

---

Subject: Re: Efficiently perform histogram reverse indices like procedure on a string array?

Posted by [Craig Markwardt](#) on Thu, 26 Jul 2012 21:33:51 GMT

[View Forum Message](#) <> [Reply to Message](#)

---

On Wednesday, July 25, 2012 7:39:33 PM UTC-4, Bogdanovist wrote:

> I have an array of a data structure, one tag of which is a string identifier indicating which location the data belongs to. There are many thousands of data points, but only about a dozen or so unique locations.

>

> I make frequent use of the HISTOGRAM function with the reverse\_indices in order to carve up data by some identifier, most commonly the time. In this case, I want to divide out the data by site efficiently. I can't use HISTOGRAM on strings, so I need some other approach. There are plenty of ways this can be done, but I'd like some views on the better and most efficient ways to do it.

>

> Take an example, say we have a simple string array

>

> foo=[ 'a','b','c','b','b','a','a','c'];

>

> To determine the list of unique strings we could do

>

> sfoo = foo[sort(foo)]

> print,sfoo[uniq(sfoo)]

>

> We can then repeatedly use WHERE to find the indices in the data array(s) corresponding to each site.

>

> Is there a quicker/better way to do this? Repeatedly calling WHERE seems inefficient (certainly HISTOGRAM is way faster when it is usable)

I prefer to do it slightly differently than your other suggestions.

I locate the breakpoints between different runs of strings like this,

```
ibreaks = where(sfoo[1:*] NE sfoo, ct)
```

This gives the interior breakpoints. In your case, ibreaks = [2,5], which is the point where 'a' changes to 'b', and 'b' changes to 'c'. Usually I add this little bit of extra post-processing,

```
if ct EQ 0 then begin
```

```
  ibreaks = [0, n_elements(sfoo)]
```

```
endif else begin
```

```
  ibreaks = [0, ibreaks+1, n_elements(sfoo)]
```

```
endelse
```

You need that little extra 'if' statement to handle the case where you have only one unique string,

so there are no breaks at all.

The start of the  $i$ th run is indexed by `ibreaks[i]`, and the end of the  $i$ th run is indexed by `ibreaks[i+1]-1`, where  $i$  goes from 0 through `n_elements(sfoo)-1`.

I.e. the  $i$ th run is given by `sfoo[ibreaks[i]:ibreaks[i+1]-1]`. Of course you can index back into the original array once you've done this.

Craig

---