete my work by meeting deadline.

Dear Jeremy,

I hope this email finds you well. I'm writing to let you know that I will not be able to come into the office for the next two days.

My 9-year-old son Peter is feeling upset with me because he feels like I haven't been giving him enough attention lately. In order to make it up to him, I've decided to take him on a little trip to help us reconnect.

I understand that this may cause some inconvenience, but I promise to make up for any lost time by working extra hard to meet any deadlines when I return.

Thank you for your understanding.

Best regards,
[Your Name]

Give a step-by-step analysis of the worst-case time complexity of the bubble sort algorithm, with Python code examples, but write every sentence in the speaking style of a fast-talkin' wise guy from a 1940's gangster movie.

User    What is unusual about this image?

Source: Barnorama

GPT-4    The unusual thing about this image is that a man is ironing clothes on an ironing board attached to the roof of a moving taxi.

# How do we represent the meaning of a word?

Definition: **meaning** (Webster dictionary)

- the idea that is represented by a word, phrase, etc.

- the idea that a person wants to express by using words, signs, etc.

- the idea that is expressed in a work of writing, art, etc.

**Commonest linguistic way of thinking of meaning:**

signifier (symbol) ⟺ signified (idea or thing)

= denotational semantics

tree ⟺ {🌳, 🌲, 🌴, …}

# How do we have usable meaning in a computer?

**Previously commonest NLP solution:** Use, e.g., WordNet, a thesaurus containing lists of **synonym sets** and **hypernyms** ("is a" relationships)

*e.g., synonym sets containing "good":*

```python
from nltk.corpus import wordnet as wn
poses = { 'n':'noun', 'v':'verb', 's':'adj (s)', 'a':'adj', 'r':'adv'}
for synset in wn.synsets("good"):
    print("{}: {}".format(poses[synset.pos()],
            ", ".join([l.name() for l in synset.lemmas()])))
```

```
noun: good
noun: good, goodness
noun: good, goodness
noun: commodity, trade_good, good
adj: good
adj (sat): full, good
adj: good
adj (sat): estimable, good, honorable, respectable
adj (sat): beneficial, good
adj (sat): good
adj (sat): good, just, upright
…
adverb: well, good
adverb: thoroughly, soundly, good
```

*e.g., hypernyms of "panda":*

```python
from nltk.corpus import wordnet as wn
panda = wn.synset("panda.n.01")
hyper = lambda s: s.hypernyms()
list(panda.closure(hyper))
```

```
[Synset('procyonid.n.01'),
Synset('carnivore.n.01'),
Synset('placental.n.01'),
Synset('mammal.n.01'),
Synset('vertebrate.n.01'),
Synset('chordate.n.01'),
Synset('animal.n.01'),
Synset('organism.n.01'),
Synset('living_thing.n.01'),
Synset('whole.n.02'),
Synset('object.n.01'),
Synset('physical_entity.n.01'),
Synset('entity.n.01')]
```
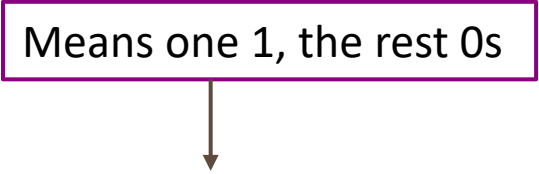
17

# Problems with resources like WordNet

- A useful resource but missing nuance:
  - e.g., "proficient" is listed as a synonym for "good"
    This is only correct in some contexts
  - Also, WordNet list offensive synonyms in some synonym sets without any coverage of the connotations or appropriateness of words
- Missing new meanings of words:
  - e.g., wicked, badass, nifty, wizard, genius, ninja, bombest
  - Impossible to keep up-to-date!
- Subjective
- Requires human labor to create and adapt
- Can't be used to accurately compute word similarity (see following slides)

# Representing words as discrete symbols

In traditional NLP, we regard words as discrete symbols:

hotel, conference, motel – a localist representation

Means one 1, the rest 0s

Such symbols for words can be represented by one-hot vectors:

motel = [0 0 0 0 0 0 0 0 0 0 1 0 0 0]

hotel = [0 0 0 0 0 0 0 1 0 0 0 0 0 0]

Vector dimension = number of words in vocabulary (e.g., 500,000+)

# Problem with words as discrete symbols

**Example:** in web search, if a user searches for "Seattle motel", we would like to match documents containing "Seattle hotel"

But:

$$\text{motel} = [0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 0]$$
$$\text{hotel} = [0\ 0\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 0\ 0\ 0\ 0]$$

These two vectors are orthogonal

There is no natural notion of **similarity** for one-hot vectors!

**Solution:**

- Could try to rely on WordNet's list of synonyms to get similarity?
  - But it is well-known to fail badly: incompleteness, etc.
- **Instead: learn to encode similarity in the vectors themselves**

# Representing words by their context

- Distributional semantics: **A word's meaning is given by the words that frequently appear close-by**

  - *"You shall know a word by the company it keeps"* (J. R. Firth 1957: 11)

  - One of the most successful ideas of modern statistical NLP!

- When a word *w* appears in a text, its **context** is the set of words that appear nearby (within a fixed-size window).

- We use the many contexts of *w* to build up a representation of *w*

…government debt problems turning into   banking   crises as happened in 2009…

…saying that Europe needs unified   banking   regulation to replace the hodgepodge…

…India has just given its   banking   system a shot in the arm…

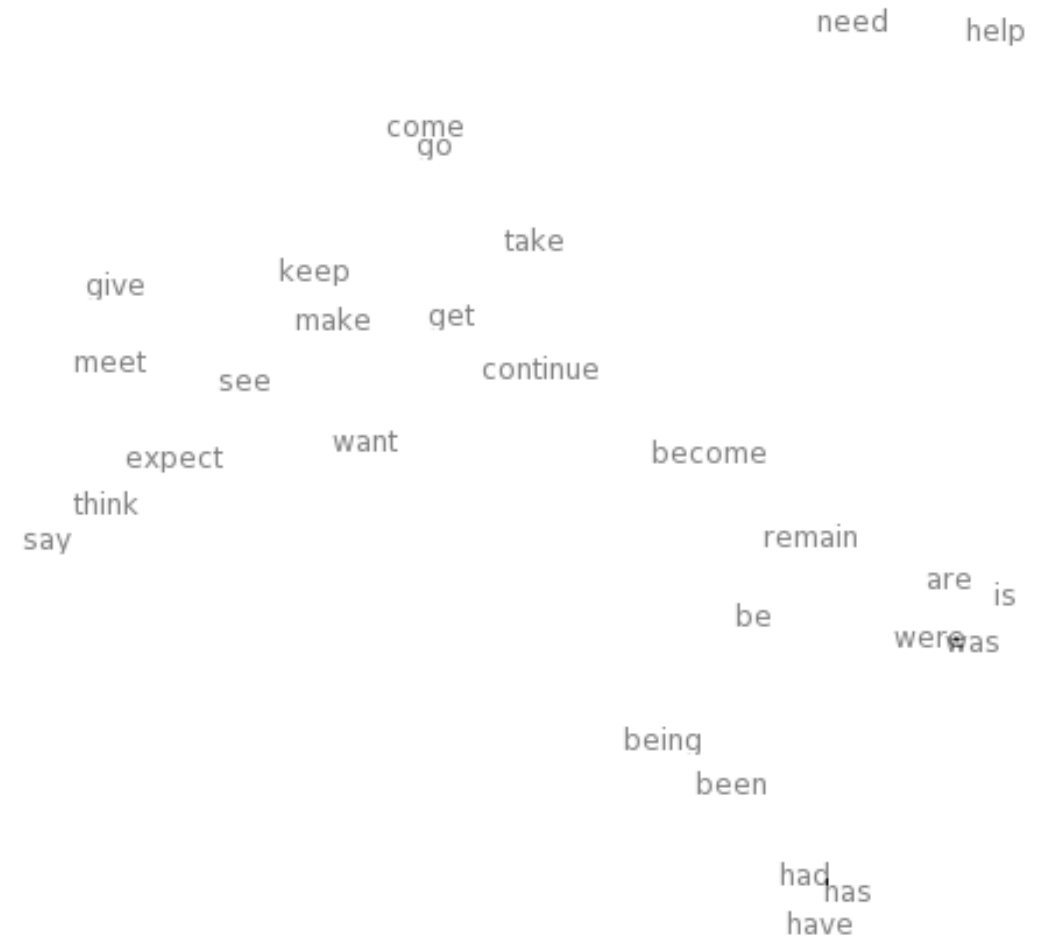These context words will represent *banking*

# Word vectors

We will build a dense vector for each word, chosen so that it is similar to vectors of words that appear in similar contexts, measuring similarity as the vector dot (scalar) product

$$
banking = \begin{pmatrix} 0.286 \\ 0.792 \\ -0.177 \\ -0.107 \\ 0.109 \\ -0.542 \\ 0.349 \\ 0.271 \end{pmatrix}
\qquad
monetary = \begin{pmatrix} 0.413 \\ 0.582 \\ -0.007 \\ 0.247 \\ 0.216 \\ -0.718 \\ 0.147 \\ 0.051 \end{pmatrix}
$$

Note: word vectors are also called (word) embeddings or (neural) word representations
They are a distributed representation

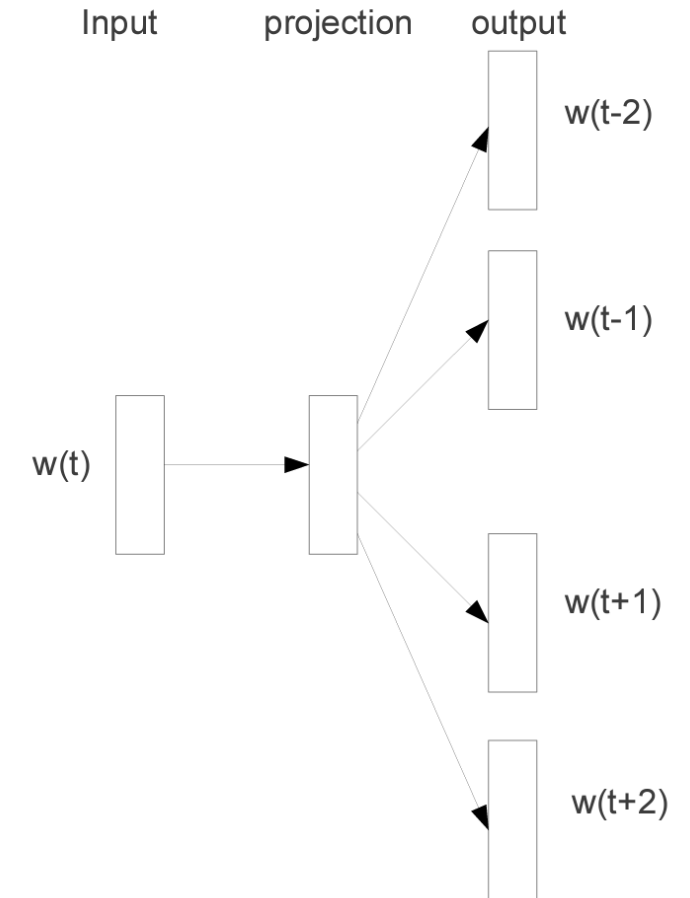# Word meaning as a neural word vector – visualization

$$expect = \begin{pmatrix} 0.286 \\ 0.792 \\ -0.177 \\ -0.107 \\ 0.109 \\ -0.542 \\ 0.349 \\ 0.271 \\ 0.487 \end{pmatrix}$$

need    help

come
go

take

give    keep

make    get

meet        continue
see

want            become

expect

think

say                         remain

are    is

be

were    was

being

been

had    has
have

# 3. Word2vec: Overview

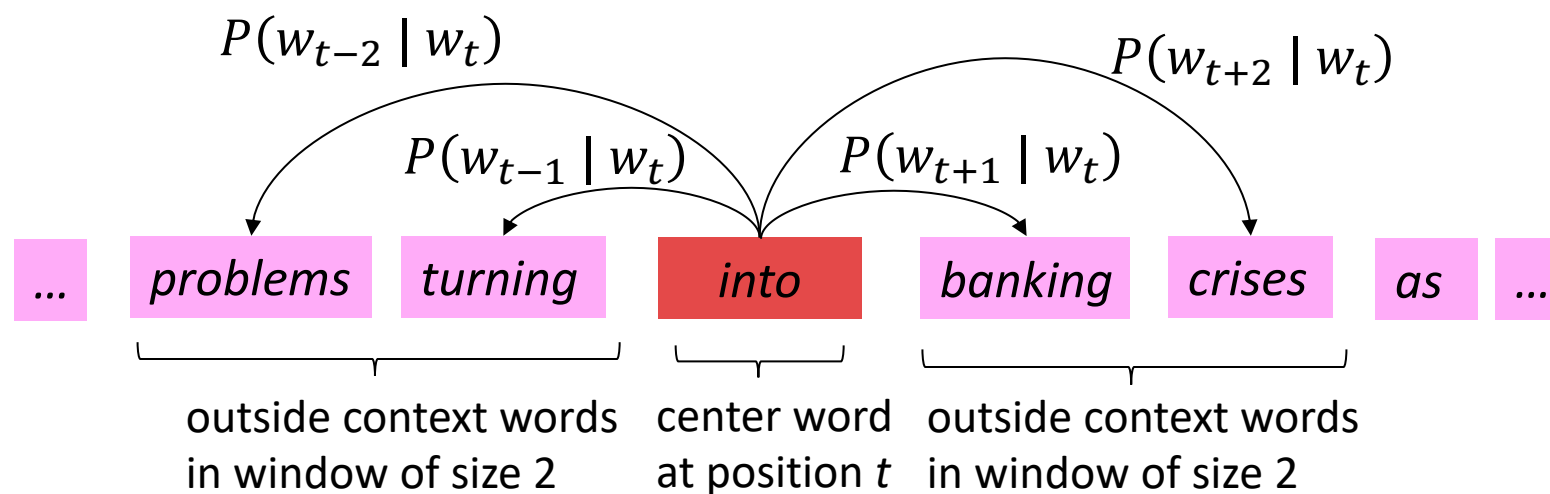Word2vec is a framework for learning word vectors
(Mikolov et al. 2013)

Idea:

- We have a large corpus ("body") of text: a long list of words
- Every word in a fixed vocabulary is represented by a vector
- Go through each position $t$ in the text, which has a center word $c$ and context ("outside") words $o$
- Use the similarity of the word vectors for $c$ and $o$ to calculate the probability of $o$ given $c$ (or vice versa)
- Keep adjusting the word vectors to maximize this probability

Input     projection     output

w(t)

w(t-2)

w(t-1)

w(t+1)

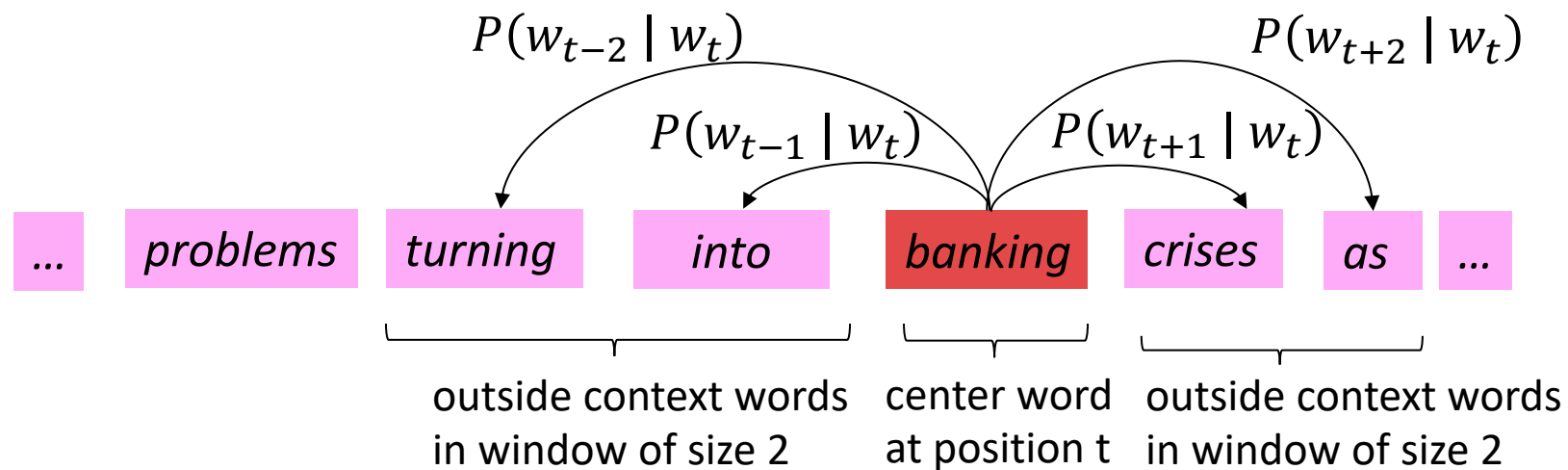w(t+2)

Skip-gram model
(Mikolov et al. 2013)

# Word2Vec Overview

Example windows and process for computing $P(w_{t+j} \mid w_t)$

# Word2Vec Overview

Example windows and process for computing $P(w_{t+j} \mid w_t)$

# Word2vec: objective function

For each position $t = 1, \ldots, T$, predict context words within a window of fixed size $m$, given center word $w_t$. Data likelihood:

Likelihood = $L(\theta) = \displaystyle\prod_{t=1}^{T} \prod_{\substack{-m \leq j \leq m \\ j \neq 0}} P(w_{t+j} \mid w_t; \theta)$

$\theta$ is all variables to be optimized

sometimes called a *cost* or *loss* function

The objective function $J(\theta)$ is the (average) negative log likelihood:

$$J(\theta) = -\frac{1}{T}\log L(\theta) = -\frac{1}{T}\sum_{t=1}^{T} \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P(w_{t+j} \mid w_t; \theta)$$

Minimizing objective function $\Longleftrightarrow$ Maximizing predictive accuracy

# Word2vec: objective function

- We want to minimize the objective function:

$$J(\theta) = -\frac{1}{T}\sum_{t=1}^{T}\sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P\left(w_{t+j} \mid w_t; \theta\right)$$

- **Question:** How to calculate $P\left(w_{t+j} \mid w_t; \theta\right)$ ?

- **Answer:** We will *use two* vectors per word *w*:
  - $v_w$ when *w* is a center word
  - $u_w$ when *w* is a context word

- Then for a center word *c* and a context word *o*:

$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

# Word2vec: prediction function

② Exponentiation makes anything positive

① Dot product compares similarity of *o* and *c*.
$u^T v = u.v = \sum_{i=1}^{n} u_i v_i$
Larger dot product = larger probability

$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

③ Normalize over entire vocabulary to give probability distribution

- This is an example of the **softmax function** $\mathbb{R}^n \to (0,1)^n$   Open region

$$\text{softmax}(x_i) = \frac{\exp(x_i)}{\sum_{j=1}^{n} \exp(x_j)} = p_i$$

- The softmax function maps arbitrary values $x_i$ to a probability distribution $p_i$
  - "max" because amplifies probability of largest $x_i$
  - "soft" because still assigns some probability to smaller $x_i$
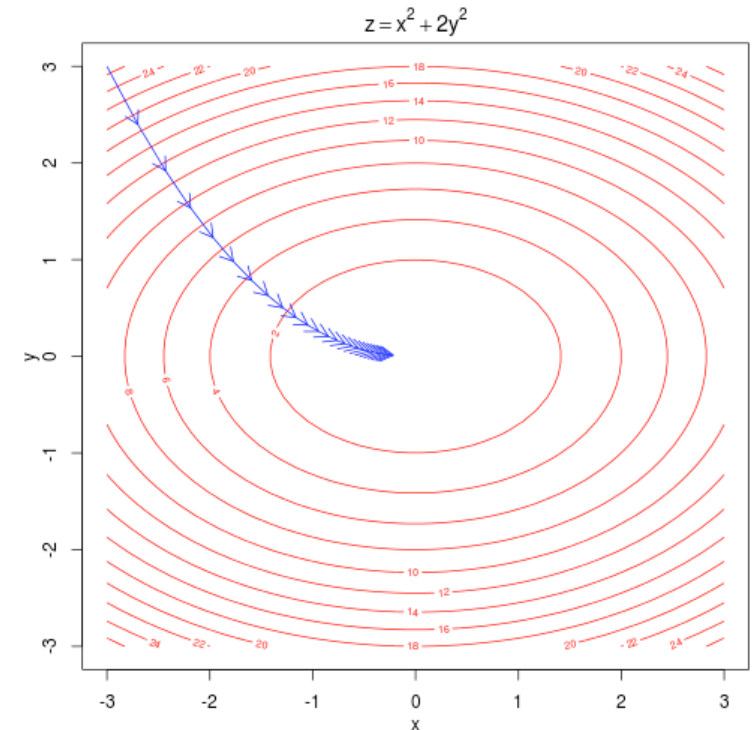  - Frequently used in Deep Learning

But sort of a weird name because it returns a distribution!

# To train the model: Optimize value of parameters to minimize loss

To train a model, we gradually adjust parameters to minimize a loss

- Recall: $\theta$ represents **all** the model parameters, in one long vector

- In our case, with $d$-dimensional vectors and $V$-many words, we have →

- Remember: every word has two vectors

$$\theta = \begin{bmatrix} v_{aardvark} \\ v_a \\ \vdots \\ v_{zebra} \\ u_{aardvark} \\ u_a \\ \vdots \\ u_{zebra} \end{bmatrix} \in \mathbb{R}^{2dV}$$



$z = x^2 + 2y^2$

- We optimize these parameters by walking down the gradient (see right figure)

- We compute **all** vector gradients!

31

# Interactive Session!

- $L(\theta) = \prod_{t=1}^{T} \prod_{\substack{-m \leq j \leq m \\ j \neq 0}} P\left(w_{t+j} \mid w_t; \theta\right)$

- For a center word $c$ and a context word $o$: $P(o|c) = \dfrac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$

**4.**

Objective Function

Maximize $J'(\theta) = \prod\limits_{t=1}^{T} \prod\limits_{\substack{-m \leq j \leq m \\ j \neq 0}} p(w'_{t+j} | w_t; \theta)$

Or minimize **ave.**
neg. log
likelihood $J(\theta) = -\frac{1}{T} \sum\limits_{t=1}^{T} \sum\limits_{\substack{-m \leq j \leq m \\ j \neq 0}} \log p(w'_{t+j} | w_t)$

[negate to minimize;
log is monotone]

Text
length

window
size

where

$p(o|c) = \dfrac{\exp(u_o^T v_c)}{\sum\limits_{w=1}^{V} \exp(u_w^T v_c)}$

word IDs

We now take derivatives to work out minimum

Each word type
(vocab entry)
has two word
representations:
as center word
and context word

$$\frac{\partial}{\partial v_c} \log \frac{\exp(u_o^T v_c)}{\sum\limits_{w=1}^{V} \exp(u_w v_c)}$$

$$= \underbrace{\frac{\partial}{\partial v_c} \log \exp(u_o^T v_c)}_{①} - \underbrace{\frac{\partial}{\partial v_c} \log \sum_{w=1}^{V} \exp(u_w^T v_c)}_{②}$$

$$① \quad \frac{\partial}{\partial v_c} \log \exp(u_o^T v_c) = \frac{\partial}{\partial v_c} u_o^T v_c = u_o$$

$\underbrace{\phantom{\frac{\partial}{\partial v_c}}}_{inverses}$

$\uparrow$

Vector!
Not high
school
single
variable
calculus

You can do things one variable at a time, and this may be helpful when things get gnarly.

$$\forall j \quad \frac{\partial}{\partial (v_c)_j} u_o^T v_c = \frac{\partial}{\partial (v_c)_j} \sum_{i=1}^{d} (u_o)_i (v_c)_i$$

$$= (u_o)_j$$

Each term is zero except when $i = j$

② $\dfrac{\partial}{\partial v_c} \underbrace{\log \underbrace{\sum\limits_{w=1}^{V} \exp\left(u_w^T v_c\right)}_{z = g(v_c)}}_{f}$

$= \underbrace{\dfrac{1}{\sum\limits_{w=1}^{V} \exp\left(u_w^T v_c\right)}}_{z} \cdot \dfrac{\partial}{\partial v_c} \sum\limits_{x=1}^{V} \exp\left(u_x^T v_c\right)$

— Important to change index

$\underbrace{\dfrac{\partial}{\partial v_c} f(\overbrace{g(v_c)}^{z}) = \dfrac{\partial f}{\partial z}}_{} \cdot \underbrace{\dfrac{\partial z}{\partial v_c}}_{}$   Use chain rule

$= \dfrac{1}{\sum\limits_{w=1}^{V} \exp\left(u_w^T v_c\right)} \cdot \left( \sum\limits_{x=1}^{V} \dfrac{\partial}{\partial v_c} \underbrace{\exp\left(u_x^T v_c\right)}_{\substack{f \quad z = g(v_c)}} \right)$

Move deriv inside sum

$\left( \sum\limits_{x=1}^{V} \exp\left(u_x^T v_c\right) \dfrac{\partial}{\partial v_c} u_x^T v_c \right)$   Chain rule

$\left( \sum\limits_{x=1}^{V} \exp\left(u_x^T v_c\right) u_x \right)$

38

$$\frac{\partial}{\partial v_c} \log\left(p(o|c)\right) = u_o - \frac{1}{\sum\limits_{w=1}^{V} \exp\left(u_w^T v_c\right)} \cdot \left(\sum\limits_{x=1}^{V} \exp\left(u_x^T v_c\right) u_x\right)$$

$$= u_o - \sum\limits_{x=1}^{V} \frac{\exp\left(u_x^T v_c\right)}{\sum\limits_{w=1}^{V} \exp\left(u_w^T v_c\right)} u_x$$

Distribute term across sum

$$= u_o - \sum\limits_{x=1}^{V} p(x|c) u_x$$

This an expectation: average over all context vectors weighted by their probability

$$= \text{observed} - \text{expected}$$

This is just the derivatives for the center vector parameters
Also need derivatives for output vector parameters
(they're similar)
Then we have derivative w.r.t. all parameters and can minimize