

1 Prototype Implementation: Matplottoy

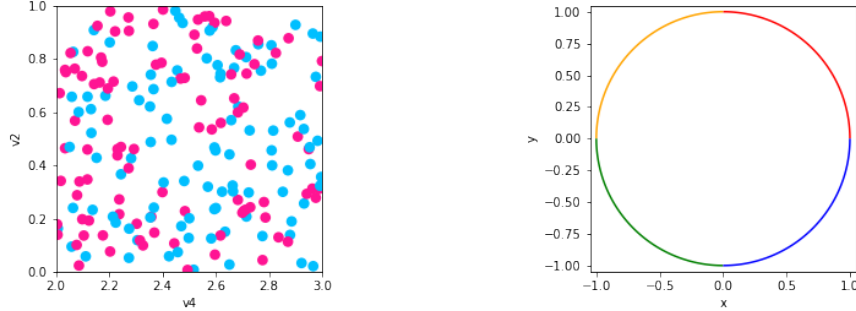


Figure 1: Scatter plot and line plot implemented using prototype artists and data models, building on Matplotlib rendering. **update with bar**

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 ν that composite to different marks Q with different continuities ξ . We make use of the Matplotlib figure and axes artists [4, 5] 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 artist that construct the scatter and line graphics. These artists are implemented as the equivalence class A' with the aesthetic configurations factored out into a dictionary that specifies the visual bundle V . The equivalence classes A' map well to Python classes since the functional aspects- ν , \hat{Q} , and ξ - are completely reusable in a consistent composition, while the visual values in V are what change between different artists belonging to the same class A' . The object is an abstraction of a data bundle E with a specified section τ . Implementing H and ρ are out of scope for this prototype because they are part of the rendering process. We also did not im

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 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 .

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 we can use to query back into data to select a subset as discussed in section ???. The ν functions are then applied to the data to generate the visual section

20 μ that here is the object . We allow for fixed visual parameter, such as setting a constant
 21 color for all sections, via setting in as the data fiber F name in the dictionary.

22 The visual object is then passed into the function that is \hat{Q} . This assembly function is
 23 responsible for generating a representation of the glyph such that it could be serialized to
 24 recreate a static version of the graphic. Although could be implemented outside the class
 25 such that it returns an object the artist could then parse to set attributes, the attributes are
 26 directly set here to reduce indirection. This artist is not optimized because we prioritized
 27 demonstrating the separability of ν and \hat{Q} . The last step in the artist function is handing
 28 itself off to the renderer. The extra arguments in are artifacts of how these objects are
 29 currently implemented in Matplotlib.

30 The artist builds on artists because collections are optimized to efficiently draw a
 31 sequence of primitive point and area marks. In this prototype, the scatter marker shape
 32 is fixed as a circle, and the only visual fiber components are x and y position, size, and
 33 the facecolor of the marker. We only show the function here because the are identical
 34 the prototype artist shown in ???. The method repackages the data as a fiber component
 35 indexed table of vertices, as described in section ???. Even though the is fiber indexed, each
 36 vertex at an index k has corresponding values in section $\tau(k_i)$. This means that all the data
 37 on one vertex maps to one glyph. To ensure the integrity of the section, must be atomic.
 38 This means that the values cannot change after the method is called in draw until a new
 39 call in draw. We put this constraint on the return of the pythonview method so that we
 40 do not risk race conditions. Using triangulization provides a common interface to the data,
 41 one we will reuse for all the artists presented here.

42 This table is converted to a table of visual variables and is then passed into . In , the μ
 43 is used to individually construct the vector path of each circular marker with center (\mathbf{x}, \mathbf{y})
 44 and size \mathbf{x} and set the colors of each circle. This is done via the object. As mentioned in
 45 sections ?? and ??, this assembly function could as easily be implemented such that it was
 46 fed one $\tau(k)$ at a time.

47 The main difference between the and objects is in th3 function because line has different
 48 continuity from scatter and is represented by a different type of graphical mark.

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

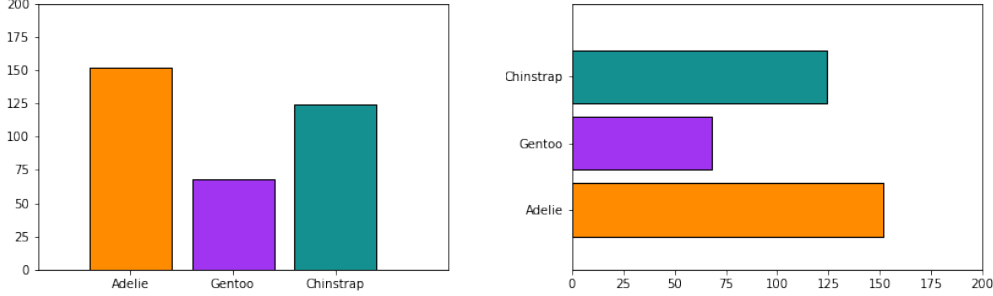


Figure 2: Frequency of Penguin types visualized as discrete bars.

55 The bar charts in figure 2 are generated with a `BarChart` artist. They have the same required
 56 P components of (position, length). In addition, an additional parameter is set, `orientation`
 57 which controls whether the bars are arranged vertically or horizontally and only applies
 58 holistically to the graphic and never to individual data parameters. This highlights another
 59 advantage of this model, that it encourages explicit differentiation between parameters in
 60 V and graphic parameters that specify \hat{Q} .

61 The method is the same as the one in the artist example ?? so is omitted here. This bar
 62 underpins a rudimentary version of the complex operator such that we can specify multiple
 63 components to map to the same visual components without requiring the components to
 64 be individually specified in the dictionary. The implementation of this support is slightly
 65 more complicated than the simple case, but not in a way that is useful for understanding the
 66 framework. The function constructs bars and sets their edge color to black. Defaults are
 67 provided for 'width' and 'floor' to make this function more reusable. Typically the defaults
 68 are used for the type of chart shown in figure 2, but these visual variables are often set when
 69 building composite versions of this chart type as discussed in section ??.

70 1.2 Encoders ν

71 As mentioned above, the encoding dictionary is specified by the visual fiber component, the
 72 corresponding data fiber component, and the mapping function. The visual parameter serves
 73 as the dictionary key because the visual representation is constructed from the encoding
 74 applied to the data $\mu = \nu \circ \tau$. For the scatter plot, the mappings for the visual fiber
 75 components $P = (x, y, facecolors, s)$ are defined as where the position (x, y) ν transformers
 76 are identity functions. The size s transformer is not acting on a component of F , instead it
 77 is a ν that returns a constant value. While size could be embedded inside the function, it is
 78 added to the transformers to illustrate user configured visual parameters that could either
 79 be constant or mapped to a component in F . The identity and constant ν are explicitly
 80 implemented here to demonstrate their implicit role in the visual pipeline, but they could
 81 be optimized away. More complex encoders can be implemented as callable classes, such as
 82 where one can validate that the output of the ν is a valid element of the P component the
 83 ν function is targeting. Creating a callable class also provides a simple way to swap out the
 84 specific (data, value) mapping without having to reimplement the validation or conversion
 85 logic. A test for equivariance can be implemented trivially but is currently factored out of the
 86 artist for clarity. 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 encoder, but for minimal overhead, the equivariant it is implemented as part of the library tests.

1.3 Data E

The data input into the 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 that returns an atomic object that encapsulates τ . The object returned by the view must be key valued pairs of where each section is a component as defined in equation ?? . To support specifying the fiber bundle, we define a data class[1]

that asks the user to specify how K is triangulated and the attributes of F . Python dataclasses are a good abstraction for the fiber bundle class because the class only stores data. The K is specified as tables because the functions expect tables that match the continuity of the graphic; scatter expects a vertex table because it is discontinuous, line expects an edge table because it is 1D continuous. The fiber informs appropriate choice of ν therefore it is a dictionary of attributes of the fiber components.

To generate the scatter plot in figure 1, we fully specify a dataset with random keys and values in a section chosen at random from the corresponding fiber component. The fiberbundle is a class level attribute since all instances of VertexSimplex come from the same fiberbundle. The view method returns a dictionary where the key is a fiber component name and the value is a list of values in the fiber component. The table is built one call to at a time, guaranteeing that all the fiber component values are over the same k . Table has a method as it is a method on Python dictionaries. In contrast, the line in is defined as the functions on each edge.

Unlike scatter, the line method returns the functions on the edge evaluated on the interval $[0,1]$. By default these means each returns a list of 1000 x and y points and the associated color. As with scatter, builds a table by calling for each k . Unlike scatter, the line table is a list where each item contains a list of points. This bookkeeping of which data is on an edge is used by the functions to bind segments to their visual properties.

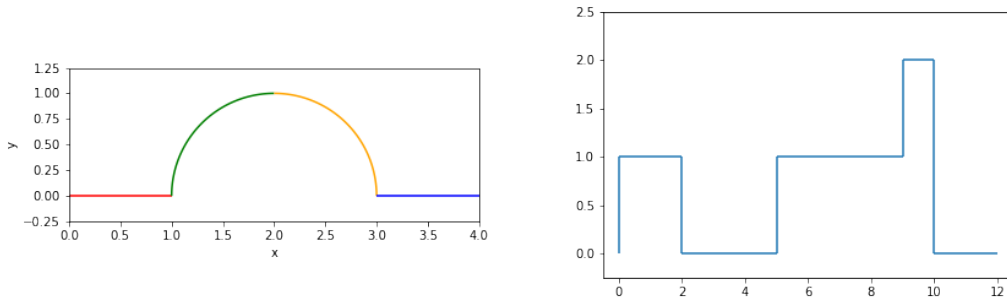


Figure 3: Continuous and discontinuous lines as defined by different data models, but generated with the same A'

The graphics in figure 3 are made using the artist and the data source

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[2, 3] since it is multivariate and has a varying number of penguins. We use a version of the data packaged as a pandas dataframe[6, 7] since that is a very commonly used Python labeled 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.

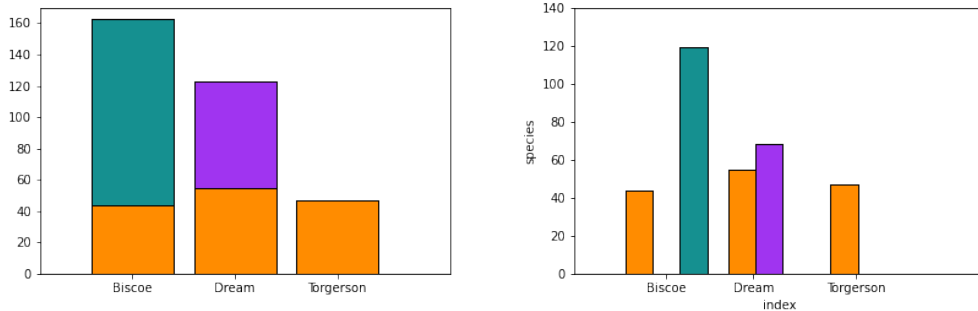


Figure 4: Penguin count disaggregated by island and species

For example, the artist that makes figure 4 reuses but does not reuse the assemble function because the composition of elements forces fundamental 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 separate 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 transformer 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

148 where \mathcal{D} is the same dictionary for both stacked and grouped version, as is the \mathcal{O} object .
 149 The only difference between the two versions is the **stacked** flag, and the only difference
 150 between figures 2 is the **orientation** argument. By decomposing the architecture into data,
 151 visual encoding, and assembly steps, we are able to build components that are more flexible
 152 and also more self contained than the existing code base.

153 1.5 Summary

154 In general, the way in which we implemented the artist is as follows:

155 ν encoder function converting data to library normalized form, can also be method

156 μ dictionary of the form \mathcal{D} curried until `Artist.draw`

157 A' objects that take in as input E and V for example \mathcal{D}, \mathcal{O} ,

158 A method which is A' parameterized by V

159 \hat{Q} method that arranges components of V into glyphs

160 with the data representation implemented as

161 E object with `method`

162 F named components of the data and their types

163 K determines how the data is bound together when returned by view

164 τ method that returns component sections

165 While very rough, this API demonstrates that the ideas presented in the math framework
 166 are implementable. In choosing a functional approach, if not implementation, we provide a
 167 framework for library developers to build reusable encoder ν and assembly \hat{Q} . We argue that
 168 if these functions are built such that they are equivariant with respect to monoid actions
 169 and the graphic topology is a deformation retraction of the data topology, then the artist
 170 by definition will be a structure and property preserving map from data to graphic.