

ENCODING AND TRANSFORMING DATA

LECTURE #8
SECTION 1
JUNE 15



INSTITUTE FOR APPLIED
COMPUTATIONAL SCIENCE
AT HARVARD UNIVERSITY



UNIVERSITY *of*
RWANDA

CATEGORICAL DATA

CATEGORICAL DATA:

So far our covariates have been numerical data, data that are recorded as real numbers and these values can be ordered, e.g. $\text{Age}=50$ is greater than $\text{Age}=20$.

But many important real-life measurements are not numerical, e.g. if $\text{nationality} = \text{"USA"}, \text{"Rwanda"}, \dots$ Why can't we simply encode the values as numbers?

$$\text{"USA"} = 0, \text{"Rwanda"} = 1, \dots$$

Encoding categories as numbers is misleading, since numbers imply ordering:

$$\text{"Rwanda"} = 1 > 0 = \text{"USA"} \dots$$

Instead, we encode these categorical variables as binary vectors:

One-hot Encoding

<u>Nationality</u>	"USA"	"Rwanda"	"China"	"Mexico"
"USA"	[1 , 0 , 0 , 0]			
"Rwanda"		[0 , 1 , 0 , 0]		
"China"			[1 , 0 , 0]	
"Mexico"				[0 , 0 , 0 , 1]

TEXT DATA

REPRESENTING TEXT DATA:

Sometimes our covariates contain entire blocks of text, e.g.

diagnostic-notes-1 = "The Patient has fever and chill and has not eaten."

diagnostic-notes-2 = "The patient has no fever and had no food."

We can encode each word using one-hot encoding and add-up all the onehot vectors. This simply counts the number of times each word occurs in the text. This is called the **count vector**:

Count Vector Representations

	the	Patient	has	fever	and	chills	not	eaten	no	food
diagnostic-notes-1	[1	1	2	1	2	1	1	1	0	0]
diagnostic-notes-2	[1	1	1	1	1	0	0	0	1	1]

We need to worry about:

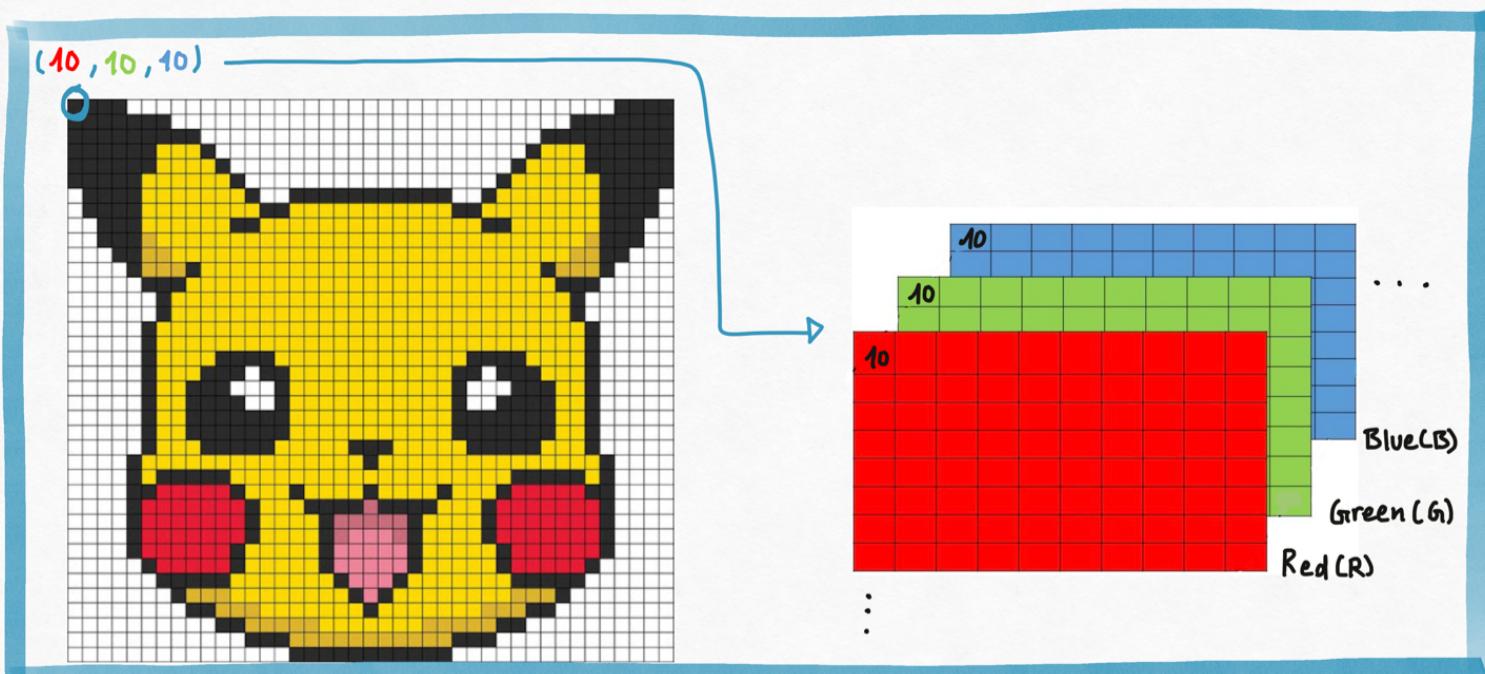
1. punctuation
2. conjugation
3. words with little meaning: "the", "an", "and", ...
4. ambiguity: "the river bank", "the bank account".

IMAGE DATA

REPRESENTING IMAGE DATA:

Images are represented as a grid of numbers. Each cell in the grid is called a pixel and consists of three numbers: (R, G, B) for how much red, green and blue color to display for the pixel.

So each image with $L \times H$ pixels is represented by 3 $L \times H$ matrices, one for each colour channel.



FLATTENED IMAGE VECTORS:

The representation of images as three color matrices is inconvenient to work with, all of our models take vectors as inputs.

In order to apply our models to image data, we flatten the matrices into a single vector.

