

Scaling up reproducible research for single cell transcriptomics using MetaNeighbor (Anticipated results)

MetaNeighbor anticipated results

Because MetaNeighbor is non-parametric, there is no fine-tuning to be done for any of the protocols presented here. Over time, we have identified two sources of potential error: bad highly variable gene selection and coding or formatting errors, which can be easily diagnosed by looking at AUROC heatmaps. We expect AUROCs to correctly represent global relationships between cell types and generalize across studies.

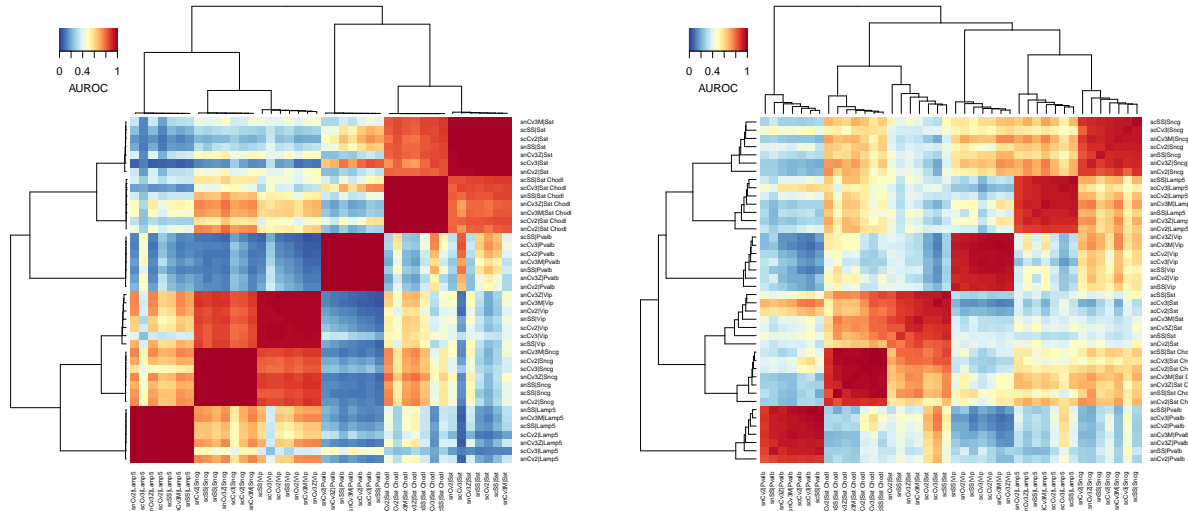
Bad gene set selection

The most common problem is to forget to select a set of highly variable genes, which is expected to dampen the impact of technical variability on neighbor voting. First, we present an example of a correct analysis, where we load the BICCN GABAergic neurons, select highly variable genes, and compute cluster similarities (see Protocol 1 for more details).

```
library(MetaNeighbor)
biccn_data = readRDS("biccn_gaba.rds")
biccn_hvgs = variableGenes(biccn_data, exp_labels = biccn_data$study_id)

# GOOD
aurocs = MetaNeighborUS(var_genes = biccn_hvgs,
                        dat = biccn_data,
                        study_id = biccn_data$study_id,
                        cell_type = biccn_data$joint_subclass_label,
                        fast_version = TRUE)
plotHeatmap(aurocs, cex = 0.5)

# BAD
aurocs = MetaNeighborUS(var_genes = sample(rownames(biccn_data), length(biccn_hvgs)),
                        dat = biccn_data,
                        study_id = biccn_data$study_id,
                        cell_type = biccn_data$joint_subclass_label,
                        fast_version = TRUE)
plotHeatmap(aurocs, cex = 0.5)
```



We recognize strong replicability structure, evidenced by the presence of dark red blocks. When we repeat the analysis with random genes, the replicability structure is still present, but we recognize two signatures of bad gene set selection: (a) AUROCs are low overall (shift to light red and orange), (b) within red blocks, there is a clear gradient structure. In our experience, there are 3 scenarios that lead to bad gene selection: errors in gene symbol conversion, errors when genes are stored as factors in R (that are implicitly converted to numerical values during indexing), forgetting to select highly variable genes altogether.

No overlap between datasets

The second problem occurs when there is no overlap between datasets. We illustrate this problem with the data from Protocol 2, where we expect all cell types from the Tasic dataset to be present in the pre-trained BICCN model. According to our expectations, all cell types have strong hits with BICCN clusters and we see a hierarchical structure that is consistent with prior biological knowledge: lighter red blocks corresponding to MGE and CGE-derived inhibitory neurons. We compare with the same block of code, where we mistakenly keep non-neurons from the BICCN taxonomy instead of inhibitory neurons. The lack of biological overlap can be deduced from 3 factors: (a) low AUROC values overall, (b) almost no strong hits (contrary to expectations), (c) lack of expected hierarchical structure (MGE and CGE derived interneurons).

```
library(scRNAseq)
tasic = TasicBrainData(ensembl = FALSE)
tasic$study_id = "tasic"

biccn_subclasses = read.table("pretrained_biccn_subclasses.txt", check.names = FALSE)
global_aurocs = MetaNeighborUS(
  trained_model = biccn_subclasses, dat = tasic,
  study_id = tasic$study_id, cell_type = tasic$primary_type,
  fast_version = TRUE
)
gabaergic_tasic = splitTestClusters(global_aurocs, k = 4)[[2]]

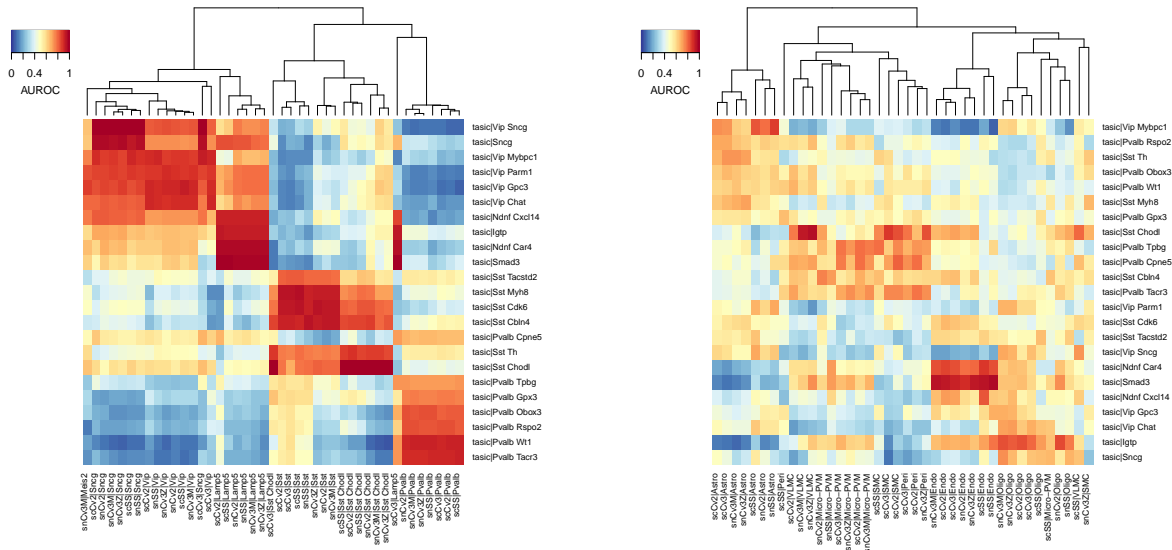
# GOOD
gabaergic_biccn = splitTrainClusters(global_aurocs[gabaergic_tasic,], k = 4)[[4]]
keep_cell = makeClusterName(tasic$study_id, tasic$primary_type) %in% gabaergic_tasic
tasic_subdata = tasic[, keep_cell]
aurocs = MetaNeighborUS(
```

```

trained_model = biccn_subclasses[, gabaergic_biccn],
dat = tasic_subdata, study_id = tasic_subdata$study_id,
cell_type = tasic_subdata$primary_type, fast_version = TRUE
)
plotHeatmapPretrained(aurocs, cex = 0.7)

# BAD: non-neurons instead of GABAergic neurons
gabaergic_biccn = splitTrainClusters(global_aurocs, k = 5)[[1]]
keep_cell = makeClusterName(tasic$study_id, tasic$primary_type) %in% gabaergic_tasic
tasic_subdata = tasic[, keep_cell]
aurocs = MetaNeighborUS(
  trained_model = biccn_subclasses[, gabaergic_biccn],
  dat = tasic_subdata, study_id = tasic_subdata$study_id,
  cell_type = tasic_subdata$primary_type, fast_version = TRUE
)
plotHeatmapPretrained(aurocs, cex = 0.7)

```



Pretrained MetaNeighbor: bad name formatting

The third problem we have encountered is a mistake that occurs when loading pre-trained model and forgetting to specify “check.names = FALSE”, which is essential to preserve correct formatting of cell type names. Below, we present an example of correct code based on data from Protocol 2. We obtain the expected replicability structure, with evidence of strong hits across all cell types (see Protocol 2 for further details and analyses). When we forget “check.names = FALSE”, MetaNeighbor is unable to correctly recognize dataset names and cell type names in the pre-trained model, the similarity computations become meaningless, leading to AUROC values that are essentially 0.5. This problem is easy to diagnose and fix, but can be very confusing when it occurs.

```

# GOOD
biccn_subclasses = read.table("pretrained_biccn_subclasses.txt", check.names = FALSE)
aurocs = MetaNeighborUS(
  trained_model = biccn_subclasses, dat = tasic,
  study_id = tasic$study_id, cell_type = tasic$primary_type,
  fast_version = TRUE
)

```

Figure 1: Heatmap of AUROC values for 100 cell lines across 100 genes. The heatmap is color-coded by AUROC, with red indicating high values (0.8-1.0) and blue indicating low values (0.0-0.4). Hierarchical clustering is shown on the top and left. A dendrogram on the right shows the relationship between cell lines, with a color scale for AUROC values. The cell lines are grouped into two main clusters: a large cluster of 80 cell lines and a smaller cluster of 20 cell lines. The genes are grouped into two main clusters: a large cluster of 80 genes and a smaller cluster of 20 genes. The heatmap shows that the large cluster of cell lines has high AUROC values for most genes, while the smaller cluster of cell lines has low AUROC values for most genes.

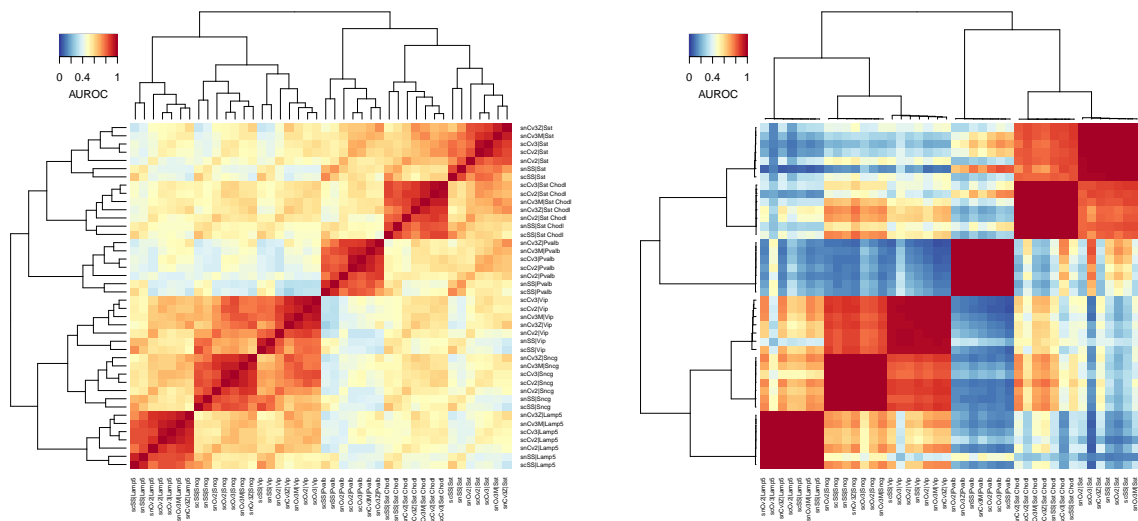
In their computation, MetaNeighbor’s AUROCs are directly related to Spearman correlations. More precisely, all computations are based on average Spearman correlations between cells from two cell types, but include an additional prediction that alleviate batch effects, while keeping an interpretability power that is comparable to correlations (where AUROC = 0 maps to correlation = -1, AUROC = 0.5 maps to correlation = 0, AUROC = 1 maps to correlation = 1).

```
cell_types = as.factor(
  makeClusterName(biccn_data$study_id, biccn_data$joint_subclass_label)
)
normalization_factor = Matrix::colSums(assay(biccn_data)) / 1000000
```

```
cpm = assay(biccn_data)
cpm@x = cpm@x / rep.int(normalization_factor, diff(cpm@p))
cpm = as.matrix(cpm[biccn_hvgs,])

centroids = sapply(levels(cell_types), function(ct) {
  matrixStats::rowMeans2(log2(cpm+1), cols = cell_types == ct)
})

centroid_cor = cor(centroids, method = "spearman")
aurocs = MetaNeighborUS(var_genes = biccn_hvgs,
                        dat = biccn_data,
                        study_id = biccn_data$study_id,
                        cell_type = biccn_data$joint_subclass_label,
                        fast_version = TRUE)
plotHeatmap((1+centroid_cor)/2, cex = 0.5)
plotHeatmap(aurocs, cex = 0.5)
```



Compared to correlations, AUROCs have one additional “parameter”: the outgroup used for the prediction task. In Protocol 1, we illustrated how the outgroup can be controlled and interpreted. A deviation from $\text{AUROC} = 1$ can thus be interpreted as a combination of two factors: lack of similarity between cell types and choice of outgroup (difficulty of prediction task). If the outgroup is well controlled, AUROC values will generalize across studies and fundamentally indicate the quality of the clustering. For a given cell type in a given background (for example Sst cells in an unbiased sample of MOp interneurons), the similarity to Sst cells in another dataset (or any other interneuron type for that matter) should be in the range of similarity observed within the BICCN.

As a robust alternative to centroid correlations, MetaNeighbor AUROCs can be applied to simple pre-processing tasks, such as identifying and selecting cell types that overlap between datasets before applying a merging framework. However, beyond the purely applicative viewpoint, we believe that MetaNeighbor-style AUROCs are a stepping stone toward a generalizable formalization of cell type similarity.