

1 Prototype: Matplottoy

To evaluate the feasibility of the model described in ??, we built prototypes of a `point`, `line`, and `bar` artist. We make use of the Matplotlib figure and axes artists [1, 2] so that we can initially focus on the data to graphic transformations and exploit the Matplotlib transform stack to convert data coordinates into screen coordinates. While the artist is specified in a fully functional manner in ??, we implement the prototype in a heavily object oriented manner. We do so mostly to more easily manage function inputs, especially parameters that are passed through to methods that are structurally functional.

Building on the current Matplotlib artists which construct an internal representation of the graphic, acts as an equivalence class artist A' as described in ??. The visual bundle V is specified as the dictionary of the form where parameter is a component in P , variable is a component in F , and the ν encoders are passed in as functions or callable objects. The data bundle E is passed in as a object. By binding data and transforms to A' inside , the method is a fully specified artist A as defined in ??. The data is fetched in section τ via a method on the data because the input to the artist is a section on E . The method takes the attribute because it provides the region in graphic coordinates S that can be used to query back into data to select a subset as described in ??. To ensure the integrity of the section, must be atomic, which means that the values cannot change after the method is called in draw until a new call to draw[3]. We put this constraint on the return of the method so that we do not risk race conditions.

The ν functions are then applied to the data, as describe in ??, to generate the visual section μ that here is the object . The conversion from data to visual space is simplified here to directly show that it is the encoding ν applied to the component. The function that is \hat{Q} , as defined in ??, is responsible for generating a representation such that it could be serialized to recreate a static version of the graphic. This artist is not optimized because we prioritized demonstrating the separability of ν and \hat{Q} . The last step in the artist function is handing itself off to the renderer. The extra arguments in are artifacts of how these objects are currently implemented.

29 1.1 Scatter, Line, and Bar Artists

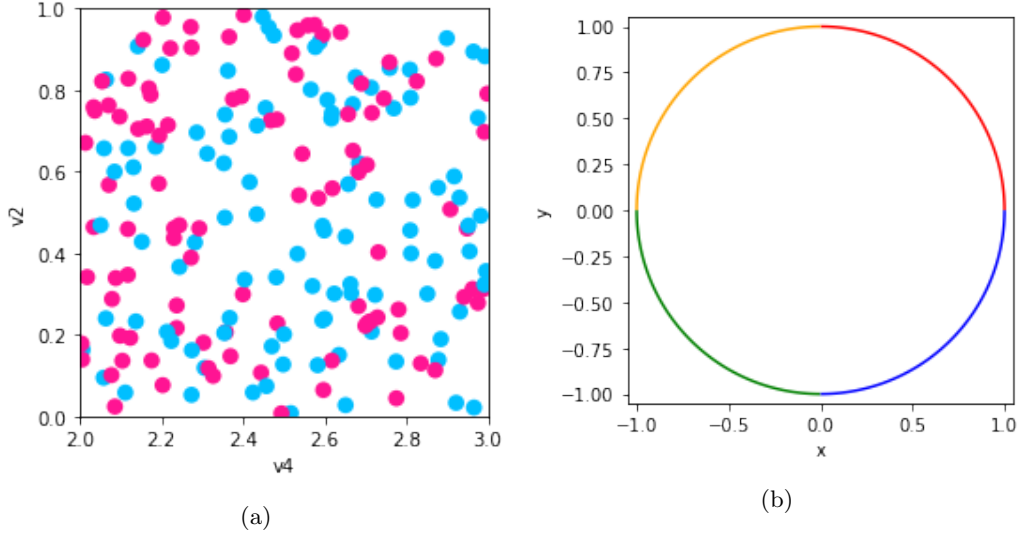


Figure 1: Scatter plot and line plot implemented using prototype artists and data models, building on Matplotlib rendering.

30 The figure in [Figure 1a](#) is described by `??`. This is implemented via a `Scatter` object where the
 31 scatter marker shape is fixed as a circle, and the visual fiber components are x and y position
 32 and the facecolor and size of the marker. We only show the `Scatter` function here because the `Scatter`
 33 are identical the prototype artist.

34 The `Scatter` method repackages the data as a fiber component indexed table of vertices. Even
 35 though the `Scatter` is fiber indexed, each vertex at an index k has corresponding values in section
 36 $\tau(k_i)$. This means that all the data on one vertex maps to one glyph. In `Scatter`, the μ components
 37 are used to construct the vector path of each circular marker with center (x, y) and size x
 38 and set the colors of each circle. This is done via the `Scatter` object. To generate [Figure 1b](#), the
 39 `Scatter` method returns a table of edges. Each edge consists of (x, y) points sampled along the
 40 line defined by the edge and information such as the color of the edge. As with `Scatter`, the data
 41 is then converted into visual variables. In `Scatter`, described by `??`, this visual representation is
 42 composed into a set of line segments, where each segment is the array generated by `Scatter`. Then
 43 the colors of each line segment are set. The colors are guaranteed to correspond to the correct
 44 segment because of the atomicity constraint on the view. The current implementation of
 45 line in Matplotlib does not have this functionality because it has no notion of an atomic
 46 view τ , instead the user keeps track of which values are input into a single call of \hat{Q} by
 47 creating a loop where every color is a new call to the line plotting function.

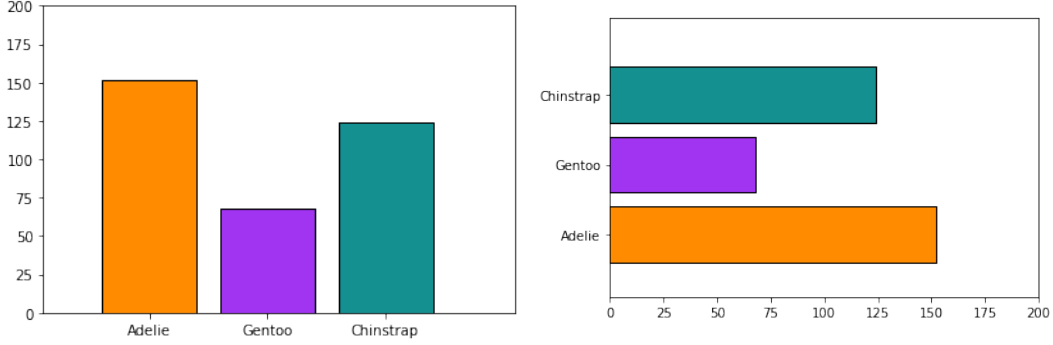


Figure 2: Frequency of Penguin types visualized as discrete bars.

The bar charts in figure 2 are generated with a `artist`. The artist has required visual parameters P of (position, length), and an additional parameter `orientation` which controls whether the bars are arranged vertically or horizontally. This parameter only applies holistically to the graphic and never to individual data parameters, and highlights how the model encourages explicit differentiation between parameters in V and graphic parameters applied directly to \hat{Q} .

As with `scatter` and `line`, `bar` function constructs bars and sets their properties, face and edge colors. The `bar` function converts the input position and length to the coordinates of a rectangle of the given width. Defaults are provided for 'width' and 'floor' to make this function more reusable. Typically the defaults are used for the type of chart shown in figure 2, but these visual variables are often set when building composite versions of this chart type as discussed in section 1.2.

1.2 Visual Encoders

The visual parameter serves as the dictionary key because the visual representation is constructed from the encoding applied to the data $\mu = \nu \circ \tau$. For the scatter plot, the mappings for the visual fiber components $P = (x, y, facecolors, s)$ are defined as where id is an identity ν , id maps into P without corresponding τ to set a constant visual value, and id is a custom ν implemented as a class for reusability. As described in ??, a test for equivariance can be implemented trivially. It is currently factored out of the artist for clarity.

1.3 Data Model

The data input into the `Data` will often be a wrapper class around an existing data structure. This wrapper object must specify the fiber components F and connectivity K and have a method that returns an atomic object that encapsulates τ . To support specifying the fiber bundle, we define a `Data` class[4]

that asks the user to specify the the properties of F and the K connectivity as either discrete vertices or continuous data along edges. To generate the scatter plot and the line plot, the distinction is in the `plot` method that is the section. The discrete `plot` method returns a record of discrete points, akin to a row in a table, while the line returns a sampling of points along an edge k . In both cases, the `plot` method packages the data into a data structure that

77 the artist can unpack via data component name, akin to a table with column names when
 78 K is 0 or 1 D.

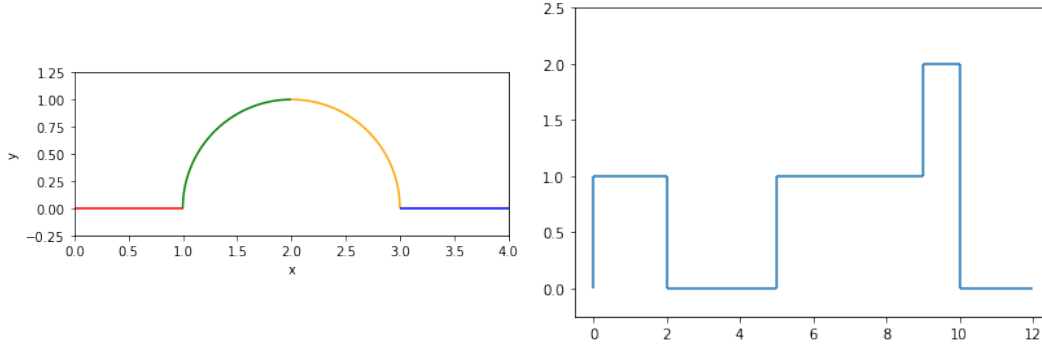


Figure 3: Continuous and discontinuous lines as defined via the same data model, and generated with the same A'

79 The graphics in figure [Figure 3](#) are made using the `artist` and the `data` source where
 80 if told that the data is connected, the data source will check for that connectivity by con-
 81 structing an adjacency matrix. The multicolored line is a connected graph of edges with
 82 each edge function evaluated on 100 samples, which is an arbitrary choice made to display
 83 a smooth curve. The axes can also be used to choose an appropriate number of samples. In
 84 contrast, the stair chart is discontinuous and only needs to be evaluated at the edges of the
 85 interval such that one advantage of this model is it helps differentiate graphics that have
 86 different artists from graphics that have the same artist but make different assumptions
 87 about the source data.

88 1.4 Case Study: Penguins

89 Building on the `artist` in [subsection 1.1](#), we implement grouped bar charts as these do not
 90 exist out of the box in the current version of Matplotlib. Instead, grouped bar charts are
 91 often achieved via looping over an implementation of `bar`, which forces the user to keep track
 92 which values are mapped to a single visual element and how that is achieved. For this case
 93 study, we use the Palmer Penguins dataset[5, 6], packaged as a pandas dataframe[7] since
 94 that is a very commonly used Python labeled data structure. The wrapper is very thin
 95 because there is explicitly only one section. Since the aim for this wrapper is to be very
 96 generic, here the fiber is set by querying the dataframe for its metadata. There are a list of
 97 column names and the datatype of the values in each column; this is the minimal amount
 98 of information the model requires to verify constraints. The pandas indexer is a key valued
 99 set of discrete vertices, so there is no need to repackage for the data interface.

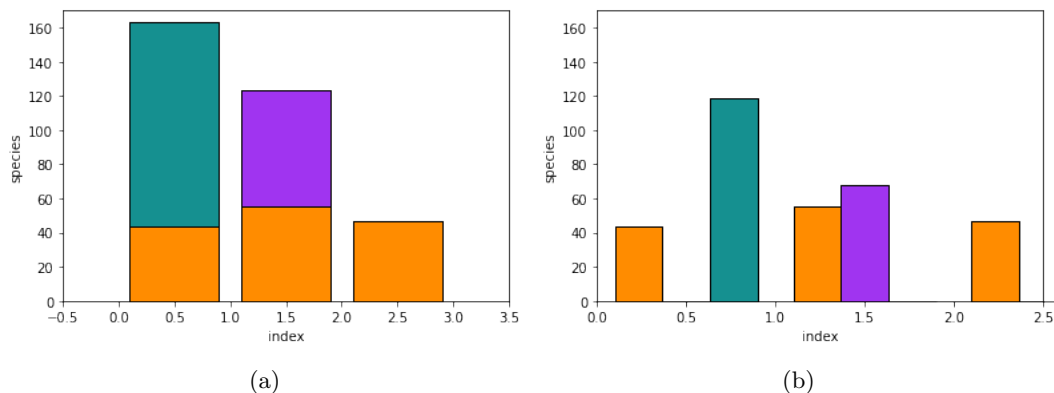


Figure 4: Penguin count disaggregated by island and species

The stacked and grouped bar charts in figure 4 are both composites of artists such that the difference between and is specific to the ways in which the are stitched together. These two artists have identical constructors and methods. As with , the orientation is set in the constructor. In both these artists, we separate the transforms that are applied to only one component (column) from transforms applied to multiple components (columns). We encode the one component case as dictionary, and the multiple components case as . This convention allows us to, for example, map one column to position, but multiple to length.

Since all the visual transformation is passed through to , the method does not do any visual transformations. In the is used to adjust the for every subsequent bar chart, since a stacked bar chart is bar chart area marks concatenated together in the length parameter. In contrast, does not even need the view, but instead keeps track of the relative position of each group of bars in the visual only variable .

Since the only difference between these two glyphs is in the composition of , they take in the exact same transform specification dictionaries. The dictionary dictates the position of the group, in this case by island the penguins are found on.

describes the group, and takes a list of dictionaries where each dictionary is the aesthetics of each group. That and are required parameters is enforced in the creation of the artist. These means that these two artists have identical function signatures

but differ in assembly \hat{Q} . By decomposing the architecture into data, visual encoding, and assembly steps, we are able to build components that are more flexible and also more self contained than the existing code base. While very rough, this API demonstrates that the ideas presented in the math framework are implementable. For example, the function that maps most closely to A is functional, with state only being necessary for bookkeeping the many inputs that the function requires. In choosing a functional approach, if not implementation, we provide a framework for library developers to build reusable encoder ν assembly \hat{Q} and artists A . We argue that if these functions are built such that they are equivariant with respect to monoid actions and the graphic topology is a deformation retraction of the data topology, then the artist by definition will be a structure and property preserving map from data to graphic.