

CS 489/698: Introduction to Natural Language Processing

Lecture 4: Classification

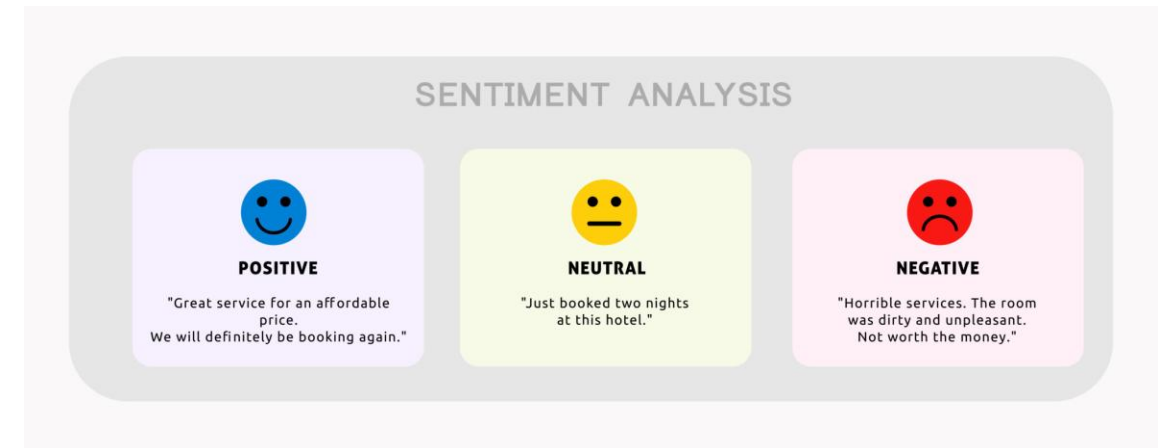
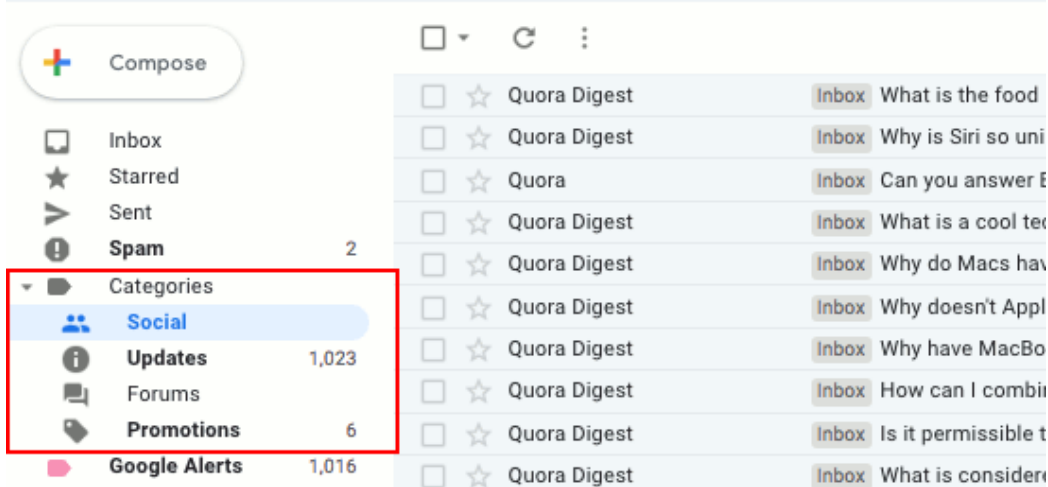
Instructor: Freda Shi

rhs@uwaterloo.ca

January 14th, 2026

Classification

Simplest user facing NLP application



Outline of Today's Lecture

Modeling approaches for classification

- Rule-based classification
- Naïve Bayes
- Logistic regression
- Neural networks

Rule-Based Classifier

Sentiment classification of sentence s (classes: positive, negative)

- If s contains words in [*good, excellent, extraordinary, ...*] return positive
- If s contains words in [*bad, terrible, awful, ...*] return negative



Nice interpretability



Can be very accurate (with carefully refined rules)



Rules are difficult to define



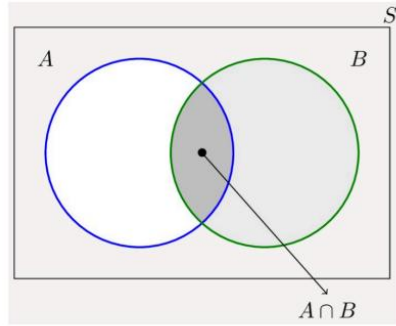
System can be very complicated



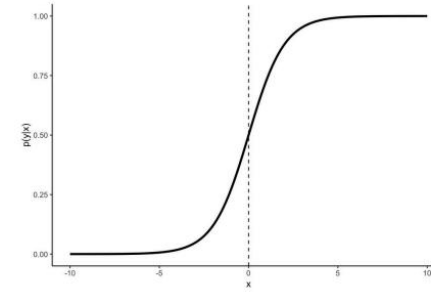
Hardly generalizable

VADER (Valence Aware Dictionary and sEntiment Reasoner) is a lexicon and rule-based sentiment analysis tool that is *specifically attuned to sentiments expressed in social media*. It is fully open-sourced under the [\[MIT License\]](#) (we sincerely appreciate all attributions and readily accept most contributions, but please don't hold us liable).

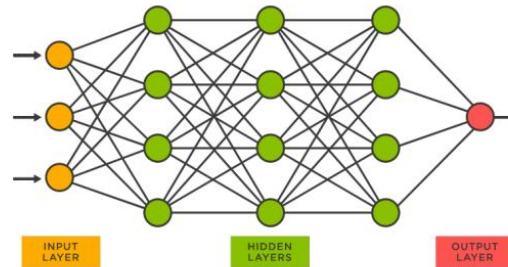
Statistical Classifiers



Naïve Bayes



Logistic Regression



Neural Networks

Statistical Classifiers: General Formulation

Data: A set of labeled sentences $\mathcal{D} = \{(s_1, y_1), (s_2, y_2), \dots, (s_n, y_n)\}$.

- s_i : a sentence (or a piece of text).
- y_i : label, usually represented as an integer.

Deliverable: A classifier that takes an arbitrary sentence s as input, and predicts the label y .

$$\text{classify}(s) = \underset{y}{\operatorname{argmax}} \text{ score}(s, y; \Theta)$$

Θ : parameters of
score function

Modeling, Learning, and Inference

Inference: solve arg max

Modeling: define score function

$$\text{classify}(s) = \underset{y}{\operatorname{argmax}} \text{score}(s, y; \Theta)$$

Learning: choose parameter

Naïve Bayes: Modeling

Probabilistic model: $\text{score}(s, y; \Theta) := P(y \mid s)$ estimated from the data.

How to estimate $P(y \mid s)$ to ensure sufficient generalizability?

The Bayes rule:

$$P(y \mid s) = \frac{P(s \mid y)P(y)}{P(s)} \propto P(s \mid y)P(y)$$

Unigram assumption:

$$P(s \mid y) = P(w_1 \mid y)P(w_2 \mid y) \dots P(w_k \mid y) \quad \text{if } s = w_1 w_2 \dots w_k$$

Θ in $\text{score}(s, y; \Theta)$: look-up table that stores $P(w \mid y)$ and $P(y)$.

Naïve Bayes: Learning

$$P(s \mid y) = P(w_1 \mid y)P(w_2 \mid y) \dots P(w_k \mid y) \quad \text{if } s = w_1 w_2 \dots w_k$$

Each $P(w \mid y)$ and $P(y)$ can be estimated by counting from the corpora.

$$P(w \mid y) = \frac{\text{count}(w, y)}{\sum_{w' \in V} \text{count}(w', y)} \quad P(y) = \frac{\text{count}(y)}{n}$$

This is the estimation that maximizes dataset probability.

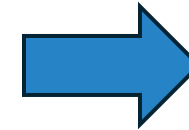
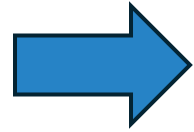
Issue: an unseen word in certain class will lead to a zero probability!

Solution: smoothing.

$$P(w \mid y) = \frac{\text{count}(w, y) + \alpha}{\sum_{w' \in V} \text{count}(w', y) + \alpha|V|}$$

“Bag-of-words” Features

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!



it	6
I	5
the	4
to	3
and	3
seen	2
yet	1
would	1
sweet	1
genre	1
fairy	1
humor	1
have	1
great	1
...	

Naïve Bayes: Inference

$$\text{classify}(s) = \underset{y}{\operatorname{argmax}} \text{ score}(s, y; \Theta) = \underset{y}{\operatorname{argmax}} P(y \mid s)$$

$$P(y \mid s) \propto P(y) P(s \mid y) = P(y) P(w_1 \mid y) P(w_2 \mid y) \dots P(w_k \mid y)$$

Consult the look-up tables for values of $P(y)$ and $P(w \mid y)$.

Naïve Bayes: Example

Type	Category	Sentence
Training	-	just plain boring
Training	-	entirely predictable and lacks energy
Training	-	no surprises and very few laughs
Training	+	very powerful
Training	+	the most fun film of the summer
Testing	?	predictable with no fun

- Scoring the examples

$$P(+)P(\text{predictable} \mid +)P(\text{no} \mid +)P(\text{fun} \mid +) = 3.2 \times 10^{-5}$$

$$P(-)P(\text{predictable} \mid -)P(\text{no} \mid -)P(\text{fun} \mid -) = 6.1 \times 10^{-5}$$

- Prior from training set:

$$P(+) = \frac{2}{5}, P(-) = \frac{3}{5}$$

- $|V| = 20$ from training set

- Conditional probabilities ($\alpha = 1$):

$$P(\text{predictable} \mid +) = \frac{0 + 1}{9 + 20} = \frac{1}{29}$$

$$P(\text{no} \mid +) = \frac{0 + 1}{9 + 20} = \frac{1}{29}$$

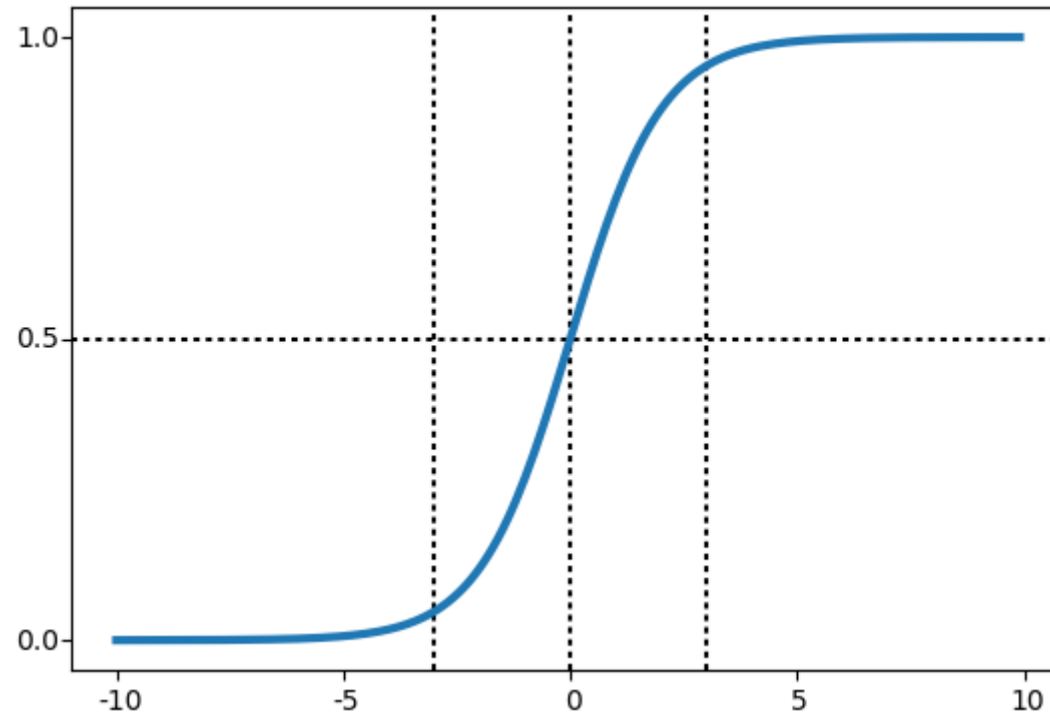
$$P(\text{fun} \mid +) = \frac{1 + 1}{9 + 20} = \frac{2}{29}$$

$$P(\text{predictable} \mid -) = \frac{1 + 1}{14 + 20} = \frac{2}{34}$$

$$P(\text{no} \mid -) = \frac{1 + 1}{14 + 20} = \frac{2}{34}$$

$$P(\text{fun} \mid -) = \frac{0 + 1}{14 + 20} = \frac{1}{34}$$

Recap: Logistic Function



Logistic function: $\sigma(x) = \frac{1}{1+e^{-x}}$, $\sigma: \mathbb{R} \rightarrow (0, 1)$.

“Native” for binary classification.

Logistic Regression: Modeling

Suppose we can represent sentence s with a vector \mathbf{x} .

Probabilistic model:

$$\text{score}(s, y = 1; \boldsymbol{\theta}) := P(y = 1 \mid s) = \sigma(\boldsymbol{\theta} \cdot \mathbf{x} + b)$$

$$\text{score}(s, y = 0; \boldsymbol{\theta}) := P(y = 0 \mid s) = 1 - \sigma(\boldsymbol{\theta} \cdot \mathbf{x} + b)$$

$\boldsymbol{\theta} \cdot \mathbf{x} + b$ can be written as $\boldsymbol{\theta} \cdot \mathbf{x}$ if \mathbf{x} has a constant dimension.

Logistic Regression: Learning

Objective: maximizing the dataset probability, under the assumption that each example is sampled independently.

$$\hat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} \prod_{i=1}^n \sigma(\boldsymbol{\theta} \cdot \mathbf{x}_i)^{y_i} (1 - \sigma(\boldsymbol{\theta} \cdot \mathbf{x}_i))^{1-y_i}$$

For better numerical representation, we usually take the negative logarithm of the probability as the loss and minimize it:

$$Loss(\boldsymbol{\theta}) = - \sum_{i=1}^n y_i \log \sigma(\boldsymbol{\theta} \cdot \mathbf{x}_i) + (1 - y_i) \log(1 - \sigma(\boldsymbol{\theta} \cdot \mathbf{x}_i))$$

Gradient descent:

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} Loss(\boldsymbol{\theta})$$

Logistic Regression: Inference

Compute whether class 0 or 1 has larger probability.

What if there are more than 2 classes?

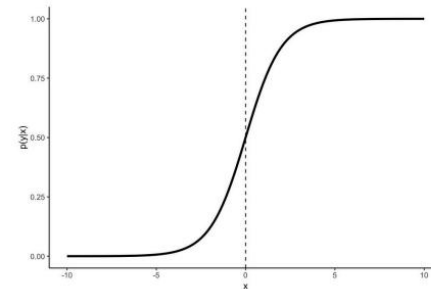
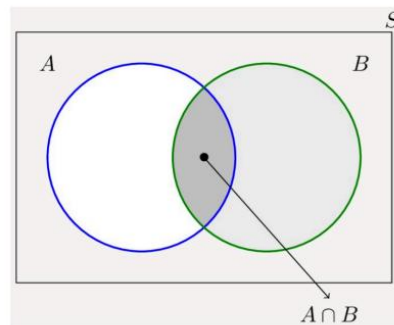
- 1 vs. 1 for $\frac{C(C-1)}{2}$ class pairs and do voting; or
- 1 vs. all for C classes and do argmax.
- Probabilistic interpretation (over classes) no longer hold for multi-class classification with logistic regression.

Generative vs. Discriminative Model for Classification

- Generative models: $P(s, y)$ is accessible when modeling $P(y | s)$.
- Discriminative models: $P(y | s)$ is directly modeled.

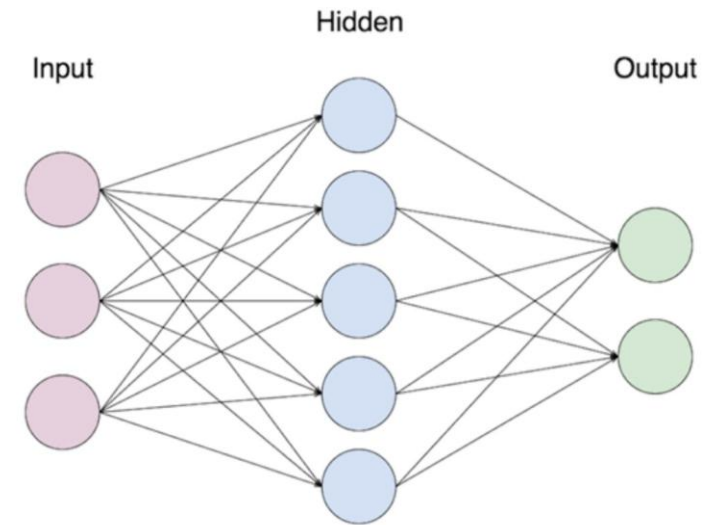
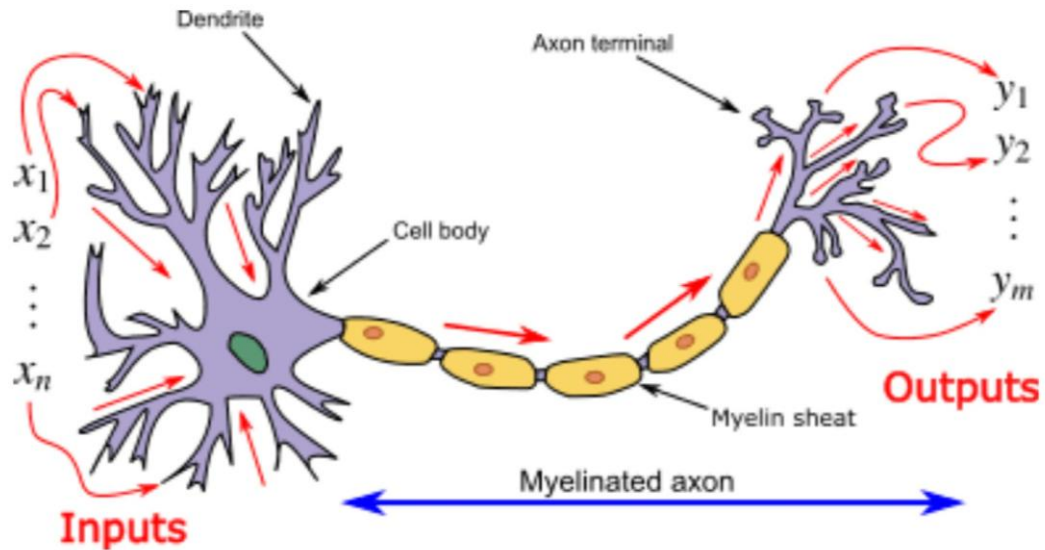
What are the differences?

Key difference: if you can generate a new data example once you have the model.



Neural-Network Classifier

- A neural network is a function
 - It has inputs and outputs.
 - “Neural modeling” now is better thought of as dense representation learning.



Neural Implementation of Scoring Function

Inference: solve arg max

Modeling: define score function

$$\text{classify}(s) = \underset{y}{\operatorname{argmax}} \operatorname{score}(s, y; \Theta)$$

Learning: choose parameter

With a neural network--based function f_{θ} , we input \mathbf{x} and collect a vector $\hat{\mathbf{y}} = f_{\mathbf{w}}(\mathbf{x})$; $\operatorname{score}(\mathbf{x}, y; \boldsymbol{\theta})$ is defined by selecting the corresponding entry in $\hat{\mathbf{y}}$.

Common Neural Network Notations

\mathbf{u} : a vector

u_i : entry i in the vector

\mathbf{W} : a matrix

$w_{i,j}$: entry (i,j) in the matrix

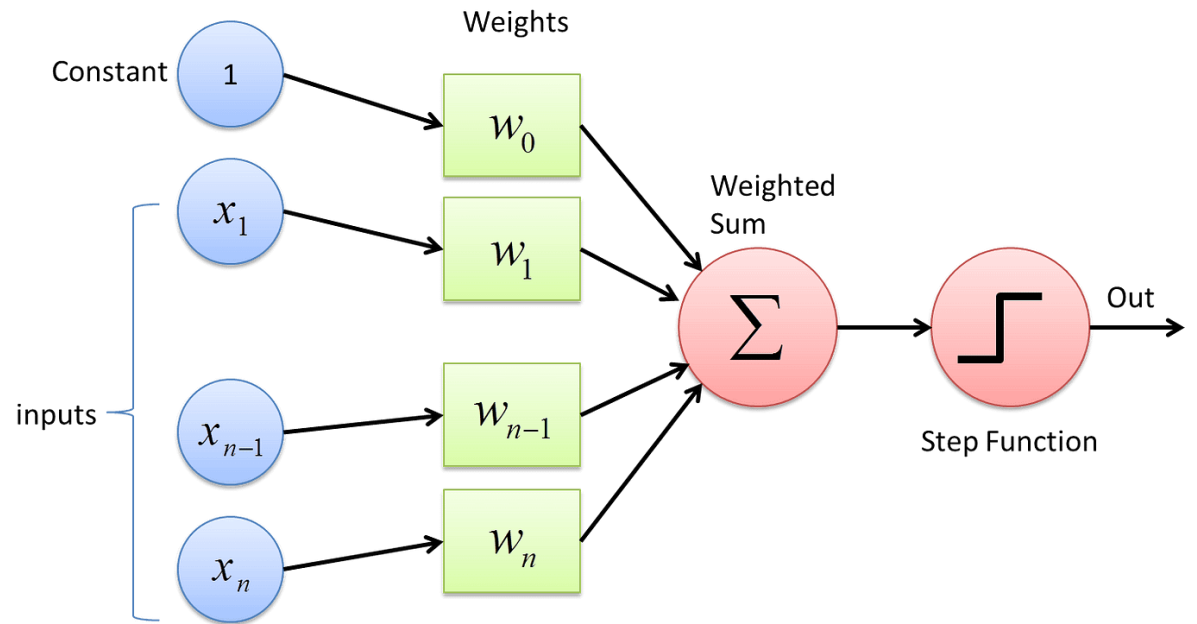
Perceptron

$$\text{perceptron}(\mathbf{x}; \boldsymbol{\theta}) = \text{step}(\boldsymbol{\theta}^T \mathbf{x})$$

activation
function

affine
transform

$$\text{step}(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$$



Perceptron: Learning

$$\text{perceptron}(\mathbf{x}; \boldsymbol{\theta}) = \text{step}(\boldsymbol{\theta}^T \mathbf{x})$$

Predict the label: $\hat{y}^{(i)} = \text{perceptron}(\mathbf{x}^{(i)}; \boldsymbol{\theta})$

Update weights: $\boldsymbol{\theta} = \boldsymbol{\theta} + \eta \cdot (y^{(i)} - \hat{y}^{(i)})\mathbf{x}^{(i)}$

$\hat{y}^{(i)}$	$y^{(i)}$	What happens
0	0	nothing
1	1	nothing

Perceptron: Learning

$$\text{perceptron}(\mathbf{x}; \boldsymbol{\theta}) = \text{step}(\boldsymbol{\theta}^T \mathbf{x})$$

Predict the label: $\hat{y}^{(i)} = \text{perceptron}(\mathbf{x}^{(i)}; \boldsymbol{\theta})$

Update weights: $\boldsymbol{\theta} = \boldsymbol{\theta} + \eta \cdot (y^{(i)} - \hat{y}^{(i)})\mathbf{x}^{(i)}$

$y^{(i)}$	$\hat{y}^{(i)}$	What happens
0	0	nothing
1	1	nothing
0	1	$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \mathbf{x}^{(i)}$

Perceptron: Learning

$$\text{perceptron}(\mathbf{x}; \boldsymbol{\theta}) = \text{step}(\boldsymbol{\theta}^T \mathbf{x})$$

Predict the label: $\hat{y}^{(i)} = \text{perceptron}(\mathbf{x}^{(i)}; \boldsymbol{\theta})$

Update weights: $\boldsymbol{\theta} = \boldsymbol{\theta} + \eta \cdot (y^{(i)} - \hat{y}^{(i)})\mathbf{x}^{(i)}$

$y^{(i)}$	$\hat{y}^{(i)}$	What happens
0	0	nothing
1	1	nothing
0	1	$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \mathbf{x}^{(i)}$
1	0	$\boldsymbol{\theta} = \boldsymbol{\theta} + \eta \mathbf{x}^{(i)}$

Neural Layer: Generalized Perceptron

A neural layer = affine transformation + nonlinearity

$$\text{perceptron}(\mathbf{x}) = \text{step}(\boldsymbol{\theta}^T \mathbf{x} + b)$$

activation function

affine transform

$$\text{neural_layer}(\mathbf{x}) = g(\mathbf{W}^T \mathbf{x} + \mathbf{b})$$

Output is a vector (results from multiple independent perceptrons).

Can have other activation functions for nonlinearity.

Stacking Neural Layers

Multiple neural layers can be stacked together.

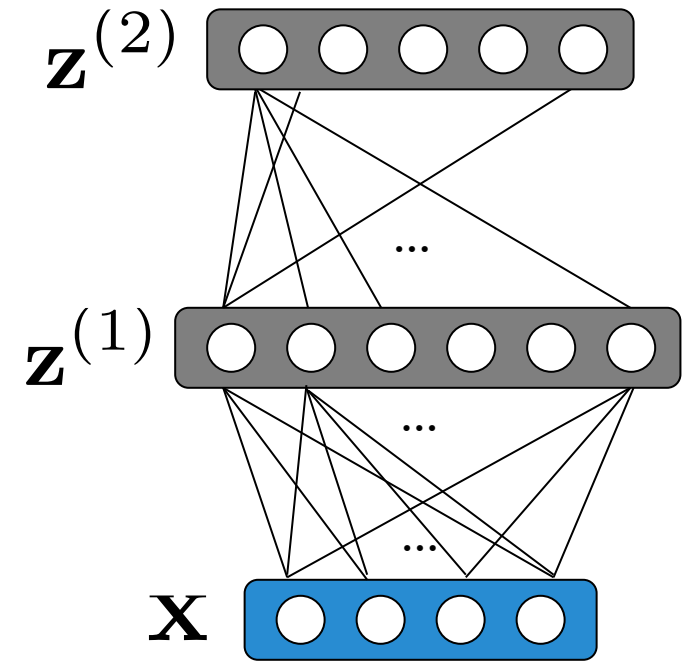
hidden units $\mathbf{z}^{(1)} = g \left(\mathbf{W}^{(0)} \mathbf{x} + \mathbf{b}^{(0)} \right)$

layer index

$$\mathbf{z}^{(2)} = g \left(\mathbf{W}^{(1)} \mathbf{z}^{(1)} + \mathbf{b}^{(1)} \right)$$

...

- Use output of one layer as input to the next
- *Feed-forward* and/or *fully-connected* layers
- Also called *multi-layer perceptron (MLP)*



Nonlinearities (activation function)

$$\mathbf{z}^{(1)} = g \left(\mathbf{W}^{(0)} \mathbf{x} + \mathbf{b}^{(0)} \right)$$

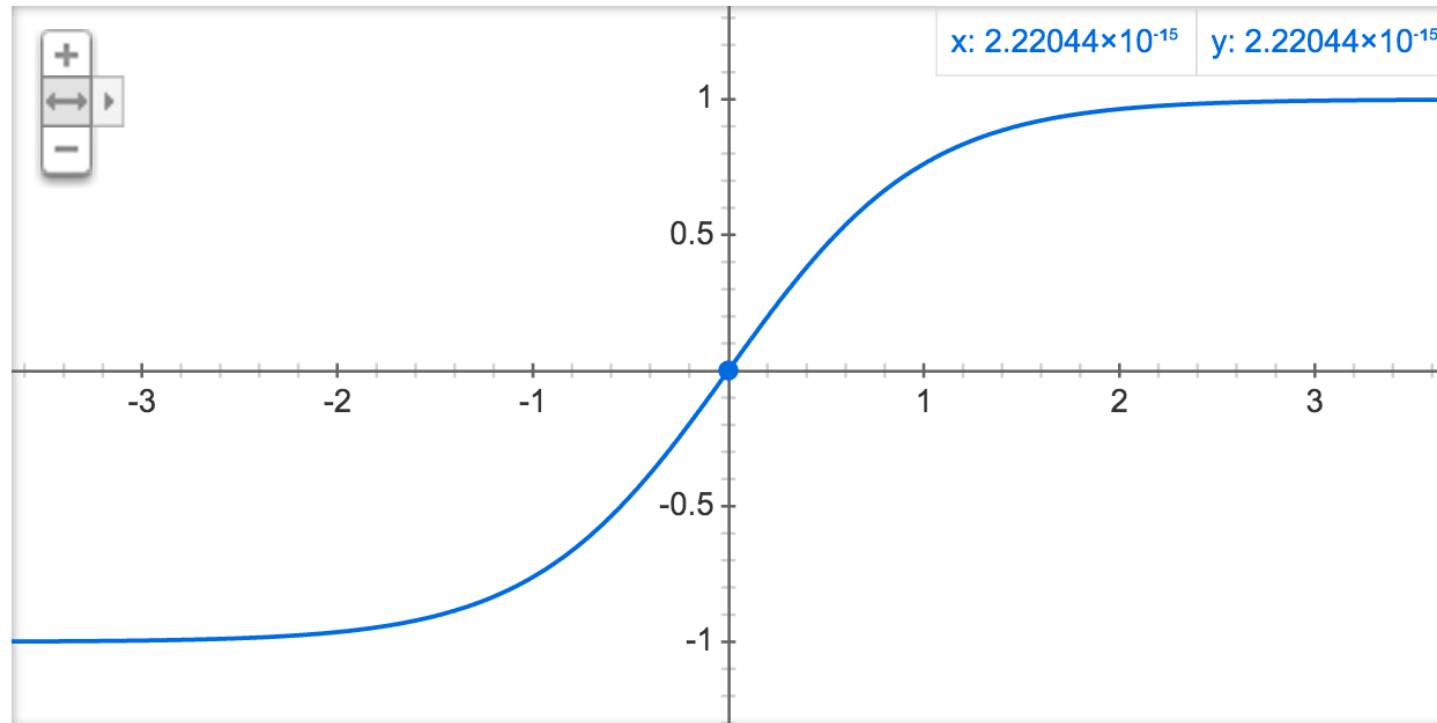
- g can be applied to each entry in a vector in an element-wise manner
- Common activation functions: tanh, sigmoid, and ReLU

Why nonlinearities?

- Otherwise stacking neural layers results in a simple affine transform.

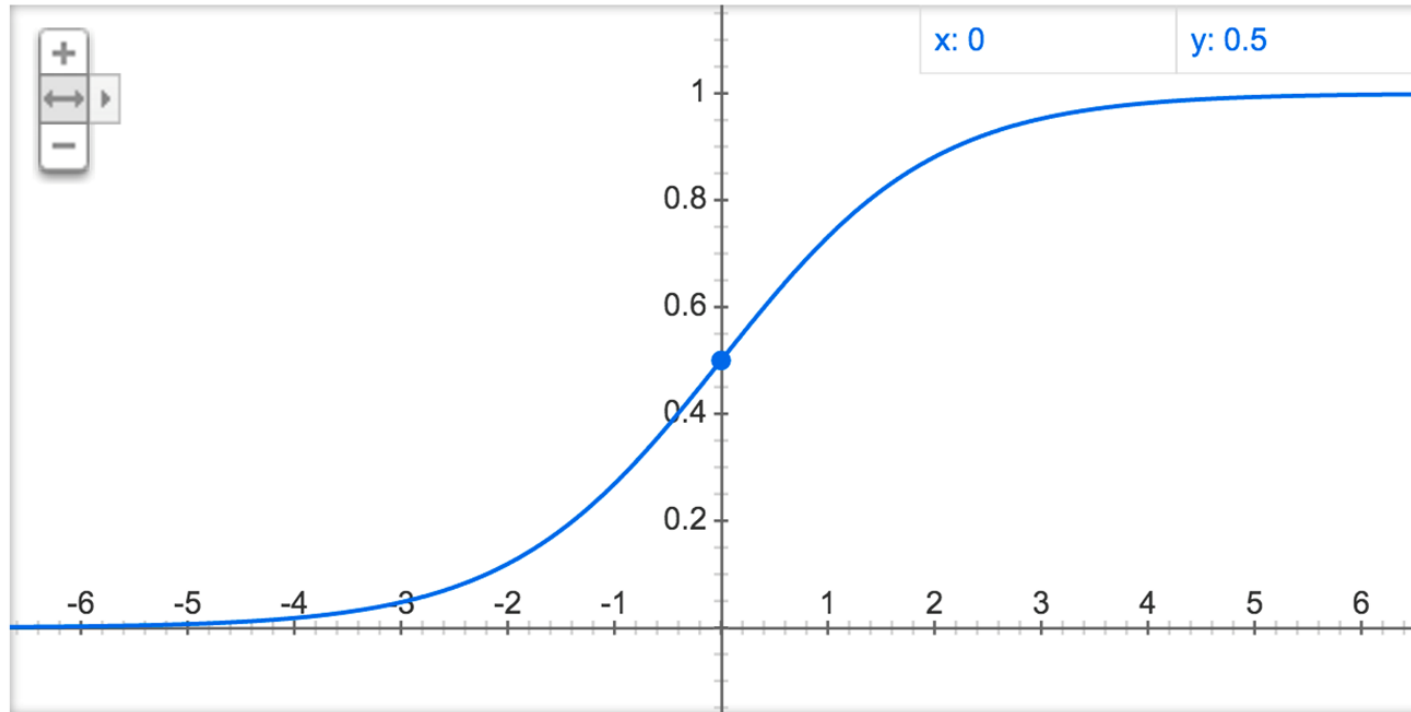
Nonlinearities: tanh

$$y = \tanh(x) = \frac{e^x + e^{-x}}{e^x - e^{-x}}$$



Nonlinearities: sigmoid

$$y = \sigma(x) = \frac{1}{1 + e^{-x}}$$



Sigmoid and tanh

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

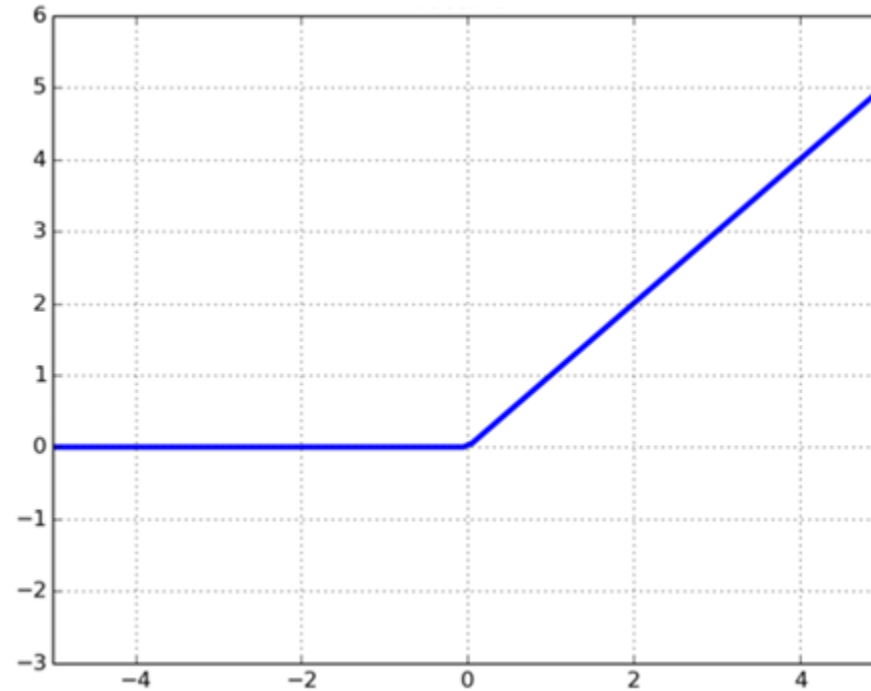
divide both sides by e^x

$$= \frac{1 - e^{-2x}}{1 + e^{-2x}}$$

$$= \sigma(2x) - (1 - \sigma(2x)) = 2\sigma(2x) - 1$$

Nonlinearity: Rectified Linear Unit (ReLU)

$$y = \text{ReLU}(x) = \max\{0, x\}$$



Sentiment Classification with Neural Network

$$\mathbf{z}^{(1)} = g \left(\mathbf{W}^{(0)} \mathbf{x} + \mathbf{b}^{(0)} \right)$$
$$\mathbf{s} = \mathbf{W}^{(1)} \mathbf{z}^{(1)} + \mathbf{b}^{(1)}$$

We empirically don't pass the final layer into an activation function

How can we get \mathbf{x} for a sentence?

- Average word embeddings
- More complicated neural network structures

Neural Network: Training

$$\mathbf{s} = \begin{bmatrix} \text{score}(\mathbf{x}, 0; \mathbf{w}) \\ \text{score}(\mathbf{x}, 1; \mathbf{w}) \end{bmatrix} \xRightarrow{\text{softmax}} \mathbf{p} = \begin{bmatrix} \frac{e^{\text{score}(\mathbf{x}, 0; \mathbf{w})}}{Z} \\ \frac{e^{\text{score}(\mathbf{x}, 1; \mathbf{w})}}{Z} \end{bmatrix}$$

$$Z = e^{\text{score}(\mathbf{x}, 0; \mathbf{w})} + e^{\text{score}(\mathbf{x}, 1; \mathbf{w})}$$

Maximize the probability (minimize negative log probability) of gold standard label

$$\text{loss} = -\log P(y \mid \mathbf{x}) = -\log \left(\frac{e^{\text{score}(\mathbf{x}, y; \mathbf{w})}}{Z} \right)$$

Also called *cross entropy* (between \mathbf{p} and the 1-hot gold standard distribution) loss.

Backpropagation

- Chain rule: suppose $y = f(x)$, $z = g(y)$, then $\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx}$

$$\mathbf{z}^{(1)} = g\left(\mathbf{W}^{(0)}\mathbf{x} + \mathbf{b}^{(0)}\right)$$

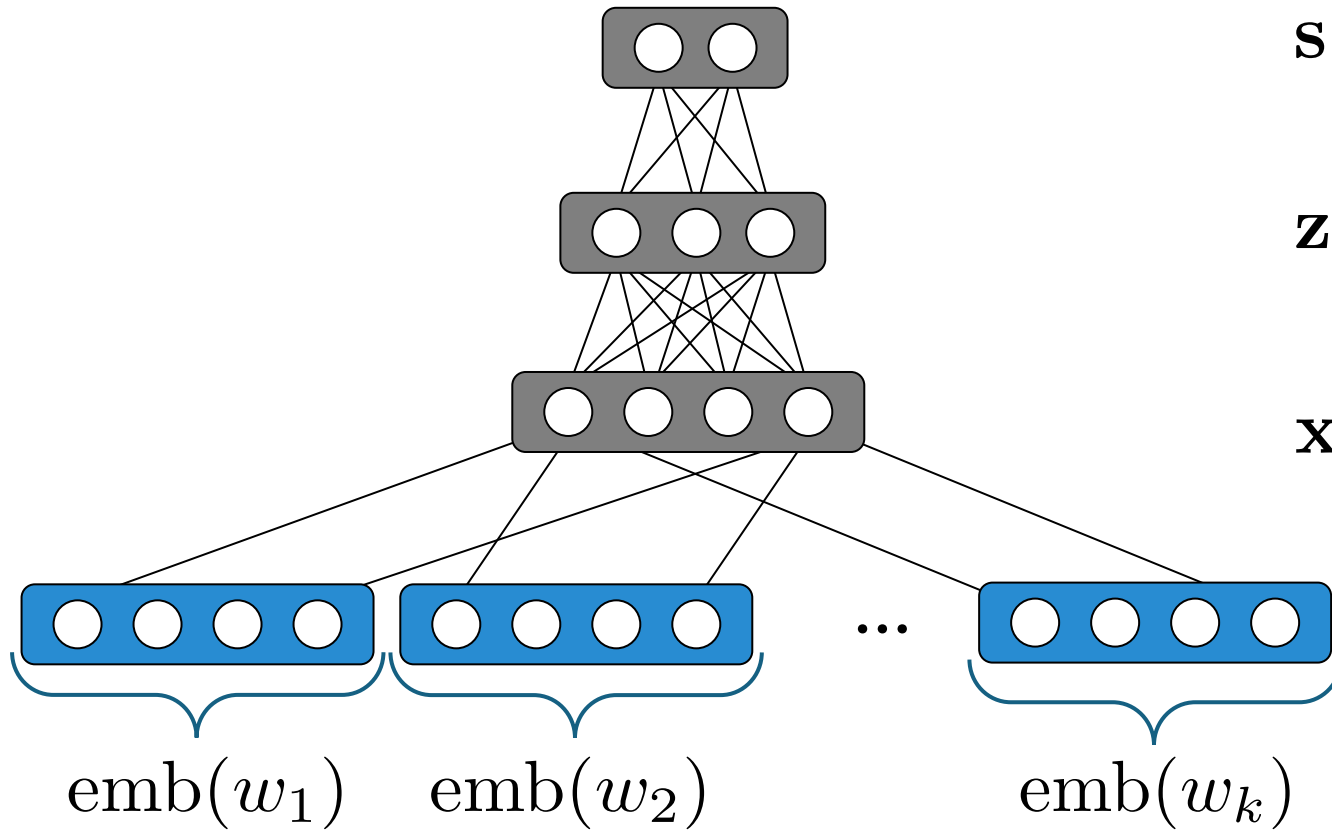
$$\mathbf{s} = \mathbf{W}^{(1)}\mathbf{z}^{(1)} + \mathbf{b}^{(1)}$$

$$\text{loss}(\mathbf{s}; \mathbf{x}, y) = -\log\left(\frac{e^{s_y}}{Z}\right)$$

Now we have $\frac{\partial \text{loss}}{\partial \mathbf{s}}$, how should we update $\mathbf{W}^{(0)}$?

$$\frac{\partial \text{loss}}{\partial \mathbf{W}^{(0)}} = \frac{\partial \text{loss}}{\partial \mathbf{s}} \frac{\partial \mathbf{s}}{\partial \mathbf{z}^{(1)}} \frac{\partial \mathbf{z}^{(1)}}{\partial \mathbf{W}^{(0)}}$$

Visualization of an Average Word Vector Classifier



$$\mathbf{s} = \mathbf{W}^{(1)}\mathbf{z}^{(1)} + \mathbf{b}^{(1)}$$

$$\mathbf{z}^{(1)} = g\left(\mathbf{W}^{(0)}\mathbf{x} + \mathbf{b}^{(0)}\right)$$

$$\mathbf{x} = \frac{1}{k} \sum_{i=1}^k \text{emb}(w_k)$$

Next

Common neural network architectures for NLP (for more powerful models).