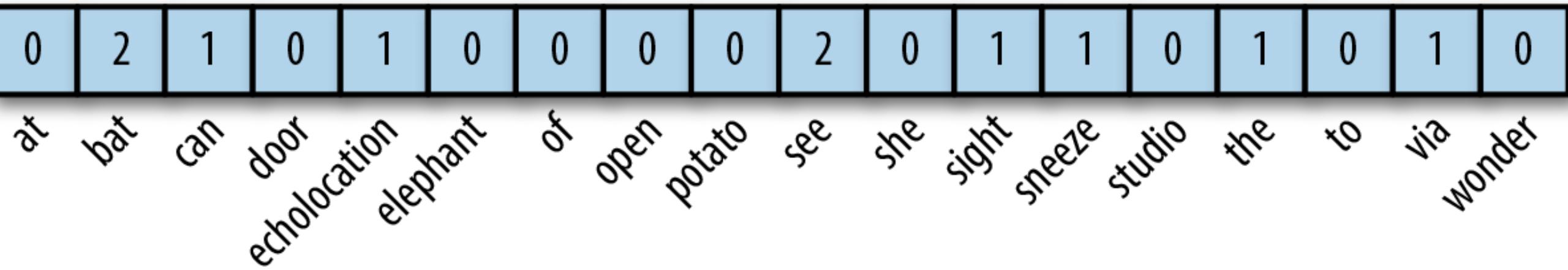
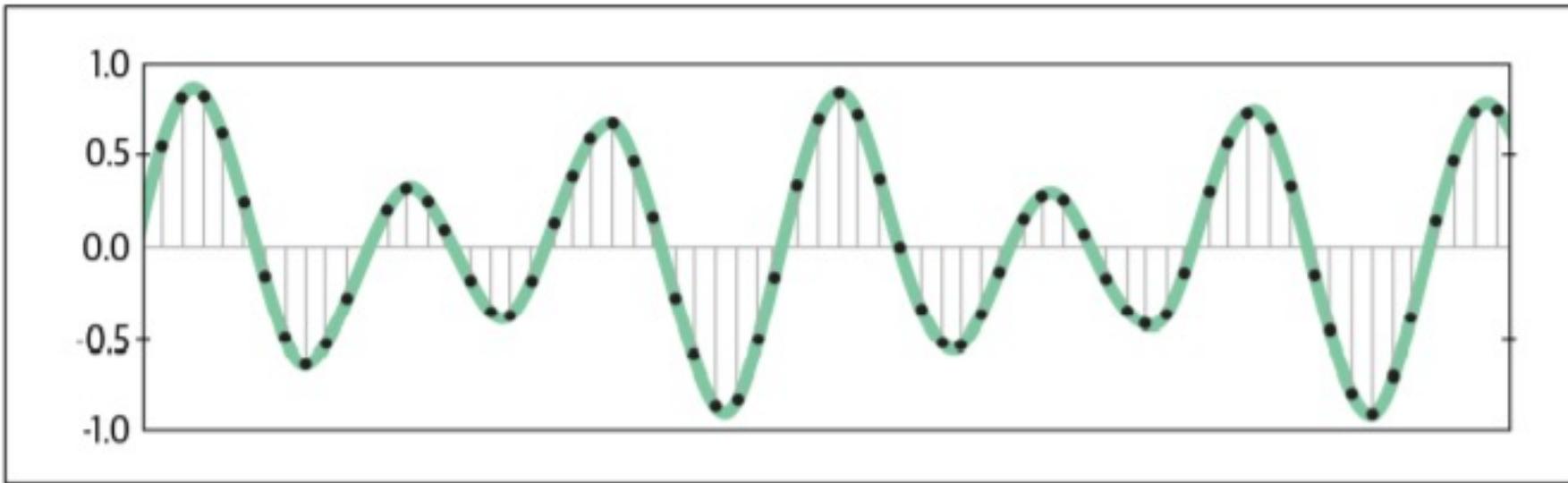


# Topic 3

## Text Representation



# Speech



```
[-1274, -1252, -1160, -986, -792, -692, -614, -429, -286, -134, -57, -41,
-169, -456, -450, -541, -761, -1067, -1231, -1047, -952, -645, -489, -448,
-397, -212, 193, 114, -17, -110, 128, 261, 198, 390, 461, 772, 948, 1451,
1974, 2624, 3793, 4968, 5939, 6057, 6581, 7302, 7640, 7223, 6119, 5461,
4820, 4353, 3611, 2740, 2004, 1349, 1178, 1085, 901, 301, -262, -499,
-488, -707, -1406, -1997, -2377, -2494, -2605, -2675, -2627, -2500, -2148,
-1648, -970, -364, 13, 260, 494, 788, 1011, 938, 717, 507, 323, 324, 325,
350, 103, -113, 64, 176, 93, -249, -461, -606, -909, -1159, -1307, -1544]
```

Figure 3-4. Speech signal represented by a numerical vector

# Imaging



What We See

```
08 02 22 97 38 15 00 40 00 75 04 05 07 78 52 12 50 77 91 08  
49 49 99 40 17 81 18 57 60 87 17 40 98 43 69 48 04 56 62 00  
81 49 31 73 55 79 14 29 93 71 40 67 53 88 30 03 49 13 36 65  
52 70 95 23 04 60 11 42 69 24 68 56 01 32 56 71 37 02 36 91  
22 31 16 71 51 67 63 89 41 92 36 54 22 40 40 28 66 33 13 80  
24 47 32 60 99 03 45 02 44 75 33 53 78 36 84 20 35 17 12 50  
32 98 81 28 64 23 67 10 26 38 40 67 59 54 70 66 18 38 44 70  
67 26 20 68 02 62 12 20 95 63 94 39 63 08 40 91 66 49 94 21  
24 55 58 05 66 73 99 26 97 17 78 78 96 83 14 88 34 89 43 72  
21 36 23 09 75 00 76 44 20 45 35 14 00 41 33 97 34 31 33 95  
78 17 53 28 22 75 31 67 15 94 03 80 04 42 16 14 09 53 56 92  
16 39 05 42 96 35 31 47 55 58 88 24 00 17 54 24 36 29 85 57  
86 56 00 48 35 71 89 07 05 44 44 37 44 60 21 58 51 54 17 58  
19 80 81 68 05 94 47 69 28 73 92 13 86 52 17 77 04 89 55 40  
04 52 08 83 97 35 99 16 07 97 57 32 16 26 26 79 33 27 98 66  
88 36 68 87 57 62 20 72 03 46 33 67 46 55 12 32 63 93 53 69  
04 42 16 73 38 25 39 11 24 94 72 18 08 46 29 32 40 62 76 36  
20 69 36 41 72 30 23 88 34 62 99 69 82 47 59 85 74 04 36 16  
20 73 35 29 78 31 90 01 74 31 49 71 48 86 81 16 23 57 05 54  
01 70 54 71 83 51 54 69 16 92 33 48 61 43 52 01 89 19 67 48
```

What Computers See

Figure 3-2. How we see an image versus how computers see it [1]

# What can we do with a string?

- Many methods cannot process strings or plain text
- Need to work on a numeric feature space
- Feature extraction, also known as *vectorization*
  - Character-level representation
  - Word-level representation
  - Document-level representation

# Vector Representation

- Mathematical model that represents text units as vectors
  - What is a vector?
- Different strategies that capture the linguistic properties of a text to different degrees (but the strategies we cover will fall under vector space models - VSM)

# One Hot Encoding

apple	eat	horse
1	0	0
0	1	0
0	0	1

“The horse eats the apple.”

preprocessing

“horse eat apple.”

vectorization

$[[0\ 0\ 1]\ [0\ 1\ 0]\ [1\ 0\ 0]]$

# One Hot Encoding

```
#Get one hot representation for any string based on this vocabulary.  
#If the word exists in the vocabulary, its representation is returned.  
#If not, a list of zeroes is returned for that word.  
def get_onehot_vector(somestring):  
    onehot_encoded = []  
    for word in somestring.split():  
        temp = [0]*len(vocab)  
        if word in vocab:  
            temp[vocab[word]-1] = 1 # -1 is to take care of the fact indexing  
            onehot_encoded.append(temp)  
    return onehot_encoded
```

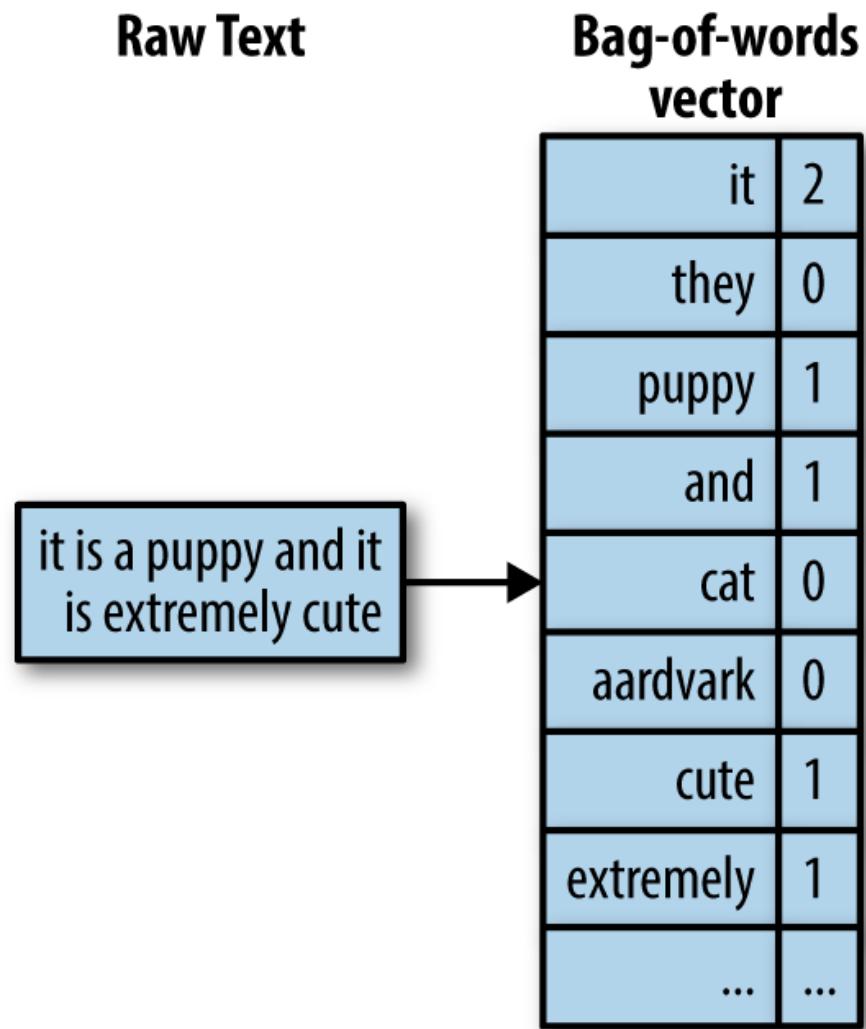
if word not present, it's 0  
drawback : if sth not in vocab, it can't be represented

What happens if we add new documents?

# Drawbacks

- All words equidistant, i.e. normalization extra important
- Relationships between tokens lost
- Scalability
- Sparsity
- Out of vocabulary words (OOV)

# Bag-of-Words

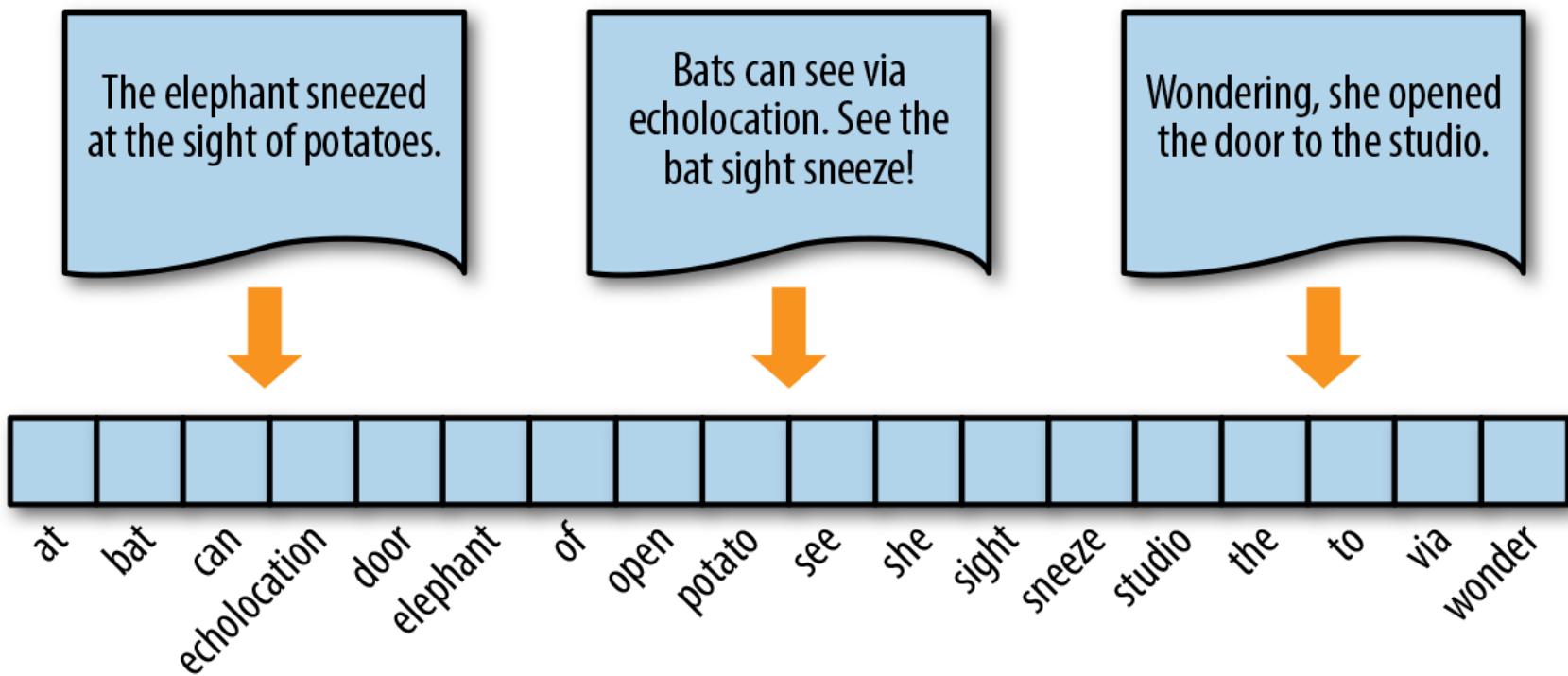


- Represent text as the bag (multiset) of its words
  - Keep multiplicity
  - Commonly used for classification problems

Text book dh3 GitHub have example  
of how to build vocabulary

# Bag-of-Words

- What numeric values should be put into this vector?



# Drawbacks

"dog likes man" and "man likes dogs"  
means the same thing

- The size of the vector still increases with the size of the vocabulary
  - Sparsity
- OOV words
- Disregarding grammar, word order
- Still can't capture information about words that mean the same thing

# Bag of N-grams representation

*"The cat sat on the mat"*

deal with frequency  
of  $n$  words instead of  
1 word

We could break this up into:

- **1-grams:** The, cat, sat, on, the, mat
- **2-grams:** The cat, cat sat, sat on, on the, the mat
- **3-grams:** The cat sat, cat sat on, sat on the, on the mat
- etc.

What will happen as N increases for an N-gram representation?

What will happen as N increases for an N-gram representation?

- We still have sparsity and OOV problems

# Term frequency (TF) Representation

- TF representation weights words proportionally to their frequency
  - Most frequent words not always most informative
- 
- $\text{TF}(t, d) = \frac{\text{\# of occurrences of term } t \text{ in } d}{\text{total \# of terms in } d}$

$t$  = term,  $d$  = document

# TF-IDF

- Inverse-Document-Frequency
- Basic idea: Common words may not add anything to our understanding of a specific text. Conversely, if a word occurs less frequently we might want to give it a larger weight in our representation

$$\text{IDF}(t) = \log \left( \frac{\text{total documents}}{\# \text{ documents with } t} \right)$$

$$\boxed{\text{TF-IDF} = \text{TF} * \text{IDF}}$$

# Activity

- Calculate the TF-IDF of term  $t$  for the following toy corpus:
- **Document 1:** "The cat sat on the mat."
- **Document 2:** "The mat was red."
- **Document 3:** "The cat liked the mat."

$$\frac{1}{6}$$

$$\text{cat} = \text{TF-IDF} = \frac{1}{6} \cdot \log\left(\frac{3}{2}\right)$$

$$\text{mat} = \frac{1}{6} \cdot \log\left(\frac{3}{3}\right) = 0$$

Revisit - Learn change it after so TF-IDF of mat not  
tobe 0

# Distributed Representations

- Encode text along a continuous scale with a distributed representation
- Example: Word2Vec

The elephant sneezed  
at the sight of potatoes.



-0.0225403

-0.0212964

0.02708783

0.0049877

0.0492694

-0.03268785

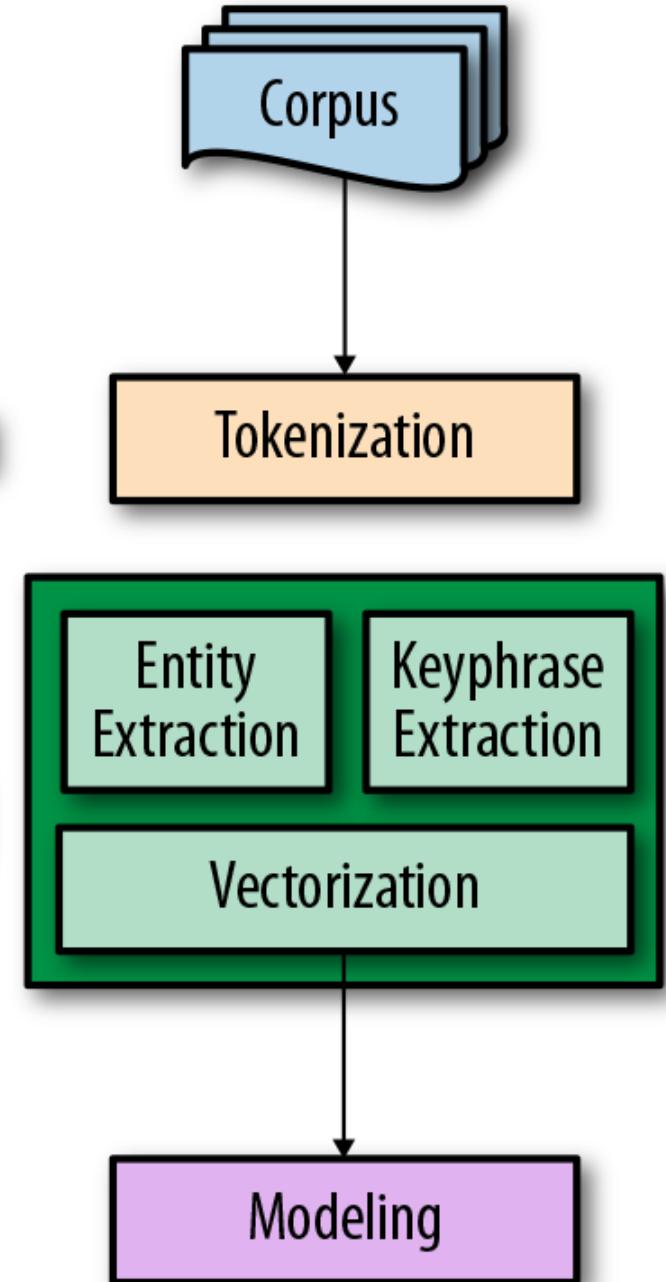
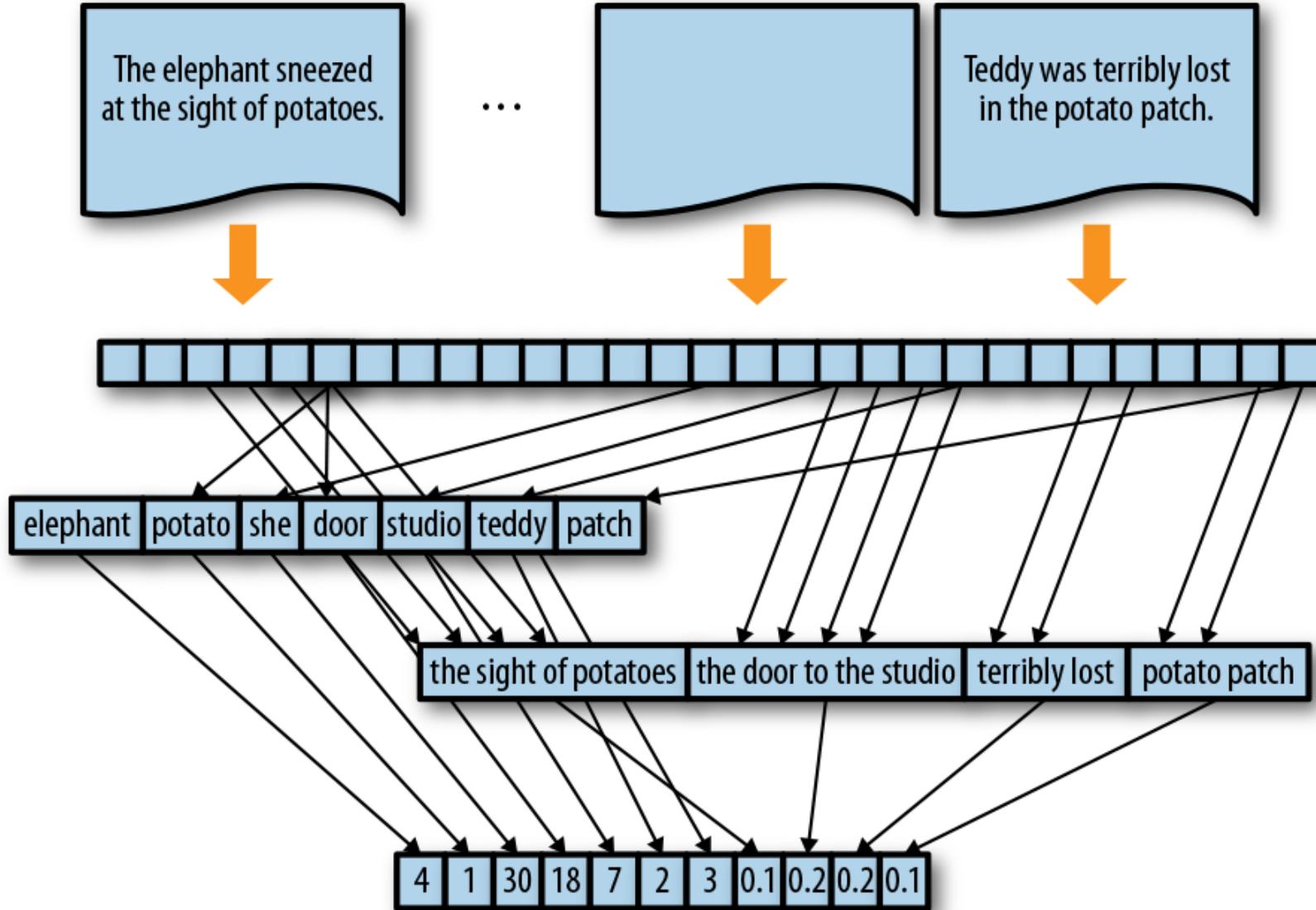
-0.0320941

Bats can see via  
echolocation. See the  
bat sight sneeze!



Wondering, she opened  
the door to the studio.





# Representing text as tensors

- What is a tensor?

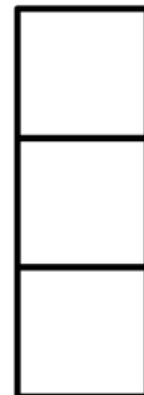
**NumPy Array:** [1 2 3]

**PyTorch Tensor:** tensor([1, 2, 3])

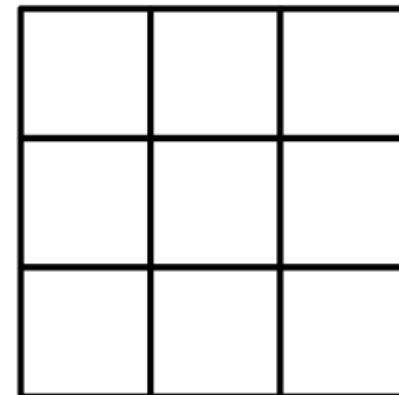
Scalar  
Rank 0 Tensor



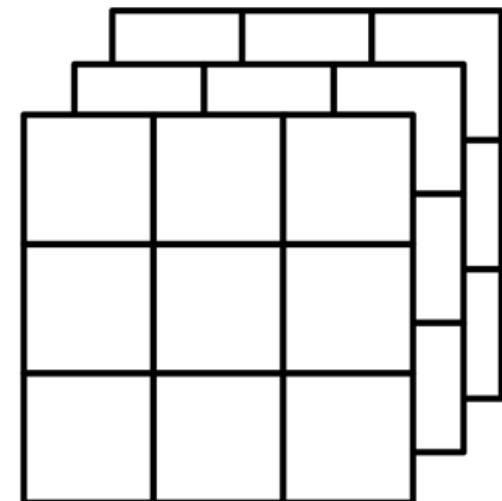
Scalar  
Rank 1 Tensor



Scalar  
Rank 2 Tensor



Rank 3 Tensor



# Why use tensors?

- More efficient
- That is what many architectures expect

# CUDA Tensors

- Using GPU instead of CPU
- Allocate the tensor on the GPU's memory

```
torch.cuda.is_available()
```

```
torch.device()
```

```
x = torch.rand(3, 3).to(device)
```

```
describe(x)
```

Next time:

- Crash course in neural networks and autoencoders

1, CountVectorizer : Create dictionary

$\rightarrow \text{Count\_vec} = \text{CountVectorizer()}$

$\text{Count\_vec}.fit\text{-transform}(\text{text})$

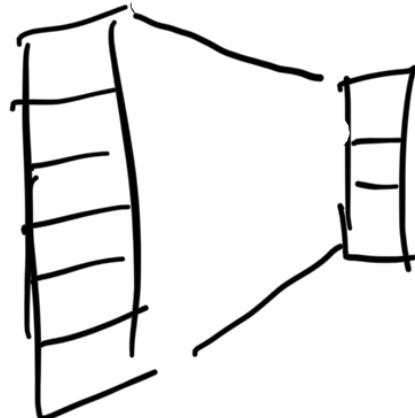
2, Create dict for series of words, not only single word

$\text{Count\_vec} = \text{CountVectorizer(ngram\_range} = (1, 3))$

3, TF-IDF  
Create dictionary for words but with normalize

(\*) If words doesn't exist, can't be represented (as not exist in  
dictionary)

- In Bag of word : Dog eats man and man eats dog is the same
- In Bag of n-gram : They are different as it conta the representation of phrase in addition to representation of words
- In TF-IDF , not exist = 1 at not exist = 0 but also include the intensity



$6 \times 3$

$$\begin{matrix} & \times & \begin{matrix} & & \\ & & \\ & & \\ & & \\ & & \\ & & \end{matrix} & = & (6 \times 1) \times & = \\ \begin{matrix} & & \\ & & \\ & & \end{matrix} & \times & \begin{matrix} & & \\ & & \\ & & \\ & & \\ & & \\ & & \end{matrix} & = & (6 \times 1) \times & = \end{matrix}$$