
HyperTools: A Python toolbox for visualizing and manipulating high-dimensional data

Andrew C. Heusser[†], Kirsten Ziman[†], Lucy L. W. Owen, and Jeremy R. Manning[‡],

**Department of Psychological and Brain Sciences,
Dartmouth College, Hanover, NH 03755**

[†] Denotes equal contribution

[‡] Address correspondence to jeremy.r.manning@dartmouth.edu

Abstract

Data visualizations can reveal trends and patterns that are not otherwise obvious from the raw data or summary statistics. While visualizing low-dimensional data is relatively straightforward (for example, plotting the change in a variable over time as (x, y) coordinates on a graph), it is not always obvious how to visualize high-dimensional datasets in a similarly intuitive way. Here we present `HyperTools`, a Python toolbox for visualizing and manipulating large, high-dimensional datasets. Our primary approach is to use dimensionality reduction techniques [14, 17] to embed high-dimensional datasets in a lower-dimensional space, and plot the data using a simple (yet powerful) API with many options for data manipulation (e.g. hyperalignment [8], clustering, normalizing, etc.) and plot styling. The toolbox is designed around the notion of *data trajectories* and *point clouds*. Just as the position of an object moving through space can be visualized as a 3D trajectory, `HyperTools` uses dimensionality reduction algorithms to create similar 2D and 3D trajectories for time series of high-dimensional observations. The trajectories may be plotted as interactive static plots or visualized as animations. These same dimensionality reduction and alignment algorithms can also reveal structure in static datasets (e.g. collections of observations or attributes). We present several examples showcasing how using our toolbox to explore data through trajectories and low-dimensional embeddings can reveal deep insights into datasets across a wide variety of domains.

Introduction

“To deal with hyper-planes in a fourteen dimensional space, visualize a 3D space and say ‘fourteen’ to yourself very loudly. Everyone does it.”

—Geoffrey Hinton [9]

1

2

3

4

5

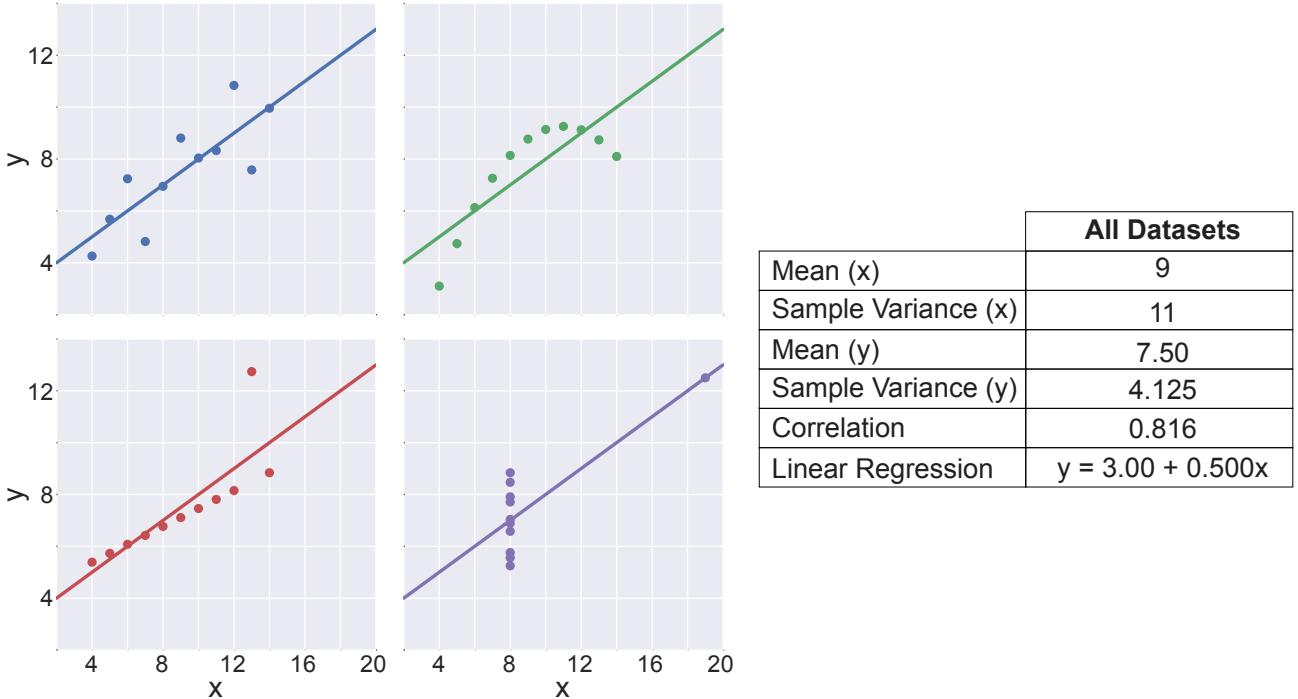


Figure 1. Anscombe’s Quartet. Each dataset shares the same descriptive statistics but exhibits a unique shape.

The **HyperTools** toolbox is designed to reveal geometric structure in high-dimensional data through visualizations and manipulations. Modern data visualizations date back to at least the 16th century, when early data pioneers began to develop the sorts of accurate maps and diagrams we might still recognize today [5, 19]. Visualizations can reveal deep insights and intuitions about geometric structure and patterns in complex datasets by capitalizing on the human visual system’s ability to quickly and efficiently extract meaning and structure from highly complex visual information [20]. This is perhaps especially true of high-dimensional datasets, where different dimensions or features may interact in complex ways that may not be immediately obvious through conventional summary statistics.

As an illustration of the potential for summary statistics to mislead in the absence of visualization, consider the classic example, Anscombe’s Quartet [1] (Fig. 1). Anscombe’s Quartet comprises four datasets that share a common statistical profile. Because the datasets are exactly equal along several common summary measures (mean, variance, trend lines), at first glance they seem highly similar. However, plotting the datasets and comparing them visually reveals that they differ substantially in structure. Whereas low-dimensional datasets like those in Anscombe’s Quartet can be easily plotted, it is not always obvious how to visualize high-dimensional datasets (e.g. with greater than 3 dimensions) in a similarly intuitive way.

A number of techniques, collectively referred to as *dimensionality reduction algorithms* have been developed over the past half-century to map high-dimensional data onto lower-dimensional representations that may be more easily manipulated and visualized. Some well-known examples include Principal Components Analysis (PCA) [14], Probabilistic Principal Components Anal-

ysis (PPCA) [17], Independent Components Analysis (ICA) [4, 11], Multidimensional Scaling (MDS) [18], and t -Distributed Stochastic Neighbor Embedding (t-SNE) [21]. While the details of these algorithms differ, they each provide a means of obtaining a low-dimensional representation of the original high-dimensional dataset that preserves many of the geometric properties (e.g. the overall covariance structure of the data, data grouping, etc.) to the extent possible within a low-dimensional space. In the **HyperTools** toolbox, we leverage these dimensionality reduction algorithms (Fig. 2a) to aide in visualizing high-dimensional data.

A second class of algorithms leveraged in our toolbox provide techniques for manipulating and aligning different high-dimensional datasets (Fig. 2b). These algorithms draw inspiration from the Procrustean transformation [16], which computes the affine transformations (i.e. translation, reflection, rotation, and scaling) that bring one trajectory into alignment with another (in terms of minimizing the mean squared Euclidean distances between the corresponding points). The hyperalignment algorithm [8] and the Shared Response Model (SRM) [3] extend this technique to find a common set of transformations that bring many (more than two) high-dimensional trajectories into common alignment. Our **HyperTools** toolbox leverages these alignment algorithms to allow users to manipulate and compare high-dimensional data, even when the dimensions (features) of original observations are different (e.g. brain patterns from different people, observations from different modalities, etc.).

Taken together, the **HyperTools** toolbox provides a set of powerful functions for visualizing and manipulating high-dimensional data using dimensionality reduction and data alignment algorithms. The toolbox is designed with ease of use as a primary goal, such that complex visualizations and analyses may often be carried out with a single line of code. Another major goal is to enable users to easily produce visually appealing publication-quality plots, also often with only a single line of code. Our toolbox is open-source and is distributed with the MIT License.

In the next section we provide a detailed overview of the components of the **HyperTools** toolbox and describe how the codebase is organized. We then describe a series of analyses of datasets from a wide array of domains to highlight many of the main functions of the toolbox.

Materials and Methods

Overview

The **HyperTools** toolbox is written in Python and can be downloaded from our GitHub page or with pip:

```
pip install hypertools
```

 (1)

HyperTools depends on the following open-source software packages: **Matplotlib** [10] for plotting functionality, **Seaborn** [23] for plot styling, **scikit-learn** [15] for data manipulation (dimensionality reduction, clustering, etc.), and PPCA for inferring missing data using PPCA. The toolbox also includes a port of the hyperalignment algorithm [8] from the **PyMVPA** library, as well as the shared response model from the **BrainIAK** toolbox, as an alternative data alignment

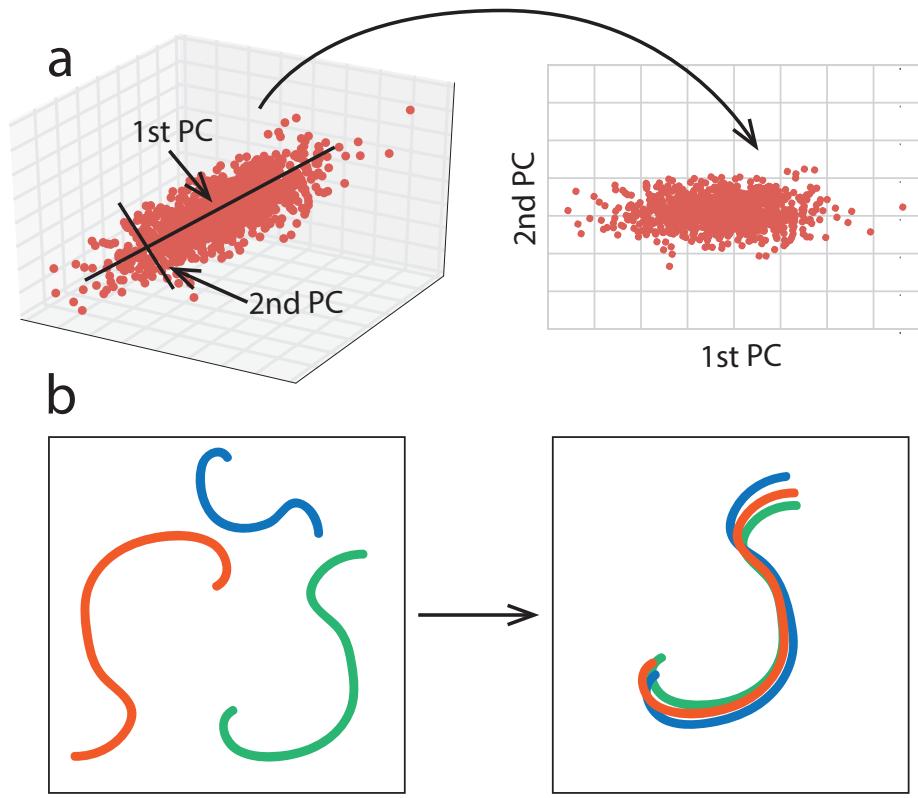


Figure 2. Visualizing and manipulating high-dimensional data. **a.** HyperTools uses dimensionality reduction algorithms to project high-dimensional data onto 2D and 3D plots. As shown in the panel, the dimensionality reduction algorithm PCA may be used to find the axes that explain the most variance in the original data (left panel). The data may then be projected onto a small number of those axes to facilitate plotting (right panel). **b.** Another important feature of HyperTools concerns aligning datasets with different fundamental coordinate systems. The left panel displays three trajectories with similar geometries but different coordinate systems, and the right panel displays how those trajectories may be aligned (via linear transformations) into a common space using hyperalignment.

Filename	Description
plot/plot.py	Main plotting function: parses arguments, dispatches to <code>static.py</code> and <code>animate.py</code>
plot/static.py	Handles all static plot logic
plot/animate.py	Handles all animated plot logic
tools/align.py	Aligns the coordinate space of a list of matrices using hyperalignment
tools/cluster.py	Parcellates observations into discrete clusters using k -means clustering
tools/describe_pca.py	Analyzes and plots how many principal components are needed to capture the covariance structure of the data
tools/missing_inds.py	Find nans in data and returns indices
tools/normalize.py	z -scores rows/columns of matrices
tools/df2mat.py	Converts <code>Pandas</code> dataframes to <code>Numpy</code> arrays
tools/procrustes.py	Aligns the coordinate spaces of two arrays
tools/reduce.py	Reduces the dimensionality of one or more arrays using PCA and PPCA
_externals/srm.py	Implements the Shared Response Model (alternative alignment algorithm)
_shared/helpers.py	Collection of helper functions used across many files

Table 1. `HyperTools` code organization. The table lists the main files and functions that comprise the toolbox. We provide a feature complete description of the API on the project’s GitHub page and in the documentation included with the toolbox download.

technique. In addition to providing a simple interface to several functions from these libraries, `HyperTools` adds a number of custom arguments to facilitate data visualization and manipulation of high-dimensional data. Table 1 provides summary of the `HyperTools` code base. In the remainder of this section, we provide descriptions of the primary toolbox functions, but we have not provided an exhaustive list. A feature-complete description of the API may be found on the project’s GitHub page and in the documentation included with the toolbox download.

Nearly all of the `HyperTools` functions may be accessed through the main `plot` function. This design enables complex data analysis, data manipulation, and plotting to be carried out in a single function call. There are two general types of plots supported by the toolbox: *static plots* and *animated plots*.

Static Plots

Accessing the `HyperTools` plot functionality entails first loading the to-be-analyzed dataset into the Python workspace and converting it to a `Numpy` array [22] or a `Pandas` dataframe [13]. The format of the data should be samples (S) by features (F). Once the dataset conforms to this

62
63
64
65
66
67
68
69
70
71

72

format, simply import the library and call the plot function:

```
import hypertools as hyp
```

(2)

```
hyp.plot(data)
```

(3)

By default (i.e. with no additional arguments specified), this function will perform dimensionality reduction (using PCA), convert the $S \times F$ `data` matrix into an $S \times 3$ matrix, and then create an interactive 3D line plot that can be explored by using the mouse to rotate the plot. If there are `nans` present in the dataset, these missing values will be automatically interpolated using PPCA [17]. If $F < 3$, a 2D plot is created instead of a 3D plot. This simple interface to plotting is deceptively powerful: with a single command, the toolbox automatically fills in missing data and determines whether to create a 2D or 3D plot (reducing the dimensionality of the observations as needed).

`HyperTools` can also accommodate lists of Numpy arrays or Pandas dataframes (only single-level indexed dataframes are currently supported):

```
hyp.plot([array1, array2, array3])
```

(4)

When passed a list of arrays, `HyperTools` will plot each array in a distinct color. Colors and styling can be customized in several ways. Like `Matplotlib`, `HyperTools` can parse format strings passed as positional arguments. For example:

```
hyp.plot(array1, 'k-')
```

(5)

```
hyp.plot([array1, array2, array3], ['bo', 'r--', 'g:'])
```

(6)

Line colors may also be specified via the `color` (or `colors`) keyword argument:

```
hyp.plot(array1, color='g')
```

(7)

```
hyp.plot([array1, array2, array3], colors=['b', '#FF0000', (.3, .5, .4)])
```

(8)

Colors may be defined using format strings, hex codes, RGB values, or a mix of these formats. Rather than specifying the colors of each data array, colors may instead be specified for each individual sample by providing labels for each sample:

```
hyp.plot(data, group=group_labels),
```

(9)

where `group_labels` is a list of length S (number of samples). (Lists of group label arrays are also supported, e.g. if the data are passed in as a list of arrays; the list of labels must be of the same length as the list of data arrays.)

`HyperTools` parses this function call by sub-dividing each data matrix into new lists defined by each unique label in `group_labels`. For example, if each sample label is a string from the set ('a', 'b', 'c') then each of these unique labels will be assigned a unique color, and the datapoints assigned to each label will be assigned that label's color.

In addition to specifying string labels for each sample, `HyperTools` also supports numerical labeling. If `group_labels` is a list of numbers, `HyperTools` will bin the range covered by those

numerical values (excluding `nans`, `Nones`, and `infs`) into n evenly spaced bins (default: $n = 100$) and map these values onto a color palette. The color palette used for this mapping may also be customized using the `palette` keyword:

```
hyp.plot([array1, array2], palette='husl')
```

 (10)

All Matplotlib and Seaborn color palettes and plot styles are supported by HyperTools.

In addition to specifying group-level labels (which are used to determine the colors of each sample), each sample may also be labeled with an additional text label that may be shown (as text) on the plot. The `label` and `labels` keyword arguments allow the user to define a list of strings (or a list of lists of strings) to be displayed next to each sample datapoint with an arrow pointing to it. Each list of labels must be of length S (number of samples). (The `None` value may be used to specify “blank” labels, which will not show up on the plot.)

By default, all datapoint labels are shown if the `label` or `labels` keyword is specified. However, HyperTools also supports a “data exploration” mode whereby the datapoint labels will only be shown when the mouse pointer hovers over the corresponding datapoint:

```
hyp.plot(data, labels=['a', None, 'a', 'b'], explore=True)
```

 (11)

This plotting mode is especially useful when there are many datapoints, or when the data labels are long. If `explore` is set to `True` and no labels are specified, the labels will be auto-generated as an index and the PCA coordinate (e.g. '45: (3.0, 4.0, 5.0)'). Note: at the time of this writing, the `labels` and `explore` arguments are only supported for 3D static plots.

Animated Plots

Animated 3D plots are especially useful for visualizing high-dimensional timeseries data. To create an animation, simply toggle the `animate` keyword:

```
hyp.plot(data, animate=True)
```

 (12)

This will create a 3D animated representation of the data, where the animation occurs over the rows of the `data` matrix. As with static plots, the user may pass a list of data matrices to plot multiple datasets on a single plot, and format strings and keyword arguments may be used to customize the plot appearance. Each frame of the animation displays a portion of the total data trajectory enclosed in a cube (Fig. 3). In successive frames, the displayed portion of the data trajectory advances by a small amount, and the camera angle rotates around the cube, providing visual access to different aspects of the data as the animation progresses.

The formats of animated HyperTools plots may be customized using the following keyword arguments: `duration` specifies the animation duration in seconds, `tail_duration` specifies the duration of the trailing tail in seconds, `rotations` specifies the number of camera rotations around the data (over the course of the entire animation), `zoom` will zoom the camera in (positive number) or out (negative number) from the data, and setting `chemtrails=True` will plot a low opacity version of samples prior to the currently active trajectory so that the full structure and history of the data may be visualized. For a complete list of animation-specific arguments,

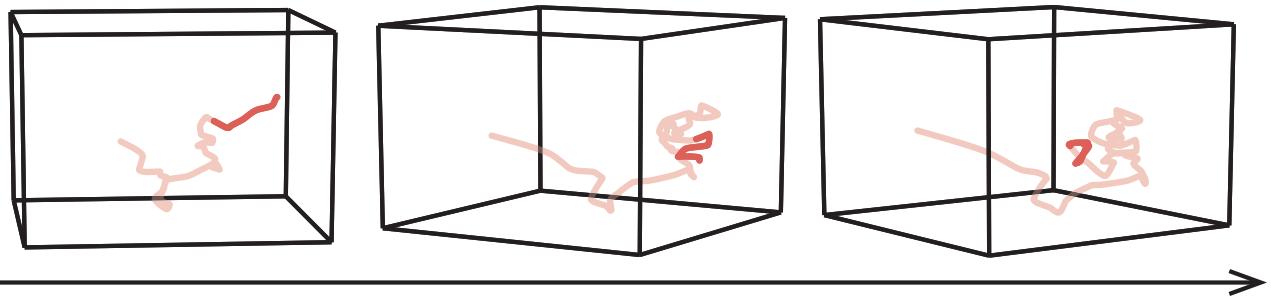


Figure 3. Frames from an animated plot. Three frames (with time increasing moving from left to right) from an example animation are displayed in each panel.

please see the API documentation. Both animated and static plots can be saved by including the `save_path` argument (with the file extension included):

```
hyp.plot(data, save_path='path/to/the/file.pdf')
```

(13)

```
hyp.plot(data, animate=True, save_path='path/to/the/file.mp4'),
```

(14)

where Snippet 13 saves a static plot to a resolution-independent PDF, and Snippet 14 saves an animation as an MP4 movie. Note that saving animated plots requires that `ffmpeg` is installed on your computer; see the API documentation for more details.

107
108
109

110

Reduce

When passed high-dimensional data, the `plot` function uses PPCA to fill in missing data and PCA to project the data onto 3 dimensions. We provide API access to the `reduce` function that underlies these transformations. At its core, the `reduce` function is primarily a wrapper for `scikit-learn`'s PCA function and the PPCA package. `HyperTools` extends the functionality of these tools by providing an easier-to-use syntax and adding support for lists of matrices. The function may be used as follows:

```
reduced_data = hyp.tools.reduce(data, ndims=3)
```

(15)

Because dimensionality reduction results in information loss (relative to the original dataset), it is important to consider how accurately a low-dimensional projection of the data reflects the original high-dimensional dataset. `HyperTools` includes a function that plots the correlation between the covariance matrices of the reduced and full datasets, as function of the number of principal components in the reduced dataset (Fig. 4):

```
fig, ax, data = hyp.tools.describe_pca(data)
```

(16)

The `describe_pca` function computes these correlations iteratively (i.e. starting with one principal component, then two, then three, etc.) until a local maximum is detected. The resulting plot provides insights into the increase in explanatory power (in terms of the across-sample covariance) associated with each new principal component.

111
112
113
114

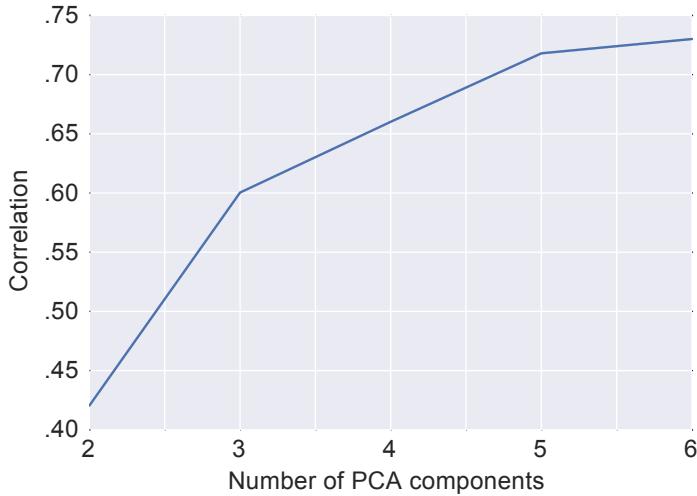


Figure 4. Covariance preserved as a function of the number of principal components. For an example dataset (Example 2, *Results*) the panel displays the correlation between the upper triangle of the across-sample covariance matrices of the reduced versus original data, as a function of the number of principal components.

Align

115

Two or more datasets may share geometrical structure, but reside in different coordinate systems. Hyperalignment is a method that aligns the representational spaces over a list of datasets, effectively co-registering them to a common space [8]. Using linear transformations, hyperalignment find a common space that minimizes the distance between two or more datasets (Fig. 2b). Aligning them to a common space allows one to visualize commonalities between the two different kinds of data. The `align` function accepts a list of arrays as input and returns a hyperaligned list of arrays in a common geometric space:

```
hyperaligned_list = hyp.tools.align([array1, array2, array3]) (17)
```

In addition to supporting alignment via the hyperalignment algorithm proposed by [8], we have also added support for alignment via the Shared Response Model [3], which was ported from the `BrainIAK` toolbox:

```
SRM_aligned_list = hyp.tools.align([array1, array2, array3], method='SRM') (18)
```

Cluster

116

Some datasets exhibit *clustering* tendencies, whereby the data may be divided into discrete groups of similar or related samples (i.e. samples that are comprised of similar features). When these discrete groups are unlabeled or unknown, clustering algorithms provide heuristics for recovering these clusters of similar samples automatically. `HyperTools` incorporates the *k*-means clustering algorithm [7] to facilitate automatic data clustering. Given a pre-chosen number of

clusters, k , the `cluster` keyword argument to the `plot` function uses k -means clustering to automatically assign each observation to a cluster, and then colors each observation's point according to its cluster membership:

```
hyp.plot(data, n_clusters=k) (19)
```

We also expose the k -means clustering algorithm directly through the `cluster` function:

```
cluster_labels = hyp.tools.cluster(data, n_clusters=k) (20)
```

The `cluster` function wraps the `scikit-learn` implementation of k -means clustering and extends it to work with lists of data matrices.

Results

Example 1: Visualizing hypercubes in 3D

To illustrate how a user might visualize high-dimensional data with `HyperTools`, we start by examining four synthetic datasets with unique, known structures. We generated datasets of one cube (3 dimensions) and three hypercubes of increasing dimensionality (4, 5 and 6 dimensions), each comprised of 100 points along each of their respective edges. We then used `HyperTools` to project the hypercubes into 3 dimensional space (using PCA) and visualize the result.

Figure 5 illustrates how projecting hypercubes of different dimensionalities into 3 dimensional space distorts some aspects of their shapes, while preserving others. In the original (high-dimensional) data, all edges of each respective cube are of equal length, and each vertex comprises n adjacent edges converging orthogonally (where n is the dimensionality of the hypercube). However, in Fig. 5, some edges appear longer than others, and some vertices appear to form acute and obtuse angles.

Despite these differences, many of the underlying structural components are accurately reflected in the visualization. Namely, the visualization of each n -dimensional cube correctly depict 2^n vertices, $2^{n-1} * n$ edges, and n edges converging at each vertex. Each edge is also reliably reconstructed as a straight (rather than curved) line segment. The visualizations also depict increasing complexity with increasing dimensionality.

Example 2: Dimensionality reduction and clustering with various types of mushrooms

In this section, we highlight the dimensionality reduction and clustering capabilities of `HyperTools`. We retrieved the 'mushroom classification' dataset from the Kaggle database. The dataset contains annotated descriptive features of 8,124 mushrooms spanning 23 mushroom species from the *Audubon Society Field Guide to North American Mushrooms* [12]. Each observation comprises a list of 22 descriptive features (e.g. cap shape, cap surface, habitat, etc.) along with a tag identifying each mushroom exemplar as poisonous or non-poisonous (features for five example mushrooms are shown in Tab. 2).

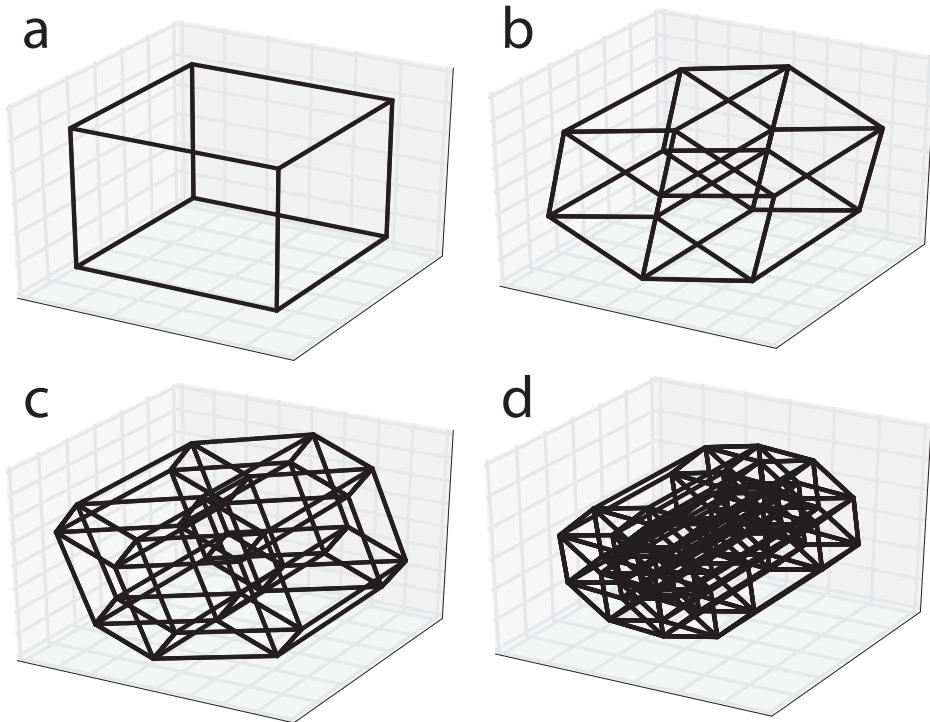


Figure 5. Hypercubes with increasing dimensionality. Each dataset comprises 100 evenly spaced points along each edge of the corresponding cube with dimensionality **a.** 3, **b.** 4, **c.** 5, and **d.** 6.

	class	cap-shape	cap-surface	cap-color	bruises	odor	...	habitat
0	p	x	s	n	t	p	...	u
1	e	x	s	y	t	a	...	g
2	e	b	s	w	t	l	...	m
3	p	x	y	w	t	p	...	u
4	e	x	s	g	f	n	...	g

Table 2. Example of mushrooms dataset. The dataset contains annotated features (columns) of each mushroom (row), along with labels indicating whether each mushroom is poisonous or non-poisonous (not shown).

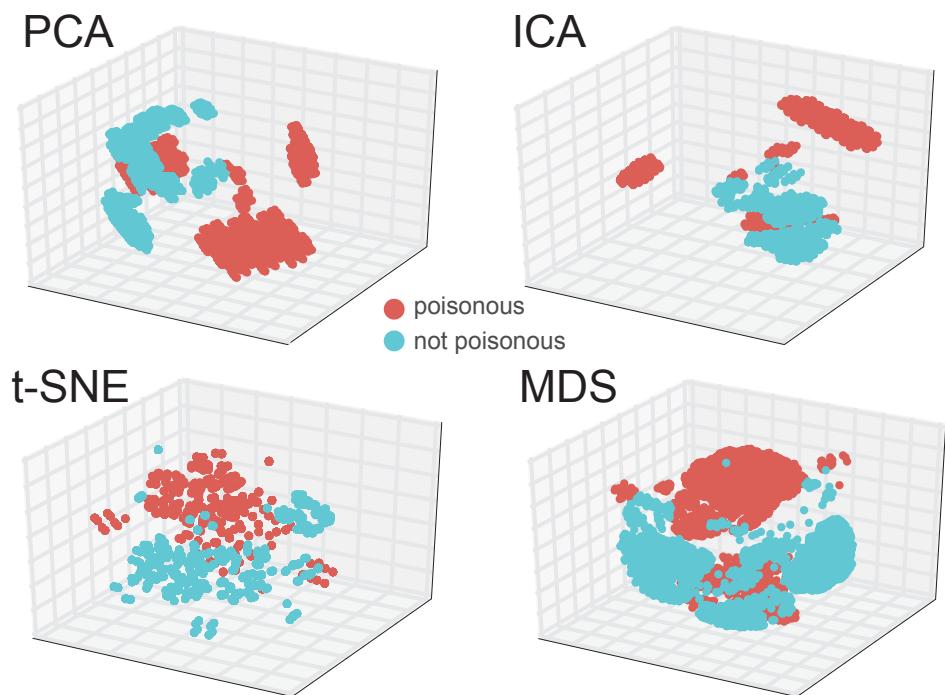


Figure 6. Three-dimensional embeddings of the mushrooms dataset using several dimensionality reduction techniques. Each point represents a sample (mushroom). Red dots indicate poisonous mushrooms and blue indicate non-poisonous mushrooms.

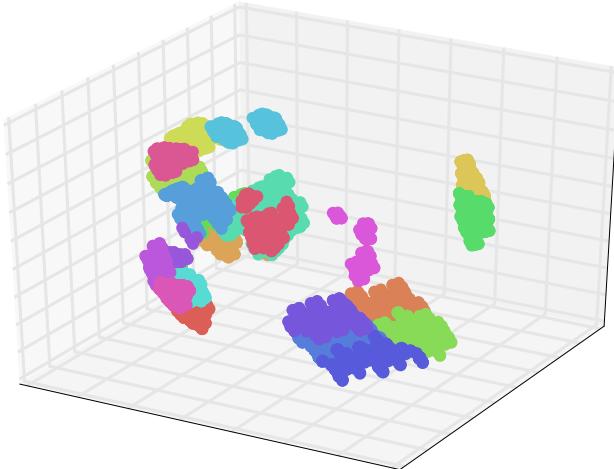


Figure 7. Mushrooms dataset, colored by k -means cluster.

Because the mushroom features are provided as character strings, they must be transformed into numerical vectors to plot them. When passed a `Pandas` dataframe with columns containing text, `HyperTools` automatically converts the data into a binary matrix, where each column reflects one of the unique values of one of the features. The underlying function for converting dataframes into matrices may also be called directly:

```
matrix = hyp.tools.df2mat(dataframe)                                (21)
```

Plotting the resulting matrix with `HyperTools` reveals a striking clustered structure. Overall, ¹⁴⁶ the samples appear to cluster by whether or not they are poisonous, but they also appear to ¹⁴⁷ group into sub-clusters (Fig. 6). By default, `HyperTools` uses PCA for dimensionality reduction, ¹⁴⁸ but different dimensionality reduction techniques can reveal distinct geometrical properties of a ¹⁴⁹ dataset. To highlight this, we plotted the transformed binary matrix using several dimensionality ¹⁵⁰ reduction techniques (PCA, ICA, t-SNE, and MDS) to visualize their effects on clustering (Fig. 6). ¹⁵¹ Each technique produces a unique low-dimensional projection of the data, highlighting distinct ¹⁵² structural aspects. ¹⁵³

To highlight the sub-clustering structure in this dataset, we use the `n_clusters` argument to ¹⁵⁴ plot:

```
hyp.plot(mushrooms_data, n_clusters = 23)                            (22)
```

This command relies on k -means clustering to empirically derive cluster labels, and then plot ¹⁵⁴ each cluster in a different color (Fig. 7). ¹⁵⁵

Example 3: Exploring factors that influence educational outcomes. ¹⁵⁶

Next, we analyzed an education dataset containing, for each of 480 students around the world, ¹⁵⁷ performance ratings (high, medium, and low performance), demographic descriptors (e.g. gender, ¹⁵⁸ nationality, place of birth, etc.) as well as classroom behaviors (number of times the student ¹⁵⁹

	gender	Nationality	PlaceofBirth	StageID	GradeID	...	Class
0	M	KW	KuwaIT	lowerlevel	G-04	...	M
1	M	KW	KuwaIT	lowerlevel	G-04	...	M
2	M	KW	KuwaIT	lowerlevel	G-04	...	L
3	M	KW	KuwaIT	lowerlevel	G-04	...	L
4	M	KW	KuwaIT	lowerlevel	G-04	...	M

Table 3. Example features in education dataset. The dataset contained categorical and numerical features, as well as student performance labels.

raised their hand, days absent from class, number of times the student visited online resources, etc.) and others (features for five example students are displayed in Tab. 3; for full list of features and to download the data, see the Kaggle database). Given a dataframe with the student features, `HyperTools` automatically converts this into a binary data matrix (as described above) for visualization.

In contrast to the mushroom dataset, where the samples formed clear clusters, the distribution of samples in this dataset appear to form a single contiguous mass. Further, coloring each sample (student) by their performance rating reveals a striking correspondence between the student’s attributes and performance ratings (Fig. 8a). For example, as the student attributes vary along the first principal component, the student performance ratings appear to transition smoothly from low, to medium, to high (Fig. 8b). To highlight this pattern, we fit a linear regression model whose output variable was student performance and the input variables were the first three principal components (Fig. 8c). In this way, the regression model’s outputs provide a continuous estimate of student performance, whereas the original data contained only discrete (categorized) estimates. In Figure 8d, each dot from Panel a has been re-colored according to the regression model’s performance predictions, resulting in a smooth gradient from low to high performance.

Example 4: Exploring linguistic data from presidential nominees’ Twitter posts.

Whereas the above examples illustrate how simple numerical and categorical features are processed by `HyperTools` to reveal geometric patterns in the data, we can use a similar approach to extract and visualize more complex features. For example, topic models [2] may be used to derive a vector representation of each document in a corpus according to its linguistic properties. Specifically, topic models identify “themes” that are reflected in varying amounts by different documents in the corpus, where each theme (*topic*) is defined formally as a distribution over words in the vocabulary. In other words, a neuroscience-themed topic might heavily weight words like *neuron* and *brain*, whereas a sports-themed topic might heavily weight words like *running* and *athlete*. (Fitting a topic model to a text corpus reveals what the specific topics are and how much each document reflects each topic.) Once we have derived topic vectors for each document in the corpus, we can use `HyperTools` to visualize the full corpus to potentially gain insights into its geometric structure.

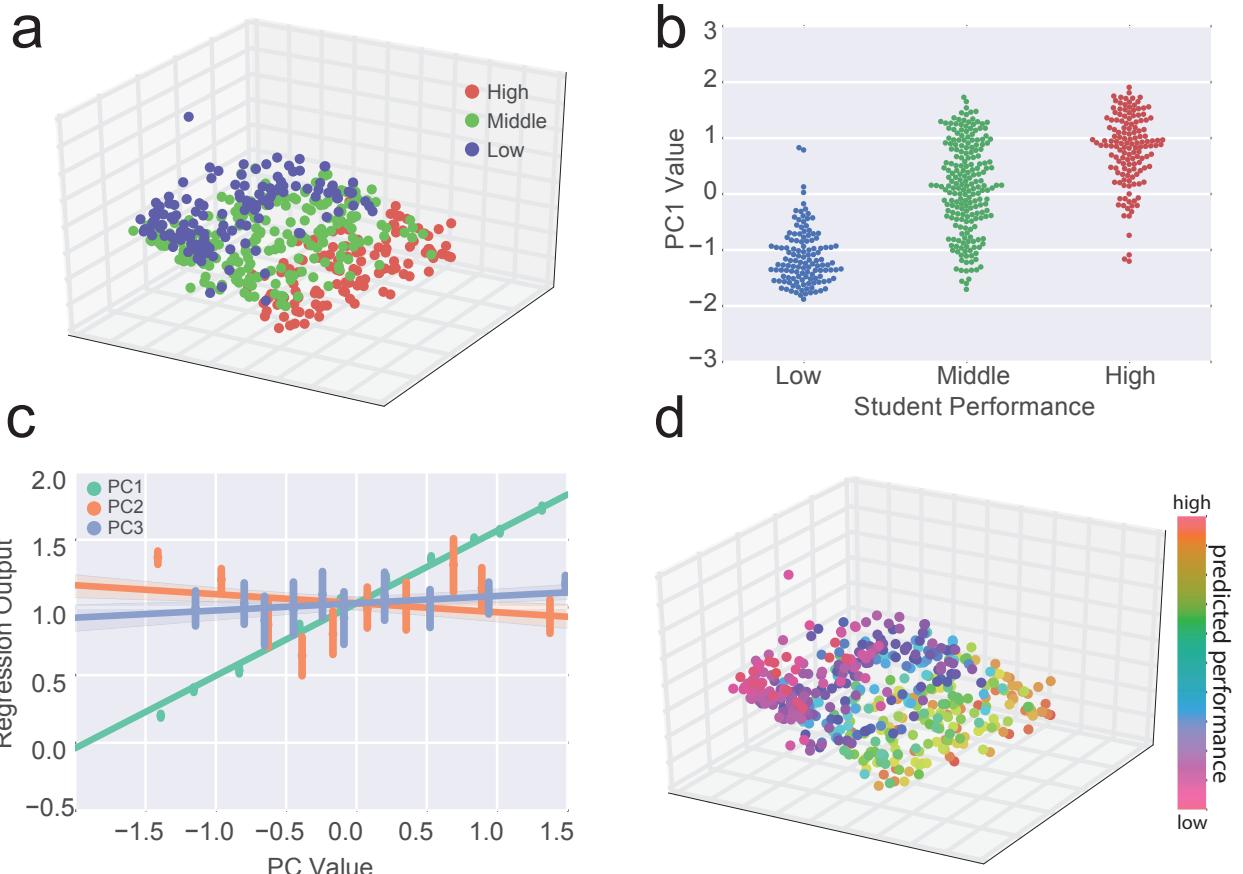


Figure 8. Relationship between student attributes and performance. **a.** Each point represents the feature vector associated with a student, and the points are colored by student performance (red: high, green: middle, blue: low). **b.** Swarm plot of the first principal component split by student performance (coloring same as above). **c.** Predicted student performance from linear regression of each PC on student performance. **d.** Same as (a), but points are colored by linear regression predictions of student performance by PC1 value (red to violet gradient represents high to low predicted performance).

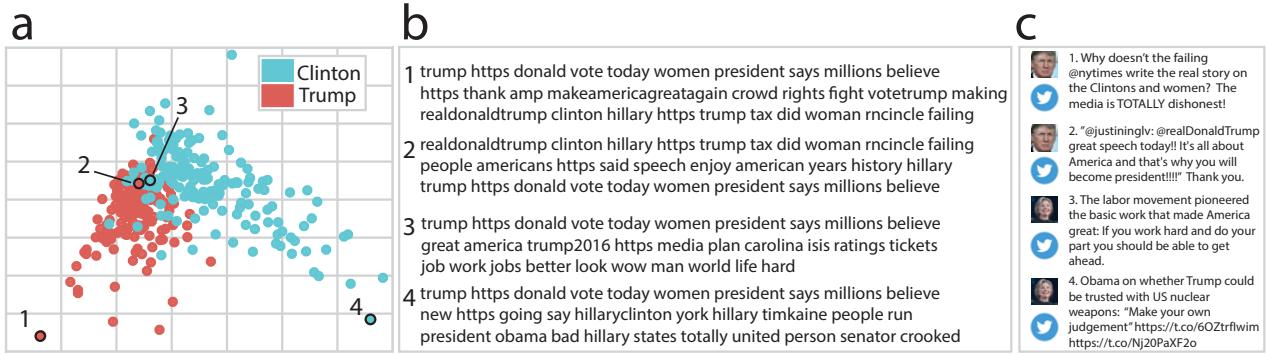


Figure 9. Topic models of political Twitter data. a. Two-dimensional representation of Clinton's (blue) and Trump's (red) day-by-day tweet content. **b.** Top ten words from each of the top three topics on selected days. **c.** Representative tweets from the selected days.

As an example of this approach, we next turn to an analysis of Twitter data ("tweets") from the Twitter accounts of Hillary Clinton (@HillaryClinton) and Donald Trump (@realDonaldTrump) over the course of their 2016 political campaigns. The dataset, sourced from FiveThirtyEight, contains 6,444 tweets sent from the candidates' primary Twitter accounts between April 17, 2016 and September 26, 2016.

We began our analysis by fitting a 20-topic topic model to the entire collection of tweets from both candidates, yielding a single topic vector for each tweet. Separately for each candidate, we next computed daily average topic vectors over the six month interval covered by the dataset, and we used **HyperTools** to visualize the resulting day-by-day Twitter topics.

Plotting the candidates' tweet content in two dimensions reveals that Clinton's and Trump's tweets were primarily about different topics, resulting in a V-like topic cloud (Fig. 9a). We leveraged this structure revealed by **HyperTools** to select several days of interest to examine further. Specifically, we examined (1) a day of Trump tweets whose topic coordinates were especially Trump-like (i.e. at the end of the Trump side of the V), (2) a day of Trump tweets whose topic coordinates fell at the intersection of the V, (3) a day of Clinton tweets whose topic coordinates fell at the intersection of the V, and (4) a day of Clinton tweets that fell at the end of the Clinton side of the V. For example, we wondered whether the candidates' tweets that fell at the extreme ends of the V might be especially representative of each candidates' unique features, whereas tweets that fell at the intersection of the V might express points of similarity between the candidates. Figure 9b displays the top 10 words from each of the top three topics for each of these days of interest, and Figure 9c provides representative tweets from each day. Strikingly (perhaps), the most Trump-like tweets appear to disparage Clinton, the most Clinton-like tweets appear to disparage Trump, and the overlapping tweets appear to praise America's greatness.

Example 5: Cyclical increases in global temperatures over time.

In addition to generating static point cloud plots, **HyperTools** may be used to generate trajectory plots to illustrate dynamic patterns in the data. To highlight this feature, we used a global

temperatures dataset which we acquired from Berkeley Earth. The Berkeley Earth averaging
216 method takes temperature observations from a large array of weather monitoring stations
217 throughout the world and produces a time-varying estimate of the underlying global temperature
218 field across all of the Earth’s land areas. This temperature field may then be sampled to obtain
219 location-specific temperature estimates.
220

To visualize how the global temperature field changes over time, we acquired monthly average
temperature estimates for 20 cities throughout the world (Fig. 10c) over the 138 year interval
from 1875–2013. We used `HyperTools` to plot the resulting temperature trajectory (Fig. 10a,b).
To visualize systematic changes over time, we plotted the month-by-month trajectory for each
year in a different color using the `group` keyword argument to `plot`:

```
hyp.plot(data, group = years, palette =' RdBu_r' )
```

 (23)

Two general trends were revealed by plotting the temperature data in this way. First, the
221 month-by-month temperatures within a year are cyclical (e.g. reflecting the changing seasons),
222 which appears in the trajectory as a “figure 8” (this trend is most visible in Fig. 10b). Second,
223 there has been a systematic shift in global temperatures over the 138 year period we examined.
224 This appears as a systematic shift in the position of the trajectory over time (Fig. 10a), and can
225 also be seen by directly plotting the temperatures over time (Fig. 10d).
226

Example 6: Visualizing the correspondence between neural trajectories and a movie stimulus

 227
228

In addition to providing plotting tools for visualizing complex data, `HyperTools` also provides
229 tools for aligning trajectories from different sources (see *Align*). For example, suppose we have
230 brain recordings from different people who all watched the same movie. The general shapes of
231 different people’s brain data trajectories (showing how everyone’s brain responses changed over
232 time while watching the movie), as well as the movie trajectory (showing how the movie itself
233 changed over time), might all share similar properties (e.g. reflecting the covariance structure of
234 the movie and how people responded to it). However, different people’s brains may have reflected
235 those similar responses differently, and the dimensions of “brain space” and “movie space” are not
236 directly comparable. As described in *Materials and Methods*, the `HyperTools` toolbox provides
237 an easy-to-use interface for aligning datasets. In this section example we demonstrate some uses
238 of the `align` function using a previously published fMRI dataset [8], available for download here.
239 The dataset comprises voxel responses from ventral temporal cortex, from each of 11 people,
240 as they watched the feature-length film *Raiders of the Lost Ark*. The data were processed and
241 hyperaligned as described in the original manuscript [8].
242

Figure 11a displays the trajectory plots for the averaged hyperaligned brain responses from
243 two groups of participants in the original experiment (six in group 1, the remaining five in group
244 2). The trajectories appear similar in their overall shape (indicating that the two groups of
245 participants had roughly similar brain responses to the movie), but the alignment is imperfect
246 (indicating that understanding individual differences between people’s responses might be an
247 interesting future direction to explore).
248

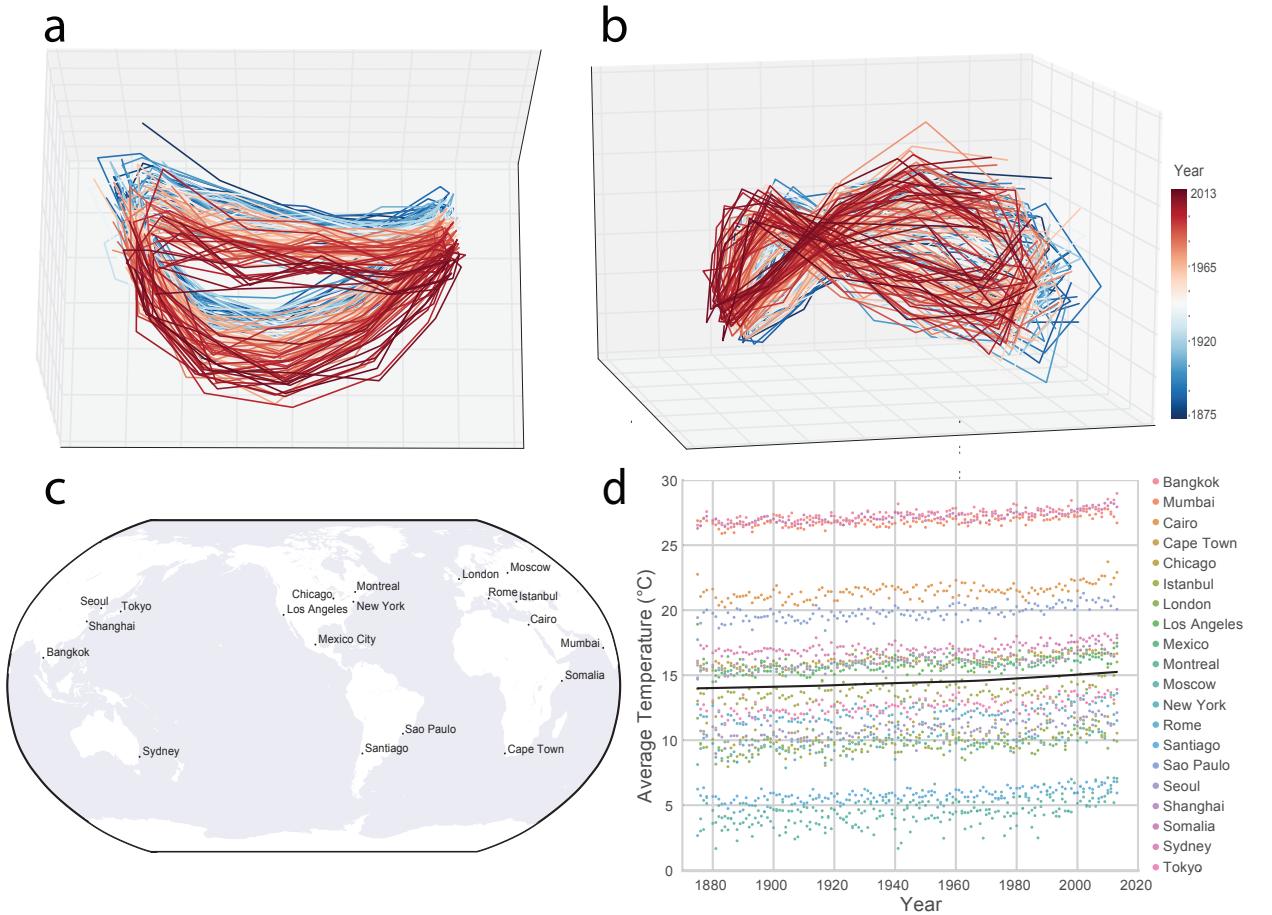


Figure 10. Global temperatures from 1875–2013. **a.** and **b.** The global temperatures dataset plotted using PCA dimensionality reduction in two views. The line colors change over time (from the earliest time point in blue to the most recent time point in red). The view on the left shows the temporal progression in one of the dimensions while the view on right highlights the cyclical nature of the dataset. **c.** Locations of the 20 cities in the dataset. **d.** Yearly mean temperatures colored by location and black LOWESS line fit to the full dataset.

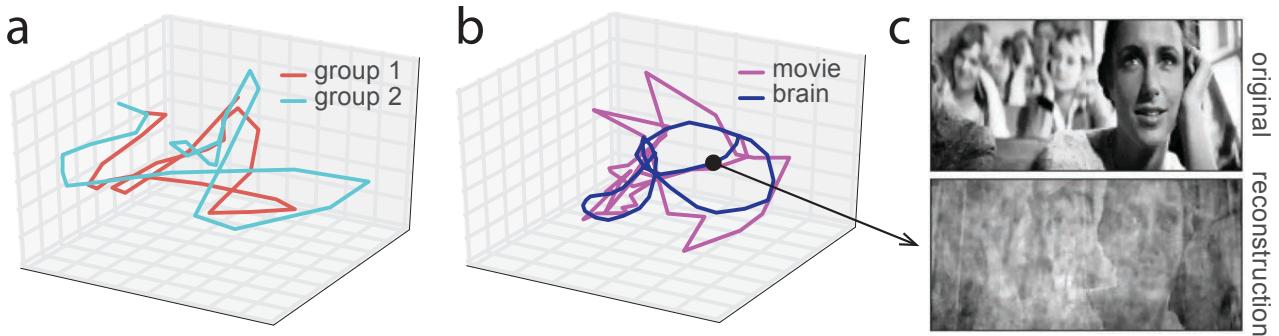


Figure 11. Brain/movie trajectories during movie viewing. **a.** Group-averaged trajectory of brain activity from ventral visual cortex split into two randomly-selected groups of subjects (group 1: $n = 6$, group 2: $n = 5$) watching the same movie. **b.** Group-averaged trajectory of brain activity from ventral visual cortex and trajectory of movie (pixel intensities over time) hyperaligned to a common space. **c.** Movie frame reconstructed from ventral visual brain activity that is aligned to movie space.

We next demonstrate how **HyperTools** may be used to visualize the correspondence between datasets with different coordinate systems—specifically, time-varying brain responses to the movie and the time-varying pixel intensities of the movie frames. To align these spaces, we first preprocessed the movie frames to convert the movie into the $S \times F$ matrix format required by **HyperTools** (here S is the number of movie frames and F is the number of pixels per frame). We downsampled the movie frames from 540×960 RGB pixels at 30 FPS to 108×192 grayscale pixels at 1 FPS. We then re-shaped each downsampled frame into a 20,736-dimensional vector.

We next averaged the (hyperaligned) brain responses from the 11 experimental participants to obtain a single brain response matrix. We used piecewise cubic interpolation [6] to re-sample this averaged brain response matrix from the original data acquisition rate (one image acquired every 2.5 s) to the downsampled movie frame rate (one image per second). We used the `reduce` function to project both the movie and brain data onto 6,641 dimensions (i.e. the number of voxels in the original brain data) and shifted the time labels of the brain matrix backwards by 5 s to account for the hemodynamic response. We then used the `procrustes` function to align the brain and movie data:

$$\text{brain_aligned_to_movie} = \text{hyp.tools.procrustes}(\text{movie_data}, \text{brain_data}) \quad (24)$$

The resulting aligned brain data matrix may then be plotted in the same space as the movie data matrix (Fig. 11b). This visualization can provide insights into the similarities and differences between the geometric structure of the original movie and the structure of the brain responses to the movie.

In addition to facilitating visual comparisons of the geometries of the movie and brain data, the aligned data may also be compared in the “native” data space. For example, each coordinate of “movie space” corresponds to an image, which may be displayed and examined. Aligning the brain data to this movie space (using the `procrustes` function) means that each brain pattern now corresponds to a coordinate in movie space, and therefore the corresponding image may

also be displayed and examined (Fig. 11c). This provides a means of viewing the original movie
265 through the “lens” of the brain responses to that movie. This general approach could also be
266 carried out in a cross-validated way (i.e. using one portion of the data to compute the Procrustean
267 transformation from brain space to movie space, and then applying that transformation to the
268 held-out brain data). We plan to explore this form of alignment-based decoding in future work.
269

Discussion

Visualizing high-dimensional data via low-dimensional embeddings provides an intuitive means
271 of exploring the geometric and statistical properties of complex datasets. This can help to guide
272 analysis decisions and facilitate hypothesis generation and testing. Returning briefly to the
273 example of Anscombe’s quartet we discussed in the *Introduction* (Fig. 1), striking differences
274 between datasets with very different geometries may be overlooked when solely considering their
275 summary statistics, and this principle can be extended to high-dimensional data as well. Our
276 **HyperTools** toolbox aims to assist in high-dimensional data visualization by providing a simple
277 (yet powerful) set of plotting functions and data manipulation tools.
278

We have provided brief examples of how our toolbox may be used to examine data from a
279 wide array of domains: geometry (Example 1), biological data (Example 2), educational and
280 sociological data (Example 3), political and linguistic data (Example 4), and neuroscientific data
281 (Example 5). We chose these particular examples to showcase a broad sampling of the types of
282 visualizations and analyses our toolbox supports, but they are not intended to indicate that our
283 toolbox may be used in only these ways or in these domains.
284

We hope that **HyperTools** will prove useful in analyzing and visualizing complex data from
285 a wide array of domains. We have released the toolbox under an open-source license to facilitate
286 transparency and widespread adoption. We also hope that users will contribute to the toolbox
287 by providing feedback and suggestions, and by sharing their own extensions and applications
288 with the community.
289

Acknowledgments

We are grateful for useful discussions with Luke J. Chang and Matthijs van der Meer. We are
290 also grateful for the help of J. Swaroop Guntapalli in implementing our `align` function. Our
work was supported in part by NSF EPSCoR Award Number 1632738. The content is solely
the responsibility of the authors and does not necessarily represent the official views of our
supporting organizations.

References

- [1] F. J. Anscombe. Graphs in statistical analysis. *American Statistitian*, 27(1):17–21, 1973.
- [2] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent dirichlet allocation. *Journal of Machine Learning Research*, 3:993 – 1022, 2003.

-
- [3] P.-H. C. Chen, J. Chen, Y. Yeshurun, U. Hasson, J. Haxby, and P. J. Ramadge. A Reduced-Dimension fMRI Shared Response Model. In C. Cortes and N. D. Lawrence and D. D. Lee and M. Sugiyama and R. Garnett, editor, *Advances in Neural Information Processing Systems 28*, pages 460–468. Curran Associates, Inc., 2015.
 - [4] P. Comon, C. Jutten, and J. Herault. Blind separation of sources, part II: Problems statement. *Signal Processing*, 24(1):11 – 20, 1991.
 - [5] M. Friendly. A brief history of data visualization. In C. Chen and W. Härdle and A Unwin, editor, *Handbook of Computational Statistics: Data Visualization*, volume III. Springer-Verlag, Heidelberg, 2006.
 - [6] F. N. Fritsch and R. E. Carlson. Monotone piecewise cubic interpolation. *SIAM J. Numer. Anal.*, 17(2):238–246, 1980.
 - [7] J. A. Hartigan and M. A. Wong. Algorithm AS 136: A k-means clustering algorithm. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 28(1):100–108, 1979.
 - [8] J. V. Haxby, J. S. Guntupalli, A. C. Connolly, Y. O. Halchenko, B. R. Conroy, M. I. Gobbini, M. Hanke, and P. J. Ramadge. A common, high-dimensional model of the representational space in human ventral temporal cortex. *Neuron*, 72:404–416, 2011.
 - [9] G. Hinton. Neural networks for machine learning. *Coursera*, <https://www.class-central.com/mooc/398/coursera-neural-networks-for-machine-learning>, 2012.
 - [10] J. D. Hunter. Matplotlib: A 2D graphics environment. *Computing In Science & Engineering*, 9(3):90–95, 2007.
 - [11] C. Jutten and J. Herault. Blind separation of sources, part I: An adaptive algorithm based on neuromimetic architecture. *Signal Processing*, 41(1):1 – 10, 1991.
 - [12] G. Lincoff and National Audubon Society. *The Audubon Society Field Guide to North American Mushrooms*. A Chanticleer Press edition. Knopf, 1981.
 - [13] W. McKinney. Data structures for statistical computing in Python. In *Proceedings of the 9th Python in Science Conference*, pages 51–56, 2010.
 - [14] K. Pearson. On lines and planes of closest fit to systems of points in space. *Philosophical Magazine*, 2:559–572, 1901.
 - [15] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
 - [16] P. Schönemann. A generalized solution of the orthogonal Procrustes problem. *Psychometrika*, 31:1–10, 1966.

-
- [17] M. E. Tipping and C. M. Bishop. Probabilistic principal component analysis. *Journal of Royal Statistical Society, Series B*, 61(3):611–622, 1999.
 - [18] W. S. Torgerson. *Theory and methods of scaling*. Wiley, New York, 1958.
 - [19] E. R. Tufte and P. Graves-Morris. *The visual display of quantitative information*, volume 2. Graphics press, Cheshire, CT, 1983.
 - [20] S. Uddenberg, G. Newman, and B. Scholl. Perceptual averaging of scientific data: Implications of ensemble representations for the perception of patterns in graphs. *Journal of Vision*, 16(12):1081, 2016.
 - [21] L. J. P. van der Maaten and G. E. Hinton. Visualizing high-dimensional data using t-SNE. *Journal of Machine Learning Research*, 9:2579–2605, 2008.
 - [22] S. van der Walt, S. C. Colbert, and G. Varoquaux. The NumPy array: A structure for efficient numerical computation. *Computing in Science & Engineering*, 13:22–30, 2011.
 - [23] M. Waskom, O. Botvinnik, D. Okane, P. Hobson, David, Y. Halchenko, S. Lukauskas, J. B. Cole, J. Warmenhoven, J. de Ruiter, S. Hoyer, J. Vanderplas, S. Villalba, G. Kunter, E. Quintero, M. Martin, A. Miles, K. Meyer, T. Augspurger, T. Yarkoni, P. Bachant, M. Williams, C. Evans, C. Fitzgerald, Brian, D. Wehner, G. Hitz, E. Ziegler, A. Qalieh, and A. Lee. Seaborn: v0.7.1 (June 2016), June 2016.