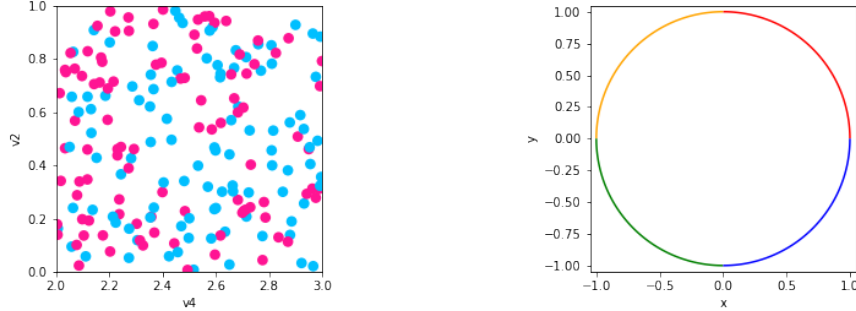# 1 Prototype Implementation: Matplottoy



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

To prototype our model, we implemented the artist classes for the scatter and line plots shown in figure 1 because they differ in every attribute: different visual channels $\nu$ that composite to different marks $Q$ with different continuities $\xi$ We make use of the Matplotlib figure and axes artists [3, 4] so that we can initially focus on the data to graphic transformations.

To generate the images in figure 1, we instantiate  artists that will contain the new primitive objects we implemented based on our topology model.

We then add the $A'=$ and $A'=$ artists that construct the scatter and line graphics. The arguments to the artist are the data $E=$ that is to be plotted and the aesthetic configuration $\nu=$. We implement the artists as equivalence classes $A'$ because it would be impractical to implement a new artist for every aesthetic setting, such as one artist for red lines and another for green.

## 1.1 Artist Class $A'$

The artist is the piece of the matplotlib architecture that constructs an internal representation of the graphic that the render then uses to draw the graphic. In the prototype artist, is a dictionary of the form  where parameter is a component in $P$, variable is a component in $F$, and the $\nu$ 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 $A$.

The data is fetched in section $\tau$ via a  method on the data because the input to the artist is a section on $E$. The return view object has a  method to support querying for components that are not in $F$ which we exploit to support parameters in the visual fiber that are not bound to fiber components in $F$. The $\nu$ functions are then applied to the data to generate the $\mu=$ input to $Q$. An explicit $\xi$ is not implemented since that would mean copying a single $\mu$ on $k$ to all the associated $s$, as illustrated in figure **??**, and that is unnecessary overhead for these scatter and line plots. In $\hat{Q}=$ the artist generates instructions for the render by

setting the attributes that are related to the graphic. These are the settings that would have to be serialized in order to recreate a static version of the graphic. Although could be implemented outside the class such that it returns an object the artist could then parse to set attributes, the attributes are directly set here to reduce indirection. The $\nu$ functions could be evaluated in this function to avoid passing over $K$ twice but are not done so here to demonstrate the seperability of $\nu$ and $\hat{Q}$ The last step in the artist function is handing itself off to the renderer.

The artist builds on artists because collections are optimized to efficiently draw a sequence of primitive point and area marks. In this prototype, the scatter marker shape is fixed as a circle, and the only visual fiber components are x and y position, size, and the facecolor of the marker. The method repackages the data as a fiber component indexed table of vertices, as described in section **??**; even though the is fiber indexed, each vertex at an index $k$ has corresponding values in section $\tau(k_i)$ such that all the data on one vertex maps to one marker. To ensure the integrity of the section, must be atomic, meaning that the values cannot change after the method is called in draw until a new call in draw. This table is converted to a table of visual variables. It is then passed into , where it is used to individually construct the vector path of each circular marker with center `(x,y)` and size `x` and set the colors of each circle. Since returns a $\tau$ all these operations could be applied on a section on one $k$ or a subset of $K$.

The only difference between the and objects is in the and function because line has different continuity from scatter and is represented by a different type of graphical mark.

In the artist, returns a table of edges. Each edge consists of (x,y) points sampled along the line defined by the edge and information such as the color of the edge. As with , the data is then converted into visual variables. In , this visual representation is composed into a set of line segments and then the colors of each line segment are set. The colors are guaranteed to correspond to the correct segment because of the atomicity constraint on view.

## 1.2   Encoders $\nu$

As mentioned above, the encoding dictionary is specified by the visual fiber component, the corresponding data fiber component, and the mapping function. 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 the position *(x,y)* $\nu$ transformers are identity functions. The size $s$ transformer is not acting on a component of $F$, instead it is a $\nu$ that returns a constant value. While size could be embedded inside the function, it is added to the transformers to illustrate user configured visual parameters that could either be constant or mapped to a component in $F$. The identity and constant $\nu$ are explicitly implemented here to demonstrate their implicit role in the visual pipeline, but they could be optimized away. More complex encoders can be implemented as callable classes, such as where can validate that the output of the $\nu$ is a valid element of the $P$ component the $\nu$ function is targeting. Creating a callable class also provides a simple way to swap out the specific (data, value) mapping without having to reimplement the validation or conversion logic.

A test for equivariance can be implemented trivially such that it is independent of data or encoder. In this example, checks for equivariance of permutation group actions by applying the encoder to a set of values, shuffling values, and checking that (value, encoding) pairs remain the same. This equivariance test can be implemented as part of the artist or

74 encoder, but for minimal overhead, the equivariant it is implemented as part of the library
75 tests.

## 1.3 Data $E$

77 The data input into the will often be a wrapper class around an existing data structure,
78 but must meet the following criteria:

79     1. specify the fiber components $F$ and connectivity $K$

80     2. have a that returns an atomic object that encapsulates $\tau$

81     3. the view object must have that returns a fiber component

82 To support specifying the fiber bundle, we define an optional class
83     that asks the user to specify how $K$ is triangulated and the attributes of $F$. The functions
84 expect tables that match the continuity of the graphic; scatter expects a vertex table because
85 it is discontinuous, line expects an edge table because it is 1D continuous. The fiber informs
86 appropriate choice of $\nu$ therefore it is a dictionary of attributes of the fiber components.
87 I've basically stripped this out of the artists above so should I just ditch this section?
88     To generate the scatter plot in figure 1, we fully specify a dataset with random keys
89 and values in a section chosen at random form the corresponding fiber component. The
90 fiberbundle is a class level attribute since all instances of VertexSimplex come from the
91 same fiberbundle. The view method returns a dictionary where the key is a fiber component
92 name and the value is a list of values in the fiber component. The table is built one call to
93 at a time, guaranteeing that all the fiber component values are over the same $k$. Table has
94 a method as it is a method on Python dictionaries. In contrast, the line in is defined as
95 the functions on each edge.
96     Unlike scatter, the line method returns the functions on the edge evaluated on the
97 interval [0,1]. By default these means each returns a list of 1000 x and y points and the
98 associated color. As with scatter, builds a table by calling for each $k$Unlike scatter, the
99 line table is a list where each item contains a list of points. This bookkeeping of which data
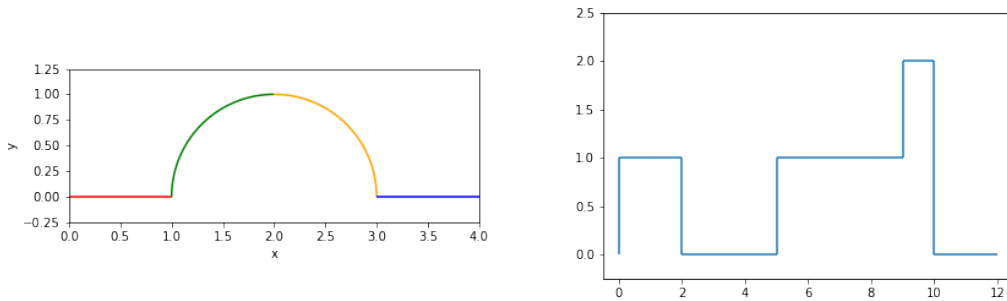100 is on an edge is used by the functions to bind segments to their visual properties.



Figure 2: Continuous and discontinuous lines as defined by different data models, but generated with the same $A'=$

101     The graphics in figure 2 are made using the artist and the data source

3

where if told that the data is connected, the data source will check for that connectivity by constructing an adjacency matrix. The multicolored line is a connected graph of edges with each edge function evaluated on 1000 samples while the stair chart is discontinuous and only needs to be evaluated at the edges of the interval such that one advantage of this model is it helps differentiate graphics that have different artists from graphics that have the same artist but make different assumptions about the source data.

## 1.4   Case Study: Penguins

For this case study, we use the Palmer Penguins dataset[1, 2] since it is multivariate and has a varying number of penguins. We use a version of the data packaged as a pandas dataframe[5, 6] since that is a very commonly used Python labled data structure. The wrapper is very thin since here there is explicitly only one section. The pandas indexer is a key valued set of discrete vertices, so there is no need to repackage for triangulation. As with the previous examples, there is no need to implement an explicit  method since the object has a get method.
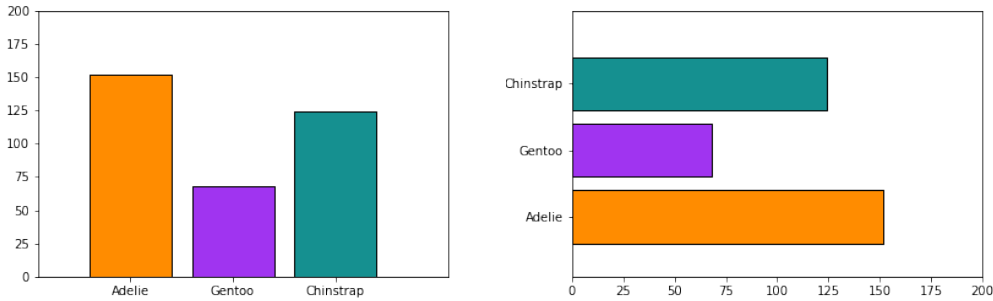


Figure 3: Frequency of Penguin types visualized as discrete bars.

The bar charts in figure 3 are generated with a  artist. The have the same required $P$ components of (position, length). In  of an additional parameter is set, `orientation` which only applies holistically to the graphic and never to individual data parameters. Explicitly differentiate between parameters in $V$ and ones that are only in $\hat{Q}$ is another way this model allows for cleaner separation of roles in the code.

The  method identical to the ones above, but here the visual transformations are factored out into a separate function. The  function sets some defaults, constructs bars, and sets their edge color to black. The  function is somewhat factored out because this is an operation that may be used by other bar making functions that may not be able to make use of bars assemble or draw.
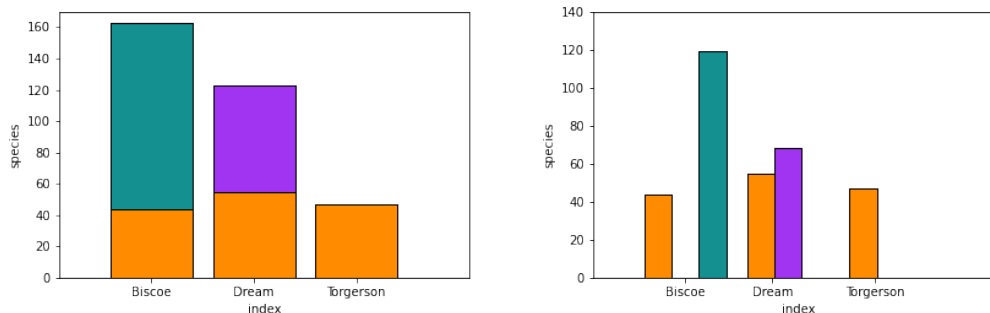
Figure 4: Penguin count disaggregated by island and species

For example, the  artist that makes figure 4 reuses  but does not reuse the assemble function because the composition of elements forces fundemental differences in glyph construction. As demonstrated in the , the composite bar chart has orientation and whether it is stacked or not. While the stacked bar chart and the grouped bar chart could be seperate artists, as demonstrated they share so much overlapping code that it is far less redundant to implement them together. looking at the mess that is this code, I'm a) not convinced these should be combined b) no longer convinced this provides anything over just bar if it isn't rewritten to use bar more

In the , a utility function is used for conversions, but the length transforms are held until after assembly because the length is computed by adding the current length to the previous and many transforms are not distributable such that $\nu(x_0 + x_1 + x_2) = \nu(x_0) + \nu(x_1) + \nu(x_2)$. Inside , the glyphs are either shifted vertically (`stacked`) or horizontally (`grouped`) such that the positions are recorded and added to with the next group. This function allows multiple columns to be mapped to a visual parameter, but it must be equal numbers of columns such as in this example where for each column contributing to a segment of the bar there is a corresponding column of colors for this segment. The reason the multibar can work with such a trasnformer is because it is relying on the data model to do most of the bookkeeping of which values get mapped to which bars. This also yields a much simpler function call to the artist

where  is the same dictionary for both stacked and grouped version, as is the  object . The only difference between the two versions is the `stacked` flag, and the only difference between figures 3 is the `orientation` argument. 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.

This API may want to be redesigned such that there's a way to clearly couple the columns when doing multindex broadcasting

5