

The script contains relevant code used in the study. Finding the adaptive needles in a population-structured haystack is a task worthy of a New Zealand kōwhiri. Citations for all software and scripts are in the original article. The code below is free to be used, just cite our work (Salloom et al. 2022).

Notes on terminology:

- "all", generally referring to a dataset containing all SNPs for a set of populations
- "neutral", referring to a dataset containing outlier SNPs. It does not mean that the SNPs are truly evolutionary neutral
- "NZ", referring to a dataset containing outlier SNPs. It does not mean that the SNPs are adaptive.
- "NZ-wide" and "all conserved" refer to the same set of populations (all populations sampled across New Zealand)
- "NI" refers to North Island
- "SI" refers to South Island
- "Southern" refers to South and Sub-antarctic Islands/ Please refer to the original study for more information on terminology.

SNP calling using stacks and the same parameters used in the first submission of man 3 (m=3, M=2, n=4)

calling SNPs only for the North Island populations

```
In [ ]: cd ~/nesi/nobackup/uoa089538/programs/Stacks_0neq
mkdir samples_all
cd samples_all
cp AdaptationInStr/outcutedReps/Stacks/samples_all/* ./
cd ..
mkdir NI
nano cstacks_NI.sh

In [ ]: #!/bin/bash -e
#SBATCH --job-name=cstacksNI
#SBATCH --account=uoa089538
#SBATCH --partition=long
#SBATCH --time=80:00:00
#SBATCH --mem=24G
#SBATCH --cpus-per-task=16
#SBATCH --output=cstacks_NI%j.out

module purge
module load Stacks/2.2-gimkl-2017a

b=0
for sample in $(ls)
do
    files=$(ls)
    OneqEC01
    OneqEC02
    OneqEC03
    OneqEC04
    OneqEC05
    OneqEC06
    OneqEC07
    OneqEC08
    OneqEC09
    OneqEC10
    OneqEC11
    OneqEC12
    OneqEC13
    OneqEC14
    OneqEC15
    OneqEC16
    OneqEC17
    OneqEC18
    OneqEC19
    OneqEC20
    OneqEC21
    OneqEC22
    OneqEC23
    OneqEC24
    OneqEC25
    OneqEC26
    OneqEC27
    OneqEC28
    OneqEC29
    OneqEC30
    OneqEC31
    OneqEC32
    OneqEC33
    OneqEC34
    OneqEC35
    OneqEC36
    OneqEC37
    OneqEC38
    OneqEC39
    OneqEC40
    OneqEC41
    OneqEC42
    OneqEC43
    OneqEC44
    OneqEC45
    OneqEC46
    OneqEC47
    OneqEC48
    OneqEC49
    OneqEC50
    OneqEC51
    OneqEC52
    OneqEC53
    OneqEC54
    OneqEC55
    OneqEC56
    OneqEC57
    OneqEC58
    OneqEC59
    OneqEC60
    OneqEC61
    OneqEC62
    OneqEC63
    OneqEC64
    OneqEC65
    OneqEC66
    OneqEC67
    OneqEC68
    OneqEC69
    OneqEC70
    OneqEC71
    OneqEC72
    OneqEC73
    OneqEC74
    OneqEC75
    OneqEC76
    OneqEC77
    OneqEC78
    OneqEC79
    OneqEC80
    OneqEC81
    OneqEC82
    OneqEC83
    OneqEC84
    OneqEC85
    OneqEC86
    OneqEC87
    OneqEC88
    OneqEC89
    OneqEC90
    OneqEC91
    OneqEC92
    OneqEC93
    OneqEC94
    OneqEC95
    OneqEC96
    OneqEC97
    OneqEC98
    OneqEC99
    OneqEC100
    OneqEC101
    OneqEC102
    OneqEC103
    OneqEC104
    OneqEC105
    OneqEC106
    OneqEC107
    OneqEC108
    OneqEC109
    OneqEC110
    OneqEC111
    OneqEC112
    OneqEC113
    OneqEC114
    OneqEC115
    OneqEC116
    OneqEC117
    OneqEC118
    OneqEC119
    OneqEC120
    OneqEC121
    OneqEC122
    OneqEC123
    OneqEC124
    OneqEC125
    OneqEC126
    OneqEC127
    OneqEC128
    OneqEC129
    OneqEC130
    OneqEC131
    OneqEC132
    OneqEC133
    OneqEC134
    OneqEC135
    OneqEC136
    OneqEC137
    OneqEC138
    OneqEC139
    OneqEC140
    OneqEC141
    OneqEC142
    OneqEC143
    OneqEC144
    OneqEC145
    OneqEC146
    OneqEC147
    OneqEC148
    OneqEC149
    OneqEC150
    OneqEC151
    OneqEC152
    OneqEC153
    OneqEC154
    OneqEC155
    OneqEC156
    OneqEC157
    OneqEC158
    OneqEC159
    OneqEC160
    OneqEC161
    OneqEC162
    OneqEC163
    OneqEC164
    OneqEC165
    OneqEC166
    OneqEC167
    OneqEC168
    OneqEC169
    OneqEC170
    OneqEC171
    OneqEC172
    OneqEC173
    OneqEC174
    OneqEC175
    OneqEC176
    OneqEC177
    OneqEC178
    OneqEC179
    OneqEC180
    OneqEC181
    OneqEC182
    OneqEC183
    OneqEC184
    OneqEC185
    OneqEC186
    OneqEC187
    OneqEC188
    OneqEC189
    OneqEC190
    OneqEC191
    OneqEC192
    OneqEC193
    OneqEC194
    OneqEC195
    OneqEC196
    OneqEC197
    OneqEC198
    OneqEC199
    OneqEC200
    OneqEC201
    OneqEC202
    OneqEC203
    OneqEC204
    OneqEC205
    OneqEC206
    OneqEC207
    OneqEC208
    OneqEC209
    OneqEC210
    OneqEC211
    OneqEC212
    OneqEC213
    OneqEC214
    OneqEC215
    OneqEC216
    OneqEC217
    OneqEC218
    OneqEC219
    OneqEC220
    OneqEC221
    OneqEC222
    OneqEC223
    OneqEC224
    OneqEC225
    OneqEC226
    OneqEC227
    OneqEC228
    OneqEC229
    OneqEC230
    OneqEC231
    OneqEC232
    OneqEC233
    OneqEC234
    OneqEC235
    OneqEC236
    OneqEC237
    OneqEC238
    OneqEC239
    OneqEC240
    OneqEC241
    OneqEC242
    OneqEC243
    OneqEC244
    OneqEC245
    OneqEC246
    OneqEC247
    OneqEC248
    OneqEC249
    OneqEC250
    OneqEC251
    OneqEC252
    OneqEC253
    OneqEC254
    OneqEC255
    OneqEC256
    OneqEC257
    OneqEC258
    OneqEC259
    OneqEC260
    OneqEC261
    OneqEC262
    OneqEC263
    OneqEC264
    OneqEC265
    OneqEC266
    OneqEC267
    OneqEC268
    OneqEC269
    OneqEC270
    OneqEC271
    OneqEC272
    OneqEC273
    OneqEC274
    OneqEC275
    OneqEC276
    OneqEC277
    OneqEC278
    OneqEC279
    OneqEC280
    OneqEC281
    OneqEC282
    OneqEC283
    OneqEC284
    OneqEC285
    OneqEC286
    OneqEC287
    OneqEC288
    OneqEC289
    OneqEC290
    OneqEC291
    OneqEC292
    OneqEC293
    OneqEC294
    OneqEC295
    OneqEC296
    OneqEC297
    OneqEC298
    OneqEC299
    OneqEC300
    OneqEC301
    OneqEC302
    OneqEC303
    OneqEC304
    OneqEC305
    OneqEC306
    OneqEC307
    OneqEC308
    OneqEC309
    OneqEC310
    OneqEC311
    OneqEC312
    OneqEC313
    OneqEC314
    OneqEC315
    OneqEC316
    OneqEC317
    OneqEC318
    OneqEC319
    OneqEC320
    OneqEC321
    OneqEC322
    OneqEC323
    OneqEC324
    OneqEC325
    OneqEC326
    OneqEC327
    OneqEC328
    OneqEC329
    OneqEC330
    OneqEC331
    OneqEC332
    OneqEC333
    OneqEC334
    OneqEC335
    OneqEC336
    OneqEC337
    OneqEC338
    OneqEC339
    OneqEC340
    OneqEC341
    OneqEC342
    OneqEC343
    OneqEC344
    OneqEC345
    OneqEC346
    OneqEC347
    OneqEC348
    OneqEC349
    OneqEC350
    OneqEC351
    OneqEC352
    OneqEC353
    OneqEC354
    OneqEC355
    OneqEC356
    OneqEC357
    OneqEC358
    OneqEC359
    OneqEC360
    OneqEC361
    OneqEC362
    OneqEC363
    OneqEC364
    OneqEC365
    OneqEC366
    OneqEC367
    OneqEC368
    OneqEC369
    OneqEC370
    OneqEC371
    OneqEC372
    OneqEC373
    OneqEC374
    OneqEC375
    OneqEC376
    OneqEC377
    OneqEC378
    OneqEC379
    OneqEC380
    OneqEC381
    OneqEC382
    OneqEC383
    OneqEC384
    OneqEC385
    OneqEC386
    OneqEC387
    OneqEC388
    OneqEC389
    OneqEC390
    OneqEC391
    OneqEC392
    OneqEC393
    OneqEC394
    OneqEC395
    OneqEC396
    OneqEC397
    OneqEC398
    OneqEC399
    OneqEC400
    OneqEC401
    OneqEC402
    OneqEC403
    OneqEC404
    OneqEC405
    OneqEC406
    OneqEC407
    OneqEC408
    OneqEC409
    OneqEC410
    OneqEC411
    OneqEC412
    OneqEC413
    OneqEC414
    OneqEC415
    OneqEC416
    OneqEC417
    OneqEC418
    OneqEC419
    OneqEC420
    OneqEC421
    OneqEC422
    OneqEC423
    OneqEC424
    OneqEC425
    OneqEC426
    OneqEC427
    OneqEC428
    OneqEC429
    OneqEC430
    OneqEC431
    OneqEC432
    OneqEC433
    OneqEC434
    OneqEC435
    OneqEC436
    OneqEC437
    OneqEC438
    OneqEC439
    OneqEC440
    OneqEC441
    OneqEC442
    OneqEC443
    OneqEC444
    OneqEC445
    OneqEC446
    OneqEC447
    OneqEC448
    OneqEC449
    OneqEC450
    OneqEC451
    OneqEC452
    OneqEC453
    OneqEC454
    OneqEC455
    OneqEC456
    OneqEC457
    OneqEC458
    OneqEC459
    OneqEC460
    OneqEC461
    OneqEC462
    OneqEC463
    OneqEC464
    OneqEC465
    OneqEC466
    OneqEC467
    OneqEC468
    OneqEC469
    OneqEC470
    OneqEC471
    OneqEC472
    OneqEC473
    OneqEC474
    OneqEC475
    OneqEC476
    OneqEC477
    OneqEC478
    OneqEC479
    OneqEC480
    OneqEC481
    OneqEC482
    OneqEC483
    OneqEC484
    OneqEC485
    OneqEC486
    OneqEC487
    OneqEC488
    OneqEC489
    OneqEC490
    OneqEC491
    OneqEC492
    OneqEC493
    OneqEC494
    OneqEC495
    OneqEC496
    OneqEC497
    OneqEC498
    OneqEC499
    OneqEC500
    OneqEC501
    OneqEC502
    OneqEC503
    OneqEC504
    OneqEC505
    OneqEC506
    OneqEC507
    OneqEC508
    OneqEC509
    OneqEC510
    OneqEC511
    OneqEC512
    OneqEC513
    OneqEC514
    OneqEC515
    OneqEC516
    OneqEC517
    OneqEC518
    OneqEC519
    OneqEC520
    OneqEC521
    OneqEC522
    OneqEC523
    OneqEC524
    OneqEC525
    OneqEC526
    OneqEC527
    OneqEC528
    OneqEC529
    OneqEC530
    OneqEC531
    OneqEC532
    OneqEC533
    OneqEC534
    OneqEC535
    OneqEC536
    OneqEC537
    OneqEC538
    OneqEC539
    OneqEC540
    OneqEC541
    OneqEC542
    OneqEC543
    OneqEC544
    OneqEC545
    OneqEC546
    OneqEC547
    OneqEC548
    OneqEC549
    OneqEC550
    OneqEC551
    OneqEC552
    OneqEC553
    OneqEC554
    OneqEC555
    OneqEC556
    OneqEC557
    OneqEC558
    OneqEC559
    OneqEC560
    OneqEC561
    OneqEC562
    OneqEC563
    OneqEC564
    OneqEC565
    OneqEC566
    OneqEC567
    OneqEC568
    OneqEC569
    OneqEC570
    OneqEC571
    OneqEC572
    OneqEC573
    OneqEC574
    OneqEC575
    OneqEC576
    OneqEC577
    OneqEC578
    OneqEC579
    OneqEC580
    OneqEC581
    OneqEC582
    OneqEC583
    OneqEC584
    OneqEC585
    OneqEC586
    OneqEC587
    OneqEC588
    OneqEC589
    OneqEC590
    OneqEC591
    OneqEC592
    OneqEC593
    OneqEC594
    OneqEC595
    OneqEC596
    OneqEC597
    OneqEC598
    OneqEC599
    OneqEC600
    OneqEC601
    OneqEC602
    OneqEC603
    OneqEC604
    OneqEC605
    OneqEC606
    OneqEC607
    OneqEC608
    OneqEC609
    OneqEC610
    OneqEC611
    OneqEC612
    OneqEC613
    OneqEC614
    OneqEC615
    OneqEC616
    OneqEC617
    OneqEC618
    OneqEC619
    OneqEC620
    OneqEC621
    OneqEC622
    OneqEC623
    OneqEC624
    OneqEC625
    OneqEC626
    OneqEC627
    OneqEC628
    OneqEC629
    OneqEC630
    OneqEC631
    OneqEC632
    OneqEC633
    OneqEC634
    OneqEC635
    OneqEC636
    OneqEC637
    OneqEC638
    OneqEC639
    OneqEC640
    OneqEC641
    OneqEC642
    OneqEC643
    OneqEC644
    OneqEC645
    OneqEC646
    OneqEC647
    OneqEC648
    OneqEC649
    OneqEC650
    OneqEC651
    OneqEC652
    OneqEC653
    OneqEC654
    OneqEC655
    OneqEC656
    OneqEC657
    OneqEC658
    OneqEC659
    OneqEC660
    OneqEC661
    OneqEC662
    OneqEC663
    OneqEC664
    OneqEC665
    OneqEC666
    OneqEC667
    OneqEC668
    OneqEC669
    OneqEC670
    OneqEC671
    OneqEC672
    OneqEC673
    OneqEC674
    OneqEC675
    OneqEC676
    OneqEC677
    OneqEC678
    OneqEC679
    OneqEC680
    OneqEC681
    OneqEC682
    OneqEC683
    OneqEC684
    OneqEC685
    OneqEC686
    OneqEC687
    OneqEC688
    OneqEC689
    OneqEC690
    OneqEC691
    OneqEC692
    OneqEC693
    OneqEC694
    OneqEC695
    OneqEC696
    OneqEC697
    OneqEC698
    OneqEC699
    OneqEC700
    OneqEC701
    OneqEC702
    OneqEC703
    OneqEC704
    OneqEC705
    OneqEC706
    OneqEC707
    OneqEC708
    OneqEC709
    OneqEC710
    OneqEC711
    OneqEC712
    OneqEC713
    OneqEC714
    OneqEC715
    OneqEC716
    OneqEC717
    OneqEC718
    OneqEC719
    OneqEC720
    OneqEC721
    OneqEC722
    OneqEC723
    OneqEC724
    OneqEC725
    OneqEC726
    OneqEC727
    OneqEC728
    OneqEC729
    OneqEC730
    OneqEC731
    OneqEC732
    OneqEC733
    OneqEC734
    OneqEC735
    OneqEC736
    OneqEC737
    OneqEC738
    OneqEC739
    OneqEC740
    OneqEC741
    OneqEC742
    OneqEC743
    OneqEC744
    OneqEC745
    OneqEC746
    OneqEC747
    OneqEC748
    OneqEC749
    OneqEC750
    OneqEC751
    OneqEC752
    OneqEC753
    OneqEC754
    OneqEC755
    OneqEC756
    OneqEC757
    OneqEC758
    OneqEC759
    OneqEC760
    OneqEC761
    OneqEC762
    OneqEC763
    OneqEC764
    OneqEC765
    OneqEC766
    OneqEC767
    OneqEC768
    OneqEC769
    OneqEC770
    OneqEC771
    OneqEC772
    OneqEC773
    OneqEC774
    OneqEC775
    OneqEC776
    OneqEC777
    OneqEC778
    OneqEC779
    OneqEC780
    OneqEC781
    OneqEC782
    OneqEC783
    OneqEC784
    OneqEC785
    OneqEC786
    OneqEC787
    OneqEC788
    OneqEC789
    OneqEC790
    OneqEC791
    OneqEC792
    OneqEC793
    OneqEC794
    OneqEC795
    OneqEC796
    OneqEC797
    OneqEC798
    OneqEC799
    OneqEC800
    OneqEC801
    OneqEC802
    OneqEC803
    OneqEC804
    OneqEC805
    OneqEC806
    OneqEC807
    OneqEC808
    OneqEC809
    OneqEC810
    OneqEC811
    OneqEC812
    OneqEC813
    OneqEC814
    OneqEC815
    OneqEC816
    OneqEC817
    OneqEC818
    OneqEC819
    OneqEC820
    OneqEC821
    OneqEC822
    OneqEC823
    OneqEC824
    OneqEC825
    OneqEC826
    OneqEC827
    OneqEC828
    OneqEC829
    OneqEC830
    OneqEC831
    OneqEC832
    OneqEC833
    OneqEC834
    OneqEC835
    OneqEC836
    OneqEC837
    OneqEC838
    OneqEC839
    OneqEC840
    OneqEC841
    OneqEC842
    OneqEC843
    OneqEC844
    OneqEC845
    OneqEC846
    OneqEC847
    OneqEC848
    OneqEC849
    OneqEC850
    OneqEC851
    OneqEC852
    OneqEC853
    OneqEC854
    OneqEC855
    OneqEC856
    OneqEC857
    OneqEC858
    OneqEC859
    OneqEC860
    OneqEC861
    OneqEC862
    OneqEC863
    OneqEC864
    OneqEC865
    OneqEC866
    OneqEC867
    OneqEC868
    OneqEC869
    OneqEC870
    OneqEC871
    OneqEC872
    OneqEC873
    OneqEC874
    OneqEC875
    OneqEC876
    OneqEC877
    OneqEC878
    OneqEC879
    OneqEC880
    OneqEC881
    OneqEC882
    OneqEC883
    OneqEC884
    OneqEC885
    OneqEC886
    OneqEC887
    OneqEC888
    OneqEC889
    OneqEC890
    OneqEC891
    OneqEC892
    OneqEC893
    OneqEC894
    OneqEC895
    OneqEC896
    OneqEC897
    OneqEC898
    OneqEC899
    OneqEC900
    OneqEC901
    OneqEC902
    OneqEC903
    OneqEC904
    OneqEC905
    OneqEC906
    OneqEC907
    OneqEC908
    OneqEC909
    OneqEC910
    OneqEC911
    OneqEC912
    OneqEC913
    OneqEC914
    OneqEC915
    OneqEC916
    OneqEC917
    OneqEC918
    OneqEC919
    OneqEC920
    OneqEC921
    OneqEC922
    OneqEC923
    OneqEC924
    OneqEC925
    OneqEC926
    OneqEC927
    OneqEC928
    OneqEC929
    OneqEC930
    OneqEC931
    OneqEC932
    OneqEC933
    OneqEC934
    OneqEC935
    OneqEC936
    OneqEC937
    OneqEC938
    OneqEC939
    OneqEC940
    OneqEC941
    OneqEC942
    OneqEC943
    OneqEC944
    OneqEC945
    OneqEC946
    OneqEC947
    OneqEC948
    OneqEC949
    OneqEC950
    OneqEC951
    OneqEC952
    OneqEC953
    OneqEC954
    OneqEC955
    OneqEC956
    OneqEC957
    OneqEC958
    OneqEC959
    OneqEC960
    OneqEC961
    OneqEC962
    OneqEC963
    OneqEC964
    OneqEC965
    OneqEC966
    OneqEC967
    OneqEC968
    OneqEC969
    OneqEC970
    OneqEC971
    OneqEC972
    OneqEC973
    OneqEC974
    OneqEC975
    OneqEC976
    OneqEC977
    OneqEC978
    OneqEC979
    OneqEC980
    OneqEC981
    OneqEC982
    OneqEC983
    OneqEC984
    OneqEC985
    OneqEC986
    OneqEC987
    OneqEC988
    OneqEC989
    OneqEC990
    OneqEC991
    OneqEC992
    OneqEC993
    OneqEC994
    OneqEC995
    OneqEC996
    OneqEC997
    OneqEC998
    OneqEC999
    OneqEC1000
    OneqEC1001
    OneqEC1002
    OneqEC1003
    OneqEC1004
    OneqEC1005
    OneqEC1006
    OneqEC1007
    OneqEC1008
    OneqEC1009
    OneqEC1010
    OneqEC1011
    OneqEC1012
    OneqEC1013
    OneqEC1014
    OneqEC1015
    OneqEC1016
    OneqEC1017
    OneqEC1018
    OneqEC1019
    OneqEC1020
    OneqEC1021
    OneqEC1022
    OneqEC1023
    OneqEC1024
    OneqEC1025
    OneqEC1026
    OneqEC1027
    OneqEC1028
    OneqEC1029
    OneqEC1030
    OneqEC1031
    OneqEC1032
    OneqEC1033
    OneqEC1034
    OneqEC1035
    OneqEC1036
    OneqEC1037
    OneqEC1038
    OneqEC1039
    OneqEC1040
    OneqEC1041
    OneqEC1042
    OneqEC1043
    OneqEC1044
    OneqEC1045
    OneqEC1046
    OneqEC1047
    OneqEC1048
    OneqEC1049
    OneqEC1050
    OneqEC1051
    OneqEC1052
    OneqEC1053
    OneqEC1054
    OneqEC1055
    OneqEC1056
    OneqEC1057
    OneqEC1058
    OneqEC1059
    OneqEC1060
    OneqEC1061
    OneqEC1062
    OneqEC1063
    OneqEC1064
    OneqEC1065
    OneqEC1066
    OneqEC1067
    OneqEC1068
    OneqEC1069
    OneqEC1070
    OneqEC1071
    OneqEC1072
    OneqEC1073
    OneqEC1074
    OneqEC1075
    OneqEC1076
    OneqEC1077
    OneqEC1078
    OneqEC1079
    OneqEC1080
    OneqEC1081
    OneqEC1082
    OneqEC1083
    OneqEC1084
    OneqEC1085
    OneqEC1086
    OneqEC1087
    OneqEC1088
    OneqEC1089
    OneqEC1090
    OneqEC1091
    OneqEC1092
    OneqEC1093
    OneqEC1094
    OneqEC1095
    OneqEC1096
    OneqEC1097
    OneqEC1098
    OneqEC1099
    OneqEC1100
    OneqEC1101
    OneqEC1102
    OneqEC1103
    OneqEC1104
    OneqEC1105
    OneqEC1106
    OneqEC1107
    OneqEC1108
    OneqEC1109
    OneqEC1110
    OneqEC1111
    OneqEC1112
    OneqEC1113
    OneqEC1114
    OneqEC1115
    OneqEC1116
    OneqEC1117
    OneqEC1118
    OneqEC1119
    OneqEC1120
    OneqEC1121
    OneqEC1122
    OneqEC1123
    OneqEC1124
    OneqEC1125
    OneqEC1126
    OneqEC1127
    OneqEC1128
    OneqEC1129
    OneqEC1130
    OneqEC1131
    OneqEC1132
    OneqEC1133
    OneqEC1134
    OneqEC1135
    OneqEC1136
    OneqEC1137
    OneqEC1138
    OneqEC1139
    OneqEC1140
    OneqEC1141
    OneqEC1142
    OneqEC1143
    OneqEC1144
    OneqEC1145
    OneqEC1146
    OneqEC1147
    OneqEC1148
    OneqEC1149
    OneqEC1150
    OneqEC1151
    OneqEC1152
    OneqEC1153
    OneqEC1154
    OneqEC1155
    OneqEC1156
    OneqEC1157
    OneqEC1158
    OneqEC1159
    OneqEC1160
    OneqEC1161
    OneqEC1162
    OneqEC1163
    OneqEC1164
    OneqEC1165
    OneqEC1166
    OneqEC1167
    OneqEC1168
    OneqEC1169
    OneqEC1170
    OneqEC1171
    OneqEC1172
    OneqEC1173
    OneqEC1174
    OneqEC1175
    OneqEC1176
    OneqEC1177
    OneqEC1178
    OneqEC1179
    OneqEC1180
    OneqEC1181
    OneqEC1182
    OneqEC1183
    OneqEC1184
    OneqEC1185
    OneqEC1186
    OneqEC1187
    OneqEC1188
    OneqEC1189
    OneqEC1190
    OneqEC1191
    OneqEC1192
    OneqEC1193
    OneqEC1194
    OneqEC1195
    OneqEC1196
    OneqEC1197
    OneqEC1198
    OneqEC1199
    OneqEC1200
    OneqEC1201
    OneqEC1202
    OneqEC1203
    OneqEC1204
    OneqEC1205
    OneqEC1206
    OneqEC1207
    OneqEC1208
    OneqEC1209
    OneqEC1210
    OneqEC1211
    OneqEC1212
    OneqEC1213
    OneqEC1214
    OneqEC1215
    OneqEC1216
    OneqEC1217
    OneqEC1218
    OneqEC1219
    OneqEC1220
    OneqEC1221
    OneqEC1222
    OneqEC1223
    OneqEC1224
    OneqEC1225
    OneqEC1226
    OneqEC1227
    OneqEC1228
    OneqEC1229
    OneqEC1230
    OneqEC1231
    OneqEC1232
    OneqEC1233
    OneqEC1234
    OneqEC1235
    OneqEC1236
    OneqEC1237
    OneqEC1238
    OneqEC1239
    OneqEC1240
    OneqEC1241
    OneqEC1242
    OneqEC1243
    OneqEC1244
    OneqEC1245
    OneqEC1246
    OneqEC1247
    OneqEC1248
    OneqEC1249
    OneqEC1250
    OneqEC1251
    OneqEC1252
    OneqEC1253
    OneqEC1254
    OneqEC1255
    OneqEC1256
    OneqEC1257
    OneqEC1258
    OneqEC1259
    OneqEC1260
    OneqEC1261
    OneqEC1262
    OneqEC1263
    OneqEC1264
    OneqEC1265
    OneqEC1266
    OneqEC1267
    OneqEC1268
    OneqEC1269
    OneqEC1270
    OneqEC1271
    OneqEC1272
    OneqEC1273
    OneqEC1274
    OneqEC1275
    OneqEC1276
    OneqEC1277
    OneqEC1278
    OneqEC1279
    OneqEC1280
    OneqEC1281
    OneqEC1282
    OneqEC1283
    OneqEC1284
    OneqEC1285
    OneqEC1286
    OneqEC1287
    OneqEC1288
    OneqEC1289
    OneqEC1290
    OneqEC1291
    OneqEC1292
    OneqEC1293
    OneqEC1294
    OneqEC1295
    OneqEC1296
    OneqEC1297
    OneqEC1298
    OneqEC1299
    OneqEC1300
    OneqEC1301
    OneqEC1302
    OneqEC1303
    OneqEC1304
    OneqEC1305
    OneqEC1306
    OneqEC1307
    OneqEC1308
    OneqEC1309
    OneqEC1310
    OneqEC1311
    OneqEC1312
    OneqEC1313
    OneqEC1314
    OneqEC1315
    OneqEC1316
    OneqEC1317
    OneqEC1318
    OneqEC1319
    OneqEC1320
    OneqEC1321
    OneqEC1322
    OneqEC1323
    OneqEC1324
    OneqEC1325
    OneqEC1326
    OneqEC1327
    OneqEC1328
    OneqEC1329
    OneqEC1330
    OneqEC1331
    OneqEC1332
    OneqEC1333
    OneqEC1334
    OneqEC1335
    OneqEC1336
    OneqEC1337
    OneqEC1338
    OneqEC1339
    OneqEC1340
    OneqEC1341
    OneqEC1342
    OneqEC1343
    OneqEC1344
    OneqEC1345
    OneqEC1346
    OneqEC1347
    OneqEC1348
    OneqEC1349
    OneqEC1350
    OneqEC1351
    OneqEC1352
    OneqEC1353
    OneqEC1354
    OneqEC1355
    OneqEC1356
    OneqEC1357
    OneqEC1358
    OneqEC1359
    OneqEC1360
    OneqEC1361
    OneqEC1362
    OneqEC1363
    OneqEC1364
    OneqEC1365
    OneqEC1366
    OneqEC1367
    OneqEC1368
    OneqEC1369
    OneqEC1370
    OneqEC1371
    OneqEC1372
    OneqEC1373
    OneqEC1374
    OneqEC1375
    OneqEC1376
    OneqEC1377
    OneqEC1378
    OneqEC1379
    OneqEC1380
    OneqEC1381
    OneqEC1382
    OneqEC1383
    OneqEC1384
    OneqEC1385
    OneqEC1386
    OneqEC1387
    OneqEC1388
    OneqEC1389
    OneqEC1390
    OneqEC1391
    OneqEC1392
    OneqEC1393
    OneqEC1394
    OneqEC1395
    OneqEC1396
    OneqEC1397
    OneqEC1398
    OneqEC1399
    OneqEC1400
    OneqEC1401
    OneqEC1402
    OneqEC1403
    OneqEC1404
    OneqEC1405
    OneqEC1406
    OneqEC1407
    OneqEC1408
    OneqEC1409
    OneqEC1410
    OneqEC1411
    OneqEC1412
    OneqEC1413
    OneqEC1414
    OneqEC1415
    OneqEC1416
    OneqEC1417
    OneqEC1418
    OneqEC1419
    OneqEC1420
    OneqEC1421
    OneqEC1422
    OneqEC1423
    OneqEC14
```



```
In [ ] : cd ..//nemo
          nano keep_NI.txt
OneyGC01
OneyGC02
OneyGC03
OneyGC04
OneyGC05
OneyGC06
OneyGC07
OneyGC08
OneyGC09
OneyNC02
OneyNC03
OneyNC04
OneyNC05
OneyNC06
OneyNC07
OneyNC08
OneyNC09
OneyNC10
OneyNC11
OneyRU03
OneyRU04
OneyRU05
OneyRU06
OneyRU07
OneyRU08
OneyRU09
OneyRU13
OneyRU14
OneyRU15
OneyRU16
OneyRU17
OneyRU18
OneyRU19
OneyTI09
OneyTI10
OneyTI11
OneyTI12
OneyTI13
OneyTI14
OneyTI15
OneyTI16
OneyTI17
OneyTI18
OneyTI21
OneyTI22

# Filtering for specific individuals (keep a list of individuals saved in a file, with one individual name
per row and no header - names must match vcf file)
vcftools --vcf ../Conserved.vcf --keep keep_NI.txt --recode --recode-INFO-all --out NI.out

#Filtering out alleles with count smaller than 1 within a clade
vcftools --vcf NI.out.recode.vcf --mac 1 --recode --recode-INFO-all --out NI_mac1
```

BayeScan

From the directory where BayeScan was installed, sbatch the script below.

What is needed:

- BayeScan installed locally (<http://cmpg.unibe.ch/software/BayeScan/download.html>)
- VCF file converted to bayescan format (done using PGD spider, code at the end)
- The slurm script below to run BayeScan
- The R script below to plot results and identify outliers (the cut-off used was a q-value <= 0.05. The q-value is in one of BayeScan's output file, the one ending with _fst.txt)
- The outputs from the Bayesian run (_fst.txt, which will have the results and sel, which can be used to check the run for convergence)
- A csv file with ID info from the vcf file (loci_info_vstYcfile.csv), for matching the locus ID of the outliers found with BayeScan with the locus ID of the vcf file

```
In [ ] : less bayescenv.sh

#!/bin/bash -e
#SBATCH --job-name=bayescnev
#SBATCH --account=uoa09538
#SBATCH --time=72:00:00
#SBATCH --cpus-per-task=16
#SBATCH --ntasks=1
#SBATCH --mem=24G
#SBATCH --output=$slurm_out
#SBATCH --error=$slurm_err

./bayescan_2.1 allConservedToBayescan.txt -o all_ConservednotNeutral -out freq
```

```
In [ ] : # This file is used to plot figures for the software BayeScan in R.
```

```
# Copyright (C)
#
# This program
# it under the
# the Free Sof
```

```
# but WITHOUT ANY WARRANTY; without even the implied warranty of
# MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the
# GNU General Public License for more details.

# You should have received a copy of the GNU General Public License
# along with this program. If not, see <http://www.gnu.org/licenses/>.

# Arguments:
```

```
# - the q-value
# - size is the
# - pos is the
```

- ```
- highlight the optional list of outlier indices to display in red
name.highlighted allows to write the indices of highlighted markers instead of using a point like the o
other markers
- add.txt adds the indices of the outlier markers

Output:
- This function returns different parameters in a list
- o.outliers: the list of outliers
- nb.outliers: the number of outliers

Typical usage:
- load this file into R (file/source R code)
- in R, go to the directory where "output_fst.txt" is (file/change current dir)
- at the R prompt, type:
- plot_bayescan("output_fst.txt", g, FDR=0.05)
if you save the output in a variable, you can recall the different results:
results=plot_bayescan("output_fst.txt", g, FDR=0.05)
results$outliers
results$nb_outliers
```

```
#
plotting posterior distribution is very easy in R with the output of BayesCan:
First load the output file ".sel" produced by BayesCan
mydata=read.table("bi.sel",colClasses="numeric")
choose the parameter you want to plot by setting for example:
parameter="Fst1"
then this line will make the plot for:
plot(density(mydata[[parameter]]),xlab=parameter,main=paste(parameter,"posterior distribution"))
you can plot population specific Fst coefficient by setting
parameter="Fst1"
if you have non-codominant data you can plot posterior for FIS coefficients in each population:
parameter="Fis1"
if you test for selection, you can plot the posterior for alpha coefficient for selection:
parameter="alpha1"
you also have access to the likelihood with:
parameter="logL"
if you have the package "boa" installed, you can very easily obtain Highest Probability
Density Interval (HPDI) for your parameter of interest (example for the 95% interval):
> boa.hpdi(mydata[[parameter]][,0.05])

plot_bayescan<-function(res,FDR=0.05,size1,pos=0.35,highlight=NULL,name_highlighted=F,add_text=T)
{
 if (is.character(res))
 res=read.table(res)

 colfstat=5
 colq=colfstat-2

 highlight_rows=which(is.element(as.numeric(row.names(res)),highlight))
 non_highlight_rows=setdiff(1:nrow(res),highlight_rows)

 outliers=as.integer(row.names(res[res[,colq]<=FDR,]))

 ok.outliers=TRUE
 if (sum(res[,colq]<=FDR]==0)
 ok.outliers=FALSE;
}
```

```
res[,col]=0-.0001,colq]=0-.0001

Plot
plot(log10(res[,colq]),res[,colfstat],xlim=rev(range(log10(res[,colq]))),xlab="log10(q value)",ylab=names(
points(log10(res[,colq]=FDR,],col="red",cex=size)

if (name.highlighted) {
 if (length(highlight_rows)>0) {
 text(log10(res[highlight_rows,colq]),res[highlight_rows,colfstat],row.names(res[highlight_
rows,]),col="red",cex=size*1.2,font=2)
 }
} else {
 points(log10(res[highlight_rows,colq]),res[highlight_rows,colfstat],col="red",pch=19,cex=size)
 # add names of loci over p and vertical line
 if (ok.outliers & add.text) {
 text(log10(res[,colq]=FDR,],colq])~pos"(round(runif(nrow(res[,colq]=FDR,],1,2
)))^2-3),res[,colq]=FDR,],colfstat],row.names(res[,colq]=FDR,],cex=size)
 }
}
lines(c(log10(FDR),log10(FDR)),c(-1,1),lwd=2)
return(list("outliers"=outliers,"nb_outliers"=length(outliers)))
}

My data -----
#1. Read the function above into R
#Plot the Fst results
fst_results = plot_bayescan("all_ConservednotNeutral_fst.txt", 0, FDR=0.05)
plot_bayescan("all_ConservednotNeutral_fst.txt")
#Identify outliers
outliers = fst_results$outliers
length(outliers)
fst_results$nb_outliers #this just gives you the length of the outliers
write.table(outliers, "outliers.csv", sep = ",")

plotting posterior distribution of the parameters
#read the file with the .sel extension in R
my_sel = read.table("all_ConservednotNeutral.sel", colClasses = "numeric")
#define which parameter to plot
parameter = "Fst1" #this is the Fst in population 1
plot(density(my_sel[[parameter]]), xlab = parameter, main = paste(parameter, "posterior distribution"))
#obtain Highest Probability Density Interval (HPDI) for your parameter of interest (example for the 95% in
library(boa)
boa.hpd(my_sel[[parameter]], 0.05)

#plotting the above for all populations
fsts = c("Fst1", "Fst2", "Fst3", "Fst4", "Fst5", "Fst6", "Fst7", "Fst8", "Fst9", "Fst10", "Fst11", "Fst12", "Fst13", "F
st14", "Fst15", "Fst16")
for (i in fsts) {
 parameter = i
 ploth = plot(density(my_sel[[parameter]]), xlab = parameter, main = paste(parameter, "posterior distrib
ution"))
 plot2 = boa.hpd(my_sel[[parameter]], 0.05)
 print(plot1)
 print(plot2)
}

if you test for selection, you can plot the posterior for alpha coefficient for selection:
parameters="alpha1" ##there was no alpha in my_sel output file, only in the fst.txt output file.

you also have access to the likelihood with:
parameters="logL"
parameter = "logL"
plot(density(my_sel[[parameter]]), xlab = parameter, main = paste(parameter, "posterior distribution"))
boa.hpd(my_sel[[parameter]], 0.05)

##Customised code and plots for finding outliers in the Bayesian output
library(ggplot2)
results = read.table("all_ConservednotNeutral_fst.txt", header = TRUE)
#save a table with the q-values for each locus
locus = 1:10007
qvalues = results$qval
df = as.data.frame(cbind(locus, qvalues))
colnames(df) = c("locus", "qvalue")
#write table(df, "all_loci_qvalues.csv", sep = ",", row.names = FALSE)
ggplot(results, aes(x=fst, y=alpha)) + geom_point()
ggplot(results, aes(x=fst, y=alpha)) + geom_text(label = row.names(results))
ggplot(results, aes(x=fst, y=alpha, colour = qval)) + geom_point() + scale_colour_continuous()

##Classic Bayesian plot
plot_bayescan("all_ConservednotNeutral_fst.txt", 1, FDR = 0.05)
plot(results, aes(x = -log10(qval), y = fst)) +
 geom_point() +
 geom_segment(aes(y=min(fst), yend=max(fst)+0.2, x=1.3, xend=1.3), color="black") +
 labs(x = "-log10(qvalue)", y = "Fst") +
 theme_classic()

ggplot(results, aes(x=qval, y=alpha)) + geom_point()

which(results$alpha<0)
length(which(results$alpha<0))

#Distributions (all loci)
hist(results$qval)
hist(results$prob)
hist(results$logBP.Q)
hist(results$alpha)
hist(results$fst)

#Match Bayesian outliers with the loci ID in the vcf file
#read a csv file with ID info from the vcf file
all_loci = read.csv("loci_info_vcfStyle.csv", header = T)
#read a csv file with the outlier position detected by Bayescan (in the MOD dataset (in which I added line
numbers under the ID column in the vcf file), the outlier position matches the line number, which matches
the SNP ID in the vcf)
bayscanOut = read.csv("outliers.csv", header = T)

head(all_loci)
head(bayscanOut)

#id which loci are outliers for Bayescan and mark them with TRUE
all_loci$bayscanOUT = all_loci$ID %in% bayscanOut$outliers

#add that column to the table with vcf information
write.table(all_loci, "loci_info_vcfStyle.csv", sep = ",")

#make a list with the real vcf ID of the Bayesian outliers
only_outs = all_loci[all_loci$bayscanOUT==TRUE,]
dim(only_outs)

outs_id = only_outs$CHROM
length(outs_id)

write.table(outs_id, "vcfID_bayscanOutliers.csv", sep = ",")

PCAdapt

PCAdapt runs in R. Create a directory to save all tables and plots from PCAdapt.

What is needed:

- PCAdapt, qvalue, tidyverse and ggplot2 installed. For PCAdapt tutorial: https://bcn-uga.github.io/pcadapt/articles/pcadapt.html
- Filtered VCF file
- The R script below to run analyses, plot results and identify outliers (the cut-off used was a q-value < 0.05).
- A csv file with ID info from the vcf file (loci_info_vcfstyle.csv), for matching the locus ID of the outliers found with PCAdapt with the locus ID of the vcf file

In []:
#Load packages
library(PCAdapt)
library(qvalue)
library(ggplot2)

#Load data (in vcf format)
path.to.file <- "C:/.../files_vcf-env/popmap/all_conserved/conserved_MOD_vcf"

#read popmap input file (this file has information on individual ID, population ID, group or clade ID (e.g.,
#North America Central and Southern))
#read = read.csv("lid_pop_env.csv", header = T)

#You should use the read.pcadapt function to convert your genotype file to the pcadapt format.
#pcadapt files should have individuals in columns, SNPs in lines, and missing values should be
#recoded with 9's.
filename <- read.pcadapt(path.to.file, type = "vcf")

#Perform PCA to choose the number of K (should be done with a large enough number of PCs)
#The function will perform a PCA and also compute test statistics and p-values based on the
#correlation between SNPs and the principal components
#note is assumed to be diploid by default (but use ploidy = 2 or ploidy =1 to specify it)
x = PCAdapt(filename, K = 50)

#to check the proportion of variance of each PC, look at x$singular.values. This is the squared root of th
e proportion of
variance explained by each PC (https://bcn-uga.github.io/pcadapt/articles/pcadapt.html). Thus, to find th
e proportion of
variance, do
proportion_of_variance = (x$singular.values)^2
proportion_of_variance
class(proportion_of_variance)

#make a scree plot to display in decreasing order the percentage of variance explained by each PC
par(bg="white")
plot(x, option = "screeplot") + theme_classic() + scale_colour_manual (values = "black") + theme(text = el
ement_text(size=20))
plot(x, option = "screeplot", K = 10) + theme_classic() + scale_colour_manual (values = "black") #reduce t
he number of PCs plotted to only the 10 first

#Make scores plots to display population structure (PC1 vs.2 is default, but can be changed to
#display other PCs as well, and should be done to investigate until which PC there is population
#structure in the dataset)
poplist.int = info$pop #create a vector with the populations of each individual
grouplist.int = info$group # create a vector with the group of each individual
plot(x, option = "scores", pop = poplist.int)
plot(x, option = "scores", pop = grouplist.int)
plot(x, option = "scores", i = 3, j = 4, pop = poplist.int) # i sets the PC in the X axis and j sets the P
C in the Y axis
plot(x, option = "scores", i = 3, j = 4, pop = grouplist.int)
plot(x, option = "scores", i = 5, j = 6, pop = poplist.int)
plot(x, option = "scores", i = 5, j = 6, pop = grouplist.int)
plot(x, option = "scores", i = 7, j = 8, pop = poplist.int)
plot(x, option = "scores", i = 7, j = 8, pop = grouplist.int)
plot(x, option = "scores", i = 9, j = 10, pop = poplist.int)
plot(x, option = "scores", i = 9, j = 10, pop = grouplist.int)
plot(x, option = "scores", i = 11, j = 12, pop = poplist.int)
plot(x, option = "scores", i = 11, j = 12, pop = grouplist.int)
plot(x, option = "scores", i = 13, j = 14, pop = poplist.int)
plot(x, option = "scores", i = 13, j = 14, pop = grouplist.int)
plot(x, option = "scores", i = 15, j = 16, pop = poplist.int)
plot(x, option = "scores", i = 15, j = 16, pop = grouplist.int)

#CUSTOM PLOTS
dat = as.data.frame(x$scores)
#Defining a different colour for each population
#Specify the factor levels in the order you want
poplist.int = factor(info$pop, levels = c("RU", "TI", "NC", "EC", "WE", "CP", "CR", "OU", "AK", "CU", "ST")
"AF", "CA", "AN", "BO", "CH")
#Create a custom color scale (called pops.scale)
cols.pops = c("#c80000", "#f70000", "#990066", "#ff33cc", "#ff9900", "#ffc000", "#006600", "#000066", "#3399
93", "#33ff33", "#99cc66", "#3399ff", "#33ccff", "#9900cc", "#669999", "#000099")
cols.pops.pops = levels(poplist.int)
cols.pops.scale = scale_colour_manual(name = "poplist.int", values = cols.pops)
#plot data using the custom colour scale
#PC1 vs PC2
ggplot(dat, aes(x = dat[,1], y = dat[,2], colour = poplist.int)) +
 geom_point(size = 3, alpha = 1, position=position_dodge2(width=.01, padding = .2, preserve = "total")) +
 labs(x = "PC1",
 y = "PC2",
 title = "PC1 versus PC2, complete dataset") +
 cols.pops.scale +
 theme_classic()
#PC9 vs PC10
ggplot(dat, aes(x = dat[,9], y = dat[,10], colour = poplist.int)) +
 geom_point(size = 6) +
 labs(x = "PC9",
 y = "PC10",
 title = "PC9 versus PC10, complete dataset") +
 cols.pops.scale +
 theme_classic()
#Defining a different colour for each group/clade
#Specify the factor levels in the order you want
grouplist.int = factor(grouplist.int, levels = c("North_Is", "Central_NZ", "South_Is", "Subantarctics"))
#Create a custom color scale (called cols.group.scale)
cols.group = c("#ff9933", "#ff0000", "#3300ff", "#33ff00", "#9900cc", "#669999", "#000099")
names(cols.group) = levels(grouplist.int)
cols.group.scale = scale_colour_manual(name = "grouplist.int", values = cols.group)
#plot data using the custom colour scale
#PC1 vs PC2
ggplot(dat, aes(x = dat[,1], y = dat[,2], colour = grouplist.int)) +
 geom_point(size = 6) +
 labs(x = "PC1",
 y = "PC2",
 title = "PC1 versus PC2, complete dataset") +
 cols.group.scale +
 theme_classic()
#PC9 vs PC10
ggplot(dat, aes(x = dat[,9], y = dat[,10], colour = grouplist.int)) +
 geom_point(size = 6) +
 labs(x = "PC9",
 y = "PC10",
 title = "PC9 versus PC10, complete dataset") +
 cols.group.scale +
 theme_classic()
#Computing the test statistic based on PCA
#For a given SNP, the test statistic is based on the z-scores obtained when regressing SNPs
#against the K principal components. The test statistic for detecting outlier SNPs is the
#Mahalanobis distance, which is a multi-dimensional approach that measures how distant
#is a point from the mean. Once divided by a constant (lambda) taking the genomic inflation factor, the
#rescaled squared distances Q2/lambda should have a chi-squared distribution with K degrees of
#freedom under the assumption that there are no outliers.
#If K=9 or K=3 corresponds to the optimal choice of the number of PCs, do:
x <- readapt(filename, K = 9)
x3 <- readapt(filename, K = 3)
#In addition to the number K of principal components to work with, the user can also set the
#parameter min.maf that corresponds to a threshold of minor allele frequency. By default,
#the parameter min.maf is set to 5%. P-values of SNPs with a minor allele frequency smaller
#than the threshold are not computed (NA is returned). My dataset was filtered for MAF, so min.maf did not
#make a difference in this analysis.
summary(x)
summary(x3)
#We assume in the following that there are n individuals and L markers.
#stat is a vector of size L containing scaled Mahalanobis distances by default.
#pvalues is a vector containing L p-values.
#maf is a vector of size L containing minor allele frequencies.
#qval is a numerical value corresponding to the genomic inflation factor estimated from stat.
#ch2.stat is a vector of size L containing the rescaled statistics stat/gif that follow a chi-squared di
stribution with K degrees of freedom.
#scores is a (n,K) matrix corresponding to the projections of the individuals onto each PC.
#loadings is a (L,K) matrix containing the correlations between each genetic marker and each PC.
#singular.values is a vector containing the K ordered squared root of the proportion of variance explaine
d by this PC.
#zscores is a (L,K) matrix containing the z-scores.
#All of these elements are accessible using the $ symbol. For example, the p-values are contained in x$p
```

```
#Histograms of
#A histogram of
#and that the e
```

```
#The presence or outliers is also visible when plotting a histogram or the test statistic D).
plot(x, option = "stat.distribution")
plot(x3, option = "stat.distribution")

#Choosing a cutoff for outlier detection
#To provide a list of outliers, use the R package qvalue, transforming the
#p-values into q-values.
#For a given alpha (real valued number between 0 and 1), SNPs with q-values less than alpha will be
#considered as outliers with an expected false discovery rate bounded by alpha. The false discovery
```

```
To provide a list of candidate SNPs with an expected false discovery rate lower than 5%:
qval <- qvalue(x3spvalues)$qvalues
alpha <- 0.05
outliers <- which(qval < alpha)
length(outliers)

qval3 <- qvalue(x3spvalues)$qvalues
qvalues3= qvalue(x3spvalues)
alpha3 <- 0.05
outliers3 <- which(qval3 < alpha3)
```

```

hist(qval13, xlab = "q-values", main = NULL, breaks = 50, col = "orange")
write.table(outliers, 'outliers_qval095.csv', sep = ",")
#save outliers for each loci in a table
locus = 1:10987
qval13
df = cbind(locus, qval13)
class(df)
df = as.data.frame(df)
colnames(df) = c('locus', 'qvalue')
write.table(df, "all_loci_qvalues.csv", sep = ",", row.names = FALSE)

#manhattan plot with q-values
qval_plot = -log10(qval13)
locus = c(1:length(qval_plot))
qval_plot = data.frame(qval_plot, locus)
length(qval_plot)
plot(qval_plot,
 qseg=segment(aes(x=0, xend=9275, y=1.3, yend=1.3), color="red", size = 1) + theme_classic()
 #adds a red line to q-value = 0.05 (-log10(0.05) = 1.3)
 ,
 npar
 ggplot(qval_plot, aes(x = locus, y = qval_plot)) +
 geom_point() +
 qseg=segment(aes(x=0, xend = length(x$3pvalues), y = 1.3, yend=1.3), color = "red", size = 1) +
 theme_classic() +
 labs(title = "Manhattan plot", y = "-log10(qvalue)", x = "locus")
 axis.title = element_text(size = 12, family = "arial")
 #adds a red line to q-value = 0.05 (-log10(0.05) = 1.3)
)

#distribution and summary
hist(values$S1fdr)
summary(values$S1fdr)

#To know which principal components are actually the most correlated with the outlier SNPs, use
#the function get.pc()
snp_pc <- get.pc(x3, outliers)
write.table(snp_pc, "outliers_pcs.tsv", sep = "\t")
snp3_pc <- get.pc(x3, outliers)
write.table(snp3_pc, "outliers_pcs3.tsv", sep = "\t")

#####
#Match PCAdapt outliers with the loci ID in the vcf file
#read a PCAdapt with ID info from the vcf file
all_loci = read.csv("loci_info_vcfStyle.csv", header = T)
#read a csv file with the outlier position detected by PCAdapt (in the complete dataset, the outlier position matches the line number, which matches the SNP ID in the vcf)
pcadapt_out = read.csv("outliers_qval095.csv", header = T)

head(all_loci)
head(pcadapt_out)

#ID which loci are outliers for pcardapt and mark them with TRUE
all_loci$pcadaptOUT = all_loci$ID %in% pcardapt_out$locus

#add that column to the table with vcf information
write.table(all_loci, "loci_info_vcfStyle.csv", sep = ",")

#make a list with the real vcf ID of the PCAdapt outliers
only_outs = all_loci[all_loci$pcadaptOUT==TRUE,]
dim(only_outs)

outs_id = only_outs$CHROM
length(outs_id)

write.table(outs_id, "vcfID_pcardaptOutliers.csv", sep = ",")

```

```
pca = read.csv("pcadapt_all_loci_qvalues.csv", header = T)
bay = read.csv("bayescan_all_loci_qvalues.csv", header = T)

#Combine the q-values of pcadapt and bayescan
head(pca)
head(bay)

#create a function to calculate the geometric mean of q-values
gmean = function(x, y) {
 sqrt(x*y)
}
```

```
#combine the pcadapt and bayescan outputs using the geometric mean of q-values (with the function above)
combgq = gmean(pca$Qvalue, bay$Qvalue)
head(combgq)

#Make a dataframe with all results (combined and uncombined)
df = as.data.frame(cbind(pca$Qvalue, pca$Qvalue, bay$Qvalue, combgq))
head(df)
colnames(df) = c("locus", "PCAdapt_qval", "Bayescan_qval", "Combined_qval")
summary(df)
```

- ```
#id outliers using the combined q-value
df$combined_out = ifelse(df$combined_qval <= 0.05, TRUE, FALSE)

#id outliers using the intersection of methods
#pcadapt outliers
df$pcadapt_out = ifelse(df$PCAdapt_qval <= 0.05, TRUE, FALSE)
#bayescan outliers
df$bay_out = ifelse(df$Bayescan_qval <= 0.05, TRUE, FALSE)
```

```

function dfSpca_out == TRUE & dfSbay_out == TRUE, TRUE, FALSE)
summary(df)

write.table(df, "outliers_combined_quals.csv", sep = ",", row.names = FALSE)

```

BayesScEnv

BayesScEnv was run with real environmental variables and with random variables. What is needed for both cases

- BayesScEnv installed locally
- The filtered VCF file converted to Bayescan format (which is also used by BayesScEnv). This was done with PGD spider, code at the end.

Real Environmental Variables: What is needed

- A file for each environmental variable (example giver for mean air temperature below, each column has the value for each population. Specifications can be found in BayesScEnv manual). More information on obtaining and extracting the environmental variables is also given below.
- The slurm script below to run BayesScEnv
- The R script below to plot results and identify outliers (the cut-off used was a q-value <= 0.05. The q-value is in one of BayesScEnv's output file, the one ending with _fstxt)
- The outputs from the BayesScEnv runs (two for each environmental variable: _fstxt, which will have the results and _sel, which can be used to check the run for convergence)

```

In [ ]: less mean_air.txt
0 - 1.87538249      - 1.37885565      - 1.140623761      - 0.988321799      - 1.678099623      - 0.291675506      - 0.42891364
1      0.452218721      - 0.716593907      - 0.187538249      - 1.199136564      - 1.188097937      - 1.498130364      -
6.679646409      - 1.498130364      - 0.428913641

```

```

In [ ]: #!/bin/bash -e

#SBATCH --job-name=bayescenv
#SBATCH --account=uoa09538
#SBATCH --time=48:00:00
#SBATCH --cpus-per-task=16
#SBATCH --ntasks=1
#SBATCH --mem=24G
#SBATCH --output=slurm_out
#SBATCH --error=slurm_err

#./bayescenv allTobayescenv.txt -env env_file_bayescenv.txt -o all_out -nbp 20 -pilot 5000 -thin 10 - n 5000 -burn 5000 -out_pilot_out_freq
./bayescenv allTobayescenv.txt -env meanAir.txt -o all_out -nbp 20 -pilot 5000 -thin 10 - n 5000 -burn 5000 -out_pilot_out_freq
./bayescenv allTobayescenv.txt -env meanLat.txt -o all_out -nbp 20 -pilot 5000 -thin 10 - n 5000 -burn 5000 -out_pilot_out_freq
./bayescenv allTobayescenv.txt -env meanLong.txt -o all_out -nbp 20 -pilot 5000 -thin 10 - n 5000 -burn 5000 -out_pilot_out_freq
./bayescenv allTobayescenv.txt -env meanMslp.txt -o all_out -nbp 20 -pilot 5000 -thin 10 - n 5000 -burn 5000 -out_pilot_out_freq
./bayescenv allTobayescenv.txt -env meanKwsrs.txt -o all_out -nbp 20 -pilot 5000 -thin 10 - n 5000 -burn 5000 -out_pilot_out_freq
./bayescenv allTobayescenv.txt -env meanSalt.txt -o all_out -nbp 20 -pilot 5000 -thin 10 - n 5000 -burn 5000 -out_pilot_out_freq
./bayescenv allTobayescenv.txt -env meanOtem.txt -o all_out -nbp 20 -pilot 5000 -thin 10 - n 5000 -burn 5000 -out_pilot_out_freq
./bayescenv allTobayescenv.txt -env meanPrate.txt -o all_out -nbp 20 -pilot 5000 -thin 10 - n 5000 -burn 5000 -out_pilot_out_freq
./bayescenv allTobayescenv.txt -env meanPres.txt -o all_out -nbp 20 -pilot 5000 -thin 10 - n 5000 -burn 5000 -out_pilot_out_freq
./bayescenv allTobayescenv.txt -env meanSalt.txt -o all_out -nbp 20 -pilot 5000 -thin 10 - n 5000 -burn 5000 -out_pilot_out_freq
./bayescenv allTobayescenv.txt -env meanShq.txt -o all_out -nbp 20 -pilot 5000 -thin 10 - n 5000 -burn 5000 -out_pilot_out_freq
./bayescenv allTobayescenv.txt -env meanSst.txt -o all_out -nbp 20 -pilot 5000 -thin 10 - n 5000 -burn 5000 -out_pilot_out_freq
./bayescenv allTobayescenv.txt -env meanTcdc.txt -o all_out -nbp 20 -pilot 5000 -thin 10 - n 5000 -burn 5000 -out_pilot_out_freq
./bayescenv allTobayescenv.txt -env meanUwnd.txt -o all_out -nbp 20 -pilot 5000 -thin 10 - n 5000 -burn 5000 -out_pilot_out_freq
./bayescenv allTobayescenv.txt -env meanWmd.txt -o all_out -nbp 20 -pilot 5000 -thin 10 - n 5000 -burn 5000 -out_pilot_out_freq

```



```
In [ ]: #Load packages
library(coda)
library(tibble)
library(ggplot2)

# 1. Analysing convergence and results for each environmental variable separately

#lat
filename = "meanLat_out_sel"
chain <- read.table(filename, header = TRUE)
chain <- mcmc(chain, thin = 10) #Adapt thin to its actual value (10 is the default)
heidel.diag(chain) #To test for convergence
effectiveSize(chain) #To compute effective sample size (good above 1000)
autocorr.diag(chain) #To look for auto-correlation (good below 0.7)
par(mar=c(1,1,1,1))
plot(chain) #To plot the "trace" and the posterior distribution

tablename = "meanLat_out_fst.txt"
dat = read.table(tablename, header = TRUE)
head(dat)
colnames(dat)
dat = add_column(dat, locus = 1:nrow(dat), before = "PEP_g")
dat
dim(dat)
#fst
ggplot(dat, aes(x = locus, y=fst)) + geom_point() + labs(title = tablename)
#fst distribution
d1st(d1stfst)
gg
ggplot(dat, aes(x = locus, y=g)) + geom_point() + labs(title = tablename)
#alpha
ggplot(dat, aes(x = locus, y=log10(qval.g)) + geom_point() + labs(title = tablename)
#-log10(PEP.g)
ggplot(dat, aes(x = locus, y=-log10(PEP.g))) + geom_point() + labs(title = tablename)
#-log10(qval.g)
ggplot(dat, aes(x = locus, y=-log10(qval.g)) + geom_point() + labs(title = tablename)
#-log10(PEP_alpha)
ggplot(dat, aes(x = locus, y=-log10(PEP_alpha)) + geom_point() + labs(title = tablename)
#-log10(qval_alpha)
ggplot(dat, aes(x = locus, y=-log10(qval_alpha)) + geom_point() + labs(title = tablename)
#long
ggplot(dat, aes(x = locus, y=log10(qval.g), y=fst)) + geom_point() + labs(title = tablename)
geom_segment(aes(y=min(fst), yend=max(fst)+0.2, x=1.3, xend=1.3), color="black") + theme_classic()
#PEP versus -log10(qval_alpha)
ggplot(dat, aes(x = -log10(qval_alpha), y=fst)) + geom_point() + labs(title = tablename)
geom_segment(aes(y=min(fst), yend=max(fst)+0.2, x=1.3, xend=1.3), color="black") + theme_classic()
#Rise and repeat. Change the variables below and re-run the code above
#long
filename = "meanLong_out_fst.txt"
tablename = "meanLong_out_fst.txt"
#lat
filename = "meanLat_out_sel"
tablename = "meanLat_out_fst.txt"
#mslp
filename = "meanMslp_out_sel"
tablename = "meanMslp_out_fst.txt"
#meanNsws
filename = "meanNsws_out_sel"
tablename = "meanNsws_out_fst.txt"
#meanOsalT
filename = "meanOsalT_out_sel"
tablename = "meanOsalT_out_fst.txt"
#meanOtem
filename = "meanOtem_out_sel"
tablename = "meanOtem_out_fst.txt"
#meanPrate
filename = "meanPrate_out_sel"
tablename = "meanPrate_out_fst.txt"
#meanPres
filename = "meanPres_out_sel"
tablename = "meanPres_out_fst.txt"
#meanSalt
filename = "meanSalt_out_sel"
tablename = "meanSalt_out_fst.txt"
#meanSshg
filename = "meanSshg_out_sel"
tablename = "meanSshg_out_fst.txt"
#meanSst
filename = "meanSst_out_sel"
tablename = "meanSst_out_fst.txt"
#meanTcdc
filename = "meanTcdc_out_sel"
tablename = "meanTcdc_out_fst.txt"
#meanUwnd
filename = "meanUwnd_out_sel"
tablename = "meanUwnd_out_fst.txt"
#meanVwnd
filename = "meanVwnd_out_sel"
tablename = "meanVwnd_out_fst.txt"
#2. Counting loci with sigifican alpha and q' (-log10(qval) = 1.3)

env_variable = "meanLat_out_fst.txt"
df = read.table(env_variable, header= TRUE)
g = which (df$qval.g < 0.05)
length(g)
alpha = which (df$qval_alpha < 0.05)
length(alpha)
print(g)
print(alpha)
#Rise and repeat. Change the variables below and re-run the code above
env_variable = "meanLong_out_fst.txt"
env_variable = "meanMslp_out_fst.txt"
env_variable = "meanNsws_out_fst.txt"
env_variable = "meanOsalT_out_fst.txt"
env_variable = "meanOtem_out_fst.txt"
env_variable = "meanSalt_out_fst.txt"
env_variable = "meanSshg_out_fst.txt"
env_variable = "meanSst_out_fst.txt"
env_variable = "meanTcdc_out_fst.txt"
env_variable = "meanUwnd_out_fst.txt"
env_variable = "meanVwnd_out_fst.txt"

#3. (optional) extract all outliers and make a data frame with it
na.pad <- function(x,len){
  1:len
}
makePaddedDataFrame <- function(1,...){
  maxLen <- max(sapply(1,length))
  data.frame(lapply(1,na.pad,len=maxlen),...)
}

library(magickfor)
file_names=dir(pattern="*_out_fst.txt")
magic_for(print, silent = TRUE)

for (f in file_names) {
  df = read.table(f, header = TRUE)
  g = which (df$qval.g < 0.05)
  length(g)
  alpha = which (df$qval_alpha < 0.05)
  length(alpha)
  print(g)
  print(alpha)
}
results = magic_result_as_vector()
outliers_realvar = makePaddedDataFrame(results)
write.table(outliers_realvar, 'outliers_realvar.csv', sep = ",",)

#
library(dplyr)
library(magickfor)
filelist=dir(pattern="*_out_fst.txt")
magic_for(print, silent = TRUE)

for (file in filelist) {
  df = read.table(file, header = TRUE)
  select (df,
    colnames=df) <- c(file, paste(file, 'qval.g', sep = ""), paste(file, 'qval_alpha', sep = ""))
  results = magic_result_as_vector()
  head(results)
  results
}

#4. Store the qvalues for alpha and g in a different object for each environmental variable and make a data
aframe with all results (a different env variable in each column for g and alpha)
this will be the most helpful table to define outliers later on
data_all <- data.frame("1:10987,1" #change the 10987 to reflect the number of loci in the dataset
head(data_all)
data_all <- select(data_all,1)
colnames(data_all) = 'loci'
head(data_all)
lat = read.table("meanLat_out_fst.txt", header = T)
data_all$latG = lat$qval.g
data_all$latAlpha = lat$qval_alpha
#long data
long = read.table("meanLong_out_fst.txt", header = T)
data_all$longG = long$qval.g
data_all$longAlpha = long$qval_alpha
#mslp data
mslp = read.table("meanMslp_out_fst.txt", header = T)
data_all$mslpG = mslp$qval.g
data_all$mslpAlpha = mslp$qval_alpha
#nsws data
nsws = read.table("meanNsws_out_fst.txt", header = T)
data_all$nswsG = nsws$qval.g
data_all$nswsAlpha = nsws$qval_alpha
#osalT data
osalT = read.table("meanOsalT_out_fst.txt", header = T)
data_all$osalTG = osalT$qval.g
data_all$osalTAlpha = osalT$qval_alpha
#otemp data
otemp = read.table("meanOtem_out_fst.txt", header = T)
data_all$otempG = otemp$qval.g
data_all$otempAlpha = otemp$qval_alpha
#prate data
prate = read.table("meanPrate_out_fst.txt", header = T)
data_all$prateG = prate$qval.g
data_all$prateAlpha = prate$qval_alpha
#pres data
pres = read.table("meanPres_out_fst.txt", header = T)
data_all$presG = pres$qval.g
data_all$presAlpha = pres$qval_alpha
#salt data
salt = read.table("meanSalt_out_fst.txt", header = T)
data_all$saltG = salt$qval.g
data_all$saltAlpha = salt$qval_alpha
#sshg data
sshg = read.table("meanSshg_out_fst.txt", header = T)
data_all$sshgG = sshg$qval.g
data_all$sshgAlpha = sshg$qval_alpha
#sst data
sst = read.table("meanSst_out_fst.txt", header = T)
data_all$ssTG = sst$qval.g
data_all$ssTAlpha = sst$qval_alpha
#tdc data
tdc = read.table("meanTcdc_out_fst.txt", header = T)
data_all$tdcG = tdc$qval.g
data_all$tdcAlpha = tdc$qval_alpha
#uwnd data
uwnd = read.table("meanUwnd_out_fst.txt", header = T)
data_all$uwndG = uwnd$qval.g
data_all$uwndAlpha = uwnd$qval_alpha
#vwnd data
vwnd = read.table("meanVwnd_out_fst.txt", header = T)
data_all$vwndG = vwnd$qval.g
data_all$vwndAlpha = vwnd$qval_alpha
#Air data
air = read.table("meanAir_out_fst.txt", header = T)
data_all$airG = air$qval.g
data_all$airAlpha = air$qval_alpha

head(data_all)
dim(data_all)
getwd()
write.table(data_all, "summary_qval_allRealEnv.csv", sep = ",", row.names = FALSE)

#5. Define outliers using the summary_qval_allRealEnv.csv created above
data_all <- read.csv("summary_qval_allRealEnv.csv", header = T)
head(data_all)
str(data_all)
#lat outs
length(which(data_all$latAlpha < 0.05)) #1391
length(which(data_all$latG < 0.05)) #0
#long outs
length(which(data_all$longAlpha < 0.05)) #1437
length(which(data_all$longG < 0.05)) #0
#mslp
length(which(data_all$mslpAlpha < 0.05)) #1486
length(which(data_all$mslpG < 0.05)) #0
#nswr
length(which(data_all$nswsAlpha < 0.05)) #2086
length(which(data_all$nswsG < 0.05)) #2784
#osalT
length(which(data_all$osalAlpha < 0.05)) #140
length(which(data_all$osalTG < 0.05)) #3064
#otemp
length(which(data_all$otempAlpha < 0.05)) #477
length(which(data_all$otempG < 0.05)) #485
#prate
length(which(data_all$prateAlpha < 0.05)) #1221
length(which(data_all$prateG < 0.05)) #0
#pres
length(which(data_all$presAlpha < 0.05)) #865
length(which(data_all$presG < 0.05)) #0
#salt
length(which(data_all$saltAlpha < 0.05)) #333
length(which(data_all$saltG < 0.05)) #1252
#sshg
length(which(data_all$sshgAlpha < 0.05)) #1533
length(which(data_all$sshgG < 0.05)) #0
#sst
length(which(data_all$ssTG < 0.05)) #1177
length(which(data_all$ssTG < 0.05)) #0
#tdc
length(which(data_all$tdcAlpha < 0.05)) #1473
length(which(data_all$tdcG < 0.05)) #0
#uwnd
length(which(data_all$uwndAlpha < 0.05)) #186
length(which(data_all$uwndG < 0.05)) #2872
#vwnd
length(which(data_all$vwndAlpha < 0.05)) #1736
length(which(data_all$vwndG < 0.05)) #0
#air
length(which(data_all$airAlpha < 0.05)) #474
length(which(data_all$airG < 0.05)) #12

#allAlpha
library(tidyverse)
allAlpha = data_all %>% select (1, 3, 5, 7, 9, 11, 13, 15, 17, 19, 21, 23, 25, 27, 29, 31)
allAlpha$locID = allAlpha$loci
allAlpha = %>%
  gather(loci, 'qval_alpha', 1:16) %>%
  filter (qval_alpha < 0.05)
head(allAlpha)
colnames(allAlpha) = c('loci', 'env_variable', 'qval_alpha')
head(allAlpha)
head(allAlpha)
#duplicate alpha = (duplicated(allAlpha$loci))==TRUE)
length(which(duplicated_alpha==FALSE)) #1736 Unique outliers
summary(duplicated_alpha)
summary(allAlpha)

#allG
library(tidyverse)
allG = data_all %>% select (1, 2, 4, 6, 8, 10, 12, 14, 16, 18, 20, 22, 24, 26, 28, 30)
allG$locID = allG$loci
allG = allG %>%
  gather(loci, 'qval_g', 1:16) %>%
  filter (qval_g < 0.05)
head(allG)
colnames(allG) = c('loci', 'env_variable', 'qval_g')
head(allG)
duplicated_g = (duplicated(allG$loci))==TRUE)
length(which(duplicated_g==FALSE)) #3093 Unique outliers
summary(duplicated_g)
summary(allG)


```

Random Variables: There is needed

- An R code to generate 10 random variables to be used as environmental files (will be reused in the RDA analyses)
- A population map (ld.pop.csv)

- The slurm script below to run BayesEnv, in a directory that has a subdirectory called env_all
- The R script below to identify outliers (the cut-off used was a q-value <= 0.05. The q-value is in one of BayesEnv's output file, the one ending with .fst.txt)

- The outputs from the BayesEnv run on each random variable (focus on the .fst.txt, which will have the results. The .sel will also be generated, and can be used to check the run for convergence)

```
In [ ]: #generate a set of 16 random values to be the environmental variables, but do that 100 times
nrSamples = 100 #how many times you want to generate the variables
e <- list(nodes="vector",length=nrSamples) # "e" is a list to store the different variables
for (i in 1:nrSamples) {
  e[[i]] <- rnorm(n = 16, mean = 0, sd = 1) #for R script generates 16 random variables normally distributed for
  #nrSamples length (100 in this case)
}
e = as.data.frame(e) #convert the list e to a data frame

#For BayesEnv -----
#save each of the columns of e (each of the 100 sets of 16 random variables) to a different env file
for (i in 1:ncol(e)) {
  filename = paste(i, "env", sep = "") #define each of the 100 file names
  write.table(e[[i]], filename, sep = ",", row.names = FALSE, col.names = FALSE) #save each of the co
  lnames of the dataframe e in a different file
}

#For RDA -----
#save it as a dataframe with information on population ID
dim(e)
colnames(e) = c(1:100)
colnames(e) = paste("X", colnames(e), sep = "")
class(colnames(e))
class(e)
esppos = as.character(rownames(e))
head(e)
dim(e)
#read individual, population and group information
pop.info= read.csv("ld.pop.csv")
#organize environmental data based on population information
env_all = merge(pop.info, e, by.x = 2, by.y = 0, all.x = TRUE) # The order of the data frames to merge 1
s important -> pop.info should come first, as e will be replicated for each pop
dim(env_all)
env_all$ind <- as.character(env_all$ind) # Make individual names characters (not factors)
env_all$pop <- as.character(env_all$pop)
env_all$group <- as.character(env_all$group)
write.table(env_all, "all_arbitraryEnv.csv", sep = "")

In [ ]: #!/bin/bash -e
#SBATCH --job-name=BayesEnv
#SBATCH --account=usd0858
#SBATCH --time=36:00:00
#SBATCH --cpus-per-task=8
#SBATCH --mem=1600
#SBATCH --array=1-100
#SBATCH --output=slurm%j.%a.out

./BayesEnv all7BayesEnv.txt \
  -env ./env_all/${SLURM_ARRAY_TASK_ID}.env \
  -o ./env_all/${SLURM_ARRAY_TASK_ID}.out \
  -np 20 -pilot 5000 -thin 10 - n 5000 -burn 50000 \
  -out_pilot -out_freq

In [ ]: file_names=dir(pattern="*_out_sel")

#Checking if the MCMC run converged
library(coda)
for (f in file_names) {
  chain <- read.table(f, header = TRUE)
  chain <- mcmc(chain, thin = 10) #Adapt thin to its actual value (10 is the default)
  heidel.diag(chain) #To test for convergence
  effectiveSize(chain) #To compute effective sample size (good above 1000)
  autocorr.diag(chain) #To look for auto-correlation (good below 0.7)
  par(mar=c(1,1,1,1))
  plot(chain) #To plot the "trace" and the posterior distribution
}

#####
# SUMMARIZE DATA INCLUDING RANDOM LOCI, RANDOM ENV, AND q-value
#to read each data.frame and combine them into one big dataframe
library(tidyverse)
#sewdl("env_all")
#If you want to create a column in each of the data frames from the results with the locus ID, the loop be
#will do it and save each data frame as a new table in the working directory
#Create a column in all dataframes for merging with rbind later (locus ID)
list.files()
file_names=dir(pattern="*_out_fst.txt")
loci=1:10987
for (f in file_names) {
  df1 = read.table(f, header = T)
  df1$loc = 1:10987
  write.table(df1, paste(f, ".MOD.csv", sep = ""), sep = ",")
}
allframes = lapply(1:100, function(x){read.csv(paste(x,"out_fst.txt.MOD.csv",sep=''))})
sapply(allframes,dim) #all data frames have all the loci 10987
df1 = bind_rows(allframes) #if you bind them like this, each run will take 13094 rows, so each variable
will be organized under the corresponding column
df1$env = rep(1:100, each = 10987)
summary(df1)
#res(r, times = 1, length.out = NA, each = 1)

#Select only the columns of qval, each = 1) across env variables, a
#transform q-values into 0 (larger than 0.05) and 1 (smaller than 0.05)
#g
gdf = bigdf %>%
  select("loci", "env", "qval.g")
gdf = gdf %>% spread(key = "env", value = "qval.g")
str(gdf)
head(gdf)
summary(gdf)
write.table(gdf, "../summary2021_qvaluesRandomEnv_BayesEnv.csv", sep = ",", row.names = F)


```

RDA

RDA implemented on the R package vegan was run with real environmental variables and with random variables. In addition, the test-statistic from the RDA is z-scores (coming from RDA loadings) and it was converted to q-values with a custom script. What is needed for running with both real environmental and random variables:

- The vegan and qvalue packages in R (in addition to other packages - check R script below)
- A population map
- An environmental file (a csv file with individual, population and environmental data information for each individual sample. In the case of the random variables, the environmental variables were replaced by the same variables used in the random runs with BayesEnv - code above) The RDA does not run with missing data, so imputation was done in the following

- NZ-wide was imputed using the most common genotype within clades
- Within-clade datasets were imputed using the most common genotype across all individuals

```
In [ ]: #Imputation
library(adegenet)
library(vcfR)

#-----NZ-all-----
#Read the data
readVcf file with input
vcr = read.vcfR("input.vcf")
vcr
#convert data to a genind object
gen = vcfR2genind(vcr)
gen = genInd2df(gen)
dim(gen)
sum(is.na(gen))
rownames(gen)
ak = 1 - 9
#AN - 10 - 17
#H - 18 - 26
#BD - 27 - 36
#CA - 37 - 46
#CH - 47 - 55
#CP - 56 - 65
#CR - 66 - 74
#CU - 72 - 76
#DU - 77 - 96
#EC - 97 - 105
#NC - 106 - 115
#RU - 116 - 129
#ST - 130 - 137
#TI - 138 - 149
#WE - 150 - 167

#Proportion of missing data per population (NA in population/total NA)
sum(is.na(gen[1:9,]))/sum(is.na(gen)) # AK 0.0727051
sum(is.na(gen[10:17,]))/sum(is.na(gen)) # AN 0.08951131
sum(is.na(gen[18:26,]))/sum(is.na(gen)) # AU 0.0360852
sum(is.na(gen[27:36,]))/sum(is.na(gen)) # BD 0.04648177
sum(is.na(gen[37:46,]))/sum(is.na(gen)) # CA 0.04503713
sum(is.na(gen[47:55,]))/sum(is.na(gen)) # CH 0.0548761
sum(is.na(gen[56:65,]))/sum(is.na(gen)) # CP 0.06970757
sum(is.na(gen[66:74,]))/sum(is.na(gen)) # CR 0.07930848
sum(is.na(gen[72:76,]))/sum(is.na(gen)) # CU 0.0358696
sum(is.na(gen[77:96,]))/sum(is.na(gen)) # DU 0.06858979
sum(is.na(gen[97:105,]))/sum(is.na(gen)) # EC 0.06887273
sum(is.na(gen[106:115,]))/sum(is.na(gen)) # NC 0.04338956
sum(is.na(gen[116:129,]))/sum(is.na(gen)) # RU 0.07376897
sum(is.na(gen[130:137,]))/sum(is.na(gen)) # ST 0.05992728
sum(is.na(gen[138:149,]))/sum(is.na(gen)) # TI 0.03850882
sum(is.na(gen[150:167,]))/sum(is.na(gen)) # WE 0.0851832

#having a look at AK population first
ak = gen[1:9,]
is.na(ak[,1])
aki =
for (i in 1:ncol(ak)) {
  replace (ak[,i], is.na(ak[,i]), as.numeric(names(which.max(table(ak[,i])))))
} #Error in x[!list] <- values : replacement has length zero
aki = apply(ak, 2, function(x) replace (x, is.na(x), as.numeric(names(which.max(table(x))))) #Error in x
[!list] <- values : replacement has length zero

#This error means that within the AK population, there are columns (SNPs) that are completely missing, so
it is impossible to find a value to input them (there is no replacement, or the replacement has zero leng
th).
#Thus, instead of finding by population for the RDA, I will try to input by clade (N, C, S)
#Because the three inputed clades
#The above was done for the second submiss
ion
RDA requires no missing data. Input missing values using the most common genotype at each SNP across all
individuals
#Read the data
#Read Vcf file with vcfR
vcr = read.vcfR("input.vcf")
vcr
#convert data to a genind object
gen = vcfR2genind(vcr)
gen = genInd2df(gen)
dim(gen)
sum(is.na(gen))
sum(is.na(gen[, 2, function(x) replace(x, is.na(x), as.numeric(names(which.max(table(x)))))])
sum(is.na(gen)) # No NAs

#Read environmental data in
env_all <- read.csv("../envdata_all.csv")
str(env_all)
env_all$ind <- as.character(env_all$ind) # Make individual names characters (not factors)
env_all$pop <- as.character(env_all$pop)
env_all$group <- as.character(env_all$group)

#Confirm genotypes and environment data are in the same order
identical(rownames(gen.imp), env_all[,1])
#Re-order to match
rownames(gen.imp)
gen.imp = gen.imp[order(rownames(gen.imp)),]
rownames(gen.imp)
identical(rownames(gen.imp), env_all[,1])
class(gen.imp)
gen.imp = as.data.frame(gen.imp)
str(gen.imp)
write.table(gen.imp, "gen.imp.csv", sep = ",") #this file will also be used to run RDA with the 100 random
variables

#-----Within-clade datasets-----
#This was done previous to revision, for the first submission. The above was done for the second submiss
ion
RDA requires no missing data. Input missing values using the most common genotype at each SNP across all
individuals
#Read the data
#Read Vcf file with vcfR
vcr = read.vcfR("input.vcf")
vcr
#convert data to a genind object
gen = vcfR2genind(vcr)
gen = genInd2df(gen)
dim(gen)
sum(is.na(gen))
sum(is.na(gen[, 2, function(x) replace(x, is.na(x), as.numeric(names(which.max(table(x)))))])
sum(is.na(gen)) # No NAs

#Read environmental data in
env_all <- read.csv("../envdata_all.csv")
str(env_all)
env_all$ind <- as.character(env_all$ind) # Make individual names characters (not factors)
env_all$pop <- as.character(env_all$pop)
env_all$group <- as.character(env_all$group)

#Confirm genotypes and environment data are in the same order
identical(rownames(gen.imp), env_all[,1])
#Re-order to match
rownames(gen.imp)
gen.imp = gen.imp[order(rownames(gen.imp)),]
rownames(gen.imp)
identical(rownames(gen.imp), env_all[,1])
class(gen.imp)
gen.imp = as.data.frame(gen.imp)
str(gen.imp)
write.table(gen.imp, "gen.imp.csv", sep = ",") #this file will also be used to run RDA with the 100 random
variables


```

RDA with real environmental variables


```
In [ ] : # Load packages
# -----
library(psychn) # Used to investigate correlations among predictors
library(vegan) # Used to run RDA
library(aegeanet)
library(vcfR)
library(usdm) # calculates variance inflation factors
library(robuust) # for robust mahalanobis distance calculation
library(qvalue) # FDR calculation

##Read datasets
##Environmental dataset
env.all <- read.csv("../envdata.all.csv")
str(env.all)
env.all$ind <- as.character(env.all$ind) # Make individual names characters (not factors)
env.all$pop <- as.character(env.all$pop)
env.all$group <- as.character(env.all$group)
##Genomic dataset (imputed)
gen.imp <- read.csv("gen.imp.csv", header = TRUE)
# Real predictors are often in different units (e.g. temp., forest cover, soil pH)
# In this case it is best to scale the predictors (which we'll do here as well):
env <- scale(env.all[,4:18]) # default is center and scale = T; this subtracts the mean from each value, then
env divides by the std. dev.
env <- as.data.frame(env)
env <- cbind.data.frame(env.all$ind, env.all$pop, env.all$group, env)
# "ind", "pop", "group", "lat", "lon", "mean_rslp", "mean_rswrs", "mean_otpem", "mean_pratre",
"mean_pres", "mean_calt", "mean_osalt", "mean_sshg", "mean_tcdc", "mean_uwnd", "mean_wvnd", "me
head(env)

##SUBSET THE DATA TO INCLUDE NON-CORRELATED ENV VARIABLES
pred <- subset(env, select=c(ind, pop, group, lat, lon, mean_pratre, mean_pres, mean_sshg))
rownames(pred) = c("lat", "lon", "mean_pratre", "mean_pres", "mean_sshg")
#RUN RDA
rda <- rda1 = rda(formula = gen.imp ~ lat + lon + mean_pratre + mean_pres + mean_sshg, data = pred, scale =
T)
#RDA SUMMARY OF RESULTS.
chiton.rda
plot(chiton.rda) #histogram of the loadings of the RDA axes
load.rda <- scores(chiton.rda, choices=c(1:5), display="species") # Species scores for the constrained ax
s
write.table(load.rda, "rda_loading.csv", sep = ",")
##TABLE OF PROPORTION EXPLAINED BY RDA AND PCA
summary(eigenvals(chiton.rda, model = "constrained")) #Proportion explained by the constrained axes (RDA
(canonical axis), or by the env variables)
summary(eigenvals(chiton.rda, model = "unconstrained")) #Proportion explained by the unconstrained axes (P
CA, or the residual variance left from the RDA)
#PLOT RDA
plot(chiton.rda)

##EXTRACT DATA TO CUSTOMIZE PLOTS
rda1.triplot <- rda.biplot
#Extract the data to be plotted from the RDA object
smry <- summary(chiton.rda)
df1 <- data.frame(smry$sites[,1:5], pred$pop) #get RDA1 and RDA2 for all populations
df2 <- data.frame(smry$species[,1:5]) #get the loadings of RDA1 and RDA2 for each locus
df3 <- data.frame(smry$biplot)
rownames(df3)
colnames(df3) <- c("lat", "lon", "mean_pratre", "mean_pres", "mean_sshg")
write.table(df1, "df1_plotRDA_ind.csv", sep = ",")
write.table(df2, "df2_plotRDA_loci.csv", sep = ",")
write.table(df3, "df3_plotRDA_env.csv", sep = ",")

##CUSTOMIZED RDA PLOTS
library(ggplot2)
#Set up colours
pred.int <- factor(c(env$pop, levels = c("R0", "T1", "N0", "EC", "ME", "CP", "CR", "DU", "AK", "CU", "ST",
"AU", "CA", "AN", "BO", "CH")),
#Create a custom color scale (called cols_pop_scale)
cols.pop <- c("ccc9000", "#ff9900", "#990066", "#ff33cc", "#ff9900", "#ffcc99", "#003399", "#006699", "#3399
33", "#33ff33", "#99cc66", "#3399ff", "#33ccff", "#99cccc", "#669999", "#000099")
names(cols.pop) <- levels(pred.int)
cols.pop_scale <- scale_colour_manual(name = "poplist.int", values = cols.pop)
#make a base plot with the individuals
rda.biplot <- ggplot(df3, aes(x=0, yend=RDA1, y=RDA2, colour = poplist.int)) +
geom_point(size = 4) +
geom_hline(yintercept=0, linetype="dotted") +
geom_vline(xintercept=0, linetype="dotted") +
coord_fixed()
cols.pop_scale +
theme_classic()
rda.biplot
#Add the SNPs in the RDA space
rda.biplot <- rda.biplot + geom_point(data = df2, aes (x =RDA1, y =RDA2), colour = "grey", size = 1, shape=3)
#Add the arrows indicating the direction of each environmental variable in the RDA space
rda.triplot <- rda.biplot +
geom_segment(data=df3, aes(x=0, yend=RDA1, y=RDA2),
color="black", arrow=arrow(length=unit(0.01,"npc")), arrow.fill = NULL, lineend = "butt", l
geom_text(data=df3,
aes(x=RDA1, y=RDA2, label=rownames(df3),
hjust=-0.5*(1-sign(RDA1)), vjust=0.5*(1-sign(RDA2))),
color="black", size=5)
rda.triplot +
theme(text = element_text(size=20))
#Have the arrows without labels
no.arrow.labels <- rda.biplot +
geom_segment(data=df3, aes(x=0, yend=RDA1, y=RDA2),
color="black", size = 1, arrow=arrow(length=unit(0.02,"npc")), arrow.fill = NULL, lineend =
"butt", linejoin = "round") +
theme(text = element_text(size=20))
no.arrow.labels
##READING OUTLIER LOCI (treating loadings as z-scores and converting them to p-values and q-values)
#Read the RDA loadings in
load.rda <- read.csv("rda_loading.csv", header = T)
head(load.rda)
class(load.rda)
load.rda <- as.data.frame(load.rda)
load.rda$locus <- rownames(load.rda)
colnames(load.rda)
head(rownames(load.rda))
dim(load.rda)
#group loadings from all axes together
library(tidy)
load.rda <- gather(load.rda, axis, loading, -locus) #same as using reshape2 function melt(): mload.rda =
melt(load.rda, id = "locus")
tail(mload.rda)
#mean-center and scale loadings to variance = 1
mload.rda$scaled_loading <- scale(mload.rda$loading)
head(mload.rda)
var(mload.rda$scaled_loading)
hist(mload.rda$scaled_loading)
hist(mload.rda$scaled_loading)
#transform loadings into p-values
mload.rda$pvalues <- (1-pnorm(abs(mload.rda$scaled_loading)))^2
head(mload.rda)
hist(mload.rda$pvalues)
#transform p-values into q-values
library(qvalue)
qvalue.rda <- qvalue(mload.rda$pvalues)
summary(qvalue.rda)
hist(qvalue.rda$pvalues)
mload.rda$qvalues <- qvalue.rda$qvalues
head(mload.rda)
summary(mload.rda$qvalues)
write.table(mload.rda, "SNPs_loadings_pvalues_qvalues.csv", sep = ",")
#find outliers based on q-values (cut off = 0.05)
outlier.loci <- function(x, y) { # x is a data frame with the q-values in a row named qvalues; y is the qv
#size cutoff for outliers
outs <- x$qvalues <= y #find loci with qvalue smaller than cutoff
outliers <- x[outs, ] #subset the dataframe to include only outliers
if (nrow(outliers) > 0) { #if there are outliers
outliers_unique <- outliers[duplicated(outliers$locus),] # include only unique outliers in the dataframe
print(dim(outliers_unique)) #print outliers to the console
print(head(outliers_unique))
hist(outliers_unique$qvalues, breaks = 10)#plot the distribution of the q-values of outliers
write.table(outliers_unique, "outliers_qvalue_rda.csv", sep = ",")# write a table to the working dire
tory with the defined outliers (not duplicated)
} else { print ("No outlier locus was found with the selected q-value cutoff")} #if there are no outliers
}
outlier.loci(mload.rda, 0.05)
#No outlier locus was found with the selected q-value cutoff"
##HOW TO RUN THE CODE BELOW IF OUTLIERS WERE IDENTIFIED ABOVE
#Add the correlation of each SNP with the env variables
outliers.unique <- read.csv("outliers_qvalue_rda.csv")
#FOO = matrix (nrow = nrow(outliers.unique), ncol = 5)
colnames(foo) <- c("A", "B", "C", "D", "E")
pred.1 <- pred[,4:8]
for (i in 1:length(outliers.unique$locus)) {
nam <- outliers.unique[,i,2] #get the nameID of the outlier loci
snp.gen <- gen.imp[,nam] #get its genotypes from the gen.imp dataset (genetic data imputed used in RDA)
foo[,i] <- apply(pred.1,2,function(x) cor(x,snp.gen)) #apply a correlation function to the genetic and t
he env data for this loci
}
full <- cbind.data.frame(outliers.unique, foo)
head(full)
#see which of the predictors each candidate SNP is most strongly correlated with
for (z in 1:length(full$locus)) {
bar <- full[,z]
full[,12] <- names(which.max(abs(bar[,7:11]))) # gives the variable
full[,13] <- max(abs(bar[,7:11])) # gives the correlation
}
colnames(full)[12] <- "predictor"
colnames(full)[13] <- "correlation"
head(full)
write.table(full, "outliers_qvalue_correlation_fullTable.csv", sep = ",")
#The code below tests the RDA model for significance (can be run when it did not identify outliers as wel
#TESTS OF SIGNIFICANCE FOR THE RDA MODEL AND MULTICOLLINEARITY AMONG ENV VARIABLES USED (run this last, it
takes time)
#TEST THE RDA FULL MODEL FOR SIGNIFICANCE USING F-STATISTICS
signif.full <- anova.cca(chiton.rda, parallel=getOption("mc.cores")) # default is permutation=999
signif.full
#TEST EACH CONSTRAINED AXIS (RDAs) FOR SIGNIFICANCE AS DONE WITH THE FULL MODEL, USING F-STATISTICS (THIS
IS THE ONE THAT TAKES TIME)
signif.axis <- anova.cca(chiton.rda, by="axis", parallel=getOption("mc.cores")) #will take a couple of hou
rs
signif.axis
RDA with 100 random variables
```

In hindsight, I do not think it makes any difference (maybe in the speed of the runs), but I broke the 100 variables into 20 RDA runs with 5 variables each. 20 folders for the output were needed. What was needed:

- the gen.imp.csv file (imputed file created from vcf, same used in the RDA run with real environmental variables)
- the csv file created with all random variables (same random variables as used in BayesEnv)
- the R script below, which was run in the background with nohup in Nesi using:

```
In [ ] : nohup Rscript script.R &
```

The output of the nohup run (nohup.out) will have the results of the RDA run. If the function to identify outliers return any outlier with a q-value cut-off 0.05, it will be printed in this outfile.

```
In [ ] : nano nohup.out
```

```
In [ ] : #script.R to run the RDA
# Load packages
# -----
library(psychn) # Used to investigate correlations among predictors
library(vegan) # Used to run RDA
library(aegeanet)
library(vcfR)
library(usdm) # calculates variance inflation factors
library(robuust) # for robust mahalanobis distance calculation
library(qvalue) # FDR calculation

##Read genetic data (imputed from a previous run of the code)
gen.imp <- read.csv("gen.imp.csv", header = TRUE)
##Read environmental data
env <- read.csv("../all_arbitrary_env.csv")
colnames(env) <- as.character(colnames(env))
#mean-center and scale environmental data
senv <- scale(env[,1:183])
env <- data.frame(env[,1:183], senv)

#Subset the env file with 100 environmental variables - this can be used for RDA without formula, but I am
just keeping because all the code was developed with pred (after the RDA with formula)
pred <- subset(env, select=c(ind, pop, group, lat, lon, X1, X2, X3, X4, X5)) #CHANGE THE LINE FOR EACH OF THE 20 RUN
S (run 2 should use the line below, run 3 the line below that, and so on)
#pred <- subset(env, select=c(ind, pop, group, X6, X7, X8, X9, X10))
#pred <- subset(env, select=c(ind, pop, group, X11, X12, X13, X14, X15))
#pred <- subset(env, select=c(ind, pop, group, X16, X17, X18, X19, X20))
#pred <- subset(env, select=c(ind, pop, group, X21, X22, X23, X24, X25))
#pred <- subset(env, select=c(ind, pop, group, X26, X27, X28, X29, X30))
#pred <- subset(env, select=c(ind, pop, group, X31, X32, X33, X34, X35))
#pred <- subset(env, select=c(ind, pop, group, X36, X37, X38, X39, X40))
#pred <- subset(env, select=c(ind, pop, group, X41, X42, X43, X44, X45))
#pred <- subset(env, select=c(ind, pop, group, X46, X47, X48, X49, X50))
#pred <- subset(env, select=c(ind, pop, group, X51, X52, X53, X54, X55))
#pred <- subset(env, select=c(ind, pop, group, X56, X57, X58, X59, X60))
#pred <- subset(env, select=c(ind, pop, group, X61, X62, X63, X64, X65))
#pred <- subset(env, select=c(ind, pop, group, X66, X67, X68, X69, X70))
#pred <- subset(env, select=c(ind, pop, group, X71, X72, X73, X74, X75))
#pred <- subset(env, select=c(ind, pop, group, X76, X77, X78, X79, X80))
#pred <- subset(env, select=c(ind, pop, group, X81, X82, X83, X84, X85))
#pred <- subset(env, select=c(ind, pop, group, X86, X87, X88, X89, X90))
#pred <- subset(env, select=c(ind, pop, group, X91, X92, X93, X94, X95))
#pred <- subset(env, select=c(ind, pop, group, X96, X97, X98, X99, X100))

#RDA with formula (allows calculating anova.cca per axis later on)
chiton.rda <- rda(formula = gen.imp ~ X1 + X2 + X3 + X4 + X5, data = env, scale = T) #CHANGE THE LINE FOR E
ACH OF THE 20 RUNS (run 2 should use the line below, run 3 the line below that, and so on)
#chiton.rda <- rda(formula = gen.imp ~ X6 + X7 + X8 + X9 + X10, data = env, scale = T)
#chiton.rda <- rda(formula = gen.imp ~ X11 + X12 + X13 + X14 + X15, data = env, scale = T)
#chiton.rda <- rda(formula = gen.imp ~ X16 + X17 + X18 + X19 + X20, data = env, scale = T)
#chiton.rda <- rda(formula = gen.imp ~ X21 + X22 + X23 + X24 + X25, data = env, scale = T)
#chiton.rda <- rda(formula = gen.imp ~ X26 + X27 + X28 + X29 + X30, data = env, scale = T)
#chiton.rda <- rda(formula = gen.imp ~ X31 + X32 + X33 + X34 + X35, data = env, scale = T)
#chiton.rda <- rda(formula = gen.imp ~ X36 + X37 + X38 + X39 + X40, data = env, scale = T)
#chiton.rda <- rda(formula = gen.imp ~ X41 + X42 + X43 + X44 + X45, data = env, scale = T)
#chiton.rda <- rda(formula = gen.imp ~ X46 + X47 + X48 + X49 + X50, data = env, scale = T)
#chiton.rda <- rda(formula = gen.imp ~ X51 + X52 + X53 + X54 + X55, data = env, scale = T)
#chiton.rda <- rda(formula = gen.imp ~ X56 + X57 + X58 + X59 + X60, data = env, scale = T)
#chiton.rda <- rda(formula = gen.imp ~ X61 + X62 + X63 + X64 + X65, data = env, scale = T)
#chiton.rda <- rda(formula = gen.imp ~ X66 + X67 + X68 + X69 + X70, data = env, scale = T)
#chiton.rda <- rda(formula = gen.imp ~ X71 + X72 + X73 + X74 + X75, data = env, scale = T)
#chiton.rda <- rda(formula = gen.imp ~ X76 + X77 + X78 + X79 + X80, data = env, scale = T)
#chiton.rda <- rda(formula = gen.imp ~ X81 + X82 + X83 + X84 + X85, data = env, scale = T)
#chiton.rda <- rda(formula = gen.imp ~ X86 + X87 + X88 + X89 + X90, data = env, scale = T)
#chiton.rda <- rda(formula = gen.imp ~ X91 + X92 + X93 + X94 + X95, data = env, scale = T)
#chiton.rda <- rda(formula = gen.imp ~ X96 + X97 + X98 + X99 + X100, data = env, scale = T)

print(chiton.rda)
sum.constrained <- print(summary(eigenvals(chiton.rda, model = "constrained"))) #Proportion explained by th
e constrained axes (RDA (canonical axis), or by the env variables)
sum.unconstrained <- print(summary(eigenvals(chiton.rda, model = "unconstrained"))) #Proportion explained b
y the unconstrained axes (PCA, or the residual variance left from the RDA)
hists <- sapply(chiton.rda) #HISTOGRAM OF THE EIGENVALUES OF THE RDAs
#make a plot with the loadings for each RDA and each SNP
load.rda <- scores(chiton.rda, choices=c(1:5), display="species") # Species scores for the constrained ax
es
write.table(load.rda, "rda_loading.csv", sep = ",")

#Read the RDA loadings in
load.rda
head(load.rda)
class(load.rda)
load.rda <- as.data.frame(load.rda)
load.rda$locus <- rownames(load.rda)
colnames(load.rda)
head(rownames(load.rda))
dim(load.rda)
#group loadings from all axes together
library(tidy)
load.rda <- gather(load.rda, axis, loading, -locus) #same as using reshape2 function melt(): mload.rda =
melt(load.rda, id = "locus")
tail(mload.rda)
#mean-center and scale loadings to variance = 1
mload.rda$scaled_loading <- scale(mload.rda$loading)
head(mload.rda)
var(mload.rda$scaled_loading)
hist(mload.rda$scaled_loading)
hist(mload.rda$scaled_loading)
#transform loadings into p-values
mload.rda$pvalues <- (1-pnorm(abs(mload.rda$scaled_loading)))^2
head(mload.rda)
hist(mload.rda$pvalues)
#transform p-values into q-values
library(qvalue)
qvalue.rda <- qvalue(mload.rda$pvalues)
summary(qvalue.rda)
hist(qvalue.rda$pvalues)
mload.rda$qvalues <- qvalue.rda$qvalues
head(mload.rda)
summary(mload.rda$qvalues)
write.table(mload.rda, "SNPs_loadings_pvalues_qvalues.csv", sep = ",")
#find outliers based on q-values (cut off = 0.05)
outlier.loci <- function(x, y) { # x is a data frame with the q-values in a row named qvalues; y is the qv
#size cutoff for outliers
outs <- x$qvalues <= y #find loci with qvalue smaller than cutoff
outliers <- x[outs, ] #subset the dataframe to include only outliers
if (nrow(outliers) > 0) { #if there are outliers
outliers_unique <- outliers[duplicated(outliers$locus),] # include only unique outliers in the dataframe
print(dim(outliers_unique)) #print outliers to the console
print(head(outliers_unique))
hist(outliers_unique$qvalues) #plot the distribution of the q-values of outliers
write.table(outliers_unique, "outliers_qvalue_rda.csv", sep = ",") # write a table to the working dire
ctory with the defined outliers (not duplicated)
} else { print ("No outlier locus was found with the env variables")
}
outlier.loci(mload.rda, 0.05)

#Add the correlation of each SNP with the env variables
outliers.unique <- read.csv("outliers_qvalue_rda.csv")
head(outliers.unique)
foo <- matrix (nrow = nrow(outliers.unique), ncol = 5)
colnames(foo) <- c("A", "B", "C", "D", "E")
pred.1 <- pred[,4:8]
for (i in 1:length(outliers.unique$locus)) {
nam <- outliers.unique[,i,2] #get the nameID of the outlier loci
snp.gen <- gen.imp[,nam] #get its genotypes from the gen.imp dataset (genetic data imputed used in RDA)
foo[,i] <- apply(pred.1,2,function(x) cor(x,snp.gen)) #apply a correlation function to the genetic and t
he env data for this loci
}
full <- cbind.data.frame(outliers.unique, foo)
head(full)

#see which of the predictors each candidate SNP is most strongly correlated with
for (z in 1:length(full$locus)) {
bar <- full[,z]
full[,12] <- names(which.max(abs(bar[,7:11]))) # gives the variable
full[,13] <- max(abs(bar[,7:11])) # gives the correlation
}
colnames(full)[12] <- "predictor"
colnames(full)[13] <- "correlation"
head(full)
write.table(full, "outliers_qvalue_correlation_fullTable.csv", sep = ",")
#TESTS OF SIGNIFICANCE FOR THE RDA MODEL AND MULTICOLLINEARITY AMONG ENV VARIABLES USED (run this
last, it takes time)
#TEST THE RDA FULL MODEL FOR SIGNIFICANCE USING F-STATISTICS
signif.full <- anova.cca(chiton.rda, parallel=getOption("mc.cores")) # default is permutation=999
signif.full
#TEST EACH CONSTRAINED AXIS (RDAs) FOR SIGNIFICANCE AS DONE WITH THE FULL MODEL, USING F-STATISTICS (THIS
IS THE ONE THAT TAKES TIME)
signif.axis <- anova.cca(chiton.rda, by="axis", parallel=getOption("mc.cores")) #will take a couple of ho
urs
signif.axis
```

For the NZ-wide dataset, the RDA results with Random variables showed two random variables correlated with the same six outliers (output from nohup, run the code below in R)

```
In [ ] : ##Running the new imputed dataset (imputed by clade) with the random environmental variables identified onl
y six outliers
4, 12, 19, 36, and 44.
#ALL of these outliers are homozygous ref to N and S (00 in the plot and dataset), and homozygous alt to C
(11 in the plot and dataset), for both not-imputed and imputed datasets
```

```
gen.imp <- read.csv("../gen.imp.csv", header = TRUE)

plot(gen.imp[,2])
plot(gen.imp[,4])
plot(gen.imp[,12])
plot(gen.imp[,19])
plot(gen.imp[,36])
plot(gen.imp[,44])

plot(gen[,2])
plot(gen[,4])
plot(gen[,12])
plot(gen[,19])
plot(gen[,36])
plot(gen[,44])

#-----Random Variables-----
##Read files with RDA q-values and Bayescenv q-values for the 100 random env variables
rda <- read.csv("all_random_qvalues2021.csv")
bay <- read.csv("summary_qval_allRealEnv.csv", header = T)

#working and tidying with the RDA data frame
head(rda)
library(tidyverse)
library(dplyr)
#select the columns of interest (locus ID (which is CHROM in this case), qvalue and RDA/env) and transform
rda into wide form
rda <- rda %>% select (1, 2, 6)
head(rda)
rda <- rda %>% spread (axis, qvalues)
#remove unnecessary rows (the q-values of CHROM) and re-order table to increasing loci id/number (as to be abl
e to use locus = 1:10987, matching rownumber and also bayescenv)
head(rda$locus)
library(stringr)
Cid <- rda$locus
Cid <- str_split_fixed(Cid, "-", n = 2)
Cid <- Cid[,1,3]
Cid <- str_split_fixed(Cid, "X", n = 2)
Cid <- Cid[,2]
Cid <- as.numeric(Cid)
Cid

#Add Cid (Chrom ID without X and position) as a column in the data frame rda
head(rda$CHROM)
Cid <- rda$CHROM
#re-order the data frame so CHROM is in ascending order
rda <- rda[order(rda$CHROM),]
head(rda$CHROM)
head(rda)
#re-write the 'locus' column with loci ID matching rownumber and bayescenv
rda$locus <- 1:10987 #change this range to the number of loci in the dataset
head(rda)

#working and tidying the Bayescenv dataframe (we are only interested in g)
head(bay)
#select columns with q-values for g in each real env variable
bay <- bay %>% select (1, 2, 4, 6, 8, 10, 12, 14, 16, 18, 20, 22, 24, 26, 28, 30)
#we want to combine the q-values of Bayescenv and RDA for each env variable, but RDA does not include vari
ables correlated with latitude (only latitude)
#so we will combine the q-value of bayescenv-latitude with the q-value of all the other variables correlat
ed with latitude
#when interpreting the results, this will mean that gene-env correlations that are true for latitude, coul
d be true for any of the latitude-correlated envs
#to combine the q-values of these, we will use the geometric mean ( the nth root of the product of n numbe
rs)
#The variables used in RDA and Bayescenv were Lat, Long, Prate, Pres, and Sshg. All others are correlated
with latitude

broad.lat <- (bay$lat0+bay$slp0+bay$snwrs0+bay$ssalt0+bay$stc0+bay$stemp0+bay$slalt0+bay$ssst0+bay$stc0+bay$wnd0+ba
y$wndm0+bay$airs0)/(1/11)

#now remake the bayescenv (bay) dataframe to contain the latitude "broad sense" just calculated, and the ot
her env variables also included in the RDA
bay$lat <- broad.lat
bay <- bay %>% select (1, 2, 3, 8, 9, 11)
head(bay)

#Combine the q-values of rda and bayescenv (don't forget that lat is now "broad sense")
#create a function to do the geometric mean as above
gmean <- function(x, y) {
sqrt(x*y)}
#Select the RDA columns to be used (leave loci ID out)
colnames(rda)
dim(rda)
rda <- rda %>% select (2:6)
#select the bayescenv columns to be used (leave loci ID out)
colnames(bay)
dim(bay)
bay <- bay %>% select (2:6)

#Create a variable from 1:5 to iterate from (the 5 columns, or env variables, of each dataset)
envs <- 1:5
#create a list to store the results of each env variable
result_list <- list()
#make a for loop to calculate the geometric mean
for (i in envs) {
n <- gmean (rda[,i], bay[,i])
result_list [[i]] <- n
}
summary(result_list) #check that the list stored the results
summary(result_list[[1]]) #check that there are values of geometric mean for env variable 1
#convert the list with all results to a data frame, containing a locus per row and an env variable per col
column
df <- data.frame(matrix(unlist(result_list), nrow=10987, byrow=F)) #change the nrow to the number of loci i
n the dataset
#locus <- 1:10987 had the column with locus id (which is just the row number because both datasets were
ordered) -- change the nrow to the number of loci in the dataset
head(df)
#save the table
write.table(df, "combined_qvalues_real_RDA_Bayescenv.csv", sep = ",", row.names = F)

#-----Random Variables-----
##Read files with RDA q-values and Bayescenv q-values for the 100 random env variables
rda <- read.csv("all_random_qvalues2021.csv")
bay <- read.csv("summary_qval_allRealEnv.csv", header = T)

head(rda)
head(bay)
tail(rda)
tail(bay)
#calculate the geometric mean of q-values for random env 1 (combined q-value). The geometric mean is the
nth root of the product of n numbers
gmean <- sqrt(rda$X1*bay$X1)
gmean1
#create a function to do the geometric mean as above
gmean <- function(x, y) {
sqrt(x*y)}
}
#test the function and compare with the result of gmean1
summary(rda$X1, bay$X1)
#now automate calculating the geometric mean between RDA and Bayescenv q-values of each random env
library(tidyverse)
#select the RDA dataframe to be used (leave loci ID out)
colnames(rda)
dim(rda)
rda <- rda %>% select (2:181)
#select the Bayescenv dataframe to be used (leave loci ID out)
colnames(bay)
dim(bay)
bay <- bay %>% select (2:181)
#Create a variable from 1:100 to iterate from (the 100 columns, or random env variables, of each dataset)
envs <- 1:100
#create a list to store the results of each random variable
result_list <- list()
#make a for loop to calculate the geometric mean
for (i in envs) {
n <- gmean (rda[,i], bay[,i])
result_list [[i]] <- n
}
summary(result_list) #check that the list stored the results
summary(result_list[[1]]) #check that there are values of geometric mean for random variable 1
#convert the list with all results to a data frame, containing a locus per row and a random env variable p
er column
df <- data.frame(matrix(unlist(result_list), nrow=10987, byrow=F)) #change the nrow to the number of loci i
n the dataset
#locus <- 1:10987 had the column with locus id (which is just the row number because both datasets were
ordered) -- change the nrow to the number of loci in the dataset
head(df)
#save the table
write.table(df, "combined_qvalues_random_RDA_Bayescenv.csv", sep = ",", row.names = F)

#How many 5% random associations, combined qvalue?
df.bin <- df[,1:180]
df.bin <- ifelse(df.bin <= 0.05, 1, 0)
df.bin <- as.data.frame(df.bin)
head(df.bin)
df <- df %>% bintotal = apply(df.bin, 1, sum)
head(df.bin)
length(which(df.bin$total>=5)) #RDA-Bayescenv combined 5% associations
summary(df.bin)

##Looking into RDA and Bayescenv results separately
rda.bin <- ifelse (rda <= 0.05, 1, 0)
head(rda.bin)
rda.bin <- as.data.frame(rda.bin)
rda.bin$total <- apply (rda.bin, 1, sum)
head(rda.bin)
length(which(rda.bin$total>=5)) #RDA 5% random associations
summary(rda.bin)

bay.bin <- ifelse(bay <= 0.05, 1, 0)
bay.bin <- as.data.frame(bay.bin)
head(bay.bin)
bay.bin$total <- apply(bay.bin, 1, sum)
head(bay.bin)
length(which(bay.bin$total>= 5)) #722 Bayescenv 5% random associations
summary(bay.bin)
```

Checking back with the runs with real environmental variables to see if any association was found after c

ombining q-values

```
head(real)
real <- read.csv("../real/combined_qvalues_real_RDA_Bayescenv.csv")
real <- read[,1:5]
real.bin <- ifelse(real <= 0.05, 1, 0)
real.bin <- as.data.frame(real.bin)
head(real.bin)
real.bin$total <- apply(real.bin, 1, sum)
length(which(real.bin$total==1))
summary(real.bin)
```

Checking gene-environment associations with real environmental variables and random variables (loci associated with at least 5% of the random variables will no longer be considered true associations)

What is needed

- The file with the combined q-value of BayeScEnv and RDA for the real environmental variables (code for combining q-values provided above)
- The file with the combined q-values of BayeScEnv and RDA for the random variables (code for combining q-values provided above)
- A csv file with ID info from the vcf file (loci_info_vcfStyle.csv), for matching the locus ID of the loci associated with real environmental variables with the locus ID of the vcf file.


```
In [ ]: #Southern (this example is provided with the Southern dataset, since the NZ-wide did not have any associations in the RDA)
library(tidyverse)

#read file with combined q-values of Bayescenv and RDA for real and random variables
real = read.csv("real/combined_qvalues_real_RDA_Bayescenv.csv")
random = read.csv("arbitrary/Random.combinedqVals.RDABayescenv.csv")

head(real)
head(random)

summary(real)
summary(random)

#Define which loci are outliers (qvalue < 0.05) in real and random datasets using T and F
envs = 5 #number of environmental variables
loci = nrow(real) #number of loci

real_out = matrix(ncol=envs, nrow=loci) #create a matrix to store the results of a for loop

for (col in 1:envs) { #to each column (each environmental variable)
  random_out[,col] = ifelse (real[,col] < 0.05, TRUE, FALSE) #call qvalues < 0.05 of TRUE and larger of FALSE
}

head(real_out)
class(real_out)
real_out = as.data.frame(real_out) #Transform the matrix of results into a dataframe
real_out$loci = 1:loci #add the locus ID into the results dataframe
colnames(real_out) = colnames(real) #rename the columns using the env variable names
summary(real_out) # There are 111 outliers in LAT*, 1 in LONG, 5 in PRATE, 135 IN PRESS and 23 IN SSHg.

#Find the total number of outliers and which ones they are for vcf_style table
total_random_outlier = ifelse(random_out[,1:100]==TRUE, TRUE, FALSE)
ifelse(real_out$long==TRUE , TRUE,
       ifelse(real_out$press==TRUE, TRUE,
              ifelse(real_out$sshg==TRUE, TRUE, FALSE))))

real_out$total_real = total_real
head(real_out)
length(real_out$total_real)
write.table (real_out, "real/combinedqVals.vcfStyleSummary.csv", sep = ",", row.names = FALSE)

#Random variables
envs = 100 #number of random variables
loci = nrow(random) #number of loci

random_out = matrix(ncol=envs, nrow=loci) #create a matrix to store the results of a for loop

for (col in 1:envs) { #to each column (each environmental variable)
  random_out[,col] = ifelse (random[,col] < 0.05, TRUE, FALSE) #call qvalues < 0.05 of TRUE and larger of FALSE
}

head(random_out)
class(random_out)
random_out = as.data.frame(random_out) #transform the matrix of results into a dataframe
random_out$loci = 1:loci #add the locus ID into the results dataframe
summary(random_out) #There are some associations between random env variables and some of the loci when th
e q-values of Bayescenv and RDA are combined

write.table (random_out, "Arbitrary/combedQvals.vcfStyleSummary.csv", sep = ",", row.names = FALSE)

#Find the total number of random outliers and which ones they are for vcf_style table (associated with onl
y one random variable)
random_out = ifelse (random_out[,1:100]==TRUE, TRUE, 0)
total_random = apply(random_out[,1,sum])
class(total_random)
length(total_random)
class(random_out)
as.data.frame(ifelse(total_random==0, FALSE, TRUE))
total_random = as.data.frame(total_random)
random_out$total_random = total_random
dim(random_out)
length(total_random)
head(random_out)
summary(random_out) #2088 random associations (with at least one loci, not the 5K)

#High confidence outliers (associated only with real envs and less than 5 random envs)
# Transform random dataset from T anf F to 0 and 1
#random
random_out = ifelse (random_out[,1:100]==TRUE, 1, 0)
sum(random_out[,1:100]) #1193 outliers in total (regardless of random variable - that means one locus can
be associated with more than one environmental variable)
random_out = as.data.frame(random_out)
head(random_out)
#sum to find out which loci have sum > 5 (outliers in at least 5 out of the 100 runs with random variabl
es)
random_out$sum = apply(random_out[, 1, sum])
summary(random_out)
random_out$loci = 1:nrow(random_out)#add loci ID
length(which loci were associated with at least 5 random variables (TRUE if sum column is equal to or lar
ger than 5))
random_out$out = ifelse (random_out$sum == 5, TRUE, FALSE) ####THIS IS THE LINE TO COPY TO THE VCF STYLE SUM
MARY
head(random_out)
summary(random_out) #268 of the loci were associated with at least five random environmental variables (th
e 5K)
write.table (random_out, "Arbitrary/combedQvals.vcfStyleSummary.csv", sep = ",", row.names = FALSE)

#Which of the real variables is high confidence (not associated with five or more random env variab
les?)
head(real_out)
class(real_out)
head(random_out)
class(random_out)
library(tidyverse)
%>%
  random_out$loci = random_out %>% filter (out == TRUE) %>% select (loci)
random_out$loci = random_out$loci[,1] #random_out$loci is now a vector including loci id of loci associa
te d with at least 5 random env variables

#Latitude* (broad sense)
lat = real_out %>% filter (lat == TRUE) %>% select (loci)
lat = lat[,1:]#lat is now a vector containing loci ID of loci associated with latitude
high_confidence_lat = lat %>% filter (high_confidence_lat==TRUE) %>% select (lat) #the 81 loci are only
associated with latitude broad sense (no duplication with other real env variables)
high_confidence_prate = as.vector(high_confidence_lat$lat)

#Longitude
long = real_out %>% filter (long == TRUE) %>% select (loci)
long = long[,1:]#long is now a vector containing loci ID of loci associated with longitude
high_confidence_long = long %>% filter (high_confidence_long==TRUE) %>% select (loci) #the 81 loci are only
associated with longitude broad sense (no duplication with other real env variables)
high_confidence_sshg = as.vector(high_confidence_long$long)

#prate
prate = real_out %>% filter (prate == TRUE) %>% select (loci)
prate = prate[,1:]#prate is now a vector containing loci ID of loci associated with latitude
high_confidence_press = press %>% filter (high_confidence_press==TRUE) %>% select (press) #Locus ID 322
is associated with latitude broad sense (no duplication with other real env variables)
df_prate = data.frame(prate, high_confidence_press)
high_confidence_prate = df_prate %>% filter (high_confidence_press==TRUE) %>% select (prate) #Locus ID 322
is associated with latitude broad sense (no duplication with other real env variables)
high_confidence_press = as.vector(high_confidence_press$prate)

#ssshg
ssshg = real_out %>% filter (ssshg == TRUE) %>% select (loci)
ssshg = ssshg[,1:]#ssshg is now a vector containing loci ID of loci associated with latitude
high_confidence_ssshg = ssshg %>% filter (high_confidence_ssshg==TRUE) %>% select (ssshg) #Locus ID 5134
- same as long
high_confidence_ssshg = as.vector(high_confidence_ssshg$ssshg)

#make a vector of T and F to add to the vcf_style file
lat_high = ifelse (real_out$loci %in% high_confidence_lat, TRUE, FALSE)
summary(lat_high)
prate_high = ifelse (real_out$loci %in% high_confidence_prate, TRUE, FALSE)
press_high = ifelse (real_out$loci %in% high_confidence_press, TRUE, FALSE)
ssshg_high = ifelse (real_out$loci %in% high_confidence_ssshg, TRUE, FALSE)
#combine the vectors above in a table
highConfidenceSummaryTable = data.frame(lat_high, long_high, prate_high, press_high, ssshg_high, real_out$loci)
head(highConfidenceSummaryTable)
#convert the vectors above in a table
highConfidenceSummaryTable = c('lat_high','long_high','prate_high','press_high','ssshg_high','lo
ci')

#make a vector with the combined result (total high confidence associations)
total_high = ifelse (highConfidenceSummaryTable$lat_high== TRUE, TRUE,
                     ifelse (highConfidenceSummaryTable$prate_high == TRUE, TRUE,
                             ifelse (highConfidenceSummaryTable$press_high== TRUE, TRUE,
                                     ifelse (highConfidenceSummaryTable$ssshg_high==TRUE, TRUE, FALSE
                                             )))))
highConfidenceSummaryTable$total_high = total_high
summary(highConfidenceSummaryTable)
write.table (highConfidenceSummaryTable, "highconf_geneEnvAssociations.csv", sep = ",", row.names = FALSE)

#Which are the high-confidence loci id?
lat_id = highConfidenceSummaryTable %>% filter (lat_high==TRUE) %>% select (loci)
lat_id = lat_id$loci
lat_id
#
prate_id = highConfidenceSummaryTable %>% filter (prate_high==TRUE) %>% select (loci)
prate_id = prate_id$loci
prate_id
#
press_id = highConfidenceSummaryTable %>% filter (press_high==TRUE) %>% select (loci)
press_id = press_id$loci
press_id
#
ssshg_id = highConfidenceSummaryTable %>% filter (ssshg_high==TRUE) %>% select (loci)
ssshg_id = ssshg_id$loci
ssshg_id

#Check for duplicates -> are the high-confidence loci associated with more than one environmental variabl
e?
head(highConfidenceSummaryTable)
lat_high and prate_high
highConfidenceSummaryTable %>% filter (lat_high == TRUE & prate_high == TRUE) %>% select (loci) #0
#lat_high and press_high
highConfidenceSummaryTable %>% filter (lat_high == TRUE & press_high == TRUE) %>% select (loci) #0
#lat_high and ssshg_high
highConfidenceSummaryTable %>% filter (lat_high == TRUE & ssshg_high == TRUE) %>% select (loci) #0
#prate_high and press_high
highConfidenceSummaryTable %>% filter (prate_high == TRUE & press_high == TRUE) %>% select (loci) #0
#prate_high and ssshg_high
highConfidenceSummaryTable %>% filter (prate_high == TRUE & ssshg_high == TRUE) %>% select (loci) #0
#press_high and ssshg_high
highConfidenceSummaryTable %>% filter (press_high == TRUE & ssshg_high == TRUE) %>% select (loci) #0

#Find CHROM id of high-confidence associations
chromid = read.csv(".././CHROM_LOCI_ID_Sonly.csv", header = T)
head(chromid)
lat_chrom = chromid[lat_id, 1]
prate_chrom = chromid[prate_id, 1]
press_chrom = chromid[press_id, 1]
ssshg_chrom = chromid[ssshg_id, 1]

#Make a table with all high-confidence loci IDs to use for grepping the fasta sequence and blast it agains
t Acanthopleura granulata genome
headf = cbind(highConfidenceSummaryTable, chromid)
hic = hiccdf %>% filter (total_high == TRUE) %>% select (CHROM)
head(hic)
chrom = hic$CHROM

for (i in 1:200) { #change 200 by nrow(chrom), that is, the number of CHROM ids to search the fasta sequenc
ce for
  text = paste ("grep -A 1 -w CLOCUS:", chrom[i], " populations.loci.fa > loci", chrom[i], ".fasta", sep
    )
  print(text)
} #copy and paste output in notepad, than save as .sh to run in bash
```

Extract fasta sequences from loci associated with the environment (high-confidence) and Blast against the Acanthopleura granulata genome (custom Blast database)

What is needed

- the .sh file from the R code above, which contain grep commands to find the CHROM id and fasta sequence using bash
- the fasta output from the populations module on Stacks (it is unfiltered, but it is fine, since it will be used to search for specific loci IDs)
- Acanthopleura granulata assembly and transcripts (downloaded fasta from Dryad
<https://doi.org/10.5061/dryad.wstg9q2k9>)
- Blast installed locally for load the module in NeSI)

```
In [ ]: cd Stacks/ #Go to the Stacks directory that contains the fasta output from the populations module
run seqprep.sh #Run the file with the fasta sequences to find, which was generated with the R script above

cd Agranulata/ #Go to the directory with A. granulata .fasta assembly and transcripts
cp ../Sonly_outlocus_RDA_Bayescnv_HighConfidence.fasta ./copy_the_fasta_file_with_the_target_loci_sequences_output_from_previous_step
load Blast module
Module spider/load Blast
Run Blast (can also be done in a slurm script) with A. granulata scaffolds (example below) and transcript
ones (output detected of blastn)
blast query Sonly_outlocus_RDA_Bayescnv_HighConfidence.fasta -subject Acanthopleura_granulata-scaffold
-f.sas -outfast '7 qseqid sseqid length qlen slen qstart qend start send evalua' -out SonlyHighConf_Agransc
affold.BlastResult.txt -export_search_strategy SonlyHighConf_Agranscaffold.BlastSearchStrategy.txt
Note the results of blasting against A. granulata sequences - did it hit something?
less SonlyHighConf_Agranscaffold.BlastResult.txt
Look A. granulata flanking sequences from the hits and blast in the NCBI website. If it hits something, I
got for the G term and ProtDB
```

Demographic Analyses

These analyses were done comparing datasets with all loci, only outliers and only neutral. The dataset with all loci can be split into the other two datasets using VCFtools

1. Make a list of loci to be kept in the outliers list with combined q-value of BayeScan and PCAdapt (> 0.05) and the true/false for each method and combination will be helpful for that, and maybe the outliers_combined_qvals.csv (from combining q-values of BayeScan and PCAdapt) will be helpful too.
2. Copy the VCF file with all loci in a directory containing the lists of loci to be kept in each subset (the ID column in the VCF must be present, used line numbers and made it match the IDs I was using to identify outliers. Another option is to copy the CHROM ID under the chromosome in the vcf and use the CHROM ID to identify outliers)
3. Load VCFtools module and run the VCFtools fork below

```
In [ ]: module load VCFtools/0.1.15-GCC-9.2.0-Pperl-5.30.1

#NZ-wide
vcftools --vrf MOD_input.vcf --snps outliers.txt --recode --recode-INFO-all --out outliers.out
#Southern
vcftools --vrf MOD_input.vcf --snps Soulliers.txt --recode --recode-INFO-all --out Soulliers.out
#NZ
vcftools --vrf NIMOD_input.vcf --snps NIoutliers.txt --recode --recode-INFO-all --out NIoutliers.out
vcftools --vrf NIMOD_input.vcf --snps NIneutral.txt --recode --recode-INFO-all --out NIneutral.out
#F
vcftools --vrf SIMOD_input.vcf --snps SIoutliers.txt --recode --recode-INFO-all --out SIoutliers.out
vcftools --vrf SIMOD_input.v
```



```
In [ ] : #!Rscript /git/Hub/shenglin-Liu/vcf2fs/01nb/master/vcf2fs.R
# Functions for creating site frequency spectra (SFS) from a VCF file.
# Functions for manipulating and plotting SFS.
# By Shenglin Liu, Feb 12, 2020.

## Only for diploids.
# Convention for file name extension (not mandatory):
# f (dadi), txt (dadi SNP data), fsc (fastSimcoal).
# Format of dadi:
# strictly two lines;
# dimension vector for the first line;
# SFS for the second line;
# no information for fold, mask or population names.
# Format of dadi SNP data:
# refIn refOut Allele1 Aaa Aro Allele2 Aaa Aro
# Aaa Aaa 49 56 7 0
## Frequent variables in the script:
# gt: a genotype object created from a VCF file and a popmap file.
# sfs: a table or array object representing SFS.
# vcf2gt(f,vcf,f.popmap)
# choose.pops(gt,pops)
# read.dadi(fsc,fsc.pops)
# viewMissing(gt)
# filter.gt(gt,filter.indi=NA,filter.snp=NA)
# vcf2fs.raw(gt,pops)
# vcf2fs.impute(gt,pops,sampleSizes)
# project.sfs(sfs,sampleSizes)
# write.sfs(sfs,f,output)
# read.dadi(fsc,fsc.pops=NULL)
# fold.sfs(sfs)
# sampleSizes(sfs)
# read.dadi(fsc,fsc.pops=NULL)
# read.dadi(fsc,fsc.pops=NULL)
## Read VCF file and popmap file, and create a gt (genotype) object.
# It will generate a named list of two elements, i.e., a popmap vector and a genotype matrix.
# The FORMAT column of the VCF file should start with GT tag;
# if a individual will be filtered out when NA, no filter=NA.
# The popmap file lists the population IDs of individuals;
# two columns; tab delimited; individual ID as 1st column, pop ID as 2nd column;
# 1st column can be arbitrary; 2nd column should be an integer or string vector;
# minimum sample size per population is 2.
vcf2gt<-function(f,vcf,f.popmap)
{
  #dadi<-getOptions("dadi")
  options(warn=1)
  #Read VCF file and popmap file.
  vcf.f<-as.matrix(read.table(f.vcf,sep="\t",stringsAsFactors=F)[,c(1:9)])
  popmap<-read.table(f.popmap,sep="\t",stringsAsFactors=F)[,2]
  nrow.vcf<-nrow(vcf.f)
  ncol.vcf<-ncol(vcf.f)
  # Parse genotypes.
  chrom<-substr(nrow(vcf.f),1,1)
  chrom2<-substr(nrow(vcf.f),3,3)
  chrom<-matrix(as.integer(chrom1)+as.integer(chrom2),nrow.vcf,ncol.vcf)
  options(warn=oldw)
  list(popmap=popmap,genotype=chrom)
}
## Choose populations (subsetting the gt object by populations).
choose.pops<-function(gt,pops)
{
  # pops: a character or integer vector; IDs of the chosen populations.
  popmap<-gt$popmap
  chrom<-gt$genotype
  # Choose populations.
  index<-sapply(pops,function(x){which(popmap==x)})
  index<-sort(unlist(index))
  popmap<-popmap[index]
  chrom<-chrom[index]
  list(popmap=popmap,genotype=chrom)
}
## Calculate sample sizes of the populations.
size<-function(gt)
{
  popmap<-gt$popmap
  chrom<-gt$genotype
  pops<-unique(popmap)
  sampleSizes<-sapply(pops,function(x){sum(popmap==x)})
  names(sampleSizes)<-pops
  sampleSizes
}
## Calculate minimum sample sizes of the populations accounting for the missing values.
minSampleSize<-function(gt)
{
  # pops: a character or integer vector; IDs of populations to be included in the SFS.
  popmap<-gt$popmap
  chrom<-gt$genotype
  pops<-unique(popmap)
  sampleSizes<-sapply(pops,function(x){min(rowSums(!is.na(chrom[,popmap==x]))))})
  names(sampleSizes)<-pops
  sampleSizes
}
# View the distribution of the missing values.
# It will output a SFS based on row count without accounting for the missing values.
viewMissing<-function(gt)
{
  # pops: a character or integer vector; IDs of populations to be included in the SFS.
  popmap<-gt$popmap
  chrom<-gt$genotype
  # Filter individuals according to missing values.
  if(!is.na(filter.indi))
  {
    missingVal<-is.na(chrom)
    delete.indi<-which(colSums(missingVal)>filter.indi)
    if(length(delete.indi)>0)
    {
      chrom<-chrom[, -delete.indi]
      popmap<-popmap[ -delete.indi]
      rm(missingVal)
    }
  }
  # Filter SFS according to missing values.
  if(!is.na(filter.snp))
  {
    missingVal<-is.na(chrom)
    delete.snp<-which(rowSums(missingVal)>filter.snp)
    if(length(delete.snp)>0)
    {
      chrom<-chrom[ -delete.snp,]
      rm(missingVal)
    }
  }
  list(popmap=popmap,genotype=chrom)
}
## Generate a SFS (table object) from the gt object.
# It will output a SFS based on row count without accounting for the missing values.
# pops: a character or integer vector; IDs of populations to be included in the SFS.
gt2sfs.raw<-function(gt,pops)
{
  # pops: a character or integer vector; IDs of populations to be included in the SFS.
  popmap<-gt$popmap
  chrom<-gt$genotype
  nrow.vcf<-nrow(chrom)
  ncol.vcf<-ncol(chrom)
  n.pops<-length(pops)
  # Number of chromosomes.
  ns.chr<-sapply(pops,function(x){sum(popmap==x)})*2
  # SFS based on row count.
  cnt<-matrix(0,nrow.vcf,n.pops)
  for(i in 1:n.pops)
  {
    index<-which(popmap==pops[i])
    cnt[,i]<-rowSums(chrom[,index],na.rm=T)
    ext<-c(ext,list(0:ns.chr[i]))
  }
  ext<-as.matrix(expand.grid(ext))
  cnt<-data.frame(rbind(cnt,ext))
  sfs<-table(cnt-1)
  names(dimnames(sfs.raw))<-pops
  sfs.raw<-table(cnt-1)
}
## Generate a SFS (table object) from the gt object.
# It will output a SFS based on row count without accounting for the missing values.
# pops: a character or integer vector; IDs of populations to be included in the SFS.
gt2sfs.impute<-function(gt,pops,sampleSizes)
{
  # pops: a character or integer vector; IDs of populations to be included in the SFS.
  popmap<-gt$popmap
  chrom<-gt$genotype/2
  nrow.vcf<-nrow(chrom)
  ncol.vcf<-ncol(chrom)
  n.pops<-length(pops)
  # Number of chromosomes.
  ns.chr<-sampleSizes*2
  # SFS based on count imputed through binomial distribution.
  cnt<-matrix(0,nrow.vcf,n.pops)
  ext<-list()
  for(i in 1:n.pops)
  {
    index<-which(popmap==pops[i])
    p<-rowMeans(chrom[,index],na.rm=T)
    cnt[,i]<-rbinom(nrow.vcf,ns.chr[i],p)
    ext<-c(ext,list(0:ns.chr[i]))
  }
  ext<-as.matrix(expand.grid(ext))
  cnt<-data.frame(rbind(cnt,ext))
  sfs<-table(cnt-1)
  names(dimnames(sfs.imp))<-pops
  sfs.imp<-table(cnt-1)
}
## Project a SFS to smaller sample sizes according to hypergeometric distribution.
# sampleSizes: an integer vector; downsized number of individuals in each population.
project.sfs<-function(sfs,sampleSizes)
{
  # project<-function(a,m,n)
  {
    b<-numeric(m+1)
    for(i in 0:m){
      for(j in 0:(n-m+i))
      {
        b[i+j]<-b[i+j]+choose(m,i)*choose(m-j,m-j)/choose(n,j)*a[j+1]
      }
    }
    dims<-dim(sfs)
    l.dims<-length(dims)
    sfs.pro<-sfs
    if(l.dims==1)
    {
      ms<-sampleSizes[1]^2
      ns<-dims[1]-1
      sfs.pro<-array(project(sfs.pro,m,n))
    }
    else
    {
      for(k in 1:l.dims)
      {
        ms<-sampleSizes[k]^2
        ns<-dims[k]-1
        sfs.pro<-apply(sfs.pro,c(1:l.dims)[ -k],function(x){project(x,m,n)})
        perm<-integer(1,dim)
        perm[k]<-1
        perm[ -k]<-2:l.dims
        sfs.pro<-aperm(sfs.pro,perm)
      }
    }
    names(dims(sfs.pro))<-sapply(sampleSizes,function(x){x:0})
    names(dimnames(sfs.pro))<-names(dimnames(sfs))
    sfs.pro<-as.table(sfs.pro)
    sfs.pro
  }
}
## Write a SFS to a file in dadi format.
write.dadi<-function(f,output)
{
  sink(f,output)
  fcat<-fcat(sfs)
  cat(fcat)
  cat(aperm(sfs))
  sink()
}
## Read a SFS from a dadi format file.
read.dadi<-function(f,dadi,popnames=NULL)
{
  # dadi<-scan(f,dadi,nlines=1,comment.char="#",quiet=T)
  sfs<-scan(f,dadi,sfsc=1,nlines=1,comment.char="#",quiet=T)
  sfsc<-aperm(array(sfs,dims))
  sfs<-as.table(sfs)
  dimnames(sfs)<-lapply(dims-1,function(x){x:0})
  names(dimnames(sfs))<-popnames
  sfs
}
## Fold a SFS.
fold.sfs<-function(sfs)
{
  sfs[[]<-sfs$rev(sfs)
  dims<-dim(sfs)
  cnt.pool<-rowSums(expand.grid(lapply(dims-1,function(x){x:0})))
  index<-cnt.pool<-sum(dims-1)/2
  sfs[index]<-0
  index<-cnt.pool<-sum(dims-1)/2
  sfs[index]<-sfs[index]/2
  sfs
}
## Resampling from a SFS for the purpose of bootstrapping.
sample.sfs<-function(sfs)
{
  sfs<-sum(sfs)
  sfs[[]<-rpois(length(sfs),sfs)
  sfs[[]<-n.site<-sum(sfs)
}
## Transform the gt object to a headered data frame of dadi SNP data format.
# The function will output a headered data frame as such:
# refIn refOut Allele1 Aaa Aro Allele2 Aaa Aro
# Aaa Aaa 49 56 7 0
# pops: a character or integer vector; IDs of populations to be included in the SFS.
gt2snps<-function(gt,pops)
{
  # pops: a character or integer vector; IDs of populations to be included in the SFS.
  popmap<-gt$popmap
  chrom<-gt$genotype
  nrow.vcf<-nrow(chrom)
  ncol.vcf<-ncol(chrom)
  n.pops<-length(pops)
  index<-which(popmap==pops[i])
  ref<-integer(0)
  alt<-integer(0)
  for(i in 1:n.pops)
  {
    index<-which(popmap==pops[i])
    ms<-rowMeans(chrom[,index],na.rm=T)
    ref<-cbind(ref,rowSums(!is.na(chrom[,index]))^2-2*ms)
    alt<-cbind(alt,temp)
  }
  snps<-data.frame("AAA","AAA","A","ref","T",alt)
  names(snps)<-c("refIn","refOut","Allele1","Allele2",pops,"Allele2",pops)
}
## Transform dadi SNP data format to SFS (table object).
# It will output a SFS based on row count.
# snp: a headered data frame as such:
# refIn refOut Allele1 Aaa Aro Allele2 Aaa Aro
# Aaa Aaa 49 56 7 0
# pops: a character vector; IDs of populations to be included in the SFS.
# sampleSizes: an integer vector of the same length as pops;
# if NA, calculate number of diploid individuals of populations;
# if NA, calculate from the SNP data.
snp2sfs.raw<-function(snp,pops,sampleSizes=NA)
{
  i.ref<-sapply(pops,function(x){which(names(snp)==x)[1])
  i.alt<-sapply(pops,function(x){which(names(snp)==x)[2])
  names(i.ref)<-names(i.alt)<-pops
}
# Number of chromosomes.
if(is.na(sampleSizes[1]))
{
  ns.chr<-sapply(pops,function(x){max(snp[,i.ref[x]]*snp[,i.alt[x])})
}
else
{
  ns.chr<-sampleSizes*2
}
# SFS based on row count.
ext<-list()
for(n.chr in ns.chr){ext<-c(ext,list(0:n.chr))}
ext<-as.matrix(expand.grid(ext))
cnt<-data.frame(rbind(as.matrix(snp),ext))
sfs.raw<-table(cnt-1)
names(dimnames(sfs.raw))<-pops
sfs.raw
}
## Plot a SFS (barplot for 1D, image for 2D).
plot.sfs<-function(sfs,cr=c(1,20000),colScheme=function(x)rainbow(x,s=0.8),...)
{
  pops<-names(dimnames(sfs))
  dims<-dim(sfs)
  n.pops<-length(dims)
  sfs[[]<-sfs[length(sfs)]<-0
  # Color coding for number of SNPs.
  cr<-log10(cr)
  nbc<-101
  breaks<-10*log(seq(r[1],r[2],length.out=nbc))
  if(n.pops==1)
  {
    barplot(sfs,names.arg=c(0,rep(NA,dims-2),dims-1),main=pops,...)
  }
  if(n.pops==2)
  {
    image(sfs,axes=F,col=colScheme(nbc-1),breaks=breaks,
  xlab=pops[1],ylab=pops[2],...)
    axis(1,at=c(0,1),labels=c(0,dims[1]-1))
    axis(2,at=c(0,1),labels=c(0,dims[2]-1))
    box()
  }
  # (n.pops>2)stop("Cannot plot SFS with dimension higher than 2!")
}
## Plot pairwise 2D-SFS among populations.
# Reads in from SFS files in dadi format.
# Each file should have names of population pairs separated by certain symbol.
# pops: names of the populations as in the file names; order matters.
# ext.o: extension name of the files for the expected SFS; if NA, not plotted.
# ext.s: a two-column number vector of separating sites when projections = 0, 0.5, 1, 1.5, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 101, 102, 103, 104, 105, 106, 107, 108, 109, 110, 111, 112, 113, 114, 115, 116, 117, 118, 119, 120, 121, 122, 123, 124, 125, 126, 127, 128, 129, 130, 131, 132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 142, 143, 144, 145, 146, 147, 148, 149, 150, 151, 152, 153, 154, 155, 156, 157, 158, 159, 160, 161, 162, 163, 164, 165, 166, 167, 168, 169, 170, 171, 172, 173, 174, 175, 176, 177, 178, 179, 180, 181, 182, 183, 184, 185, 186, 187, 188, 189, 190, 191, 192, 193, 194, 195, 196, 197, 198, 199, 200, 201, 202, 203, 204, 205, 206, 207, 208, 209, 210, 211, 212, 213, 214, 215, 216, 217, 218, 219, 220, 221, 222, 223, 224, 225, 226, 227, 228, 229, 230, 231, 232, 233, 234, 235, 236, 237, 238, 239, 240, 241, 242, 243, 244, 245, 246, 247, 248, 249, 250, 251, 252, 253, 254, 255, 256, 257, 258, 259, 260, 261, 262, 263, 264, 265, 266, 267, 268, 269, 270, 271, 272, 273, 274, 275, 276, 277, 278, 279, 280, 281, 282, 283, 284, 285, 286, 287, 288, 289, 290, 291, 292, 293, 294, 295, 296, 297, 298, 299, 300, 301, 302, 303, 304, 305, 306, 307, 308, 309, 310, 311, 312, 313, 314, 315, 316, 317, 318, 319, 320, 321, 322, 323, 324, 325, 326, 327, 328, 329, 330, 331, 332, 333, 334, 335, 336, 337, 338, 339, 340, 341, 342, 343, 344, 345, 346, 347, 348, 349, 350, 351, 352, 353, 354, 355, 356, 357, 358, 359, 360, 361, 362, 363, 364, 365, 366, 367, 368, 369, 370, 371, 372, 373, 374, 375, 376, 377, 378, 379, 380, 381, 382, 383, 384, 385, 386, 387, 388, 389, 390, 391, 392, 393, 394, 395, 396, 397, 398, 399, 400, 401, 402, 403, 404, 405, 406, 407, 408, 409, 410, 411, 412, 413, 414, 415, 416, 417, 418, 419, 420, 421, 422, 423, 424, 425, 426, 427, 428, 429, 430, 431, 432, 433, 434, 435, 436, 437, 438, 439, 440, 441, 442, 443, 444, 445, 446, 447, 448, 449, 450, 451, 452, 453, 454, 455, 456, 457, 458, 459, 460, 461, 462, 463, 464, 465, 466, 467, 468, 469, 470, 471, 472, 473, 474, 475, 476, 477, 478, 479, 480, 481, 482, 483, 484, 485, 486, 487, 488, 489, 490, 491, 492, 493, 494, 495, 496, 497, 498, 499, 500, 501, 502, 503, 504, 505, 506, 507, 508, 509, 510, 511, 512, 513, 514, 515, 516, 517, 518, 519, 520, 521, 522, 523, 524, 525, 526, 527, 528, 529, 530, 531, 532, 533, 534, 535, 536, 537, 538, 539, 540, 541, 542, 543, 544, 545, 546, 547, 548, 549, 550, 551, 552, 553, 554, 555, 556, 557, 558, 559, 560, 561, 562, 563, 564, 565, 566, 567, 568, 569, 570, 571, 572, 573, 574, 575, 576, 577, 578, 579, 580, 581, 582, 583, 584, 585, 586, 587, 588, 589, 590, 591, 592, 593, 594, 595, 596, 597, 598, 599, 600, 601, 602, 603, 604, 605, 606, 607, 608, 609, 610, 611, 612, 613, 614, 615, 616, 617, 618, 619, 620, 621, 622, 623, 624, 625, 626, 627, 628, 629, 630, 631, 632, 633, 634, 635, 636, 637, 638, 639, 640, 641, 642, 643, 644, 645, 646, 647, 648, 649, 650, 651, 652, 653, 654, 655, 656, 657, 658, 659, 660, 661, 662, 663, 664, 665, 666, 667, 668, 669, 670, 671, 672, 673, 674, 675, 676, 677, 678, 679, 680, 681, 682, 683, 684, 685, 686, 687, 688, 689, 690, 691, 692, 693, 694, 695, 696, 697, 698, 699, 700, 701, 702, 703, 704, 705, 706, 707, 708, 709, 710, 711, 712, 713, 714, 715, 716, 717, 718, 719, 720, 721, 722, 723, 724, 725, 726, 727, 728, 729, 730, 731, 732, 733, 734, 735, 736, 737, 738, 739, 740, 741, 742, 743, 744, 745, 746, 747, 748, 749, 750, 751, 752, 753, 754, 755, 756, 757, 758, 759, 760, 761, 762, 763, 764, 765, 766, 767, 768, 769, 770, 771, 772, 773, 774, 775, 776, 777, 778, 779, 780, 781, 782, 783, 784, 785, 786, 787, 788, 789, 790, 791, 792, 793, 794, 795, 796, 797, 798, 799, 800, 801, 802, 803, 804, 805, 806, 807, 808, 809, 810, 811, 812, 813, 814, 815, 816, 817, 818, 819, 820, 821, 822, 823, 824, 825, 826, 827, 828, 829, 830, 831, 832, 833, 834, 835, 836, 837, 838, 839, 840, 841, 842, 843, 844, 845, 846, 847, 848, 849, 850, 851, 852, 853, 854, 855, 856, 857, 858, 859, 860, 861, 862, 863, 864, 865, 866, 867, 868, 869, 870, 871, 872, 873, 874, 875, 876, 877, 878, 879, 880, 881, 882, 883, 884, 885, 886, 887, 888, 889, 890, 891, 892, 893, 894, 895, 896, 897, 898, 899, 900, 901, 902, 903, 904, 905, 906, 907, 908, 909, 910, 911, 912, 913, 914, 915, 916, 917, 918, 919, 920, 921, 922, 923, 924, 925, 926, 927, 928, 929, 930, 931, 932, 933, 934, 935, 936, 937, 938, 939, 940, 941, 942, 943, 944, 945, 946, 947, 948, 949, 950, 951, 952, 953, 954, 955, 956, 957, 958, 959, 960, 961, 962, 963, 964, 965, 966, 967, 968, 969, 970, 971, 972, 973, 974, 975, 976, 977, 978, 979, 980, 981, 982, 983, 984, 985, 986, 987, 988, 989, 990, 991, 992, 993, 994, 995, 996, 997, 998, 999, 1000, 1001, 1002, 1003, 1004, 1005, 1006, 1007, 1008, 1009, 1010, 1011, 1012, 1013, 1014, 1015, 1016, 1017, 1018, 1019, 1020, 1021, 1022, 1023, 1024, 1025, 1026, 1027, 1028, 1029, 1030, 1031, 1032, 1033, 1034, 1035, 1036, 1037, 1038, 1039, 1040, 1041, 1042, 1043, 1044, 1045, 1046, 1047, 1048, 1049, 1050, 1051, 1052, 1053, 1054, 1055, 1056, 1057, 1058, 1059, 1060, 1061, 1062, 1063, 1064, 1065, 1066, 1067, 1068, 1069, 1070, 1071, 1072, 1073, 1074, 1075, 1076, 1077, 1078, 1079, 1080, 1081, 1082, 1083, 1084, 1085, 1086, 1087, 1088, 1089, 1090, 1091, 1092, 1093, 1094, 1095, 1096, 1097, 1098, 1099, 1100, 1101, 1102, 1103, 1104, 1105, 1106, 1107, 1108, 1109, 1110, 1111, 1112, 1113, 1114, 1115, 1116, 1117, 1118, 1119, 1120, 1121, 1122, 1123, 1124, 1125, 1126, 1127, 1128, 1129, 1130, 1131, 1132, 1133, 1134, 1135, 1136, 1137, 1138, 1139, 1140, 1141, 1142, 1143, 1144, 1145, 1146, 1147, 1148, 1149, 1150, 1151, 1152, 1153, 1154, 1155, 1156, 1157, 1158, 1159, 1160, 1161, 1162, 1163, 1164, 1165, 1166, 1167, 1168, 1169, 1170, 1171, 1172, 1173, 1174, 1175, 1176, 1177, 1178, 1179, 1180, 1181, 1182, 1183, 1184, 1185, 1186, 1187, 1188, 1189, 1190, 1191, 1192, 1193, 1194, 1195, 1196, 1197, 1198, 1199, 1200, 1201, 1202, 1203, 1204, 1205, 1206, 1207, 1208, 1209, 1210, 1211, 1212, 1213, 1214, 1215, 1216, 1217, 1218, 1219, 1220, 1221, 1222, 1223, 1224, 1225, 1226, 1227, 1228, 1229, 1230, 1231, 1232, 1233, 1234, 1235, 1236, 1237, 1238, 1239, 1240, 1241, 1242, 1243, 1244, 1245, 1246, 1247, 1248, 1249, 1250, 1251, 1252, 1253, 1254, 1255, 1256, 1257, 1258, 1259, 1260, 1261, 1262, 1263, 1264, 1265, 1266, 1267, 1268, 1269, 1270, 1271, 1272, 1273, 1274, 1275, 1276, 1277, 1278, 1279, 1280, 1281, 1282, 1283, 1284, 1285, 1286, 1287, 1288, 1289, 1290, 1291, 1292, 1293, 1294, 1295, 1296, 1297, 1298, 1299, 1300, 1301, 1302, 1303, 1304, 1305, 1306, 1307, 1308, 1309, 1310, 1311, 1312, 1313, 1314, 1315, 1316, 1317, 1318, 1319, 1320, 1321, 1322, 1323, 1324, 1325, 1326, 1327, 1328, 1329, 1330, 1331, 1332, 1333, 1334, 1335, 1336, 1337, 1338, 1339, 1340, 1341, 1342, 1343, 1344, 1345, 1346, 1347, 1348, 1349, 1350, 1351, 1352, 1353, 1354, 1355, 1356, 1357, 1358, 1359, 1360, 1361, 1362, 1363, 1364, 1365, 1366, 1367, 1368, 1369, 1370, 1371, 1372, 1373, 1374, 1375, 1376, 1377, 1378, 1379, 1380, 1381, 1382, 1383, 1384, 1385, 1386, 1387, 1388, 1389, 1390, 1391, 1392, 1393, 1394, 1395, 1396, 1397, 1398, 1399, 1400, 1401, 1402, 1403, 1404, 1405, 1406, 1407, 1408, 1409, 1410, 1411, 1412, 1413, 1414, 1415, 1416, 1417, 1418, 1419, 1420, 1421, 1422, 1423, 1424, 1425, 1426, 1427, 1428, 1429, 1430, 1431, 1432, 1433, 1434, 1435, 1436, 1437, 1438, 1439, 1440, 1441, 1442, 1443, 1444, 1445, 1446, 1447, 1448, 1449, 1450, 1451, 1452, 1453, 1454, 1455, 1456, 1457, 1458, 1459, 1460, 1461, 1462, 1463, 1464, 1465, 1466, 1467, 1468, 1469, 1470, 1471, 1472, 1473, 1474, 1475, 1476, 1477, 1478, 1479, 1480, 1481, 1482, 1483, 1484, 1485, 1486, 1487, 1488, 1489, 1490, 1491, 1492, 1493, 1494, 1495, 1496, 1497, 1498, 1499, 1500, 1501, 1502, 1503, 1504, 1505, 1506, 1507, 1508, 1509, 1510, 1511, 1512, 1513, 1514, 1515, 1516, 1517, 1518, 1519, 1520, 1521, 1522, 1523, 1524, 1525, 1526, 1527, 1528, 1529, 1530, 1531, 1532, 1533, 1534, 1535, 1536, 1537, 1538, 1539, 1540, 1541, 1542, 1543, 1544, 1545, 1546, 1547, 1548, 1549, 1550, 1551, 1552, 1553, 1554, 1555, 1556, 1557, 1558, 155
```



```
In [ ] : less example.vcf2bayescan.spid
# spid-file generated: Tue Sep 17 16:13:54 NZST 2019

# VCF Parser questions
PARSER_FORMAT=VCF

# Only output SNPs with a phred-scaled quality of at least:
VCF_PARSER_QUAL_QUESTION=
# Select population definition file:
VCF_PARSER_POP_FILE_QUESTION=./pop_def_conserved.tsv
# What is the ploidy of the data?
VCF_PARSER_PLOIDY_QUESTION=
# Do you want to include a file with population definitions?
VCF_PARSER_POP_QUESTION=TRUE
# Output genotypes as missing if the phred-scale genotype quality is below:
VCF_PARSER_GTQUAL_QUESTION=
# Do you want to include non-polymorphic SNPs?
VCF_PARSER_MONOMORPHIC_QUESTION=
# Only output following individuals (ind1, ind2, ind4, ...):
VCF_PARSER_IND_QUESTION=
# Only input following regions (refSeqName:start:end, multiple regions: whitespace separated):
VCF_PARSER_REGION_QUESTION=
# Output genotypes as missing if the read depth of a position for the sample is below:
VCF_PARSER_READ_QUESTION=
# Take most likely genotype if "PL" or "GL" is given in the genotype field?
VCF_PARSER_PL_QUESTION=
# Do you want to exclude loci with only missing data?
VCF_PARSER_EXC_MISSING_LOCI_QUESTION=

# GESTE / BayeScan writer questions
PARSER_FORMAT=GESTE_BAYE_SCAN

# Specify which data type should be included in the GESTE / BayeScan file (GESTE / BayeScan can only anal
# yze one data type per file):
GESTE_BAYE_SCAN_WRITER_DATA_TYPE_QUESTION=SNP
```

```
In [ ] : #population definition file (one individual per line, no header, second column containing population defin
# ition)
head NZ-wide.popsmap.txt

OnegAK10      AK
OnegAK11      AK
OnegAK12      AK
OnegAK13      AK
OnegAK14      AK
OnegAK15      AK
OnegAK16      AK
OnegAK17      AK
OnegAK18      AK
OnegAN01      AN
```

```
In [ ] : #Command to run PGDSpider (make sure to change the -inputfile to the name of the vcf input to be converted,
# the -outputfile to the desired name of the converted output file (with .txt at the end) and
# the -spid to the corresponding spid file for this dataset (the spid file contains information on the popul
# ation definition file, so it is important that it matched the correct vcf file)

less southern_split_vcf2bayescan.sh
#!/bin/bash -e
#SBATCH --job-name=PGDSpider
#SBATCH --account=ne060308
#SBATCH --time=00:15:00
#SBATCH --mem=12G
#SBATCH --output=PS%j.out
#SBATCH --error=PS%j.err

module load Java/1.8.0_144
module load SAMtools/1.9-GCC-7.4.0
module load BamTools/2.5.1-gimkl-2018b

java -Xmx2024m -Xms12m -jar PGDSpider2-cli.jar -inputfile ./southern_split.vcf -inputformat VCF -outputfil
# e ./southern_splitBayescan.txt -outputformat GESTE_BAYE_SCAN -spid southern_split_vcf2bayescan.spid
```

Environmental files with ncdf4 in R

Note: the input are downloaded netCDF version 4 files from the database (in this case NOAA)

```
In [ ] : #Extracting environmental information from netCDF version 4 files

#Load packages
library(chron)
library(RColorBrewer)
library(lattice)
library(ncdf4)

#check working directory
getwd()

#Read a NetCDF file using the ncdf package.
#Also, set the values for some temporary variables: ncname is the name of the ncCDF file, while dname is t
# he name of the variable that will be read in
ncname = "FILE NAME.nc"
dname = "sst" #sea surface temperature

#Open the netcdf file
ncin = nc_open(ncname)
#print basic info about the file
print(ncin)

#Read the longitude and latitude
lon = ncvar_get(ncin, "lon")
nlon = dim(lon)
head(lon)

lat = ncvar_get(ncin, "lat")
nlat = dim(lat)
head(lat)

#confirm the dimensions of the data
print(c(nlon, nlat))

t = ncvar_get(ncin, "time")
tunits = ncatt_get(ncin, "time", "units")
t
# get the number of values and list them
nt = dim(t)

#Print the time units string.
tunits

#Read the variable and its longname, units and fill value (.FillValue) attributes
tmp.array = ncvar_get(ncin, dname)
dname = ncatt_get(ncin, dname, "long.name")
units = ncatt_get(ncin, dname, "units")
fillvalue = ncatt_get(ncin, dname, "FillValue")

#the dimension of the array containing the values of the variables can be confirmed
dim(tmp.array)

#Read a set of global attributes (metadata)
title = ncatt_get(ncin, 0, "title")
institution = ncatt_get(ncin, 0, "institution")
datasource = ncatt_get(ncin, 0, "source")
references = ncatt_get(ncin, 0, "references")
history = ncatt_get(ncin, 0, "history")
Conventions = ncatt_get(ncin, 0, "Conventions")
#List the attribute values
title$value
institution$value
datasource$value
references$value
history$value
Conventions$value

#We are now done with the input dataset, so close it
nc_close(ncin)

#Convert the time variable from [time-since] units into [real] time values
# split the time units string into fields
tmp.slice = tmp.array[, , m]
dim(tmp.slice)
lonlat = expand.grid(lon, lat)
tmp.vec = as.vector(tmp.slice)
length(tmp.vec)
tmp.df01 = data.frame(cbind(lonlat, tmp.vec))
names(tmp.df01) = c("lon", "lat", paste(dname, as.character(m), sep = "_"))
head(na.omit(tmp.df01), 20)
#Save the data frame as a .csv file, using na.omit() to drop the observations with missing data
csvfile = "FILENAME.TO.SAVE.csv"
write.table(na.omit(tmp.df01), csvfile, row.names = FALSE, sep = ",")

#Convert the whole array to a data frame, and calculate minimum, maximum and the annual mean
tmp.vec.long = as.vector(tmp.array)
tmp.mat = matrix(tmp.vec.long, nrow = nlon * nlat, ncol = nt)
dim(tmp.mat)
head(na.omit(tmp.mat))
lonlat = expand.grid(lon, lat)
tmp.df02 = data.frame(cbind(lonlat, tmp.mat))
#CHECK WHAT ARE EACH COLUMNS OF MY DATAFRAME BEFRE USING THE NAMES AS BELOW, WHICH WERE FOR ONLY 12 COLUMN
5 AKA NOT 365
names(tmp.df02) = c("lon", "lat", 1:365)
options(width = 110)
head(na.omit(tmp.df02, 365))

#Get annual mean, mtwa and mtco values and add them the second data frame
tmp.df02mtwa <- apply(tmp.df02[3:367], 1, max) # mtwa
tmp.df02mtco <- apply(tmp.df02[3:367], 1, min) # mtco
tmp.df02smat <- apply(tmp.df02[3:367], 1, mean) # annual (i.e. row) means
head(na.omit(tmp.df02))

dim(na.omit(tmp.df02))

#Write the second data frame out as a .csv file, dropping NAs.
csvfile = "FILENAME.TO.SAVE.csv"
write.table(na.omit(tmp.df02), csvfile, row.names = FALSE, sep = ",")

#Create a third data frame, with only non-missing values
tmp.df03 <- na.omit(tmp.df02)
head(tmp.df03)

write.table(na.omit(tmp.df03), "FILENAME.TO.SAVE.csv", row.names = FALSE, sep = ",")
```