

FUNCTIONAL PEARLS

Composing data visualizations

TOMAS PETRICEK
University of Kent, UK
(*e-mail*: t.petricek@kent.ac.uk)

1 Introduction

Let's say we want to create the two charts in Figure 1. The chart on the left is a bar chart that shows two different values for each bar. The chart on the right consists of two line charts that share the X axis and highlight two parts of the timeline with two different colors.

There is a plenty of libraries that can draw bar charts and line charts, but adding those extra features will only be possible if the author already thought about your exact scenario. For example, Google Charts supports the left chart (it is called Dual-X Bar Chart) but there is no way for adding a background, or sharing an axis between charts. The alternative is to use a more low-level library such as D3. In D3 you construct the chart piece by piece, but then you have to tediously transform your values to coordinates in pixels yourself. For scientific plots, you could use an implementation of Grammar of Graphics such as ggplot2, where a chart is a mapping from data to geometric objects (such as points, bars, and lines) and their visual properties (X and Y coordinate, shape and color). However, the range of charts that can be created using this systematic approach is still somewhat limited.

What would an elegant functional approach to data visualization look like? A functional programmer would want a domain-specific language that has a small number of primitives; allow us to define high level abstractions such as a bar chart and its basic building blocks are expressed in terms of domain values such as the exchange rate, rather than pixels.

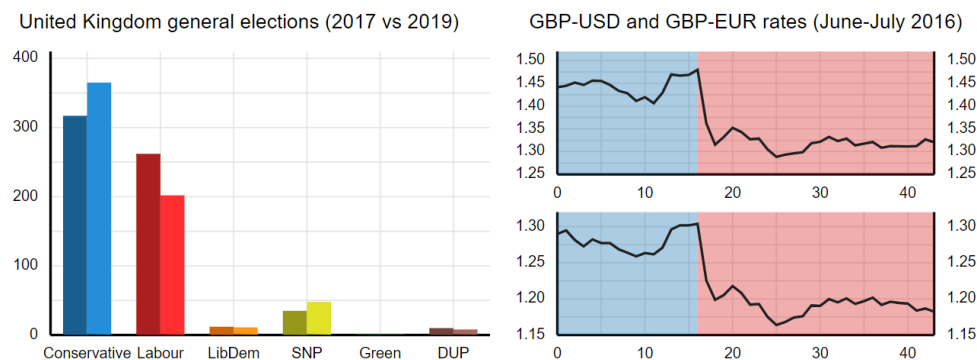


Fig. 1. Two charts about the UK politics: Comparison of election results from 2017 and 2019 (left) and GBP-USD exchange rate with highlighted areas before and after the 23 June 2016 Brexit vote.

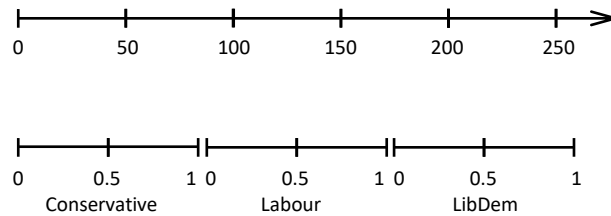


Fig. 2. On a continuous scale (above), an exact position is determined by a number. On a categorical scale (below), an exact position is determined by the category and a numerical ratio from 0 to 1.

As is often the case with domain-specific languages, finding the right primitives is more of an art than science. For this reason, I present my answer – a library named *Compost* – as a functional pearl. I hope to convince the reader that *Compost* is elegant and I illustrate this with a wide range of examples. *Compost* has a number of specific desirable properties:

- Charts are composed from a small number of primitive building blocks using a small number of combinators. In particular, concepts such as bar charts, line charts or charts with aligned axes are expressed in terms of more basic concepts.
- The primitives are specified in domain terms. When drawing a line, the value of an Y coordinate is an exchange rate of 1.36 USD/GBP, not 137 pixels from the bottom.
- Most common chart types can be easily captured as high level abstractions, but there is an elegant way of creating a majority of more interesting custom charts.
- The approach can easily be extended to creating web-based charts that involve animations or interaction with the user.

The presentation in this paper focuses on explaining the primitives and combinators of the domain-specific language. I outline the structure of an implementation, but omit the details. Filling those in requires careful thinking about geometry and projections, but there are no unexpected surprises. A complete F# implementation, including the examples used in this paper, is available at: <http://github.com/compostjs>.

2 Composing simple visualizations

I will introduce individual features of the *Compost* library gradually. The first important aspect of *Compost* is that properties of shapes are defined in terms of domain-specific values. I first explain what this means and then use domain-specific values to specify the core part of the UK election results bar chart.

2.1 Working with domain-specific values

In the election results chart in Figure 1 (left), the X values are categorical values representing the political parties such as *Conservative* or *Labour*. The Y values are numerical values representing the number of seats won such as 365 MPs. When creating data visualizations, those are the values that the user needs to specify. This is akin to most high-level charting libraries such as Google Charts, but in contrast with more flexible libraries like D3.

v	$=$	<code>cat</code> c, r	s	$=$	<code>line</code> $\gamma, [v_{x1}, v_{y1}, \dots, v_{xn}, v_{yn}]$	<code>overlay</code> $[s_1, \dots, s_n]$
		<code>cont</code> n			<code>fill</code> $\gamma, [v_{x1}, v_{y1}, \dots, v_{xn}, v_{yn}]$	<code>axis</code> $l/r/t/b$ s
					<code>text</code> γ, v_x, v_y, t	<code>padding</code> n_l, n_r, n_b, n_t, s

Fig. 3. Core primitives of the Compost domain-specific language. Values v are either categorical or continuous; a shape s is then defined as a simple recursive algebraic data type.

Our design focuses on two-dimensional charts with X and Y axes. Values mapped to those axes can be either categorical (such as different political parties or countries) or continuous (such as number of votes or exchange rates). The mapping from categorical and continuous values to exact positions on the chart is done automatically. For continuous values, this simply means applying a linear transformation. For categorical values, the mapping is more difficult.

For example, in the UK election results chart, the X axis is categorical. The library automatically divides the available space between the six categorical values (political parties). The value `Green` does not determine an exact position on the axis, but rather a range. To determine an exact position, we also need to attach a value between 0 and 1 to the categorical value. This identifies a relative position in the available range.

Figure 2 illustrates the two kinds of values using the axes from the UK election results chart. In Figure 3, we define a value v as either a continuous value `cont` n containing any number n or a categorical value `cat` c, r , consisting of a categorical value c (implemented as a string) and a ratio r between 0 and 1.

2.2 Introducing basic primitives and combinators

Now that we know how Compost represents values, we can define the basic elements of its domain-specific language. A chart is represented by the shape s defined in Figure 3. A primitive shape can be a text label, a line connecting a list of points or a filled polygon defined by a list of points. The position of points is specified by X and Y coordinates, which can be either categorical or continuous values. For text, line and polygon, we also include a parameter γ that specifies the element color.

Figure 3 also defines three combinators. The most important is `overlay`, which overlays all shapes from a given list. When doing this, Compost automatically infers the scales of X and Y axes and calculates suitable projections using a method discussed in the next section. Finally, `padding` adds padding around a specified shape and `axis` adds an axis showing the inferred scale on the left, right, top or bottom of a given shape. Using those primitives, we can construct the simple UK election results bar chart in Figure 4 as follows:

```
axisl (axisb (overlay [
  fill #0000ff, [ (cat Conservative, 0), (cont 0), (cat Conservative, 0), (cont 365),
    (cat Conservative, 1), (cont 365), (cat Conservative, 1), (cont 0) ],
  fill #ff0000, [ (cat Labour, 0), (cont 0), (cat Labour, 0), (cont 202),
    (cat Labour, 1), (cont 202), (cat Labour, 1), (cont 0) ] ]))
```

The chart specification overlays two bars of different colors and then adds axes to the bottom and left of the chart. The two bars are filled rectangles defined using four corner

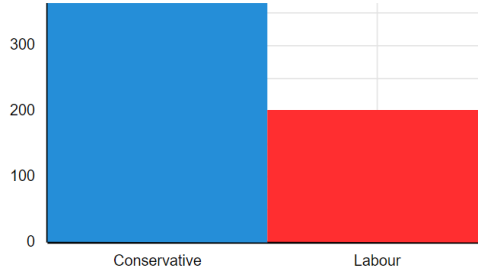


Fig. 4. Two charts about the UK politics: Comparison of election results from 2017 and 2019 (left) and GBP-USD exchange rate with highlighted areas before and after the 23 June 2016 Brexit vote.

points. The Y coordinates are specified as continuous values, while the X coordinates are categorical. For the Conservative party, two of the points have the Y coordinate set to `cont 0` (bottom of the bar) and two have the Y coordinate set to `cont 365` (top of the bar). The two X coordinates are the start and the end of the range allocated for the `Conservative` category, i.e. `cat Conservative, 0` on the left and `cat Conservative, 1` on the right.

Extending the snippet to generate a grouped bar chart that shows two results for each party as in Figure 1 is easy. Given a party p , we need to generate two rectangles, one with X coordinates `cat p, 0` and `cat p, 0.5` and the other with X coordinates `cat p, 0.5` and `cat p, 1`. In the following snippet, we use a `for` comprehension to generate the list. All remaining constructs are primitives of the Compost domain-specific language. Assuming elections is a list of election results containing a five-element tuple consisting of a party name, colors for 2017 and 2019 and results for 2017 and 2019 we create the chart using:

```
axisl (axisb (overlay [
  for party, clr17, clr19, mp17, mp19 in elections →
  padding 0, 10, 0, 10, overlay [
    fill clr17, [(cat party, 0), (cont 0), (cat party, 0.5), (cont mp17),
                (cat party, 0.5), (cont mp17), (cat party, 1), (cont 0)],
    fill clr19, [(cat party, 0.5), (cont 0), (cat party, 0.5), (cont mp19),
                (cat party, 1), (cont mp19), (cat party, 1), (cont 0)] ] ]))
```

Aside from iterating over all available parties and splitting the bar, the example also adds padding around the bars. A padding is specified in pixels rather than in terms of domain values. This is sometimes preferable over, for example, drawing a bar using a range from `0.05` to `0.5`. One remaining feature of the chart in Figure 1 that is still missing is the caption. We will add this in Section 4.

2.3 Inferring scales and projections automatically

When composing shapes using the `overlay` primitive, the user does not need to specify how to position the child elements relatively to each other. The Compost library positions the elements automatically. This is done in two steps. First, Compost infers the *scales* for X and Y axes. A scale represents the range of values that needs to fit in the space available for the chart. Second, Compost calculates a *projection*, a mapping from domain-specific values of the scale to the available screen space. A scale l can be either categorical or continuous:

$$l = \text{continuous } n_{\min}, n_{\max} \mid \text{categorical } [c_1, \dots, c_k]$$

A continuous scale is defined by a minimal and maximal value that need to be mapped to the available chart space. A categorical scale is defined by a list of individual categorical values. Note that we do not need a minimal and maximal ratios of the used categorical values as Compost will use an equal space for each category, regardless of where in this space a shape needs to appear.

Inferring scales is done by a simple recursive function that walks over the given shape and constructs two scales for the X and Y axis, using the X and Y coordinates that appear in the shape. Most of the work is done by a simple helper function that takes two scales, l_1 and l_2 , and produces a new scale that represents the union of the two:

$$\begin{aligned} \text{union } (\text{continuous } n_l, n_h) (\text{continuous } n'_l, n'_h) &= \\ \text{continuous } \min(n_l, n'_l), \max(n_h, n'_h) & \\ \text{union } (\text{categorical } [c_1, \dots, c_p]) (\text{categorical } [c'_1, \dots, c'_q]) &= \\ \text{categorical } [c_1, \dots, c_p] @ [c'_i \mid \forall i \in 1 \dots q, j. c_j = c'_i] & \end{aligned}$$

When unioning two continuous scales, the minimum and maximum of the resulting scale is the smallest and largest of the two minimums and maximums, respectively. When unioning two categorical scales, we take all values of the first scale and append all values of the second scale that do not appear in the first one. Note that this means that the order of categorical values in a scale depends on the order in which they appear in the shape. (A possible improvement to Compost would be to support ordinal values, which are categorical values with a well-defined ordering.) It is also worth noting that a categorical scale cannot be combined with a continuous scale. In other words, mixing categorical and continuous values in a single scale results in an error.

Once Compost computes scales for the given shape and all its sub-shapes, it constructs a projection function that maps domain-specific values to the available chart space. I say more about this in Section 6. For a continuous scale, the projection is a linear transformation. For categorical scale with k values, we split the available chart space into k equally sized regions and then map a categorical value $\text{cat } c, r$ to the region corresponding to c according to the ratio r .

3 Composing advanced visualizations

3.1 Creating scales with nice ranges

3.2 Controlling the composition of scales

Functional pearls

7

4 Defining abstractions

5 Interactive charts

6 Implementation structure

Controlling and nesting scales

The process of computing scales can be controlled by additional primitives in the Compost domain-specific language for describing shapes that were not discussed earlier. The following definition lists additional primitives that are also a part of the definition of s . Note that the primitives can be applied either to the X scale or to the Y scale - we write x/y to indicate that a sub-script on those primitives can be set either to x or to y (both scales can be modified by nesting two primitives):

$$s = \begin{array}{l} (\dots) \\ | \text{axis}_{x/y} s \\ | \text{roundScale}_{x/y} s \\ | \text{explicitScale}_{x/y} l, s \\ | \text{nest}_{x/y} v_{min}, v_{max}, s \end{array}$$

We first focus on the first three primitives, while the `nest` primitive will be discussed in the next section. The `axis` and `roundScale` constructs could be defined as derived constructs, but it is easier to consider them as primitives for now. The `axis` primitive simply takes the shape s and draws an axis around the contents of s , using the inferred scale of s to determine values displayed on the axis. This could be a derived construct, because axes can be rendered using lines and text labels. The `roundScale` operation takes the inferred X or Y scale of the shape s and, if it is a continuous scale, rounds its minimal and maximal values to “nice” numbers. For example, if a continuous scale has minimum 0 and maximum 66.04, the resulting scale would have maximum 70. For categorical scale, the operation does not have any effect. The `explicitScale` operation is similar, but it replaces the inferred scale with an explicitly provided scale (the type of the inferred scale has to match with the type of the explicitly given scale). For example, assuming `populationBar` is the example given earlier, we can write:

```
axisx (axisy (roundScaley
  (explicitScalex (categorical CZ,MN,UK) populationBar)))
```

This replaces the X axis with an explicitly given one that includes extra country code and also overrides the implicitly inferred order of the categorical values. We then ask Compost to automatically round the Y scale (showing population) so that we get a “nice” number as the maximum. Finally, we add axes around the shape, producing a usual labelled chart. As noted earlier, both `axis` and `roundScale` could be defined as derived – `axis` would have to infer the scales of the nested shape, calculate appropriate labels and insert suitable lines and labels; the `roundScale` operation would need to infer the scale too and then add an explicit `explicitScale` with a rounded minimal and maximal value of a continuous scale.

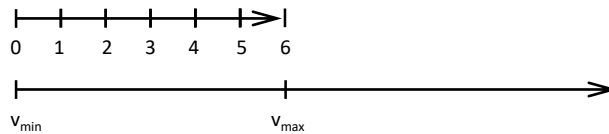


Fig. 5. A continuous scale with values from 0 to 6, nested in another scale.

Nesting of charts and scales

Finally, one more primitive that we added in the above definition and that we have not explained yet is `nestx/y`. The parameters of the primitive are two values, v_{min} and v_{max} together with a shape s . The primitive makes it possible to nest one scale inside another one. When inferring the scales, Compost infers the scale of the nested shape s . This will then take the area determined by the two points v_{min} and v_{max} in the scale determined by other shapes that are overlaid with the `nest` shape. The two points are also used to determine the scales of the overall shape, but the scales of s are nested and do not have any effect on the outer scales. In fact, the scale of s does not even have to be of the same type as the scale of the outer shape. Figure 5 illustrates this with an example. Here, we are nesting a continuous scale with values ranging from 0 to 6 into another scale. The values v_{max} and v_{min} do not even have to be continuous. We can, for example, nest a line chart inside a bar of a bar chart and then the two values could be `cat One, 0` and `cat One, 1` (which are two values that define a region of a categorical scale). One area where nesting of scales is particularly useful is when combining multiple charts in a single view. For example, let's say that we have two line charts showing prices of two different stocks over the same period of time (for simplicity, we can just use UNIX timestamp as a time). The prices of the two stocks are orders of magnitude different, so they cannot fit easily into the same chart. We want to show two line charts side-by-side (one above each other). They should each have their own Y axis (price), but the X axis (time) should be shared. Assuming we have `googPrices` and `fbPrices` as two line charts without any axes, covering the same date range, we can write:

```
axisx (overlay
  (nesty (cat top-chart, 0), (cat top-chart, 1), (axisy fbPrices)),
  (nesty (cat bottom-chart, 0), (cat bottom-chart, 1), (axisy googPrices)))
```

Reading the code from the inner-most part to the outer-most, we first add separate Y axes to both of the line charts. Given the difference in the prices, the axes will have quite different values. We then nest the (continuous) Y axes using the `nest` primitive and, at the same time, implicitly define a new outer Y axis. The outer Y axis is categorical with just two values, `top-chart` and `bottom-chart`. This means that the two nested charts will take equal amount of space, one above the other (each of the charts takes the full space allocated for the category – the minimal value has ratio 0 while the maximal value has ratio 1). [TODO: We really need illustrations for those example charts, but that would be more work!] [TODO: Maybe use this chart as a motivation: <https://wellcomeopenresearch.org/articles/4-63/v1>]

7 Random ideas

- We could define a type system to catch categorical/continuous value mismatch in a single
- It would be worth thinking about ordinal values, which are categorical but can be sorted.
- There might be other kinds of scales - for example, color scale (can have meaning in scatter plot) or a scale for secondary markers (like sizes of bubbles in a bubble

chart). We could really say that every value (including bar chart colors) should be coming from some scale...

- Implementing something like `axis` or `roundScale` as an actual derived primitive is a bit tricky, because it needs to invoke a part of the normal rendering workflow (to run the automatic inference of scales on the nested shape) – this might be just implementation issue, but it could be some more basic problem.