

Detecting dispersion problems in GLMMs

Melina de Souza Leite

melina.souza-leite@ur.de

Postdoctoral Researcher

University of Regensburg

Ecological data analysis

Increased complexity and flexibility in ecological data modeling:

- Generalized linear modes (GLMs)
- Mixture models (e.g. zero-inflated GLMs)
- Hierarchical/Multilevel models, GLMMs



- But still few tools for model diagnostics
- Problem: **failing to check model assumptions**

Can you trust your model?

Dispersion problems in count data

- Example count data:
 - Species richness
 - Abundance of individuals
 - Number of success (K) within a number of trials (N)
- Modeling count data, GL(M)M distributions:
 - Poisson
 - Binomial (K/N) proportion

Problem when data has **more or less variability than expected by the distribution used for modeling:**

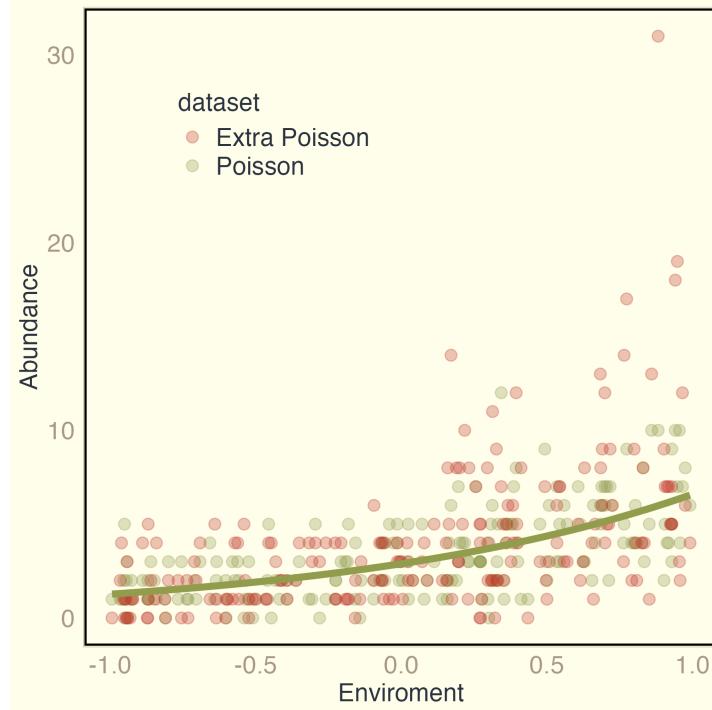
UNDER or OVERDISPERSION

GOALS

- Aware ecologists of dispersion problems with count data
- Identify and describe the 3 main causes by using model diagnostics tools with DHARMA
- Show modeling solutions for these causes

3 causes of dispersion problems

“Real” overdispersion:



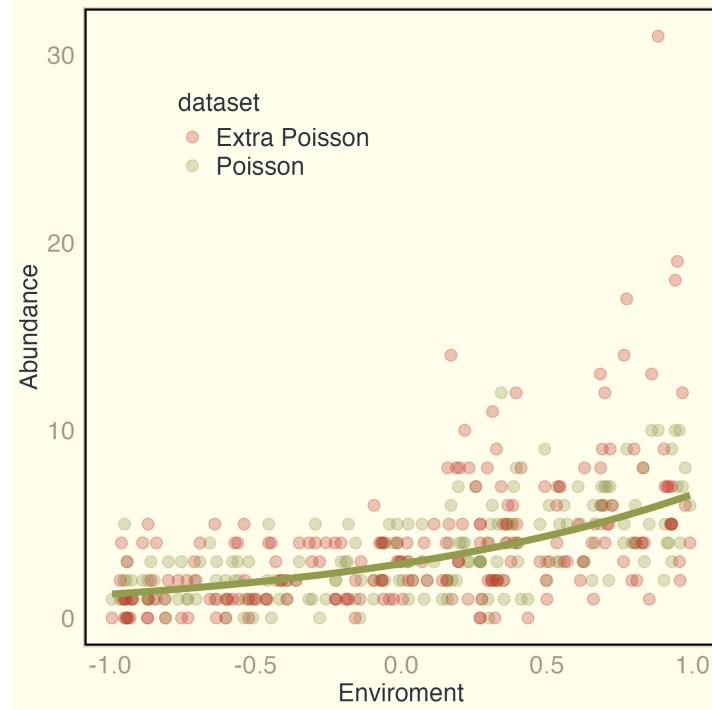
Heteroscedasticity:

Zero-inflation:

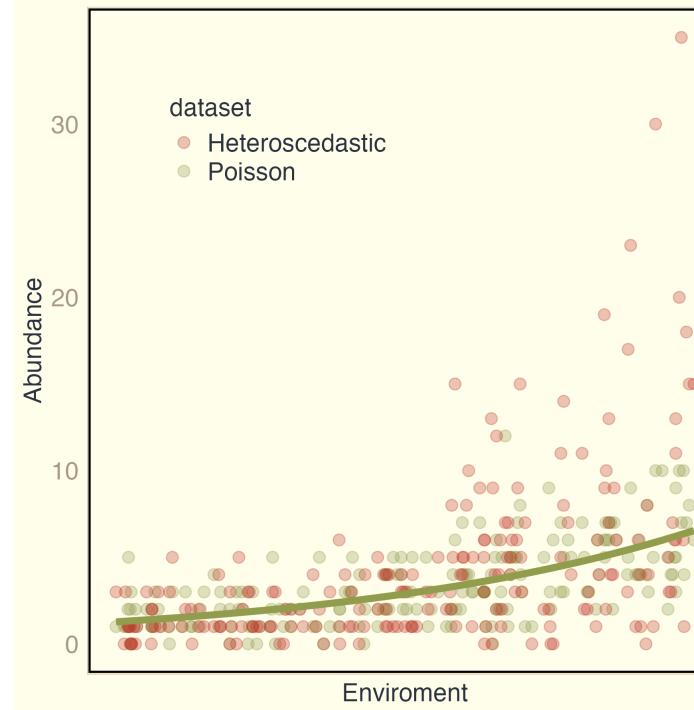
More variance than expected by the model.

3 causes of dispersion problems

“Real” overdispersion:



Heteroscedasticity:



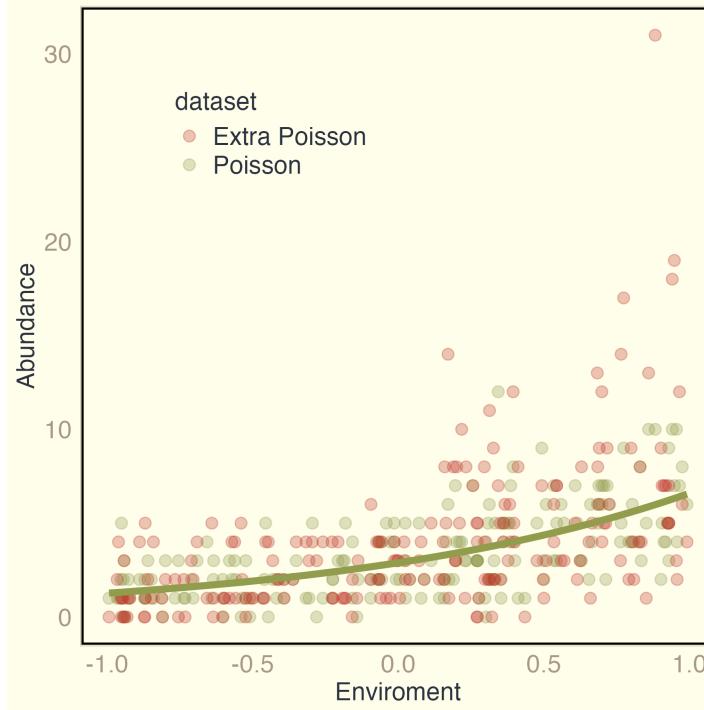
Zero-inflation:

More variance than expected by the model.

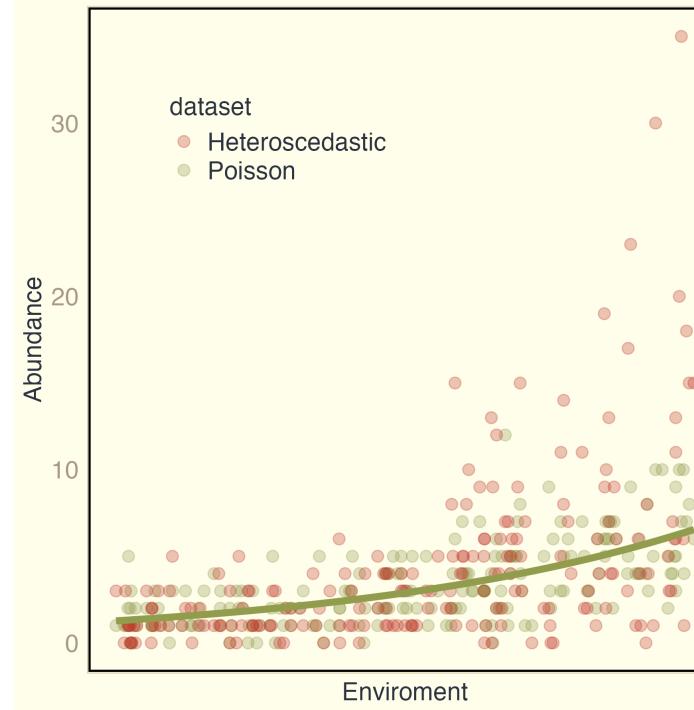
Variance increases/ decreases with a predictor.

3 causes of dispersion problems

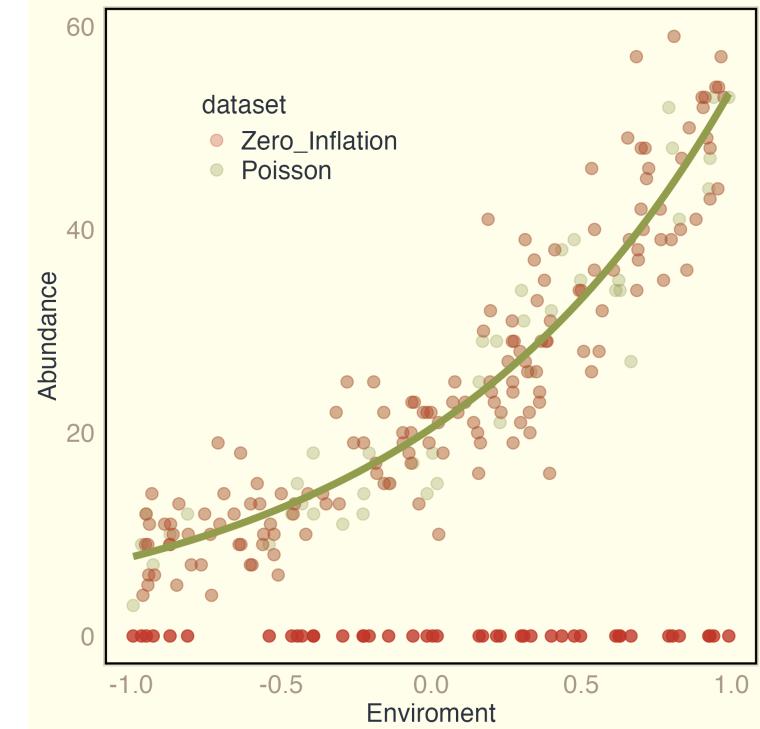
“Real” overdispersion:



Heteroscedasticity:



Zero-inflation:



More variance than expected by the model.

Variance increases/ decreases with a predictor.

More zeros than expected by the model.

Consequences of dispersion problems

- Too small standard error of estimates -> narrower confidence intervals
- Larger chance of type I error: find an effect when it doesn't exist
- Wrong estimates by ignoring other processes (e.g. zero-inflation causes) in your data-generating process.
- Missing the opportunity to learn and get more info from your data/system. Ecological meanings for modeling/understanding unexpected variability?

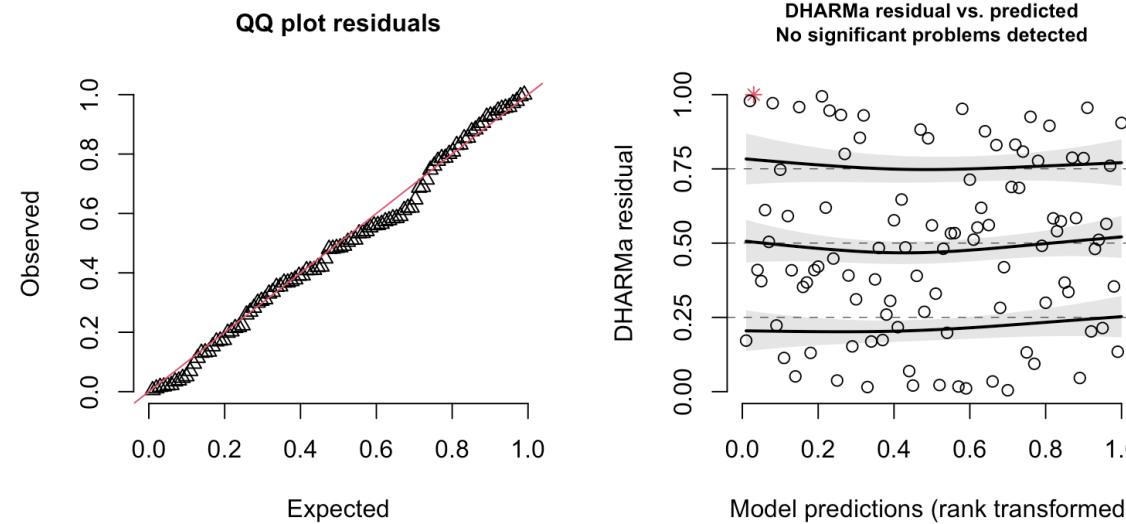


Detecting dispersion problems with DHARMa

Residual diagnostics with DHARMa

- Scaled quantile residuals -> Simulating from the model
- Residuals between 0 and 1 for ANY model complexity or distribution
- Interpreted the SAME way:

If your model is correctly specified, i.e. you have the “data-generating process”, scaled quantile residuals will present a uniform “flat” distribution between 0 and 1.



Detecting dispersion problems: DHARMa

Create DHARMa residuals

```
1 DHARMAResiduals <- simulateResiduals(model)
```

Test dispersion problems

```
1 testDispersion(DHARMAResiduals)
```

Test heteroscedasticity

```
1 plotResiduals(DHARMAResiduals,
2                   form = data$Environment1, # the predictor
3                   absoluteDeviation = T)
```

Test zero inflation

```
1 testZeroInflation(DHARMAResiduals)
```

Solving dispersion problems: glmmTMB

Overdispersion

```
1 glmmTMB(observedResponse ~ Enviroment1,  
2           family = nbinom2(), ...) # from poisson()
```

Heteroscedasticity

```
1 glmmTMB(observedResponse ~ Enviroment1,  
2           dispformula = ~ Enviroment1, # dispersion formula  
3           family = nbinom2(), ...) # needs to be negative binomial
```

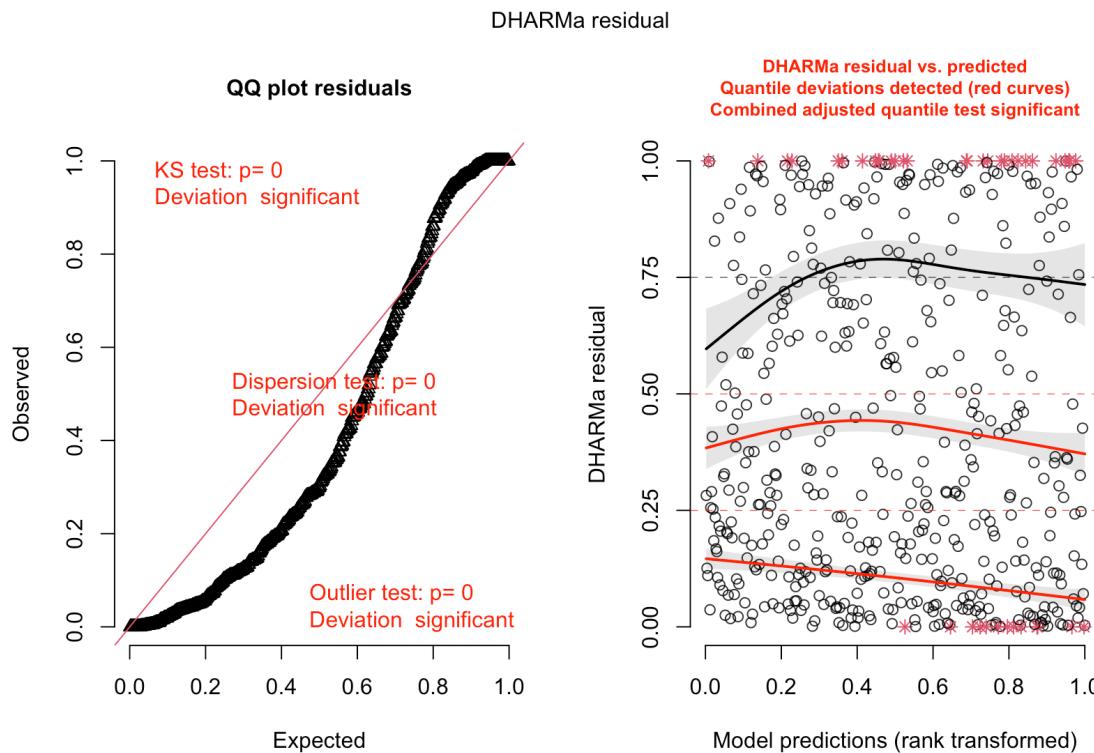
Zero-inflation

```
1 glmmTMB(observedResponse ~ Enviroment1,  
2           ziformula = ~ 1, # zero-inflation formula / can add predictor  
3           family = poisson(), ...) # can also be negative binomial
```



What is the problem?

Problem 1



```
1 testDispersion(overRes,plot=F)
```

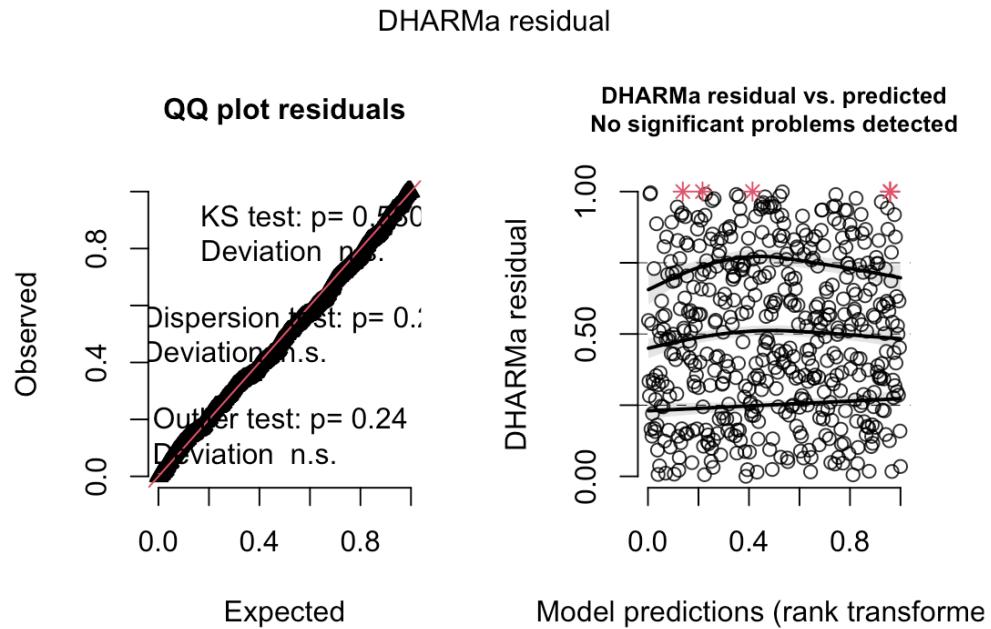
DHARMA nonparametric dispersion test via sd of residuals
fitted vs.
simulated

data: simulationOutput
dispersion = 5.2662, p-value < 2.2e-16
alternative hypothesis: two.sided

Try a negative binomial model

Modeling “real” overdispersion

```
1 glmmTMB(observedResponse ~ Environment1 + (1|group),  
2 family = nbinom2(), data = overData)
```



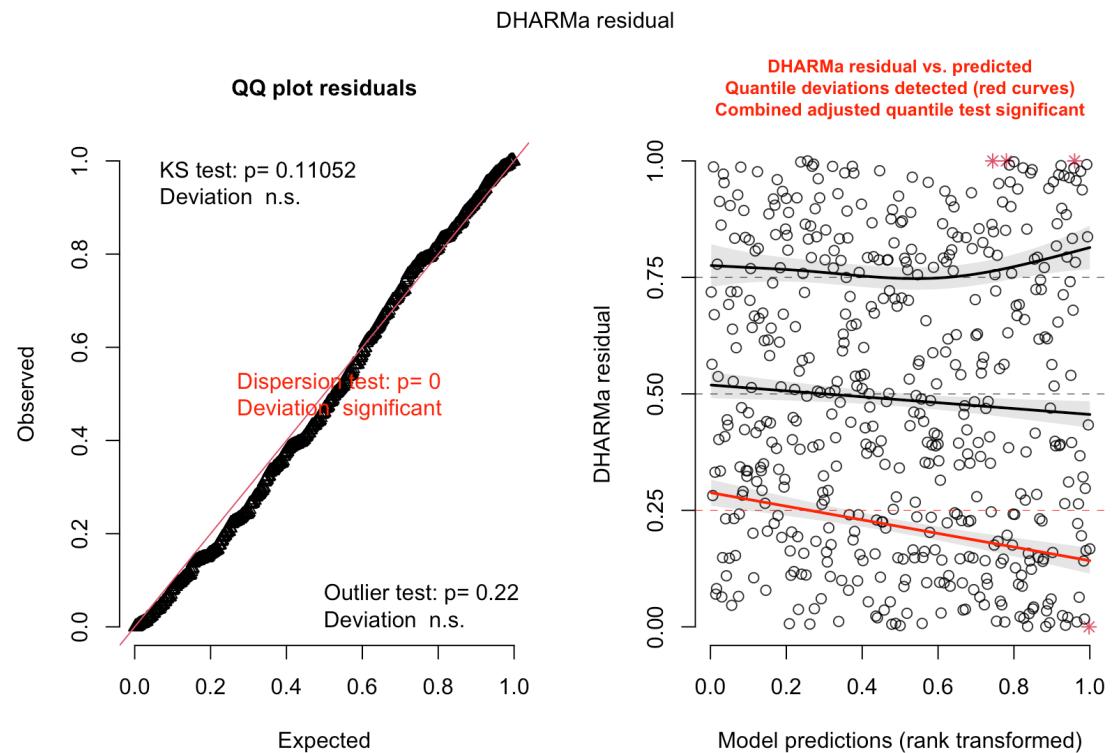
DHARMA nonparametric dispersion test via sd of residuals fitted vs.
simulated

```
data: simulationOutput  
dispersion = 1.1935, p-value = 0.224  
alternative hypothesis: two.sided
```

DHARMA zero-inflation test via comparison to
expected zeros with
simulation under H0 = fitted model

```
data: simulationOutput  
ratioObsSim = 0.97326, p-value = 0.848  
alternative hypothesis: two.sided
```

Problem 2



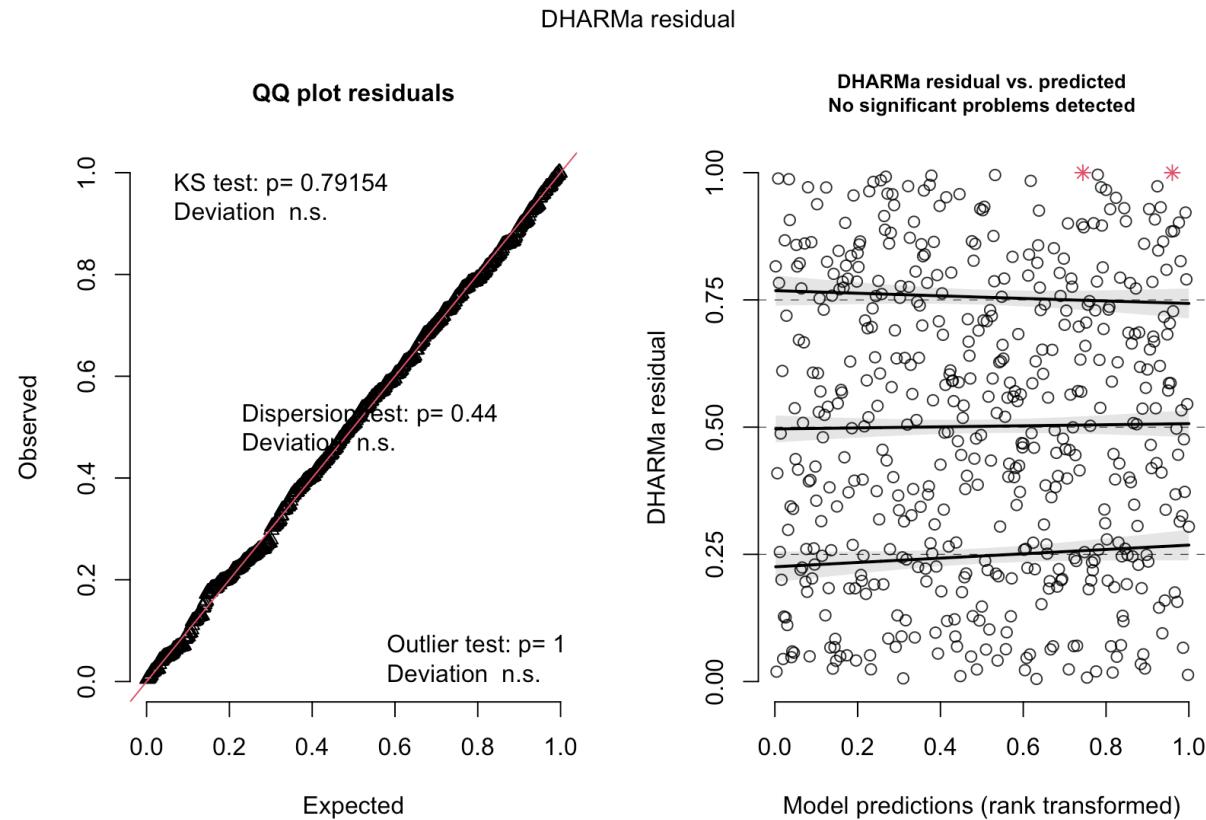
DHARMa nonparametric dispersion test via sd of residuals
fitted vs.
simulated

```
data: simulationOutput  
dispersion = 1.9199, p-value < 2.2e-16  
alternative hypothesis: two.sided
```

Add a dispersion formula

Modeling heteroscedasticity

```
1 glmmTMB(observedResponse ~ Environment1 + (1|group),  
2   dispformula = ~ Environment1,  
3   family = nbinom2(), data = heteroData)
```



DHARMA nonparametric dispersion test via sd of residuals fitted vs.
simulated

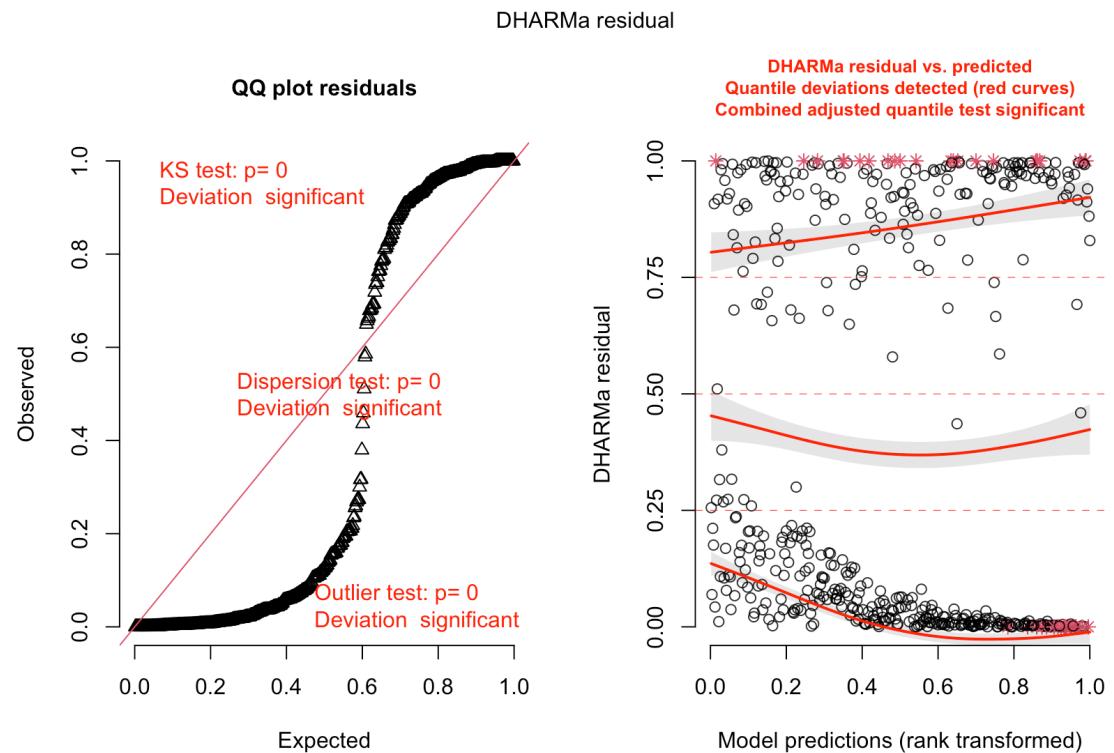
```
data: simulationOutput  
dispersion = 1.106, p-value = 0.44  
alternative hypothesis: two.sided
```

DHARMA zero-inflation test via comparison to expected zeros with
simulation under H0 = fitted model

```
data: simulationOutput  
ratioObsSim = 0.98758, p-value = 0.92  
alternative hypothesis: two.sided
```

Problem Solved!

Problem 3



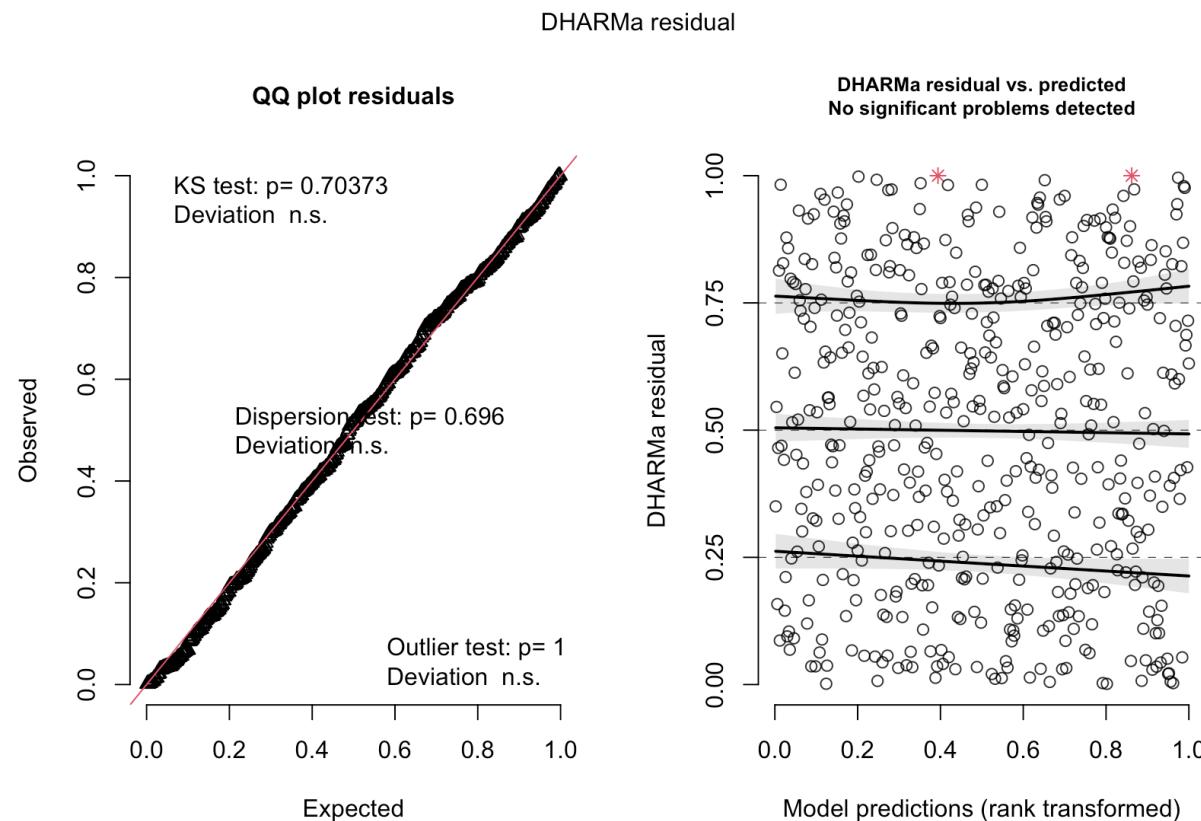
DHARMa nonparametric dispersion test via sd of residuals
fitted vs.
simulated

data: simulationOutput
dispersion = 4.9124, p-value < 2.2e-16
alternative hypothesis: two.sided

Add a zero-inflation formula

Modeling zero-inflation

```
1 glmmTMB(observedResponse ~ Environment1 + (1|group),  
2           ziformula = ~ 1,  
3           family = poisson(), data = zeroData)
```



DHARMA nonparametric dispersion test via sd of residuals fitted vs.
simulated

```
data: simulationOutput  
dispersion = 1.0414, p-value = 0.696  
alternative hypothesis: two.sided
```

DHARMA zero-inflation test via comparison to expected zeros with
simulation under H0 = fitted model

```
data: simulationOutput  
ratioObsSim = 0.99592, p-value = 1  
alternative hypothesis: two.sided
```

Problem Solved!

Solving dispersion problems

- Sometimes, residual patterns will not tell you which is the cause of overdispersion. E.g.:
 - ‘Real’ overdispersion will show significant test for zero-inflation, and vice-versa.
 - ‘Real’ overdispersion and zero-inflation may have significant heteroscedasticity/.
- Additional check: fit models addressing the potential problems and compare their fit (e.g. AIC, LRT) and residuals diagnostics.

Conclusion

- There are many causes of dispersion problems in GLMMs
- Use [DHARMA](#) residuals tools to detect them
- Address the problem with adequate models, e.g, [glmmTMB](#)

Take home message

- Models should ALWAYS be checked: residual diagnostics!
- Avoid the oversimplistic view of dispersion problems
- Detecting and addressing the causes of dispersion problems may also be informative for your system/data.



Thank you!

Vielen Dank!