

MSstatsPTM Simulation Analysis

Devon Kohler

3/7/2022

```
library(MSstatsPTM)
library(data.table)
```

```
## Warning: package 'data.table' was built under R version 4.1.2
```

```
library(tidyverse)
```

```
## -- Attaching packages ----- tidyverse 1.3.1 --
```

```
## v ggplot2 3.3.5    v purrr  0.3.4
## v tibble  3.1.6    v dplyr  1.0.8
## v tidyr   1.2.0    v stringr 1.4.0
## v readr   2.1.2    v forcats 0.5.1
```

```
## Warning: package 'tibble' was built under R version 4.1.2
```

```
## Warning: package 'tidyr' was built under R version 4.1.2
```

```
## Warning: package 'readr' was built under R version 4.1.2
```

```
## Warning: package 'dplyr' was built under R version 4.1.2
```

```
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::between()   masks data.table::between()
## x dplyr::filter()    masks stats::filter()
## x dplyr::first()     masks data.table::first()
## x dplyr::lag()       masks stats::lag()
## x dplyr::last()      masks data.table::last()
## x purrr::transpose() masks data.table::transpose()
```

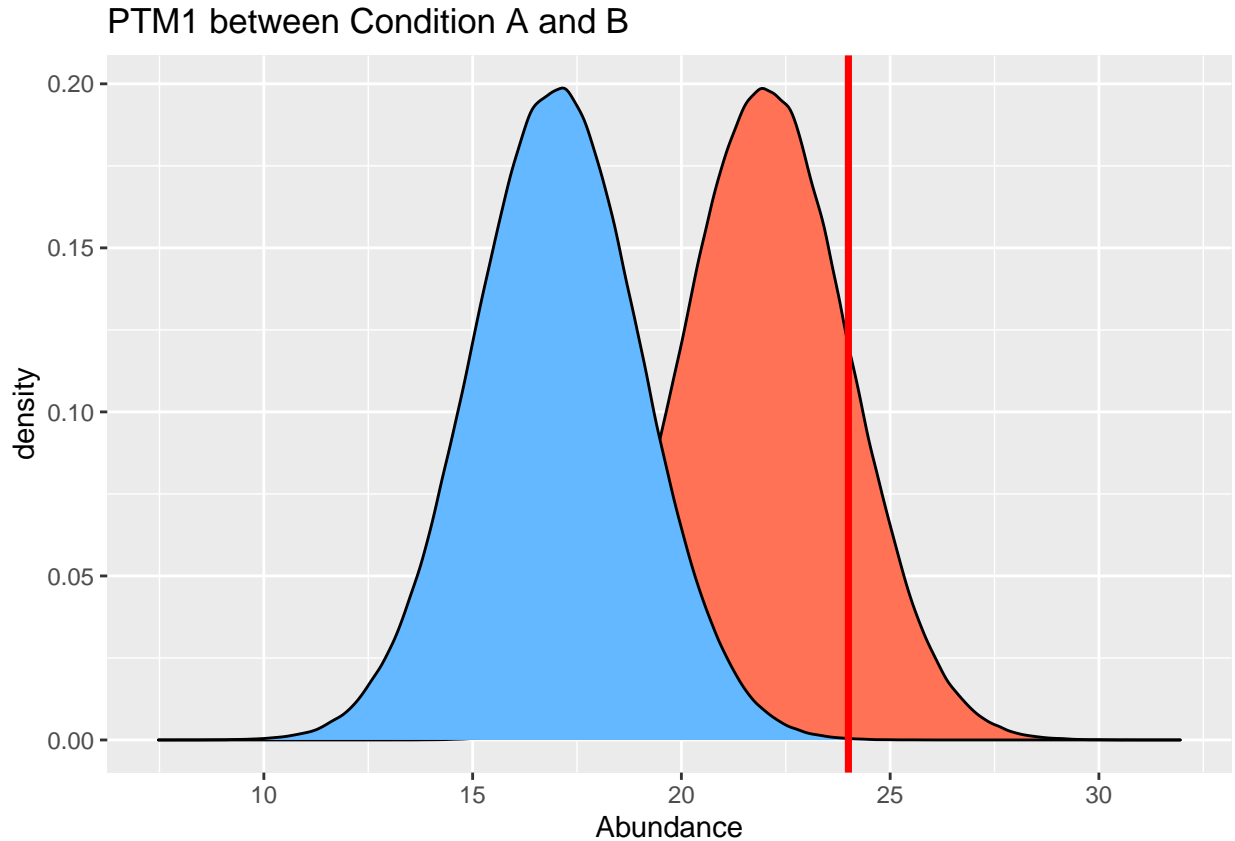
Computer Simulation

Simulation Methods

To simulate data the `PTMsimulateExperiment()` function was used. This function allows us to vary the number of conditions, replicates, number of proteins, number of sites per protein, number of spectral features per site/protein, mean log₂-abundance of PTM and PROTEIN, deviation from the mean log₂-abundance in each group, standard deviation among replicates, and standard deviation among log₂-intensities.

Three different statistical modeling methods were applied to the simulated data: MSstatsPTM, limma, and two-sample t-test. These methods were tested both with and without applying protein level adjustment. MSstatsPTM uses TMP for summarization and post modeling calculations for the adjustment. To adjust limma and t-test, the run-level data was averaged for both ptm and protein datasets and then combined. The resulting dataset was then used for limma and t-test.

Graphs can help us visualize the components of the simulation.



Here we can see a sample PTM with two conditions A and B. The red line represents a same biological replicate. We can simulate different variance's for the conditions.

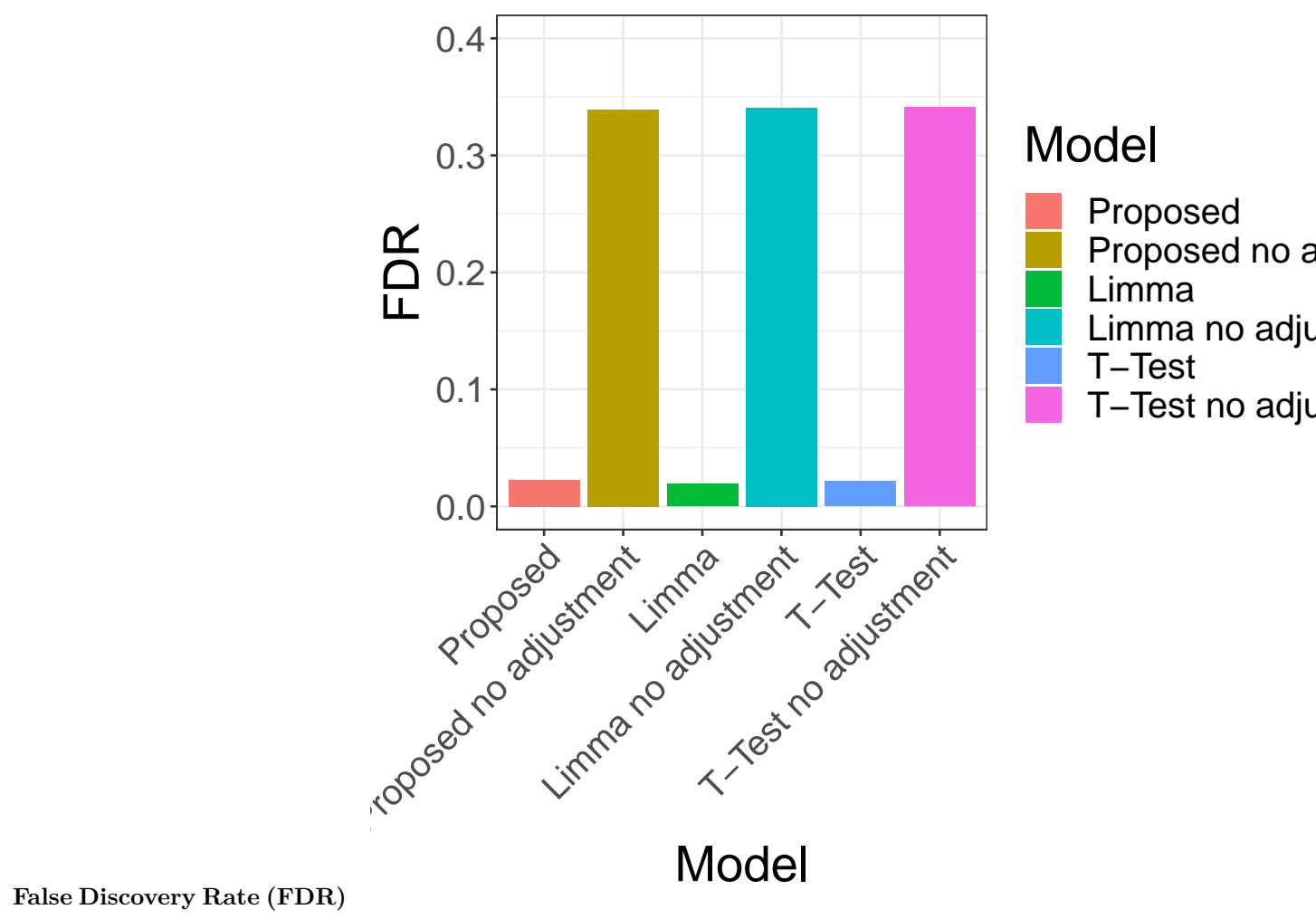
All charts are made using adjusted pvalue $< .05$ to designate a significant hit.

Simulation 1

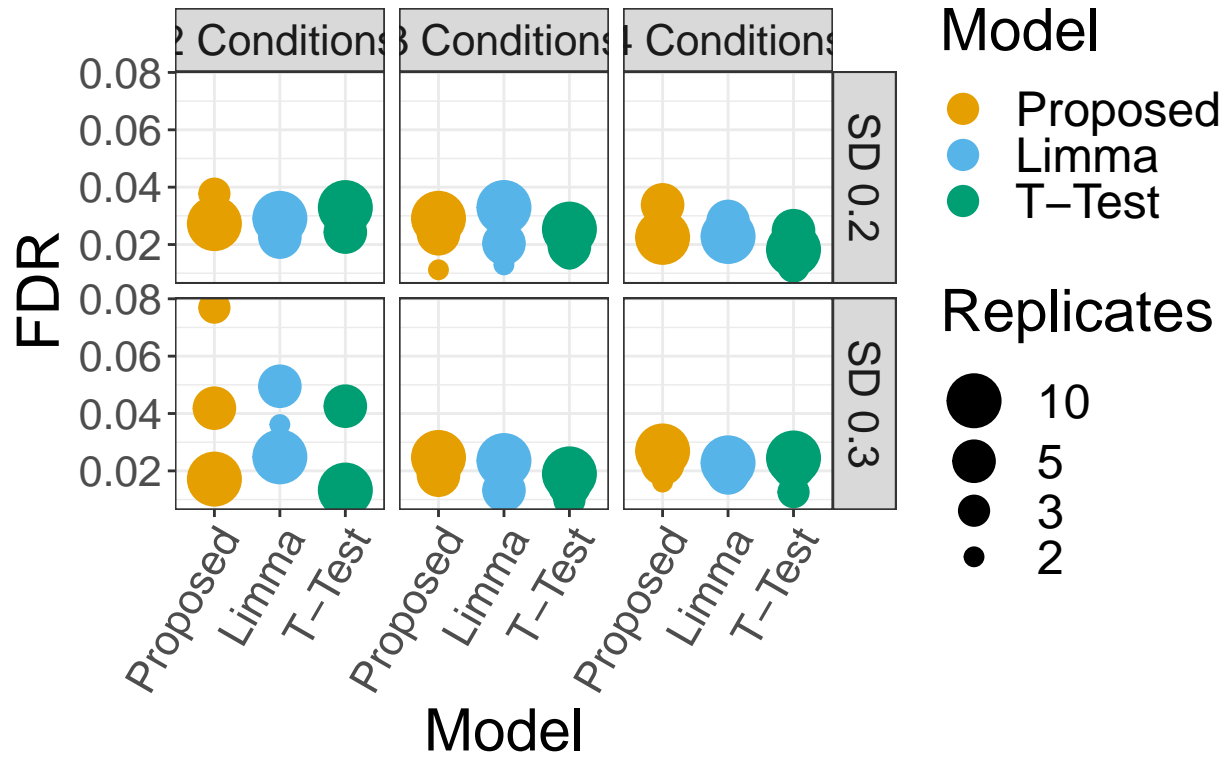
All simulations were ran with half the PTMs being differential, while the other half the difference was due to changes in global protein level. The first simulation is run with the following parameters:

- Mean of log-intensity: 25
- Number of Features: 10 (PTM), 10 (Protein)
- Standard deviations of log-intensities for modified and unmodified peptides: 0.2, 0.3
- Difference in PTM abundance between conditions: 0, 0.5, 0.75, 1
- Difference in protein abundance between conditions: 0, 0.5
- Number of replicates: 3, 4, 5
- Number of conditions: 2, 3, 4
- Number of realizations: 500

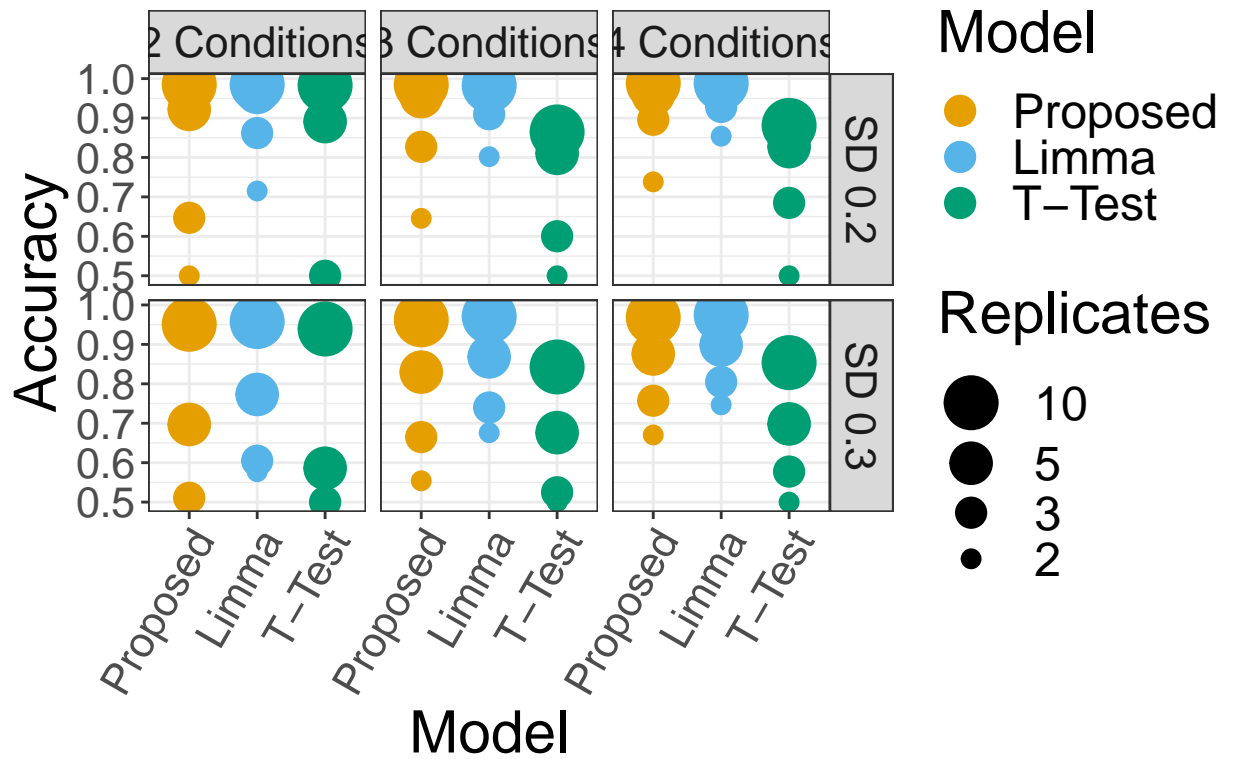
Simulation 1: FDR all models



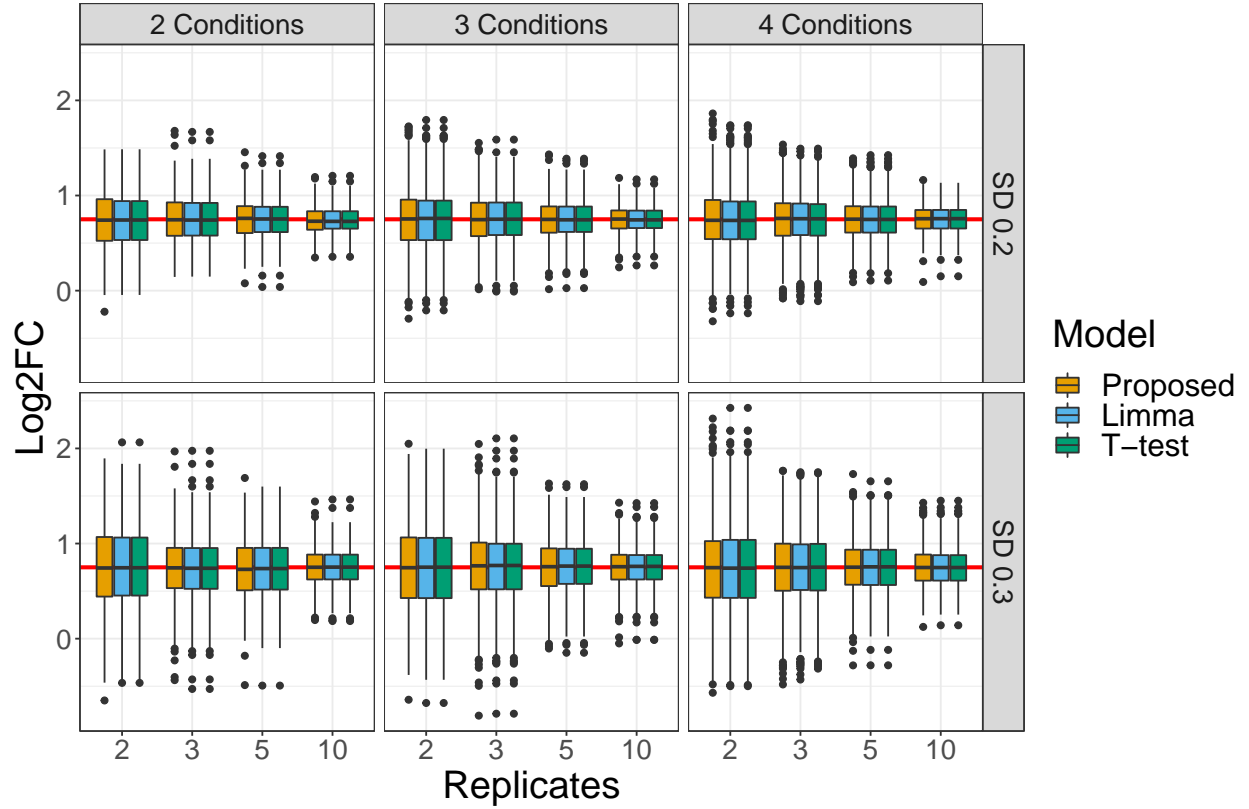
Simulation 1: FDR adjusted models



Simulation 1: Accuracy adjusted mo



Simulation 1: Fold Change Distribution

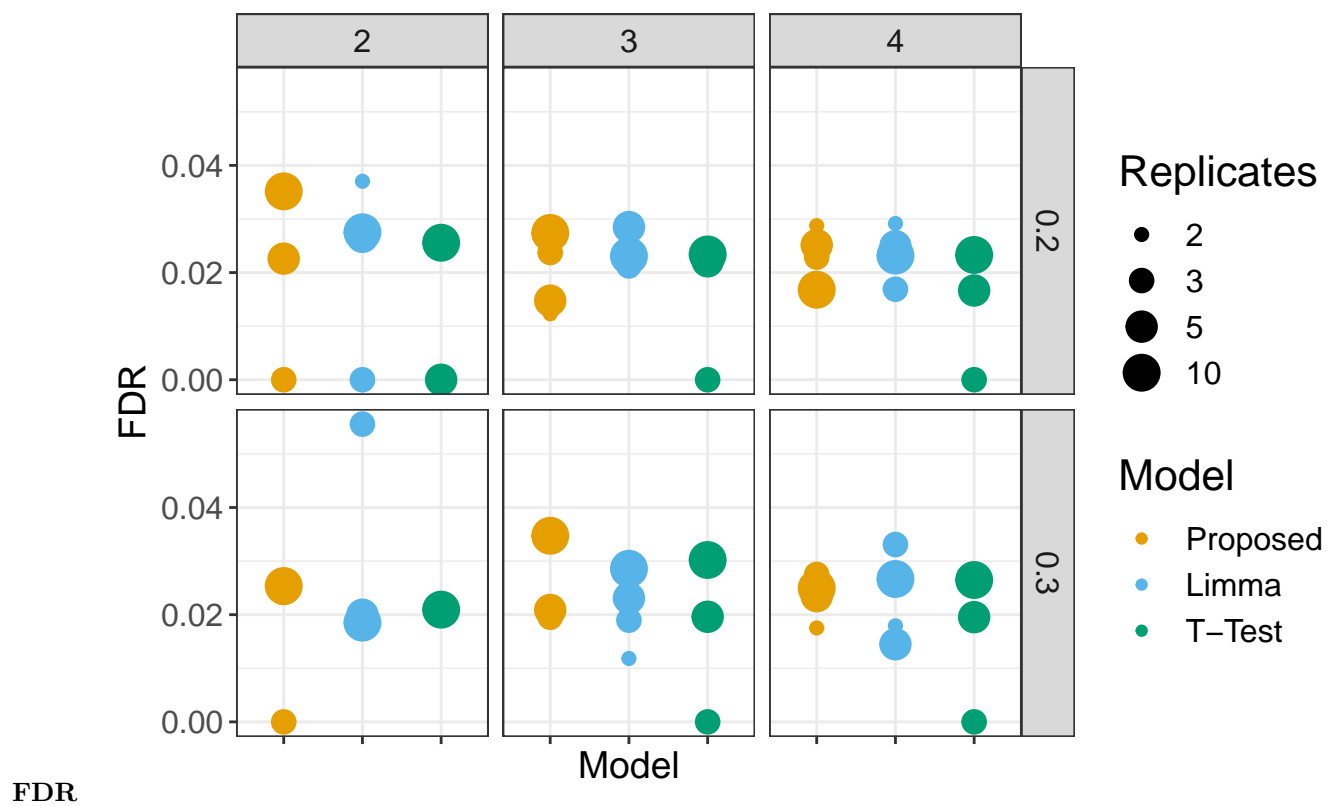


Simulation 2

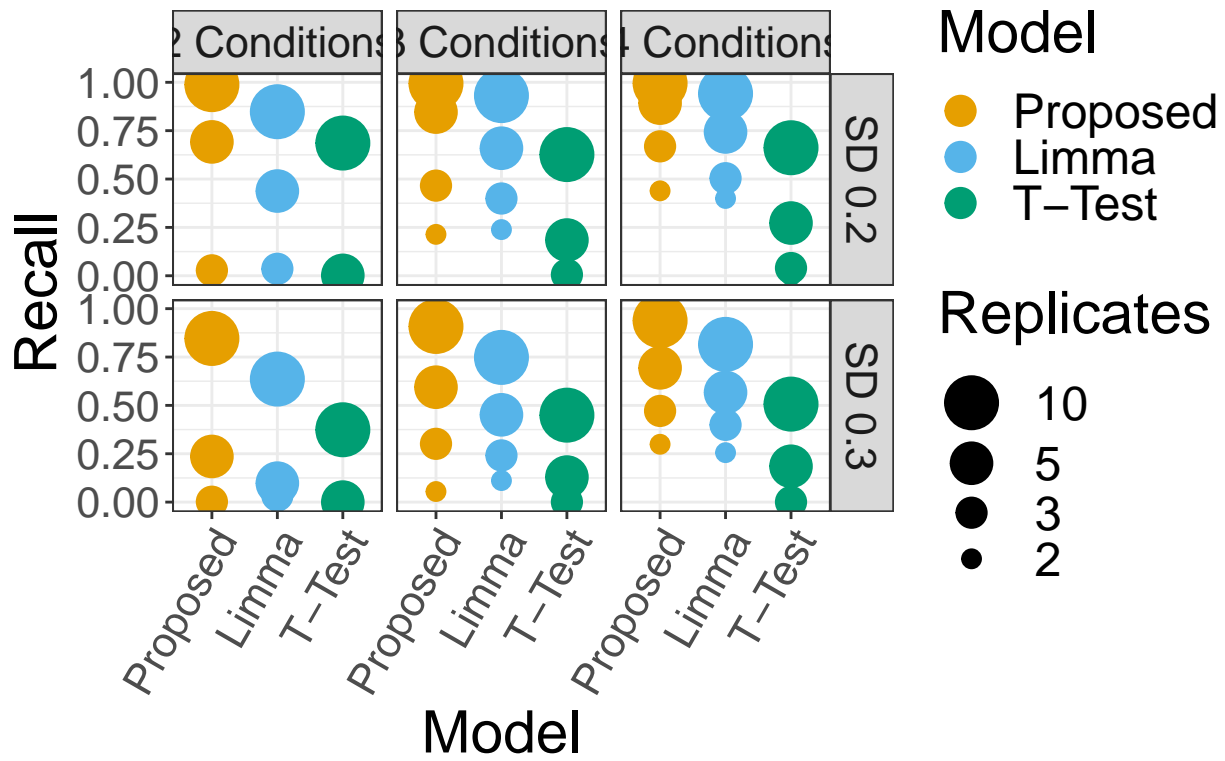
In this simulation we introduce missing values into the experiment. 20% of the features simulated were selected at random and masked with an NA value. No missing value imputation was used.

- Mean of log-intensity: 25
- Number of Features: 2 (PTM), 10 (Protein)
- Standard deviations of log-intensities for modified and unmodified peptides: 0.2, 0.3
- Difference in PTM abundance between conditions: 0, 0.75, 1.5, 2.25
- Difference in protein abundance between conditions: 0, 0.75, 1.5, 2.25
- Number of replicates: 2, 4, 6
- Number of conditions: 2, 3, 4
- Number of realizations: 1000
- 10% of features missing at random in PTM and Protein datasets

FDR adjusted models

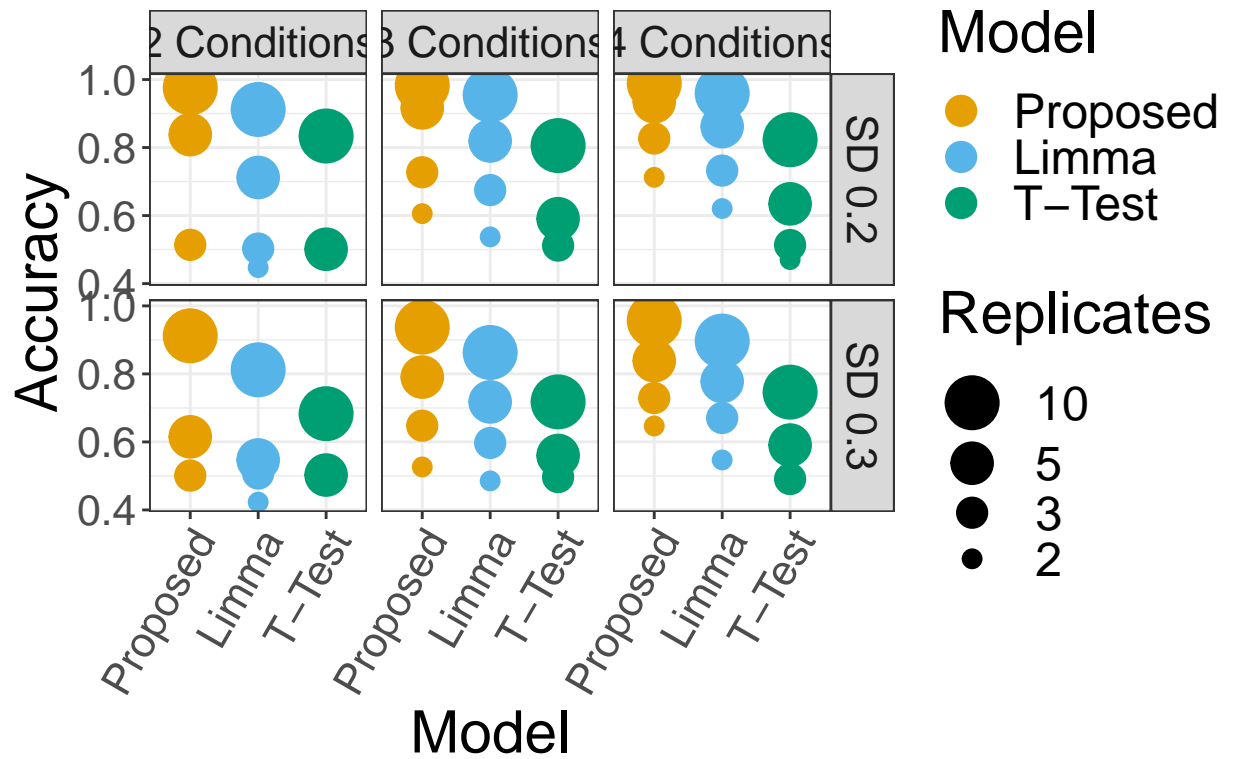


Simulation 2: Recall adjusted model



Recall (TPR)

Simulation 2: Accuracy adjusted mo



Simulation 2: Fold Change Distribution

