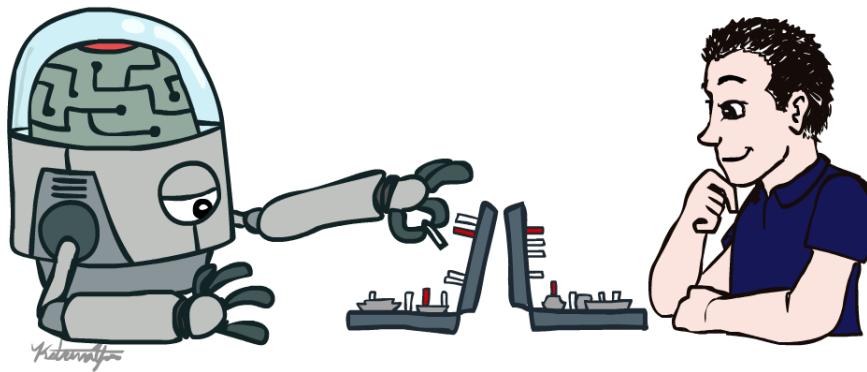


CSE 3521: Introduction to Artificial Intelligence



[Many slides are adapted from the [UC Berkeley. CS188 Intro to AI](#) at UC Berkeley and previous CSE 3521 course at OSU.]



THE OHIO STATE UNIVERSITY

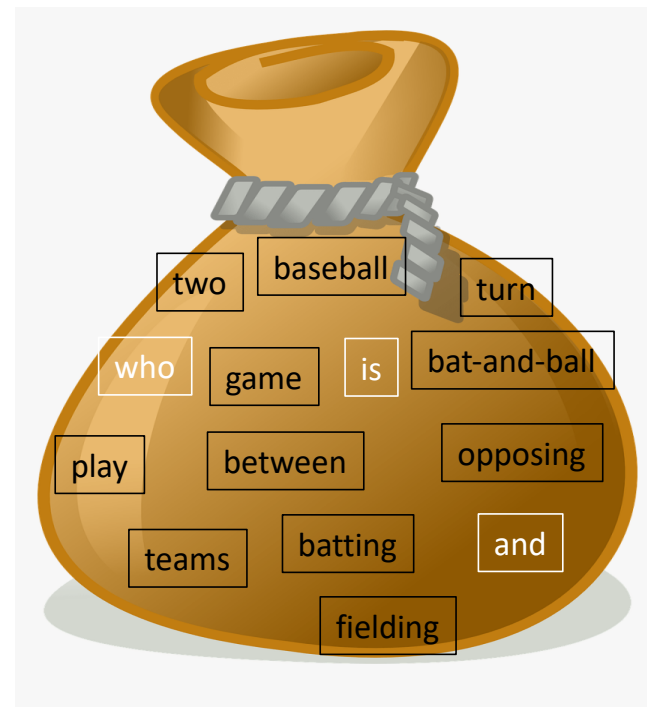
Bag of Words (BOW)

A simplified representation of sentences (documents) for classification or retrieval by **counting the occurrences** of each unique word (or phrase)

- ignore the grammar
- ignore the word order
- keep the word counts (frequencies)
- lead to a “fixed”-size vector representation

“Baseball” “is” “a” “bat-and-ball”
“game” “played” “between”
“two” “opposing” “teams”
“who” “take” “turns” “batting”
“and” “fielding” “.”

A sequence of tokens (and their indices)



Bag of Words (BOW)

Given:

- A dictionary, vocabulary, or codebook: $f(\text{token}) \rightarrow \text{index} \in \{1, \dots, D\}$ or N/A

- An all-0 vector: $\mathbf{x} = \begin{bmatrix} x[1] \\ \vdots \\ x[D] \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$

Bag of Words (BOW)

Given:

- A dictionary, vocabulary, or codebook: $f(\text{token}) \rightarrow \text{index} \in \{1, \dots, D\}$ or N/A

- An all-0 vector: $\mathbf{x} = \begin{bmatrix} x[1] \\ \vdots \\ x[D] \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$

Input: A sequence of word tokens $w[1], w[2], \dots, w[M]$

for $m = 1 : M$

if $f(w[m]) \neq \text{N/A}$

$x[f(w[m])] += 1$

end

end

Return: \mathbf{x}

Bag of Words (BOW)

Given:

- A dictionary, vocabulary, or codebook: $f(\text{token}) \rightarrow \text{index} \in \{1, \dots, D\}$ or N/A

- An all-0 vector: $\mathbf{x} = \begin{bmatrix} x[1] \\ \vdots \\ x[D] \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$

Input: A sequence of word tokens $w[1], w[2], \dots, w[M]$

for $m = 1 : M$

if $f(w[m]) \neq \text{N/A}$

$x[f(w[m])] += 1$

end

end

Return: \mathbf{x}

- **Out-of-vocabulary**
- **“Stop” words:** too frequent in all sentences/documents and less informative in differentiating them (e.g., “is”, “are”, “and”,)

Word token counts

Example: BOW for classification

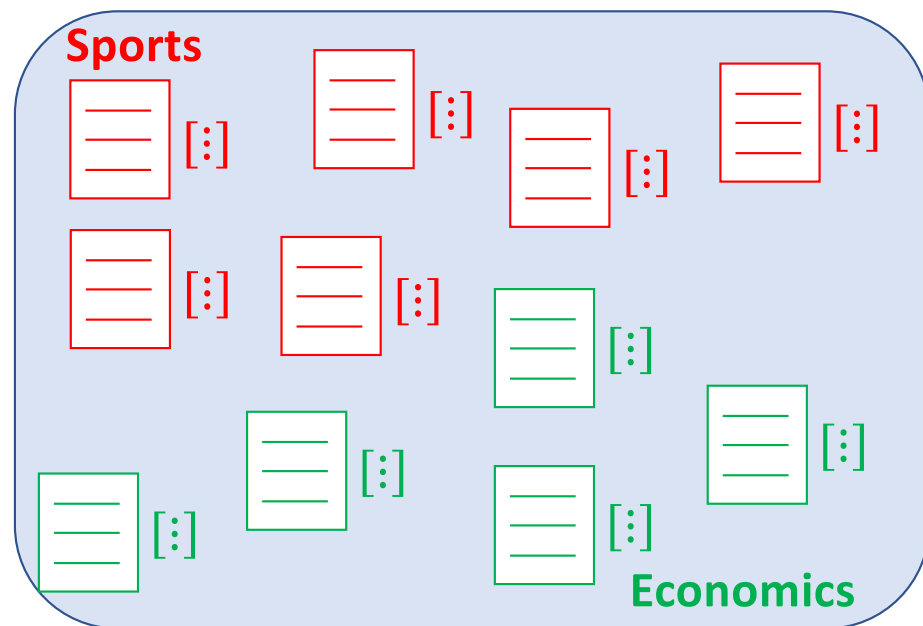
- Consider a document classification:
 - {sports, economics}
- Nearest neighbor classification (NNC)
 - Compute distance to all training documents, each of them is represented by BOW
 - Output the “label” of the nearest document

Baseball is a bat-and-ball game played between two opposing teams who take turns batting and fielding.

[:]

compute distance

Training documents



Example: BOW for classification

Distance

- Euclidean (L_2) distance:

$$\|\mathbf{x} - \mathbf{x}_n\|_2 = \left(\sum_{d=1}^D (x[d] - x_n[d])^2 \right)^{1/2}$$

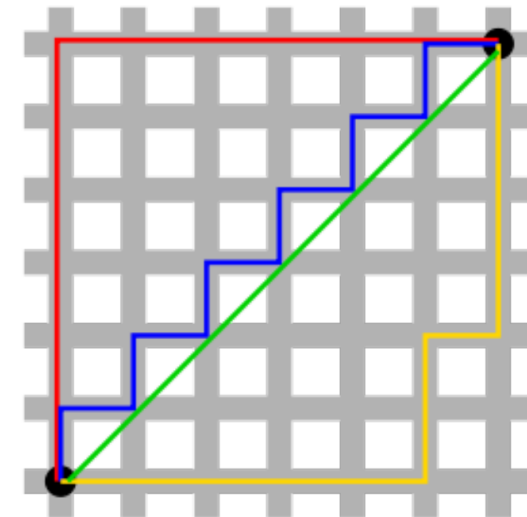
- L_1 distance:

$$\|\mathbf{x} - \mathbf{x}_n\|_1 = \sum_{d=1}^D |x[d] - x_n[d]|$$

- L_p norm:

$$\|\mathbf{x} - \mathbf{x}_n\|_p = \left(\sum_{d=1}^D |x[d] - x_n[d]|^p \right)^{1/p}$$

\mathbf{x} : 2-dimensional vectors



Green line is Euclidean distance.
Red, Blue, and Yellow lines are L_1 distance

Example: BOW for classification

Given: a distance metric **dis**, a training set = $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$, $t = \infty$ (min distance), y (predicted label)

Input: a test data instance \mathbf{x}

for $n = 1 : N$

if $\text{dis}(\mathbf{x}, \mathbf{x}_n) < t$

$y = y_n$

$t = \text{dis}(\mathbf{x}, \mathbf{x}_n)$

end

end

Return: y



Label:

$y \in \{\text{sports, economics}\}$

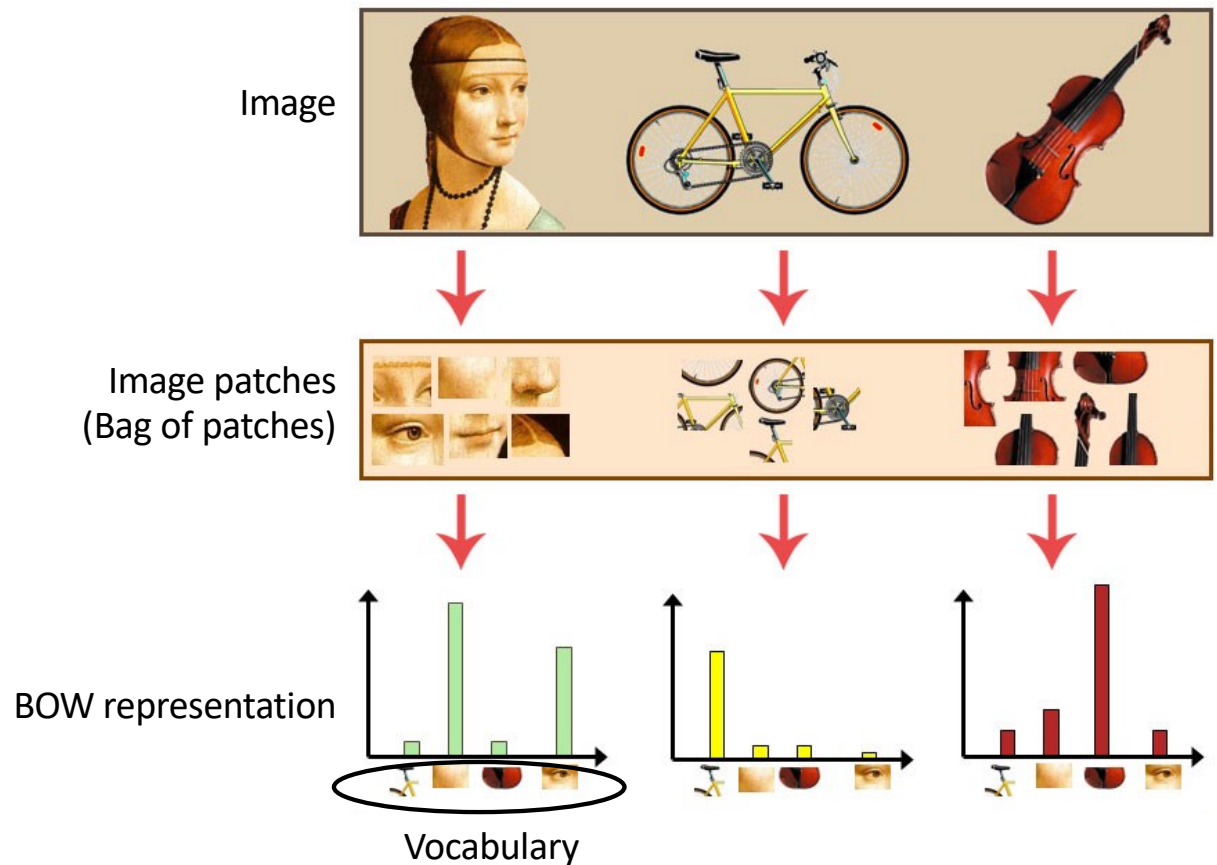
$y \in \{-1, 1\}$ or $\{0, 1\}$

$y \in \{1, \dots, C\}$

Bag of Words (BOW) for images

BOW can be applied to vision

- Vocabulary: D image patches
- Each image: a bag of patches
- BOW representation: D -dim
 - For each image patch $I[m]$, find the “nearest” patch in the vocabulary and use its index $\in \{1, \dots, D\}$ to represent $I[m]$
 - Count the time that each index shows up in the image



Bag of Words (BOW)

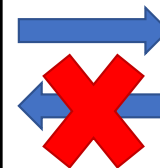
Pros:

- Simplified, easily-understandable, fixed-size

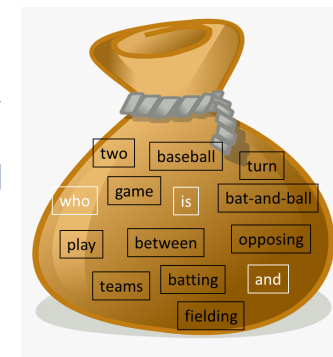
What do we lose?

- What the original sequence was

Baseball is a bat-and-ball
game played between two
opposing teams who take
turns batting and fielding.



x



Bag of Words (BOW)

Cons:

- missing the **sequential** information (e.g., “sheep follow wolfs” vs. “wolfs follow sheep”)
- not **normalized** (e.g., comparing long vs. short documents)
- words treated as **independent** (e.g., no synonyms and antonyms)
- highly sparse, **high dimensional**

N-gram vocabulary

Cons 1: Missing the sequential information

- One solution: **N-gram vocabulary**
- 1-gram (unigram): “sheep”, “follow”, “wolves”
- 2-gram (bigram): “sheep-follow”, “follow-wolves”
- 3-gram: “sheep-follow-wolves”
- ...
- Size of the vocabulary: $D + D^2 + D^3$
- “Sheep follow wolves” becomes:

	...	
“sheep”	1	}
	...	
“follow”	1	}
	...	
“wolves”	1	}
	...	
“sheep-follow”	1	
	...	
“follow-wolves”	1	}
	...	
	...	
“sheep-follow-wolves”	1	}
	...	

1-gram

2-gram

3-gram

Dataset-independent normalization

Cons 2: not **normalized**

Sheep follow wolfs.

	...
"sheep"	1
	...
"follow"	1
	...
"wolfs"	1
	...



compute
distance



Sheep follow wolfs.
Sheep follow wolfs.

	...
"sheep"	2
	...
"follow"	2
	...
"wolfs"	2
	...

Dataset-independent normalization

Cons 2: not **normalized**

- One solution: **feature normalization**
- Vector norm (i.e., length or magnitude):

$$\|\mathbf{x}\|_p = \left(\sum_{d=1}^D |x[d]|^p \right)^{1/p}$$

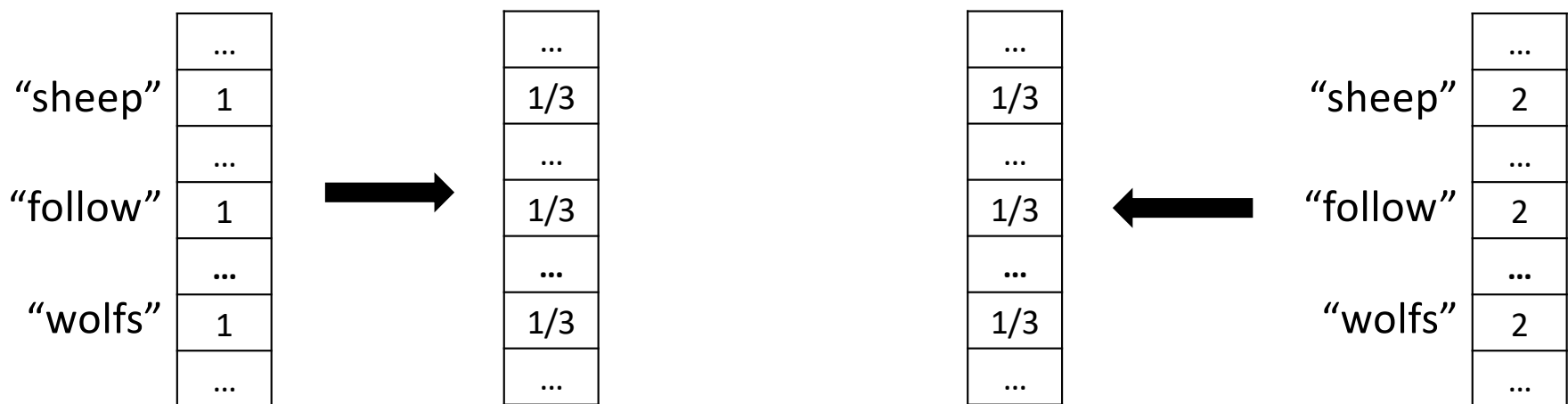
- **L_p normalization:**

$$\mathbf{z} = \frac{\mathbf{x}}{\|\mathbf{x}\|_p} = \begin{bmatrix} \frac{x[1]}{\|\mathbf{x}\|_p} \\ \vdots \\ \frac{x[D]}{\|\mathbf{x}\|_p} \end{bmatrix}$$

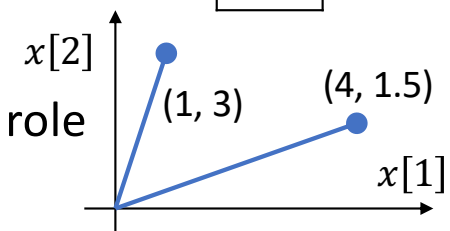
Dataset-independent normalization

Cons 2: not **normalized**

- One solution: **feature normalization**
- **L₁ normalization**: popular for counts (frequency, probability)



- **L₂ normalization**: widely used if the vector angle plays an important role



Bag of Words (BOW)

Cons:

- missing the sequential information (e.g., “sheep follow wolfs” vs. “wolfs follow sheep”)
- not normalized (e.g., comparing long vs. short documents)
- words treated as **independent** (e.g., no synonyms and antonyms)
- highly sparse, **high dimensional**

Summary

- Data and data representation (features)
 - Numerical vs. categorical variables
 - Feature extraction: from raw data to simplified, informative, non-redundant, or more interpretable representations
- Bag-of-words (BoW) representation
 - Fixed-sized representations for sentences and documents (and images)
 - Nearest neighbor classification based on distance metrics