# Christopher Keown, Ph.D.
## UC San Diego

Department of Cognitive Science

ckeown@ucsd.edu

# **White Christmas**

Weather forecasting is perhaps the most familiar domain of predictive modeling. Short-term forecasts are generally accurate, but what about longer-term prediction? What places will wake up to a snowy Christmas this year? And can you tell one month in advance?

# Publicly available datasets

- *https://github.com/awesomedata/awesome-public-datasets*
- *https://medium.com/datadriveninvestor/the-50-best-public-datasets-for-machine-learning-d80e9f030279*
- *https://aws.amazon.com/opendata/public-datasets/*
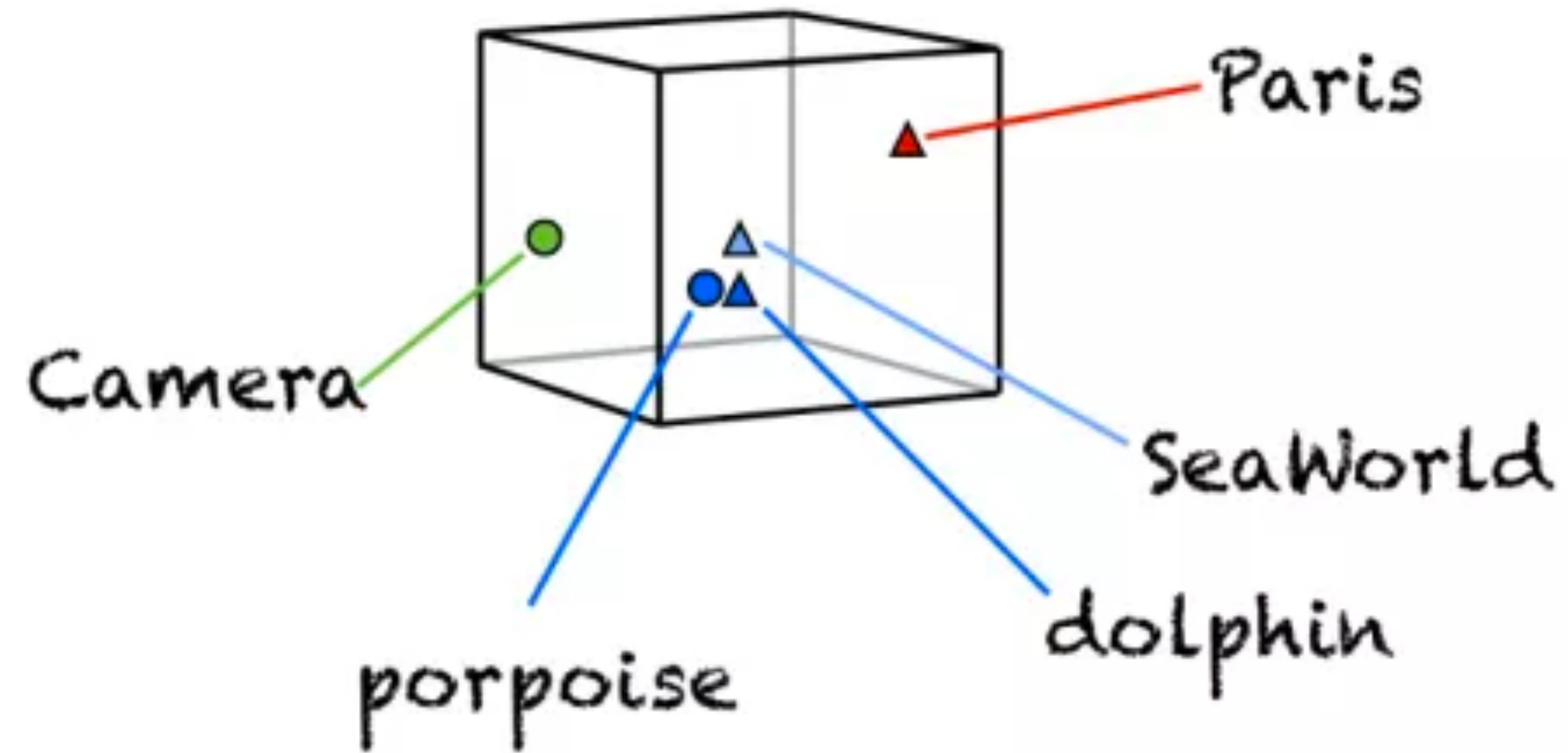- *https://registry.opendata.aws/*
- *https://www.data.gov/*

# Types of data - quantitive vs. categorical

- ***Quantitative data*** *consists of numerical values, like height and weight.*
- ***Categorical data*** *consists of labels describing the properties of the objects under investigation, like gender, hair color, and occupation*
  - *Categorical data doesn't have an order to it*
  - *Does it make any sense to talk about the maximum or minimum hair color? What is the interpretation of my hair color minus your hair color?*
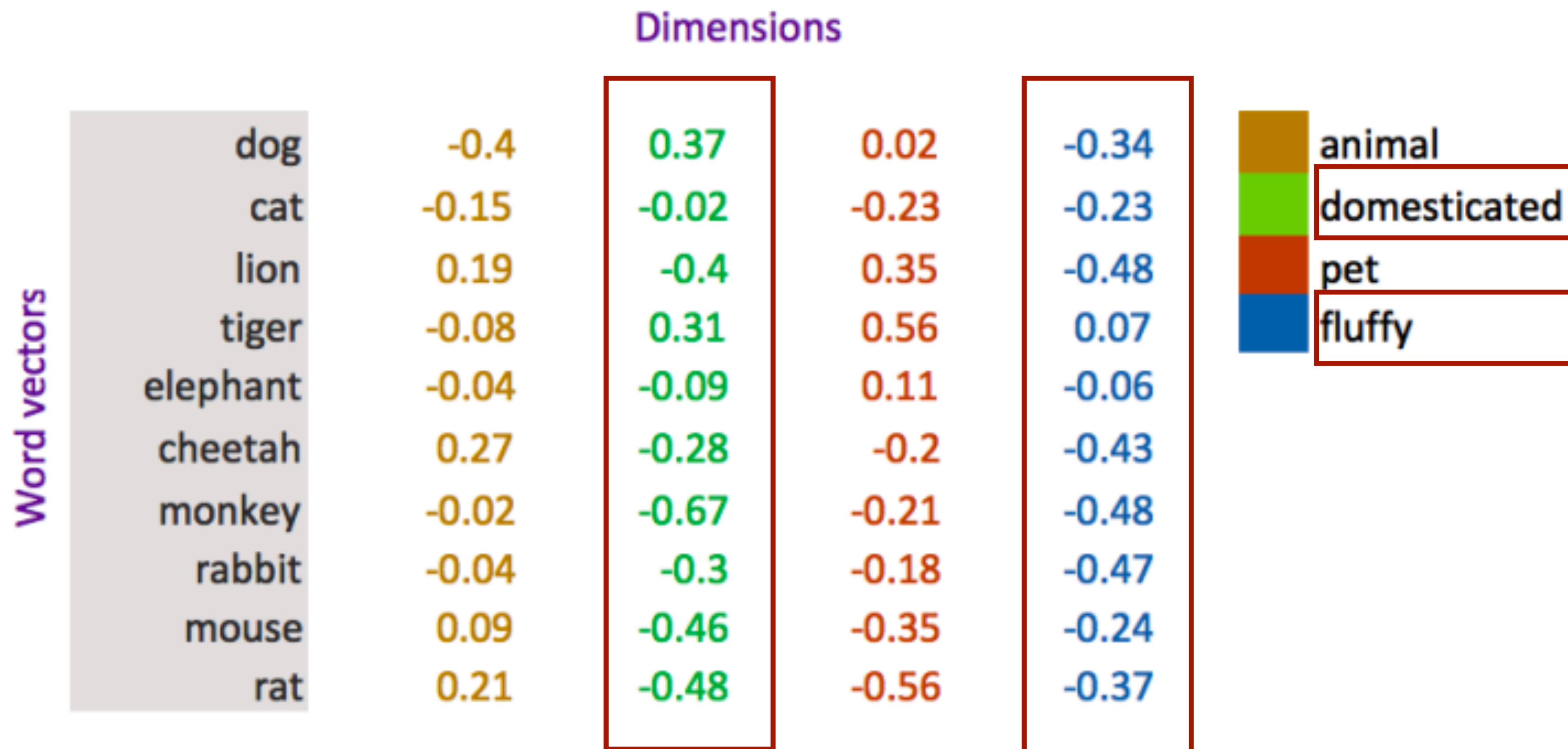
# Quantitive or Categorical?

- *Favorite ice cream flavor*
- *Money spent in marketing budget*
- *GPA*
- *Letter grade*
- *Images*
- *Words*
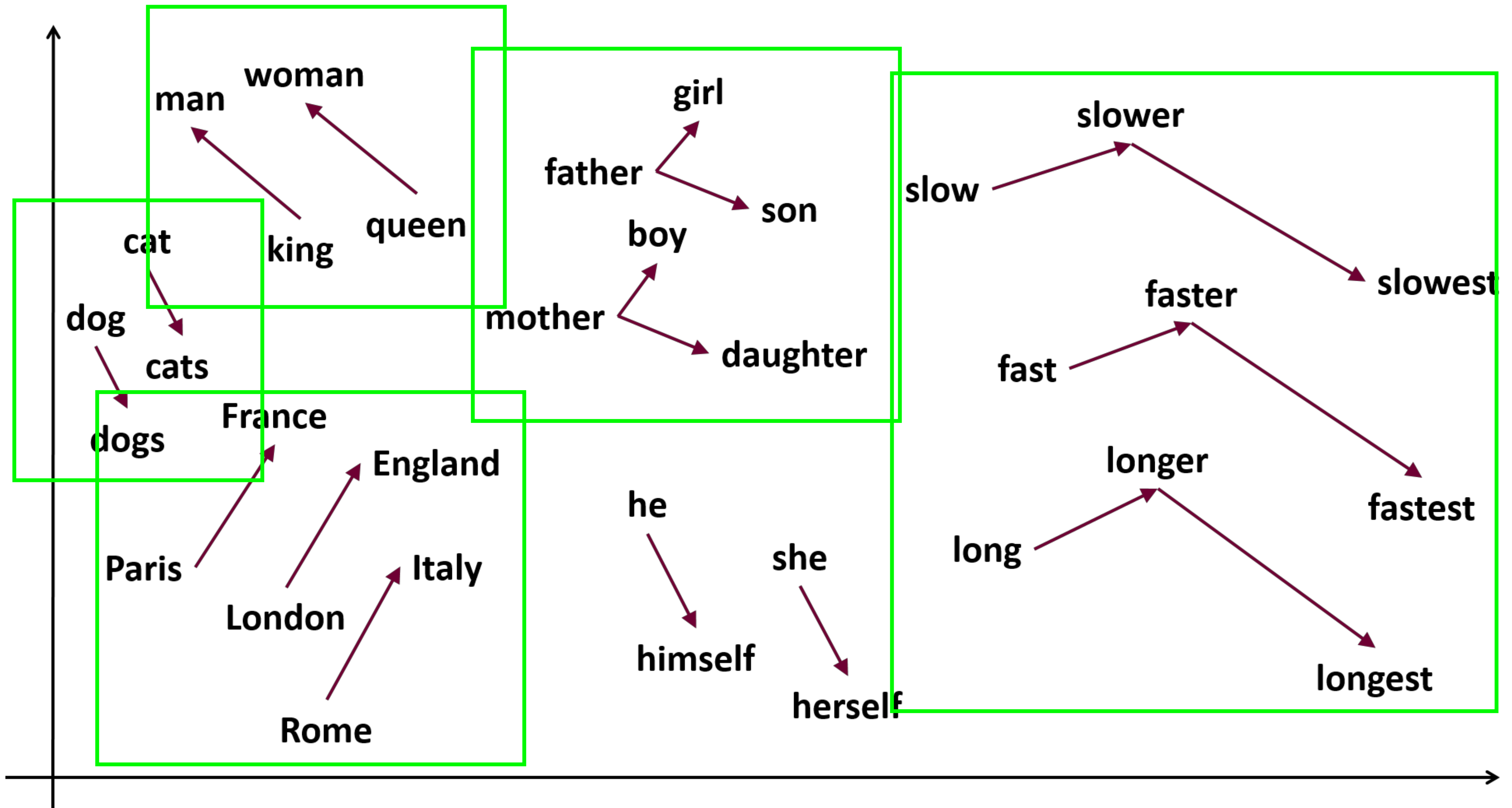
Can we turn word into numbers?

# How can we represents words of animals as vectors?

**Dimensions**

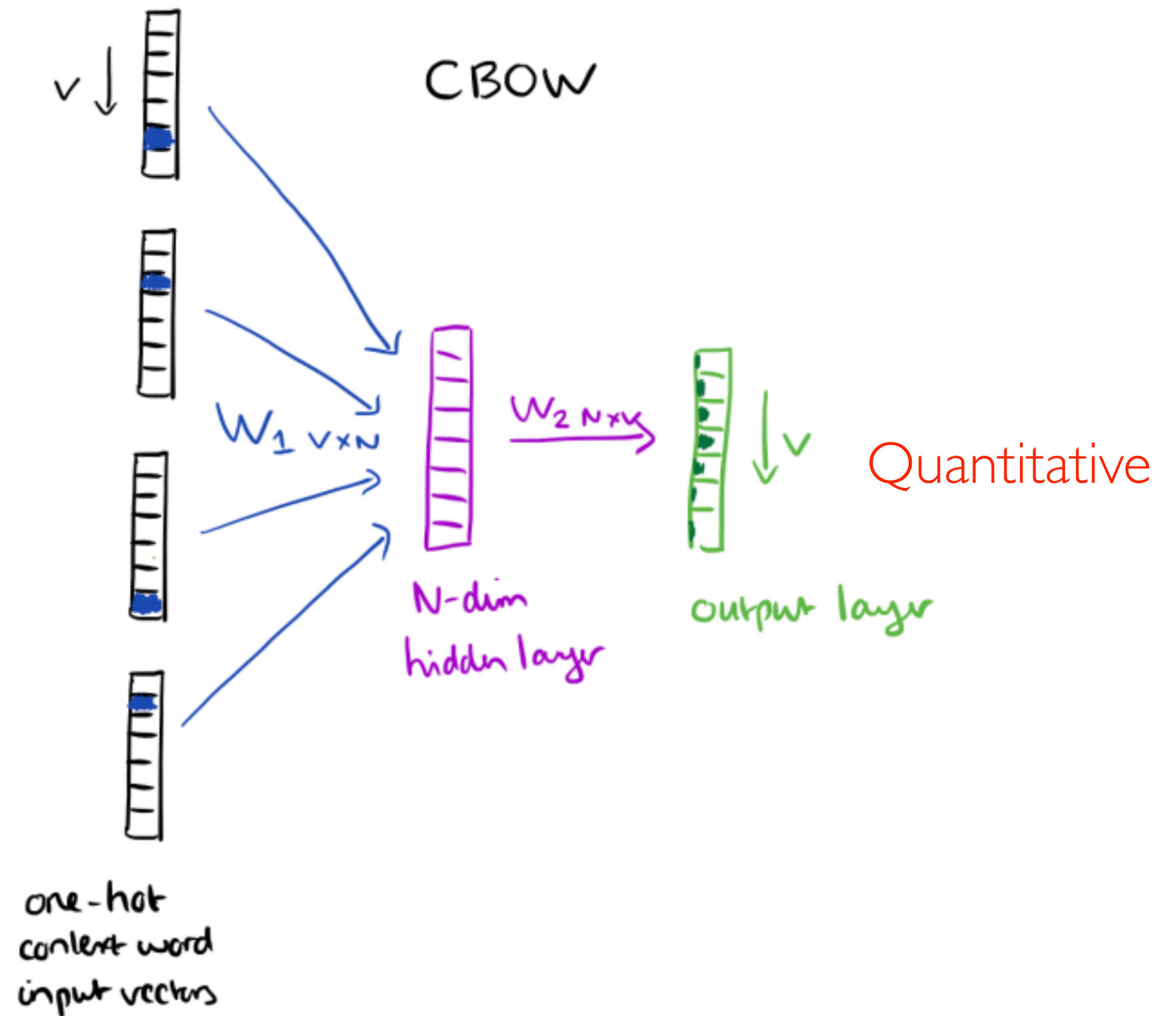| Word vectors | | | | | | |
|---|---|---|---|---|---|---|
| dog | -0.4 | 0.37 | 0.02 | -0.34 | | animal |
| cat | -0.15 | -0.02 | -0.23 | -0.23 | | domesticated |
| lion | 0.19 | -0.4 | 0.35 | -0.48 | | pet |
| tiger | -0.08 | 0.31 | 0.56 | 0.07 | | fluffy |
| elephant | -0.04 | -0.09 | 0.11 | -0.06 | | |
| cheetah | 0.27 | -0.28 | -0.2 | -0.43 | | |
| monkey | -0.02 | -0.67 | -0.21 | -0.48 | | |
| rabbit | -0.04 | -0.3 | -0.18 | -0.47 | | |
| mouse | 0.09 | -0.46 | -0.35 | -0.24 | | |
| rat | 0.21 | -0.48 | -0.56 | -0.37 | | |

# Numerical properties preserved

# Word embedding - automate the process!

**Word embedding** is the collective name for a set of language modeling and feature learning techniques in natural language processing (NLP) where words or phrases from the vocabulary are mapped to vectors of real numbers (Wikipedia).

# Categorical

## One-hot encoding

|         | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---------|---|---|---|---|---|---|---|---|---|
| man     | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| woman   | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| boy     | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| girl    | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| prince  | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| princess| 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| queen   | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| king    | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| monarch | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |



CBOW

$W_1 \; V \times N$

$W_2 \; N \times V$

N-dim hidden layer

output layer

one-hot context word input vectors

Quantitative

# WORD 2 VEC

WINDOW

THE QUICK BROWN FOX JUMPS OVER THE LAZY DOG

CLASSIFIERS

UC San Diego

# Types of data - big vs. little

# Types of data - big vs. little

- *There are difficulties in working with large data sets.*
  - *The analysis cycle time slows as data size grows (slow to iterate)*
  - *Large data sets are complex to visualize*
- *Simple models do not require massive data to fit or evaluate*

# The big vs. little data approach

What are voter preferences about the demographic presidential campaign pool?

# The big vs. little data approach

Which approach do you think will be more accurate?

**Take away:** The right data set is the one most directly relevant to the tasks at hand, not necessarily the biggest one.

UC San Diego

# Types of questions - classification

**Classification:** *Often we seek to assign a label to an item from a discrete set of possibilities. Such problems as predicting the winner of a particular sporting contest (team A or team B?) or deciding the genre of a given movie (comedy, drama, or animation?) are classification problems, since each entail selecting a label from the possible choices.*

# Types of questions - regression

**Regression:** *Another common task is to forecast a given numerical quantity. Predicting a person's weight or how much snow we will get this year is a regression problem, where we forecast the future value of a numerical function in terms of previous values and other relevant features.*

# Regression or classification?

- *Will the price of a particular stock be higher or lower tomorrow?*
- *What will the price of a particular stock be tomorrow?*
- *Is this person a good risk to sell an insurance policy to?*
- *How long do we expect this person to live?*

# Types of data - structured vs. unstructured

- **Structured data** - *data sets that are structured, like the tables in a database or spread- sheet program.*

# CSV and TSV



Lightweight but not good for hierarchical data

# XML-formatted data

```xml
<?xml version="1.0" encoding="UTF-8"?>
<customers>
    <customer>
        <customer_id>1</customer_id>
        <first_name>John</first_name>
        <last_name>Doe</last_name>
        <email>john.doe@example.com</email>
    </customer>
    <customer>
        <customer_id>2</customer_id>
        <first_name>Sam</first_name>
        <last_name>Smith</last_name>
        <email>sam.smith@example.com</email>
    </customer>
    <customer>
        <customer_id>3</customer_id>
        <first_name>Jane</first_name>
        <last_name>Doe</last_name>
        <email>jane.doe@example.com</email>
    </customer>
</customers>
```

Not good for big data

# SQL - Structured Query Language

| | EmployeeId | First Name | Last Name | Department Name |
|---|---|---|---|---|
| 1 | 1 | Ken | Sanchez | Executive |
| 2 | 2 | Terri | Duffy | Engineering |
| 3 | 3 | Roberto | Tamburello | Engineering |
| 4 | 4 | Rob | Walters | Engineering |
| 5 | 5 | Gail | Erickson | Engineering |
| 6 | 6 | Jossef | Goldberg | Engineering |
| 7 | 7 | Dylan | Miller | Support |
| 8 | 8 | Diane | Margheim | Support |
| 9 | 9 | Gigi | Matthew | Support |
| 10 | 10 | Michael | Raheem | Support |

Results    Messages

▼ DATA        RESULTS        SQL          Calculations    Row Limit 500    ☐ Totals

```sql
SELECT
    products.brand    AS "products.brand",
    products.category    AS "products.category",
    COUNT(DISTINCT products.id ) AS "products.count"
FROM public.order_items   AS order_items
LEFT JOIN public.inventory_items   AS inventory_items ON order_items
    .inventory_item_id = inventory_items.id
LEFT JOIN public.products   AS products ON inventory_items.product_id = products.id

GROUP BY 1,2
ORDER BY 3 DESC
LIMIT 500
```

Open in SQL Runner      Explain in SQL Runner

UC San Diego

# JSON

```
{
    "title": "Example Schema",
    "type": "object",                    Sample JSON Schema
    "properties": {
            "firstName": {
                    "type": "string"
            },
            "lastName": {
                    "type": "string"
            },
            "age": {
                    "description": "Age in years",
                    "type": "integer",
                    "minimum": 0
            }
    },
    "required": ["firstName", "lastName"]
}
```

Good for hierarchical data

# API - Application Programming Interface



**YOUR SYSTEMS**

Data

Applications

**API PORTAL**

**Your API Storefront**

**Developer Community**

**Apps**

# Types of data - structured vs. unstructured

**Unstructured data** - *Some datasets record information about the state of the world, but in a more heterogeneous way. Perhaps it is a large text corpus with images and links like Wikipedia, or the complicated mix of notes and test results appearing in personal medical records.*

UC San Diego

# Data Scraping

# Data Scraping

```html
<!DOCTYPE html>
<html>
    <head>
    </head>
    <body>
        <h1> First Scraping </h1>
        <p> Hello World </p>
    <body>
</html>
```

# Getting stock names & prices from Bloomberg



UC San Diego

# Data Scraping

# Data Scraping

```python
import urllib2
from bs4 import BeautifulSoup

# query the website and return the html to the variable 'page'
quote_page = 'http://www.bloomberg.com/quote/SPX:IND'
page = urllib2.urlopen(quote_page)

# parse the html using beautiful soup and store in variable `soup`
soup = BeautifulSoup(page, 'html.parser')
```

# Extracting the stock name



```
# Take out the <div> of name and get its value
name_box = soup.find('h1', attrs={'class': 'name'})
```

# Extracting the stock price



```
# get the index price
price_box = soup.find('div', attrs={'class':'price'})
```

# Christopher Keown, Ph.D.
## UC San Diego

Department of Cognitive Science

ckeown@ucsd.edu