

Data Wrangling (Data Preprocessing)

Practical assessment 2

Andre Gunawan - s3885488

Required packages

```
# This is the R chunk for the required packages
library(readr)
library(dplyr)
```

```
##
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':
##
##   filter, lag
```

```
## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union
```

```
library(tidyr)
library(MVN)
```

```
## Warning: package 'MVN' was built under R version 4.0.5
```

```
## Registered S3 method overwritten by 'GGally':
##   method from
##   +.gg      ggplot2
```

```
## sROC 0.1-2 loaded
```

```
library(forecast)
```

```
## Warning: package 'forecast' was built under R version 4.0.5
```

```
## Registered S3 method overwritten by 'quantmod':
##   method          from
##   as.zoo.data.frame zoo
```

Executive Summary

This report provides a complete data preprocessing on the unemployment rates between males and females in different countries around the world. There are several stages which I do in this data preprocessing. Firstly, I imported both data into R studio using readr package separately. Secondly, I tidy up the data using pivot_longer (). After tidying up the data, I merge these two data using the inner_join () function. Thirdly, I do a check on the data types using the str (), class (), attributes () function and then convert the data types to the proper type. Moreover, I also do labeling on one of the factor variables. Fourthly, I create/mutate a new variable from the combination of data I have. Fifthly, I scan for missing values, special values, obvious errors, and outliers then applying a proper methodology to deal with them. Lastly, I did the data transformation to reduce skewness and make the data normally distributed.

Data

- Data description: The datasets that I used represent the unemployment rates between males and females in different countries around the world. The first dataset and second dataset represents the males unemployment rate and females unemployment rate respectively.
- Both of the data was taken from The World Bank website. Male unemployment rate: <https://databank.worldbank.org/reports.aspx?source=2&series=SL.UEM.TOTL.MA.ZS&country=> (https://databank.worldbank.org/reports.aspx?source=2&series=SL.UEM.TOTL.MA.ZS&country=) Female unemployment rate: <https://databank.worldbank.org/reports.aspx?source=2&series=SL.UEM.TOTL.FE.ZS&country=> (https://databank.worldbank.org/reports.aspx?source=2&series=SL.UEM.TOTL.FE.ZS&country=)
- Variable description: Both of these datasets have the same variables, the only difference is that the first data represents males and the second data represents females.
 - Country name: name of the country
 - Country code: unique code of the country
 - Year 2011-2019: represent unemployment rate of different year

```
#Import male unemployment rate using readr package.
male_unemployment_rate <- read_csv("C:/Users/andre/OneDrive/Desktop/Data Wrangling/Practical
Assessment 2/Practical Assessment 2/Male unemployment rate.csv")
```

```
##
## -- Column specification -----
## cols(
##   `Country Name` = col_character(),
##   `Country Code` = col_character(),
##   `2011 [YR2011]` = col_character(),
##   `2012 [YR2012]` = col_character(),
##   `2013 [YR2013]` = col_character(),
##   `2014 [YR2014]` = col_character(),
##   `2015 [YR2015]` = col_character(),
##   `2016 [YR2016]` = col_character(),
##   `2017 [YR2017]` = col_character(),
##   `2018 [YR2018]` = col_character(),
##   `2019 [YR2019]` = col_character()
## )
```

```
#Import female unemployment rate using readr package.
female_unemployment_rate <- read_csv("C:/Users/andre/OneDrive/Desktop/Data Wrangling/Practica
l Assessment 2/Practical Assessment 2/Female unemployment rate.csv")
```

```
##
## -- Column specification -----
## cols(
##   `Country Name` = col_character(),
##   `Country Code` = col_character(),
##   `2011 [YR2011]` = col_character(),
##   `2012 [YR2012]` = col_character(),
##   `2013 [YR2013]` = col_character(),
##   `2014 [YR2014]` = col_character(),
##   `2015 [YR2015]` = col_character(),
##   `2016 [YR2016]` = col_character(),
##   `2017 [YR2017]` = col_character(),
##   `2018 [YR2018]` = col_character(),
##   `2019 [YR2019]` = col_character()
## )
```

#Display the first 6 rows of male unemployment rate using head() function.
head(male_unemployment_rate)

Country Name <chr>	Country Code <chr>	2011 [YR2011] <chr>	2012 [YR2012] <chr>	2013 [YR2013] <chr>	2014 [YR2014] <chr>
Afghanistan	AFG	10.89999962	10.88000011	10.89000034	10.77999973
Albania	ALB	13.27000046	14.77999973	17.62999916	19.80999947
Algeria	DZA	8.399999619	9.569999695	8.289999962	8.989999771
American Samoa	ASM
Andorra	AND
Angola	AGO	7.010000229	6.929999828	6.989999771	6.909999847

6 rows | 1-7 of 11 columns

#Display the first 6 rows of female unemployment rate using head() function.
head(female_unemployment_rate)

Country Name <chr>	Country Code <chr>	2011 [YR2011] <chr>	2012 [YR2012] <chr>	2013 [YR2013] <chr>	2014 [YR2014] <chr>
Afghanistan	AFG	14.78999996	14.85999966	14.69999981	14.52999973
Albania	ALB	13.75	11.46000004	13.34000015	15.44999981
Algeria	DZA	17.12999916	17.01000023	16.27000046	15.64000034
American Samoa	ASM
Andorra	AND
Angola	AGO	7.71999979	7.809999943	7.78000021	7.710000038

6 rows | 1-7 of 11 columns

Step explanation:

- Firstly, I import both of the data using readr function.
- Secondly, I used head() function to display the first 6 observations.
- Before merging the data, I will tidy up both of the male_unemployment_rate and female_unemployment_rate dataset above and give them the names tidy_male_unemployment_rate and tidy_female_unemployment_rate. Further explanation can be found on the Tidy & Manipulate Data I section below.

Tidy & Manipulate Data I

- There are three interrelated rules which make a dataset tidy by Wickham and Grolemund (2016). In tidy data: -Each variable must have its own column. -Each observation must have its own row. -Each value must have its own cell.
- The data above is not tidy because the years 2011 - 2019 are separated in 9 variables, instead they should be combined into one variable named "Year". Therefore, I used pivot_longer () function below to combine them into a single variable named "Year".

```
#Combine 2011 - 2019 variables into one variable named "Year" for male dataset. Then, assign it to variable tidy_male_unemployment_rate.
tidy_male_unemployment_rate <- male_unemployment_rate %>% pivot_longer(cols = -c("Country Name", "Country Code"), names_to = "Year", values_to = "Unemployment rate") %>% rename(`Male unemployment rate` = `Unemployment rate`)
```

```
#Combine 2011 - 2019 variables into one variable named "Year" for female dataset. Then, assign it to variable tidy_female_unemployment_rate.
tidy_female_unemployment_rate <- female_unemployment_rate %>% pivot_longer(cols = -c("Country Name", "Country Code"), names_to = "Year", values_to = "Unemployment rate") %>% rename(`Female unemployment rate` = `Unemployment rate`)
```

```
#Display the first several rows of tidy_male_unemployment_rate using head() function.
head(tidy_male_unemployment_rate)
```

Country Name <chr>	Country Code <chr>	Year <chr>	Male unemployment rate <chr>
Afghanistan	AFG	2011 [YR2011]	10.89999962
Afghanistan	AFG	2012 [YR2012]	10.88000011
Afghanistan	AFG	2013 [YR2013]	10.89000034
Afghanistan	AFG	2014 [YR2014]	10.77999973
Afghanistan	AFG	2015 [YR2015]	10.68000031
Afghanistan	AFG	2016 [YR2016]	10.56999969
6 rows			

```
#Display the first several rows of tidy_females_unemployment_rate using head() function.
head(tidy_female_unemployment_rate)
```

Country Name <chr>	Country Code <chr>	Year <chr>	Female unemployment rate <chr>
Afghanistan	AFG	2011 [YR2011]	14.78999996
Afghanistan	AFG	2012 [YR2012]	14.85999966
Afghanistan	AFG	2013 [YR2013]	14.69999981
Afghanistan	AFG	2014 [YR2014]	14.52999973
Afghanistan	AFG	2015 [YR2015]	14.44999981
Afghanistan	AFG	2016 [YR2016]	14.32999992

6 rows

Step explanation:

- I combined 2011 - 2019 variables into one variable named Year.
- Then, I recheck the data using head(). The 2011 - 2019 variable has been combined into one variable named Year.

Data

```
#Merged the tidy_male_unemployment_rate with tidy_female_unemployment_rate. Then, named it merged_malefemale_UR.
```

```
merged_malefemale_UR <- tidy_male_unemployment_rate %>% inner_join(tidy_female_unemployment_rate, by = c("Country Code", "Year", "Country Name"))
```

```
#Display the first several rows of merged_malefemale_UR.
```

```
head(merged_malefemale_UR)
```

Country Name <chr>	Country Code <chr>	Year <chr>	Male unemployment rate <chr>	Female unemployment rate <chr>
Afghanistan	AFG	2011 [YR2011]	10.89999962	14.78999996
Afghanistan	AFG	2012 [YR2012]	10.88000011	14.85999966
Afghanistan	AFG	2013 [YR2013]	10.89000034	14.69999981
Afghanistan	AFG	2014 [YR2014]	10.77999973	14.52999973
Afghanistan	AFG	2015 [YR2015]	10.68000031	14.44999981
Afghanistan	AFG	2016 [YR2016]	10.56999969	14.32999992

6 rows

Step explanation:

- Merged data using inner_join function and give it a new name "merged_malefemale_UR".
- Use head to display the first 6 observations of the merged_malefemale_UR.

Understand

```
#Inspect the data structures before conversion
str(merged_malefemale_UR)
```

```
## tibble [2,376 x 5] (S3: tbl_df/tbl/data.frame)
## $ Country Name      : chr [1:2376] "Afghanistan" "Afghanistan" "Afghanistan" "Afgha
nistan" ...
## $ Country Code      : chr [1:2376] "AFG" "AFG" "AFG" "AFG" ...
## $ Year              : chr [1:2376] "2011 [YR2011]" "2012 [YR2012]" "2013 [YR2013]"
"2014 [YR2014]" ...
## $ Male unemployment rate : chr [1:2376] "10.89999962" "10.88000011" "10.89000034" "10.77
999973" ...
## $ Female unemployment rate: chr [1:2376] "14.78999996" "14.85999966" "14.69999981" "14.52
999973" ...
```

```
#Inspect the variable types before conversion
class(merged_malefemale_UR$`Country Name`)
```

```
## [1] "character"
```

```
class(merged_malefemale_UR$`Country Code`)
```

```
## [1] "character"
```

```
class(merged_malefemale_UR$Year)
```

```
## [1] "character"
```

```
class(merged_malefemale_UR$`Male unemployment rate`)
```

```
## [1] "character"
```

```
class(merged_malefemale_UR$`Female unemployment rate`)
```

```
## [1] "character"
```

```
#Inspect the attributes
attributes(merged_malefemale_UR)
```

```
## $names
## [1] "Country Name"          "Country Code"
## [3] "Year"                  "Male unemployment rate"
## [5] "Female unemployment rate"
##
## $row.names
##      [1]      1      2      3      4      5      6      7      8      9     10     11     12     13     14
##     [15]     15     16     17     18     19     20     21     22     23     24     25     26     27     28
##     [29]     29     30     31     32     33     34     35     36     37     38     39     40     41     42
##     [43]     43     44     45     46     47     48     49     50     51     52     53     54     55     56
##     [57]     57     58     59     60     61     62     63     64     65     66     67     68     69     70
##     [71]     71     72     73     74     75     76     77     78     79     80     81     82     83     84
##     [85]     85     86     87     88     89     90     91     92     93     94     95     96     97     98
##     [99]     99    100    101    102    103    104    105    106    107    108    109    110    111    112
##    [113]    113    114    115    116    117    118    119    120    121    122    123    124    125    126
##    [127]    127    128    129    130    131    132    133    134    135    136    137    138    139    140
##    [141]    141    142    143    144    145    146    147    148    149    150    151    152    153    154
##    [155]    155    156    157    158    159    160    161    162    163    164    165    166    167    168
##    [169]    169    170    171    172    173    174    175    176    177    178    179    180    181    182
##    [183]    183    184    185    186    187    188    189    190    191    192    193    194    195    196
##    [197]    197    198    199    200    201    202    203    204    205    206    207    208    209    210
##    [211]    211    212    213    214    215    216    217    218    219    220    221    222    223    224
##    [225]    225    226    227    228    229    230    231    232    233    234    235    236    237    238
##    [239]    239    240    241    242    243    244    245    246    247    248    249    250    251    252
##    [253]    253    254    255    256    257    258    259    260    261    262    263    264    265    266
##    [267]    267    268    269    270    271    272    273    274    275    276    277    278    279    280
##    [281]    281    282    283    284    285    286    287    288    289    290    291    292    293    294
##    [295]    295    296    297    298    299    300    301    302    303    304    305    306    307    308
##    [309]    309    310    311    312    313    314    315    316    317    318    319    320    321    322
##    [323]    323    324    325    326    327    328    329    330    331    332    333    334    335    336
##    [337]    337    338    339    340    341    342    343    344    345    346    347    348    349    350
##    [351]    351    352    353    354    355    356    357    358    359    360    361    362    363    364
##    [365]    365    366    367    368    369    370    371    372    373    374    375    376    377    378
##    [379]    379    380    381    382    383    384    385    386    387    388    389    390    391    392
##    [393]    393    394    395    396    397    398    399    400    401    402    403    404    405    406
##    [407]    407    408    409    410    411    412    413    414    415    416    417    418    419    420
##    [421]    421    422    423    424    425    426    427    428    429    430    431    432    433    434
##    [435]    435    436    437    438    439    440    441    442    443    444    445    446    447    448
##    [449]    449    450    451    452    453    454    455    456    457    458    459    460    461    462
##    [463]    463    464    465    466    467    468    469    470    471    472    473    474    475    476
##    [477]    477    478    479    480    481    482    483    484    485    486    487    488    489    490
##    [491]    491    492    493    494    495    496    497    498    499    500    501    502    503    504
##    [505]    505    506    507    508    509    510    511    512    513    514    515    516    517    518
##    [519]    519    520    521    522    523    524    525    526    527    528    529    530    531    532
##    [533]    533    534    535    536    537    538    539    540    541    542    543    544    545    546
##    [547]    547    548    549    550    551    552    553    554    555    556    557    558    559    560
##    [561]    561    562    563    564    565    566    567    568    569    570    571    572    573    574
##    [575]    575    576    577    578    579    580    581    582    583    584    585    586    587    588
##    [589]    589    590    591    592    593    594    595    596    597    598    599    600    601    602
##    [603]    603    604    605    606    607    608    609    610    611    612    613    614    615    616
##    [617]    617    618    619    620    621    622    623    624    625    626    627    628    629    630
##    [631]    631    632    633    634    635    636    637    638    639    640    641    642    643    644
##    [645]    645    646    647    648    649    650    651    652    653    654    655    656    657    658
##    [659]    659    660    661    662    663    664    665    666    667    668    669    670    671    672
##    [673]    673    674    675    676    677    678    679    680    681    682    683    684    685    686
##    [687]    687    688    689    690    691    692    693    694    695    696    697    698    699    700
##    [701]    701    702    703    704    705    706    707    708    709    710    711    712    713    714
```

##	[715]	715	716	717	718	719	720	721	722	723	724	725	726	727	728
##	[729]	729	730	731	732	733	734	735	736	737	738	739	740	741	742
##	[743]	743	744	745	746	747	748	749	750	751	752	753	754	755	756
##	[757]	757	758	759	760	761	762	763	764	765	766	767	768	769	770
##	[771]	771	772	773	774	775	776	777	778	779	780	781	782	783	784
##	[785]	785	786	787	788	789	790	791	792	793	794	795	796	797	798
##	[799]	799	800	801	802	803	804	805	806	807	808	809	810	811	812
##	[813]	813	814	815	816	817	818	819	820	821	822	823	824	825	826
##	[827]	827	828	829	830	831	832	833	834	835	836	837	838	839	840
##	[841]	841	842	843	844	845	846	847	848	849	850	851	852	853	854
##	[855]	855	856	857	858	859	860	861	862	863	864	865	866	867	868
##	[869]	869	870	871	872	873	874	875	876	877	878	879	880	881	882
##	[883]	883	884	885	886	887	888	889	890	891	892	893	894	895	896
##	[897]	897	898	899	900	901	902	903	904	905	906	907	908	909	910
##	[911]	911	912	913	914	915	916	917	918	919	920	921	922	923	924
##	[925]	925	926	927	928	929	930	931	932	933	934	935	936	937	938
##	[939]	939	940	941	942	943	944	945	946	947	948	949	950	951	952
##	[953]	953	954	955	956	957	958	959	960	961	962	963	964	965	966
##	[967]	967	968	969	970	971	972	973	974	975	976	977	978	979	980
##	[981]	981	982	983	984	985	986	987	988	989	990	991	992	993	994
##	[995]	995	996	997	998	999	1000	1001	1002	1003	1004	1005	1006	1007	1008
##	[1009]	1009	1010	1011	1012	1013	1014	1015	1016	1017	1018	1019	1020	1021	1022
##	[1023]	1023	1024	1025	1026	1027	1028	1029	1030	1031	1032	1033	1034	1035	1036
##	[1037]	1037	1038	1039	1040	1041	1042	1043	1044	1045	1046	1047	1048	1049	1050
##	[1051]	1051	1052	1053	1054	1055	1056	1057	1058	1059	1060	1061	1062	1063	1064
##	[1065]	1065	1066	1067	1068	1069	1070	1071	1072	1073	1074	1075	1076	1077	1078
##	[1079]	1079	1080	1081	1082	1083	1084	1085	1086	1087	1088	1089	1090	1091	1092
##	[1093]	1093	1094	1095	1096	1097	1098	1099	1100	1101	1102	1103	1104	1105	1106
##	[1107]	1107	1108	1109	1110	1111	1112	1113	1114	1115	1116	1117	1118	1119	1120
##	[1121]	1121	1122	1123	1124	1125	1126	1127	1128	1129	1130	1131	1132	1133	1134
##	[1135]	1135	1136	1137	1138	1139	1140	1141	1142	1143	1144	1145	1146	1147	1148
##	[1149]	1149	1150	1151	1152	1153	1154	1155	1156	1157	1158	1159	1160	1161	1162
##	[1163]	1163	1164	1165	1166	1167	1168	1169	1170	1171	1172	1173	1174	1175	1176
##	[1177]	1177	1178	1179	1180	1181	1182	1183	1184	1185	1186	1187	1188	1189	1190
##	[1191]	1191	1192	1193	1194	1195	1196	1197	1198	1199	1200	1201	1202	1203	1204
##	[1205]	1205	1206	1207	1208	1209	1210	1211	1212	1213	1214	1215	1216	1217	1218
##	[1219]	1219	1220	1221	1222	1223	1224	1225	1226	1227	1228	1229	1230	1231	1232
##	[1233]	1233	1234	1235	1236	1237	1238	1239	1240	1241	1242	1243	1244	1245	1246
##	[1247]	1247	1248	1249	1250	1251	1252	1253	1254	1255	1256	1257	1258	1259	1260
##	[1261]	1261	1262	1263	1264	1265	1266	1267	1268	1269	1270	1271	1272	1273	1274
##	[1275]	1275	1276	1277	1278	1279	1280	1281	1282	1283	1284	1285	1286	1287	1288
##	[1289]	1289	1290	1291	1292	1293	1294	1295	1296	1297	1298	1299	1300	1301	1302
##	[1303]	1303	1304	1305	1306	1307	1308	1309	1310	1311	1312	1313	1314	1315	1316
##	[1317]	1317	1318	1319	1320	1321	1322	1323	1324	1325	1326	1327	1328	1329	1330
##	[1331]	1331	1332	1333	1334	1335	1336	1337	1338	1339	1340	1341	1342	1343	1344
##	[1345]	1345	1346	1347	1348	1349	1350	1351	1352	1353	1354	1355	1356	1357	1358
##	[1359]	1359	1360	1361	1362	1363	1364	1365	1366	1367	1368	1369	1370	1371	1372
##	[1373]	1373	1374	1375	1376	1377	1378	1379	1380	1381	1382	1383	1384	1385	1386
##	[1387]	1387	1388	1389	1390	1391	1392	1393	1394	1395	1396	1397	1398	1399	1400
##	[1401]	1401	1402	1403	1404	1405	1406	1407	1408	1409	1410	1411	1412	1413	1414
##	[1415]	1415	1416	1417	1418	1419	1420	1421	1422	1423	1424	1425	1426	1427	1428
##	[1429]	1429	1430	1431	1432	1433	1434	1435	1436	1437	1438	1439	1440	1441	1442
##	[1443]	1443	1444	1445	1446	1447	1448	1449	1450	1451	1452	1453	1454	1455	1456
##	[1457]	1457	1458	1459	1460	1461	1462	1463	1464	1465	1466	1467	1468	1469	1470
##	[1471]	1471	1472	1473	1474	1475	1476	1477	1478	1479	1480	1481	1482	1483	1484
##	[1485]	1485	1486	1487	1488	1489	1490	1491	1492	1493	1494	1495	1496	1497	1498
##	[1499]	1499	1500	1501	1502	1503	1504	1505	1506	1507	1508	1509	1510	1511	1512
##	[1513]	1513	1514	1515	1516	1517	1518	1519	1520	1521	1522	1523	1524	1525	1526


```
## [1527] 1527 1528 1529 1530 1531 1532 1533 1534 1535 1536 1537 1538 1539 1540
## [1541] 1541 1542 1543 1544 1545 1546 1547 1548 1549 1550 1551 1552 1553 1554
## [1555] 1555 1556 1557 1558 1559 1560 1561 1562 1563 1564 1565 1566 1567 1568
## [1569] 1569 1570 1571 1572 1573 1574 1575 1576 1577 1578 1579 1580 1581 1582
## [1583] 1583 1584 1585 1586 1587 1588 1589 1590 1591 1592 1593 1594 1595 1596
## [1597] 1597 1598 1599 1600 1601 1602 1603 1604 1605 1606 1607 1608 1609 1610
## [1611] 1611 1612 1613 1614 1615 1616 1617 1618 1619 1620 1621 1622 1623 1624
## [1625] 1625 1626 1627 1628 1629 1630 1631 1632 1633 1634 1635 1636 1637 1638
## [1639] 1639 1640 1641 1642 1643 1644 1645 1646 1647 1648 1649 1650 1651 1652
## [1653] 1653 1654 1655 1656 1657 1658 1659 1660 1661 1662 1663 1664 1665 1666
## [1667] 1667 1668 1669 1670 1671 1672 1673 1674 1675 1676 1677 1678 1679 1680
## [1681] 1681 1682 1683 1684 1685 1686 1687 1688 1689 1690 1691 1692 1693 1694
## [1695] 1695 1696 1697 1698 1699 1700 1701 1702 1703 1704 1705 1706 1707 1708
## [1709] 1709 1710 1711 1712 1713 1714 1715 1716 1717 1718 1719 1720 1721 1722
## [1723] 1723 1724 1725 1726 1727 1728 1729 1730 1731 1732 1733 1734 1735 1736
## [1737] 1737 1738 1739 1740 1741 1742 1743 1744 1745 1746 1747 1748 1749 1750
## [1751] 1751 1752 1753 1754 1755 1756 1757 1758 1759 1760 1761 1762 1763 1764
## [1765] 1765 1766 1767 1768 1769 1770 1771 1772 1773 1774 1775 1776 1777 1778
## [1779] 1779 1780 1781 1782 1783 1784 1785 1786 1787 1788 1789 1790 1791 1792
## [1793] 1793 1794 1795 1796 1797 1798 1799 1800 1801 1802 1803 1804 1805 1806
## [1807] 1807 1808 1809 1810 1811 1812 1813 1814 1815 1816 1817 1818 1819 1820
## [1821] 1821 1822 1823 1824 1825 1826 1827 1828 1829 1830 1831 1832 1833 1834
## [1835] 1835 1836 1837 1838 1839 1840 1841 1842 1843 1844 1845 1846 1847 1848
## [1849] 1849 1850 1851 1852 1853 1854 1855 1856 1857 1858 1859 1860 1861 1862
## [1863] 1863 1864 1865 1866 1867 1868 1869 1870 1871 1872 1873 1874 1875 1876
## [1877] 1877 1878 1879 1880 1881 1882 1883 1884 1885 1886 1887 1888 1889 1890
## [1891] 1891 1892 1893 1894 1895 1896 1897 1898 1899 1900 1901 1902 1903 1904
## [1905] 1905 1906 1907 1908 1909 1910 1911 1912 1913 1914 1915 1916 1917 1918
## [1919] 1919 1920 1921 1922 1923 1924 1925 1926 1927 1928 1929 1930 1931 1932
## [1933] 1933 1934 1935 1936 1937 1938 1939 1940 1941 1942 1943 1944 1945 1946
## [1947] 1947 1948 1949 1950 1951 1952 1953 1954 1955 1956 1957 1958 1959 1960
## [1961] 1961 1962 1963 1964 1965 1966 1967 1968 1969 1970 1971 1972 1973 1974
## [1975] 1975 1976 1977 1978 1979 1980 1981 1982 1983 1984 1985 1986 1987 1988
## [1989] 1989 1990 1991 1992 1993 1994 1995 1996 1997 1998 1999 2000 2001 2002
## [2003] 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016
## [2017] 2017 2018 2019 2020 2021 2022 2023 2024 2025 2026 2027 2028 2029 2030
## [2031] 2031 2032 2033 2034 2035 2036 2037 2038 2039 2040 2041 2042 2043 2044
## [2045] 2045 2046 2047 2048 2049 2050 2051 2052 2053 2054 2055 2056 2057 2058
## [2059] 2059 2060 2061 2062 2063 2064 2065 2066 2067 2068 2069 2070 2071 2072
## [2073] 2073 2074 2075 2076 2077 2078 2079 2080 2081 2082 2083 2084 2085 2086
## [2087] 2087 2088 2089 2090 2091 2092 2093 2094 2095 2096 2097 2098 2099 2100
## [2101] 2101 2102 2103 2104 2105 2106 2107 2108 2109 2110 2111 2112 2113 2114
## [2115] 2115 2116 2117 2118 2119 2120 2121 2122 2123 2124 2125 2126 2127 2128
## [2129] 2129 2130 2131 2132 2133 2134 2135 2136 2137 2138 2139 2140 2141 2142
## [2143] 2143 2144 2145 2146 2147 2148 2149 2150 2151 2152 2153 2154 2155 2156
## [2157] 2157 2158 2159 2160 2161 2162 2163 2164 2165 2166 2167 2168 2169 2170
## [2171] 2171 2172 2173 2174 2175 2176 2177 2178 2179 2180 2181 2182 2183 2184
## [2185] 2185 2186 2187 2188 2189 2190 2191 2192 2193 2194 2195 2196 2197 2198
## [2199] 2199 2200 2201 2202 2203 2204 2205 2206 2207 2208 2209 2210 2211 2212
## [2213] 2213 2214 2215 2216 2217 2218 2219 2220 2221 2222 2223 2224 2225 2226
## [2227] 2227 2228 2229 2230 2231 2232 2233 2234 2235 2236 2237 2238 2239 2240
## [2241] 2241 2242 2243 2244 2245 2246 2247 2248 2249 2250 2251 2252 2253 2254
## [2255] 2255 2256 2257 2258 2259 2260 2261 2262 2263 2264 2265 2266 2267 2268
## [2269] 2269 2270 2271 2272 2273 2274 2275 2276 2277 2278 2279 2280 2281 2282
## [2283] 2283 2284 2285 2286 2287 2288 2289 2290 2291 2292 2293 2294 2295 2296
## [2297] 2297 2298 2299 2300 2301 2302 2303 2304 2305 2306 2307 2308 2309 2310
## [2311] 2311 2312 2313 2314 2315 2316 2317 2318 2319 2320 2321 2322 2323 2324
## [2325] 2325 2326 2327 2328 2329 2330 2331 2332 2333 2334 2335 2336 2337 2338
```

```
## [2339] 2339 2340 2341 2342 2343 2344 2345 2346 2347 2348 2349 2350 2351 2352
## [2353] 2353 2354 2355 2356 2357 2358 2359 2360 2361 2362 2363 2364 2365 2366
## [2367] 2367 2368 2369 2370 2371 2372 2373 2374 2375 2376
##
## $class
## [1] "tbl_df"      "tbl"        "data.frame"
```

```
#convert the country code variable as a factor.
merged_malefemale_UR$`Country Code` <- as.factor(merged_malefemale_UR$`Country Code`)

#convert the year variable as a factor.
merged_malefemale_UR$Year <- as.factor(merged_malefemale_UR$Year)

#convert the male unemployment rate variable as a numeric.
merged_malefemale_UR$`Male unemployment rate` <- as.numeric(merged_malefemale_UR$`Male unempl
oyment rate`)
```

```
## Warning: NAs introduced by coercion
```

```
#convert the female unemployment rate variable as a numeric.
merged_malefemale_UR$`Female unemployment rate` <- as.numeric(merged_malefemale_UR$`Female un
employment rate`)
```

```
## Warning: NAs introduced by coercion
```

```
#Labeling factor year
merged_malefemale_UR$Year <- factor(merged_malefemale_UR$Year, levels = c("2011 [YR2011]", "2
012 [YR2012]", "2013 [YR2013]", "2014 [YR2014]", "2015 [YR2015]", "2016 [YR2016]", "2017 [YR2
017]", "2018 [YR2018]", "2019 [YR2019]"),
                                labels = c("2011", "2012", "2013", "2014", "20
15", "2016", "2017", "2018", "2019"))
```

```
#rechecking the data structures after conversion.
str(merged_malefemale_UR)
```

```
## tibble [2,376 x 5] (S3: tbl_df/tbl/data.frame)
## $ Country Name      : chr [1:2376] "Afghanistan" "Afghanistan" "Afghanistan" "Afgha
nistan" ...
## $ Country Code      : Factor w/ 264 levels "ABW","AFG","AGO",...: 2 2 2 2 2 2 2 2 2
4 ...
## $ Year              : Factor w/ 9 levels "2011","2012",...: 1 2 3 4 5 6 7 8 9 1 ...
## $ Male unemployment rate : num [1:2376] 10.9 10.9 10.9 10.8 10.7 ...
## $ Female unemployment rate: num [1:2376] 14.8 14.9 14.7 14.5 14.4 ...
```

Step explanation:

- Firstly, I inspect the data using `str()`, `class()`, and `attributes()` functions.
- Secondly, I converted the data into factors and numerics. I also labelled the Year variable.
- Lastly, I recheck the data structures using `str()` function. As the result now the data

Tidy & Manipulate Data II

```
#Mutate new variable named Unemployment rate gap between genders.
```

```
merged_mutated_MF_UR <- mutate(merged_malefemale_UR, `Unemployment rate gap between genders`  
  = `Male unemployment rate` - `Female unemployment rate`)
```

```
#Display the new variable and several first observations of the merged_mutated_MF_UR dataset.  
head(merged_mutated_MF_UR)
```

Country Name <chr>	Country Code <fct>	Y... <fct>	Male unemployment rate <dbl>	Female unemployment rate <dbl>
Afghanistan	AFG	2011	10.90	14.79
Afghanistan	AFG	2012	10.88	14.86
Afghanistan	AFG	2013	10.89	14.70
Afghanistan	AFG	2014	10.78	14.53
Afghanistan	AFG	2015	10.68	14.45
Afghanistan	AFG	2016	10.57	14.33

6 rows | 1-5 of 6 columns

Step explanation:

- In this step, I use mutate() function to create a new variable named “Unemployment rate gap between genders” which represent the difference between males unemployment rate and females unemployment rate.
- Then, I used head() function to display the new mutated variable and their first 6 observations.

Scan I

```
#scan for missing values  
colSums(is.na(merged_mutated_MF_UR))
```

```
##                Country Name                Country Code  
##                0                0  
##                Year                Male unemployment rate  
##                0                279  
##                Female unemployment rate Unemployment rate gap between genders  
##                279                279
```

```
#scan for special values (nan)  
sapply(merged_mutated_MF_UR, function(x) sum( is.nan(x) ))
```

```
##          Country Name          Country Code
##          0                      0
##          Year          Male unemployment rate
##          0                      0
##          Female unemployment rate Unemployment rate gap between genders
##          0                      0
```

```
#scan for special values (infinite)
sapply(merged_mutated_MF_UR, function(x) sum( is.infinite(x) ))
```

```
##          Country Name          Country Code
##          0                      0
##          Year          Male unemployment rate
##          0                      0
##          Female unemployment rate Unemployment rate gap between genders
##          0                      0
```

```
#scan for obvious error (inconsistencies)
merged_mutated_MF_UR$`Unemployment rate gap between genders` <= 100 & merged_mutated_MF_UR$`U
nemployment rate gap between genders` >= -100
```

[illegible]

```
## [799] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [813] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [827] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [841] TRUE TRUE TRUE TRUE TRUE TRUE NA NA NA NA NA NA NA NA
## [855] NA TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [869] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [883] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [897] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [911] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE NA NA NA NA NA NA
## [925] NA NA NA TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [939] TRUE TRUE TRUE TRUE TRUE TRUE TRUE NA NA NA NA NA NA NA
## [953] NA NA TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [967] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [981] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [995] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1009] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1023] TRUE TRUE TRUE TRUE NA NA NA NA NA NA NA NA NA TRUE
## [1037] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1051] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1065] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1079] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1093] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1107] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE NA NA NA NA
## [1121] NA NA NA NA NA TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1135] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1149] TRUE TRUE TRUE TRUE NA NA NA NA NA NA NA NA NA TRUE
## [1163] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE NA NA NA NA NA NA
## [1177] NA NA NA TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1191] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1205] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1219] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1233] TRUE NA NA NA NA NA NA NA NA NA NA TRUE TRUE TRUE TRUE
## [1247] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1261] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1275] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1289] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1303] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE NA NA
## [1317] NA NA NA NA NA NA NA TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1331] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1345] TRUE TRUE TRUE TRUE TRUE TRUE NA NA NA NA NA NA NA NA
## [1359] NA TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1373] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1387] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1401] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1415] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1429] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1443] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1457] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1471] TRUE TRUE TRUE TRUE TRUE TRUE NA NA NA NA NA NA NA NA
## [1485] NA TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1499] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1513] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE NA NA NA NA NA
## [1527] NA NA NA NA TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1541] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE NA NA NA NA NA NA
## [1555] NA NA NA TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1569] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1583] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
## [1597] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
```

[illegible]

```
#check for complete values.
```

```
merged_mutated_MF_UR[complete.cases(merged_mutated_MF_UR), ]
```

Country Name <chr>	Country Code <fct>	Y... <fct>	Male unemployment rate <dbl>	Female unemployment rate <dbl>
Afghanistan	AFG	2011	10.90	14.79
Afghanistan	AFG	2012	10.88	14.86
Afghanistan	AFG	2013	10.89	14.70
Afghanistan	AFG	2014	10.78	14.53
Afghanistan	AFG	2015	10.68	14.45
Afghanistan	AFG	2016	10.57	14.33
Afghanistan	AFG	2017	10.42	14.09
Afghanistan	AFG	2018	10.29	13.92
Afghanistan	AFG	2019	10.19	13.81
Albania	ALB	2011	13.27	13.75

1-10 of 2,097 rows | 1-5 of 6 columns

Previous 1 2 3 4 5 6 ... 210 Next

```
#Omit the missing data.
```

```
scan_MF_UR <- na.omit(merged_mutated_MF_UR)
```

```
#Rechecking the missing data after using na.omit() function.
```

```
sapply(scan_MF_UR, function(x) sum( is.infinite(x) | is.nan(x) | is.na(x) ))
```

```
##           Country Name           Country Code
##                0                0
##           Year           Male unemployment rate
##                0                0
##           Female unemployment rate Unemployment rate gap between genders
##                0                0
```

Step explanation:

- Firstly, I checked the missing values, special values, and obvious errors (inconsistencies).
- Secondly, I checked the complete observations of the data.
- Lastly, I omitted the missing values using `na.omit()` function and assign it to new variable named `scan_MF_UR`.

Methodology:

For each country in this dataset, the country have either full/complete data (2011 - 2019) or missing all the data (NA for all years 2011 - 2019). Thus, replacing / imputing missing data with data from other countries would not be appropriate because it will bias the distribution of the dataset I am using. Therefore, I chose to omit the missing data (279 observations omitted) and ended up with 2097 complete observations of year 2011 - 2019 for every country.

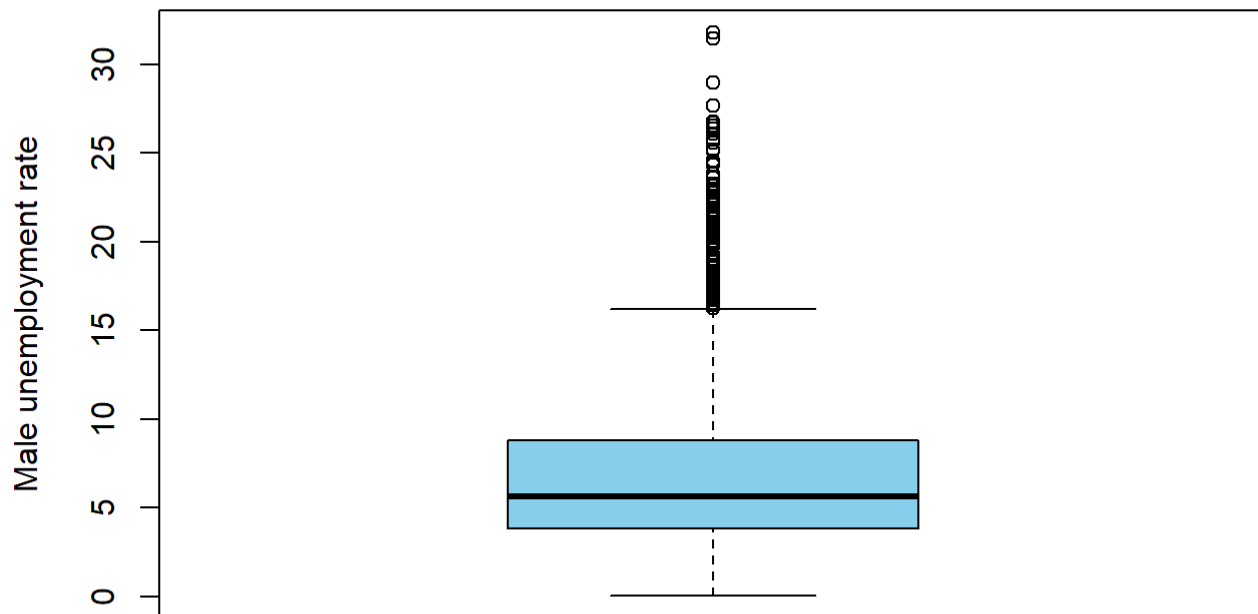
Scan II

```
#univariate outliers
```

```
#boxplot
```

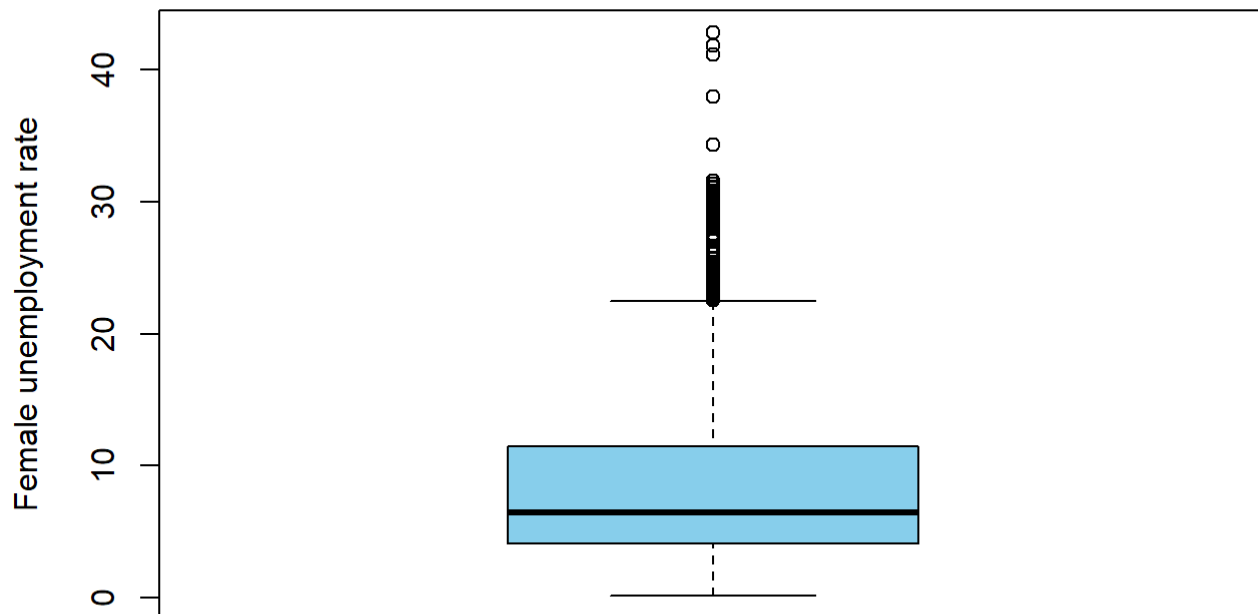
```
scan_MF_UR$`Male unemployment rate` %>% boxplot(main = "Boxplot of Male unemployment rate", ylab = "Male unemployment rate", col = "skyblue")
```

Boxplot of Male unemployment rate



```
scan_MF_UR$`Female unemployment rate` %>% boxplot(main = "Boxplot of Female unemployment rate", ylab = "Female unemployment rate", col = "skyblue")
```

Boxplot of Female unemployment rate

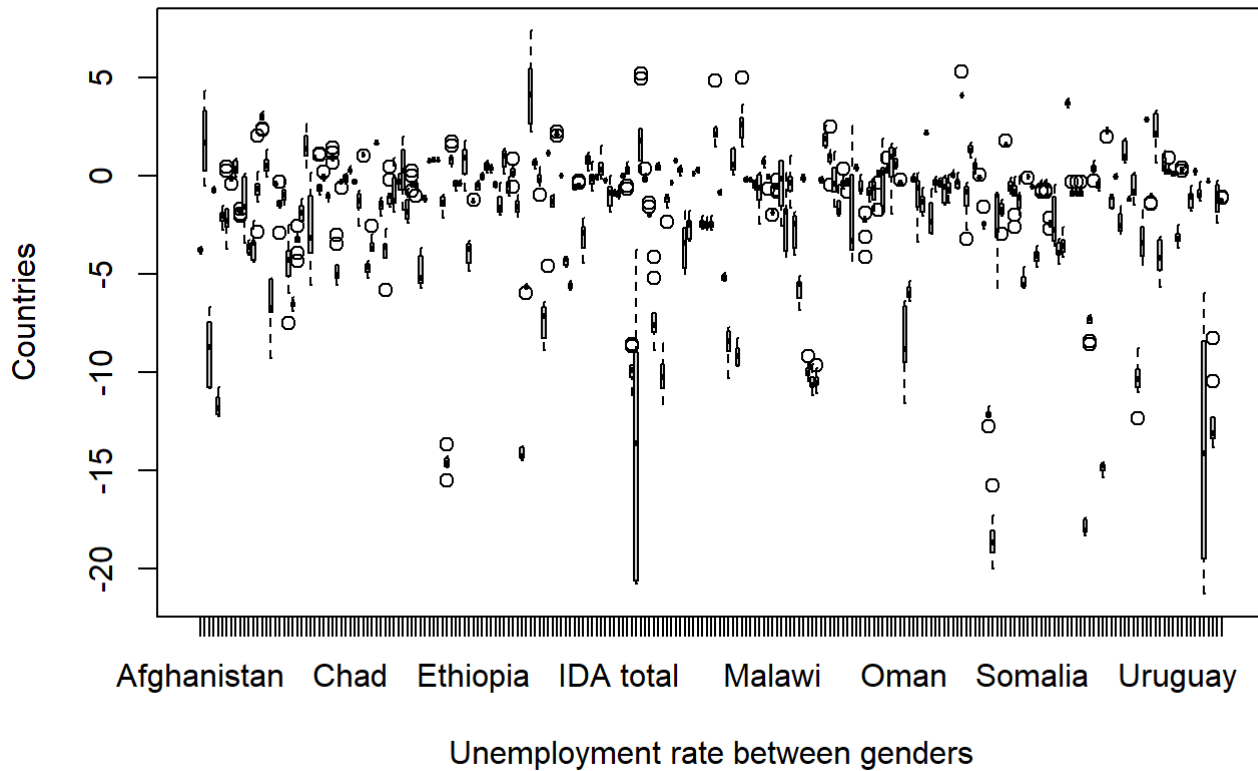


```
#multivariate outliers
```

```
#boxplot
```

```
boxplot(scan_MF_UR$`Unemployment rate gap between genders` ~ scan_MF_UR$`Country Name`, main="Country Name by unemployment rate gap between genders", ylab = "Countries", xlab = "Unemployment rate between genders")
```

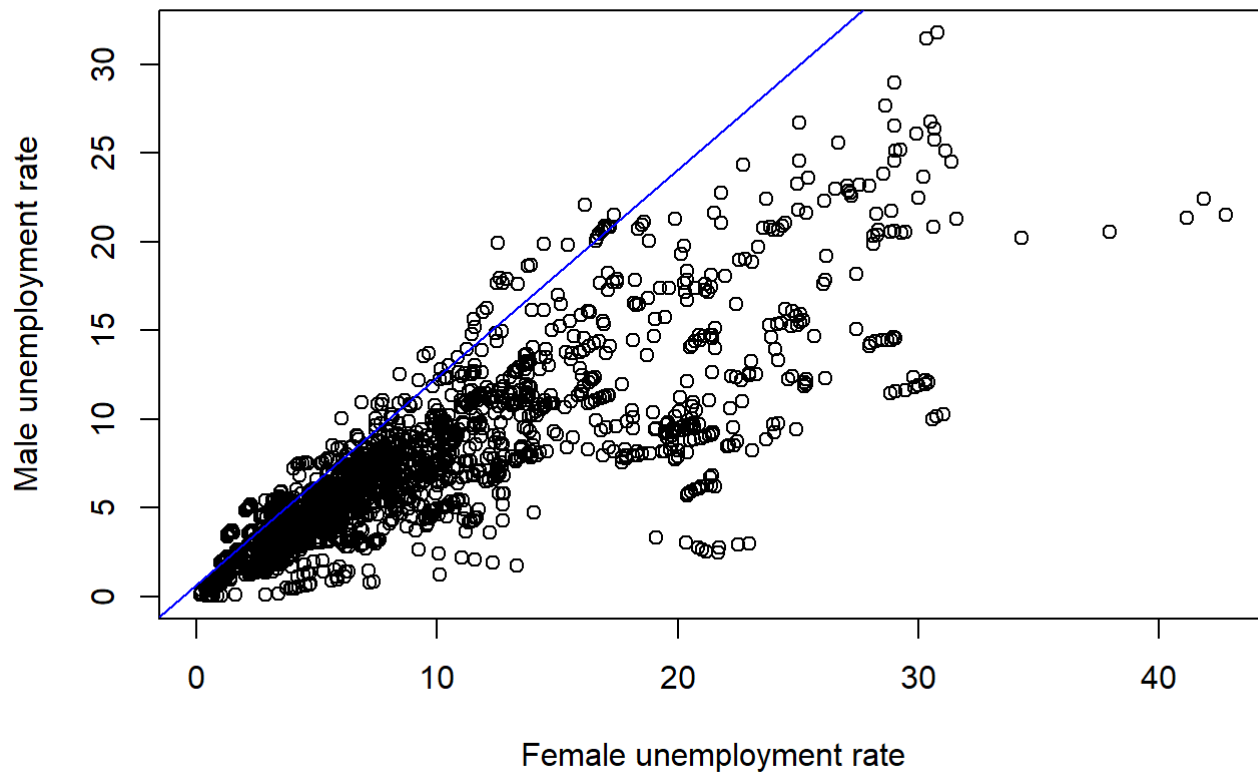
Country Name by unemployment rate gap between genders



```
#scatter plot
```

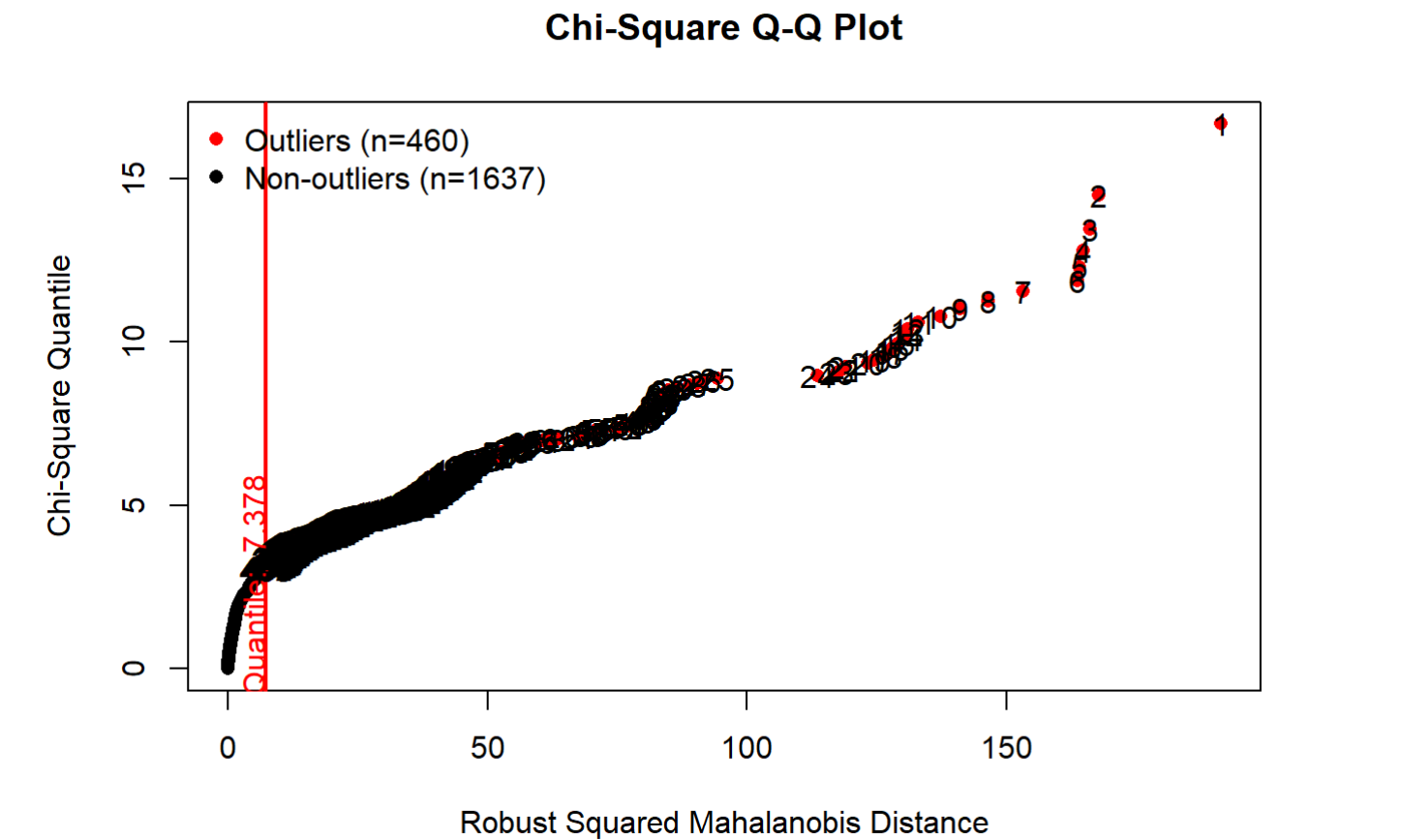
```
scan_MF_UR %>% plot(`Male unemployment rate` ~ `Female unemployment rate`, data = ., ylab= "M
ale unemployment rate", xlab = "Female unemployment rate", main = "Male and Female unemployem
ent rate distribution")
abline(lm(`Female unemployment rate` ~ `Male unemployment rate`, data = scan_MF_UR), col = "b
lue")
```

Male and Female unemployment rate distribution



```
#mahalanobis distance with QQ plots
mahalanobis <- dplyr::select(scan_MF_UR, `Male unemployment rate`, `Female unemployment rate`
`)

results <- mvn(data = mahalanobis, multivariateOutlierMethod = "quan", showOutliers = TRUE)
```



```
results$multivariateOutliers
```

Observation		Mahalanobis Distance	Outlier
<chr>		<dbl>	<chr>
1	1	191.236	TRUE
2	2	167.699	TRUE
3	3	166.127	TRUE
4	4	164.677	TRUE
5	5	163.943	TRUE
6	6	163.532	TRUE
7	7	153.155	TRUE
8	8	146.406	TRUE
9	9	141.011	TRUE
10	10	137.190	TRUE

1-10 of 460 rows

Previous123456...46Next

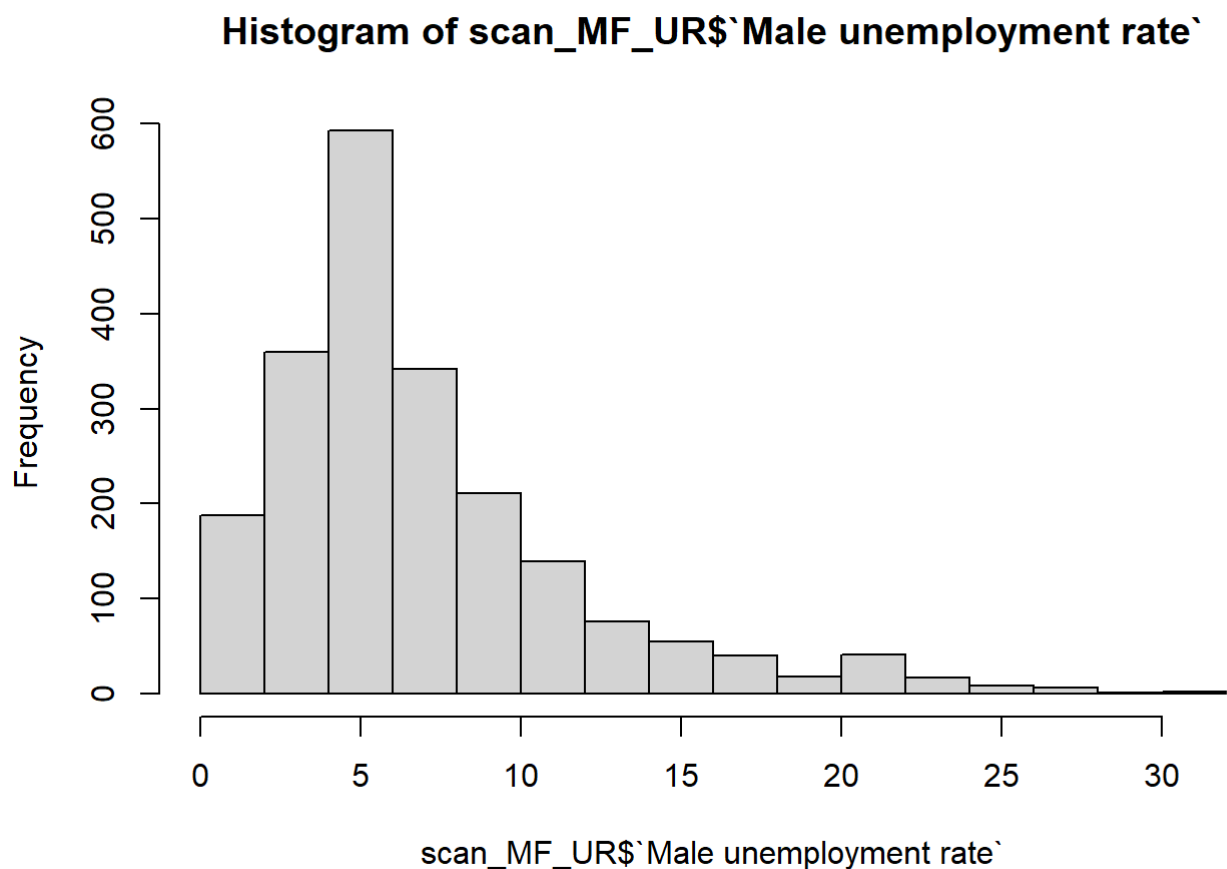
- Step explanation:
- Firstly, I check the univariate outliers using the boxplot().
 - Secondly, I check the multivariate outliers using boxplot, scatterplot, and mahalanobis distance.

Methodology:

Although the univariate boxplot and QQ plots illustrates that the data I use have many outliers. However, the multivariate boxplot and scatter plot show that the unemployment rates from one country to another can vary widely. For example, Australia has an unemployment rate of 5% compared to a country that is still in civil war like Libya, which has an unemployment rate as high as 25%. As it is known that the real world is filled with outliers, so to model the real world I believe that keeping the outliers will be beneficial for my analysis in this particular datasets.

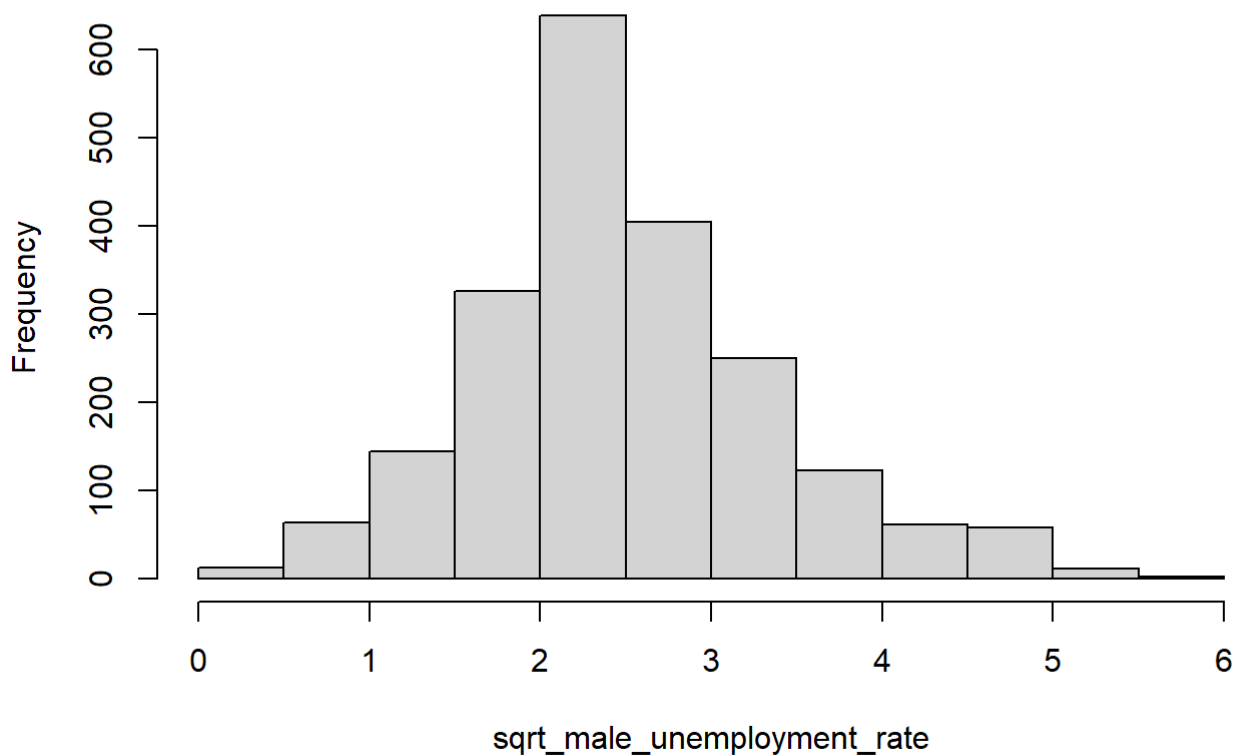
Transform

```
#Male unemployment rate before transformation  
hist(scan_MF_UR$`Male unemployment rate`)
```



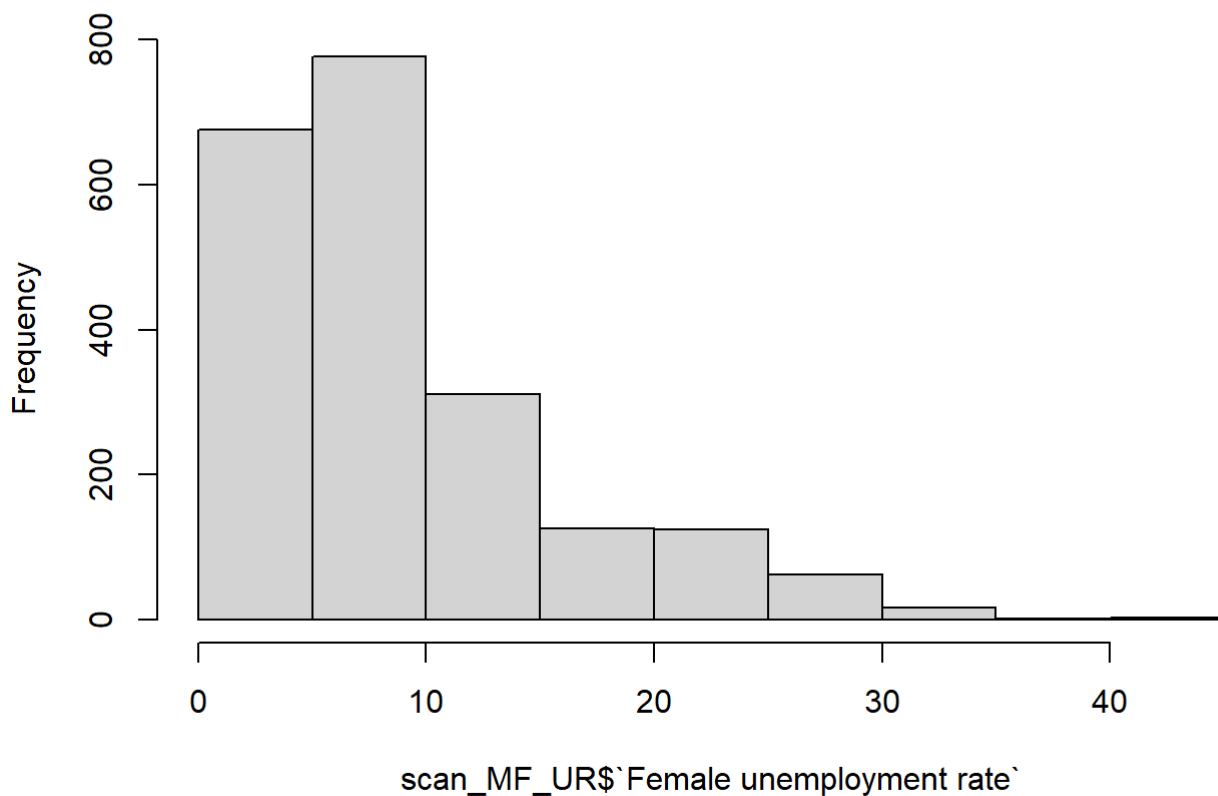
```
#Male unemployment rate after transformation  
sqrt_male_unemployment_rate <- sqrt(scan_MF_UR$`Male unemployment rate`)  
hist(sqrt_male_unemployment_rate)
```

Histogram of sqrt_male_unemployment_rate



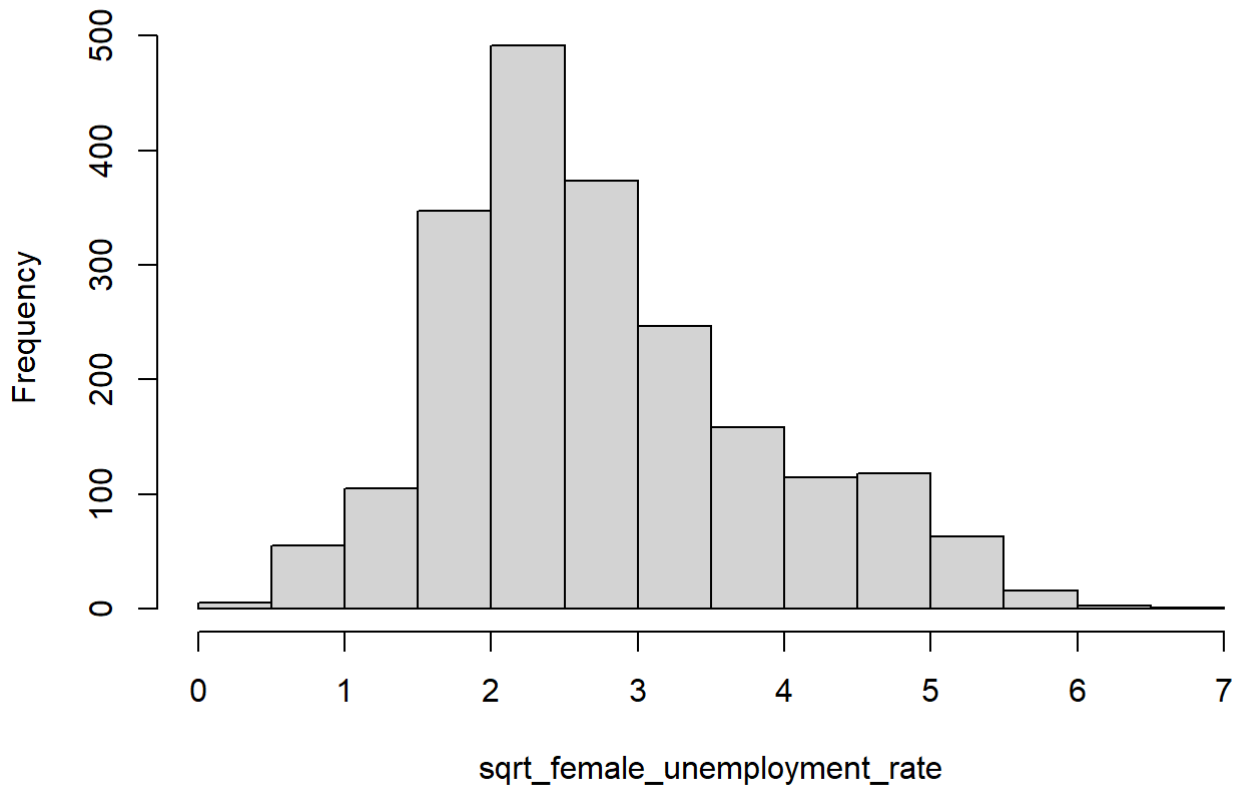
```
#Female unemployment rate before transformation.  
hist(scan_MF_UR$`Female unemployment rate`)
```

Histogram of scan_MF_UR\$`Female unemployment rate`



```
#Female unemployment rate after transformation.  
sqrt_female_unemployment_rate <- sqrt(scan_MF_UR$`Female unemployment rate`)  
hist(sqrt_female_unemployment_rate)
```

Histogram of sqrt_female_unemployment_rate



Purpose of transformation: The purpose of this transformation is to decrease the skewness of the dataset and transform the distribution into a normal distribution.

Step explanation:

- Firstly, I applied transformation to the male unemployment rate dataset using the sqrt() function.
- Secondly, I applied transformation to the female unemployment rate dataset using sqrt() function.

Decision: After trying five different data transformation (log10, log, squared, sqrt, and BoxCox), I found that square roots are the transformation that is closest to the normal distribution for variable male unemployment rate and female unemployment rate. Therefore, sqrt() function will be the one chosen/used for this transformation.

note: I couldn't put the other 4 histogram (log10, log, squared, and BoxCox) due to page limit (25 pages).