

Transforming a museum to be data-driven using R

Alice Daish
Data Scientist
@alice_data



The British Museum

Set up

Opened 1759 to all 'studious and curious persons'
1st National Public Museum in the World



Today

2nd most visited museum in the world
8 million objects
2 million years of human history



Starting point

Didn't have ...

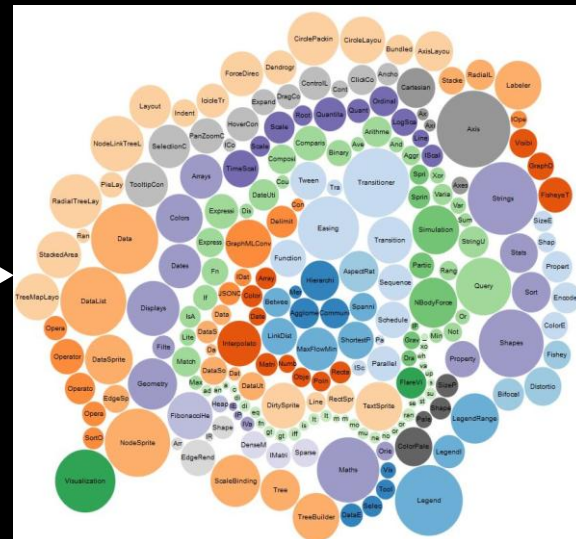
- No list of data sources
- No data access
- No databases
- No data warehouse

Did have...

- R
- Data Scientist
- Big Data: Senior Product manager
- What does “big data” mean to the museum?



Bubbles Envy



Joe Cheng to the rescue

d3.js bubble chart htmlwidget for R

This R package provides a bubble chart as seen in [this Mike Bostock example](#). It is based on [htmlwidgets](#) so it can be used from the R console, RStudio, R Markdown documents, and Shiny applications.

Installation

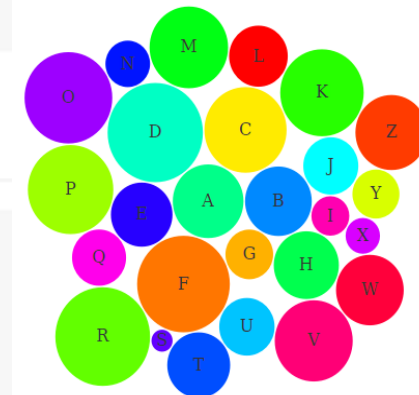
Use the **devtools** package (`install.packages("devtools")`) to install this package directly from GitHub:

```
devtools::install_github("jcheng5/bubbles")
```

Usage

```
library(bubbles)

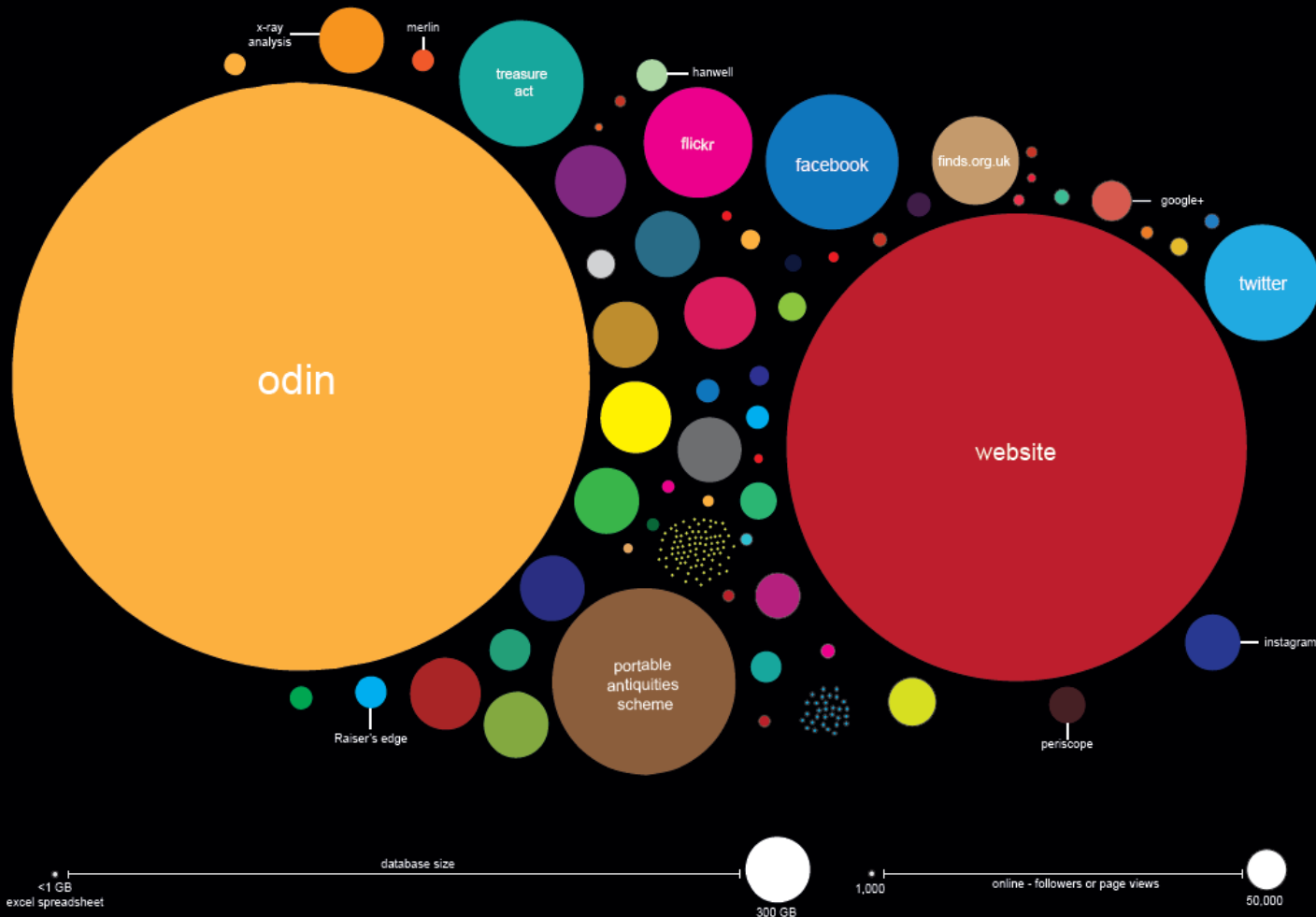
bubbles(value = runif(26), label = LETTERS,
        color = rainbow(26, alpha=NULL)[sample(26)]
)
```



The British Museum

British Museum Data Audit 2015

The big data team conducted a survey across the museum and found more than 250 data sources. This visualisation demonstrates the variety of data across the museum.



Business Problems = Data Opportunities

We don't know who our visitors are?

Online = > 9 million

Offline = 6.8 million

We don't what they do in the museum?

And we don't know the opportunities to generate revenue?

Business Problems = Data Opportunities

We don't know who our visitors are?

Online = > 9 million

Offline = 6.8 million

We don't what they do in the museum?

And we don't know the opportunities to generate revenue?

“silos” and “wrangling”

data viz

visitor movement

predictive modelling

“silos” and “wrangling”

Multiple visitor data platforms

CSV exports from external platforms

No SQL



email



online
shop



wi-fi

“silos” and “wrangling”

100's of columns



email

Multiple visitor data platforms

CSV exports from external platforms

No SQL

*print format exports
nested by timeslots*



online
shop

Split first and second name



wi-fi

No SQL = data.table

How many visitors are on multiple platforms?

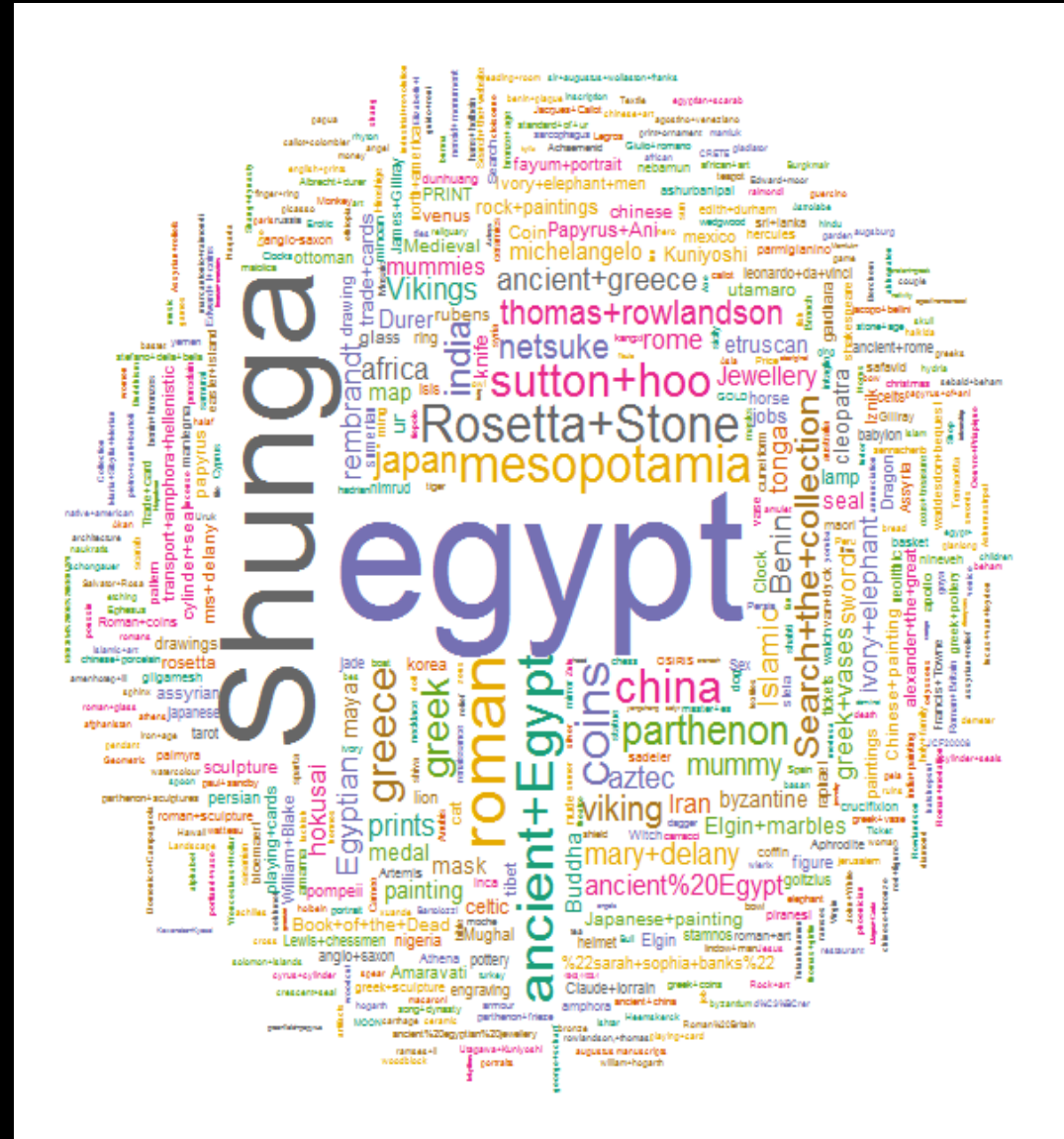
Assess the visitor data siloing

Why?

To improve engagement and access of the museum we need to examine our visitor data.



Top 500 website search



Packages : RSiteCatalyst (Adobe Analytics), WordCloud



visitor movement

*62 galleries, 3 floors,
largest covered public
square in Europe with 6.8
million visitors per year.*



Visitor Movement

Wi-Fi presence
used to sample
visitor numbers



*1st to use R to
connect to CISCO
Presence API*

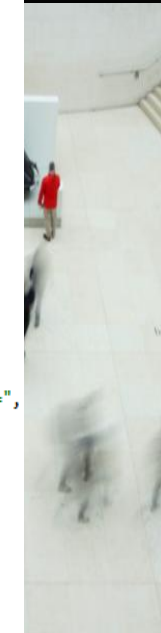


Visitor M

Wi-Fi presence
used to sample
visitor number

*1st to use R to
connect to CISCO
Presence API*

```
1 #-----CISCO API Connection to collect CISCO data-----
2 #DATE : 03/06/2016
3 #AUTHOR: Alice Daish adaish@britishmuseum.org
4
5 #----Install and Load packages
6 library(httr)
7 library(jsonlite)
8
9 #Load password and username code file
10 source("logincisco.R")
11
12 #----FIND THE LIST OF SITES-----
13 sites<-GET("https://cmxcisco.com/api/config/v1/sites",authenticate(user, password))
14 # gets the URL api content including authroization
15
16 #testing different export formats
17 str(content(sites)) #see content
18 sitelist<-content(sites, "text") #collects content as text string
19 sitelist<-fromJSON(sitelist) #convert to table format from string
20 head(sitelist) #see the top of the table
21
22 #List of site name and siteId
23 siteId<-cbind(sitelist$aesUidString,sitelist$name)
24
25 #EXAMPLE COLLECT HOURLY DATA OF ALL SITE FOR ONE DAY (14/05/2016)
26 hourdata<-matrix(NA,nrow = 97*1, ncol = 27) #blank matrix |
27 colnames(hourdata)<-c("SiteID","SiteName","Date","0","1","2","3","4","5",
28 "6","7","8","9","10","11","12","13","14","15","16",
29 "17","18","19","20","21","22","23") #Label columns
30
31 hourly<-GET(paste0("https://lnzgy2.cmxcisco.com/api/presence/v1/visitor/hourly?siteId=",
32 siteId[i,1],"&date=2016-05-14"),authenticate(user, password))
33
34 hourdata[i,1]<-siteId[i,1] #ID
35 hourdata[i,2]<-siteId[i,2] #Name of site
36 hourdata[i,3]<- "2016-05-14"
37 hourdata[i,4]<-content(hourly)$`0`
38 hourdata[i,5]<-content(hourly)$`1`
39 hourdata[i,6]<-content(hourly)$`2`
40 hourdata[i,7]<-content(hourly)$`3`
```



Visitor M

Wi-Fi presence
used to sample
visitor number

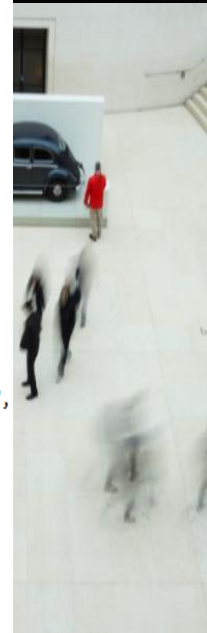
*1st to use R to
connect to CISCO
Presence API*

```
1 #-----CISCO API Connection to collect CISCO data-----
2 #DATE : 03/06/2016
3 #AUTHOR: Alice Daish adaish@britishmuseum.org
```

	9-10am	10-11am	11am-12p	12-13pm	13-14pm	14-15pm	15pm-16p	16-17pm	17-18pm	Total
5	23	155	212	254	174	165	155	141	120	1399
6	38	76	94	153	120	115	101	158	92	947
7	26	52	82	92	91	84	80	105	54	666
8	27	50	68	81	74	65	83	87	98	633
9	7	23	64	93	68	65	69	85	36	510
10	11	61	78	61	76	48	61	57	37	490
11	13	56	45	88	47	27	31	78	69	454
12	26	40	42	65	56	58	48	47	48	430
13	14	29	53	111	48	31	35	72	35	428
14	1	26	69	62	51	64	66	56	27	422
15	1	39	78	57	52	45	53	52	19	396
16	0	30	56	57	53	59	59	53	23	390
17	3	24	56	56	47	55	41	43	15	340
18	15	18	29	30	34	41	51	36	86	340
19	14	43	30	34	44	43	36	31	50	325
20	5	20	48	73	38	36	30	44	27	321
21	7	39	53	43	38	30	42	39	27	318
22	3	22	44	54	41	36	39	50	26	315
23	4	16	39	50	50	38	38	42	36	313
24	0	30	31	28	43	44	37	38	16	267
25	6	17	30	53	24	24	31	40	28	253
26	19	27	20	30	18	22	23	31	38	228
27	2	14	21	39	37	35	25	20	22	215
28	8	27	32	23	27	29	25	18	22	211
29	12	12	15	22	28	27	29	25	37	207
30	22	38	33	24	22	19	15	16	15	204
31	11	12	19	33	32	21	25	32	17	202
32	2	17	26	26	31	26	25	30	12	195
33	10	33	22	25	10	17	20	29	4	170
34	2	19	17	29	31	19	20	19	12	168
35	0	9	23	28	19	17	25	35	11	167
36	0	15	24	31	13	16	20	35	12	166
37	11	14	21	33	23	20	18	13	8	161
38	11	9	27	28	23	23	19	12	6	158
39	3	19	19	25	15	15	24	18	18	156

```
40 hourodata[i,7]<-content(hourly)$'3'
```

siteId=",
)



Predictive modelling

Can we predict ticket sales for exhibitions?

mixed effect modelling

- data wrangling
- modelling
- prediction



Predictive modelling

First initial model created – lmer()
Predicted first exhibition sales – predict()
Development continues ...



Future

Building a data pipeline including R.
Continued data wrangling of the museums data sources to find insights and value.

Who knows?

Internet of things

e.g. Toilet doors locks, Boilers, Visitor Flow Signs

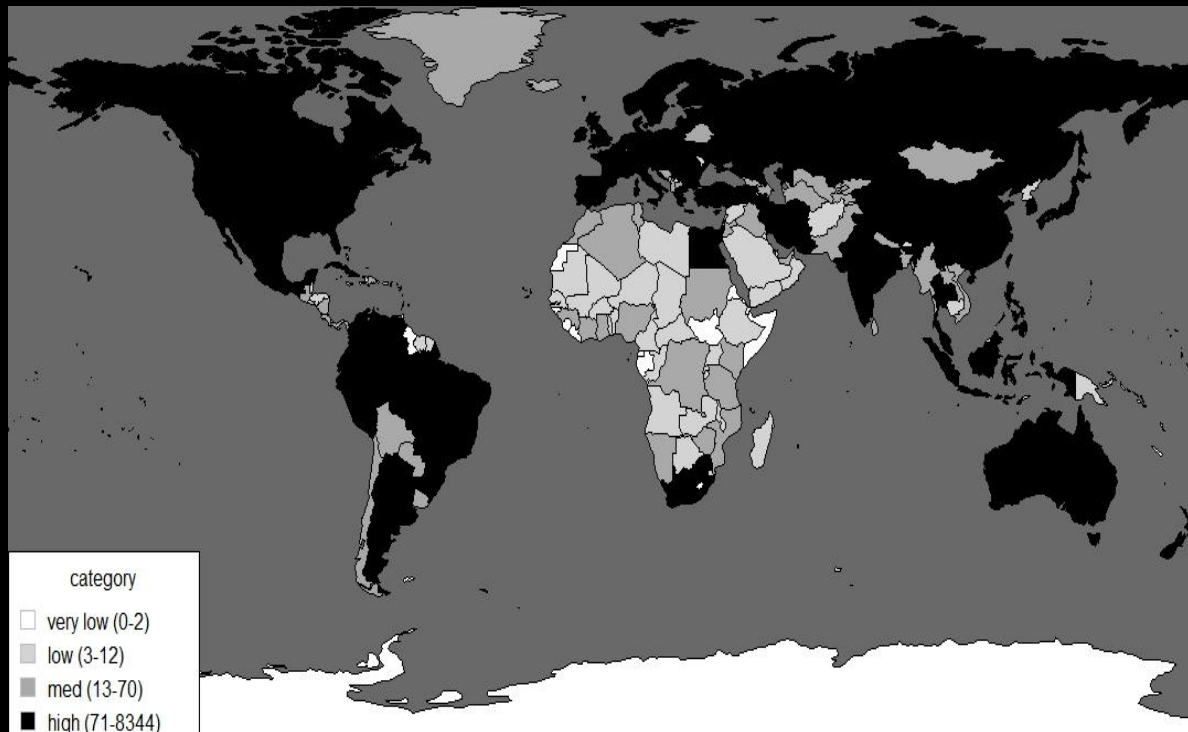
Machine Learning

e.g. Predicting Visitor Numbers, Optimization, Cognitive Services.





55,000 museums 180 countries



Packages : rworldmap



Thank you & Questions

adaish@britishmuseum.org
@alice_data

Many thanks to museum departments support and data access,
Siorna Ashby, and the R community for their continued support

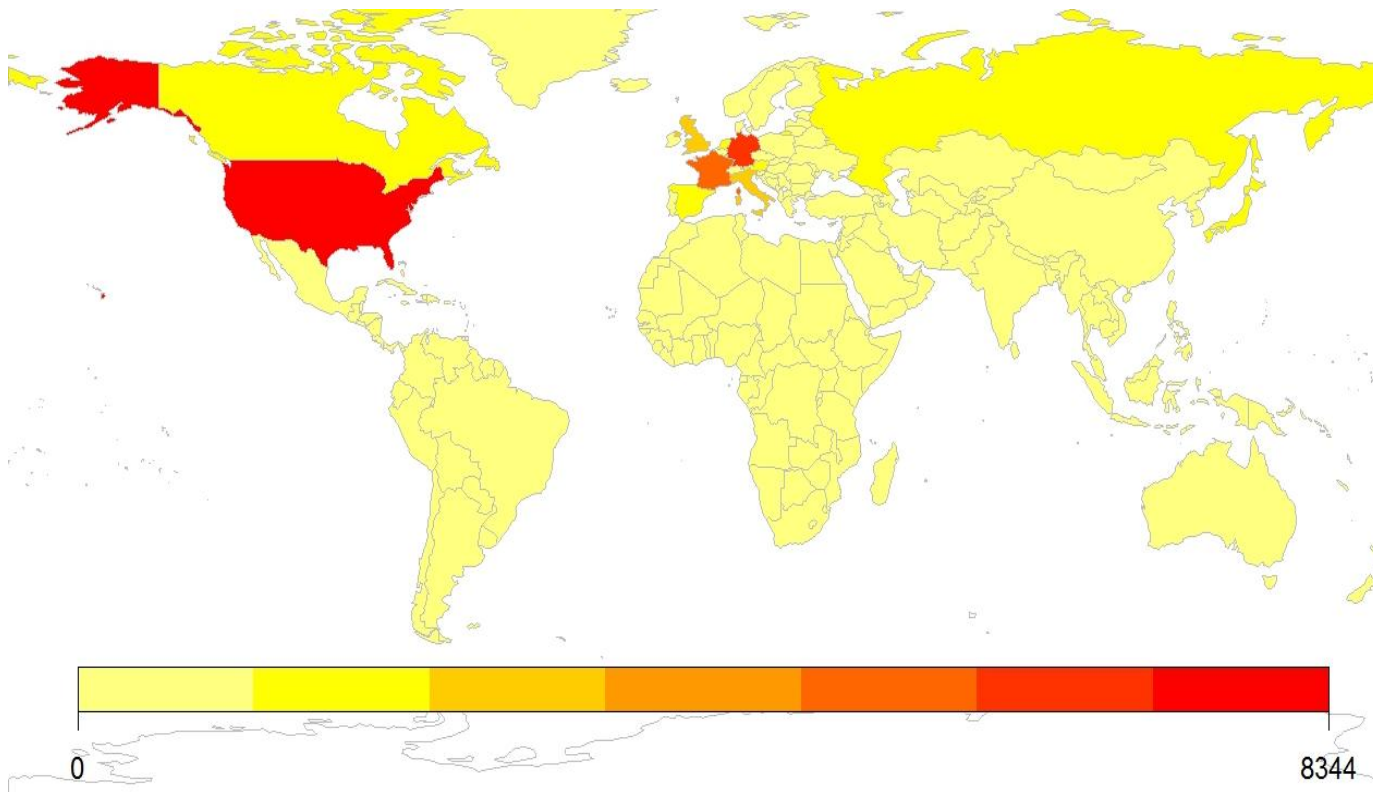
**8344
USA**

**6464
Germany**

**4842
France**

**3200
Italy**

**2943
UK**



Countries with the most museums