

# Introduction to image analysis and visual structural topic model

Michelle Torres  
*Rice University*

June 5, 2023



# Arizona Daily Star

T tucson.com

Saturday, May 30, 2020

\$2 plus tax



PHOTOS BY JOHN MINCHILLO / THE ASSOCIATED PRESS  
Demonstrators seethe with anger outside a precinct station in Minneapolis that was torched after police abandoned it. Protests over the death of a black man who died in police custody Monday flared up in the city for a fourth straight night Friday.

## Cop who knelt on man's neck arrested, charged with murder

## Assessing dire economic data, Trump lashes out at WHO, China

By Martin Crutsinger  
and Dan Sewell

THE ASSOCIATED PRESS

WASHINGTON — With new U.S. economic numbers highlighting the rough road ahead for a hoped-for rebound, President Trump on Friday took aim at the World Health Organization and China, blaming both for their roles in the pandemic's devastation.

Trump announced that the United States will end its support for WHO, charging it didn't respond adequately to the health crisis because of China's "total control" over the global organization. Trump said Chinese officials "it-

welfare net was showing signs of fraying, as protests erupted for a second day in Spain against layoffs by French carmaker Renault and Italy's chief central banker warned that "uncertainty is rife."

While some U.S. states were moving ahead with steps to reopen businesses and leisure activities needed to spur spending and restore jobs, some were finding relaxed safety measures have been followed by upicks in new cases.

Arkansas over the past week has seen a steady rise in active coronavirus cases, following moves by Gov. Asa Hutchinson to re-

**Nation+World:** Transcripts released of Michael Flynn's calls with Russian diplomat. [AIU](#)

# THE SUN

AN EDITION OF THE REGISTER

Saturday, May 30, 2020

\$2.00

FACEBOOK.COM/SBSUN TWITTER.COM/SBSUN

[sbsun.com](#)

MINNEAPOLIS PROTESTS

## EX-OFFICER FACES MURDER CHARGE



A member of the Minnesota National Guard stands watch in front of the Capitol in St. Paul on Friday. Minnesota Gov. Tim Walz announced that he asked the Minnesota National Guard to be responsible for the safety of the state Capitol.

GLEN STUBBE — STAR TRIBUNE VIA AP

EASING OF PANDEMIC ORDER

Hair salons, restaurants have state's OK to open

L.A. County given the go-ahead to allow patrons to dine on-site

By Ryan Carter  
[rcarter@sfchronicle.com](#)  
[@ryancarter on Twitter](#)

Los Angeles County has the green light from the state to allow further easing of rules in-place restrictions put in place by the coronavirus pandemic and is pushing forward with efforts to reopen a regional economy essentially shut down since March.

The state's approval came Friday after county officials said they've met the criteria to obtain a variance that allows it to move through Stage 2 of the phased-in process that eases restrictions in place since March 19, first on low-risk activities and then to high-risk.

With the state's approval, dine-in restaurants — long left with only curbside sales, takeout and delivery — can now begin allowing customers in with such rules as face masks, major cleaning and social distancing.

[OPEN > PAGE 12](#)

**NORCO '80, PART 12**

**Baldy Notch shootout led to**

# The Gazette

gazette.com

2014 PULITZER PRIZE

NATIONAL REPORTING



SERVING COLORADO SPRINGS & THE PIKES PEAK REGION SINCE 1872

SATURDAY, MAY 30, 2020 \$2.00

## INSIDE

### SPORTS



### Money where your mouth is

Doherty baseball coaches spend some of their stipend investing back into their program. **B1**

### BUSINESS



### InterQuest dining options expand

Parry's, a Denver-based restaurant chain, will open its second Colorado Springs location. **A10**

### HOME & GARDEN



# Protests erupt in U.S.

Officer charged with George Floyd's death as demonstrators nationwide march



PHOTOS BY THE ASSOCIATED PRESS

A protester yells at a member of the Minnesota National Guard on Friday in Minneapolis. Protests continued after the death of George Floyd, who died after being restrained by Minneapolis police officers on Memorial Day.

# Akron Beacon Journal

BEACONJOURNAL.COM

Saturday, May 30, 2020

Informing. Engaging. Essential. | [@beaconjournal](#) | [f facebook.com/AkronBeaconJournal](#) | [\\$2](#)

## Minneapolis smolders



A protester carries a U.S. flag upside down, a sign of distress, next to a burning building Thursday night in Minneapolis. [JULIO CORTEZ/THE ASSOCIATED PRESS]

## Local officials urging peaceful protests

DeWine, area leaders condemn killing of man in Minneapolis police custody

USA TODAY NETWORK Ohio

Gov. Mike DeWine on Friday called the death of George Floyd in the custody of Minneapolis police "horrible," but implored Ohio protesters to assemble peacefully a night after hundreds descended on the Statehouse in Columbus, shattering windows and vandalizing storefronts.

"We must not fight violence with more violence. Peaceful protest and the exercise of First Amendment rights are an important part of our civic



DeWine

## WHY IMAGES?

- Images are **powerful**: extra information + emotional activation + *see to believe* = recall and engagement

## WHY IMAGES?

- Images are **powerful**: extra information + emotional activation + *see to believe* = recall and engagement
- Images are **(kind of)** universal (e.g. compare them to spoken languages)

# WHY IMAGES?

- Images are **powerful**: extra information + emotional activation + *see to believe* = recall and engagement
- Images are **(kind of)** universal (e.g. compare them to spoken languages)
- Visuals are **frames**

# USING IMAGES IN SOCIAL SCIENCES

# USING IMAGES IN SOCIAL SCIENCES

- Studying the effect of presenting information through images
  - Labels vs. Images to signal race/ethnicity (Abrajano, Elmendorf, & Quinn 2018)
  - Visual cues and political knowledge (Prior 2014)

# USING IMAGES IN SOCIAL SCIENCES

- Studying the effect of presenting information through images
  - Labels vs. Images to signal race/ethnicity (Abrajano, Elmendorf, & Quinn 2018)
  - Visual cues and political knowledge (Prior 2014)
- Measurement
  - Election incidents in tweets (Wu and Mebane 2022)
  - Displays of emotion (Boussalis et al. 2021)
  - Rural electrification and service provision (Min 2015)

# USING IMAGES IN SOCIAL SCIENCES

- Studying the effect of presenting information through images
  - Labels vs. Images to signal race/ethnicity (Abrajano, Elmendorf, & Quinn 2018)
  - Visual cues and political knowledge (Prior 2014)
- Measurement
  - Election incidents in tweets (Wu and Mebane 2022)
  - Displays of emotion (Boussalis et al. 2021)
  - Rural electrification and service provision (Min 2015)
- Use images as a vehicle for a complex treatment
  - Masculinity/femininity (Bauer & Carpinella 2018, Bernhard 2023)
  - Police militarization (Mummolo 2018)
  - Level of conflict on attitudes towards protesters (Torres 2022)

# WHY COMPUTERS?

- Process large pools of images

# WHY COMPUTERS?

- Process large pools of images
- Increase consistency/reliability and decrease bias (\*)

# WHY COMPUTERS?

- Process large pools of images
- Increase consistency/reliability and decrease bias (\*)
- Helping humans to “see” and discover (\*)

# WHY COMPUTERS?

- Process large pools of images
- Increase consistency/reliability and decrease bias (\*)
- Helping humans to “see” and discover (\*)
- **Computer vision:** Teaching computers to see

# IMAGE AS DATA: TEACHING THE COMPUTER TO SEE

## IMAGE AS DATA: TEACHING THE COMPUTER TO SEE

A very hard task!

# IMAGE AS DATA: TEACHING THE COMPUTER TO SEE

A very hard task!

- Computers are great at following instructions reliably...



# IMAGE AS DATA: TEACHING THE COMPUTER TO SEE

A very hard task!

- Computers are great at following instructions reliably...
- ... but they are bad at inferences



# IMAGE AS DATA: TEACHING THE COMPUTER TO SEE

A very hard task!

- Computers are great at following instructions reliably...
- ... but they are bad at inferences



# IMAGE AS DATA: TEACHING THE COMPUTER TO SEE

A very hard task!

- Computers are great at following instructions reliably...
- ... but they are bad at inferences



# IMAGE AS DATA: TEACHING THE COMPUTER TO SEE

A very hard task!

- Computers are great at following instructions reliably...
- ... but they are bad at inferences

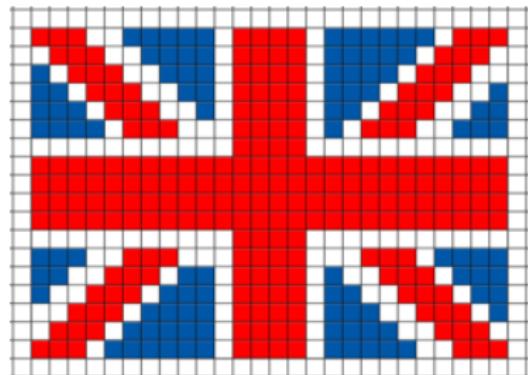


# GETTING READY

- Workshop material available at my GitHub page ([smtorres](#)).
- Google Colab notebook [here](#)
- When doing your own projects:
  - Install Keras ([here](#))
  - Install the following python libraries: numpy, scipy, cv2, matplotlib, PIL, sklearn
    - Check tutorials for OpenCV installation [here](#)
    - I suggest OpenCV 3.X and its compilation from source for full functionality

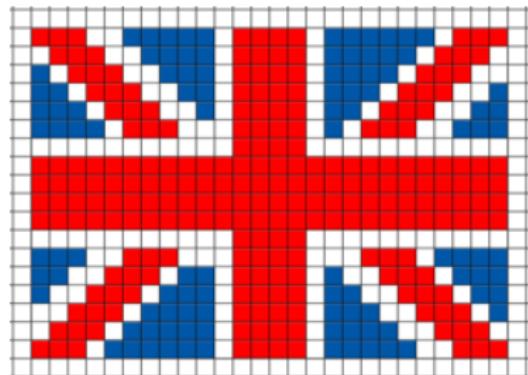
# IMAGE BASICS

- An image is a set of **pixels**:



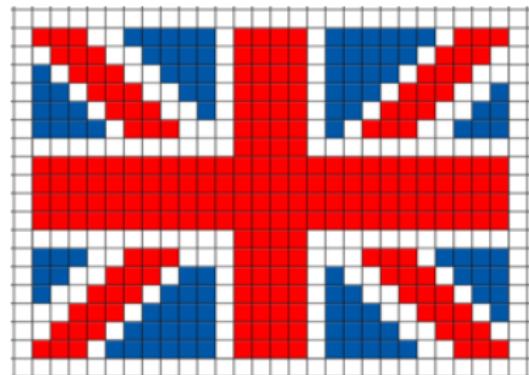
# IMAGE BASICS

- An image is a set of **pixels**:
  - Finest unit (defines **height** and **width**)



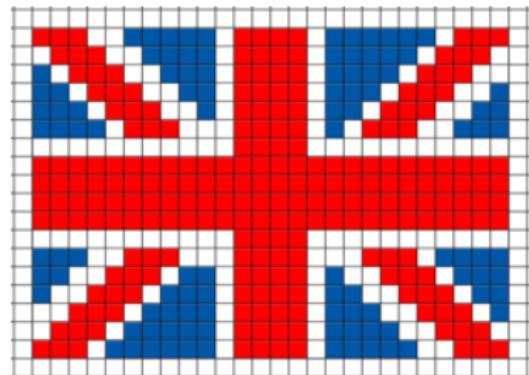
# IMAGE BASICS

- An image is a set of **pixels**:
  - Finest unit (defines **height** and **width**)
  - Grayscale: intensity of light, **Color**: color intensity per channel.



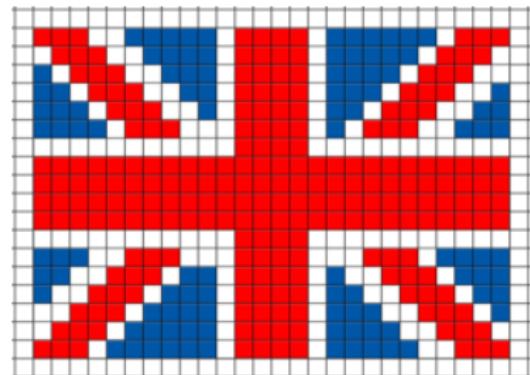
# IMAGE BASICS

- An image is a set of **pixels**:
  - Finest unit (defines **height** and **width**)
  - Grayscale: intensity of light, **Color**: color intensity per channel.
- Matrix representation



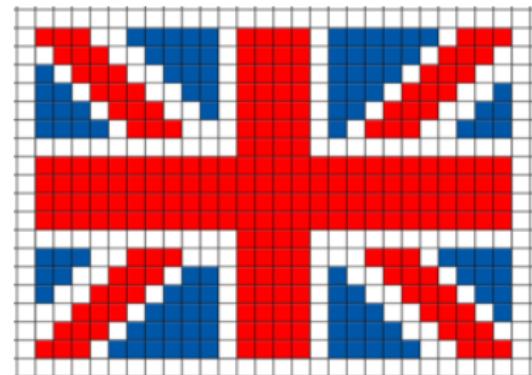
# IMAGE BASICS

- An image is a set of **pixels**:
  - Finest unit (defines **height** and **width**)
  - Grayscale: intensity of light, **Color**: color intensity per channel.
- Matrix representation
  - Grayscale: one matrix



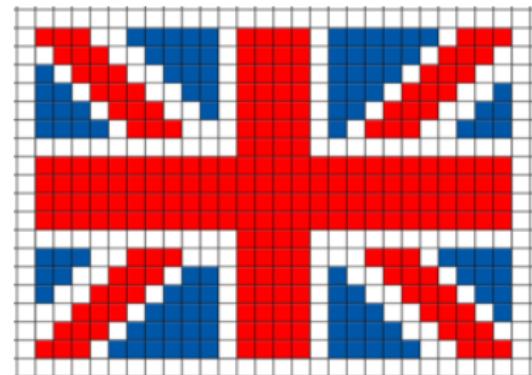
# IMAGE BASICS

- An image is a set of **pixels**:
  - Finest unit (defines **height** and **width**)
  - Grayscale: intensity of light, **Color**: color intensity per channel.
- Matrix representation
  - Grayscale: one matrix
  - Color: array with a matrix for each color channel (**Red**, **Green**, and **Blue**)



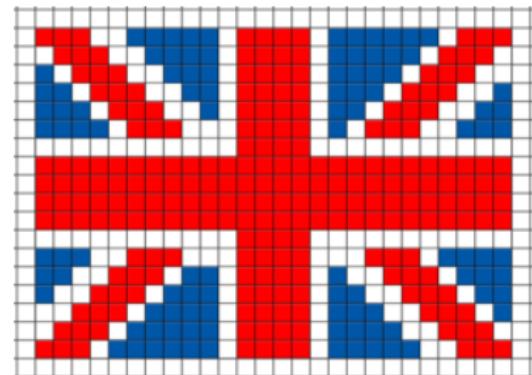
# IMAGE BASICS

- An image is a set of **pixels**:
  - Finest unit (defines **height** and **width**)
  - Grayscale: intensity of light, **Color**: color intensity per channel.
- Matrix representation
  - Grayscale: one matrix
  - Color: array with a matrix for each color channel (**Red**, **Green**, and **Blue**)
- Notice that in OpenCV:



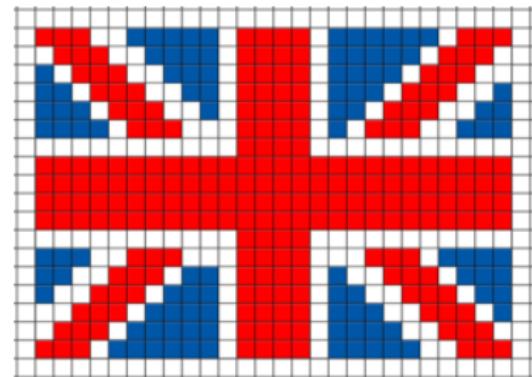
# IMAGE BASICS

- An image is a set of **pixels**:
  - Finest unit (defines **height** and **width**)
  - Grayscale: intensity of light, **Color**: color intensity per channel.
- Matrix representation
  - Grayscale: one matrix
  - Color: array with a matrix for each color channel (**Red**, **Green**, and **Blue**)
- Notice that in OpenCV:
  - Color channel specification is **BRG** instead of **RGB**



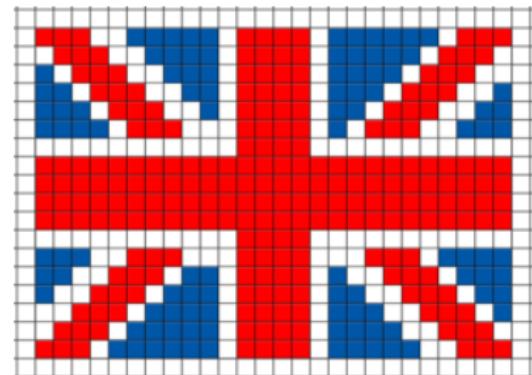
# IMAGE BASICS

- An image is a set of **pixels**:
  - Finest unit (defines **height** and **width**)
  - Grayscale: intensity of light, **Color**: color intensity per channel.
- Matrix representation
  - Grayscale: one matrix
  - Color: array with a matrix for each color channel (**Red**, **Green**, and **Blue**)
- Notice that in OpenCV:
  - Color channel specification is BRG instead of RGB
  - Origin of image is different (top left corner)



# IMAGE BASICS

- An image is a set of **pixels**:
  - Finest unit (defines **height** and **width**)
  - Grayscale: intensity of light, **Color**: color intensity per channel.
- Matrix representation
  - Grayscale: one matrix
  - Color: array with a matrix for each color channel (**Red**, **Green**, and **Blue**)
- Notice that in OpenCV:
  - Color channel specification is BRG instead of RGB
  - Origin of image is different (top left corner)
  - In numpy you specify the *y*-coordinates of an image first:  $x_2 = \text{image}[y_0:y_1, x_0:x_1]$



## DESCRIBING AN IMAGE

- Think about the “tokens” or elements that give meaning to text. What are those?

## DESCRIBING AN IMAGE

- Think about the “tokens” or elements that give meaning to text. What are those?
- The main challenge with images: a lot of pixels that mostly make sense when analyzed in clusters and not as units.

## DESCRIBING AN IMAGE

- Think about the “tokens” or elements that give meaning to text. What are those?
- The main **challenge** with images: a lot of pixels that mostly make sense when analyzed in clusters and not as units.
- Therefore, we use **image descriptors** to characterize the content on an image ***globally*** or **feature descriptors** to locally quantify ***regions*** of the image.

## DESCRIBING AN IMAGE

- Think about the “tokens” or elements that give meaning to text. What are those?
- The main **challenge** with images: a lot of pixels that mostly make sense when analyzed in clusters and not as units.
- Therefore, we use **image descriptors** to characterize the content on an image ***globally*** or **feature descriptors** to locally quantify ***regions*** of the image.
  - Color

## DESCRIBING AN IMAGE

- Think about the “tokens” or elements that give meaning to text. What are those?
- The main **challenge** with images: a lot of pixels that mostly make sense when analyzed in clusters and not as units.
- Therefore, we use **image descriptors** to characterize the content on an image ***globally*** or **feature descriptors** to locally quantify ***regions*** of the image.
  - Color
  - Texture

## DESCRIBING AN IMAGE

- Think about the “tokens” or elements that give meaning to text. What are those?
- The main **challenge** with images: a lot of pixels that mostly make sense when analyzed in clusters and not as units.
- Therefore, we use **image descriptors** to characterize the content on an image ***globally*** or **feature descriptors** to locally quantify ***regions*** of the image.
  - Color
  - Texture
  - Shape

## DESCRIBING AN IMAGE

- Think about the “tokens” or elements that give meaning to text. What are those?
- The main **challenge** with images: a lot of pixels that mostly make sense when analyzed in clusters and not as units.
- Therefore, we use **image descriptors** to characterize the content on an image ***globally*** or **feature descriptors** to locally quantify ***regions*** of the image.
  - Color
  - Texture
  - Shape
  - Pixel intensity change

## DESCRIBING AN IMAGE

- Think about the “tokens” or elements that give meaning to text. What are those?
- The main **challenge** with images: a lot of pixels that mostly make sense when analyzed in clusters and not as units.
- Therefore, we use **image descriptors** to characterize the content on an image ***globally*** or **feature descriptors** to locally quantify ***regions*** of the image.
  - Color
  - Texture
  - Shape
  - Pixel intensity change
  - Edges, objects, etc.

## DESCRIBING AN IMAGE

- Think about the “tokens” or elements that give meaning to text. What are those?
- The main **challenge** with images: a lot of pixels that mostly make sense when analyzed in clusters and not as units.
- Therefore, we use **image descriptors** to characterize the content on an image ***globally*** or **feature descriptors** to locally quantify ***regions*** of the image.
  - Color
  - Texture
  - Shape
  - Pixel intensity change
  - Edges, objects, etc.
- Feature vectors: A series of numbers used to numerically quantify the contents of an image (or regions of it) ⇒ WE USE THEM TO CREATE TOKENS!

## AND THEN WHAT DO WE DO WITH THOSE TOKENS?

- Stop thinking about images for a second
- What do you guys study or research?
- You may have:
  - Dataframes with variables in columns and observations in rows
  - Interviews with text
  - Paragraphs of speeches from MPs
  - ...
- And what do you do with that?
  - Run regressions to make inferences
  - Predict values of interest
  - Identify topics in a corpus of texts
  - ...

# WHAT WE ARE GOING TO BUILD TODAY

# WHAT WE ARE GOING TO BUILD TODAY

Speed bump!

# WHAT WE ARE GOING TO BUILD TODAY

Speed bump!

- ① Follow a Bag of Visuals Word approach in text BUT with images

# WHAT WE ARE GOING TO BUILD TODAY

Speed bump!

- ① Follow a Bag of Visuals Word approach in text BUT with images
- ② Build a Document-Term Matrix with images ⇒ Image-Visual Word Matrix

# WHAT WE ARE GOING TO BUILD TODAY

Speed bump!

- ① Follow a Bag of Visuals Word approach in text BUT with images
- ② Build a Document-Term Matrix with images ⇒ Image-Visual Word Matrix
- ③ And then...?

# WHAT WE ARE GOING TO BUILD TODAY

Speed bump!

- ① Follow a Bag of Visuals Word approach in text BUT with images
- ② Build a Document-Term Matrix with images ⇒ Image-Visual Word Matrix
- ③ And then...?

# WHAT WE ARE GOING TO BUILD TODAY

Speed bump!

- ① Follow a Bag of Visuals Word approach in text BUT with images
- ② Build a Document-Term Matrix with images ⇒ Image-Visual Word Matrix
- ③ And then...? Choose your method!

# WHAT WE ARE GOING TO BUILD TODAY

Speed bump!

- ① Follow a Bag of Visuals Word approach in text BUT with images
- ② Build a Document-Term Matrix with images ⇒ Image-Visual Word Matrix
- ③ And then...? Choose your method!

Sentiment  
Analysis

# WHAT WE ARE GOING TO BUILD TODAY

Speed bump!

- ① Follow a Bag of Visuals Word approach in text BUT with images
- ② Build a Document-Term Matrix with images ⇒ Image-Visual Word Matrix
- ③ And then...? Choose your method!

Sentiment  
Analysis

Binary  
Classification

# WHAT WE ARE GOING TO BUILD TODAY

Speed bump!

- ① Follow a Bag of Visuals Word approach in text BUT with images
- ② Build a Document-Term Matrix with images ⇒ Image-Visual Word Matrix
- ③ And then...? Choose your method!

Sentiment  
Analysis

Binary  
Classification

Clustering  
Analysis

# WHAT WE ARE GOING TO BUILD TODAY

Speed bump!

- ① Follow a Bag of Visuals Word approach in text BUT with images
- ② Build a Document-Term Matrix with images ⇒ Image-Visual Word Matrix
- ③ And then...? Choose your method!

Sentiment  
Analysis

Binary  
Classification

Clustering  
Analysis

Topic  
Modeling

# WHAT WE ARE GOING TO BUILD TODAY

Speed bump!

- ① Follow a Bag of Visuals Word approach in text BUT with images
- ② Build a Document-Term Matrix with images ⇒ Image-Visual Word Matrix
- ③ And then...? Choose your method!

Topic  
Modeling

## TODAY: A “VISUAL” STRUCTURAL TOPIC MODEL

- Structural Topic Model (Roberts et al. 2014)

## TODAY: A “VISUAL” STRUCTURAL TOPIC MODEL

- Structural Topic Model (Roberts et al. 2014)
- Tool for topic modeling of texts with document-level covariate information

## TODAY: A “VISUAL” STRUCTURAL TOPIC MODEL

- Structural Topic Model (Roberts et al. 2014)
- Tool for topic modeling of texts with document-level covariate information
- Mixture model:

## TODAY: A “VISUAL” STRUCTURAL TOPIC MODEL

- Structural Topic Model (Roberts et al. 2014)
- Tool for topic modeling of texts with document-level covariate information
- Mixture model:
  - Probability that words belong to each of the “topics” or groups of interest

## TODAY: A “VISUAL” STRUCTURAL TOPIC MODEL

- Structural Topic Model (Roberts et al. 2014)
- Tool for topic modeling of texts with document-level covariate information
- Mixture model:
  - Probability that words belong to each of the “topics” or groups of interest
  - Not a single classification outcome, but proportions of all potential topics for each document

# TODAY: A “VISUAL” STRUCTURAL TOPIC MODEL

- Structural Topic Model (Roberts et al. 2014)
- Tool for topic modeling of texts with document-level covariate information
- Mixture model:
  - Probability that words belong to each of the “topics” or groups of interest
  - Not a single classification outcome, but proportions of all potential topics for each document



# TODAY: A “VISUAL” STRUCTURAL TOPIC MODEL

- Structural Topic Model (Roberts et al. 2014)
- Tool for topic modeling of texts with document-level covariate information
- Mixture model:
  - Probability that words belong to each of the “topics” or groups of interest
  - Not a single classification outcome, but proportions of all potential topics for each document



Sky:

# TODAY: A “VISUAL” STRUCTURAL TOPIC MODEL

- Structural Topic Model (Roberts et al. 2014)
- Tool for topic modeling of texts with document-level covariate information
- Mixture model:
  - Probability that words belong to each of the “topics” or groups of interest
  - Not a single classification outcome, but proportions of all potential topics for each document



Sky:

Crowd:

# TODAY: A “VISUAL” STRUCTURAL TOPIC MODEL

- Structural Topic Model (Roberts et al. 2014)
- Tool for topic modeling of texts with document-level covariate information
- Mixture model:
  - Probability that words belong to each of the “topics” or groups of interest
  - Not a single classification outcome, but proportions of all potential topics for each document



Sky:

Crowd:

Pavement:

# TODAY: A “VISUAL” STRUCTURAL TOPIC MODEL

- Structural Topic Model (Roberts et al. 2014)
- Tool for topic modeling of texts with document-level covariate information
- Mixture model:
  - Probability that words belong to each of the “topics” or groups of interest
  - Not a single classification outcome, but proportions of all potential topics for each document



Sky:

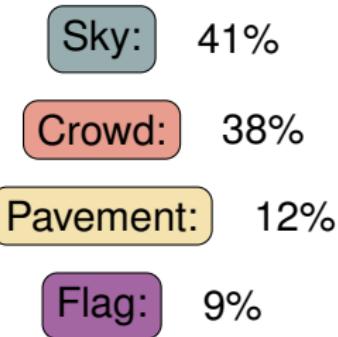
Crowd:

Pavement:

Flag:

# TODAY: A “VISUAL” STRUCTURAL TOPIC MODEL

- Structural Topic Model (Roberts et al. 2014)
- Tool for topic modeling of texts with document-level covariate information
- Mixture model:
  - Probability that words belong to each of the “topics” or groups of interest
  - Not a single classification outcome, but proportions of all potential topics for each document



## BUT FIRST: CONSTRUCTING VISUAL WORDS

## BUT FIRST: CONSTRUCTING VISUAL WORDS

- ① Identification of blocks in images

## BUT FIRST: CONSTRUCTING VISUAL WORDS

- ① Identification of blocks in images
- ② Extraction of features using a CNN

## BUT FIRST: CONSTRUCTING VISUAL WORDS

- ① Identification of blocks in images
- ② Extraction of features using a CNN
- ③ Construction of visual vocabulary based on clustering features

## BUT FIRST: CONSTRUCTING VISUAL WORDS

- ① Identification of blocks in images
- ② Extraction of features using a CNN
- ③ Construction of visual vocabulary based on clustering features
- ④ Construction of Image-Visual Word matrix

## DIVISION OF IMAGES INTO BLOCKS



(a) Original image (resized)

# DIVISION OF IMAGES INTO BLOCKS



(a) Original image (resized)



(b) Image divided into  $32 \times 32$  pixels blocks

# FEATURE EXTRACTION WITH CNNs

## FEATURE EXTRACTION WITH CNNs

- Use a **CNN** to extract features from EACH of the “mini” images composing each of the images in our corpus

## FEATURE EXTRACTION WITH CNNs

- Use a **CNN** to extract features from EACH of the “mini” images composing each of the images in our corpus
- CNN = Convolutional Neural Network

## FEATURE EXTRACTION WITH CNNs

- Use a **CNN** to extract features from EACH of the “mini” images composing each of the images in our corpus
- CNN = Convolutional Neural Network

## FEATURE EXTRACTION WITH CNNs

- Use a **CNN** to extract features from EACH of the “mini” images composing each of the images in our corpus
- CNN = Convolutional Neural Network

**SHORT PAUSE:** Brief crash course on  
Convolutional Neural Networks

## FEATURE EXTRACTION WITH CNNs

- Use a **CNN** to extract features from EACH of the “mini” images composing each of the images in our corpus
- CNN = Convolutional Neural Network

**SHORT PAUSE:** Brief crash course on  
Convolutional Neural Networks

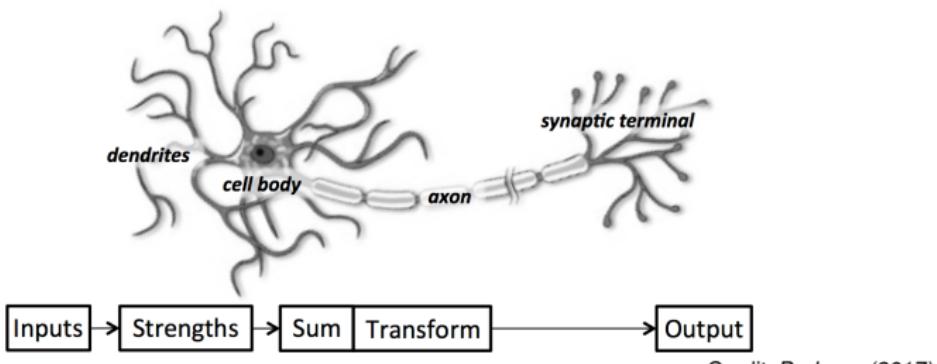
(But please hold the “mini” images thought!)

## SEEING LIKE A HUMAN

- Modern computer vision systems are meant to emulate how human brains transform sensual stimuli into conceptual understanding
- The process allows computers to set their own set of rules to classify information

# SEEING LIKE A HUMAN

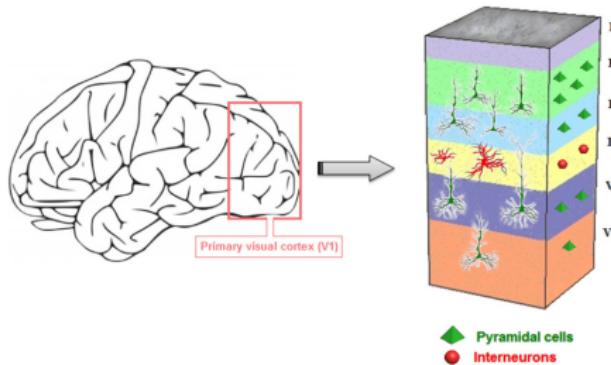
- Modern computer vision systems are meant to emulate how human brains transform sensual stimuli into conceptual understanding
- The process allows computers to set their own set of rules to classify information



Credit: Buduma (2017)

- Neurons are the core unit of brains.
- They receive information from other neurons, process the input, and send the result to other cells for further processing.

## SEEING LIKE A HUMAN, CONT.



Credit: Bachatene, Bharmauria and Molotchnikoff (2012).

- Neurons are organized into layers.
- Every layer breaks down the signal into small pieces, allowing each of its neurons to focus on a unique piece of information.
- The first layers identify basic visual patterns, intermediate layers transform patterns into shapes, and the last layers convert shapes into objects.

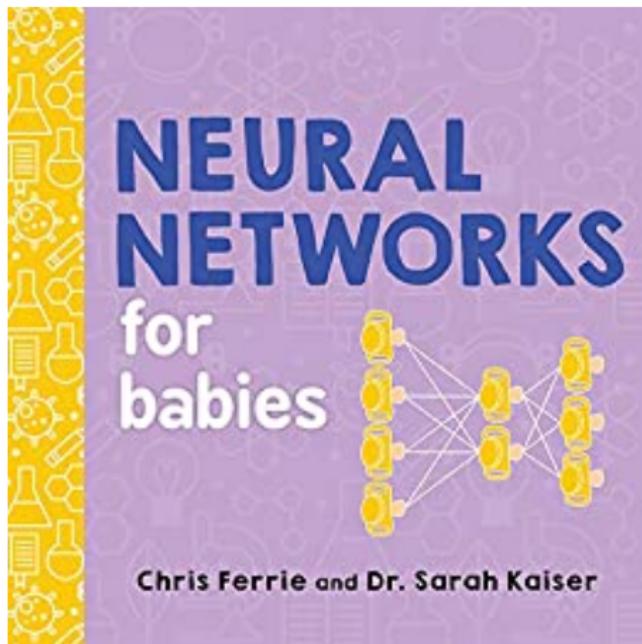
**WAIT... WHAT?**

## WAIT... WHAT?

A very sophisticated text that I've been reading a lot recently:

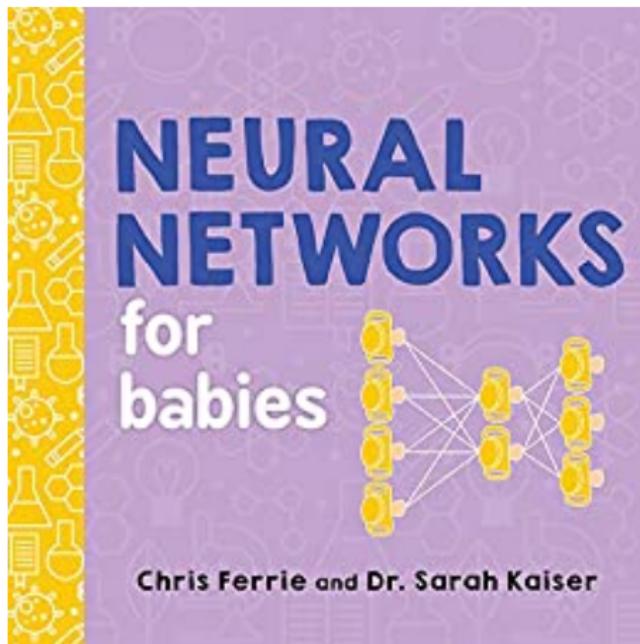
# WAIT... WHAT?

A very sophisticated text that I've been reading a lot recently:



## WAIT... WHAT?

A very sophisticated text that I've been reading a lot recently:

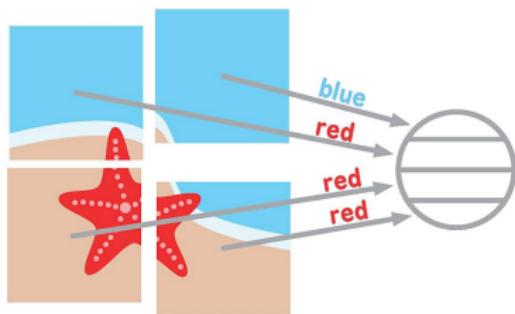


(No... I am not joking)

# THE LOGIC OF CNNs

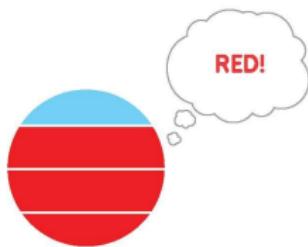


Is there a red animal in this picture?

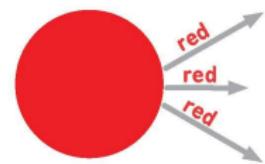


The neuron can decide based on its input.

# THE LOGIC OF CNNs



When the neuron has an answer,

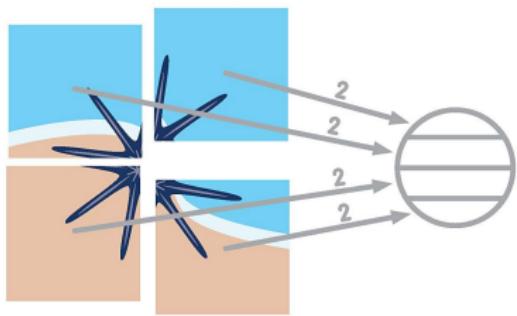


it sends its own message.

# THE LOGIC OF CNNs

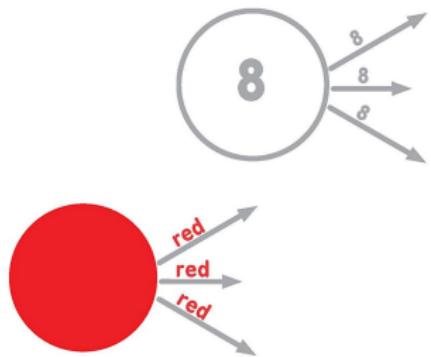


Does this animal have 8 arms?

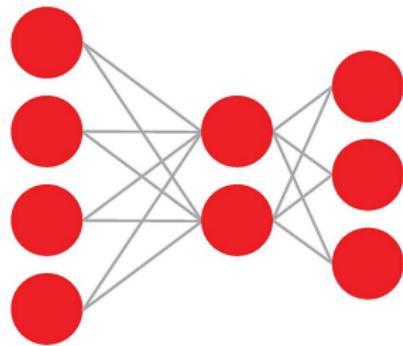


The neuron can decide based on its input.

# THE LOGIC OF CNNs



Where do the messages go?



Neurons talk to each other.  
They connect together in a network.

---

**ALMOST LIKE FINDING WALDO...**

## ALMOST LIKE FINDING WALDO...

- Ok, not that cartoonish but almost!



## ALMOST LIKE FINDING WALDO...

- Ok, not that cartoonish but almost!
- What is your approach when you want to find Waldo?



## ALMOST LIKE FINDING WALDO...

- Ok, not that cartoonish but almost!
- What is your approach when you want to find Waldo?
- Scan the image looking for particular “features”



## ALMOST LIKE FINDING WALDO...

- Ok, not that cartoonish but almost!
- What is your approach when you want to find Waldo?
- Scan the image looking for particular “features”
  - Red and white stripes



## ALMOST LIKE FINDING WALDO...

- Ok, not that cartoonish but almost!
- What is your approach when you want to find Waldo?
- Scan the image looking for particular “features”
  - Red and white stripes
  - Glasses



## ALMOST LIKE FINDING WALDO...

- Ok, not that cartoonish but almost!
- What is your approach when you want to find Waldo?
- Scan the image looking for particular “features”
  - Red and white stripes
  - Glasses
  - Hat



# ALMOST LIKE FINDING WALDO...

- Ok, not that cartoonish but almost!
- What is your approach when you want to find Waldo?
- Scan the image looking for particular “features”
  - Red and white stripes
  - Glasses
  - Hat
- There is a robot who finds him in less than 5 seconds



# FOR REAL

And it's based on CNN code (see [here](#))



# CONVOLUTIONAL NEURAL NETWORKS (CNNs)

- Convolutional Neural Networks (CNNs) is a supervised learning algorithm to classify images

## CONVOLUTIONAL NEURAL NETWORKS (CNNs)

- Convolutional Neural Networks (CNNs) is a supervised learning algorithm to classify images
- CNNs gradually learn what visual features of the image are more important in a classification task by transforming the image into multiple representations or *feature maps*.

# CONVOLUTIONAL NEURAL NETWORKS (CNNs)

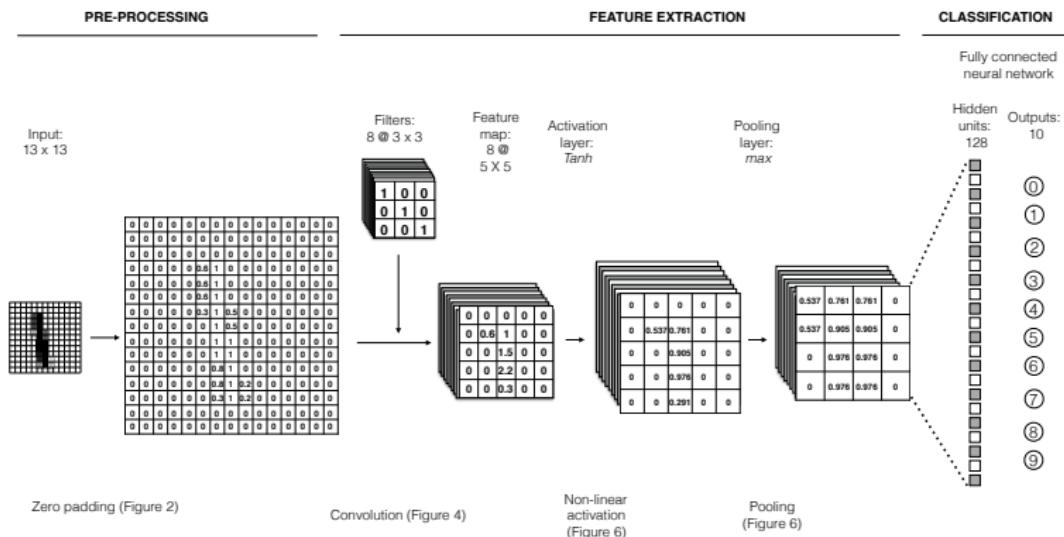
- Convolutional Neural Networks (CNNs) is a supervised learning algorithm to classify images
- CNNs gradually learn what visual features of the image are more important in a classification task by transforming the image into multiple representations or *feature maps*.
- CNNs are organized into multiple layers. Each layer contains multiple representations of the original image through maps of visual features such as edges, blobs or color combinations.

# CONVOLUTIONAL NEURAL NETWORKS (CNNs)

- Convolutional Neural Networks (CNNs) is a supervised learning algorithm to classify images
- CNNs gradually learn what visual features of the image are more important in a classification task by transforming the image into multiple representations or *feature maps*.
- CNNs are organized into multiple layers. Each layer contains multiple representations of the original image through maps of visual features such as edges, blobs or color combinations.
- The part of learning and reaching a semantic concept that humans conduct by trial and error is achieved through the training, validation and testing procedures in CNNs.

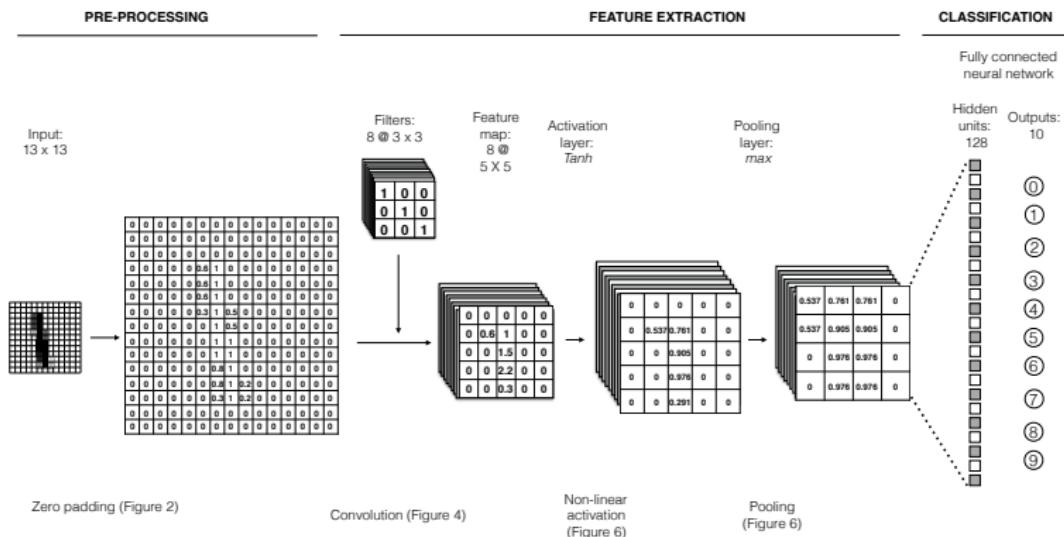
# NETWORK STRUCTURE

- **GOAL:** learn the features associated w/ outcomes



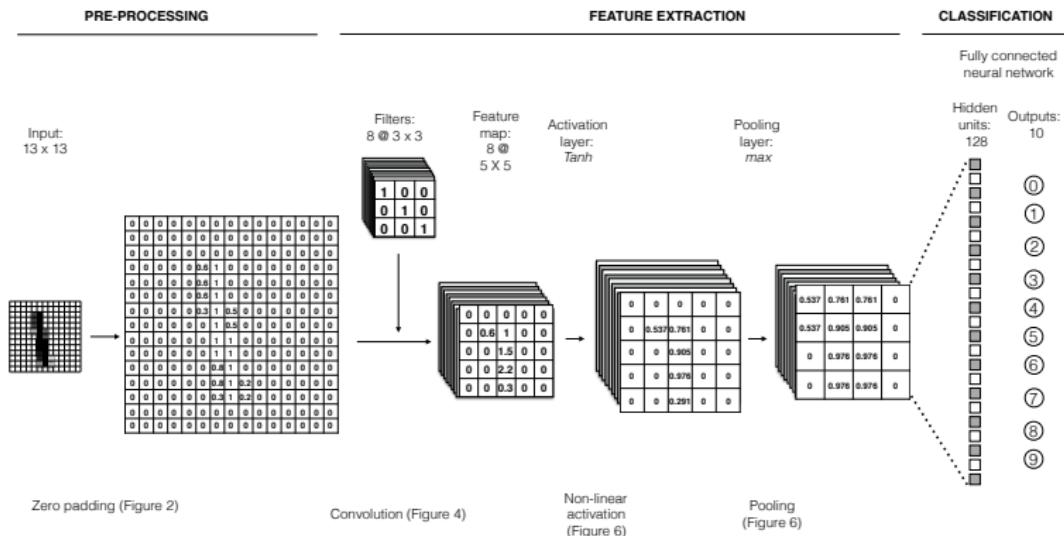
# NETWORK STRUCTURE

- **GOAL:** learn the features associated w/ outcomes
- Translation: obtain “coefficients” [weights in feature maps]



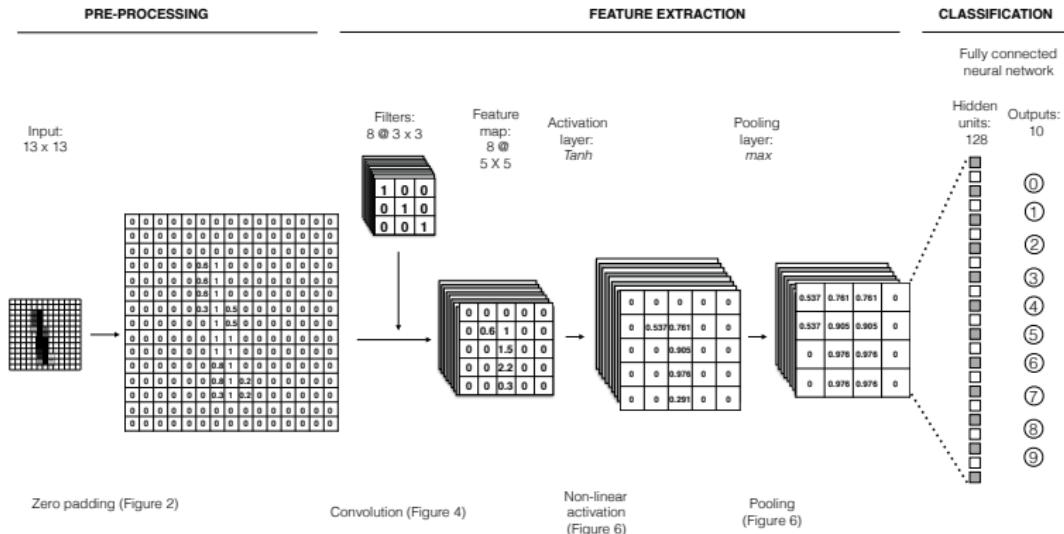
# NETWORK STRUCTURE

- **GOAL:** learn the features associated w/ outcomes
- Translation: obtain “coefficients” [weights in feature maps]
- Mainly, a data reduction technique → **Why?**

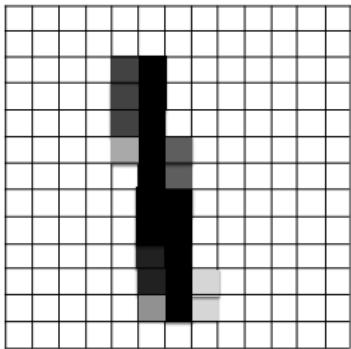


# NETWORK STRUCTURE

- **GOAL:** learn the features associated w/ outcomes
- Translation: obtain “coefficients” [weights in feature maps]
- Mainly, a data reduction technique → **Why?**
- Not a black-box! → Optimization of error



# REPRESENTING IMAGES

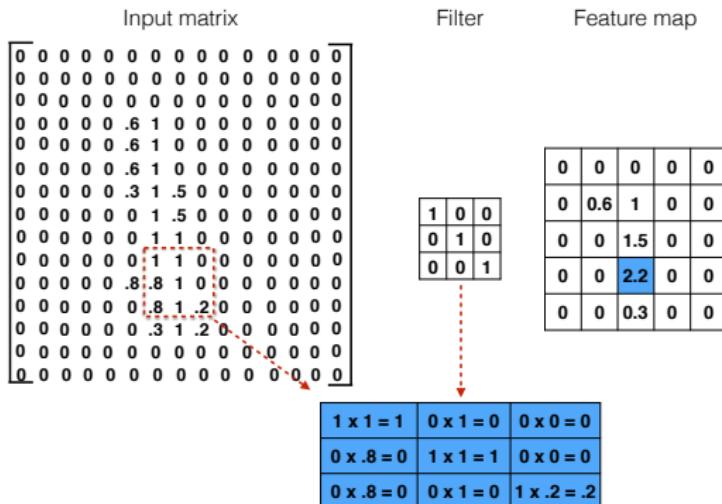


0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	.6	1	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	.6	1	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	.6	1	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	.3	1	.5	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	1	.5	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	1	1	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	1	1	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	1	.8	1	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	.8	1	.2	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	3	1	.2	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

The image is transformed into a numerical matrix, where each element represents the value of a specific pixel of the image measured as light intensity (in grayscale images) or color intensity (in color images).

## FEATURE EXTRACTION

## It's all about feature extraction!



Filters are matrixes made of *weights*, that maximize or minimize the “intensity” of a pixel. Every filter slides through each  $3 \times 3$  pixel area of the image, and computes the dot product of the region. The result is recorded on a smaller matrix to create *feature maps*. Intuitively, we want to detect whether and where a feature represented by a filter is prominent in the image.

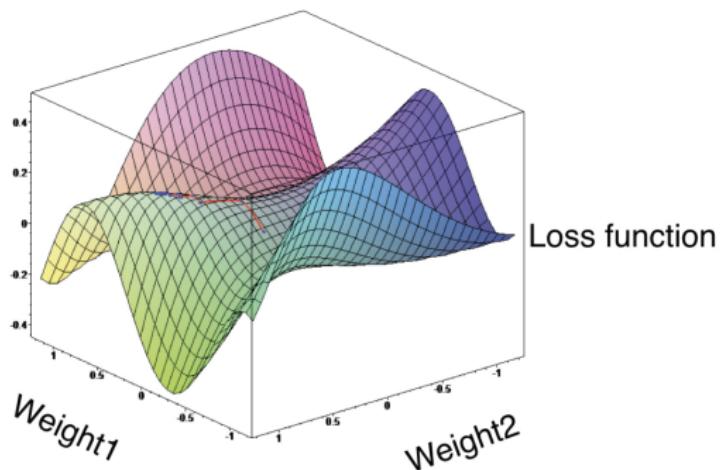
# LEARNING

- The last stage of the network involves the classification of the image. The way in which the CNN learns the features that correlate to each outcome follows a procedure called back-propagation.

[More on back-propagation](#)

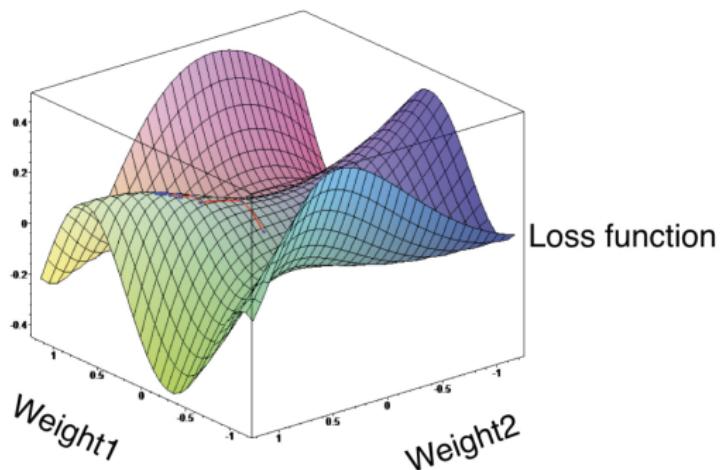
ACTUALLY, THIS SHOULD BE FAMILIAR...

## Loss function



# ACTUALLY, THIS SHOULD BE FAMILIAR...

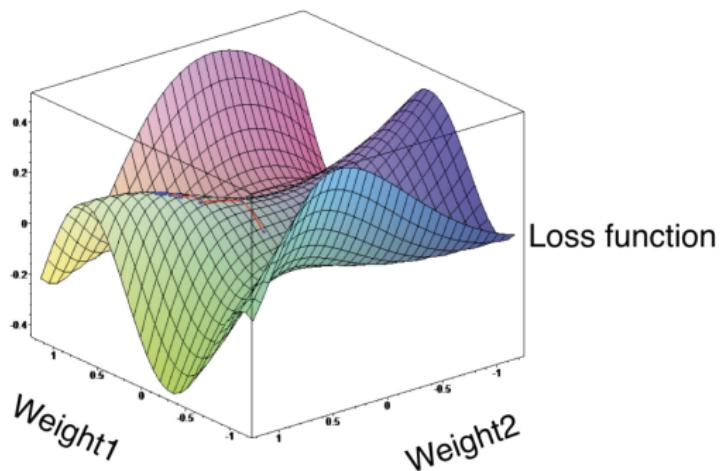
## Loss function



- Minimize multidimensional loss function →

ACTUALLY, THIS SHOULD BE FAMILIAR...

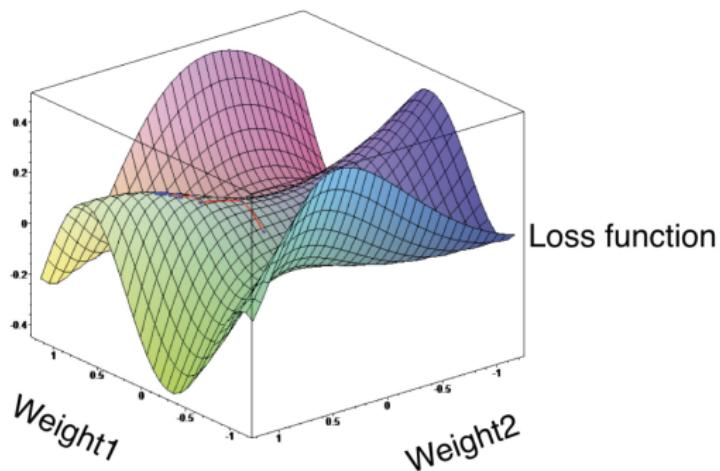
## Loss function



- Minimize multidimensional loss function →

ACTUALLY, THIS SHOULD BE FAMILIAR...

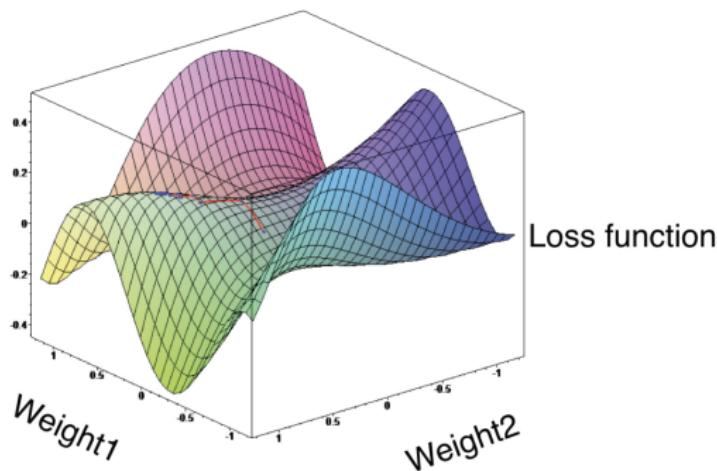
## Loss function



- Minimize multidimensional loss function → (OLS anyone?)
- By finding the minimum point [=minimum prediction error]

# ACTUALLY, THIS SHOULD BE FAMILIAR...

## Loss function



- Minimize multidimensional loss function → (OLS anyone?)
- By finding the minimum point [=minimum prediction error]
- Explore the “field” step by step

# FEATURE EXTRACTION USING CNNs

## FEATURE EXTRACTION USING CNNs

- Use pre-trained model on each block of an image

## FEATURE EXTRACTION USING CNNs

- Use pre-trained model on each block of an image
- The CNN creates feature maps of “elements/descriptors” that can be found in an image

## FEATURE EXTRACTION USING CNNs

- Use pre-trained model on each block of an image
- The CNN creates feature maps of “elements/descriptors” that can be found in an image
- Remove the dense layer (the final one) and keep an appropriate feature map vector → “Predictors”

## FEATURE EXTRACTION USING CNNs

- Use pre-trained model on each block of an image
- The CNN creates feature maps of “elements/descriptors” that can be found in an image
- Remove the dense layer (the final one) and keep an appropriate feature map vector → “Predictors”
- = Each image is described by vector of size *number of blocks* × *number of features from CNN*

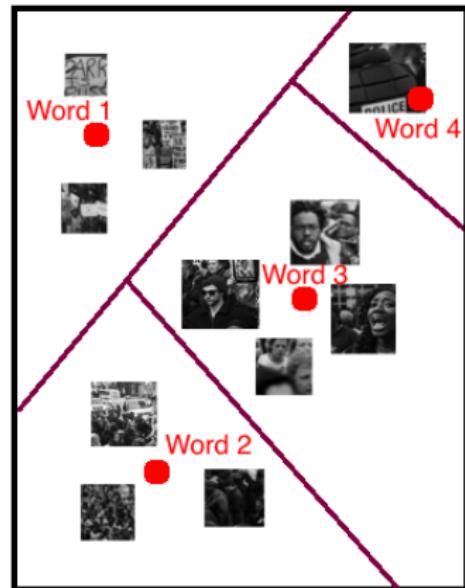
## FEATURE EXTRACTION USING CNNs

- Use pre-trained model on each block of an image
- The CNN creates feature maps of “elements/descriptors” that can be found in an image
- Remove the dense layer (the final one) and keep an appropriate feature map vector → “Predictors”
- = Each image is described by vector of size *number of blocks* × *number of features from CNN*
- In our applications, this is  $70 \times 2,048$

# CLUSTERING FEATURES TO BUILD VISUAL VOCABULARY

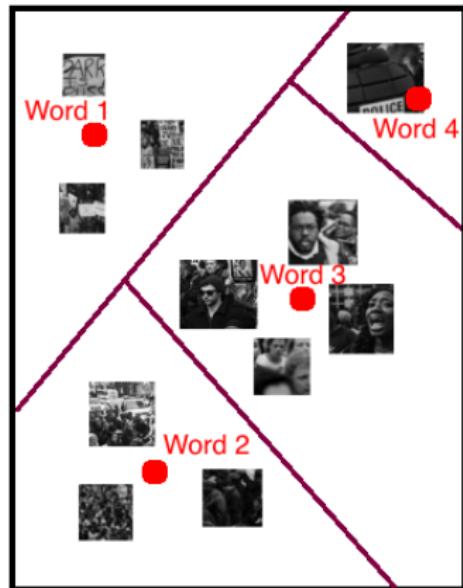
# CLUSTERING FEATURES TO BUILD VISUAL VOCABULARY

- Need for tokens → Words in columns of a DTM



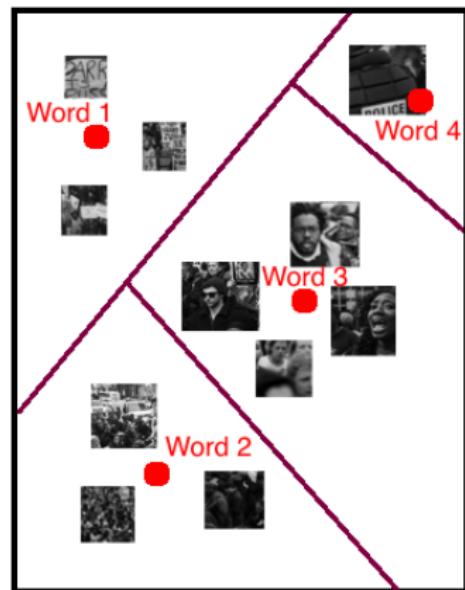
# CLUSTERING FEATURES TO BUILD VISUAL VOCABULARY

- Need for tokens → Words in columns of a DTM
- Define  $v$  clusters (= # of desired visual words)



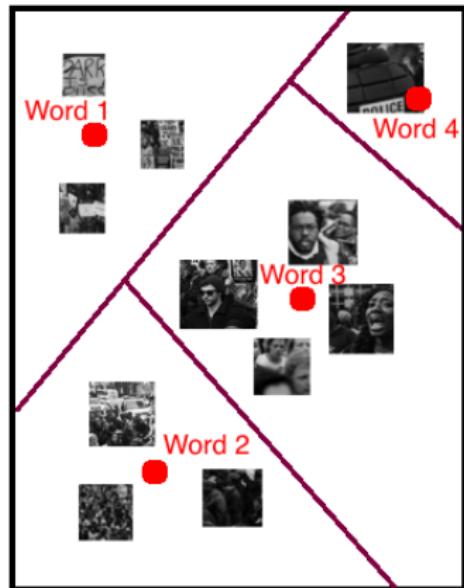
# CLUSTERING FEATURES TO BUILD VISUAL VOCABULARY

- Need for tokens → Words in columns of a DTM
- Define  $v$  clusters (= # of desired visual words)
- Cluster randomly selected sample of feature vectors



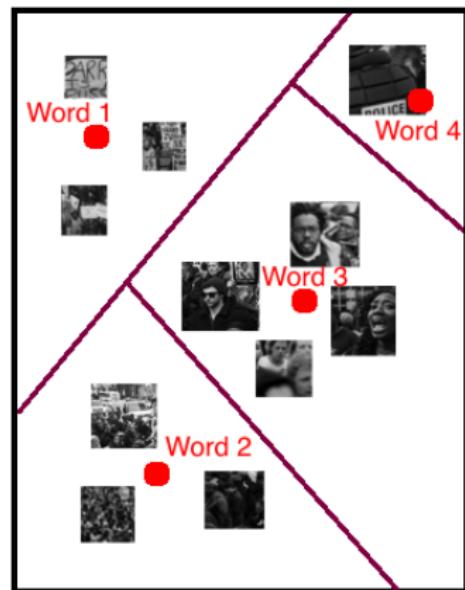
# CLUSTERING FEATURES TO BUILD VISUAL VOCABULARY

- Need for tokens → Words in columns of a DTM
- Define  $v$  clusters (= # of desired visual words)
- Cluster randomly selected sample of feature vectors
- Centroid of cluster is the “visual word”



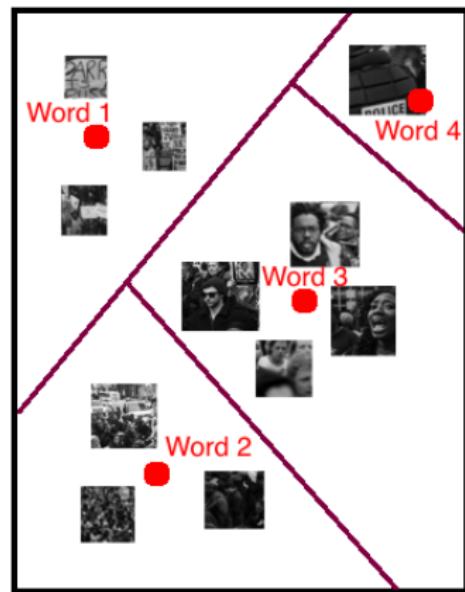
# CLUSTERING FEATURES TO BUILD VISUAL VOCABULARY

- Need for tokens → Words in columns of a DTM
- Define  $v$  clusters (= # of desired visual words)
- Cluster randomly selected sample of feature vectors
- Centroid of cluster is the “visual word”
- Why do we do this?



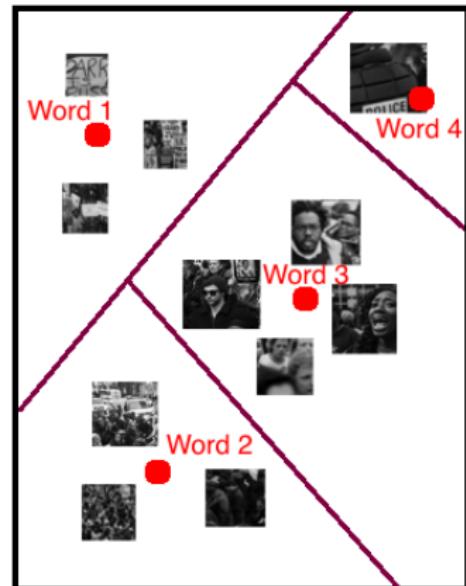
# CLUSTERING FEATURES TO BUILD VISUAL VOCABULARY

- Need for tokens → Words in columns of a DTM
- Define  $v$  clusters (= # of desired visual words)
- Cluster randomly selected sample of feature vectors
- Centroid of cluster is the “visual word”
- Why do we do this?
  - Similar features = Same concept



# CLUSTERING FEATURES TO BUILD VISUAL VOCABULARY

- Need for tokens → Words in columns of a DTM
- Define  $v$  clusters (= # of desired visual words)
- Cluster randomly selected sample of feature vectors
- Centroid of cluster is the “visual word”
- Why do we do this?
  - Similar features = Same concept
  - Reduce potential sparsity in IVWM



## VISUALIZING VISUAL WORDS

- Blocks that belong to a given cluster are similar in terms of feature vectors

## VISUALIZING VISUAL WORDS

- Blocks that belong to a given cluster are similar in terms of feature vectors
- Should look visually similar

## VISUALIZING VISUAL WORDS

- Blocks that belong to a given cluster are similar in terms of feature vectors
- Should look visually similar
- Construct “visual words” using the 16 feature vectors closest to each of the centroid of the cluster

## VISUALIZING VISUAL WORDS

- Blocks that belong to a given cluster are similar in terms of feature vectors
- Should look visually similar
- Construct “visual words” using the 16 feature vectors closest to each of the centroid of the cluster
- E.g. the most similar blocks to the “average” block representing the cluster

# VISUALIZING VISUAL WORDS

- Blocks that belong to a given cluster are similar in terms of feature vectors
- Should look visually similar
- Construct “visual words” using the 16 feature vectors closest to each of the centroid of the cluster
- E.g. the most similar blocks to the “average” block representing the cluster



# VISUALIZING VISUAL WORDS

- Blocks that belong to a given cluster are similar in terms of feature vectors
- Should look visually similar
- Construct “visual words” using the 16 feature vectors closest to each of the centroid of the cluster
- E.g. the most similar blocks to the “average” block representing the cluster



# VISUALIZING VISUAL WORDS

- Blocks that belong to a given cluster are similar in terms of feature vectors
- Should look visually similar
- Construct “visual words” using the 16 feature vectors closest to each of the centroid of the cluster
- E.g. the most similar blocks to the “average” block representing the cluster



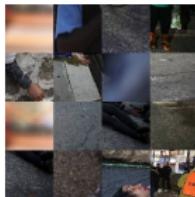
# VISUALIZING VISUAL WORDS

- Blocks that belong to a given cluster are similar in terms of feature vectors
- Should look visually similar
- Construct “visual words” using the 16 feature vectors closest to each of the centroid of the cluster
- E.g. the most similar blocks to the “average” block representing the cluster



# VISUALIZING VISUAL WORDS

- Blocks that belong to a given cluster are similar in terms of feature vectors
- Should look visually similar
- Construct “visual words” using the 16 feature vectors closest to each of the centroid of the cluster
- E.g. the most similar blocks to the “average” block representing the cluster



# **BUILDING THE IVWM TO EMULATE DTM**

## BUILDING THE IVWM TO EMULATE DTM

Count the number of times each visual word appears in an image

## BUILDING THE IVWM TO EMULATE DTM

Count the number of times each visual word appears in an image

- Also not trivial...

## BUILDING THE IVWM TO EMULATE DTM

Count the number of times each visual word appears in an image

- Also not trivial...
- Assign each feature vector to the most similar visual word in the vocabulary

## BUILDING THE IVWM TO EMULATE DTM

Count the number of times each visual word appears in an image

- Also not trivial...
- Assign each feature vector to the most similar visual word in the vocabulary
  - Compute the Euclidean distance between each feature vector and the centroids of the clusters

## BUILDING THE IVWM TO EMULATE DTM

Count the number of times each visual word appears in an image

- Also not trivial...
- Assign each feature vector to the most similar visual word in the vocabulary
  - Compute the Euclidean distance between each feature vector and the centroids of the clusters
  - Assign feature vector to visual word with shortest distance to centroid

## **ILLUSTRATING EXAMPLE: MIGRANT CARAVAN**

## ILLUSTRATING EXAMPLE: MIGRANT CARAVAN

- Groups of migrants from Central America fleeing violence in their countries and seeking refugee in the U.S.

## ILLUSTRATING EXAMPLE: MIGRANT CARAVAN

- Groups of migrants from Central America fleeing violence in their countries and seeking refugee in the U.S.
- Very polarized coverage of this phenomenon

## ILLUSTRATING EXAMPLE: MIGRANT CARAVAN

- Groups of migrants from Central America fleeing violence in their countries and seeking refugee in the U.S.
- Very polarized coverage of this phenomenon
- Emphasis on **magnitude**: threat, invasion

## ILLUSTRATING EXAMPLE: MIGRANT CARAVAN

- Groups of migrants from Central America fleeing violence in their countries and seeking refugee in the U.S.
- Very polarized coverage of this phenomenon
- Emphasis on **magnitude**: threat, invasion

“**Massive** migrant caravan on the way”

“Looks more like an **invasion** than anything”



## ILLUSTRATING EXAMPLE: MIGRANT CARAVAN

- Groups of migrants from Central America fleeing violence in their countries and seeking refugee in the U.S.
- Very polarized coverage of this phenomenon
- Emphasis on **magnitude**: threat, invasion

“**Massive** migrant caravan on the way”

“Looks more like an **invasion** than anything”



“See them as they are: **Desperate**, leaving behind whatever they had, and whomever they knew, all for a **better chance** at life”

# IDENTIFYING POLITICAL COMPONENTS IN THE DATA GENERATION PROCESS OF IMAGES

- Goal: Identify and quantify the visual framing of the magnitude of the caravan

## IDENTIFYING POLITICAL COMPONENTS IN THE DATA GENERATION PROCESS OF IMAGES

- Goal: Identify and quantify the visual framing of the magnitude of the caravan
- Structural Topic Model to identify underlying “topics”, understood as frames, in the images

# IDENTIFYING POLITICAL COMPONENTS IN THE DATA GENERATION PROCESS OF IMAGES

- **Goal:** Identify and quantify the visual framing of the **magnitude** of the caravan
- Structural Topic Model to identify underlying “topics”, understood as frames, in the images
- Visual codebook generated from  $\approx$  6,000 pictures from *Getty*

# IDENTIFYING POLITICAL COMPONENTS IN THE DATA GENERATION PROCESS OF IMAGES

- Goal: Identify and quantify the visual framing of the magnitude of the caravan
- Structural Topic Model to identify underlying “topics”, understood as frames, in the images
- Visual codebook generated from  $\approx$  6,000 pictures from *Getty*
- 500 words vocabulary

# IDENTIFYING POLITICAL COMPONENTS IN THE DATA GENERATION PROCESS OF IMAGES

- Goal: Identify and quantify the visual framing of the magnitude of the caravan
- Structural Topic Model to identify underlying “topics”, understood as frames, in the images
- Visual codebook generated from  $\approx$  6,000 pictures from *Getty*
- 500 words vocabulary
- Prevalence covariates: agency and date

# IDENTIFYING POLITICAL COMPONENTS IN THE DATA GENERATION PROCESS OF IMAGES

- Goal: Identify and quantify the visual framing of the magnitude of the caravan
- Structural Topic Model to identify underlying “topics”, understood as frames, in the images
- Visual codebook generated from  $\approx$  6,000 pictures from *Getty*
- 500 words vocabulary
- Prevalence covariates: agency and date
- Selection of 6 topics:

# IDENTIFYING POLITICAL COMPONENTS IN THE DATA GENERATION PROCESS OF IMAGES

- Goal: Identify and quantify the visual framing of the magnitude of the caravan
- Structural Topic Model to identify underlying “topics”, understood as frames, in the images
- Visual codebook generated from  $\approx$  6,000 pictures from *Getty*
- 500 words vocabulary
- Prevalence covariates: agency and date
- Selection of 6 topics:
  - Crowd

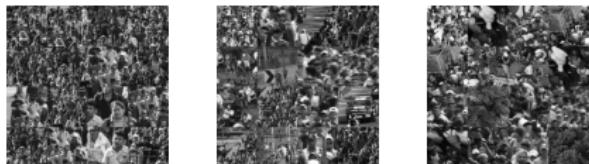
# IDENTIFYING POLITICAL COMPONENTS IN THE DATA GENERATION PROCESS OF IMAGES

- Goal: Identify and quantify the visual framing of the magnitude of the caravan
- Structural Topic Model to identify underlying “topics”, understood as frames, in the images
- Visual codebook generated from  $\approx$  6,000 pictures from *Getty*
- 500 words vocabulary
- Prevalence covariates: agency and date
- Selection of 6 topics:
  - Crowd
  - Border/Fence, Small group/Portrait, Water/Sky, Camps, Darkness

# **UNDERLYING TOPICS IN THE CARAVAN: FREX WORDS**

# UNDERLYING TOPICS IN THE CARAVAN: FREX WORDS

**Topic 1:** Crowds



**Topic 2:** Border/Fence



**Topic 3:** Water/Sky



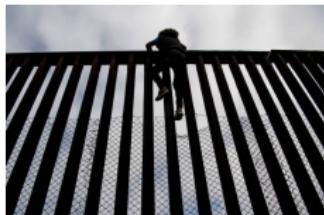
## **UNDERLYING TOPICS IN THE CARAVAN: REPRESENTATIVE IMAGES**

# UNDERLYING TOPICS IN THE CARAVAN: REPRESENTATIVE IMAGES

Crowd



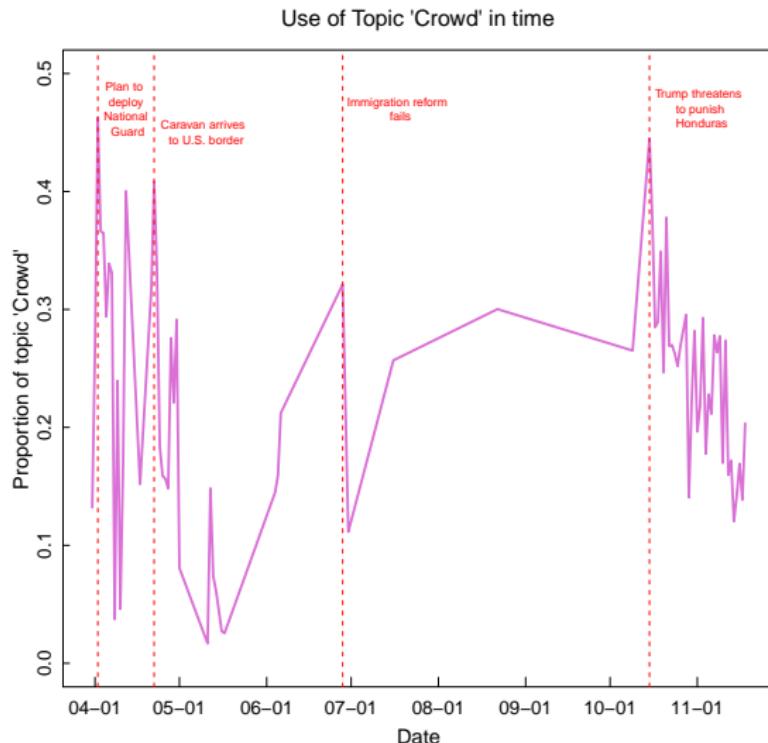
Border/  
Fence



Water/  
Sky



# CROWD TOPIC IN TIME



## VALIDATION: HIGH CORRELATION BETWEEN TOPICS AND MANUAL CODING

## VALIDATION: HIGH CORRELATION BETWEEN TOPICS AND MANUAL CODING

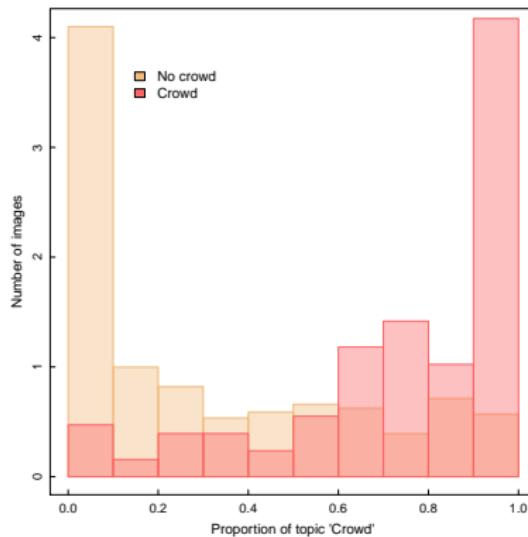
- Hand-coded sample: presence of medium/big crowd in the image (crowd=1) or no (crowd=0).

## VALIDATION: HIGH CORRELATION BETWEEN TOPICS AND MANUAL CODING

- Hand-coded sample: presence of medium/big crowd in the image ( $\text{crowd}=1$ ) or no ( $\text{crowd}=0$ ).
- Correlation with proportion topic “Crowd”: 0.58

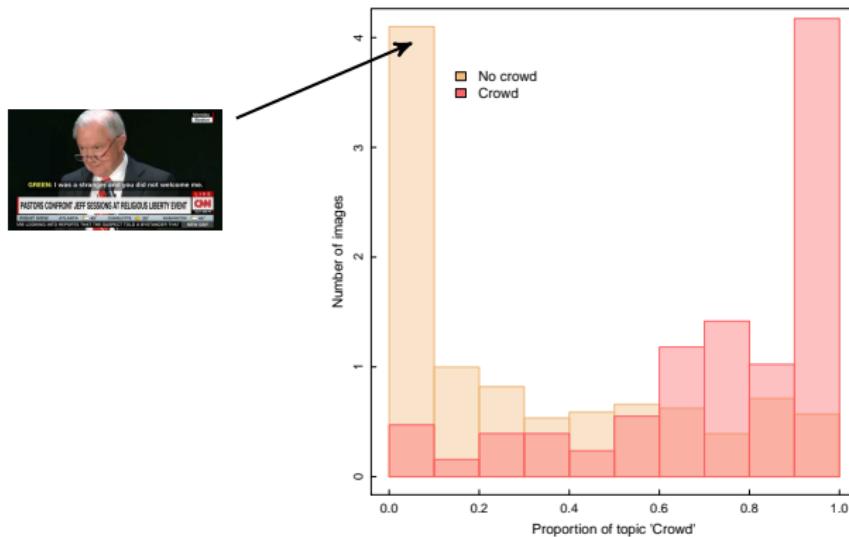
## VALIDATION: HIGH CORRELATION BETWEEN TOPICS AND MANUAL CODING

- Hand-coded sample: presence of medium/big crowd in the image (crowd=1) or no (crowd=0).
- Correlation with proportion topic “Crowd”: 0.58



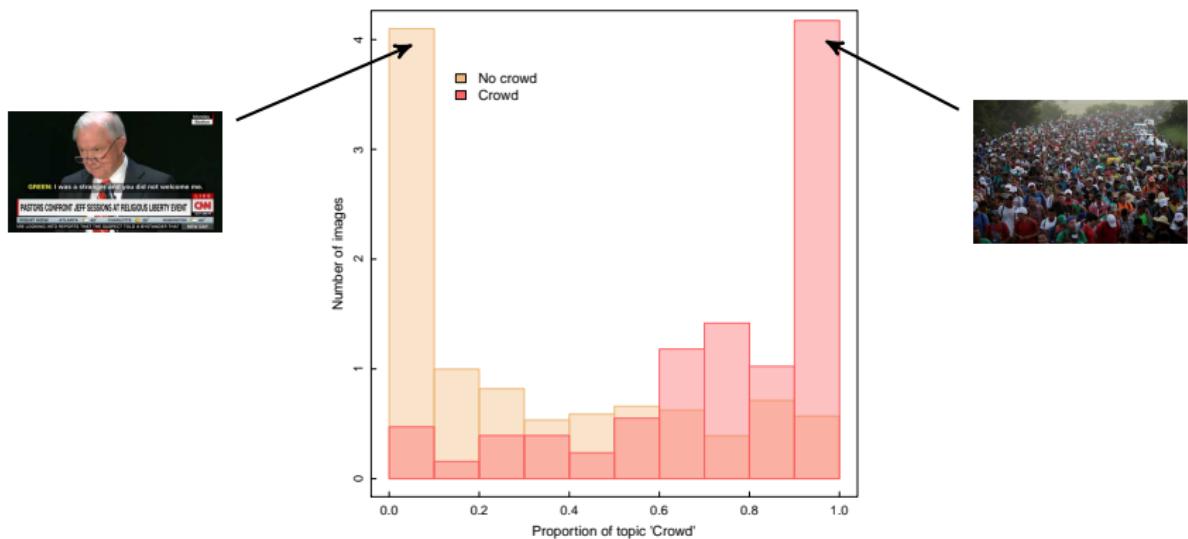
## VALIDATION: HIGH CORRELATION BETWEEN TOPICS AND MANUAL CODING

- Hand-coded sample: presence of medium/big crowd in the image (crowd=1) or no (crowd=0).
- Correlation with proportion topic “Crowd”: 0.58

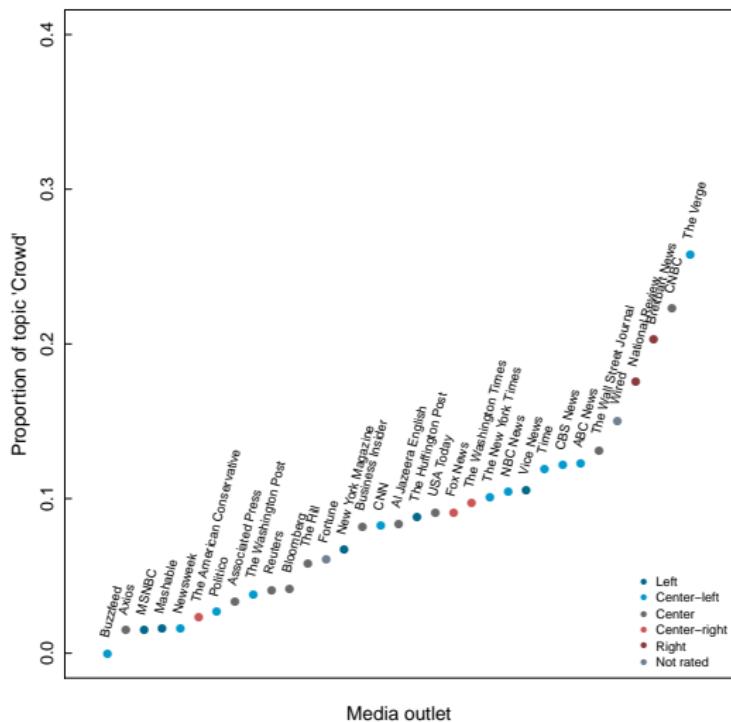


# VALIDATION: HIGH CORRELATION BETWEEN TOPICS AND MANUAL CODING

- Hand-coded sample: presence of medium/big crowd in the image (crowd=1) or no (crowd=0).
- Correlation with proportion topic “Crowd”: 0.58



# TOPIC “CROWD” BY MEDIA OUTLET

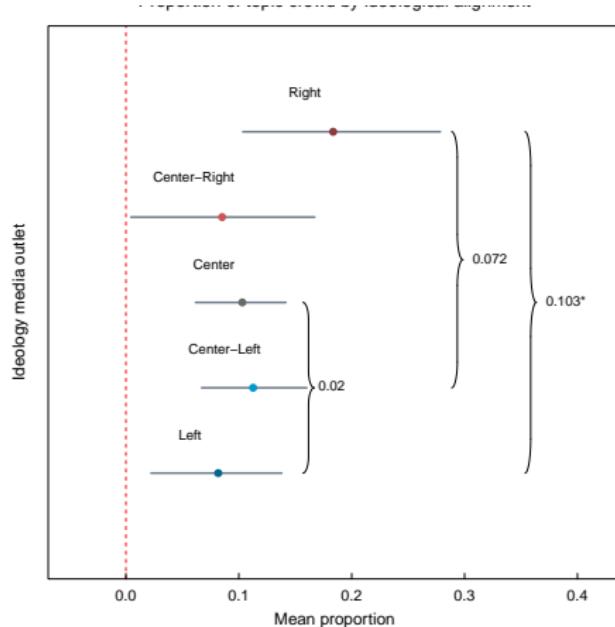


## FACTORS BEHIND THE GENERATION OF VISUAL FRAMES

- Estimate the effect of media ideology on prevalence of topic “Crowd”

# FACTORS BEHIND THE GENERATION OF VISUAL FRAMES

- Estimate the effect of media ideology on prevalence of topic “Crowd”



## THINGS TO CONSIDER

- Coherence and symbolism of visual words

## THINGS TO CONSIDER

- Coherence and symbolism of visual words
  - Some tokens are similar in terms of features but not concept

## THINGS TO CONSIDER

- Coherence and symbolism of visual words
  - Some tokens are similar in terms of features but not concept
  - Others are more synonyms than exclusive words

## THINGS TO CONSIDER

- Coherence and symbolism of visual words
  - Some tokens are similar in terms of features but not concept
  - Others are more synonyms than exclusive words
  - Some visual words are plainly “bad”

## THINGS TO CONSIDER

- Coherence and symbolism of visual words
  - Some tokens are similar in terms of features but not concept
  - Others are more synonyms than exclusive words
  - Some visual words are plainly “bad”
- Dimension reduction

## THINGS TO CONSIDER

- Coherence and symbolism of visual words
  - Some tokens are similar in terms of features but not concept
  - Others are more synonyms than exclusive words
  - Some visual words are plainly “bad”
- Dimension reduction
  - Mapping of latent concepts to low-dimensional features

## THINGS TO CONSIDER

- Coherence and symbolism of visual words
  - Some tokens are similar in terms of features but not concept
  - Others are more synonyms than exclusive words
  - Some visual words are plainly “bad”
- Dimension reduction
  - Mapping of latent concepts to low-dimensional features
  - Sensitivity of results to feature definition/extraction

## THINGS TO CONSIDER

- Coherence and symbolism of visual words
  - Some tokens are similar in terms of features but not concept
  - Others are more synonyms than exclusive words
  - Some visual words are plainly “bad”
- Dimension reduction
  - Mapping of latent concepts to low-dimensional features
  - Sensitivity of results to feature definition/extraction
- Interpretation of results

## THINGS TO CONSIDER

- Coherence and symbolism of visual words
  - Some tokens are similar in terms of features but not concept
  - Others are more synonyms than exclusive words
  - Some visual words are plainly “bad”
- Dimension reduction
  - Mapping of latent concepts to low-dimensional features
  - Sensitivity of results to feature definition/extraction
- Interpretation of results
  - “tea leaves reading”

## THINGS TO CONSIDER

- Coherence and symbolism of visual words
  - Some tokens are similar in terms of features but not concept
  - Others are more synonyms than exclusive words
  - Some visual words are plainly “bad”
- Dimension reduction
  - Mapping of latent concepts to low-dimensional features
  - Sensitivity of results to feature definition/extraction
- Interpretation of results
  - “tea leaves reading”
- Resources

## THINGS TO CONSIDER

- Coherence and symbolism of visual words
  - Some tokens are similar in terms of features but not concept
  - Others are more synonyms than exclusive words
  - Some visual words are plainly “bad”
- Dimension reduction
  - Mapping of latent concepts to low-dimensional features
  - Sensitivity of results to feature definition/extraction
- Interpretation of results
  - “tea leaves reading”
- Resources
  - Learning curve

## THINGS TO CONSIDER

- Coherence and symbolism of visual words
  - Some tokens are similar in terms of features but not concept
  - Others are more synonyms than exclusive words
  - Some visual words are plainly “bad”
- Dimension reduction
  - Mapping of latent concepts to low-dimensional features
  - Sensitivity of results to feature definition/extraction
- Interpretation of results
  - “tea leaves reading”
- Resources
  - Learning curve
  - Time?

# CONCLUSION

- Images are powerful (and abundant!) elements of frames.

## CONCLUSION

- Images are powerful (and abundant!) elements of frames.
- The BoVW can be used to detect relevant political components of visual frames ⇒ unsupervised and semi-supervised methods.

## CONCLUSION

- Images are powerful (and abundant!) elements of frames.
- The BoVW can be used to detect relevant political components of visual frames ⇒ unsupervised and semi-supervised methods.
- It is intuitive, computationally cheap and tractable

# CONCLUSION

- Images are powerful (and abundant!) elements of frames.
- The BoVW can be used to detect relevant political components of visual frames ⇒ unsupervised and semi-supervised methods.
- It is intuitive, computationally cheap and tractable
- It provides a solid input for tools and methods that social scientists are familiar with and use to answer substantive questions

# CONCLUSION

- Images are powerful (and abundant!) elements of frames.
- The BoVW can be used to detect relevant political components of visual frames ⇒ unsupervised and semi-supervised methods.
- It is intuitive, computationally cheap and tractable
- It provides a solid input for tools and methods that social scientists are familiar with and use to answer substantive questions
- The measurement and analysis of visual components are crucial to have a better understanding of political communication

# LET'S CODE!

