# Comp1204 Unix Coursework

Callum Anderson
crea1g15

February 24, 2016

# Contents

# 1 Scripts

## 1.1 Review Counting

This script searches a folder (given as an argument) for .dat files and counts the number of reviews in each file.

```bash
#!/bin/bash

for filename in $1/*.dat; do
    name=$(echo $filename | cut -d/ -f2 | cut -d. -f1)
    echo "$name $(grep Author $filename | wc -l)"
done | sort -t" " -k2nr
```

countreviews.sh

Line 4 removes the extension and directory path from the filename, so the hotel name itself can be printed next to the review count.

The for loop goes through each file, and then pipes a string of format "$hotelname $reviewcount" into sort, which prints out the sorted list of hotels.

## 1.2 Average Review Score

This script runs through all .dat files in the given folder and finds the average review score for each hotel, then prints the sorted list to the terminal.

```bash
#!/bin/bash

for filename in $1/*.dat; do
    echo $filename
    name=$(echo $filename | cut -d/ -f2 | cut -d. -f1)
    sum=0
    count=0
    for i in $(grep "<Overall>" $filename | sed "s/<Overall>//g"); do
        count=$(($count+1))
        i=`echo $i | tr -d $'\r'`
        sum=$((sum+i))
    done
    avg=`echo "scale=2; $sum/$count" | bc`
    echo "$name $avg"
done | sort -t" " -k2nr
```

averagereviews.sh

Line 8 loops through each line in the file containing "¡Overall¿", removing the 'overall' tag from it so it's just a number. Line 10 removes any linebreak characters that still remain.

## 1.3 t-Testing

```bash
#!/bin/bash

argnum=$((0))
for arg in "$@"; do

    argnum=$(($argnum+1))

    #get scores for that hotel and mean
```

```bash
 9     count[$argnum]=$((0))
10     mean[$argnum]=$((0))
11     for i in $(grep "<Overall>" $arg | sed "s/[^0-9]*//g"); do
12       count[$argnum]=$((${count[$argnum]}+1))
13       scores[${count[argnum]}]=$(echo $i)
14       mean[$argnum]=$((${mean[$argnum]} + ${scores[${count[$argnum]}]}))
15     done
16     mean[$argnum]=`echo "scale=20; ${mean[$argnum]}/${count[$argnum]}" |
           bc`
17
18
19     #get variance
20     sigma=$((0))
21     for i in "${scores[@]}"; do
22       sigma=`echo "scale=20; $sigma+($i-${mean[$argnum]})^2" | bc`
23     done
24
25     var[$argnum]=`echo "scale=20; (1/(${count[$argnum]}-1))*$sigma" | bc`
26 done
27
28 #Calculate t-stat
29 sx1x2=`echo "scale=20; sqrt((((${count[1]}-1)*${var[1]} +
        (${count[2]}-1)*${var[2]}) / (${count[1]}+${count[2]}-2))" | bc`
30 t=`echo "scale=20; (${mean[1]}-${mean[2]}) /
        ($sx1x2*sqrt((1/${count[1]})+(1/${count[2]})))" | bc`
31
32 #print all the stuff
33 echo "t: $(printf %.2f $t)"
34 sd1=$(echo "scale=2; sqrt(${var[1]})" | bc)
35 sd2=$(echo "scale=2; sqrt(${var[2]})" | bc)
36 echo "Mean $(echo $1 | cut -d/ -f2 | sed 's/\..*//'): $(printf %.2f
        ${mean[1]}), SD: $(printf %.2f $sd1)"
37 echo "Mean $(echo $2 | cut -d/ -f2 | sed 's/\..*//'): $(printf %.2f
        ${mean[2]}), SD: $(printf %.2f $sd2)"
38
39 critvalue=1.972731033408872
40
41 #bash only does int arithmetic, so multiply the t-value and
42 #crit value by big number to allow meaningful comparison
43 t=`echo "scale=0; $t*1000000000" | bc`
44 t=${t%.*}
45 critvalue=`echo "scale=0; $critvalue*1000000000" | bc`
46 critvalue=${critvalue%.*}
47 negcritvalue=$(( $critvalue*-1))
48
49 #if inside critical range, return 1
50 if [ $t -gt $critvalue ] || [ $t -lt $negcritvalue ]
51 then
52   sig=$((1))
53 else
54   sig=$((0))
55 fi
56
57 echo $sig
```

statistical_sig.sh

# 2 Hypothesis Testing

# 3 Discussion