

HANDS-ON AI I

Reading, Handling and Visualization of Datasets



Andreas Schörgenhuber
Institute for Machine Learning

Copyright Statement

This material, no matter whether in printed or electronic form, may be used for personal and non-commercial educational use only. Any reproduction of this material, no matter whether as a whole or in parts, no matter whether in printed or in electronic form, requires explicit prior acceptance of the authors.

Content of Unit 2

- Second data source: images
- Third data source: sequential data
 - Sound
 - Text

Images

A screenshot of a search engine results page. The top navigation bar includes a logo, a search bar with the word "images", and links for "Bilder", "Videos", "Maps", "News", "Mehr", "Einstellungen", and "Tools". Below the navigation are several circular buttons with labels: flower, photos, nature, wallpapers, rose, pictures, jpg, beautiful, unsplash, pic, pics, pixabay, css, stock images, shutterstock, and stock photos. The main content area displays a grid of images and their descriptions. At the top left is a large blue rose image with the text "Flower Images - Pixels - Free Stock Photos" and a link to pixels.com. Next is a bridge over water with the text "100+ Bridge Pictures | ... unplash.com". Then a nature scene with butterflies with the text "Nature Images - Pixels - Free Stock Photos" and a link to pixels.com. Following is a lake at night with a bridge, labeled "Beauty Images - Pixels - Free Stock Photo..." and a link to pixels.com. Below these are more images: a bottle in the sand labeled "5,000+ Free Bottle & Wine Images - Pixabay" and a link to pixabay.com; a sunset over the ocean labeled "900+ Sunset Images: Download HD Pictures & Photo..." and a link to unplash.com; a forest path with autumn leaves labeled "60,000+ Free Forest & Nature Images - Pixabay" and a link to pixabay.com. The bottom row includes images of a beach, peacock feathers, roses, a lotus flower, and a heart-shaped flower, each with its respective category name and a link to pixels.com.

Representation of Images

- Images are usually represented in **3 dimensions**:
 - Height
 - Length (or width)
 - Channels (often: red, green, blue)
- **Color depth** refers to the number of possible values for each channel of a pixel.
 - Color depths of **8-bit** ($2^8 = 256$ values) and **16-bit** ($2^{16} = 65536$ values) are common.
 - Higher values mean increased intensity.

0	0	0	0	0	0	0
0	255	0	0	124	0	0
0	0	0	124	0	0	0
0	0	0	0	124	0	0
0	0	0	124	0	0	0
0	0	0	0	124	0	0

RGB Model

- The RGB (red, green, blue) model is the most important model for colored images.
- 3 channels: **red, green, blue**
- Adding up the 3 color channels results in the final image.



RGB Model

- The RGB (red, green, blue) model is the most important model for colored images.
- 3 channels: **red, green, blue**
- Adding up the 3 color channels results in the final image.



RGBA Model

- The RGBA model uses a fourth channel: **alpha**.
- Alpha controls the transparency of the image.
- Higher values mean increased opaqueness ($\alpha = 0 \rightarrow$ fully transparent, $\alpha = \text{max} \rightarrow$ fully opaque).



Grayscale

- Grayscale images have only **one channel representing the brightness**.
- Color images can be converted to grayscale images (**different grayscale conversions** exist), but not vice versa (without additional information).



Data Augmentation

- The field of **computer vision** focuses on processing and understanding images.
- Strong shift towards machine learning as **a way to understand, manipulate and even generate images.**
- Example: We “show” a machine learning model a large number of images with and without Charlie.
 - Train it to distinguish between these two groups.
 - **Augment** the images to reduce the workload of getting a lot of Charlie images.
 - Feed original and/or augmented images to our model.

Data Augmentation

- In many cases, we can augment data artificially without collecting new real samples: Create “new” **artificial samples** by modifying existing samples.
- Pros:
 - Can increase the number of data points by a large factor with little effort (often on-the-fly).
 - Reduces overfitting, increases robustness of model (we will learn about this in a later unit).
- Cons:
 - Can introduce artifacts or change the task.
 - Heavily dependent on data, model and task.

Data Augmentation Techniques

- There exist a lot of data augmentation **techniques**.
- Some can be applied across different fields (e.g., adding noise), some are specific and only applicable in certain scenarios.
- Here are some common examples in the area of **image augmentation**:
 - Rotation
 - Flipping
 - Zooming/Cropping
 - Blurring
 - Noise
 - Input Dropout
 - Distortion Effects
 - Color Jittering
 - ...

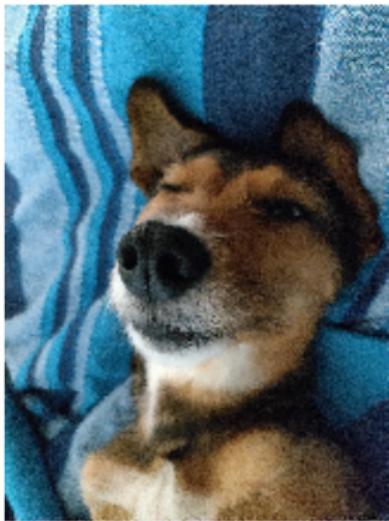
Rotation

- Rotate images around their center by a certain angle.
- It always depends on the task how much rotation is reasonable.



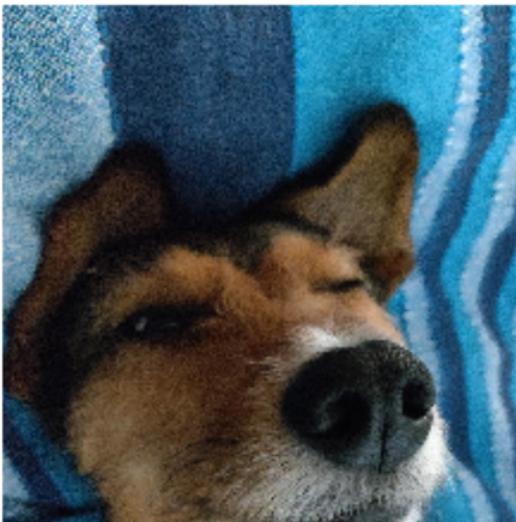
Flipping

- Mirroring the image across its vertical or horizontal axis.
- It depends on the task if flipping is reasonable.



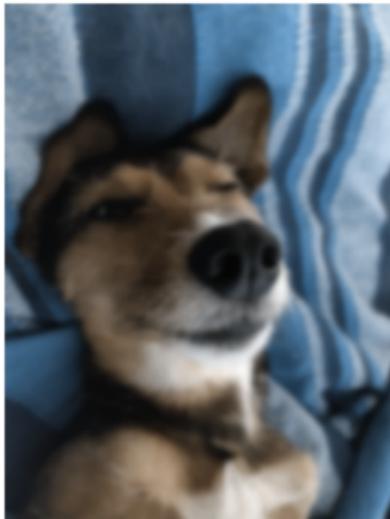
Zooming/Cropping

- Zoom into the image (different sizes of Charlie).
- Often, the same size of the image is required:
 - Selecting a smaller part and resizing it to the original resolution.



Blurring

- Usually not all images are super sharp.
 - Blurred images help the machine learning model to detect blurred Charlies as well.
- Gaussian blurring is most common blurring technique.



Noise

- Add random noise such that input feature values (pixels) are still plausible.
- Noise strength can be parameterized.



original



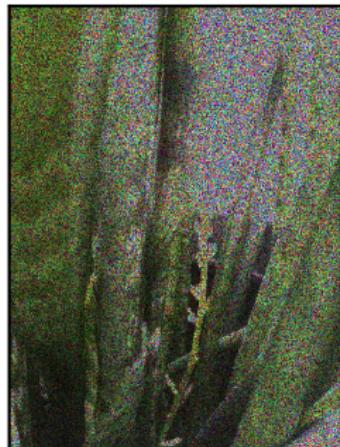
light noise



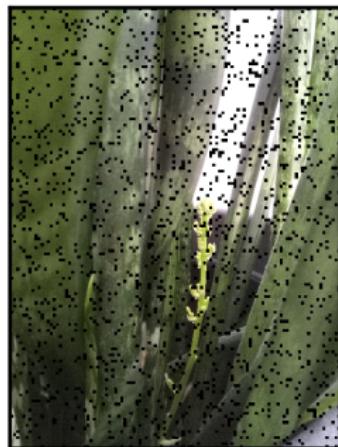
heavy noise

Input Dropout

- Temporarily remove input features to mimic data errors or hidden image content.
- Different dropout techniques exist, e.g., removing individual R, G or B values, or even entire (adjacent/neighboring) pixels over all channels.



dropout random values



adjacent pixels

Distortion Effects

- Warp image using, e.g., lens-distortion-like effects.



Color Jittering

- Modify color channels, brightness, hue, saturation or value.



Some Notes on Data Augmentation

- Very useful to get more data without having to explicitly “collect” it from the source.
- Increases variation and thereby allows for a better generalization.
- **Careless augmentations** can actually have a **negative impact**. Example:
 - Extreme zoom on Charlie’s nose.
 - Only black pixels remain with no resemblance of Charlie whatsoever.
 - Training a machine learning model with this input (while still saying this is Charlie) will yield worse results.
- Data augmentation is typical in image processing but not restricted to it.
 - Example: Add random noise to time series signal.

Image Segmentation

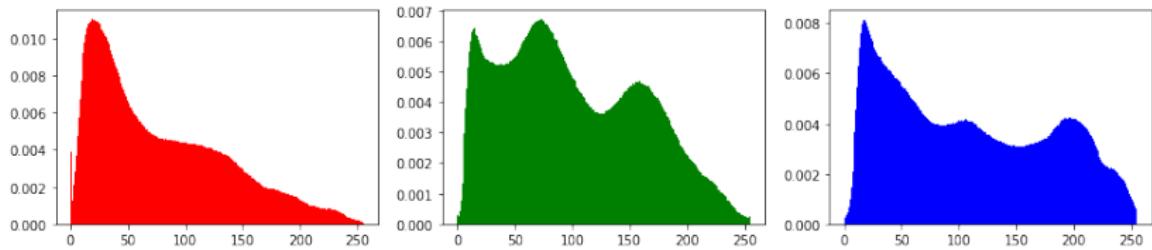
- A common problem in computer vision is to mark an object within an image, e.g., a dog with some background.



- We segment the image into **object** and **no object: image segmentation**.
- **Wide field** of research.

Simple Solution to Segment Charlie

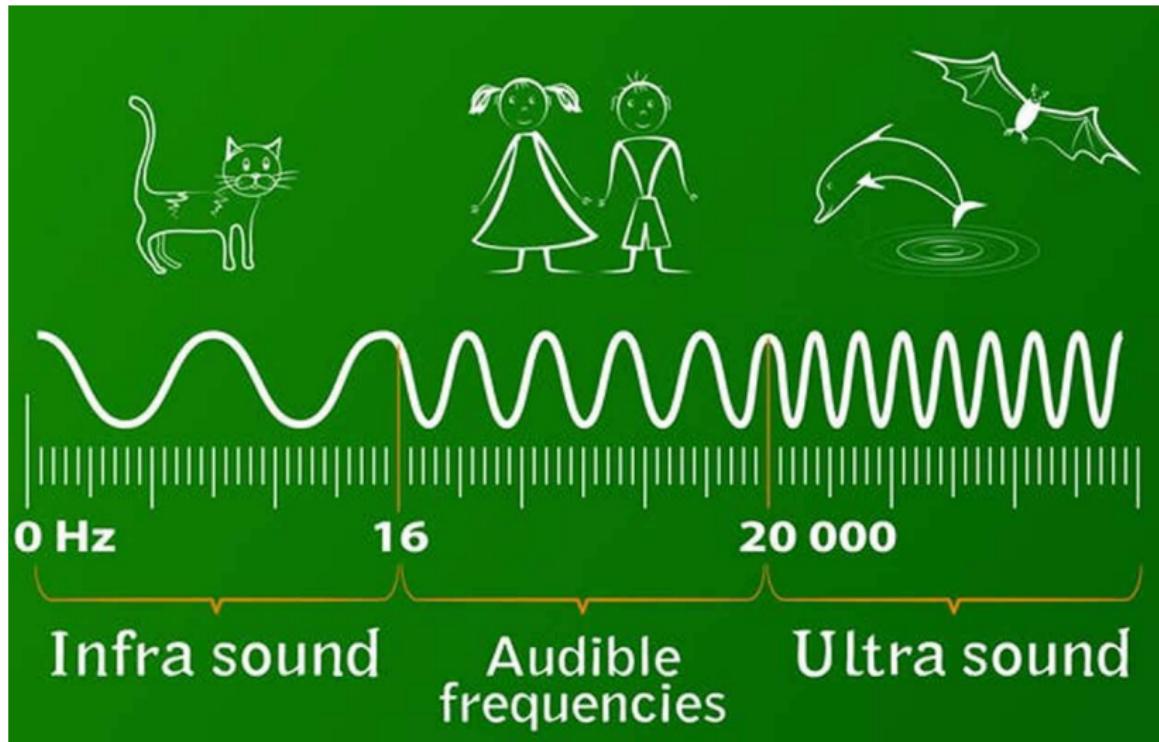
- Occurance of individual color values of an image.
- Pixels with RGB values ranging from 0 to 255.
- Setting thresholds to segment between Charlie and background.



Sequences

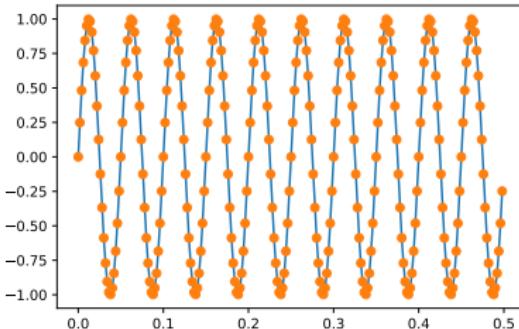
- A **sequence** is a type of data where the order of the values matters: $(a_1, a_2, a_3, a_4, \dots, a_n)$
- Some examples:
 - Text
 - Sound (or time series in general)
 - Strictly speaking: images (order of pixels is essential)
 - ...

Sequences: Sound



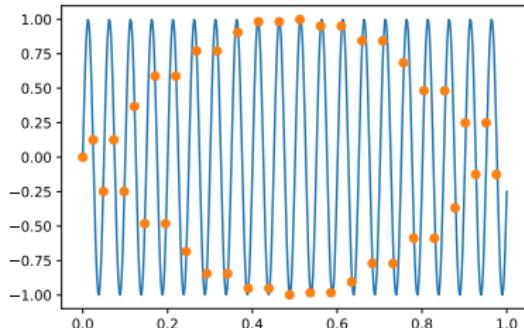
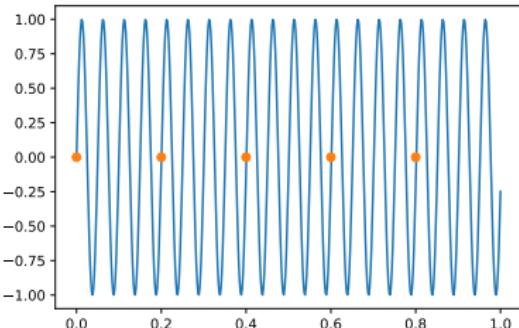
Continuous Signal

- Sound waves are a continuous signal.
- In physics, the **sound frequency** is measured in Hertz (Hz).
- In data analysis, the **sampling rate** from the continuous signal (per second) is also measured in Hz.
 - For example, a sampling rate of 48 kHz means retrieving 48000 samples per second.



Low Sampling Rates

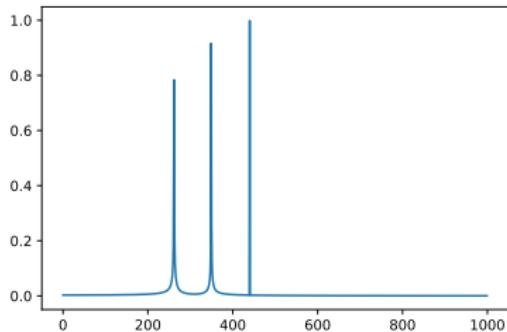
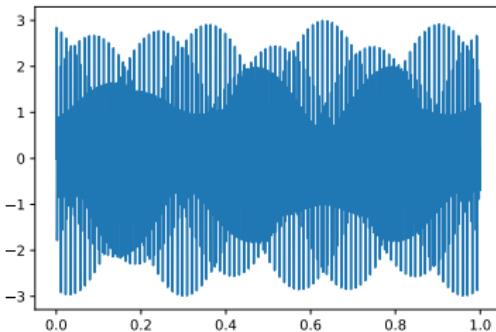
- If the sampling rate is too low, it is not possible to exactly reconstruct the signal.
- **Nyquist-Shannon theorem:** A signal is completely determined by a series of points that are $< \frac{1}{2f}$ seconds apart (f is the highest frequency in the signal), or in other words, the sampling rate f_s must be $> 2f$.



- Higher sampling rates mean higher audio quality.

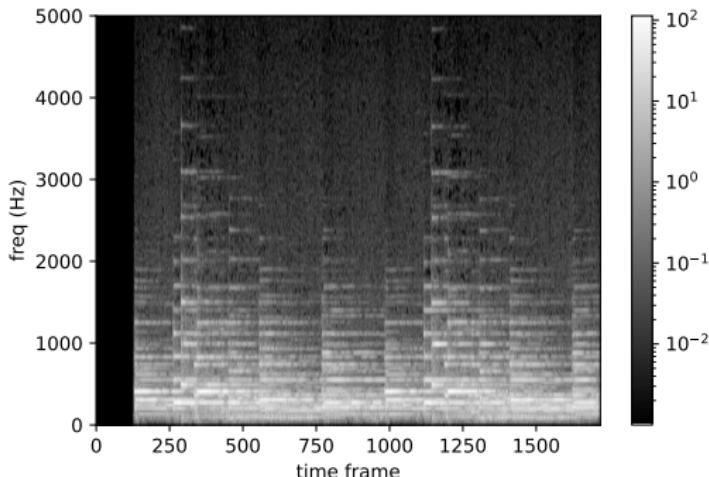
Fourier Transformation

- The Fourier transformation decomposes a signal (a function of time) into its constituent frequencies (a spectrum).
- The Fourier transformation is not limited to functions of time, but the domain of the original function is commonly referred to as the **time domain** (example left image).
- The domain of the results is typically called the **frequency domain** (example right image).



Spectrogram

- Spectrograms are visual **representations of the variation of the frequencies over time**:
 - This can be done by computing a spectrum for the first few milliseconds (time frame), another one for the next, etc.
 - Stitching these spectra together, we obtain a spectrogram.
 - It tells how much energy (=gray value) is in the signal at each time frame and frequency bin (=pixel).



Sequences: Text

Vivamus vehicula leo a justo. Quisque nec augue. Morbi mauris wisi, aliquet vitae, dignissim eget, sollicitudin molestie, ligula. In dictum enim sit amet risus. Curabitur vitae velit eu diam rhoncus hendrerit. Vivamus ut elit. Praesent mattis ipsum quis turpis. Curabitur rhoncus neque eu dui. Etiam vitae magna. Nam ullamcorper. Praesent interdum bibendum magna. Quisque auctor aliquam dolor. Morbi eu lorem et est porttitor fermentum. Nunc egestas arcu at tortor varius viverra. Fusce eu nulla ut nulla interdum consectetur. Vestibulum gravida. Morbi mattis libero sed est.

Natural Language Processing (NLP)

- NLP is a subfield of:
 - Computer science
 - Linguistics
 - Deep Learning
- Challenges in NLP involve:
 - Speech recognition
 - Language understanding
 - Language generation
 - Language translation
- Nowadays, NLP is strongly driven by modern **Deep Learning** methods. We will learn more about this in Hands-on AI II.

Text as Categorical Data – Theory

- Normally, text can be considered qualitative data without mathematical meaning.
- Mathematical operations (e.g., \sum or comparisons such as “less than”) do not make sense.
- Example: “Dog”, “Rat”, “Cat”

Text as Categorical Data – Practice

- Assume we consider n different categories ($= n$ different words), which we can also call **vocabulary** or **dictionary**.
- We could represent each category as an integer value:
 - Requires n different integer values, one for each category.
 - Example: “Dog” = 0, “Rat” = 1, “Cat” = 2
- Problem:
 - We would introduce new (probably false) information.
 - Example: In our ranking, “Dog” < “Rat” and “Rat” · 2 = “Cat”.
 - Not suitable for us.

Text as Categorical Data – Practice

- Solution: Represent a categorical feature as **one-hot-encoded** binary vector $\mathbf{v} = (v_1, \dots, v_n)$:
 - Categorical data with n different values is enumerated from $1 \dots n$ given a **fixed order**.¹
 - $v_i = 1$ if category i is true, otherwise $v_i = 0$.
 - Each element in the vector represents one category → no false information!²
- Example:
 - Possible values: “Dog”, “Rat”, “Cat” ($n = 3$)
 - Sample is “Dog” → $\mathbf{v} = (1, 0, 0)$
 - Sample is “Cat” → $\mathbf{v} = (0, 0, 1)$

¹The order itself is not important, just that it is fixed.

²Only applies if information about order in feature vector is not used.

One-Hot Encoding Example

- Vocabulary with $n = 7$ and the following fixed order:

Word	an	awesome	cat	Charlie	dog	is	rat	
Position	1		2	3	4	5	6	7

- Recall: For each word, we have a vector of size n where all numbers are zero but the element representing a specific word is set to 1.
- Example one-hot-encoded word “awesome” (position $= 2 \rightarrow v_2 = 1$, all other $v_i = 0$):
 $v_{\text{awesome}} = (0, 1, 0, 0, 0, 0, 0)$

One-Hot Encoding Example

- Entire one-hot encoded vocabulary:

	an	awesome	cat	Charlie	dog	is	rat
an	1	0	0	0	0	0	0
awesome	0	1	0	0	0	0	0
cat	0	0	1	0	0	0	0
Charlie	0	0	0	1	0	0	0
dog	0	0	0	0	1	0	0
is	0	0	0	0	0	1	0
rat	0	0	0	0	0	0	1

Word Embedding

- Problem: For large vocabularies, the usage of one-hot encoding is not ideal. For example, a vocabulary size of $n = 30\,000$ would yield vectors of the same size (= 30 000 dimensional feature vectors)!
- **Word embedding**: Represent each word by a much smaller vector but with real numbers.
- Ideally, a word embedding captures additional information compared to a simple encoding, e.g., (semantic) similarities between words.
- How to get a “good” embedding?
 - Not straightforward.
 - Often learned from a big corpus of words (in Hands-on AI II, we will actually learn an embedding ourselves).
 - Prelearned embeddings exist.

Downprojection of Word Embedding

- As always, visualizing our data can help in gaining insights.
- While embeddings are typically much smaller than simple encodings (e.g., 300 features per vector), this is still too much to visualize.
- Solution: **Reduce dimensionality** to **downproject** the embedding.
- If the embedding incorporates similarities, the downprojection of embedded words enables us to identify those.

Downprojection of Word Embedding

