

Hotel TULIP Web Server Data Analysis

Assignment 2 - SIT742 Modern Data Science

Alex Cummaudo <ca@deakin.edu.au>

Jake Renzella <jake.renzella@deakin.edu.au>

Deakin Software and Technology Innovation Laboratory

School of Information Technology

Deakin University, Australia

May 11, 2017

Executive Summary

This report summarises findings from a data exploration on the Hotel TULIP web server logs, recorded between the periods of August 2014 and August 2015. Each log contains one *request*, or *hit*, that lists fourteen attributes as described in the attached Data Dictionary spreadsheet. Publicly known client IP addresses were extracted from the MaxMind GeoIP2¹ dataset to analyse the location of requests (narrowed down to city). Additionally, user agent strings were parsed to analyse device and browser statistics using the Python user-agents library², thereby extrapolating demographics, usage trends, platform information, server performance, and security statistics from the raw logs provided in the dataset. Further details on the extraction of the data is provided in the source code attached in Appendix C, and an interactive version of this file is published on Databricks.

¹See <http://dev.maxmind.com/geoip/geoip2/>.

²See <https://pypi.python.org/pypi/user-agents>.

Contents

1	Key Findings	6
2	Introduction	7
3	Dataset	7
4	Method	7
4.1	Assumptions Made	7
4.2	Extraction Process	7
4.2.1	IP Address Source Regular Expression	7
4.3	Data Mining	8
5	Results	8
5.1	Sample Directional Network Graph	8
5.2	IP Request Sources	11
5.2.1	Internal Site Visitors	11
5.2.2	External Site Visitors	13
5.3	Top Three Countries	15
5.3.1	Hong Kong Visitors	15
5.3.2	USA Visitors	17
5.3.3	Australian Visitors	19
5.4	Platform Categories	21
5.4.1	PC Visitors	21
5.4.2	Smartphone Visitors	23
5.4.3	Tablet Visitors	25
5.4.4	Bots Visitors	27
A	Additional Tables	30
B	Data Dictionary	39
B.1	Dataset Description	39

B.2	Contact Information	39
B.3	Data Dictionary	40
C	Extrapolation Results	43

List of Figures

1	Sample frequency patterns identified in a subset of PC requests	10
2	Directional network graph visualising frequency patterns of internal visitors made on an internal IP range	12
3	Directional network graph visualising frequency patterns of external visitors made on a non-internal IP range	14
4	Directional network graph visualising frequency patterns of visitors from Hong Kong	16
5	Directional network graph visualising frequency patterns of visitors from the United States	18
6	Directional network graph visualising frequency patterns of visitors from Australia	20
7	Directional network graph visualising frequency patterns of requests made by PCs .	22
8	Directional network graph visualising frequency patterns of requests made on smart- phones	24
9	Directional network graph visualising frequency patterns of requests made on tablets	26
10	Directional network graph visualising frequency patterns of requests made by bots	28

List of Tables

1	Sample Frequency Graph	10
2	Internal Request Frequency Patterns	30
3	External Request Frequency Patterns	30
4	Hong Kong Frequency Patterns	31
5	USA Request Frequency Patterns	32
6	Australian Request Frequency Patterns	33
7	PC Request Frequency Patterns	34
8	Tablet Request Frequency Patterns	36
9	Smartphone Request Frequency Patterns	37
10	Bots Request Frequency Patterns	37
11	Dataset Description	39
12	Dataset Contact Information	39
13	Data Dictionary	40

1 Key Findings

A list of key findings in the analysis are as thus:

- Foo

2 Introduction

Browsing patterns on the Hotel TULIP Weblogs were assessed in order to gain insight on how customers navigate through the website within a typical *session*. A user session is defined as a typical visit to the website, and a collation of all the different *informational resources* that were accessed. Information resources refer to web pages that contain primary content about the hotel and its facilities, rather than multimedia and technical-related resources.

In order to assess how *different* customers do so, we extract different information based on the web log data format as prescribed in Appendix B. We contrast those users who make requests:

- internally, such as guests using the internet within the hotel’s network,
- externally, such as prospective guests browsing the website for a potential stay in the hotel,
- from users within the top three countries that visit the website (refer to Assignment 1), and
- between users on PCs, Smartphones, Tablets and Bots.

Each of these criteria were analysed against matching sessions that satisfy such criteria. Data mined using the Frequency Pattern was done so using the FPGrowth Algorithm in Apache Spark.

3 Dataset

Hotel TULIP’s web server runs Microsoft Internet Information Services (IIS) Server 7.5, and the attributes of this dataset as well as the relevant data dictionary can be found in Appendix B.3.

4 Method

4.1 Assumptions Made

4.2 Extraction Process

4.2.1 IP Address Source Regular Expression

To differentiate between private site visitors and external visitors, a regular expression was used to filter the private and public IP address ranges. The regular expression is shown below:


```

WHERE 1.c_ip REGEXP
      ' (^127\..)|(^10\..)|
      .....(^172\.1[6-9]\..)|
      .....(^172\.2[0-9]\..)|
      .....(^172\.3[0-1]\..)|
      .....(^192\.168\..) '

```

Negating this WHERE clause of the regular expression will select only public IP addresses.

4.3 Data Mining

5 Results

5.1 Sample Directional Network Graph

In our results, we have visualised the frequency patterns of users via the use of directional network graphs. In these graphs, we are able to visualise the frequency patterns of how people made requests to the website given the assumptions and extraction methods made in Section 4.

Within each graph, a *sequence* is identified as a series of multiple clicks (edges) between pages (nodes). Each sequence is coloured using the same edge colour. This sequence is also identified using a number, which is drawn on the edge label. The frequency of this pattern for the particular sequence identified is given after the forward slash on the label.

For example, a sample directional network graph, a subset of the PC requests, is given in Figure 1. This data is also presented in tabular format as Table 1.

Here we can interpret that the graph has four key sequences, as differentiated by the sequence numbers. Sequence numbers are ordered by decreasing frequency; the higher the sequence number the increased likelihood of the pattern occurring. In this graph, we see that users of PCs are most likely to move between pages in the following order:

1. from the 'Above and Beyond' page to the 'Dining' page (frequency of 293),
2. from the 'Facilities' page to the 'Dining' page (Frequency of 288),

3. from the 'Above and Beyond' page to the 'Offers' page (frequency of 286), and, with equal frequency,
4. from the 'Facilities' page to the 'Offers' page.

Table 1: Sample Frequency Graph

Sequence	From	To	Frequency
2	above and beyond	→ rooms	293
2	rooms	→ dining	293
3	facilities	→ rooms	288
3	rooms	→ dining	288
4	above and beyond	→ rooms	286
4	rooms	→ offers	286
5	facilities	→ rooms	286
5	rooms	→ offers	286

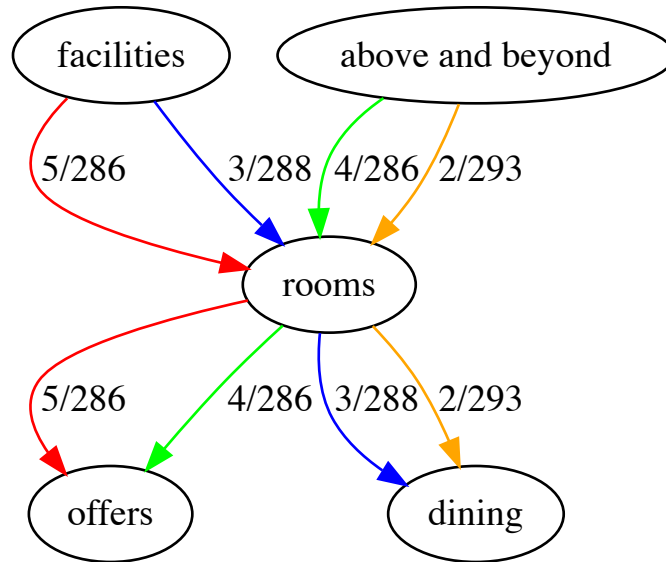


Figure 1: Sample frequency patterns identified in a subset of PC requests. Refer to Table 1 for frequency pattern interactions.

5.2 IP Request Sources

5.2.1 Internal Site Visitors

Refer to Figure 2.

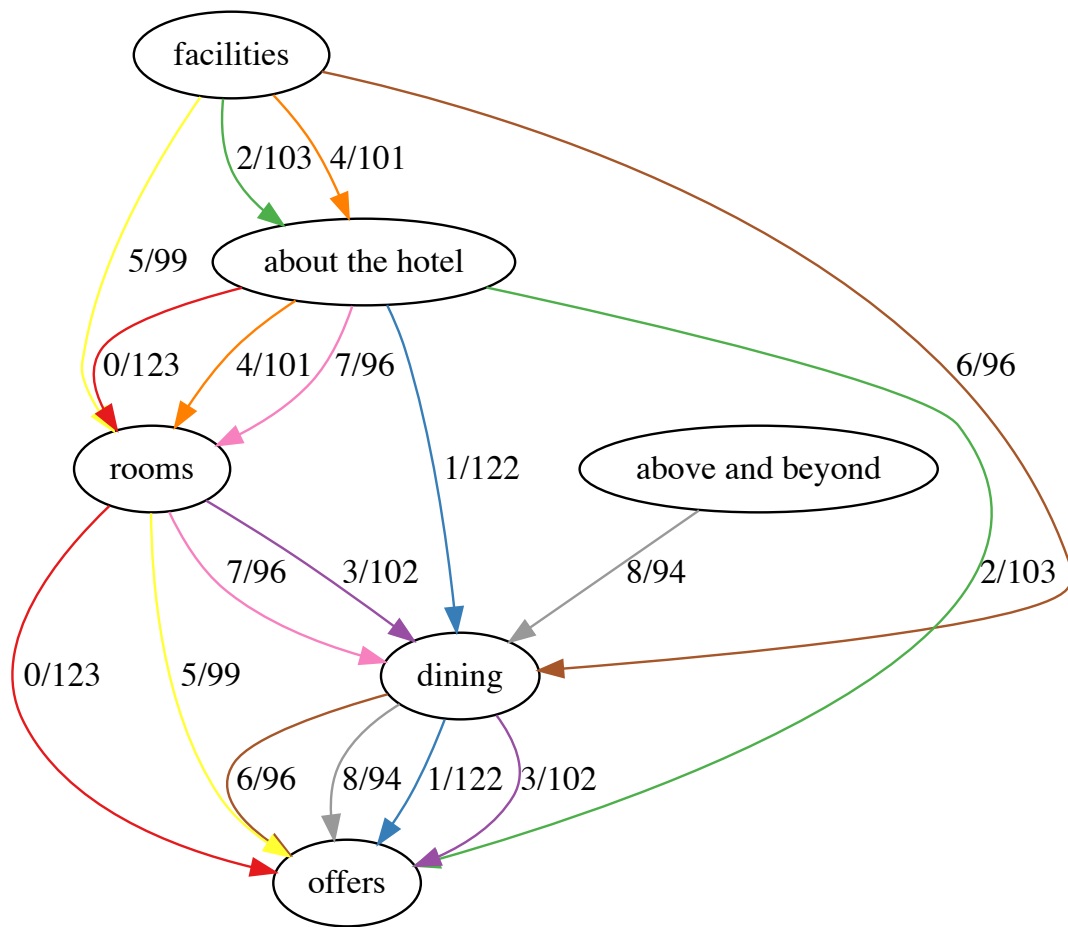


Figure 2: Directional network graph visualising frequency patterns of internal visitors made on an internal IP range. Refer to Table 2 for frequency pattern interactions.

5.2.2 External Site Visitors

Refer to Figure 3.

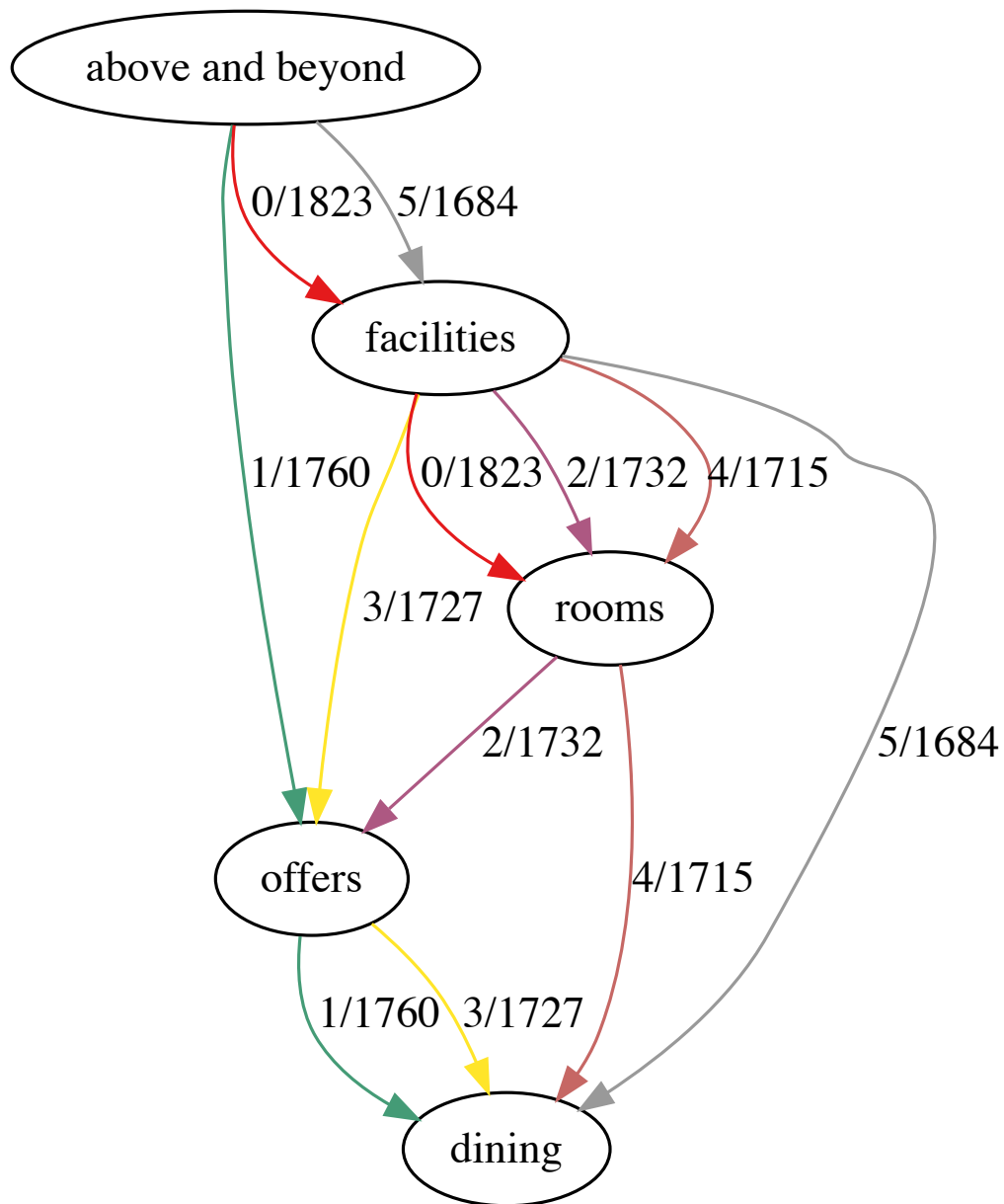


Figure 3: Directional network graph visualising frequency patterns of external visitors made on a non-internal IP range. Refer to Table 3 for frequency pattern interactions.

5.3 Top Three Countries

5.3.1 Hong Kong Visitors

Refer to Figure 4.

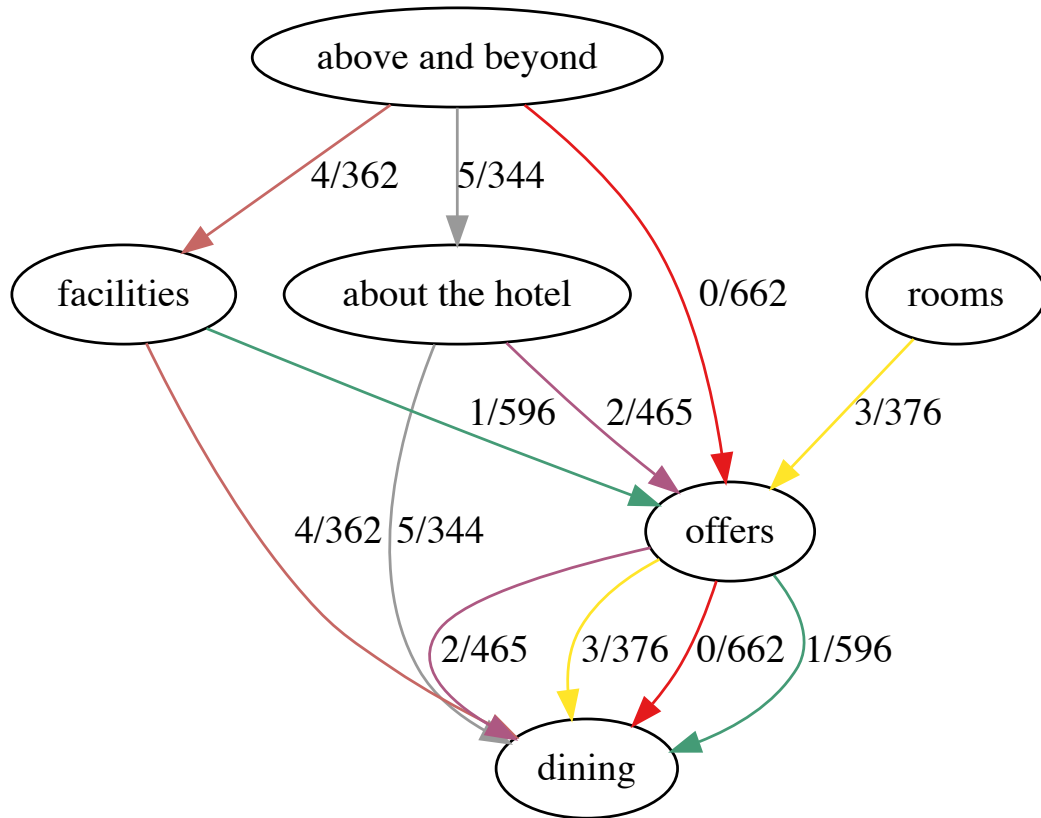


Figure 4: Directional network graph visualising frequency patterns of visitors from Hong Kong. Refer to Table 4 for frequency pattern interactions.

5.3.2 USA Visitors

Refer to Figure 5.

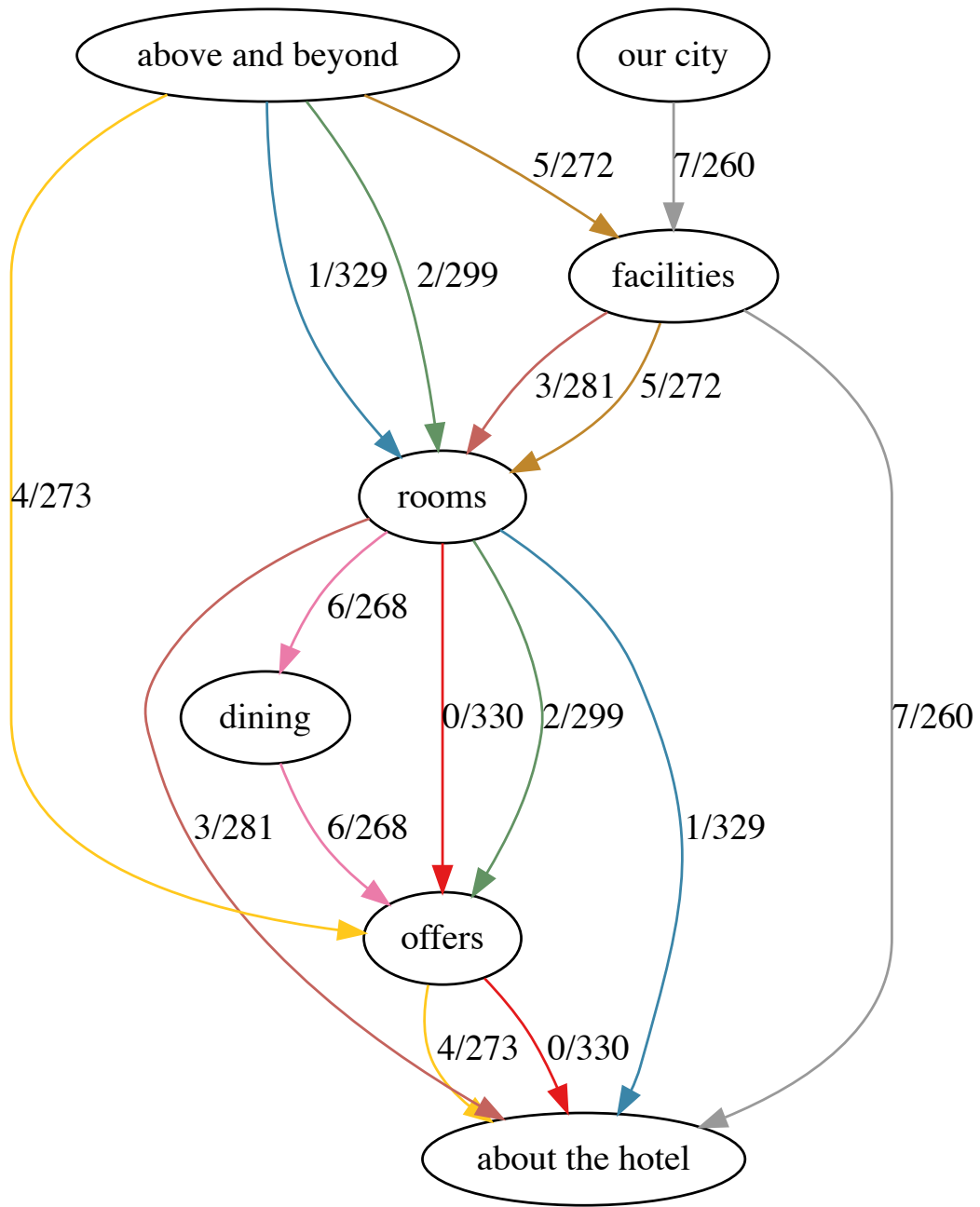


Figure 5: Directional network graph visualising frequency patterns of visitors from the United States. Refer to Table 5 for frequency pattern interactions.

5.3.3 Australian Visitors

Refer to Figure 6.

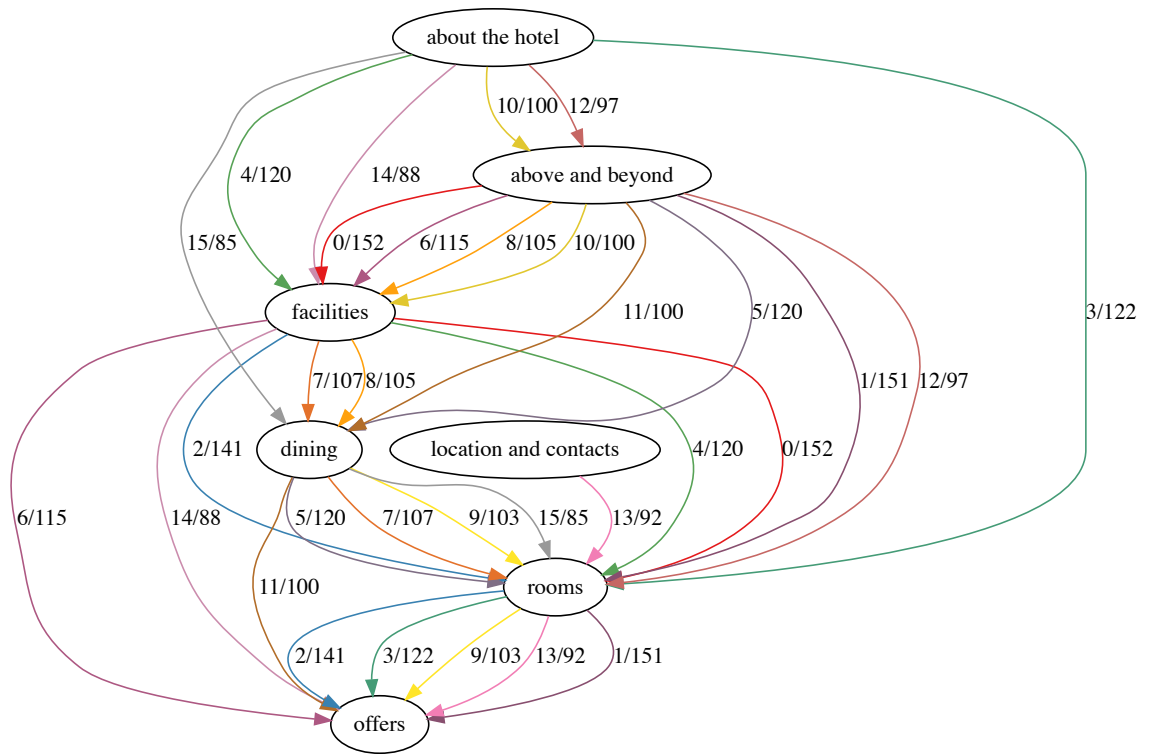


Figure 6: Directional network graph visualising frequency patterns of visitors from Australia. Refer to Table 6 for frequency pattern interactions.

5.4 Platform Categories

5.4.1 PC Visitors

Refer to Figure 7.

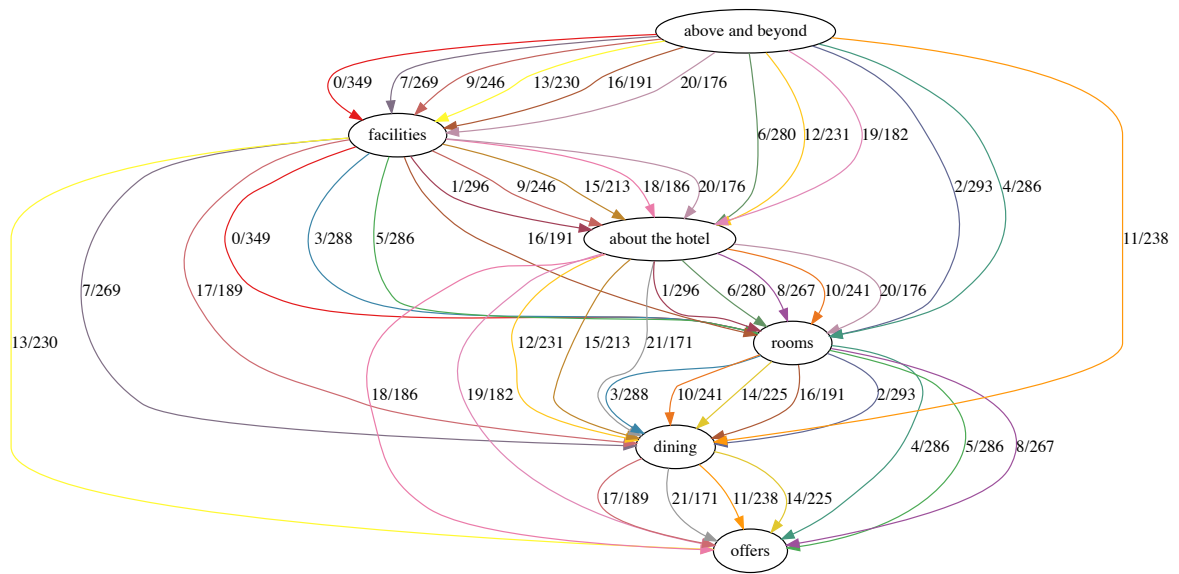


Figure 7: Directional network graph visualising frequency patterns of requests made by PCs. Refer to Table 7 for frequency pattern interactions.

5.4.2 Smartphone Visitors

Refer to Figure 8.

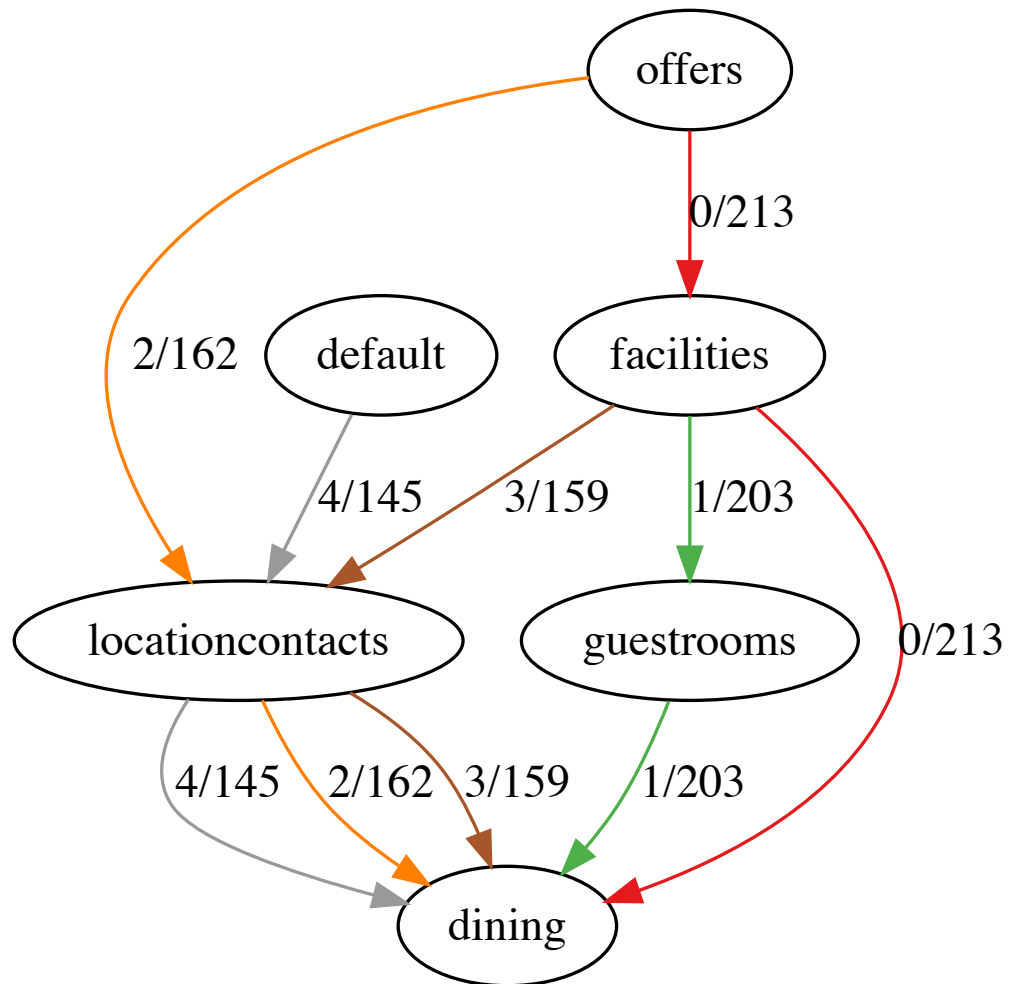


Figure 8: Directional network graph visualising frequency patterns of requests made on smart-phones. Refer to Table 9 for frequency pattern interactions.

5.4.3 Tablet Visitors

Refer to Figure 9.

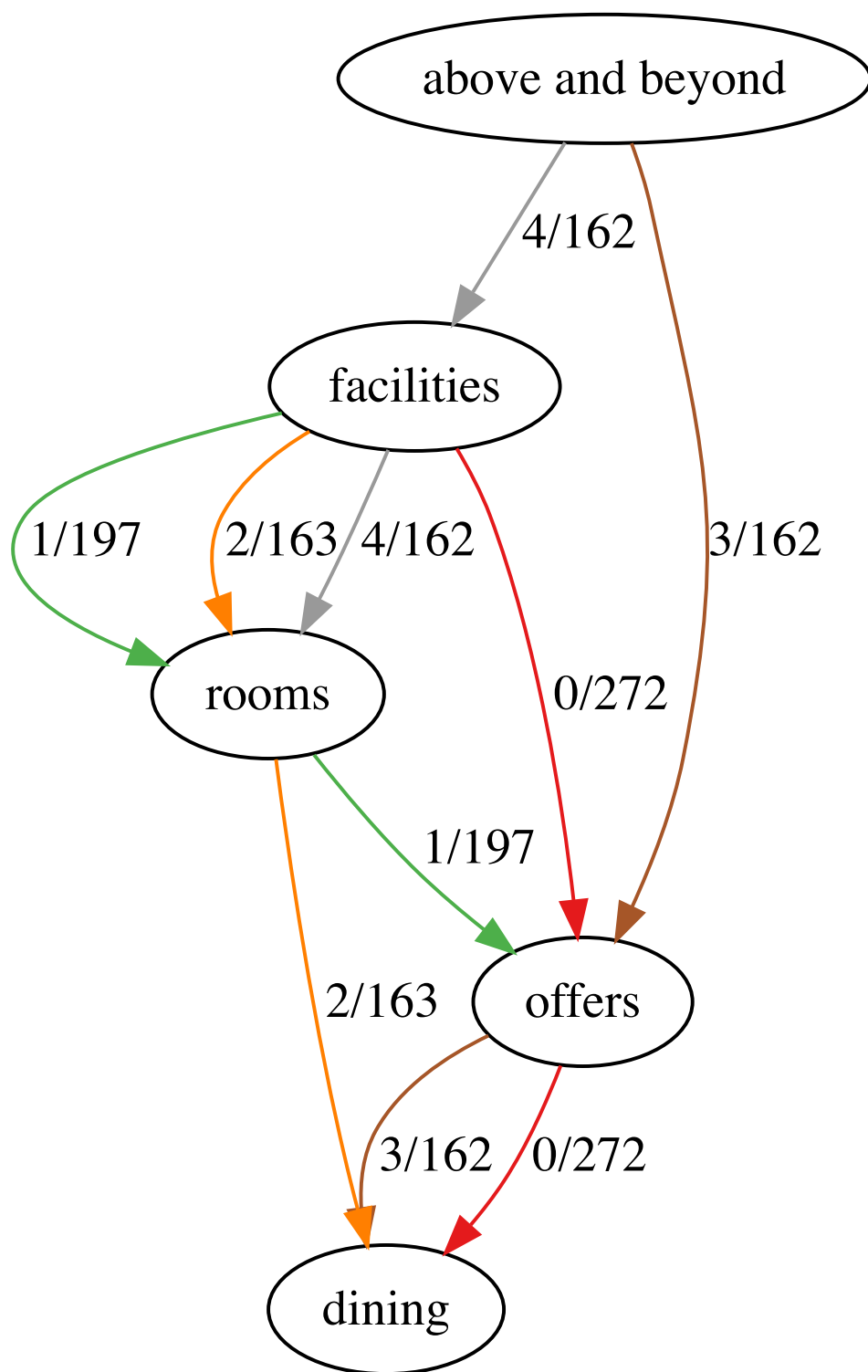


Figure 9: Directional network graph visualising frequency patterns of requests made on tablets.

Refer to Table 8 for frequency pattern interactions

5.4.4 Bots Visitors

Refer to Figure 10.

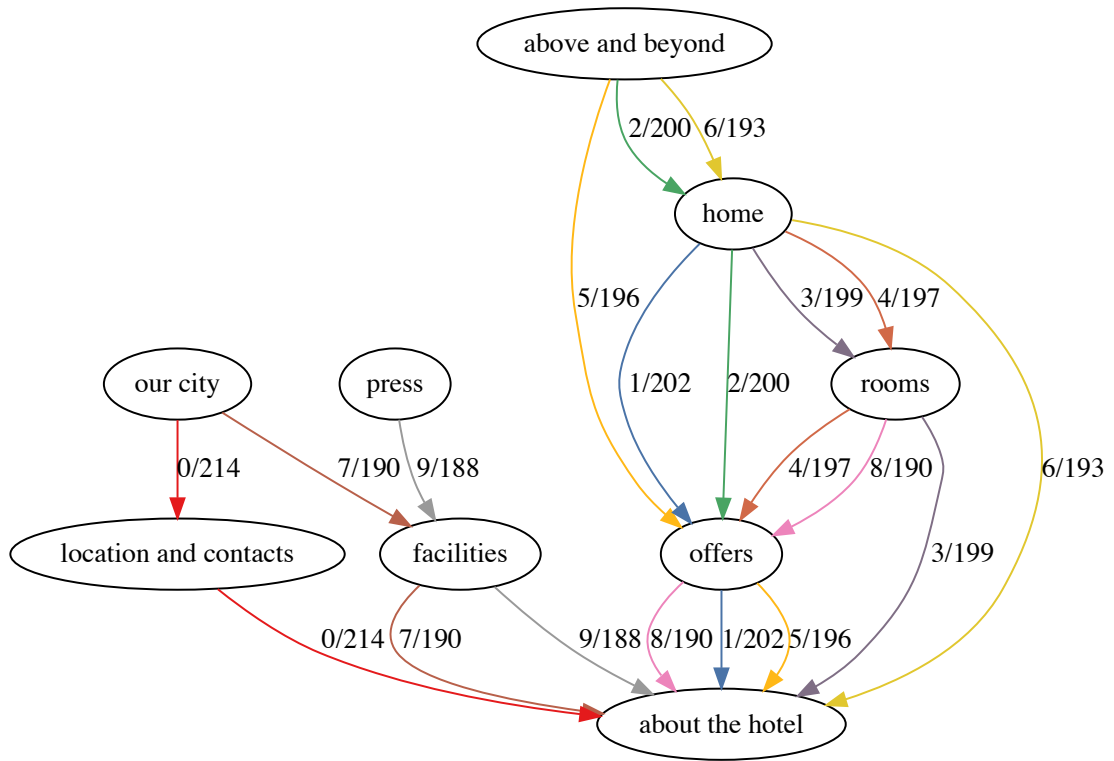


Figure 10: Directional network graph visualising frequency patterns of requests made by bots. Refer to Table 10 for frequency pattern interactions.

References

- Fielding, R., J. Gettys, J. Mogul, H. Frystyk, L. Masinter, P. Leach, and T. Berners-Lee (1999). RFC 2616, Hypertext Transfer Protocol – HTTP/1.1. Retrieved 11 May 2017, <<http://www.rfc.net/rfc2616.html>>.
- Hallam-Baker, P. and B. Behlendorf (1998). Extended Log File Format. Retrieved 11 May 2017, <<https://www.w3.org/TR/WD-logfile-960221.html>>.
- Microsoft Corporation (2003). W3C Extended Log File Format (IIS 6.0). Retrieved 11 May 2017, <<https://www.microsoft.com/technet/prodtechnol/WindowsServer2003/Library/IIS/676400bc-8969-4aa7-851a-9319490a9bbb.mspx?mfr=true>>.
- Microsoft Corporation (2017). Description of Microsoft Internet Information Services (IIS) 5.0 and 6.0 status codes. Retrieved 11 May 2017, <<https://support.microsoft.com/en-au/help/318380/description-of-microsoft-internet-information-services-iis-5.0-and-6.0-status-codes>>.
- Wolf, M. and C. Wicksteed. Date and Time Formats. <https://www.w3.org/TR/NOTE-datetime>.

A Additional Tables

Below are tables of frequency pattern results for each section identified in Section 5.

Table 2: Internal Request Frequency Patterns

Sequence	From	To	Frequency
0	about the hotel	→ rooms	123
0	rooms	→ offers	123
1	about the hotel	→ dining	122
1	dining	→ offers	122
2	facilities	→ about the hotel	103
2	about the hotel	→ offers	103
3	rooms	→ dining	102
3	dining	→ offers	102
4	facilities	→ about the hotel	101
4	about the hotel	→ rooms	101
5	facilities	→ rooms	99
5	rooms	→ offers	99
6	facilities	→ dining	96
6	dining	→ offers	96
7	about the hotel	→ rooms	96
7	rooms	→ dining	96
8	above and beyond	→ dining	94
8	dining	→ offers	94

Table 3: External Request Frequency Patterns

Sequence	From	To	Frequency
0	above and beyond	→ facilities	1823

Continued on next page...

Table 3 (*continued from Page 30*): External Request Frequency Patterns

Sequence	From	To	Frequency
0	facilities	→ rooms	1823
1	above and beyond	→ offers	1760
1	offers	→ dining	1760
2	facilities	→ rooms	1732
2	rooms	→ offers	1732
3	facilities	→ offers	1727
3	offers	→ dining	1727
4	facilities	→ rooms	1715
4	rooms	→ dining	1715
5	above and beyond	→ facilities	1684
5	facilities	→ dining	1684

Table 4: Hong Kong Frequency Patterns

Sequence	From	To	Frequency
0	above and beyond	→ offers	662
0	offers	→ dining	662
1	facilities	→ offers	596
1	offers	→ dining	596
2	about the hotel	→ offers	465
2	offers	→ dining	465
3	rooms	→ offers	376
3	offers	→ dining	376
4	above and beyond	→ facilities	362
4	facilities	→ dining	362

Continued on next page...

Table 4 (*continued from Page 31*): Hong Kong Frequency Patterns

Sequence	From	To	Frequency
5	above and beyond	→ about the hotel	344
5	about the hotel	→ dining	344

Table 5: USA Request Frequency Patterns

Sequence	From	To	Frequency
0	rooms	→ offers	330
0	offers	→ about the hotel	330
1	above and beyond	→ rooms	329
1	rooms	→ about the hotel	329
2	above and beyond	→ rooms	299
2	rooms	→ offers	299
3	facilities	→ rooms	281
3	rooms	→ about the hotel	281
4	above and beyond	→ offers	273
4	offers	→ about the hotel	273
5	above and beyond	→ facilities	272
5	facilities	→ rooms	272
6	rooms	→ dining	268
6	dining	→ offers	268
7	our city	→ facilities	260
7	facilities	→ about the hotel	260

Table 6: Australian Request Frequency Patterns

Sequence	From	To	Frequency
0	above and beyond	→ facilities	152
0	facilities	→ rooms	152
1	above and beyond	→ rooms	151
1	rooms	→ offers	151
2	facilities	→ rooms	141
2	rooms	→ offers	141
3	about the hotel	→ rooms	122
3	rooms	→ offers	122
4	about the hotel	→ facilities	120
4	facilities	→ rooms	120
5	above and beyond	→ dining	120
5	dining	→ rooms	120
6	above and beyond	→ facilities	115
6	facilities	→ offers	115
7	facilities	→ dining	107
7	dining	→ rooms	107
8	above and beyond	→ facilities	105
8	facilities	→ dining	105
9	dining	→ rooms	103
9	rooms	→ offers	103
10	about the hotel	→ above and beyond	100
10	above and beyond	→ facilities	100
11	above and beyond	→ dining	100
11	dining	→ offers	100
12	about the hotel	→ above and beyond	97
12	above and beyond	→ rooms	97

Continued on next page...

Table 6 (*continued from Page 33*): Australian Request Frequency Patterns

Sequence	From	To	Frequency
13	location and contacts	→ rooms	92
13	rooms	→ offers	92
14	about the hotel	→ facilities	88
14	facilities	→ offers	88
15	about the hotel	→ dining	85
15	dining	→ rooms	85

Table 7: PC Request Frequency Patterns

Sequence	From	To	Frequency
0	above and beyond	→ facilities	349
0	facilities	→ rooms	349
1	facilities	→ about the hotel	296
1	about the hotel	→ rooms	296
2	above and beyond	→ rooms	293
2	rooms	→ dining	293
3	facilities	→ rooms	288
3	rooms	→ dining	288
4	above and beyond	→ rooms	286
4	rooms	→ offers	286
5	facilities	→ rooms	286
5	rooms	→ offers	286
6	above and beyond	→ about the hotel	280
6	about the hotel	→ rooms	280
7	above and beyond	→ facilities	269

Continued on next page...

Table 7 (continued from Page 34): PC Request Frequency
Patterns

Sequence	From	To	Frequency
7	facilities	→ dining	269
8	about the hotel	→ rooms	267
8	rooms	→ offers	267
9	above and beyond	→ facilities	246
9	facilities	→ about the hotel	246
10	about the hotel	→ rooms	241
10	rooms	→ dining	241
11	above and beyond	→ dining	238
11	dining	→ offers	238
12	above and beyond	→ about the hotel	231
12	about the hotel	→ dining	231
13	above and beyond	→ facilities	230
13	facilities	→ offers	230
14	rooms	→ dining	225
14	dining	→ offers	225
15	facilities	→ about the hotel	213
15	about the hotel	→ dining	213
16	above and beyond	→ facilities	191
16	facilities	→ rooms	191
16	rooms	→ dining	191
17	facilities	→ dining	189
17	dining	→ offers	189
18	facilities	→ about the hotel	186
18	about the hotel	→ offers	186
19	above and beyond	→ about the hotel	182

Continued on next page...

Table 7 (continued from Page 34): PC Request Frequency Patterns

Sequence	From	To	Frequency
19	about the hotel	→ offers	182
20	above and beyond	→ facilities	176
20	facilities	→ about the hotel	176
20	about the hotel	→ rooms	176
21	about the hotel	→ dining	171
21	dining	→ offers	171

Table 8: Tablet Request Frequency Patterns

Sequence	From	To	Frequency
0	facilities	→ offers	272
0	offers	→ dining	272
1	facilities	→ rooms	197
1	rooms	→ offers	197
2	facilities	→ rooms	163
2	rooms	→ dining	163
3	above and beyond	→ offers	162
3	offers	→ dining	162
4	above and beyond	→ facilities	162
4	facilities	→ rooms	162

Table 9: Smartphone Request Frequency Patterns

Sequence	From	To	Frequency
0	offers	→ facilities	213
0	facilities	→ dining	213
1	facilities	→ guestrooms	203
1	guestrooms	→ dining	203
2	offers	→ locationcontacts	162
2	locationcontacts	→ dining	162
3	facilities	→ locationcontacts	159
3	locationcontacts	→ dining	159
4	default	→ locationcontacts	145
4	locationcontacts	→ dining	145

Table 10: Bots Request Frequency Patterns

Sequence	From	To	Frequency
0	our city	→ location and contacts	214
0	location and contacts	→ about the hotel	214
1	home	→ offers	202
1	offers	→ about the hotel	202
2	above and beyond	→ home	200
2	home	→ offers	200
3	home	→ rooms	199
3	rooms	→ about the hotel	199
4	home	→ rooms	197
4	rooms	→ offers	197
5	above and beyond	→ offers	196
5	offers	→ about the hotel	196

Continued on next page...

Table 10 (*continued from Page 37*): Bots Request Frequency
Patterns

Sequence	From	To	Frequency
6	above and beyond	→ home	193
6	home	→ about the hotel	193
7	our city	→ facilities	190
7	facilities	→ about the hotel	190
8	rooms	→ offers	190
8	offers	→ about the hotel	190
9	press	→ facilities	188
9	facilities	→ about the hotel	188

B Data Dictionary

B.1 Dataset Description

Table 11: Dataset Description

Key	Entry
Name	Hotel TULIP Web Log Dataset
Size	17.06 GB (954.7 MB compressed)
Release Date	30/4/17
Attributes	14
No. Records	73,368,256
Provider	Dr Beer Guts, CIO, Hotel TULIP (Information Technology Division) ³
Privacy	Confidential ⁴

B.2 Contact Information

Table 12: Dataset Contact Information

Key	Entry
Prepared by	Team-SIT742
Point of Contact	Alex Cummaudo <ca@deakin.edu.au> Jake Renzella <jake.renzella@deakin.edu.au>
Team Members	Alex Cummaudo <ca@deakin.edu.au> Jake Renzella <jake.renzella@deakin.edu.au>

³Download URL: <https://d2l.deakin.edu.au/d2l/1e/content/520519/topics/files/download/3482057/DirectFileTopicDownload>, <https://d2l.deakin.edu.au/d2l/1e/content/520519/viewContent/3482057/View?ou=520519>.

⁴Exclusively available for educational purposes only for the Deakin University unit SIT742. Redistribution is prohibited.

B.3 Data Dictionary

Table 13: Data Dictionary

Attribute Name	Data Type	Data Subtype	Description	Examples	Notes
date	MC	DATE - Date	Date: Date when request occurred	2014-08-01	UTC time zone; ISO 8601 date format (YYYY-MM-DD)
time	MC	DATE - Time	Time: Time when request occurred	09:51:23	UTC time zone; ISO 8601 24-hr time format (hh:mm:ss)
s-ip	CN	ADDR - Address - IP Address	Server IP Address: IP of the server generating the log responding to the request	10.130.0.12	N/A
cs-uri-stem	CN	URL - Uniform Resource Identifier (URI)	URI Stem: Stem portion of the full URI made in the client to server request	/sitecore	N/A
cs-uri-query	CN	URL - Uniform Resource Identifier (URI)	URI Query: Query portion of the full URI made in the client to server request	cmd=GetTreeview	May be empty; Query string is specifically matched as URI. See Hallam-Baker & Behlendorf (1998).

Continued on next page...

Table 13 (continued from Page 40): Data Dictionary

Attribute Name	Data Type	Data Subtype	Description	Examples	Notes
s-port	CN	ADDR - Address - Port Number	Server Port: Port for the server which request is made	80	N/A
cs-username	CN	STR - Free String	User Name: Name of authenticated user accessing the server	-	Anonymous users are represented with a hyphen.
c-ip	CN	ADDR - Address - IP Address	Client IP Address: IP of the client making the request	10.120.7.23	N/A
sc-substatus	CN	ID - Identifica- tion - Microsoft IIS Named HTTP Response Sub status Code	Protocol Sub status: The sub status error code	0	Refer to Microsoft (2017) for specific codes. A zero indicates no sub status error code.
sc-win32-status	CN	ID - Identifica- tion - Microsoft IIS Named Win- dows Code	Win32 Status: The Windows status code	0	Generally not applicable to non-Windows devices; defaults to zero.

Continued on next page...

Table 13 (*continued from Page 40*): Data Dictionary

Attribute Name	Data Type	Data Subtype	Description	Examples	Notes
time-taken	MC	DATE - Time - From zero	Duration: The time taken for the request to complete.	39	Measured in milliseconds.

C Extrapolation Results

Attached on the following pages are the results from Databricks. You may also interact with this online on Databricks.