

Permutation and scalar transformation approaches to normalisation and combining them

Sayani Gupta

1 Data generation

Observations are generated from a $\text{Gamma}(2,1)$ distribution for each combination of nx and $nfacet$ from the following sets: $nx = nfacet = \{2, 3, 5, 7, 14, 20, 31, 50\}$ to cover a wide range of levels from very low to moderately high. Each combination is being referred to as a *panel*. That is, data is being generated for each of the panels $\{nx = 2, nfacet = 2\}, \{nx = 2, nfacet = 3\}, \{nx = 2, nfacet = 5\}, \dots, \{nx = 50, nfacet = 31\}, \{nx = 50, nfacet = 50\}$. For each of the 64 panels, $ntimes = 500$ observations are drawn for each combination of the categories. That is, if we consider the panel $\{nx = 2, nfacet = 2\}$, 500 observations are generated for each of the combination of categories from the panel, namely, $\{(1, 1), (1, 2), (2, 1), (2, 2)\}$. The values of λ is set to 0.67 and values of raw wpd wpd_{raw} is obtained.

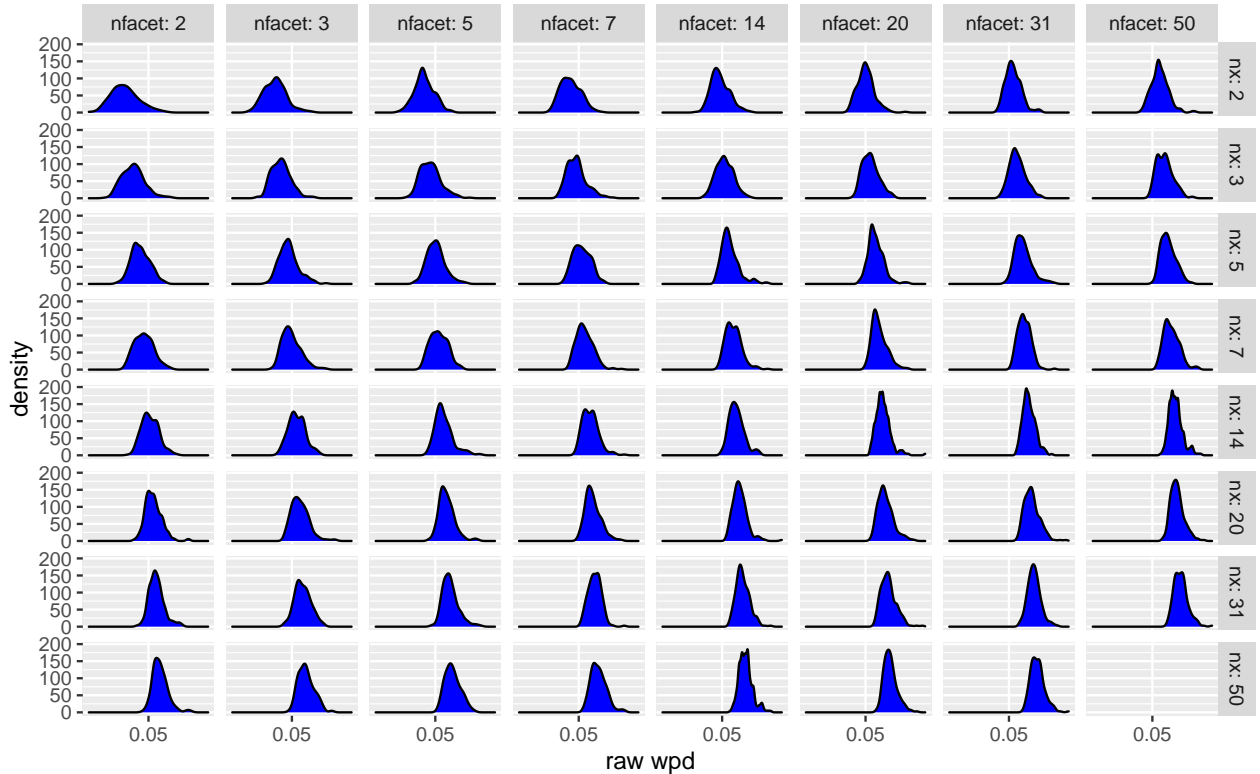


Figure 1: Distribution of raw wpd is plotted across different nx and $nfacet$ categories. Both shape and scale of the distribution changes for different nx and $nfacet$ categories.

Figure ?? shows the distribution of wpd_{raw} is plotted across different nx and $nfacet$ categories. Both shape and scale of the distribution changes for different nx and $nfacet$ categories. This is not desirable as it would

mean we would not be able to compare wpd_{raw} across different nx and $nfacet$ as each of them are drawn from distributions of different locations and scale. In Figure 2 we see how the median of wpd_{raw} varies with the total number of distances $nx_i * nfacet_j$ for a panel with $nx = nx_i$ categories and $nfacet = nfacet_j$ categories. The median increases abruptly for lower values of $nx_i * nfacet_j$ and the increase becomes more gradual for higher $nx_i * nfacet_j$.

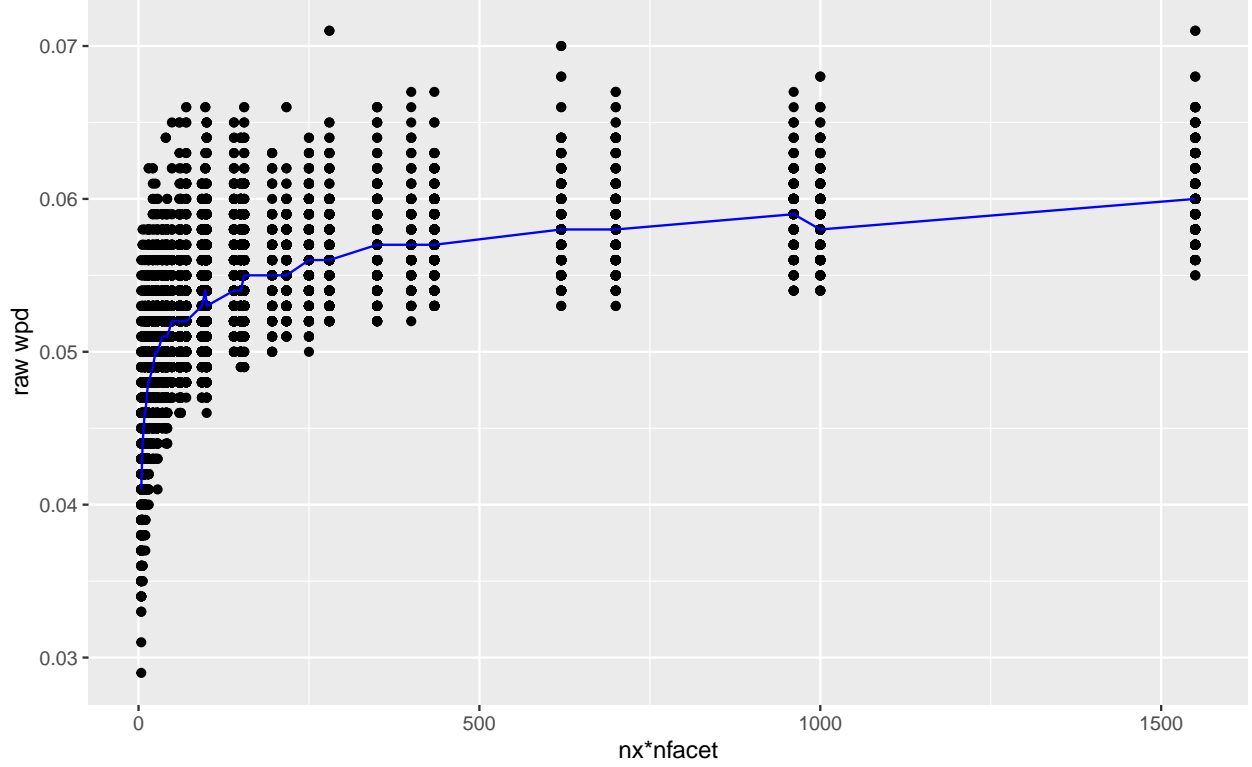


Figure 2: The raw wpd is plotted against $nxnfacet$ and the blue line represents the median of the multiple values for each $nxnfacet$ levels.

We need a transformation on raw_{wpd} which will make it independent of the values of $nx * nfacet$. Two approaches have been employed for that purpose, the first one involves fitting a model and the latter involves a permutation method to make the distribution of the transformed wpd similar across different nx and $nfacet$.

1.1 Scalar transformation approach to normalisation

1.1.1 Linear model

A log-linear model is fitted to see how the values of wpd_{raw} changes with the values of nx and $nfacet$. The model is of the form

$$y = a + b * \log(x) + e$$

, where $y = median(wpd_{raw})$ and $x = nx * nfacet$. wpd_{lm} is a transformation on wpd_{raw} which should be designed to remove the effect of $nx * nfacet$ on wpd_{raw} and thus is defined as follows: $wpd_{lm} = wpd_{raw} - a - b * \log(nx * nfacet)$

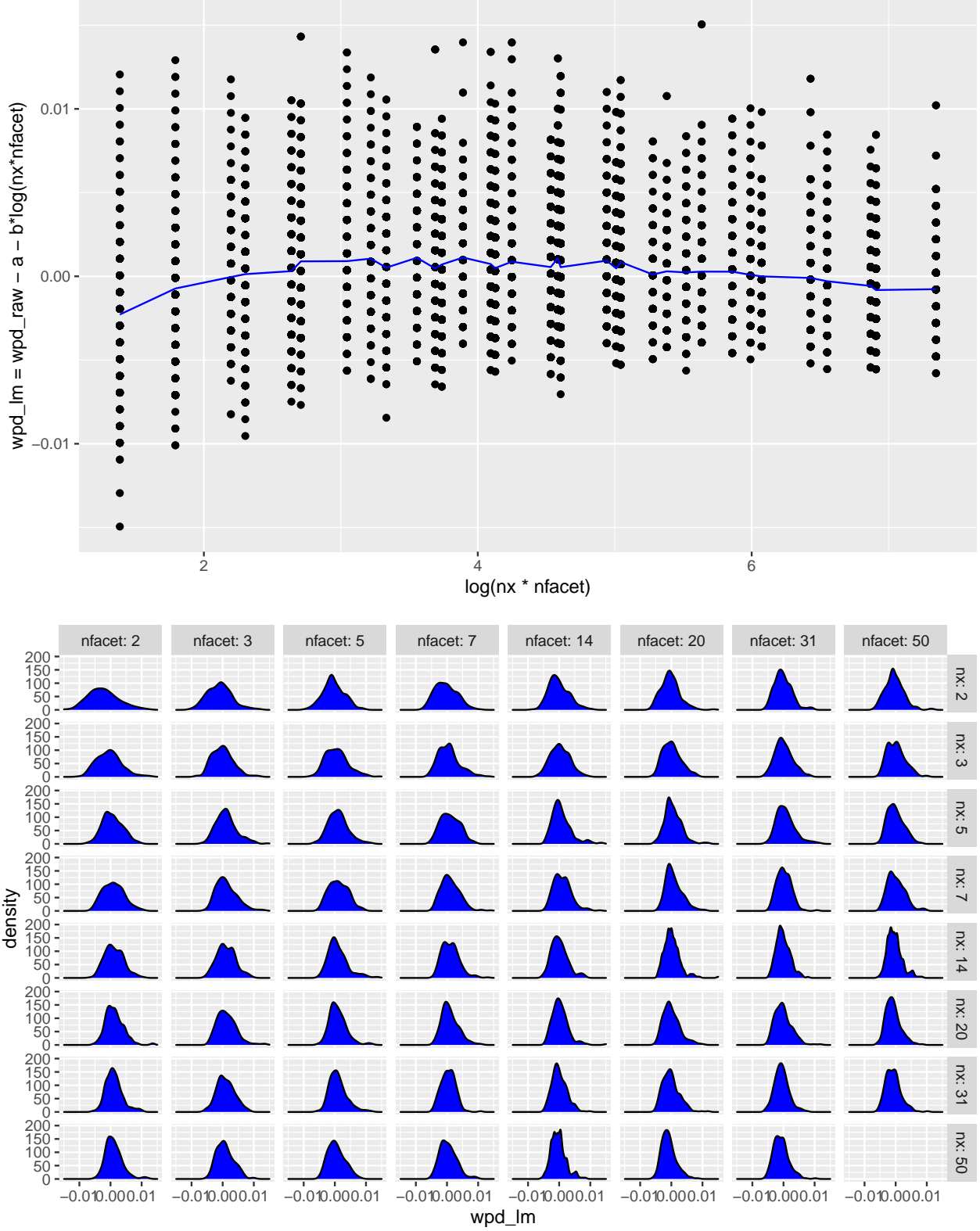


Figure 3: The distribution of wpd_{lm} is plotted. The distributions are more similar across higher nx and $nfacet$ and are different for smaller nx and $nfacet$.

1.1.2 Generalised linear model

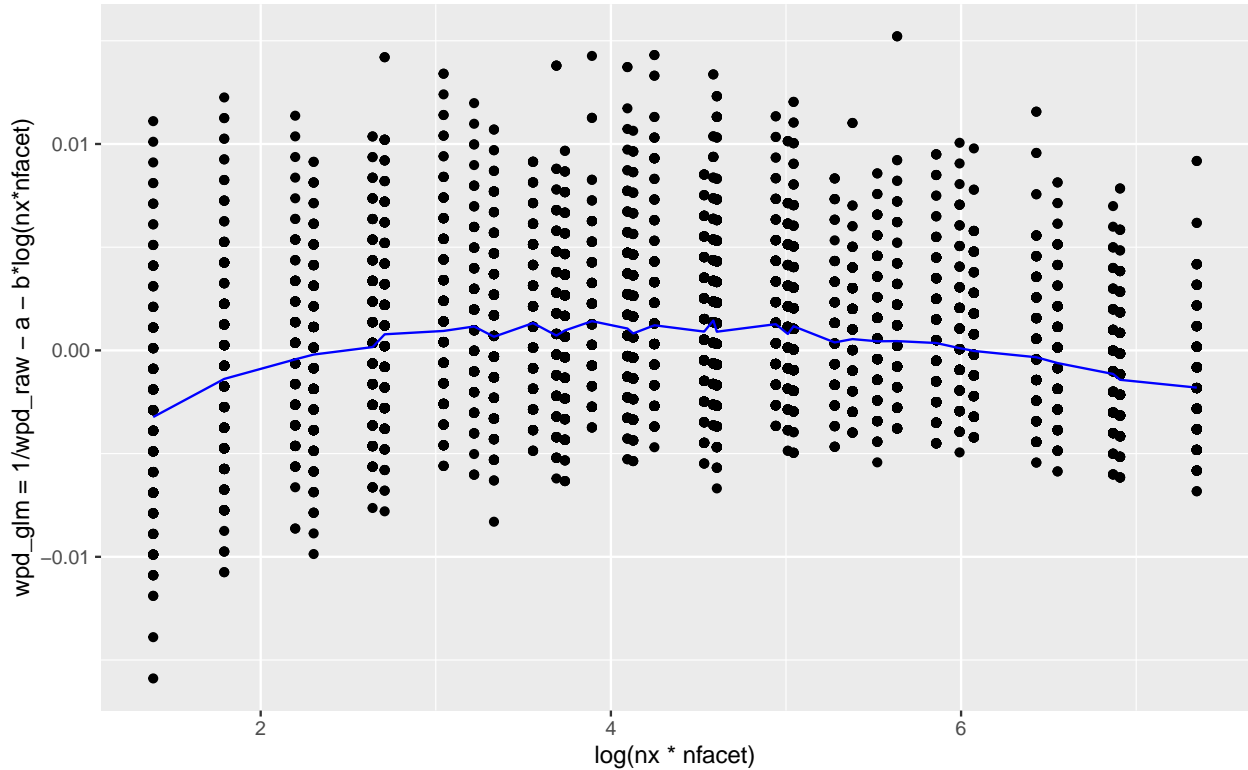
In the earlier approach, $wpd_{raw} \in R$, whereas, $0 \leq wpd_{raw} \leq 1$ since it is a JS distance. Also, JSD follows a Chi-square distribution, which could be considered as a Gamma distribution. Hence, to restrict the range of the measure between 0 and 1, and deviate from the normality assumption of the response variable in a OLS approach, we fit a generalized linear model (GLM) with the error distributed as a gamma distribution and link function as “inverse”. Hence, the model is of the form

$$1/y = a + b * \log(x) + e$$

, where $y = median(wpd_{raw})$ and $x = nx * nfacet$. Again, wpd_{glm} is a transformation on wpd_{raw} which should be designed to remove the effect of $nx * nfacet$ on wpd_{raw} and thus is defined as follows: $wpd_{glm} = 1/wpd_{raw} - a - b * \log(nx * nfacet)$.

Please note:

- subtracting the intercept brings the location of the transformed variable to 0 for either case
- heterogeneity is higher for this case after normalization as compared to the earlier one.



1.2 Permutation approach to normalisation

The simulated data for each of the panels is permuted/shuffled $nperm = 200$ times and for each of those permutations wpd_{norm} is computed as follows: $wpd_{perm} = (wpd_{raw} - mean(wpd_{raw}))/sd(wpd_{raw})$. This is done so that the distribution of the normalised measure wpd_{norm} has the same mean and standard deviation across different nx and $nfacet$.

Please note that standardizing the variable wpd_{perm} in this approach leads to $location = 0$ and $scale = 1$ for this variable.

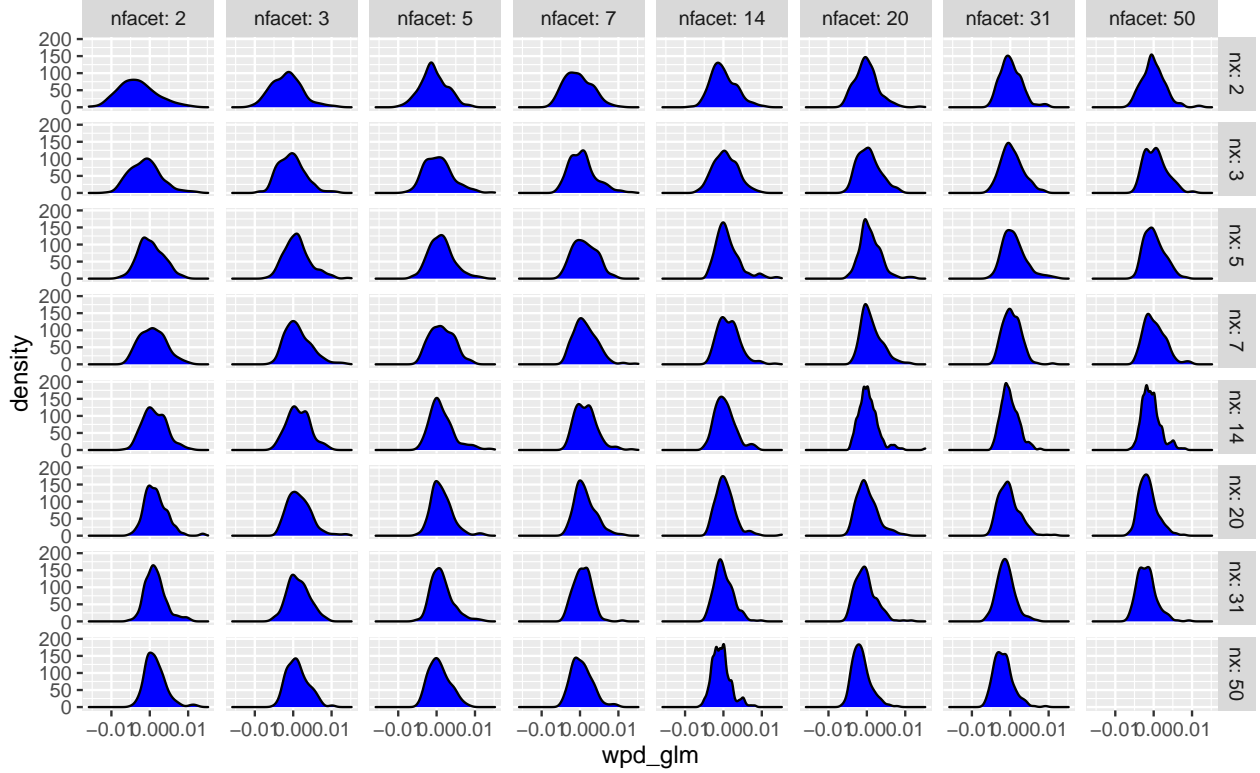


Figure 4: The distribution of wpd_{glm} is plotted. The distributions are more similar across higher nx and $nfacet$ and dissimilar for fewer nc and $nfacets$.

1.3 Bringing them both to the same scale

We see that the transformation through the modeling approach leads to very similar distribution across high nx and $nfacet$ (higher than 7) and not so much for lower nx and $nfacet$. Hence, the computational load of permutation approach could be alleviated by using the modeling approach for the higher nx and $nfacet$, however, it is important that we use the permutation approach for lower nx and $nfacet$. However, it is difficult to compare the transformed wpd from both of these approaches, since each of the variables is measured on a different scale (each of them have location 0). The transformed variables from the two approaches could be brought to the same scale so that for smaller categories, permutation approach is used and for larger categories, we can stick to modeling approach. These could be done through the following:

- Making the range of both the variables same by using min-max scaling method. In practice, however, we would only have one value of wpd_{raw} which we need to transform using the modeling approach. Hence, min-max scaling approach could not be used here.
- Standardizing the variables and expressing scores at standard deviation units. Again in practice, however, we would only have one value of wpd_{raw} which we need to transform using the modeling approach. Hence, standardizing scores could not be used here as we do not have the mean and standard deviation of a series while using transformation using modeling.
- Make the location and scale of both the approaches similar so that they could be compared. Please note that the range of values could be different in this case, however location and scale are brought to same levels.)

The measure wpd_{glm} seems to roughly follow a normal distribution except in the tails as could be seen in Figure 6 and the very method of permutation approach ensures that $wpd_{permutation}$ is also normally distributed. Further, they are brought to the similar scale and location and hence could be compared.

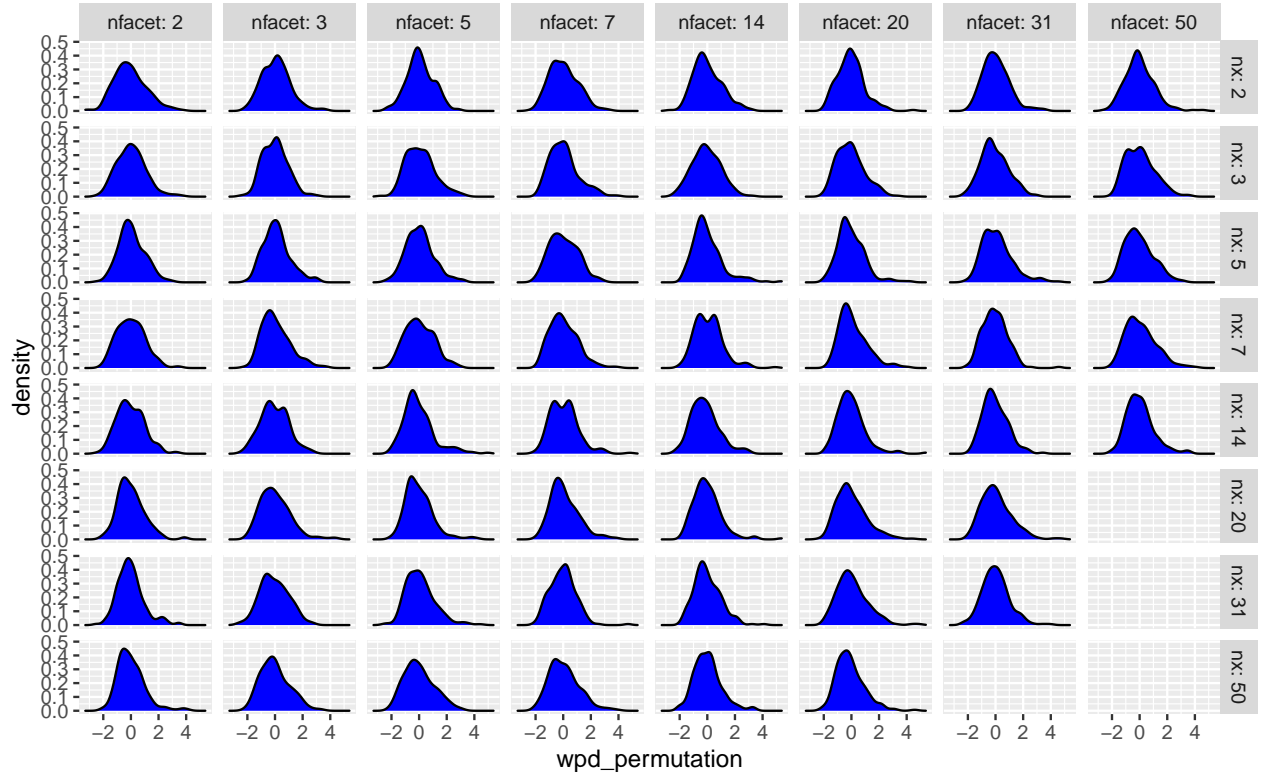


Figure 5: The distribution of $wpd_{permutation}$ is plotted. The distributions are more similar across different nx and $nfacet$ (specially for small nx and $nfacet$) but this approach has the downside of more computational time.

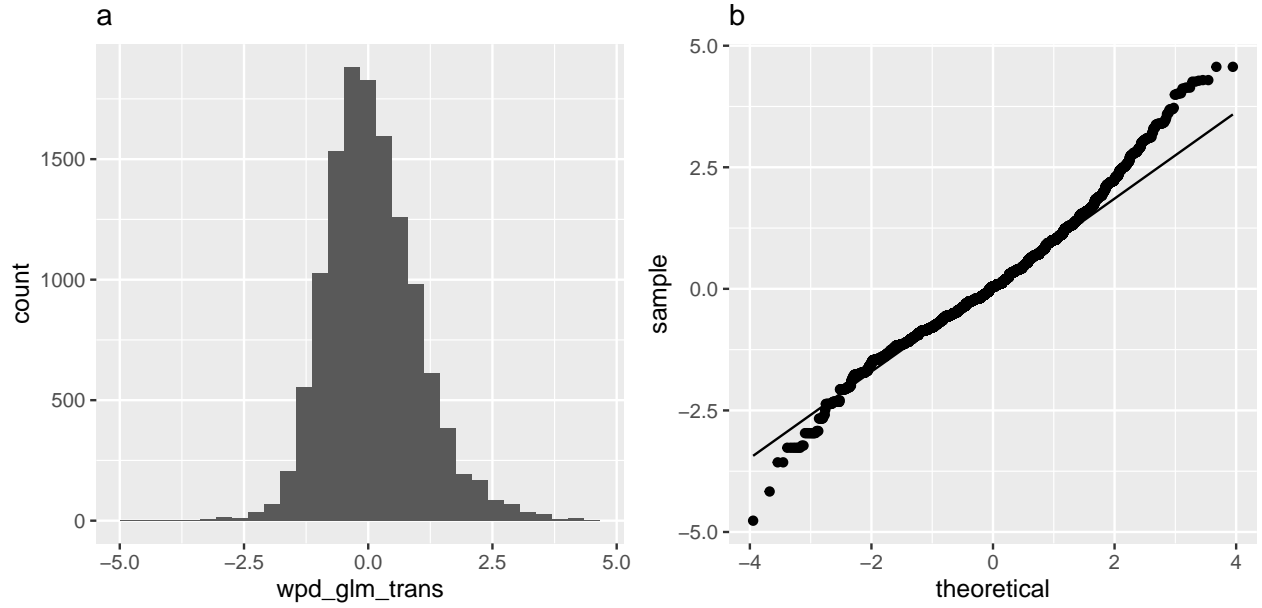


Figure 6: In panel a, the histogram of $wpd_{glm_{trans}}$ is plotted. In panel b, the QQ plot is shown with the theoretical quantiles on the x-axis and $wpd_{glm_{trans}}$ quantiles on the y-axis. The distribution looks symmetric and looks like normal except in the tails.

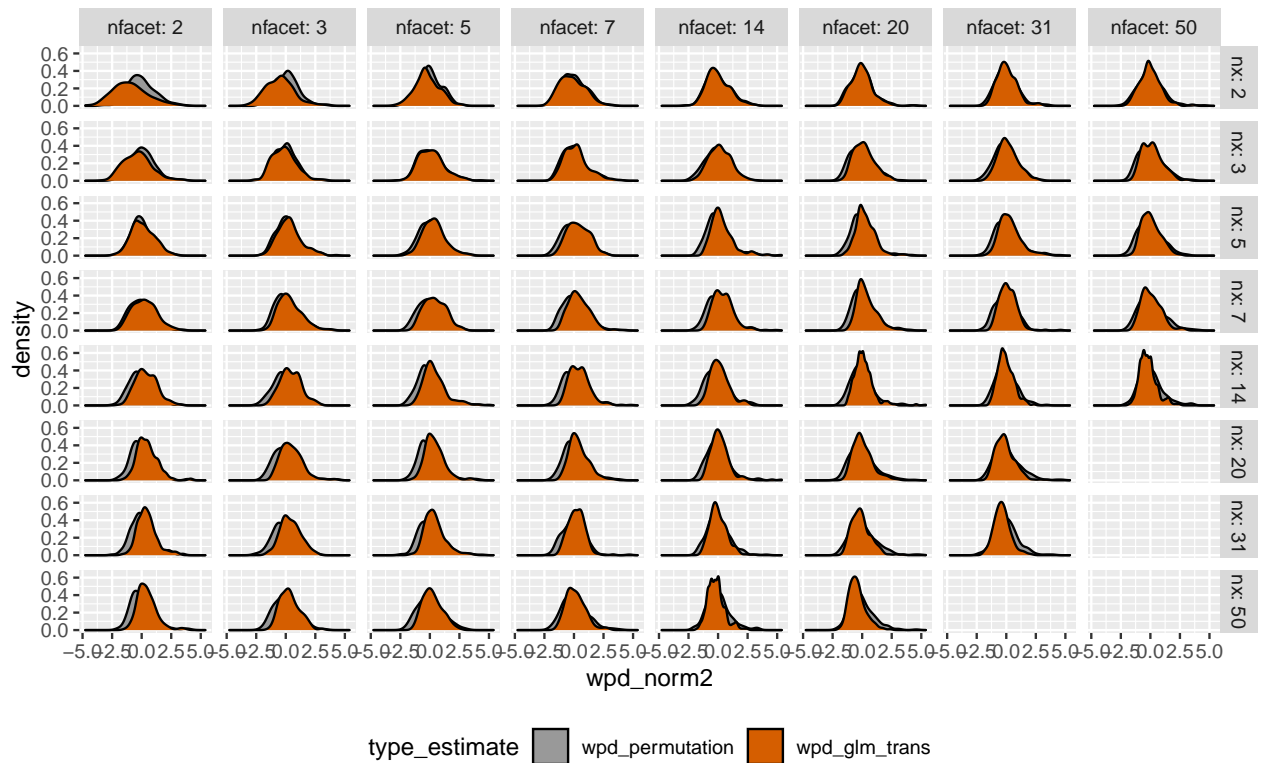


Figure 7: `wpd_permutation` and `wpd_lm` are in very different scale and could not be plotted together without transforming further. Whereas `wpd_permutation` and `wpd_glm` could be plotted together as they at least have the same location. Are they on the same scale though? Can they be compared?

Questions:

1. Does the approach of wpd_{glm} looks correct?
2. Are wpd_{glm} and $wpd_{permutation}$ be compared from this plot or both should be brought to the scale 1 for comparison?
3. Forcing same range (0,1) to both wpd_{glm} and $wpd_{permutation}$ could be obtained by using the transformation $(z - z_{min}) / (z_{max} - z_{min})$. This could be used for the permutation approach. But the modeling approach would only have one value of wpd_{raw} for a panel in practice, how to scale the values in the modeling approach so that they are within 0 and 1?

[1] 0.003051935