

SPOTlight / WesDe

Marc Elosua-Bayes

7/7/2021

```
library(SPOTlight)
```

```
## Loading required package: dplyr
##
## Attaching package: 'dplyr'
## The following objects are masked from 'package:stats':
##
##   filter, lag
## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union
```

SPOTlight objects

Load data

```
spotlight1 <- readRDS(file = "~/Downloads/spotlight_ls.rds.gz")
spotlight2 <- readRDS(file = "~/Downloads/spotlight_ls_2.rds.gz")
```

Extract info

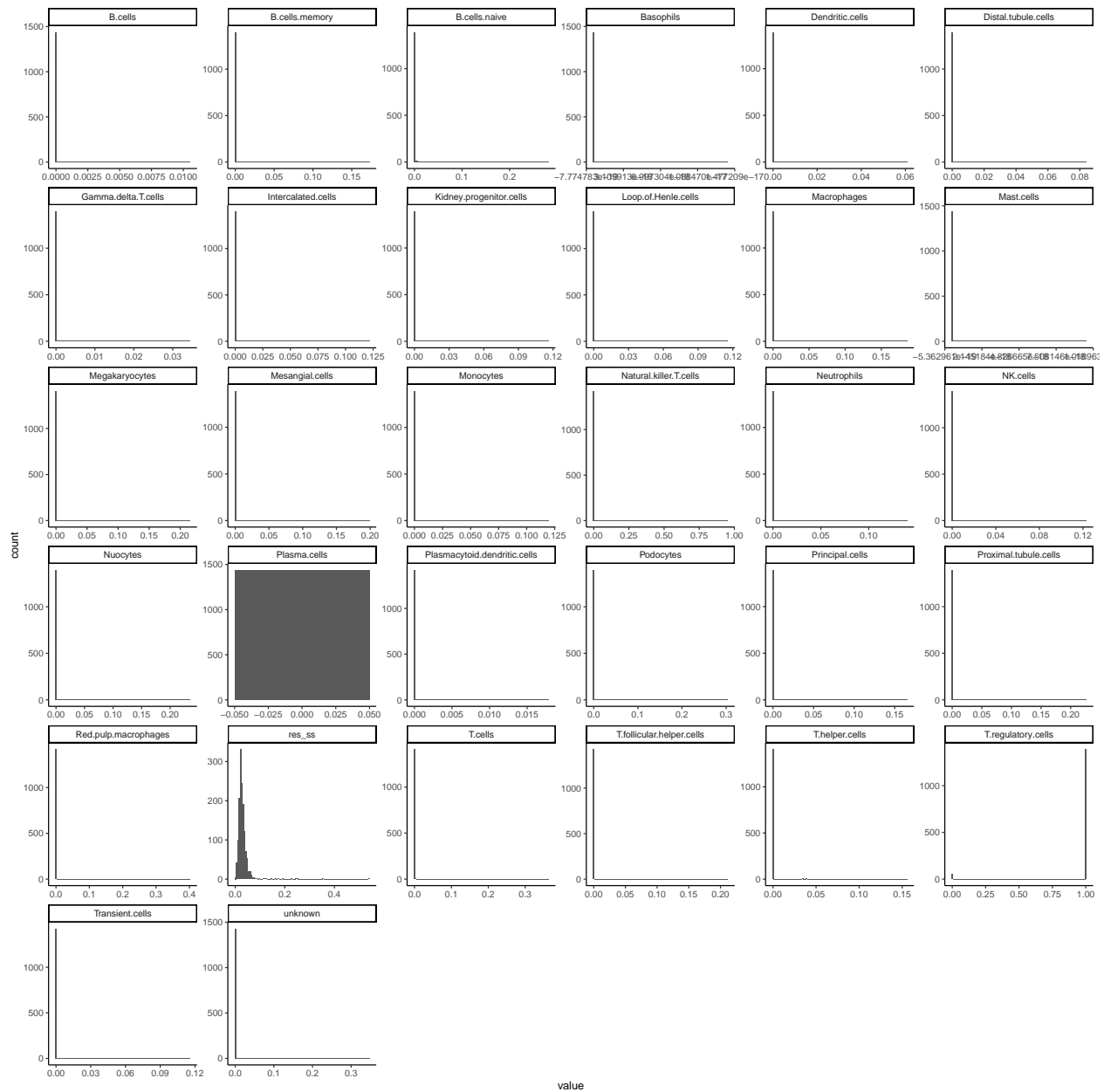
```
NMF1 <- spotlight1[[1]][[1]]
ct_labs1 <- spotlight1[[1]][[2]]
decon1 <- spotlight1[[2]]

NMF2 <- spotlight2[[1]][[1]]
ct_labs2 <- spotlight2[[1]][[2]]
decon2 <- spotlight2[[2]]
```

Explore these objects

```
colnames(decon1)
```

```
## [1] "B.cells" "B.cells.memory"
## [3] "B.cells.naive" "Basophils"
## [5] "Dendritic.cells" "Distal.tubule.cells"
## [7] "Gamma.delta.T.cells" "Intercalated.cells"
## [9] "Kidney.progenitor.cells" "Loop.of.Henle.cells"
## [11] "Macrophages" "Mast.cells"
## [13] "Megakaryocytes" "Mesangial.cells"
```

We see how regulatory T cells capture most of the signal. Lets look into this a bit further...

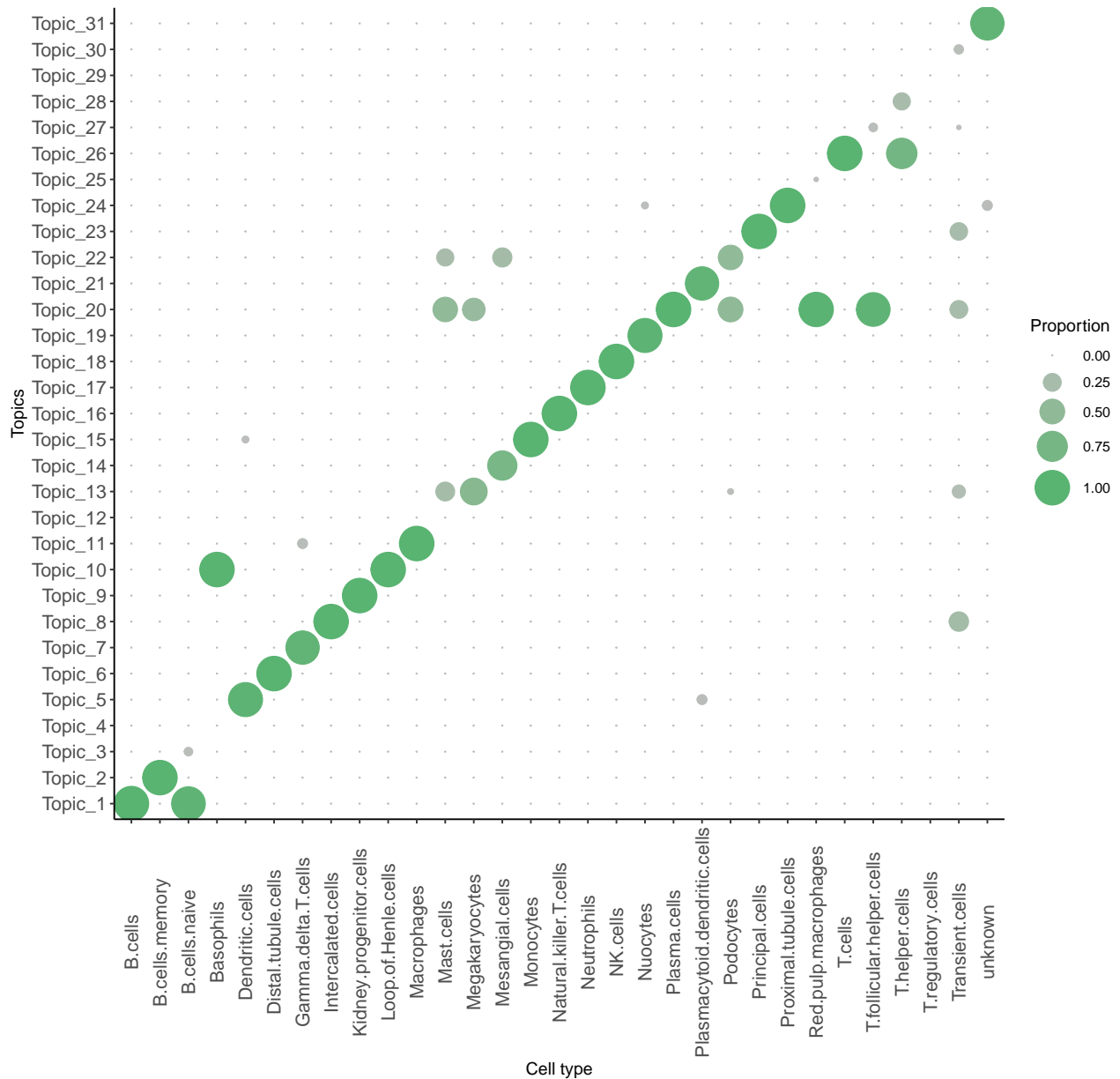
Look at the topic profiles:

```
h1 <- NMF::coef(NMF1)
rownames(h1) <- paste("Topic", 1:nrow(h1), sep = "_")
topic_profile_plts1 <- SPOTlight::dot_plot_profiles_fun(
  h = h1,
  train_cell_clust = ct_labs1)
```

Summarised topic profiles:

```
topic_profile_plts1[[2]] + ggplot2::theme(
  axis.text.x = ggplot2::element_text(angle = 90),
  axis.text = ggplot2::element_text(size = 12))
```

NMF: Topic profiles by cell type



Above we can see how T reg doesn't have a representative topic profile, something that could be happening is that Tregs are defined by a topic profile with a low weight which would make it capture a lot of signal and overshadow other signal:

Lets take a look here

```
# Prep data
h_df <- data.frame(t(h1))
colnames(h_df) <- gsub(".", " ", colnames(h_df), fixed = TRUE)
h_ds <- h_df/rowSums(h_df)
h_ds[, "clust_vr"] <- ct_labs1

# Organize data
summarise_topic_df <- h_ds %>%
  dplyr::group_by(clust_vr) %>%
```

```

dplyr::summarise_all(list(median)) %>%
tibble::remove_rownames(.) %>%
tibble::column_to_rownames("clust_vr")

summarise_topic_df[c("B.cells", "T.regulatory.cells"), ]

```

```

##           Topic_1      Topic_2      Topic_3      Topic_4
## B.cells      1.000000e+00 1.473367e-132 1.061810e-104 2.979202e-140
## T.regulatory.cells 1.012583e-116 3.700013e-123 1.408013e-149 4.858607e-153
##           Topic_5      Topic_6      Topic_7      Topic_8
## B.cells      3.264886e-128 0.000000e+00 6.417130e-152 6.401143e-283
## T.regulatory.cells 1.276929e-131 6.255418e-316 1.624871e-191 2.609962e-290
##           Topic_9 Topic_10      Topic_11      Topic_12      Topic_13
## B.cells      0          0 8.776509e-185 2.644161e-153 1.142498e-282
## T.regulatory.cells 0          0 1.402335e-155 1.004886e-153 0.000000e+00
##           Topic_14      Topic_15      Topic_16      Topic_17
## B.cells      8.761477e-283 1.584300e-171 2.928760e-269 9.280821e-233
## T.regulatory.cells 6.580530e-287 3.821942e-147 3.586926e-276 8.619301e-119
##           Topic_18      Topic_19      Topic_20      Topic_21
## B.cells      1.138082e-134 6.663898e-308 6.916919e-323 3.025777e-124
## T.regulatory.cells 7.403144e-175 0.000000e+00 0.000000e+00 5.590934e-127
##           Topic_22 Topic_23 Topic_24      Topic_25      Topic_26
## B.cells      3.754541e-292 0          0 9.352009e-229 5.974090e-131
## T.regulatory.cells 1.236199e-314 0          0 2.009710e-276 1.236381e-176
##           Topic_27      Topic_28      Topic_29      Topic_30
## B.cells      1.911883e-233 5.499163e-111 1.612499e-153 5.938242e-132
## T.regulatory.cells 1.721553e-270 2.430760e-135 1.842315e-87 9.133637e-168
##           Topic_31
## B.cells      2.435566e-207
## T.regulatory.cells 1.366724e-193

```

We can see how B cells for example is dominated by topic_1 while T regs have no defined topic and it will try to fit noise until engulfing all of it.

Profiles for all cells

```

topic_profile_plts1[[1]] + theme(axis.text.x = element_text(angle = 90),
                                axis.text = element_text(size = 12))

```



Regulatory T cells seem to have a large variability, there is no one or 2 shared topics across all cells. Consider reclustering these cells or filter the gene set used to try to obtain a cleaner signal.

lastly, look at the genes for each topic

```
basis_spotlight1 <- data.frame(NMF::basis(NMF1))

colnames(basis_spotlight1) <- glue::glue("Topic-{1:ncol(basis_spotlight1)}")

basis_spotlight1 %>%
  round(., 5) %>%
  DT::datatable(., filter = "top")
```

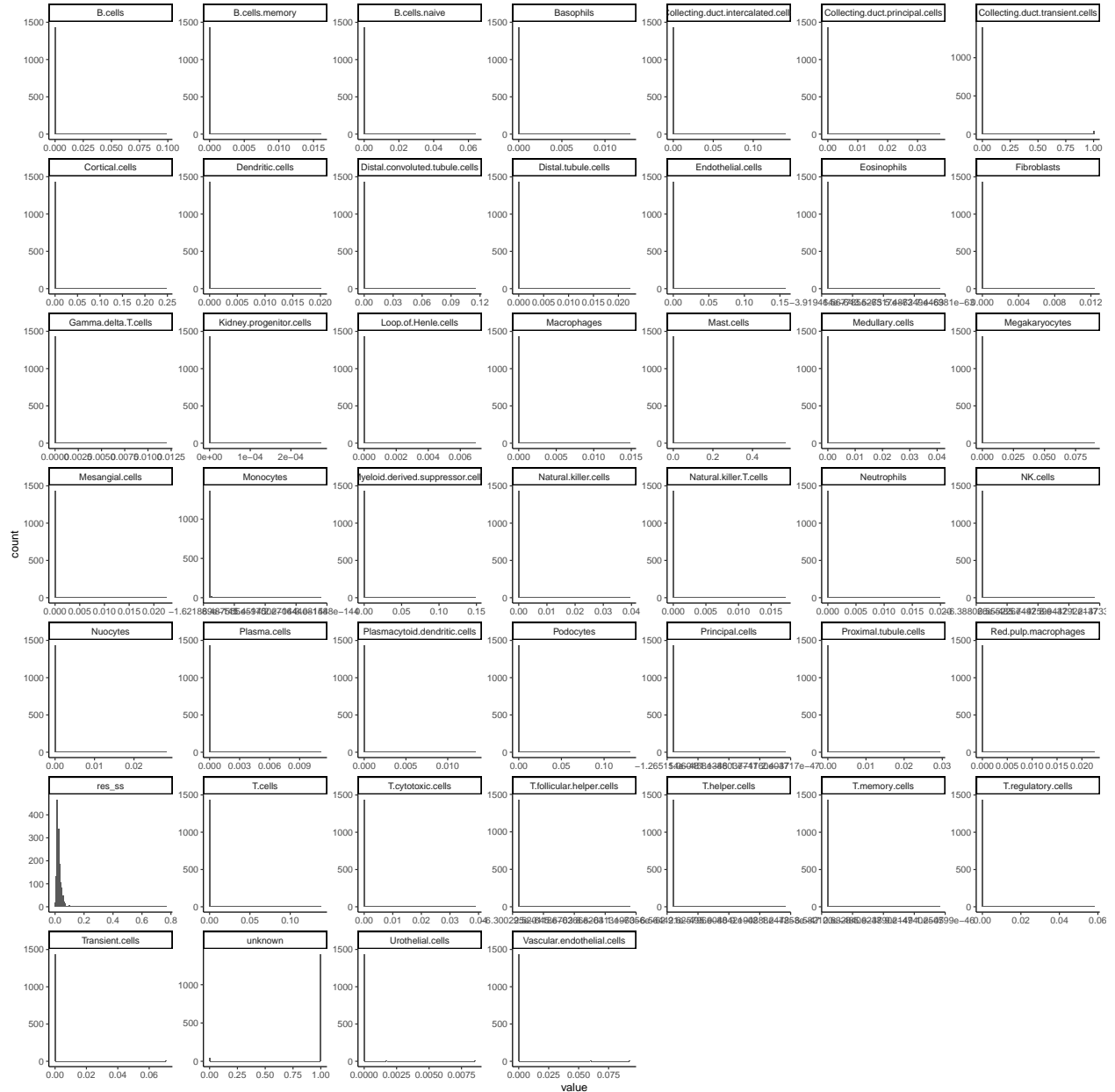
```
## Warning in instance$preRenderHook(instance): It seems your data is too big
## for client-side DataTables. You may consider server-side processing: https://
```


Look at summary statistics for each decomposed cell type

```
# Hmisc::describe(x = decon2)
```

Visualization of the predicted probabilities

```
data.frame(decon2) %>%
  tidyr::pivot_longer(cols = dplyr::everything()) %>%
  ggplot2::ggplot(ggplot2::aes(x = value)) +
  ggplot2::geom_histogram(bins = 100) +
  ggplot2::facet_wrap(. ~ name, scales = "free") +
  ggplot2::theme_classic()
```



We can see how in this case the Unknown cell type captures all the biological signal.

Look at the topic profiles:

```

h2 <- NMF::coef(NMF2)
rownames(h2) <- paste("Topic", 1:nrow(h2), sep = "_")
topic_profile_plts2 <- SPOTlight::dot_plot_profiles_fun(
  h = h2,
  train_cell_clust = ct_labs2)

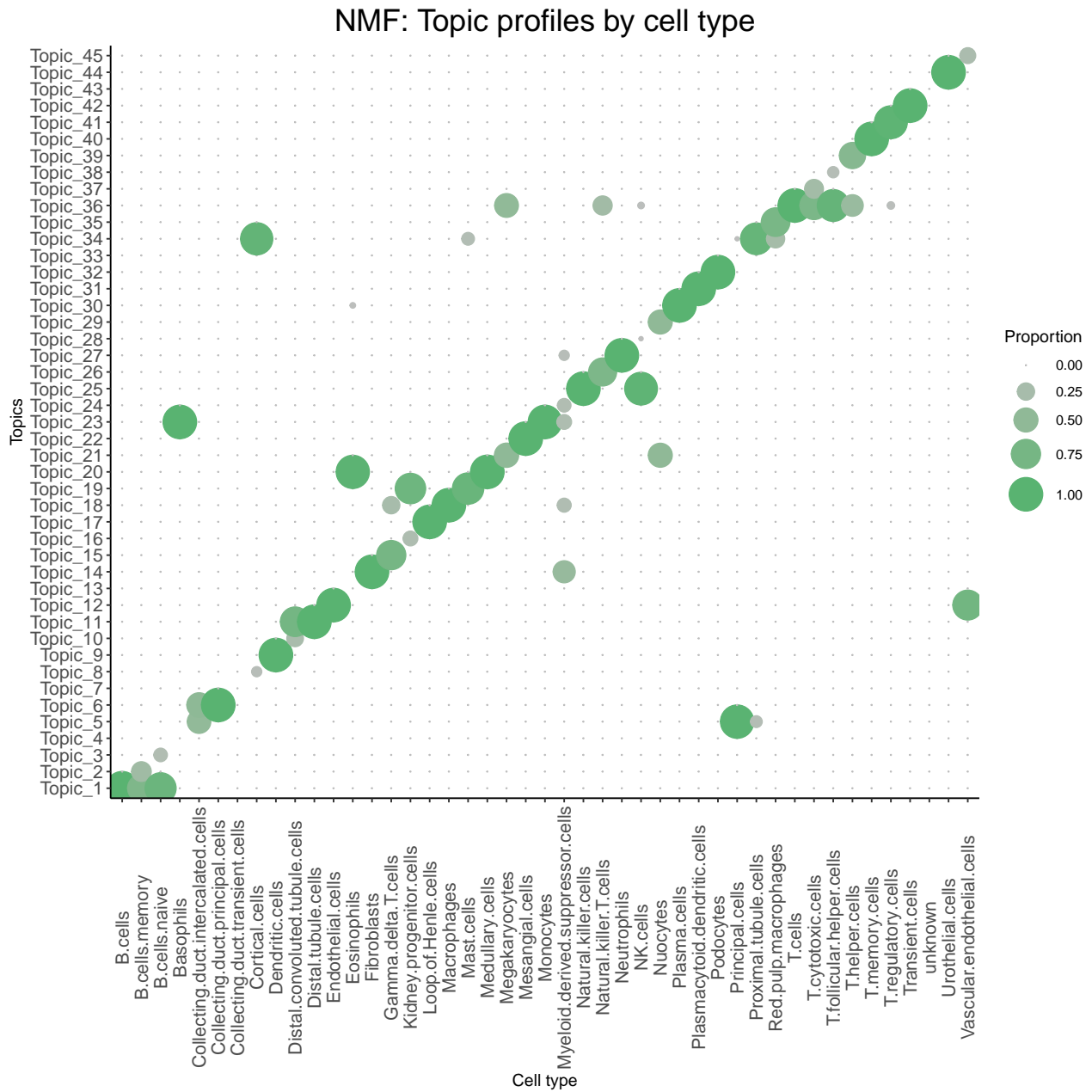
```

Summarised topic profiles:

```

topic_profile_plts2[[2]] +
  ggplot2::theme(
    axis.text.x = ggplot2::element_text(angle = 90),
    axis.text = ggplot2::element_text(size = 12))

```



We see the unknown cell type has the same issue as before with T regs... lets follow the same process as

before:

```
# Prep data
h_df <- data.frame(t(h2))
colnames(h_df) <- gsub(".", " ", colnames(h_df), fixed = TRUE)
h_ds <- h_df/rowSums(h_df)
h_ds[, "clust_vr"] <- ct_labs2

# Organize data
summarise_topic_df <- h_ds %>%
  dplyr::group_by(clust_vr) %>%
  dplyr::summarise_all(list(median)) %>%
  tibble::remove_rownames(.) %>%
  tibble::column_to_rownames("clust_vr")

summarise_topic_df[c("B.cells", "unknown"), ]

##           Topic_1      Topic_2      Topic_3      Topic_4      Topic_5
## B.cells  9.569699e-01  6.398193e-88  4.998908e-83  8.667823e-133  1.159149e-190
## unknown  1.059014e-252  8.687149e-168  9.143148e-168  6.816725e-290  1.199423e-227
##           Topic_6      Topic_7      Topic_8      Topic_9      Topic_10
## B.cells  0.000000e+00  3.836642e-208  1.802978e-98  3.952849e-129  2.048662e-126
## unknown  1.905537e-271  3.115399e-275  6.011495e-178  1.091470e-220  1.107063e-249
##           Topic_11 Topic_12      Topic_13      Topic_14      Topic_15
## B.cells  6.195972e-292          0  7.900458e-125  5.412527e-173  2.527358e-127
## unknown  1.482197e-322          0  2.343032e-239  1.702626e-300  3.127862e-247
##           Topic_16      Topic_17      Topic_18      Topic_19 Topic_20
## B.cells  6.175962e-122  2.478869e-258  1.790457e-195  8.073985e-271          0
## unknown  1.601762e-195  4.515188e-296  1.152906e-257  1.369043e-297          0
##           Topic_21      Topic_22      Topic_23      Topic_24      Topic_25
## B.cells  8.178214e-130  8.625639e-288  8.579625e-140  4.302269e-96  2.957696e-139
## unknown  5.836843e-258  6.990716e-304  2.092529e-281  2.915914e-191  2.151044e-305
##           Topic_26      Topic_27      Topic_28      Topic_29      Topic_30
## B.cells  1.935595e-102  1.685337e-234  2.224527e-95  3.292801e-84  6.459372e-124
## unknown  1.514614e-178  2.555155e-308  1.155338e-174  6.414014e-189  1.066195e-273
##           Topic_31      Topic_32      Topic_33      Topic_34      Topic_35
## B.cells  1.454772e-109  8.952524e-215  2.217329e-168  8.100327e-293  2.972661e-146
## unknown  1.364786e-183  6.618049e-280  1.320265e-237  0.000000e+00  4.387390e-216
##           Topic_36      Topic_37      Topic_38      Topic_39      Topic_40
## B.cells  2.163817e-125  1.220318e-90  6.694623e-93  7.322614e-96  2.863484e-111
## unknown  7.308671e-284  1.204820e-189  5.537398e-180  1.064330e-206  1.217565e-237
##           Topic_41      Topic_42      Topic_43      Topic_44      Topic_45
## B.cells  2.510611e-114  4.335872e-261  2.474174e-142  1.122988e-262  1.495680e-102
## unknown  5.816953e-247  2.044599e-282  1.222245e-145  8.162808e-306  1.526617e-200
```

Same as before, the B cells associated topic is dominated by topic_1 while unknowns have no defined topic and therefore it will try to fit noise to signal until engulfing all of it.

Profiles for all cells

```
topic_profile_plts2[[1]] +
  ggplot2::theme(axis.text.x = ggplot2::element_text(angle = 90),
                 axis.text = ggplot2::element_text(size = 12))
```



```
## [40] scales_1.1.1      stringi_1.6.2      farver_2.1.0
## [43] reshape2_1.4.4     doParallel_1.0.16  bslib_0.2.5.1
## [46] ellipsis_0.3.2     generics_0.1.0     vctrs_0.3.8
## [49] RColorBrewer_1.1-2 NMF_0.23.0         iterators_1.0.13
## [52] tools_4.0.4        glue_1.4.2         purrr_0.3.4
## [55] crosstalk_1.1.1    rngtools_1.5       processx_3.5.2
## [58] yaml_2.2.1         colorspace_2.0-1   cluster_2.1.1
## [61] knitr_1.33         sass_0.4.0
```