

COGS138: Neural Data Science

Lecture 16

C. Alex Simpkins, PhD

UCSD Dept. of Cognitive Science, Spring 2023

RDPRobotics, LLC

http://casimpkinsjr.radiantdolphinpress.com/pages/cogs138_sp23

rdprobotics@gmail.com | csimpkinsjr@ucsd.edu

Plan for today

- Announcements
- In class paper
- Parameterizing heterogeneous datasets II in neural data science, issues, techniques, readings
- Dimensionality reduction in NDS

Announcements

- **Deadlines upcoming this week:**
- **Tuesday** (this week should be complete) :
 - ***Previous project review 11:59pm***
 - ***Mid-quarter checkin 11:59pm***
 - ***Lecture Quiz 3 (EXTENDED) 11:59pm***
- **Wednesday (yesterday) Assignment 3 (EXTENDED) 11:59pm**
- **Saturday:**
 - Reading Quiz 3 11:59pm
- **Friday:**
 - Project proposal 11:59pm

Announcements

- We will review the mid-quarter check-in comments and grade check-in comments and get back to you
- A few asked about missed items and had other questions, some haven't been responded to yet
- If no response by Saturday evening reach out over email to Siddhant and I please

Announcements

- **github** repos
 - created,
 - invites sent,
 - please accept (time limited)
 - login and be sure you can and files are there, rename
- if you don't have an invite, there's an issue with your group record in the main list - please contact us asap
- **Procedure** : Contact Siddhant, cc me, if no response in a day, reach out to me again, I'll help

Announcements

- **Project meetings**

- You will need to sign up for a 10 min project meeting to introduce your project idea and get feedback ahead of your proposal so we can provide input early
- signup form here: <https://forms.gle/kLrLMrSs5qwC8f2P7>

Project schedule

Task due	Date due	Description
Previous project review	5/23/2023 at 11:59pm (Tuesday)	Select 2 of the 3 available, review as individuals and then come together as a group to submit your responses to the questions after a discussion. This will orient you to the class project
Project proposal	5/26/2023 at 11:59pm (Friday wk8)	Generate your question, hypothesis, initial data sets you'll be working with, etc., describe your plan, schedule, who is doing what, potential issues, suggested analysis and how it will answer your question
Data checkpoint	6/2/2023 at 11:59pm (Friday wk9)	Builds on the proposal by taking the feedback from PP above and actually getting, loading, describing your data,
EDA checkpoint	6/9/2023 at 11:59pm (Friday wk10)	Builds on the previous checkpoint, essentially most of your analysis should be done by this point
Final report	6/15/2023 at 11:59pm (Thursday Fin wk)	Due Thursday of finals week so we can grade before the Tuesday deadline, otherwise your grade may be delayed
Group evaluations	6/15/2023 at 11:59pm (Thursday Fin wk)	You will evaluate each other based on participation and performance, this will contribute to your overall final project grade 5%)

Some notes on datasets

- 2 or more discussion and motivation
- What each portion represents as an idea (prev. project review, proposal, data, eda, final checkpoints)

Remaining assignments schedule

- A4 wk8-9, A5 wk9-10
- R4 wk8-9
- LQquiz wk 8, 9, 10
- Paper this week, mostly in class or via appointment

Parameterizing heterogeneous datasets

- Definition, review
 - What do we mean by **parameterization**?
 - Reminder of what data is and stepping back to the big picture - ***representation***
 - What are **heterogeneous** datasets?
 - What are the **challenges** and solutions?
- **Tools and practice** in neural data science
 - https://nwb-overview.readthedocs.io/en/latest/tools/tools_home.html
- Examples

Parameterization vs. Hyperparameterization

- **Parameterization** - the set of parameters that define the model unknowns to be fit, typically from data
 - For example, for $y = ax + b$, what are the parameters?
 - ANN - network weights
 - Calculated/learned from data
- **Hyperparameterization** - the set of parameters for machine learning in particular that define and control the learning process and are external to the model
 - Bisection algorithm for optimization - bisection parameter
 - ANN - parameters of the learning algorithm itself
 - Heuristic, can be set by practitioner, tunable for a given problem

Parameterization vs. Hyperparameterization

- **Parameters**

- Calculated/learned from data (“the fit”)
- Internal to model
- Chosen as part of model structure either manually or algorithmically

- **Hyperparameters**

- Heuristic, can be set by practitioner, tunable for a given problem
- External from model
- Optimal parameters are not known, and are different for different problems
- Other ex.: ANN learning rate, gradient descent step size, the k in k -nearest neighbors

Stepping back: What is data?

Stepping back: What is data?

- **Data** can be of many forms
- **Data** - any representation of information that has been recorded in a fixed or dynamic state [Simpkins, 2023]
- ***“(1) Factual information (such as measurements or statistics) used as a basis for reasoning, discussion, or calculation” [Webster’s, 2023]***
- ***(2) “Information in digital form that can be transmitted or processed” [Webster’s, 2023]***

What do we mean by 'representation'

What do we mean by 'representation'

- Some sort of arbitrary symbolic link between the reality and the symbols we use to model reality, such as words, numbers, pictures, graphs, sounds, videos, smells, textures, vibrations, gestures
- Further information: See Simpkins [COGS100 lecture 10](#) on representation, and read Norman ch3: "The power of representation"

Representation defined

- Cognitive age, Norman argues started when we started using sounds, gestures and symbols to refer to objects, things and concepts - when we started generating data!
- **Representation** : The sound, gesture, symbol is not the thing itself, it stands for, refers to it
- On representation not the reality

Powers of cognition come from abstraction and representation

- Ability to represent perceptions, experiences, thoughts in some medium other than what they occurred in
- Abstracted away from irrelevant details
- “The essence of intelligence” as he states - if representation is just right, new experiences, insights, creations emerge
- **We can make symbols then use them to do our reasoning**

Representing the dimensions requires different types of data entirely

- Ultimately in neural data science we are reasoning about the brain and behavior, how it's all interconnected and the dynamics of it
- Data makes it possible to reach beyond our immediate cognitive limitations to operate on information
 - We cannot see a neuron firing when we look at each other, we measure, but then must do something with that data, related it and connect it meaningfully to other things
 - As we have been reasoning, we need massive amounts of connections to understand the patterns of it all
 - Recording it all the same way often is impossible
 - EEG vs. Behavior, text, other dimensions

Data Structures Review

Structured data

- Can be stored in database SQL
- Tables with rows and columns
- Requires a relational key
- 5-10% of all data

Semi-structured data

- Doesn't reside in a relational database
- Has organizational properties (easier to analyze)
- CSV, XML, JSON

Unstructured

- Non-tabular data
- 80% of the world's data
- Images, text, audio, videos

(Semi-)Structured Data

Data that is stored in such a way that it is easy to search and work with. These data are stored in a particular format that adheres to organization principles imposed by the file format. These are the data structures data scientists work with most often.

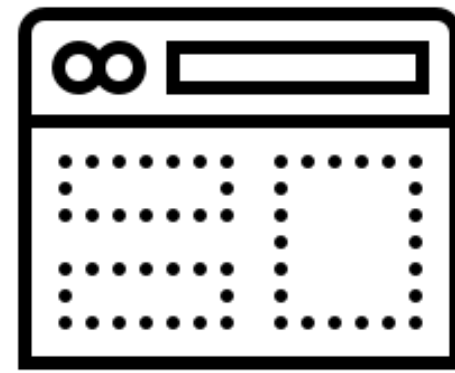
Unstructured Data

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.

Unstructured Data Types



Text files
and
documents



Websites
and
applications



Sensor
data



Image
files



Audio
files



Video
files



Email
data



Social
media
data

What are heterogeneous datasets?

- Given that ***data can represent anything that can be represented***, we can have many forms of sampling and recording systems
 - MOCAP
 - EEG/MEG
 - fMRI
 - Eye tracking
 - Text
 - Single unit recording
- What have we covered thus far for data types and forms?
- Others?

Why integrate them?

- More information can draw links that may not be clear otherwise
- Limited data source sets may not contain the necessary data for the question we want to ask
 - **Sparsity** - improved results with ***sparse*** datasets
 - **Modality** - one set might have patterns, but lack the content explaining patterns, the meaning underlying
 - **Reliability** - one dataset showing statistical significance vs. many confirming from various perspectives
 - <https://www.sciencedirect.com/science/article/pii/S1053811914003838>
 - <https://www.sciencedirect.com/science/article/pii/S1053811919300497>

Why is it a challenge to integrate them?

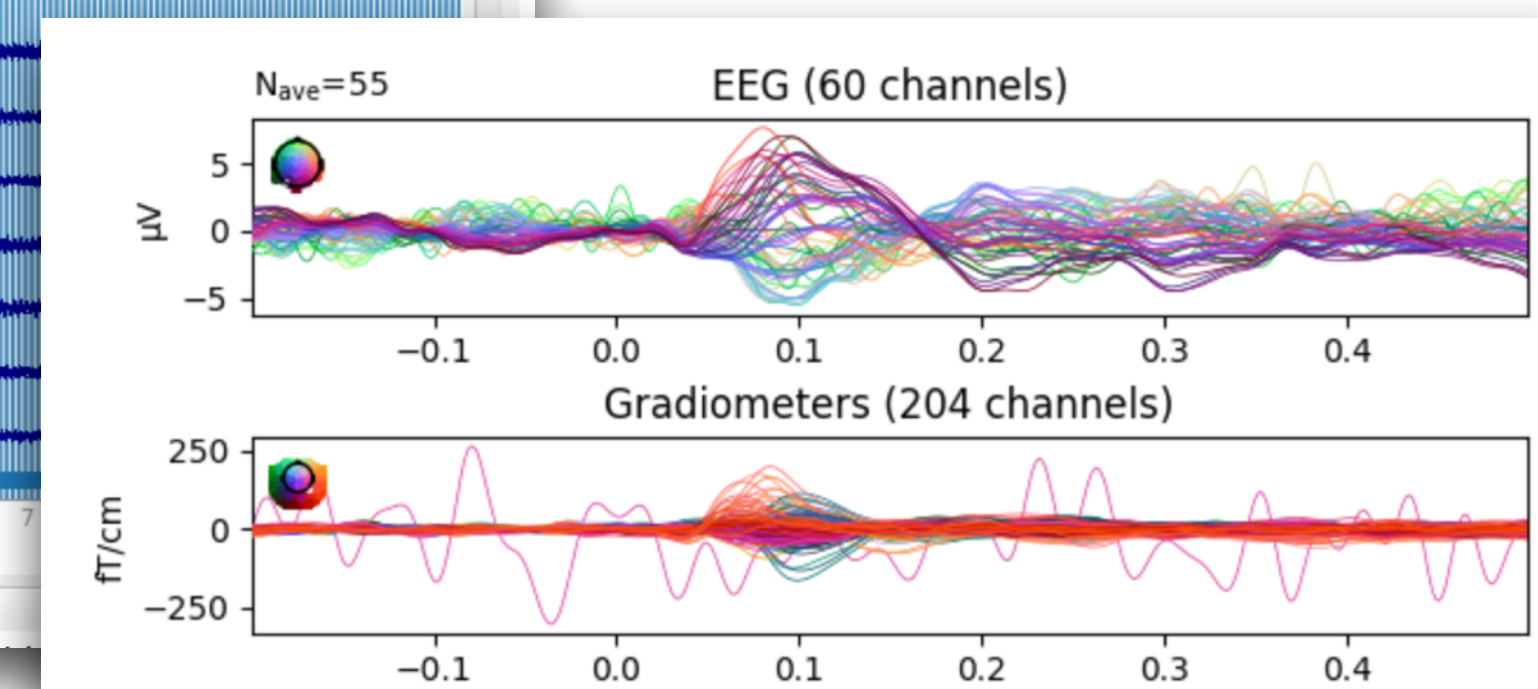
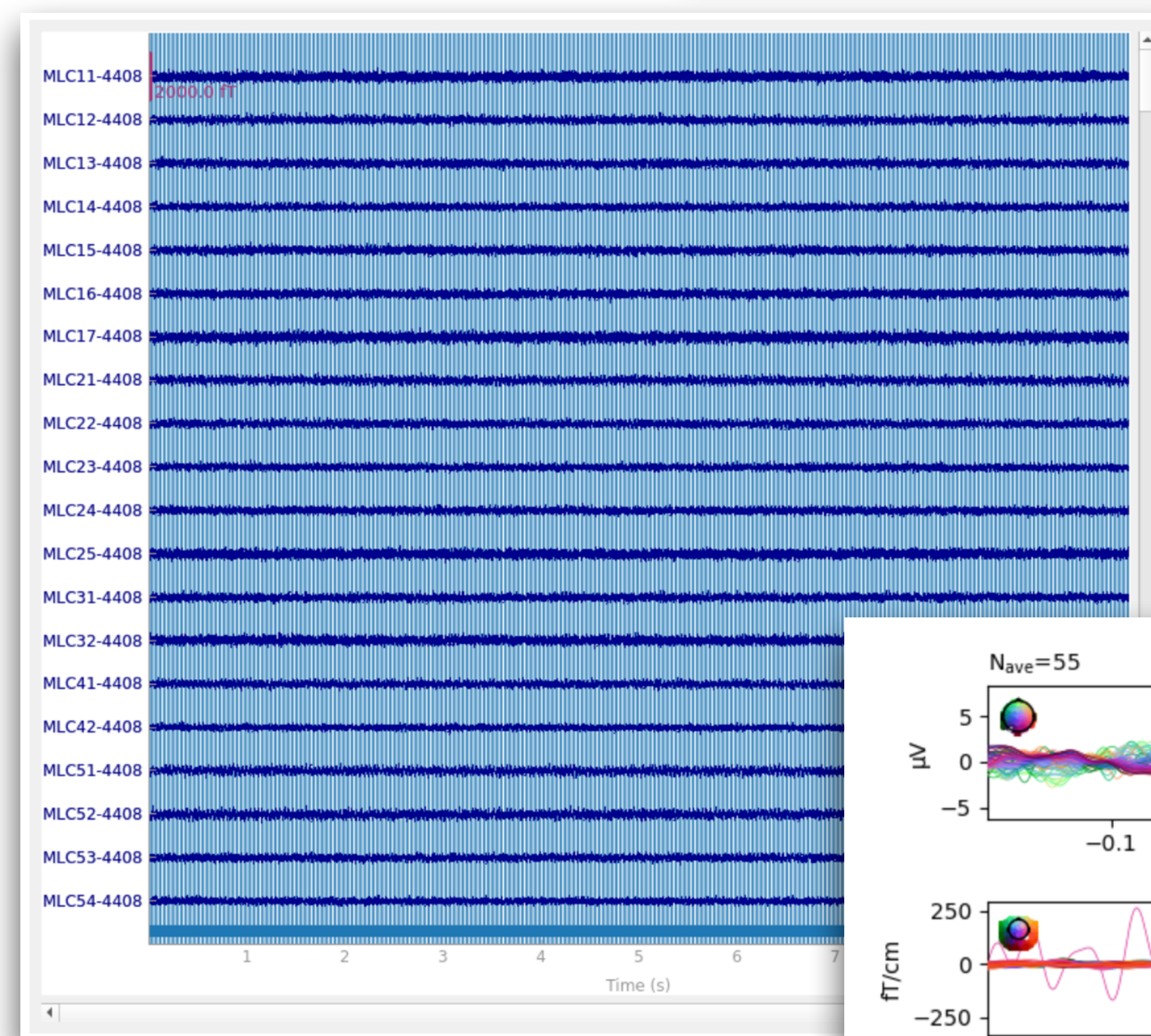
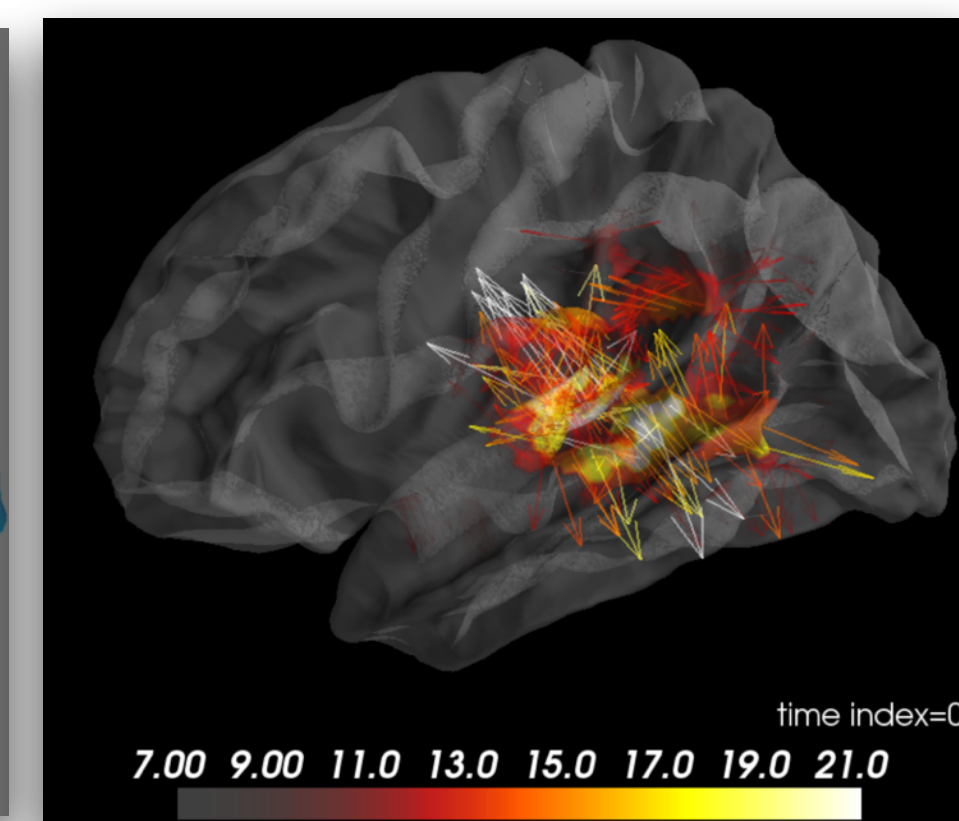
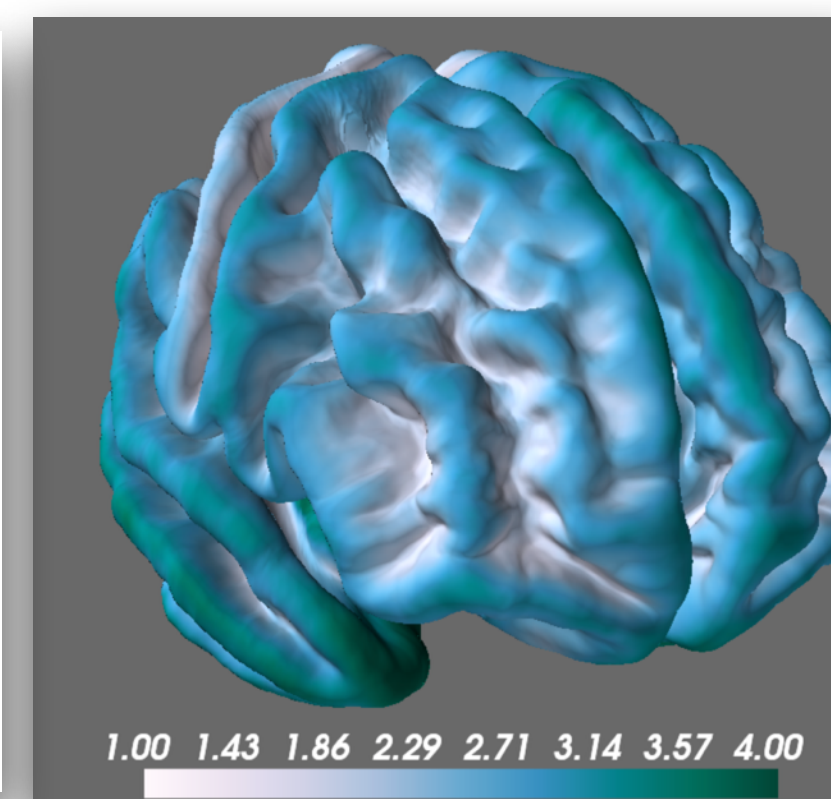
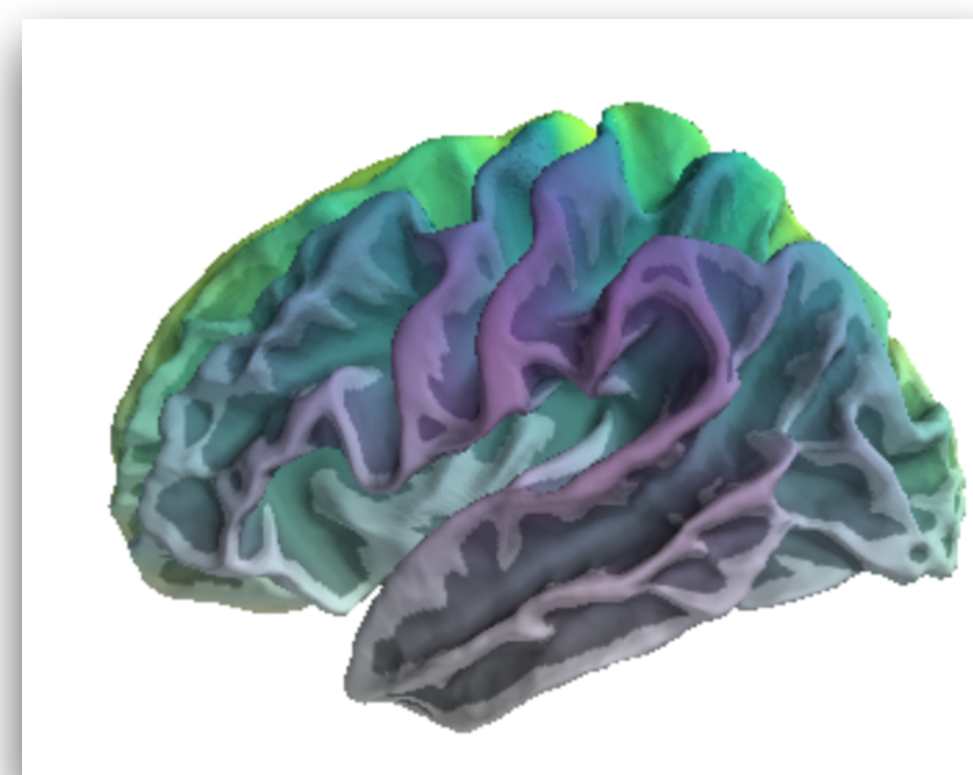
- Sampling rate mismatch
- Time/frequency/spatial domains - what is the best form of representation?
 - <https://www.sciencedirect.com/science/article/pii/S1053811919300497>
- Sample rate variability (why does this matter?)
- Sample time mismatch
- Format, software
- Missing data, data mixture/non-tabular etc
- Memory usage
- (Not an exhaustive list)

Integration strategies - *Sample rate mismatch*

- Resampling
 - Sub-sampling - (“down-sampling”) - every Mth sample, LowPass first (aliasing)
 - Super-sampling - (“up-sampling”) - padding with 0’s, then LowPass to interpolate
- Interpolation/extrapolation (what are the differences?)
 - Linear (LERP, BERP, TERP, SLERP)
 - Piecewise continuous
 - Splines, Bezier
 - Polynomial
 - Lagrange, etc

Integration strategies - Time/frequency/spatial domains

- We have data types such as structural scans of neural structure, EEG, MEG, fMRI, etc.
- How can these be synchronized spatially and temporally?
- What is an issue with spatial correlations (See A4!)?
- Mapping - coordinate, typically affine transformation
- Inverse computations - knowing locations of sensors relative to brain, can infer activation areas (localize)



Affine vs. Linear

- Can somebody explain the difference between ***linear*** and ***affine*** transformations?
- Requirements of linearity?

More on linearity vs. nonlinearity

- Power
 - **A linear system is a system whose dependent variables are related to its independent variables by a power of one**
- Linear systems have these particular properties (and they are very favorable)

- **Additive**

$$T[x_1(n) + x_2(n)] = T[x_1(n)] + T[x_2(n)]$$

- **Homogeneous**

$$T[cx(n)] = cT[x(n)]$$

- (<https://mathworld.wolfram.com/LinearSpace.html>, <https://mathworld.wolfram.com/LinearTransformation.html>)

Affine transformation

- Any transformation that preserves collinearity (i.e. points on a line remain on a line after the transformation) and ratio of distances (midpoint of a line before and after transformation remains the same)
- $y=mx+b$ is? **Affine**
- $y=mx$ is? **Linear**
- or more generally (see <https://mathworld.wolfram.com/AffineTransformation.html>, <https://mathworld.wolfram.com/AffineSpace.html>, <https://medium.com/mlait/affine-transformation-image-processing-in-tensorflow-part-1-df96256928a>)

Affine transformations in neural imaging

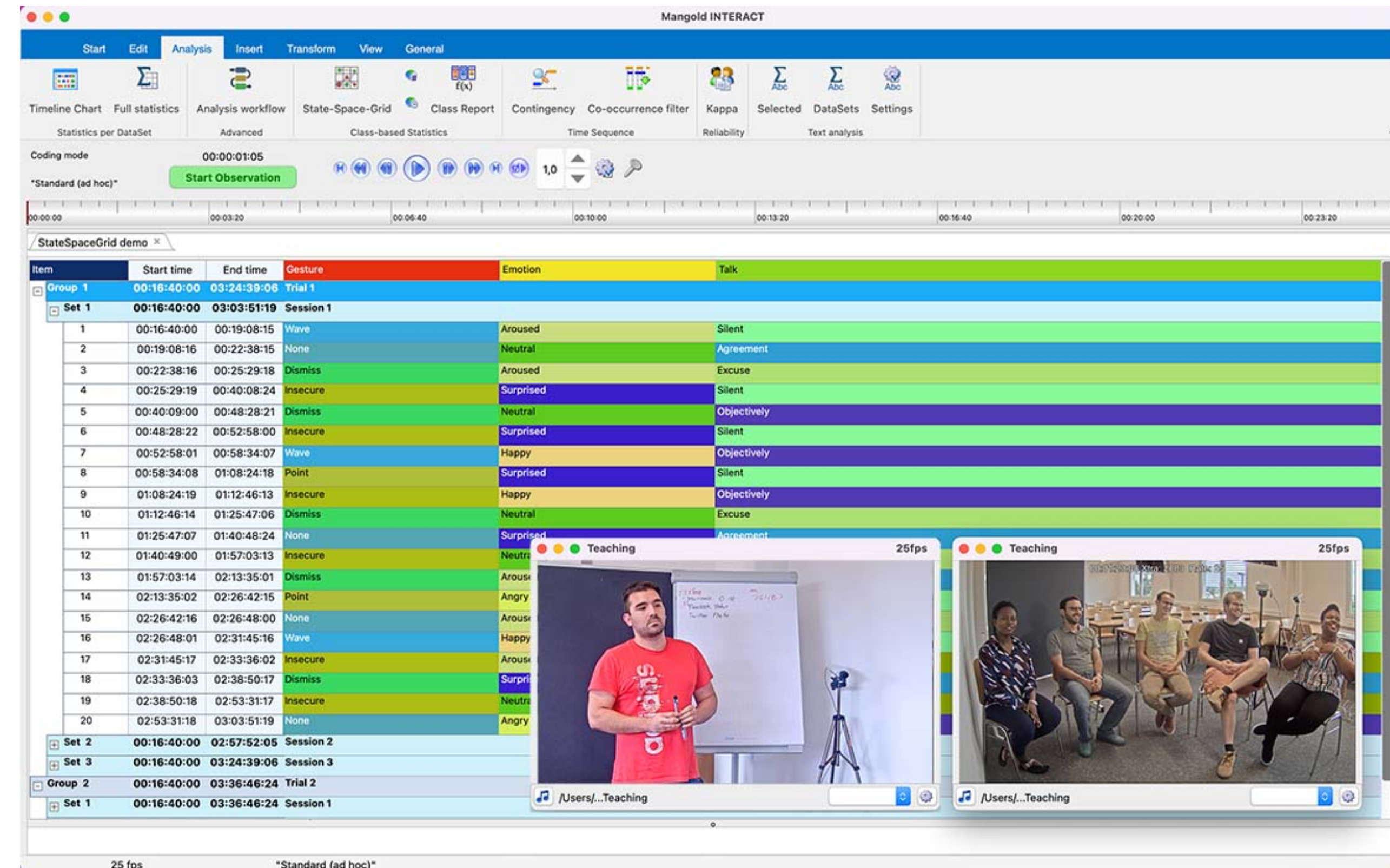
- Image processing - Correction of distortions and deformations (geometric) that occur from camera angles that are not optimal
- Brain imaging - transforming from sensor to brain coordinates, mapping different modalities, standardization for format
- Parallel lines to parallel lines
- e.g. Rotation, Translation, Scaling, Shear
- NiBabel documentation

Integration strategies - Sample time mismatch

- Super/sub sampling with filtering
- Time/sample shift to align data

Integration strategies - Software and format

- What if you have image/video data, EEG, text, audio?
- Each is in its own format, with different sample timings, not keyed events, coordinates, dataframes
- Traditional way?
- Newer way?



Integration strategies - Sample rate variability

- Do you have an accurate time measure and know the variability?
 - Yes - then you can simply interpolate and resample to create a new equally spaced set
- Inaccurate time measure, some information is lost
 - Computer timers for example do not provide accurate time measures unless they are specialized hardware
 - Can assume it's accurate if sampling much much faster than dynamics
 - Reduce sample rate (sub-sample) below estimated variability
 - Cannot use for time-critical associations

Integration strategies- missing data, mixture, non-tabular

- Addressed in earlier lectures (NANs)
- Wrangling
- Manual labor
- May need to use portions of the data
- Large sets need automated or semi-automated detection means

Integration strategies - Memory and processor usage

- Why do we need to be aware of this issue?
- Cloud computing services
- Efficient coding
- Considering data partially, in chunks, computed offline, pre-computed then processed as needed for analysis
- Variable sizes
- n-dimensions - what is the ***curse of dimensionality***?

A4: Integrating heterogeneous datasets for neuroscience

Modules for A4

- **nibabel**

- Neuroimaging in Python
- <https://nipy.org/nibabel/>

- **pysurfer**

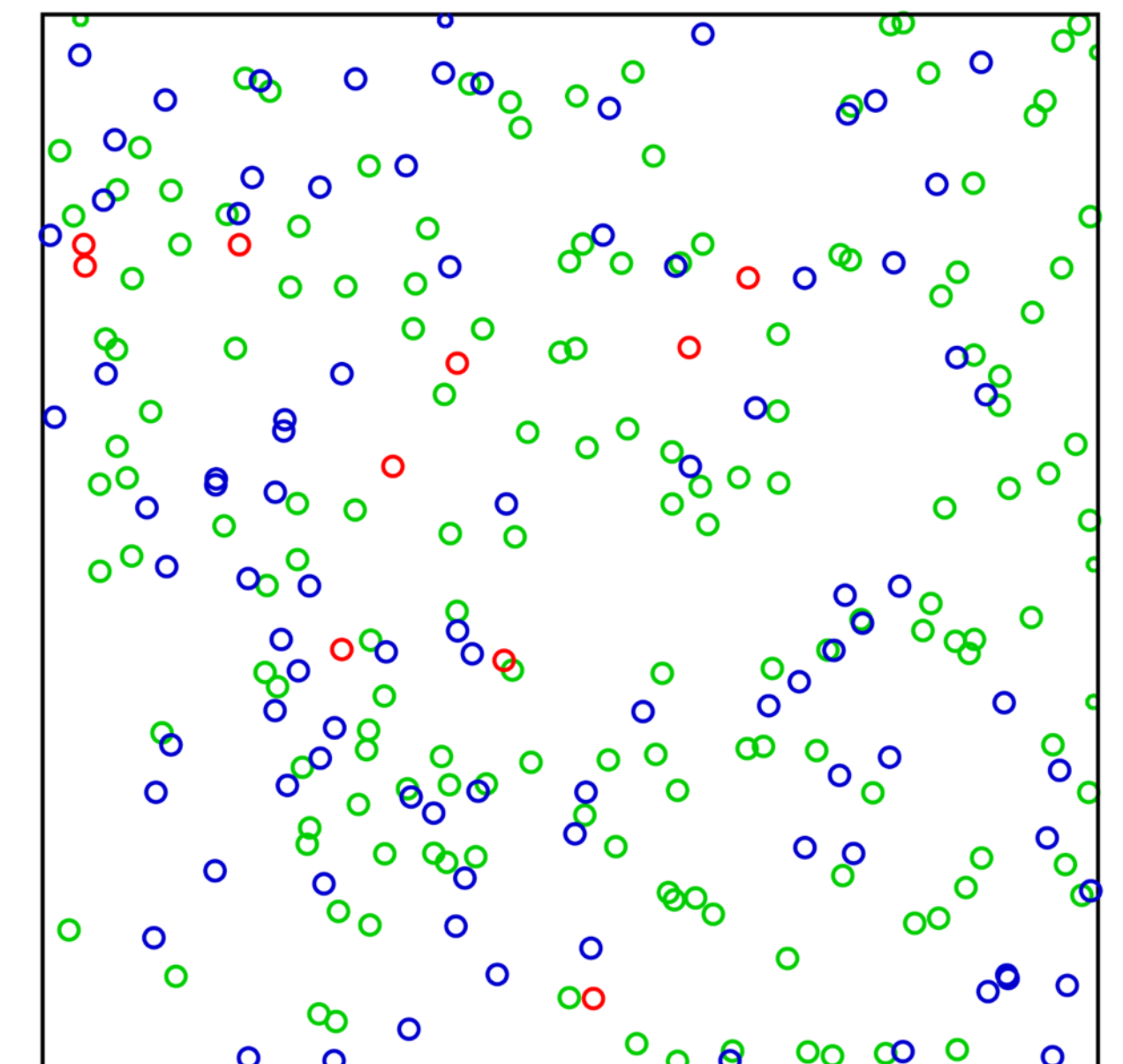
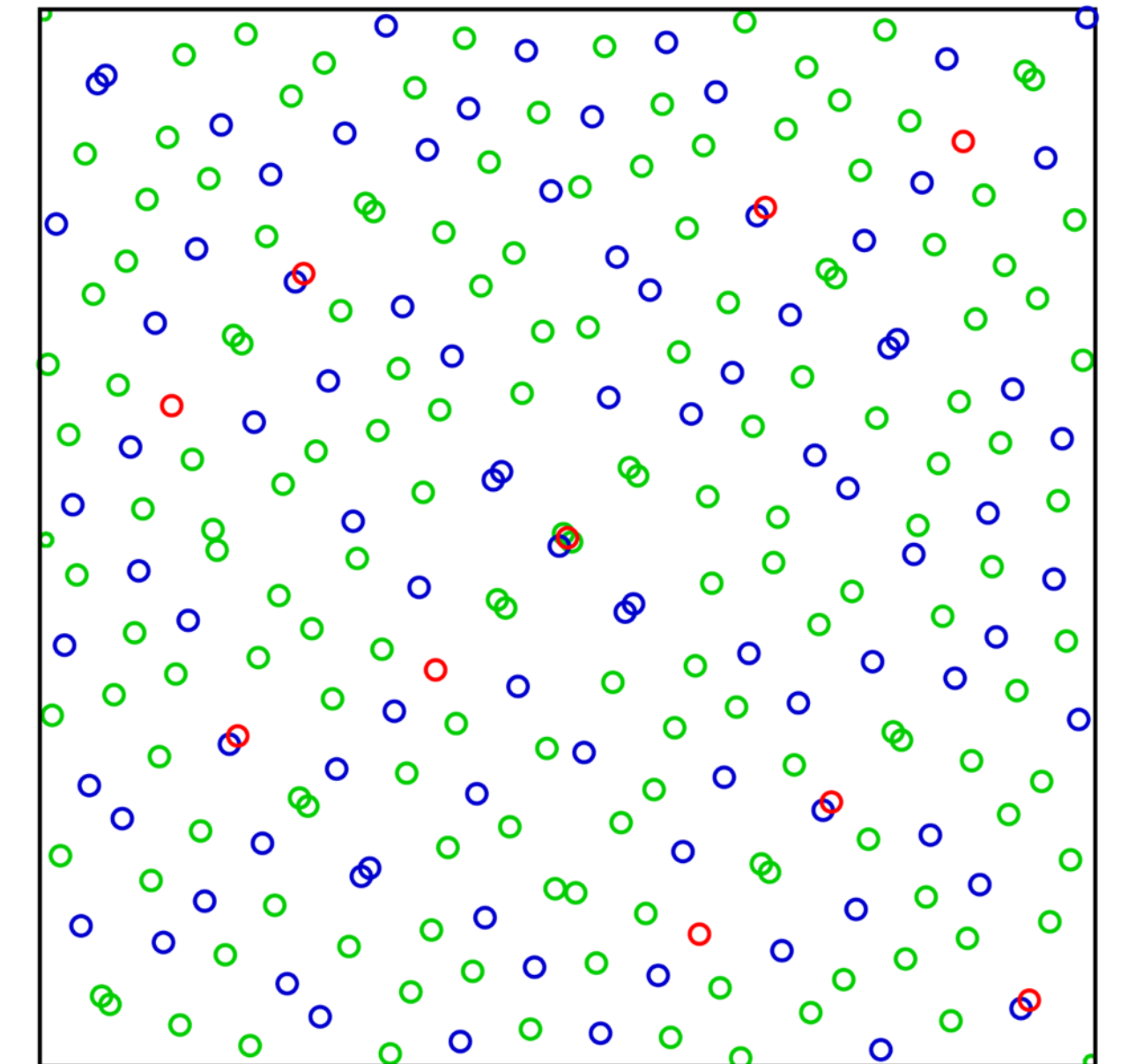
- Visualizing brain imaging data
- <https://pysurfer.github.io>

- **sobol_seq**

- Sobol sequence generator
- https://github.com/naught101/sobol_seq
- <https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.qmc.Sobol.html>

Sobol sequences

- Quasi-random low-discrepancy sequences
- https://en.wikipedia.org/wiki/Sobol_sequence
- Which one covers the space more evenly, just by eye?
 - Sobol or pseudorandom
- **Sobol sensitivity analysis** to analyze influence of parameters in computational neuroscience models
 - <https://hal.science/hal-03464025/file/root.pdf>
 - <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8184610/>
 - Model reproducibility





- PySurfer is a Python library for **visualizing cortical surface representations of neuroimaging data**.
- The package is primarily intended for use with [Freesurfer](#), but it can plot data that are drawn from a variety of sources.
- PySurfer extends [Mayavi's](#) powerful rendering engine with a high-level interface for working with MRI and MEG data.

pysurfer - installation

```
pip install pysurfer
```

Dependencies¶

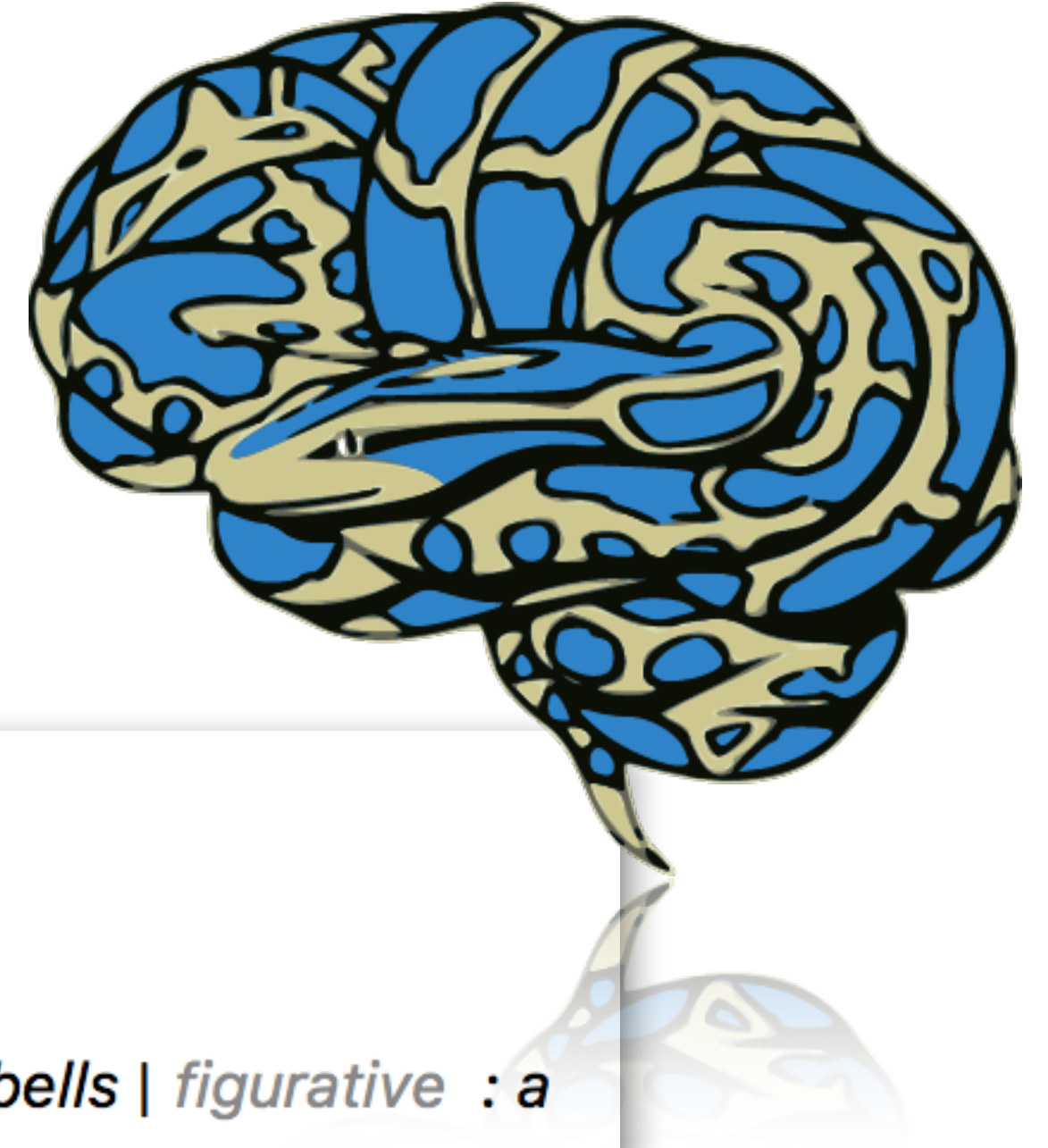
PySurfer works on Python 2.7 and 3.6+. (Older Python 3 versions will probably work, but are not tested.)

To use PySurfer, you will need to have the following Python packages:

- [numpy](#)
- [scipy](#)
- [nibabel](#)
- [mayavi](#)
- [matplotlib](#)

Some input/output functions also make use of the Python Imaging Library ([PIL](#)) and `imageio`, although they are not mandatory.

NiBabel - definition



- “Access a cacophony of neuro-imaging file formats”

- Cacophony?

cacophony | kə'käfənē |

noun (pl. **cacophonies**)

a harsh, discordant mixture of sounds: *a cacophony of deafening alarm bells* | *figurative* : *a cacophony of architectural styles* | *songs of unrelieved cacophony.*

- Read and write access to common neuroimaging file formats,
 - including: [ANALYZE](#) (plain, SPM99, SPM2 and later), [GIFTI](#), [NifTI1](#), [NifTI2](#), [CIFTI-2](#), [MINC1](#), [MINC2](#), [AFNI BRIK/HEAD](#), [ECAT](#) and Philips PAR/REC.
 - In addition, NiBabel also supports [FreeSurfer](#)'s [MGH](#), geometry, annotation and morphometry files,
 - provides some limited support for [DICOM](#)

NiBabel - Installation



```
pip install nibabel
```

NiBabel - documentation

- Coordinate systems
- Radiological vs. Neurological conventions
- Intro to DICOM

A4 - Mapping heterogeneous neural data

- How to take different neural data and map them to the human neocortex
- https://en.wikipedia.org/wiki/Human_Connectome_Project
- “A multi-modal parcellation of human cerebral cortex”
 - <https://pubmed.ncbi.nlm.nih.gov/27437579/>

In class report development (~30m)

- Define this course's intent
- Draw comparisons between this course and requirements
- How does this course build upon what came before?
- How can you use your starting point in this course to expand your understanding?