# Lesson 6:
# Sentence Classification and a little about RNNs

Partially based on slides by Jurafsky and Martin

Speech and Language Processing, 3rd Edition

# Sentiment Analysis

*Example #1: Movie Reviews*

- Unbelievably disappointing
- Full of zany characters and richly applied satire, and some great plot twists
- This is the greatest screwball comedy ever filmed
- It was pathetic. The worst part about it was the boxing scenes.

# Sentiment Analysis

*Example #2: Product Reviews*



**HP Officejet 6500A Plus e-All-in-One Color Ink-jet - Fax / copier / printer / scanner**
$89 online, $100 nearby    ★★★★☆ 377 reviews
September 2010 - Printer - HP - Inkjet - Office - Copier - Color - Scanner - Fax - 250 sh

## Reviews

**Summary** - Based on 377 reviews

| 1 star | 2 | 3 | 4 stars | 5 stars |
|--------|---|---|---------|---------|

What people are saying

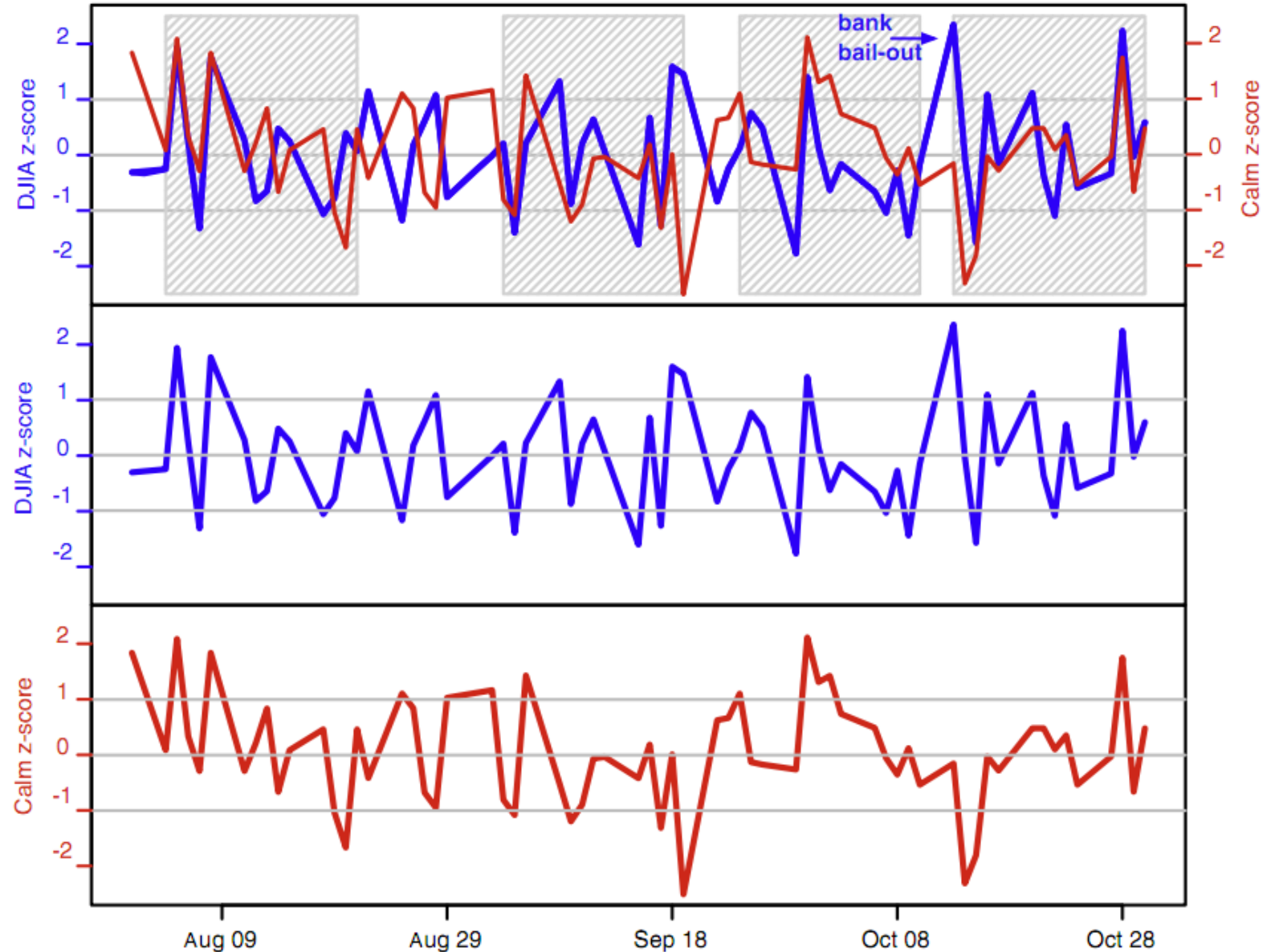| | | |
|---|---|---|
| ease of use | | "This was very easy to setup to four computers." |
| value | | "Appreciate good quality at a fair price." |
| setup | | "Overall pretty easy setup." |
| customer service | | "I DO like honest tech support people." |
| size | | "Pretty Paper weight." |
| mode | | "Photos were fair on the high quality mode." |
| colors | | "Full color prints came out with great quality." |

- A Sentiment Analysis system called CALM applied to Twitter predicts the Dow Jones Industrial Average (DJIA) 3 days later

- Such algorithms are already used by hedge funds

Johan Bollen, Huina Mao, Xiaojun Zeng. 2011. Twitter mood predicts the stock market, Journal of Computational Science 2:1, 1-8. 10.1016/j.jocs.2010.12.007.



4

# Scherer Typology of Affective States

- **Emotion**: brief organically synchronized … evaluation of a major event
  - *angry, sad, joyful, fearful, ashamed, proud, elated*
- **Mood**: diffuse non-caused low-intensity long-duration change in subjective feeling
  - *cheerful, gloomy, irritable, listless, depressed, buoyant*
- **Interpersonal stances**: affective stance toward another person in a specific interaction
  - *friendly, flirtatious, distant, cold, warm, supportive, contemptuous*
- **Attitudes**: enduring, affectively colored beliefs, dispositions towards objects or persons
  - *liking, loving, hating, valuing, desiring*
- **Personality traits**: stable personality dispositions and typical behavior tendencies
  - *nervous, anxious, reckless, morose, hostile, jealous*

# Sentiment Analysis: Definition

- Sentiment analysis is the detection of **attitudes**

  "enduring, affectively colored beliefs, dispositions towards objects or persons"

  1. **Holder (source)** of attitude
  2. **Target (aspect)** of attitude
  3. **Type** of attitude
     - From a set of types: *like, love, hate, value, desire,* etc.
     - Or (more commonly) simple weighted **polarity**: *positive, negative, neutral,* together with *strength*
  4. **Text** containing the attitude
     - Sentence or entire document

# Sentiment Analysis

- Simplest task:
  - Is the attitude of this text positive or negative?
- More complex:
  - Rank the attitude of this text from 1 to 5
- Advanced:
  - Detect the target, source, or complex attitude types

# Sentiment Classification in Movie Reviews

- Is an IMDB movie review positive or negative?

✓

when _star wars_ came out some twenty years ago
, the image of traveling throughout the stars has
become a commonplace image . […]

when han solo goes light speed , the stars change
to bright lines , going towards the viewer in lines
that converge at an invisible point .

cool .

✗

" snake eyes " is the most aggravating
kind of movie : the kind that shows so
much potential then becomes
unbelievably disappointing .

it's not just because this is a brian
depalma film , and since he's a great
director and one who's films are always
greeted with at least some fanfare .

Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan.  2002.  Thumbs up? Sentiment
Classification using Machine Learning Techniques. EMNLP-2002, 79—86.
Bo Pang and Lillian Lee.  2004.  A Sentimental Education: Sentiment Analysis Using
Subjectivity Summarization Based on Minimum Cuts.  ACL, 271-278

# A Simple Classifier

- Log-linear or Naïve Bayes classifier
- Features:
  - Tokenized words
  - Possibly mark-up (e.g., hashtags in Twitter, headers in HTMLs)
- Features are often binary
  - Indicating whether a word appeared or did not appear in the document (bag of words)
  - Often works better for text classification than word frequency

# Error Analysis:
# What makes reviews hard to classify?

- **Sarcasm:**
  - Perfume review in *Perfumes: the Guide*:
    - "If you are reading this because it is your darling fragrance, please wear it at home exclusively, and tape the windows shut."
  - On *Automobile Steering Wheel Attachable Work Surface*:
    - "You wouldn't believe how much more interesting my commute is now that I have something to do other than just stare out the window! I'm using it right now to post this review and I never"

- **Thwarted Expectations and Ordering Effects:**
  - "This film should be brilliant.  It sounds like a great plot, the actors are first grade, and the supporting cast is good as well, and Stallone is attempting to deliver a good performance. However, it **can't hold up**."

# Negation in Sentiment Analysis

- One practice is to add NOT_ to every word between negation and following punctuation:

<div align="center">

`didn't like this movie , but I`



`didn't NOT_like NOT_this NOT_movie but I`

</div>

# Negation in Sentiment Analysis

- This is a very crude solution:
  - Explicit negation is only one way to reverse polarity
    - "He **failed** to convey the importance of his message"
  - Logical structure can be more complex
    - "I do**n't** think anyone could have done a better job"
  - Negation scope:
    - "I did**n't** like the exposition, but otherwise liked the movie"

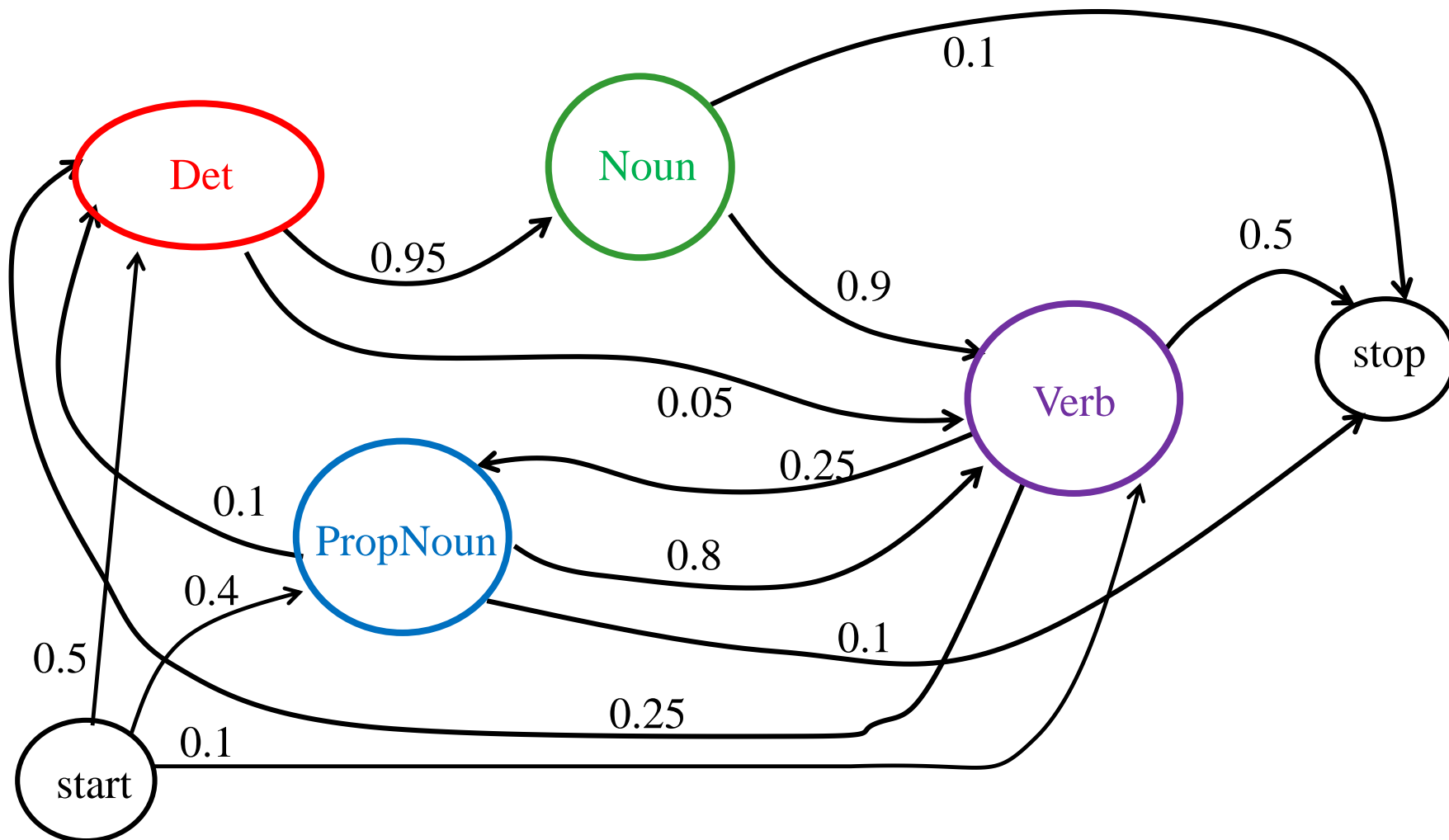- More recent approaches take the context (surrounding words) into account

# Sentiment Analysis as Sequence Labeling

- This construal of sentiment analysis attempts to capture the meaning of a word <u>in context</u> by encoding (parts of) the sentence as features

- Recently: using Recurrent Neural Networks (RNNs)

- Recall the underlying assumption in Markov (n-gram) models:
  - You only need to know the last *n* tokens you've encountered to know what's next
  - Alternative view: the probability of a sequence is the product of the probabilities of its *n-grams*

# Sentiment Analysis as Sequence Labeling

- Consider the example:
  - *How can you not see this movie?*
  - *You should not see this movie*

- How well will a bi-gram model work?
  - Similar unigrams and bigrams → similar prediction


- The problem with Markov Models: need to maintain a **state** to capture distant influences
  - The size of the space increases exponentially with the distance

# Recall: Markov Models are FSA with transition probabilities

# Recurrent Neural Networks

- Motivation:
  - Neural network model, but with a state
  - How can we borrow ideas from FSAs?

- RNNs are FSAs …
  - With a twist
  - No longer finite in the same sense
  - The state is an embedding of the history in a continuous space

# Recurrent Neural Networks

- Map from dense sequence to dense representation
  - Maps a sequence of vector inputs to a sequence of vector states

$$x_1, ..., x_n \rightarrow s_1, ..., s_n$$

  - A (parametrized) transition function $R$ does the mapping:

$$s_i = R(s_{i-1}, x_i)$$

  - $R$ is parametrized and parameters are shared across steps

$$s_4 = R(s_3, x_4) = R(R(s_2, x_3), x_4) = R(R(R(s_1, x_2), x_3), x_4)$$
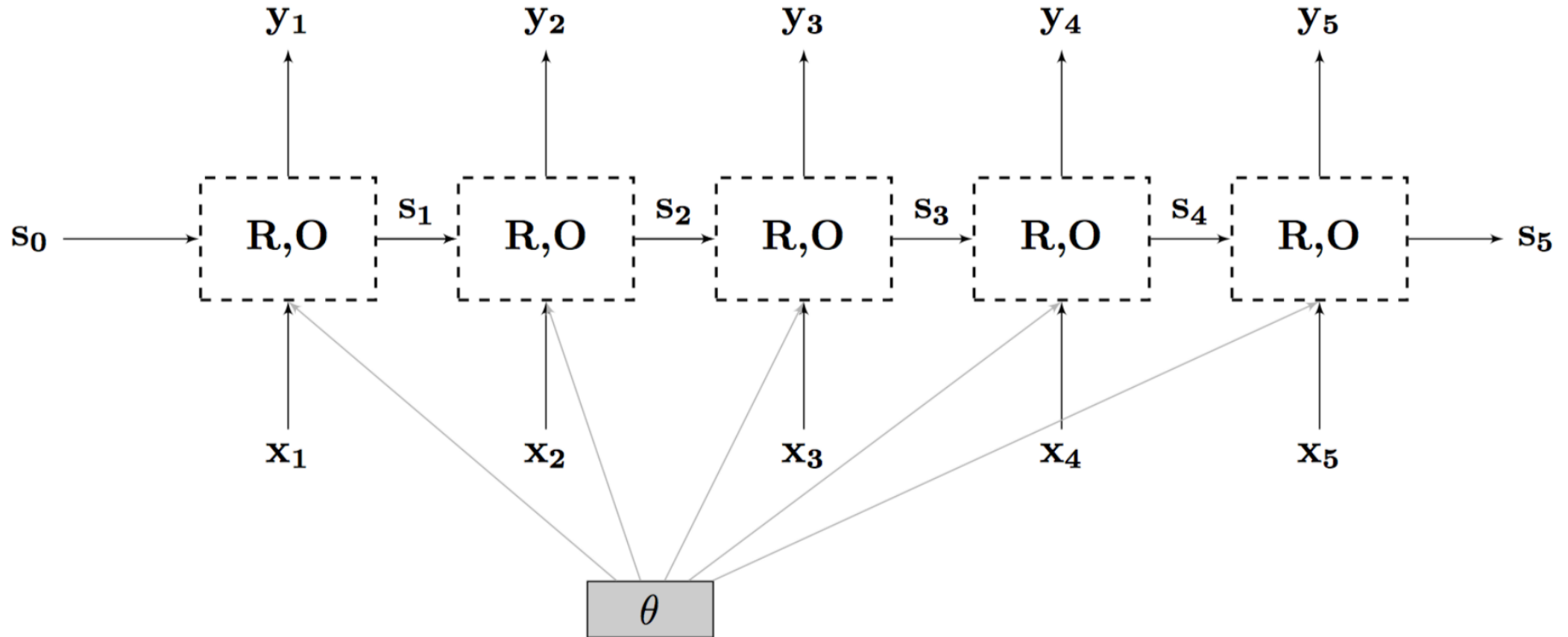
# Recurrent Neural Networks

- An output function *O* maps states to (vector) outputs, which are often viewed as a distribution over the possible labels

- Example:

$$R_{Elman}(s, x) = tanh(W[x, s] + b)$$

$$O(s_i) = softmax(W s_i + b)$$
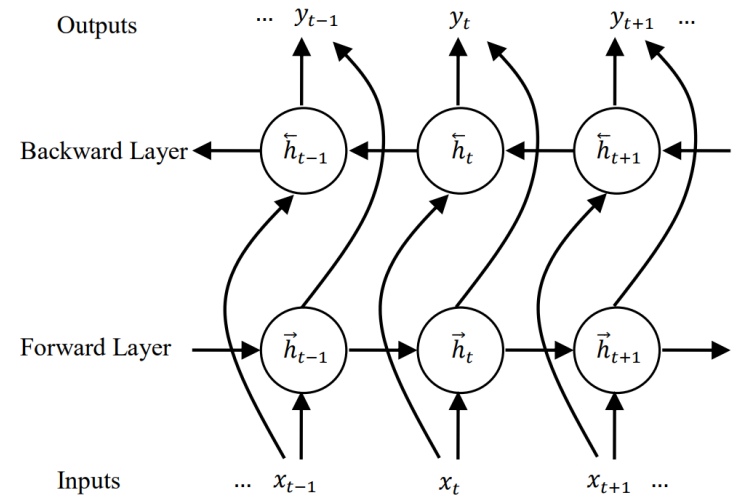
# RNNs: Graphical Representation

# Back to Sentiment Analysis

- When using RNNs for sentence classification (such as sentiment analysis), it is often practical to use Bi-RNNs

- Bi-RNNs:
  - 2 RNNs, one going back to forth and the other forth to back
  - Output function is a function of both states

- This allows the states associated with each word to encode relevant information from the words following them and preceding them

$$\overrightarrow{h}_t = \mathcal{H}(W_{x\overrightarrow{h}}x_t + W_{\overrightarrow{h}\,\overrightarrow{h}}\overrightarrow{h}_{t-1} + b_{\overrightarrow{h}}) \quad (9)$$

$$\overleftarrow{h}_t = \mathcal{H}(W_{x\overleftarrow{h}}x_t + W_{\overleftarrow{h}\,\overleftarrow{h}}\overleftarrow{h}_{t+1} + b_{\overleftarrow{h}}) \quad (10)$$

$$y_t = W_{\overrightarrow{h}y}\overrightarrow{h}_t + W_{\overleftarrow{h}y}\overleftarrow{h}_t + b_y \quad (11)$$

# Back to Sentiment Analysis

- One simple way to do sentiment analysis (or other sentence classification) with Bi-RNNs is to average the output sequence:
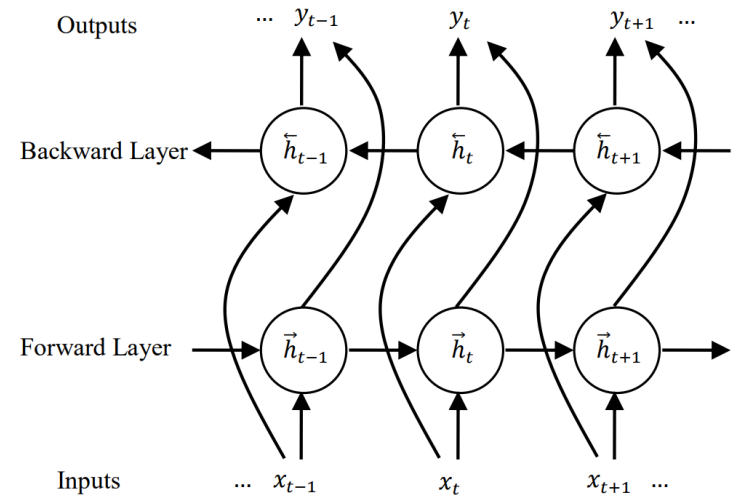
$$\overrightarrow{h}_t = \mathcal{H}(W_{x\overrightarrow{h}}x_t + W_{\overrightarrow{h}\overrightarrow{h}}\overrightarrow{h}_{t-1} + b_{\overrightarrow{h}}) \quad (9)$$

$$\overleftarrow{h}_t = \mathcal{H}(W_{x\overleftarrow{h}}x_t + W_{\overleftarrow{h}\overleftarrow{h}}\overleftarrow{h}_{t+1} + b_{\overleftarrow{h}})(10)$$

$$y_t = W_{\overrightarrow{h}y}\overrightarrow{h}_t + W_{\overleftarrow{h}y}\overleftarrow{h}_t + b_y \quad (11)$$

$$y = \frac{1}{N}\Sigma_i y_i$$

- Now train a binary logistic classifier for predicting the sentiment:

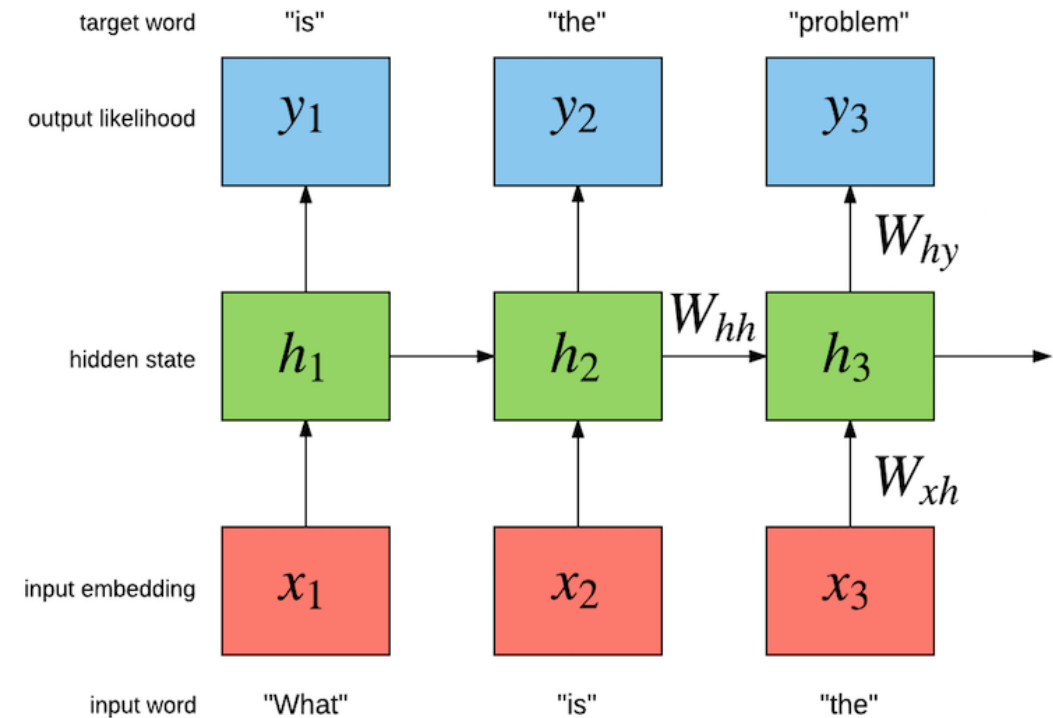$$P(+|y) = \frac{1}{1 + e^{w^t y + b}}$$

# Context in RNN models for Sentence Classification

- RNN models are able to take the context (both preceding and following) into account, as well as the linear order between the words
  - bag-of-words models cannot
- They indeed show better performance in tasks such as sentiment analysis
- However, the state sequence is not easy to interpret
- Further investigation is needed to establish what contextual and semantic aspects of sentences are captured using these techniques

# Neural Language Models

- Language models based on RNNs have shown much power in recent years, consistently surpassing n-gram models

- The basic architecture is that of a sequence recurrent neural net (RNN)

  - Input words are converted to 1-hot vectors

  - The output is passed through a softmax layer, which defines the probability of the next word

  - The loss function is often just the log probability of the next word predicted by the model



http://www.fit.vutbr.cz/~imikolov/rnnlm/thesis.pdf - Mikolov (2011) PhD Thesis

# Character-level Language Models

- Vocabulary: characters instead of words

- Advantage:
  - Small vocabulary $\rightarrow$ compact model
  - Can generalize over morphologically similar words

- However:
  - Need to learn how to spell
  - Longer range dependencies between tokens

# Character-level Language Models



[during training – green = value to increase]