

# IMRaD Report

William Pincombe

## Introduction

We are analysing data provided by Karl Berator. The data consists of the results of 8 trials where the gene expression was measured for 11 integer levels of concentration of growth factor between 0 and 10. Each trial was conducted with a treatment, either AF42 or placebo, and a cell type, either Wild Type (WT) or CT101. Two trials were carried out on each combination of treatment and cell type. Every trial was carried out on a different gene line.

There are therefore five variables in the data:

- **Gene Expression** is a continuous variable.
- **Concentration** of growth factor, measured in  $\mu\text{g/ml}$ , is a continuous variable taking integer values between 0 and 11.
- **Treatment** is a factor variable with levels AF42 and placebo.
- **Cell Type** is a factor variable with levels WT and CT101.
- **Gene Line** is a factor variable with 8 levels: CsE, bNo, JZC, fUg, jEK, Hoe, Rza and xpo.

We want to study the effect of treatment on the relationship between the concentration of growth factor and the gene expression.

## Methods

We cleaned the data for analysis by saving it in a .csv file in a long format, with a column representing each of the five variables. Some example observations of the data in this format are shown in Table 1.

Table 1: The first 5 observations of the first trial, to show how the data was structured for the analysis.

---

Concentration	Gene Expression	Cell Type	Treatment	Gene Line
0	5.51	WT	placebo	CsE
1	6.41	WT	placebo	CsE
2	5.71	WT	placebo	CsE
3	7.94	WT	placebo	CsE
4	6.87	WT	placebo	CsE
5	7.29	WT	placebo	CsE

There was one observation with a missing level of gene expression, at concentration 5 with cell type WT, treatment AF42 and gene line fUg. Since gene expression is our response variable, we excluded this observation from our data, leaving 87 observations for the analysis.

We are interested in the relationship of three predictor variables - concentration, treatment and cell type - with the response variable gene expression. However, we also must control for different gene lines used in each trial of the experiment. Therefore, we will use a mixed-effects model, with concentration, treatment and cell type as our fixed effects and gene line as a random effect.

We fit the mixed-effect model using Restricted Maximum Likelihood (REML) to avoid bias in the random effects (Morrell 1998).

We used the R statistical programming language (R Core Team 2023) for our analysis, using the tidyverse package (Wickham et al. 2019) to clean the data. We used the `lmer` command from the lme4 package (Bates et al. 2015) to fit mixed-effects models and the lmerTest package (Kuznetsova, Brockhoff, and Christensen 2017) to perform hypothesis tests on variable removal.

## Results

Looking at the scatter plot of the data in Figure 1, there appears to be a positive linear relationship between concentration and gene expression. The treatment AF42 appears to have a substantial impact on the slope of the relationship between concentration and gene expression, although this is clearly for Wild Type. Both the treatment and placebo observations appear to be higher for the same concentration in CT101 compared to Wild Type.

As a result of the relationships observed in Figure 1, we fit a model including concentration, treatment and cell type as predictors for gene expression, with interaction terms between concentration and the two categorical variables to allow the slope to vary. However, we know that each series of observations was done on a different cell type. This could have an effect on the gene expression, so we should include it in the model. However, we are not interested in this effect, we only want to remove bias from our estimates of the other coefficients, so we include it as a random effect in a Mixed-Effects Model. The fitted coefficients for the fixed effects in this model are given by Table 2.

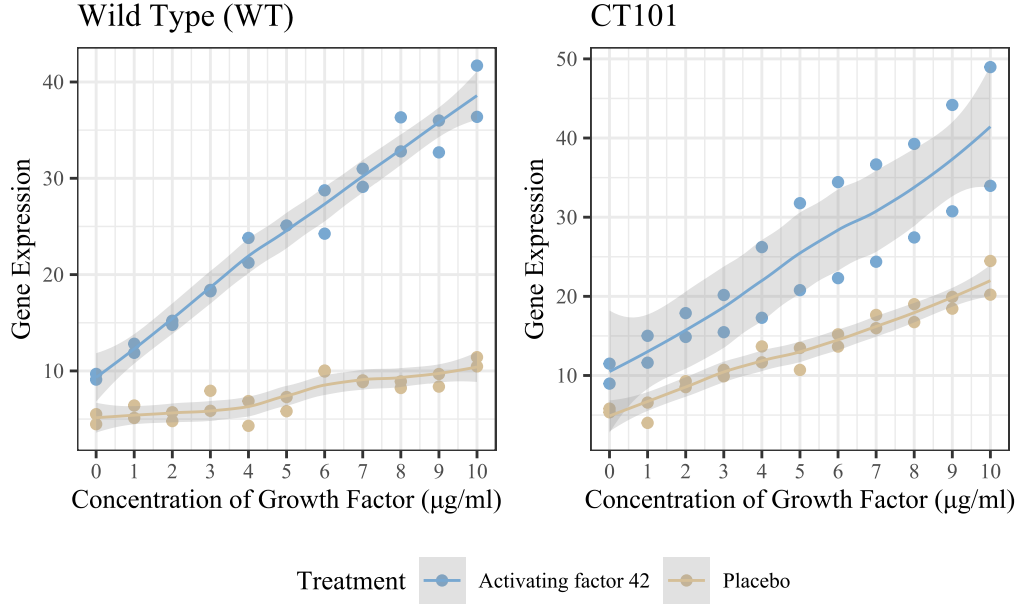


Figure 1: Scatter plot of gene expression for different concentrations of growth factor. The different cell types are shown in side-by-side plots, and the treatment is indicated by the colour.

Table 2: Summary of estimated coefficients in Mixed-Effects Model

	Estimate	Std. Error	df	t value	Pr(> t )
Intercept	9.8967	2.1084	5.6863	4.6939	0.0039
Concentration	3.2933	0.1053	76.0033	31.2840	0.0000
Treatment (placebo)	-4.8800	2.4348	5.6886	-2.0043	0.0945
cell type (WT)	-0.3213	2.4348	5.6886	-0.1320	0.8995
Concentration * Treatment (placebo)	-1.8892	0.1216	76.0033	-15.5420	0.0000
Concentration * cell type (WT)	-0.6052	0.1216	76.0033	-4.9785	0.0000

The results in Table 2 suggest that there is a strong positive relationship between concentration and gene expression. However, the effect of treatment and cell type is more ambiguous, with the slope terms returning very low standard error and high test statistics while the intercept terms have p-values greater than the standard rule of thumb of 0.05. To test the overall effect of each factor variable, we fit models excluding these variables, and then compare to the complete model.

Table 3: Model metrics for the complete mixed-effects model, and with each of the fixed effect factor variables treatment and cell type removed.

Model	AIC	BIC
Complete	387.6244	407.3516
Without cell type	408.4662	423.2616
Without Treatment	507.5614	522.3568
Without random-effects (lm)	460.0222	477.2836

The AIC and BIC of the different models is shown in Table 3. We also include a model without the random effects, a simple linear model fitted using the `lm` function. The complete mixed-effects model, with both treatment and cell type, has the lowest, hence most optimal, AIC and BIC. The removal of treatment has a much larger negative effect on both metrics than the removal of cell type.

We are also interested in whether the random effects are necessary. Firstly, we extract the values of the intercept under each level of gene line from the complete model. We then subtract the overall model intercept from each of these to show the difference in intercept for each random effect. This is shown in Table 4.

Table 4: Intercept under each level of gene line, the random effect, in the complete model. The difference to the overall intercept is also given.

	Intercept	Difference
bNo	8.197826	-1.6988968
CsE	9.281852	-0.6148705
fUg	10.077388	0.1806658
Hoe	11.965958	2.0692356
jEK	10.141254	0.2445317
JZC	12.029824	2.1331015
Rza	4.393999	-5.5027231
xpo	13.085678	3.1889558

We can test the significance of the random effect using the `ranova` command from the `lmerTest` package, which performs a likelihood ratio test. This fits a model without the random effect, and compares the log-likelihood. The results of this test are shown in Table 5, where we find that the reduced model has higher AIC and lower log likelihood. The likelihood ratio test rejects the null hypothesis that there is no difference between the log likelihoods of the models, so we take the better complete model.

Table 5: Results of a likelihood ratio test on removing the intercept random effect for gene line (GL) from the mixed-effects model.

	No. Parameters	Log Likelihood	AIC	LRT	Df	Pr(>Chisq)
Complete Model	8	-185.8122	387.6244	NA	NA	NA
Random Effect Removed	7	-224.6211	463.2421	77.61776	1	1.250408e-18

We will therefore use the complete mixed-effects model as our final model. To check that this model satisfies the regression assumptions, we plot a histogram of the residuals in Figure 2. We expect to see residuals approximately normally distributed. We observe that the residuals are close to the normal distribution, although they are somewhat right-skewed due to a few large positive values.

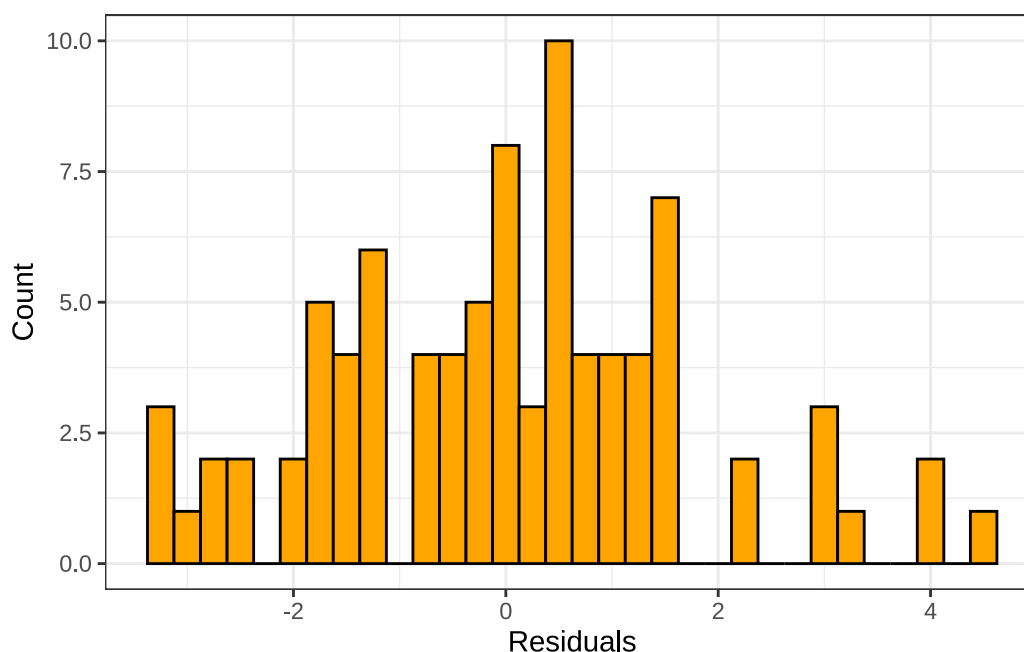


Figure 2: Histogram of the residuals of the fitted mixed-effects model.

Secondly, we plot the residuals against the fitted values in Figure 3. We expect the residuals to be roughly evenly distributed around the x-axis. There does not seem to be any substantial change in residual variance across the fitted values. However, we can again see that the residuals are somewhat right-skewed.

To check how well the final model fits the data, we can plot the fitted relationships above the scatter plot of the data from Figure 1, as shown in Figure 4.

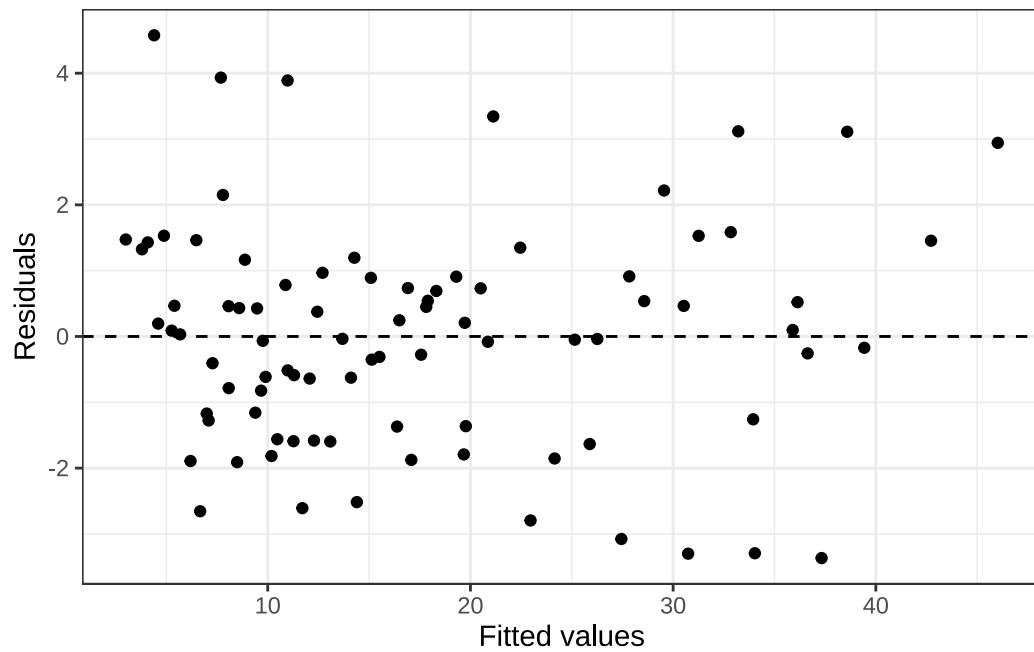


Figure 3: Plot of residuals against fitted values. The residuals satisfy the assumption of homoskedasticity, with roughly constant variance. However, the residuals are somewhat right-skewed.

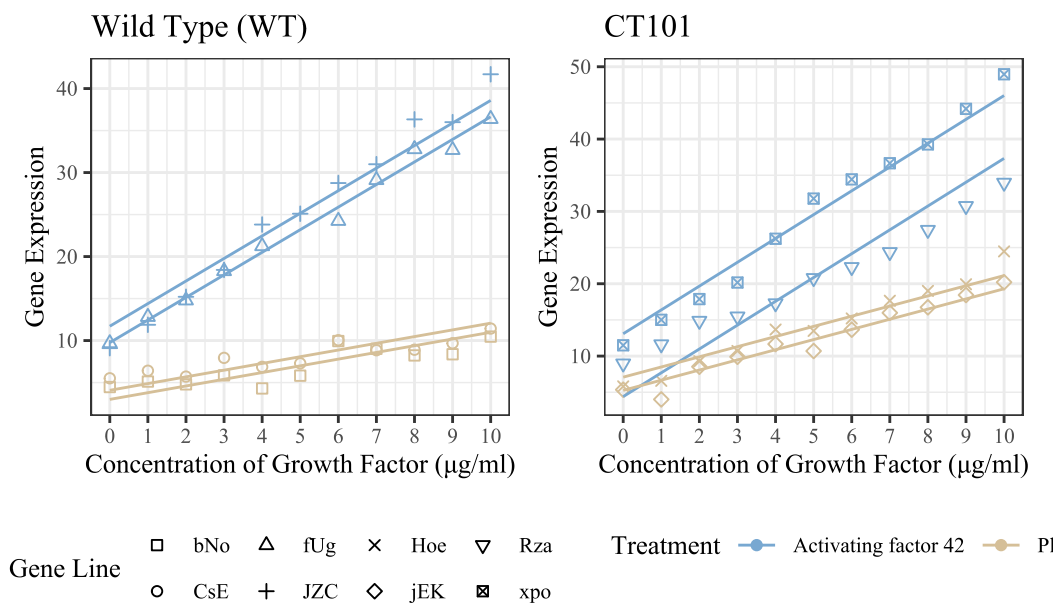


Figure 4: Scatter plot of gene expression over concentration of growth factor. The points show the observations while the lines show the fitted values from the complete mixed-effects model. The cell types are shown in separate plots. We can see here that the random-effect term has adjusted the fitted lines to more accurately account for the differences between trials on different gene lines.

## Discussion

Our model shows that there is a significant positive correlation between concentration and gene expression and that the AF42 treatment increases the slope of this relationship.

In our final model, the slope coefficient for concentration is estimated to be 1.8892 greater for the treatment than for the placebo. This indicates that we expect, for instance for the cell type CT101 that an increase in 1  $\mu\text{g}/\text{ml}$  of concentration will lead to an increase of around 1.40 in gene expression without the treatment, and an increase of 3.29 if the treatment is applied. The fact that the AIC and BIC metrics were substantially worse without treatment, as seen in Table 4, indicates the significance of this effect.

We also found that the effect of concentration was different across the two cell types, with increases in concentration having a greater effect for CT101 than for Wild Type.

We can conclude that treatment has a significant positive effect on the effect of concentration on gene expression.



## References

- Bates, Douglas, Martin Mächler, Ben Bolker, and Steve Walker. 2015. “Fitting Linear Mixed-Effects Models Using lme4.” *Journal of Statistical Software* 67 (1): 1–48. <https://doi.org/10.18637/jss.v067.i01>.
- Kuznetsova, Alexandra, Per B. Brockhoff, and Rune H. B. Christensen. 2017. “lmerTest Package: Tests in Linear Mixed Effects Models.” *Journal of Statistical Software* 82 (13): 1–26. <https://doi.org/10.18637/jss.v082.i13>.
- Morrell, Christopher H. 1998. “Likelihood Ratio Testing of Variance Components in the Linear Mixed-Effects Model Using Restricted Maximum Likelihood.” *Biometrics* 54 (4): 1560–68.
- R Core Team. 2023. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Wickham, Hadley, Mara Averick, Jennifer Bryan, Winston Chang, Lucy D’Agostino McGowan, Romain François, Garrett Golemund, et al. 2019. “Welcome to the tidyverse.” *Journal of Open Source Software* 4 (43): 1686. <https://doi.org/10.21105/joss.01686>.