

# Tutorial for running DeMixT and TmS

Wenyi Wang

In this tutorial, we use a subset of the bulk RNAseq data of prostate adenocarcinoma (PRAD) from TCGA (<https://portal.gdc.cancer.gov/>) as an example to demonstrate how to run **DeMixT**. The analysis pipeline consists of the following steps:

The analysis pipeline consists of the following steps:

- Obtaining raw read counts for the tumor and normal RNAseq data
- Loading libraries and data
- Data preprocessing
- Deconvolution using DeMixT

# Obtain raw read counts for the tumor and normal RNAseq data

The raw read counts for the tumor and normal samples from TCGA PRAD are downloaded from [TCGA data portal](#). One can also generate the raw read counts from fastq or bam files by following the [GDC mRNA Analysis Pipeline](#).

# Load libraries and data

## Load library

```
1 library(DeMixT)
```

# Load libraries and data

## Load library

```
1 library(DeMixT)
```

## Load input data

```
1 load("../docs/etc/PRAD.RData")
```

Three data are included in the `PRAD.RData`.

- `PRAD`: Read counts matrix (gene x sample) with genes as row names and sample ids as column names.
- `Normal.id`: TCGA ids of PRAD normal samples.
- `Tumor.id` TCGA ids of PRAD tumor samples.

# A glimpse of PRAD:

```
1 head(PRAD[,1:5])
2 cat('Number of genes: ', dim(PRAD)[1], '\n')
3 cat('Number of normal sample: ', length(Normal.id), '\n')
4 cat('Number of tumor sample: ', length(Tumor.id), '\n')
```

## Output:

	TCGA-CH-5761-11A	TCGA-CH-5767-11B	TCGA-EJ-7115-11A	TCGA-EJ-7123-11A	TCGA-EJ-7125-11A
TSPAN6	3876	7095	5542	2747	8465
TNMD	14	51	13	24	63
DPM1	1162	2665	1544	1974	2984
SCYL3	777	1517	1096	1231	1514
Clorf112	136	343	214	280	339
FGR	230	511	263	755	262
Number of genes:	59427				
Number of normal sample:	20				
Number of tumor sample:	30				

# Data preprocessing

Conduct data cleaning and normalization before running DeMixT.

```
1 PRAD = PRAD[, c(Normal.id, Tumor.id)]
2 selected.genes = 9000
3 cutoff_normal_range = c(0.1, 1.0)
4 cutoff_tumor_range = c(0, 2.5)
5 cutoff_step = 0.1
6 preprocessed_data = DeMixT_preprocessing(PRAD,
7                                         Normal.id,
8                                         Tumor.id,
9                                         selected.genes,
10                                        cutoff_normal_range,
11                                        cutoff_tumor_range,
12                                        cutoff_step)
13 PRAD_filter = preprocessed_data$count.matrix
14 sd_cutoff_normal = preprocessed_data$sd_cutoff_normal
15 sd_cutoff_tumor = preprocessed_data$sd_cutoff_tumor
```

# Data preprocessing

Conduct data cleaning and normalization before running DeMixT.

```
1 PRAD = PRAD[, c(Normal.id, Tumor.id)]
2 selected.genes = 9000
3 cutoff_normal_range = c(0.1, 1.0)
4 cutoff_tumor_range = c(0, 2.5)
5 cutoff_step = 0.1
6 preprocessed_data = DeMixT_preprocessing(PRAD,
7                                         Normal.id,
8                                         Tumor.id,
9                                         selected.genes,
10                                        cutoff_normal_range,
11                                        cutoff_tumor_range,
12                                        cutoff_step)
13 PRAD_filter = preprocessed_data$count.matrix
14 sd_cutoff_normal = preprocessed_data$sd_cutoff_normal
15 sd_cutoff_tumor = preprocessed_data$sd_cutoff_tumor
16 cat("Normal sd cutoff:", preprocessed_data$sd_cutoff_normal, "\n")
17 cat("Tumor sd cutoff:", preprocessed_data$sd_cutoff_tumor, "\n")
18 cat('Number of genes after filtering: ', dim(PRAD_filter)[1], '\n')
```

## Output:

```
1 Normal sd cutoff: 0.1 0.9
2 Tumor sd cutoff: 0 0.6
3 Number of genes after filtering: 9103
```

The function `DeMixT_preprocessing` identifies two intervals based on the standard deviation of the log-transformed gene expression for normal and tumor samples, respectively, within the pre-defined ranges (`cutoff_normal_range` and `cutoff_tumor_range`).

In this example, we choose to select about 9000 genes before running DeMixT with the GS (Gene Selection) method to ensure that our model-based gene selection maintains good statistical properties.

DeMixT\_preprocessing outputs a list object `preprocessed_data` containing:

- `preprocessed_data$count.matrix`: Preprocessed count matrix
- `preprocessed_data$sd_cutoff_normal`: Actual cut-off value when desired genes are selected for normal samples
- `preprocessed_data$sd_cutoff_tumor`: Actual cut-off value when desired genes are selected for tumor samples

# Deconvolution using DeMixT

- To optimize the parameters in DeMixT for input data, we recommend testing an array of combinations of number of spike-ins and number of selected genes.
- The number of CPU cores used by the DeMixT function for parallel computing is specified by the parameter `nthread`. By default, `nthread = total_number_of_cores_on_the_machine - 1`. Users can adjust `nthread` to any number between 0 and the total number of cores available on the machine.
- For reference, DeMixT takes approximately 3-4 minutes to process the PRAD data in this tutorial for each parameter combination when `nthread` is set to 55.

```

1 # Due to the random initial values and the spike-in samples used in the DeMixT function,
2 # we recommend that users set seeds to ensure reproducibility.
3 # This seed setting will be incorporated internally in DeMixT in the next update.
4
5 set.seed(1234)
6
7 data.Y = SummarizedExperiment(assays = list(counts = PRAD_filter[, Tumor.id]))
8 data.N1 <- SummarizedExperiment(assays = list(counts = PRAD_filter[, Normal.id]))
9
10 # In practice, we set the maximum number of spike-in as min(n/3, 200),
11 # where n is the number of samples.
12 nspikesin_list = c(0, 5, 10)
13 # One may set a wider range than provided below for studies other than TCGA.
14 ngene.selected_list = c(500, 1000, 1500, 2500)
15
16 for(nspikesin in nspikesin_list){
17   for(ngene.selected in ngene.selected_list){
18     name = paste("PRAD_demixt_GS_res_nspikesin", nspikesin, "ngene.selected",
19                 ngene.selected, sep = ".")

```

**Note:** We use a profiling likelihood-based method to select genes, during which we calculate confidence intervals for the model parameters using the inverse of the Hessian matrix. When the input data (e.g., gene expression levels from spatial transcriptomic data) is sparse, the Hessian matrix will contain infinite values, hence those confidence intervals can't be calculated. In this case, gene selection will be performed through differential expression analysis (identical to [DeMix\\_DE](#)). This alternative is automatically performed inside [DeMix\\_GS](#) when the above situation happens.

```

1  PiT_GS_PRAD <- c()
2  row_names <- c()
3
4  for(nspikesin in nspikesin_list){
5      for(ngene.selected in ngene.selected_list){
6          name_simplify <- paste(nspikesin, ngene.selected, sep = "_")
7          row_names <- c(row_names, name_simplify)
8
9          name = paste("PRAD_demixt_GS_res_nspikesin", nspikesin,
10                      "ngene.selected", ngene.selected, sep = "_");
11          name = paste(name, ".RData", sep = "")
12          load(name)
13          PiT_GS_PRAD <- cbind(PiT_GS_PRAD, res$pi[2, ])
14      }
15  }
16
17  colnames(PiT_GS_PRAD) <- row_names

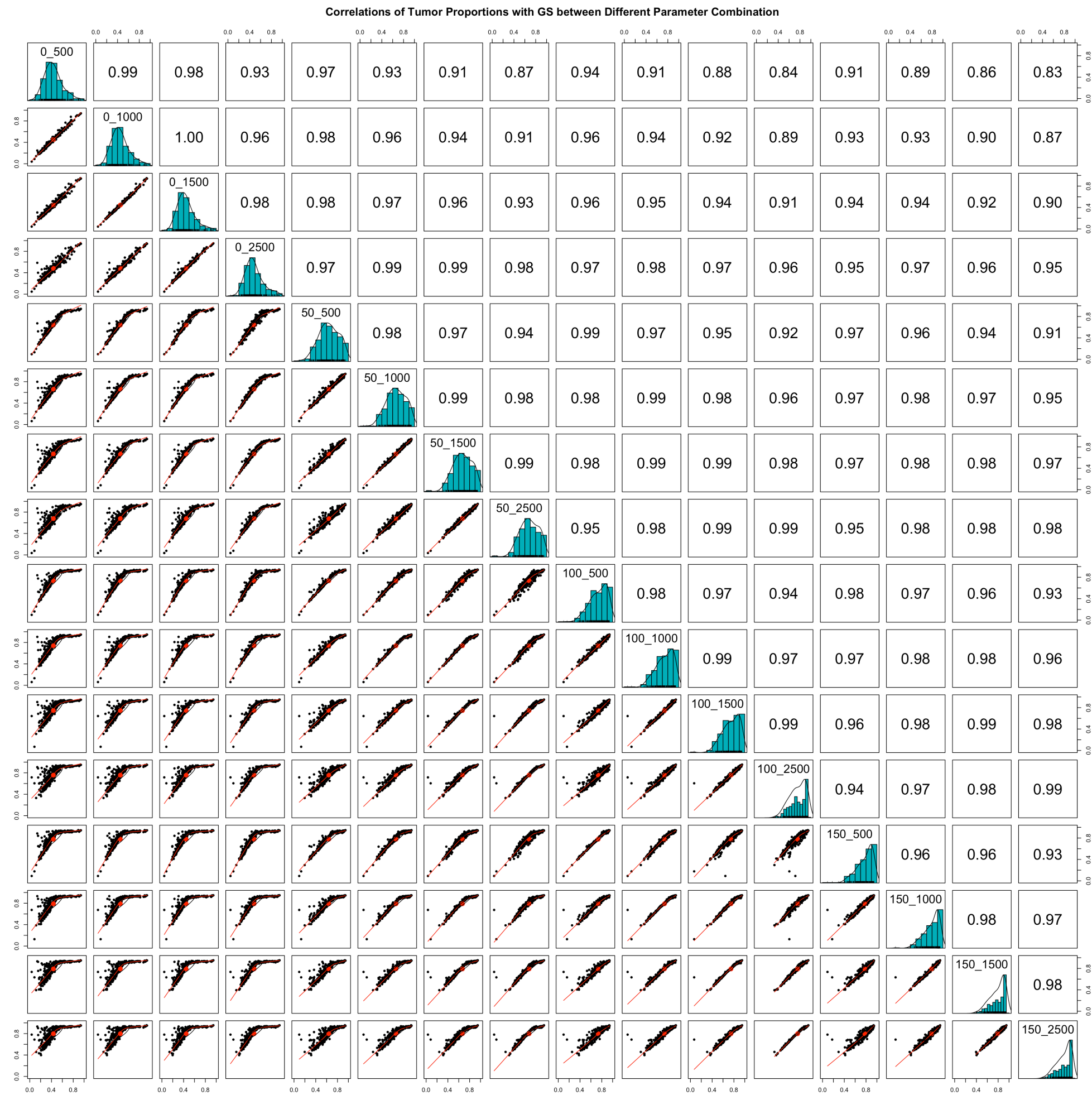
```

This step saves the deconvolution results (PiT) into a dataframe with columns named after the combination of the number of spike-ins and number of genes selected. Then one can calculate and plot the pairwise correlations of estimated tumor proportions across different parameter combinations as shown in the next slide.

```

1 pairs.panels(PiT_GS_PRAD,
2             method = "spearman",
3             # correlation method
4             hist.col = "#00AFBB",
5             density = TRUE,
6             # show density plots
7             ellipses = TRUE,
8             # show correlation ellipses
9             main = 'Correlations of Tumor
10 Proportions with GS between
11 Different Parameter Combination',
12             xlim = c(0,1),
13             ylim = c(0,1))

```



Print out the average pairwise correlation of tumor proportions across different parameter combinations.

```

1  PiT_GS_PRAD <- as.data.frame(PiT_GS_PRAD)
2  Spearman_correlations <- list()
3
4  for(entry_1 in colnames(PiT_GS_PRAD)) {
5    cor.values <- c()
6    for (entry_2 in colnames(PiT_GS_PRAD)) {
7      if (entry_1 == entry_2)
8        next
9
10     cor.values <- c(cor.values,
11                    cor(PiT_GS_PRAD[, entry_1],
12                       PiT_GS_PRAD[, entry_2],
13                       method = "spearman"))
14   }
15
16   Spearman_correlations[[entry_1]] <- mean(cor.values)
17 }
18
19 Spearman_correlations <- unlist(Spearman_correlations)

```

	num.spikein_num.selected.gene	mean.correlation
2	0_500	0.8641319
3	0_1000	0.9453534
4	0_1500	0.9401355
5	0_2500	0.9375468
6	5_500	0.9207604
7	5_1000	0.9542926
8	5_1500	0.9460006
9	5_2500	0.8992011
10	10_500	0.9237941
11	10_1000	0.9357266
12	10_1500	0.9249267
13	10_2500	0.9002124

We suggest selecting the optimal parameter combination that produces the **highest average correlation** of estimated tumor proportions.

Additionally, users are encouraged to evaluate the skewness of the PiT estimation distribution compared to a normal distribution centered around 0.5, as **Significant skewness may indicate biased estimation**.

```
1 num.spikein_num.selected.gene    mean.correlation
2 0_500                            0.8641319
3 0_1000                           0.9453534
4 0_1500                           0.9401355
5 0_2500                           0.9375468
6 5_500                            0.9207604
7 5_1000                           0.9542926
8 5_1500                           0.9460006
9 5_2500                           0.8992011
10 10_500                          0.9237941
11 10_1000                         0.9357266
12 10_1500                         0.9249267
13 10_2500                         0.9002124
```

We suggest selecting the optimal parameter combination that produces the **highest average correlation** of estimated tumor proportions.

Additionally, users are encouraged to evaluate the skewness of the PiT estimation distribution compared to a normal distribution centered around 0.5, as **Significant skewness may indicate biased estimation**.

	num.spikein_num.selected.gene	mean.correlation
1	0_500	0.8641319
2	0_1000	0.9453534
3	0_1500	0.9401355
4	0_2500	0.9375468
5	5_500	0.9207604
6	5_1000	0.9542926
7	5_1500	0.9460006
8	5_2500	0.8992011
9	10_500	0.9237941
10	10_1000	0.9357266
11	10_1500	0.9249267
12	10_2500	0.9002124

Based on these criteria, **spike-ins = 5** and **number of selected genes = 1000** are identified as the optimal parameter combination. Using these parameters, we can obtain the corresponding tumor proportions.

```
1 data.frame(sample.id=Tumor.id, PiT=PiT_GS_PRAD[['5_1000']])
2
3 sample.id          PiT
4 TCGA-2A-A8VL-01A    0.7596888
5 TCGA-2A-A8VO-01A    0.8421716
6 TCGA-2A-A8VT-01A    0.8662378
7 TCGA-2A-A8VV-01A    0.7616749
8 TCGA-2A-A8W1-01A    0.8291091
9 TCGA-2A-A8W3-01A    0.8159406
10 TCGA-CH-5737-01A    0.7314935
11 TCGA-CH-5738-01A    0.4614545
12 TCGA-CH-5739-01A    0.6349423
13 TCGA-CH-5740-01A    0.7095117
```

## List the tumor specific expression

```
1 ## Load the corresponding deconvolved gene expression
2 load("PRAD_demixt_GS_res_nspikesin_5_ngene.selected_1000.RData")
3 res$ExprT[1:5, 1:5]
4
5          TCGA-2A-A8VL-01A  TCGA-2A-A8VO-01A  TCGA-2A-A8VT-01A  TCGA-2A-A8VV-01A  TCGA-2A-A8W1-01A
6 DPM1          1710.194          1466.484          1680.4562          1644.944          1812.600
7 FUCA2          3782.990          4083.382           961.0578          4165.612          1896.901
8 GCLC          2382.106          1826.957          1527.4895          1409.707          1913.784
9 LAS1L          3329.766          2758.414          3520.9410          2834.415          2530.621
10 ENPP4          2099.591          3123.365          3173.3516          2856.371          7413.330
```

Instead of selecting using the parameter combination with the highest correlation, one can also select the parameter combination that produces estimated tumor proportions that are most biologically meaningful.

## Next,

We will provide a simple TmS tutorial which uses The estimated tumor-specific proportions (PiT) generated from DeMixT. For more details, visit <https://wwylab.github.io/TmS/articles/TmS.html>.

# TmS Calculation

Tumor-specific total mRNA expression (TmS) from bulk sequencing data, taking into account tumor transcript proportion, purity and ploidy, which are estimable through transcriptomic/genomic deconvolution.

TmS analysis pipeline consists of the following steps:

- **Step 1:** Estimate the proportion of total RNA expression ( $\pi$ ) from tumor cells using RNAseq data.
  - Achieved by using [DeMixT](#)[1]
- **Step 2:** Estimate the proportion of tumor cells and total copies of haploid genomes, i.e., tumor purity ( $\rho$ ) and tumor ploidy ( $\psi$ ), using matched DNAseq or SNP array data.
- **Step 3:** Calculate TmS, the per cell haploid genome total RNA expression for tumor, using the estimated ( $\pi$ ), ( $\rho$ ) and ( $\psi$ ):  $TmS = [\pi(1 - \rho)^2] / [(1 - \pi)\rho\psi]$ .

## Step 3: Calculate TmS using the estimated $(\pi)$ , $(\rho)$ and $(\psi)$ .

### Consensus TmS estimation

For DNA-based deconvolution methods such as [ASCAT](#) and [ABSOLUTE](#), there could be multiple tumor purity  $\rho$  and ploidy  $\psi$  pairs that have similar likelihoods. Both ASCAT and ABSOLUTE can accurately estimate the product of purity  $\rho$  and ploidy  $\psi$ ; however, they sometimes lack power to identify and separately. TmS is derived from the product of tumor ploidy and the odds of tumor purity. Hence, it is potentially more robust to ambiguity in the tumor purity and ploidy estimation, ensuring the robustness of the TmS calculation.

To calculate one final set of TmS values for a maximum number of samples, we use a consensus approach. We first calculate TmS values with tumor purity and ploidy estimates derived from both ABSOLUTE and ASCAT, and then fit a linear regression model on the log2-transformed  $TmS_{ASCAT}$  using the log2-transformed  $TmS_{ABSOLUTE}$  as a predictor variable. We remove samples with Cook's distance  $\geq 4/n$  and calculate the final

$$TmS = \sqrt{TmS_{ASCAT} \times TmS_{ABSOLUTE}}$$

The agreement between the two methods in ploidy values was low in 20% of TCGA samples. However, a large portion of these samples showed consistency in the TmS values using either ASCAT and ABSOLUTE, reducing the number of filtered TCGA samples to ~5% (264 samples) [2]. This result supports the robustness of our consensus approach.

**Input:** Tumor-specific total mRNA proportions, tumor purities, tumor ploidies

**Output:** Consensus TmS

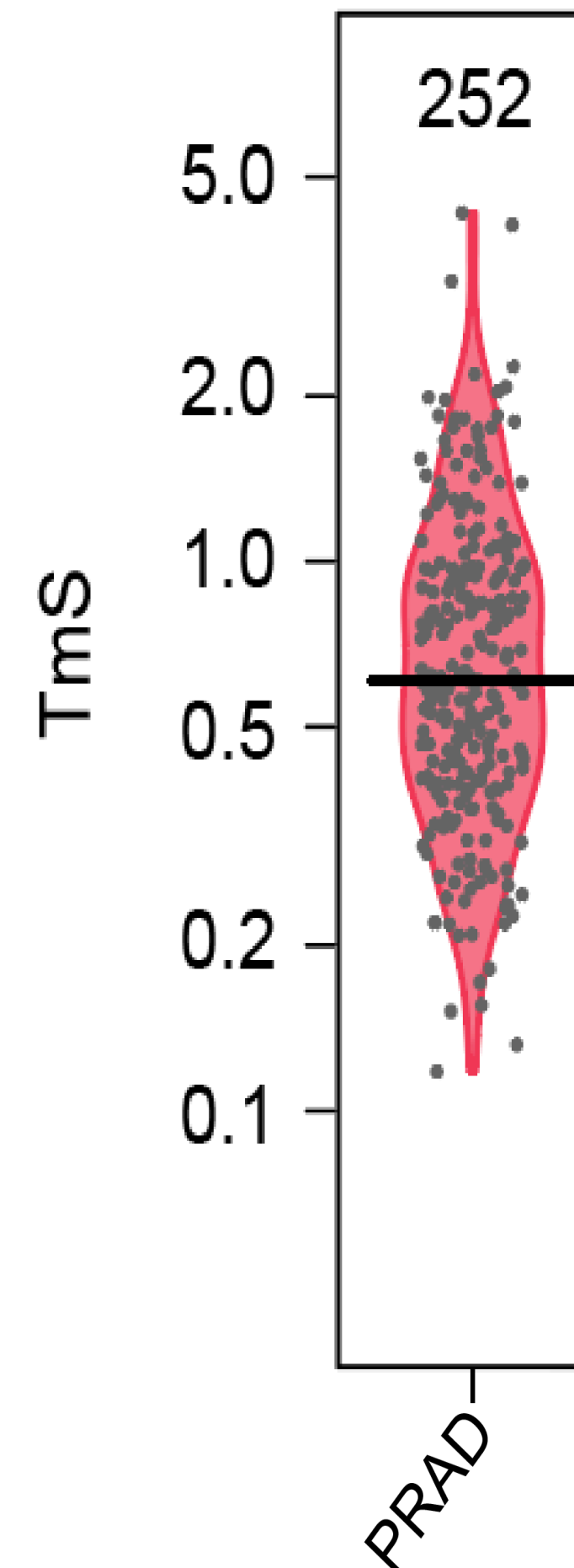
- **p**: Tumor-specific total mRNA proportions estimated by DeMixT
- **rho\_ASCAT**: tumor purity estimated by ASCAT
- **phi\_ASCAT**: tumor ploidy estimated by ASCAT
- **rho\_ABSOLUTE**: tumor purity estimated by ABSOLUTE
- **phi\_ABSOLUTE**: tumor ploidy estimated by ABSOLUTE

```

1 TmS.calculate = function(p, rho, phi){
2   return(2 * p * (1 - rho) / (phi * rho * (1 - p)))
3 }
4
5 TmS.ASCAT = TmS.calculate(p, rho_ASCAT, phi_ASCAT)
6 TmS.ABSOLUTE = TmS.calculate(p, rho_ABSOLUTE, phi_ABSOLUTE)
7
8 TmS.df = data.frame(TmS.ASCAT.log2 = log2(TmS.ASCAT),
9                    TmS.ABSOLUTE.log2 = log2(TmS.ABSOLUTE))
10
11 lm.fit = lm(TmS.df$TmS.ABSOLUTE.log2 ~ TmS.df$TmS.ASCAT.log2)
12 summary(lm.fit)
13 cooks = cooks.distance(lm.fit)
14 cooks.threshold = 4 / nrow(TmS.df)
15
16 cook.status = rep("Accept", nrow(TmS.df))
17 cook.status[cooks > cooks.threshold] = "Outlier"
18 TmS.df = TmS.df[cook.status == "Accept", ]
19

```

The estimated TmS values for TCGA PRAD tumor samples are shown in the violin plot below.



# Reference

[1] <https://github.com/wwylab/DeMixT>

[2] Cao, S. et al. Estimation of tumor cell total mRNA expression in 15 cancer types predicts disease progression. Nat Biotechnol (2022). <https://doi.org/10.1038/s41587-022-01342-x>.

