

# ACP T-SNE UMAP

## 1. Importer les données:

```
data <- read.csv('C:/Users/asus/OneDrive/Desktop/Project/data/breast_cancer.csv',
  sep=';', row.names='10_sample')

## [1] 1016 51

head(data)

##          NAT1      BIRC5      BAGO1      BCL2      BLVRA      CCNE1      CCNE1
## TCGA-3C-AAU-01A 7.100449 3.361004 3.972501 4.145609 4.765233 4.780987 2.164814
## TCGA-3C-AAU-01A 4.453040 4.501040 2.730929 1.403020 5.822480 5.291005 2.535437
## TCGA-3C-AALJ-01A 4.455574 3.164543 3.911511 4.191557 5.987255 5.229446 2.267963
## TCGA-3C-AALJ-01A 4.207961 3.920234 3.608335 3.894904 5.211594 4.614641 0.951107
## TCGA-4H-AAAK-01A 1.695378 2.950846 4.110014 3.572843 4.317859 3.772768 1.193958
## TCGA-5L-AAT0-01A 3.839200 2.630192 3.649955 3.958096 5.186200 3.130406 1.138433
##          CDC5      CDC20      CDH3      CENPF      EGFR      ERBB2      ESR1
## TCGA-3C-AAU-01A 2.633598 4.131205 1.133455 3.182165 0.680991 5.239199 3.877698
## TCGA-3C-AALJ-01A 2.734157 4.176553 0.110023 3.392215 0.954435 0.927166 0.382693
## TCGA-3C-AALJ-01A 3.379951 4.502752 0.206156 2.984169 0.821807 5.804556 4.858459
## TCGA-3C-AALK-01A 1.472956 3.896552 0.002392 2.893521 2.044761 5.958474 3.284898
## TCGA-4H-AAAK-01A 2.338953 3.473484 0.008773 2.514709 1.527022 5.942930 4.450677
## TCGA-5L-AAT0-01A 1.494393 2.720610 0.061163 2.014191 1.584887 5.579183 4.690124
##          FGFR4      FOXD1      GRS7      FOXA1      KR71      KR714      KR717
## TCGA-3C-AAU-01A 0.923478 0.342682 3.845602 1.158444 0.034380 0.080690 0.009777
## TCGA-3C-AALI-01A 4.653488 0.794085 0.780099 0.383003 2.530948 5.239199 3.755362
## TCGA-3C-AALJ-01A 0.605796 1.628313 4.487610 3.965421 0.090767 0.109529 0.310840
## TCGA-3C-AALK-01A 1.634376 2.414163 0.206156 2.984169 0.821807 5.804556 4.858459
## TCGA-4H-AAAK-01A 3.292628 1.559426 4.070559 6.456412 5.081345 5.958474 3.937990
## TCGA-5L-AAT0-01A 0.853090 1.740289 3.603108 0.353737 5.114444 0.002268 5.181969
##          MAPT      MDM2      MKI67      MMP11      MYBL2      MYC      PGR
## TCGA-3C-AAU-01A 4.571417 5.000554 3.406616 0.383997 3.991620 5.324133 3.152494
## TCGA-3C-AALJ-01A 0.767409 2.169065 3.293928 6.031836 6.328959 1.997876 0.127529
## TCGA-3C-AALJ-01A 3.530899 3.002628 2.225258 6.462122 5.776590 4.428182 1.399339
## TCGA-3C-AALK-01A 3.180234 2.629048 2.305445 5.776590 3.277669 4.540685 2.991844
## TCGA-4H-AAAK-01A 1.462092 2.608220 2.908095 1.470010 0.300921 5.853266 6.691324
## TCGA-5L-AAT0-01A 4.596884 2.475144 1.568097 7.717722 1.361077 4.059862 1.792283
##          RRM2      SFRP1      TYMS      MTA      EXO1      PTG11      MELK
## TCGA-3C-AAU-01A 4.069746 2.019955 3.102706 0.056151 1.851411 3.551624 2.541704
## TCGA-3C-AALI-01A 5.110049 1.025220 3.606615 0.230768 2.893140 3.976765 3.232838
## TCGA-3C-AALJ-01A 3.858266 0.767185 3.567142 0.847283 1.735396 4.375152 3.361018
## TCGA-3C-AALK-01A 3.293315 3.701924 3.479655 1.252874 1.311587 3.117828 2.326996
## TCGA-4H-AAAK-01A 2.302735 4.099581 3.474526 1.224086 1.119397 2.956249 2.377493
## TCGA-5L-AAT0-01A 2.092475 0.611640 0.611640 0.877212 1.120010 0.877212 1.536737
##          NDC80      KIF2C      UBE2C      ORC6      SLC39A6      PGDH
## TCGA-3C-AAU-01A 2.687899 3.422050 4.678216 1.837654 10.319962 2.991554
## TCGA-3C-AALI-01A 1.481391 3.674568 5.883007 2.762908 4.579932 2.973480
## TCGA-3C-AALJ-01A 3.615086 3.505261 5.706425 2.560333 7.747377 1.003328
## TCGA-3C-AALK-01A 1.788017 2.322912 4.378114 0.811602 7.921403 3.781136
## TCGA-4H-AAAK-01A 1.253309 1.582389 3.057444 0.621261 5.251726 2.413841
## TCGA-5L-AAT0-01A 0.916509 0.00227 0.00227 0.00227 0.00227 0.00227
##          NUF2      TMEM45B      paw50
## TCGA-3C-AAU-01A 4.571417 5.000554 3.406616 0.383997 3.991620 5.324133 3.152494
## TCGA-3C-AALI-01A 3.124620 3.946538 0.945338 0.945338 0.945338 0.945338 0.945338
## TCGA-3C-AALJ-01A 3.053335 0.281303 0.281303 0.281303 0.281303 0.281303 0.281303
## TCGA-3C-AALK-01A 1.747959 3.289543 3.289543 3.289543 3.289543 3.289543 3.289543
## TCGA-4H-AAAK-01A 1.537125 2.769603 2.769603 2.769603 2.769603 2.769603 2.769603
## TCGA-5L-AAT0-01A 0.947315 2.884545 2.884545 2.884545 2.884545 2.884545 2.884545
```

## Afficher les occurrences de chaque niveau de la colonne pam50:

```
group_sizes <- table(data$pam50)
print(group_sizes)

##          basal-like  HER2-enriched  luminal-A  luminal-B
##             199              82             543             201

library(ggplot2)
ggplot(data, aes(x = pam50)) +
  geom_bar() +
  labs(title = "Nombre d'occurrences par groupe",
    x = "pam50",
    y = "Nombre d'occurrences")

# Nombre d'occurrences par groupe
# 400
# 200
# 0
# basal-like  HER2-enriched  pam50  luminal-A  luminal-B
```

## 2.Séparer les données d'expression et les étiquettes

```
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

#Données d'expression de 50 gènes
X <- select_if(data, is.numeric)
cat('X', dim(X), '\n')

## X 1016 50

#Étiquettes correspondantes (sous-types moléculaires)
y <- data$pam50
cat('y', length(y), '\n')

## y 1016
```

## 3.Afficher les valeurs d'expression:

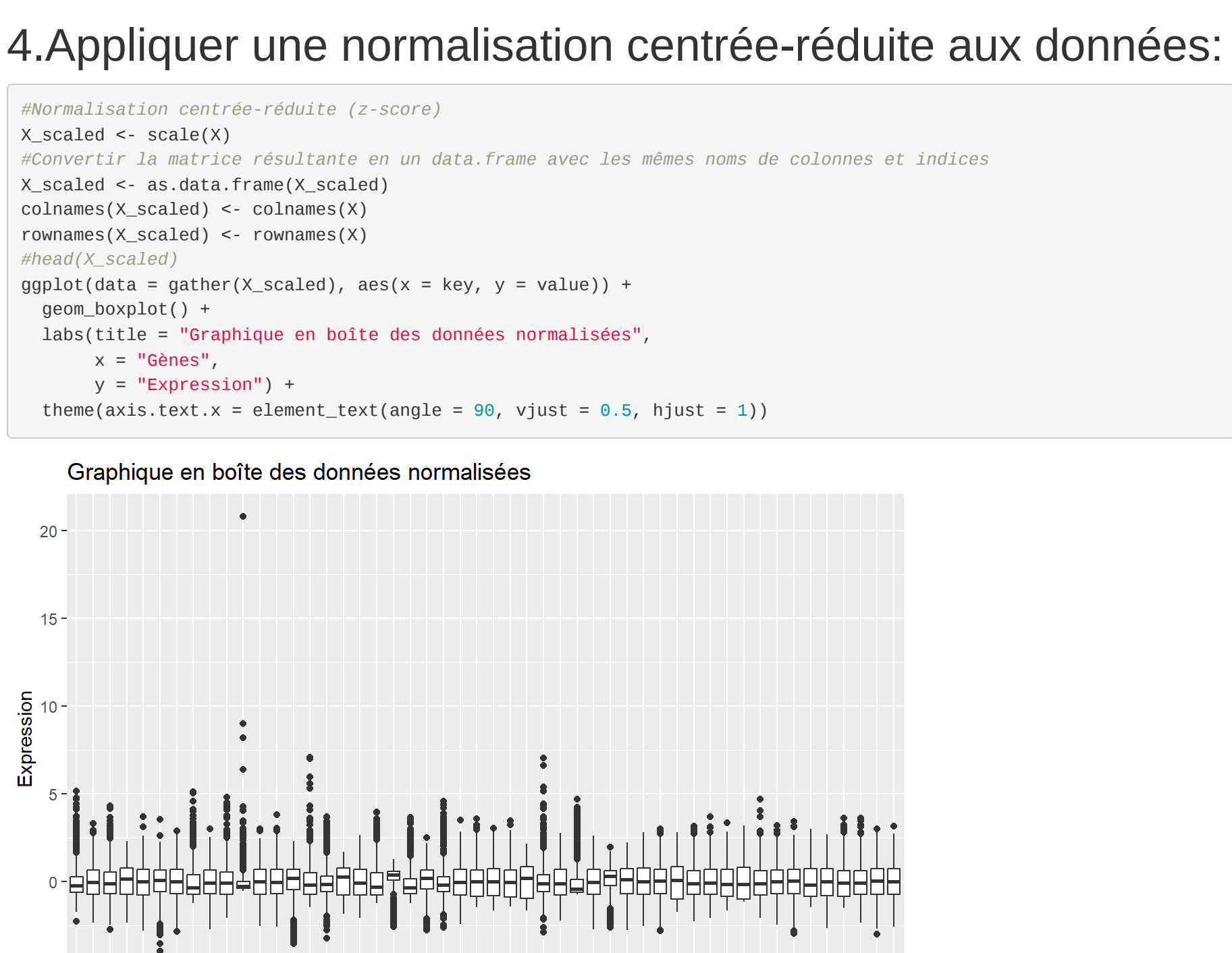
```
library(reshape2)
library(tidyrr)

##
##
## Attaching package: 'tidyr'

## The following object is masked from 'package:reshape2':
##
##   smiths

#Trier les colonnes par moyenne
sort_by_mean <- colMeans(X)
sort_by_mean <- sort_by_mean[order(sort_by_mean)]
X_sorted <- X[, names(sort_by_mean)]
#Créez le graphique en boîte avec ggplot2
ggplot(melt(X_sorted), aes(x = variable, y = value)) +
  geom_boxplot() +
  labs(title = "Graphique en boîte trié par moyenne",
    x = "Gènes",
    y = "Expression") +
  theme(axis.text.x = element_text(angle = 90, vjust = 0.5, hjust=1))

## No id variables; using all as measure variables
```



## 4.Appliquer une normalisation centrée-réduite aux données:

```
#Normalisation centrée-réduite (z-score)
X_scaled <- scale(X)
#Convertir la matrice résultante en un data.frame avec les mêmes noms de colonnes et indices
X_scaled <- as.data.frame(X_scaled)
colnames(X_scaled) <- colnames(X)
rownames(X_scaled) <- rownames(X)
#Normaliser les données
ggplot(data = gather(X_scaled, aes(x = key, y = value)) +
  geom_boxplot() +
  labs(title = "Graphique en boîte des données normalisées",
    x = "Gènes",
    y = "Expression") +
  theme(axis.text.x = element_text(angle = 90, vjust = 0.5, hjust = 1))
```



## 5. Faire une analyse en composantes principales (ACP)

### 5.1 Calcul de l'ACP:

```
#Installer et chargez les bibliothèques FactoMineR et factoextra si ce n'est pas déjà fait
#install.packages("FactoMineR")
#install.packages("factoextra")
library(FactoMineR)
library(factoextra)

## Welcome! Want to learn more? See two factoextra-related books at https://goo.gl/ve3WbA

#Réaliser l'ACP sur les données X_scaled
pca_result <- PCA(X_scaled, graph = FALSE)
#Obtenez les composantes principales (PC)
X_pca <- as.data.frame(pca_result$ind$coord)
#Renommez les colonnes des PC
colnames(X_pca) <- paste0("PC", 1:ncol(X_pca))

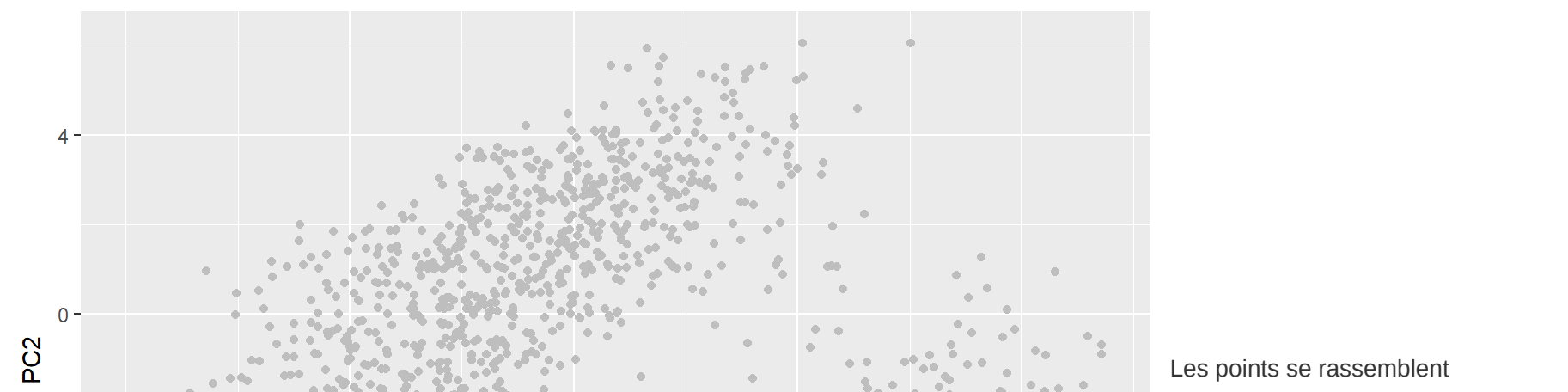
#Affichez les premières lignes du data.frame avec les composantes principales
head(X_pca)

##          PC1      PC2      PC3      PC4      PC5
## TCGA-3C-AAU-01A -0.8262713 3.4336500 2.5106320 -0.7383054 1.0699080
## TCGA-3C-AALJ-01A 2.0928393 3.9221698 -6.6983824 1.5316700 0.13862634
## TCGA-3C-AALJ-01A 0.8692634 3.7491997 0.4344589 -0.6897723 -1.0284517
## TCGA-3C-AALK-01A -1.7659214 -0.7684501 -1.3398320 2.3319832 -0.67452753
## TCGA-4H-AAAK-01A -2.0171204 -1.4094408 -1.1634389 1.3857245 -0.1715645
## TCGA-5L-AAT0-01A -4.5929983 -2.7013014 -1.1158453 0.5633696 -0.79099750
```

### 5.2 Calcul de la variance expliquée:

```
pca <- prcomp(X_scaled)
explained_variance <- tibble{
  PC = seq_along(pca$sdev),
  Explained_Variance_Ratio = pca$sdev^2 / sum(pca$sdev^2) * 100
}
# Afficher la variance expliquée
print(explained_variance)

## # A tibble: 50 x 2
##       PC Explained_Variance_Ratio
##   <int>          <dbl>
## 1     1             41.4
## 2     2             14.0
## 3     3              6.15
## 4     4              6.15
## 5     5              4.30
## 6     6              2.95
## 7     7              2.14
## 8     8              1.93
## 9     9              1.93
## 10    10             1.65
## # 40 more rows
```



## 5.3 Visualisation des deux premières composantes principales de l'ACP:

```
library(ggplot2)
#Créez un scatter plot de PC1 vs PC2
ggplot(X_pca, aes(x = PC1, y = PC2)) +
  geom_point(color = 'gray') +
  labs(title = "Scatter Plot des PC1 et PC2", x = "PC1", y = "PC2")

# Scatter Plot des PC1 et PC2
# 4
# 0
# -4
# -8
# -10
# -12
# -14
# -16
# -18
# -20
# -22
# -24
# -26
# -28
# -30
# -32
# -34
# -36
# -38
# -40
# -42
# -44
# -46
# -48
# -50
# -52
# -54
# -56
# -58
# -60
# -62
# -64
# -66
# -68
# -70
# -72
# -74
# -76
# -78
# -80
# -82
# -84
# -86
# -88
# -90
# -92
# -94
# -96
# -98
# -100
# -102
# -104
# -106
# -108
# -110
# -112
# -114
# -116
# -118
# -120
# -122
# -124
# -126
# -128
# -130
# -132
# -134
# -136
# -138
# -140
# -142
# -144
# -146
# -148
# -150
# -152
# -154
# -156
# -158
# -160
# -162
# -164
# -166
# -168
# -170
# -172
# -174
# -176
# -178
# -180
# -182
# -184
# -186
# -188
# -190
# -192
# -194
# -196
# -198
# -200
# -202
# -204
# -206
# -208
# -210
# -212
# -214
# -216
# -218
# -220
# -222
# -224
# -226
# -228
# -230
# -232
# -234
# -236
# -238
# -240
# -242
# -244
# -246
# -248
# -250
# -252
# -254
# -256
# -258
# -260
# -262
# -264
# -266
# -268
# -270
# -272
# -274
# -276
# -278
# -280
# -282
# -284
# -286
# -288
# -290
# -292
# -294
# -296
# -298
# -300
# -302
# -304
# -306
# -308
# -310
# -312
# -314
# -316
# -318
# -320
# -322
# -324
# -326
# -328
# -330
# -332
# -334
# -336
# -338
# -340
# -342
# -344
# -346
# -348
# -350
# -352
# -354
# -356
# -358
# -360
# -362
# -364
# -366
# -368
# -370
# -372
# -374
# -376
# -378
# -380
# -382
# -384
# -386
# -388
# -390
# -392
# -394
# -396
# -398
# -400
# -402
# -404
# -406
# -408
# -410
# -412
# -414
# -416
# -418
# -420
# -422
# -424
# -426
# -428
# -430
# -432
# -434
# -436
# -438
# -440
# -442
# -444
# -446
# -448
# -450
# -452
# -454
# -456
# -458
# -460
# -462
# -464
# -466
# -468
# -470
# -472
# -474
# -476
# -478
# -480
# -482
# -484
# -486
# -488
# -490
# -492
# -494
# -496
# -498
# -500
# -502
# -504
# -506
# -508
# -510
# -512
# -514
# -516
# -518
# -520
# -522
# -524
# -526
# -528
# -530
# -532
# -534
# -536
# -538
# -540
# -542
# -544
# -546
# -548
# -550
# -552
# -554
# -556
# -558
# -560
# -562
# -564
# -566
# -568
# -570
# -572
# -574
# -576
# -578
# -580
# -582
# -584
# -586
# -588
# -590
# -592
# -594
# -596
# -598
# -600
# -602
# -604
# -606
# -608
# -610
# -612
# -614
# -616
# -618
# -620
# -622
# -624
# -626
# -628
# -630
# -632
# -634
# -636
# -638
# -640
# -642
# -644
# -646
# -648
# -650
# -652
# -654
# -656
# -658
# -660
# -662
# -664
# -666
# -668
# -670
# -672
# -674
# -676
# -678
# -680
# -682
# -684
# -686
# -688
# -690
# -692
# -694
# -696
# -698
# -700
# -702
# -704
# -706
# -708
# -710
# -712
# -714
# -716
# -718
# -720
# -722
# -724
# -726
# -728
# -730
# -732
# -734
# -736
# -738
# -740
# -742
# -744
# -746
# -748
# -750
# -752
# -754
# -756
# -758
# -760
# -762
# -764
# -766
# -768
# -770
# -772
# -774
# -776
# -778
# -780
# -782
# -784
# -786
# -788
# -790
# -792
# -794
# -796
# -798
# -800
# -802
# -804
# -806
# -808
# -810
# -812
# -814
# -816
# -818
# -820
# -822
# -824
# -826
# -828
# -830
# -832
# -834
# -836
# -838
# -840
# -842
# -844
# -846
# -848
# -850
# -852
# -854
# -856
# -858
# -860
# -862
# -864
# -866
# -868
# -870
# -872
# -874
# -876
# -878
# -880
# -882
# -884
# -886
# -888
# -890
# -892
# -894
# -896
# -898
# -900
# -902
# -904
# -906
# -908
# -910
# -912
# -914
# -916
# -918
# -920
# -922
# -924
# -926
# -928
# -930
# -932
# -934
# -936
# -938
# -940
# -942
# -944
# -946
# -948
# -950
# -952
# -954
# -956
# -958
# -960
# -962
# -964
# -966
# -968
# -970
# -972
# -974
# -976
# -978
# -980
# -982
# -984
# -986
# -988
# -990
# -992
# -994
# -996
# -998
# -1000
# -1002
# -1004
# -1006
# -1008
# -1010
# -1012
# -1014
# -1016
# -1018
# -1020
# -1022
# -1024
# -1026
# -1028
# -1030
# -1032
# -1034
# -1036
# -1038
# -1040
# -1042
# -1044
# -1046
# -1048
# -1050
# -1052
# -1054
# -1056
# -1058
# -1060
# -1062
# -1064
# -1066
# -1068
# -1070
# -1072
# -1074
# -1076
# -1078
# -1080
# -1082
# -1084
# -1086
# -1088
# -1090
# -1092
# -1094
# -1096
# -1098
# -1100
# -1102
# -1104
# -1106
# -1108
# -1110
# -1112
# -1114
# -1116
# -1118
# -1120
# -1122
# -1124
# -1126
# -1128
# -1130
# -1132
# -1134
# -1136
# -1138
# -1140
# -1142
# -1144
# -1146
# -1148
# -1150
# -1152
# -1154
# -1156
# -1158
# -1160
# -1162
# -1164
# -1166
# -1168
# -1170
# -1172
# -1174
# -1176
# -1178
# -1180
# -1182
# -1184
# -1186
# -1188
# -1190
# -1192
# -1194
# -1196
# -1198
# -1200
# -1202
# -1204
# -1206
# -1208
# -1210
# -1212
# -1214
# -1216
# -1218
# -1220
# -1222
# -1224
# -1226
# -1228
# -1230
# -1232
# -1234
# -1236
# -1238
# -1240
# -1242
# -1244
# -1246
# -1248
# -1250
# -1252
# -1254
# -1256
# -1258
# -1260
# -1262
# -1264
# -1266
# -1268
# -1270
# -1272
# -1274
# -1276
# -1278
# -1280
# -1282
# -1284
# -1286
# -1288
# -1290
# -1292
# -1294
# -1296
# -1298
# -1300
# -1302
# -1304
# -1306
# -1308
# -1310
# -1312
# -1314
# -1316
# -1318
# -1320
# -1322
# -1324
# -1326
# -1328
# -1330
# -1332
# -1334
# -1336
# -1338
# -1340
# -1342
# -1344
# -1346
# -1348
# -1350
# -1352
# -1354
# -1356
# -1358
# -1360
# -1362
# -1364
# -1366
# -1368
# -1370
# -1372
# -1374
# -1376
# -1378
# -1380
# -1382
# -1384
# -1386
# -1388
# -1390
# -1392
# -1394
# -1396
# -1398
# -1400
# -1402
# -1404
# -1406
# -1408
# -1410
# -1412
# -1414
# -1416
# -1418
# -1420
# -1422
# -1424
# -1426
# -1428
# -1430
# -1432
# -1434
# -1436
# -1438
# -1440
# -1442
# -1444
# -1446
# -1448
# -1450
# -1452
# -1454
# -1456
# -1458
# -1460
# -1462
# -1464
# -1466
# -1468
# -1470
# -1472
# -1474
# -1476
# -1478
# -1480
# -1482
# -1484
# -1486
# -1488
# -1490
# -1492
# -1494
# -1496
# -1498
# -1500
# -1502
# -1504
# -1506
# -1508
# -1510
# -1512
# -1514
# -1516
# -1518
# -1520
# -1522
# -1524
# -1526
# -1528
# -1530
# -1532
# -1534
# -1536
# -1538
# -1540
# -1542
# -1544
# -1546
# -1548
# -1550
# -1552
# -1554
# -1556
# -1558
# -1560
# -1562
# -1564
# -1566
# -1568
# -1570
# -1572
# -1574
# -1576
# -1578
# -1580
# -1582
# -1584
# -1586
# -1588
# -1590
# -1592
# -1594
# -1596
# -1598
# -1600
# -1602
# -1604
# -1606
# -1608
# -1610
# -1612
# -1614
# -1616
# -1618
# -1620
# -1622
# -1624
# -1626
# -1628
# -1630
# -1632
# -1634
# -1636
# -1638
# -1640
# -1642
# -1644
# -1646
# -1648
# -1650
# -1652
# -1654
# -1656
# -1658
# -1660
# -1662
# -1664
# -1666
# -1668
# -1670
# -1672
# -1674
# -1676
# -1678
# -1680
# -1682
# -1684
# -1686
# -1688
# -1690
# -1692
# -1694
# -1696
# -1698
# -1700
# -1702
# -1704
# -1706
# -1708
# -1710
# -1712
# -1714
# -1716
# -1718
# -1720
# -1722
# -1724
# -1726
# -1728
# -1730
# -1732
# -1734
# -1736
# -1738
# -1740
# -1742
# -1744
# -1746
# -1748
# -1750
# -1752
# -1754
# -1756
# -1758
# -1760
# -1762
# -1764
# -1766
# -1768
# -1770
# -1772
# -1774
# -1776
# -1778
# -1780
# -1782
# -1784
# -1786
# -1788
# -1790
# -1792
# -1794
# -1796
# -1798
# -1800
# -1802
# -1804
# -1806
# -1808
# -1810
# -1812
# -1814
# -1816
# -1818
# -1820
# -1822
# -1824
# -1826
# -1828
# -1830
# -1832
# -1834
# -1836
# -1838
# -1840
# -1842
# -1844
# -1846
# -1848
# -1850
# -1852
# -1854
# -1856
# -1858
# -1860
# -1862
# -1864
# -1866
# -1868
# -1870
# -1872
# -1874
# -1876
# -1878
# -1880
# -1882
# -1884
# -1886
# -1888
# -1890
# -1892
# -1894
# -1896
# -1898
# -1900
# -1902
# -1904
# -1906
# -1908
# -1910
# -1912
# -1914
# -1916
# -1918
# -1920
# -1922
# -1924
# -1926
# -1928
# -1930
# -1932
# -1934
# -1936
# -1938
# -1940
# -1942
# -19
```