

Analyzing indonesian rice farms

Jonas Kernebeck, Alexander Flick, Felix Lehner

07/16/2021

Contents

1 Introduction	2
1.1 Numerical Variables	2
1.1 Categorical Variables	4
1.3 Variable selection and transformation	5
1.4 Model evaluation	5
2 First Model	6
3 Second Model	7
GAM	7
GAM Comparison	7
Final model visualization	11
4 Comparision	12
5 Conclusion	13

1 Introduction

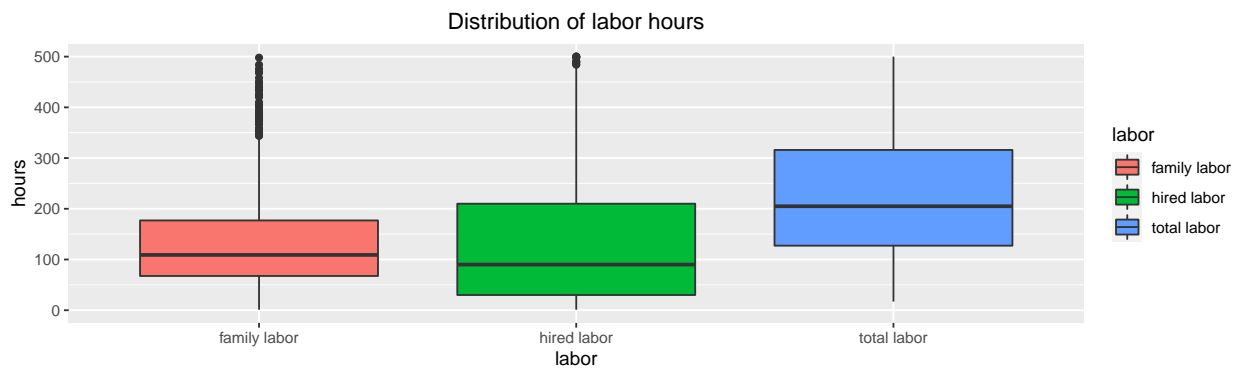
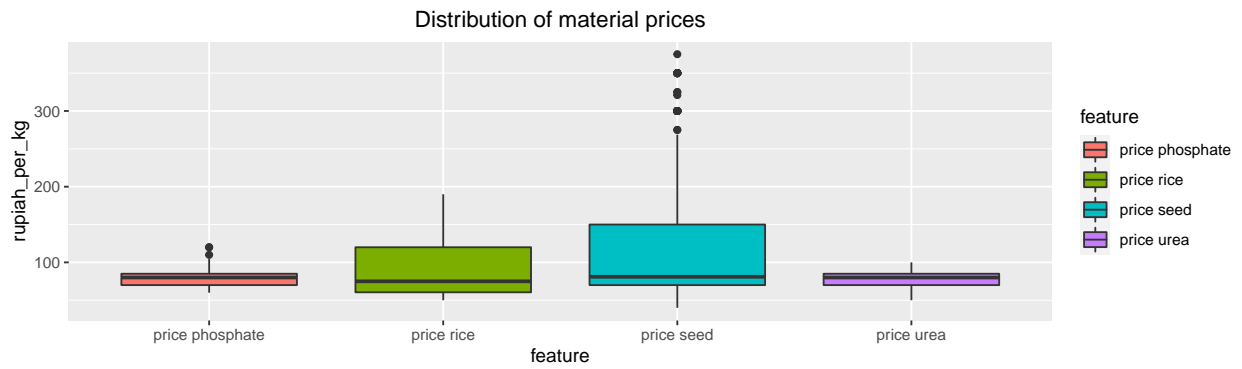
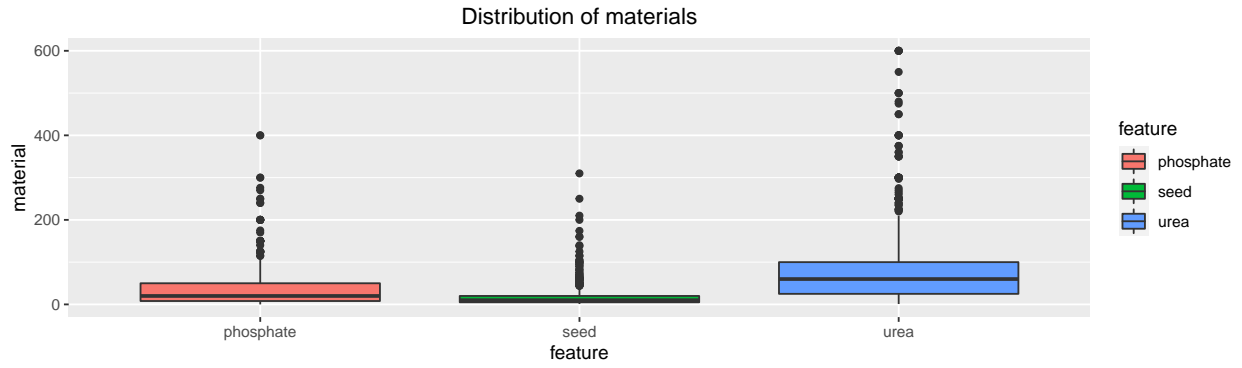
The present data set includes production data for 171 Indonesian rice farms. The dataframe contains the following variables:

variable	description	expressions
id	unique identifier for a farm	unique id
time	unique identifier for a specific growing season	1 - 6
size	total production area in hectares	0.01 - 5.322
status	status of property rights	“owner”, “share”, “mixed”
varieties	rice seed varieties	“trad”, “high”, “mixed”
bimas	bimas-status of the farmers	“no”, “yes”, “mixed”
seed	seed in kilogram	1 - 1250 kg
urea	urea in kilogram	1 - 1250 kg
phosphate	phosphate in kilogram	0 - 700 kg
pesticide	pesticide cost in Rupiah	0 - 62600 r
pseed	price of seed in Rupiah per kg	40 - 375 r/kg
purea	price of urea in Rupiah per kg	50 - 100 r/kg
pphosph	price of phosphate in Rupiah per kg	60 - 120 r/kg
hiredlabor	hired labor in hours	1 - 4536 h
famlabor	family labor in hours	1 - 1526 h
totlabor	total labor (excluding harvest labor)	1 - 4774 h
wage	labor wage in Rupiah per hour	30 - 175.35 r/h
goutput	gross output of rice in kg	42 - 20960 kg
noutput	gross output minus harvesting cost	42 - 17610 kg
price	price of rough rice in Rupiah per kg	50 - 190 r/kg
region	region of the farm	unique region

As present in the table, the data set consists of 16 numeric variables and 4 categorical variables. The target variable for the regression modeling will be *goutput*, what represents the gross output of rice in *kg* for the respective rice farm. In the following some explorative data analysis will be made to get to get a first impression of the distribution of the individual variables.

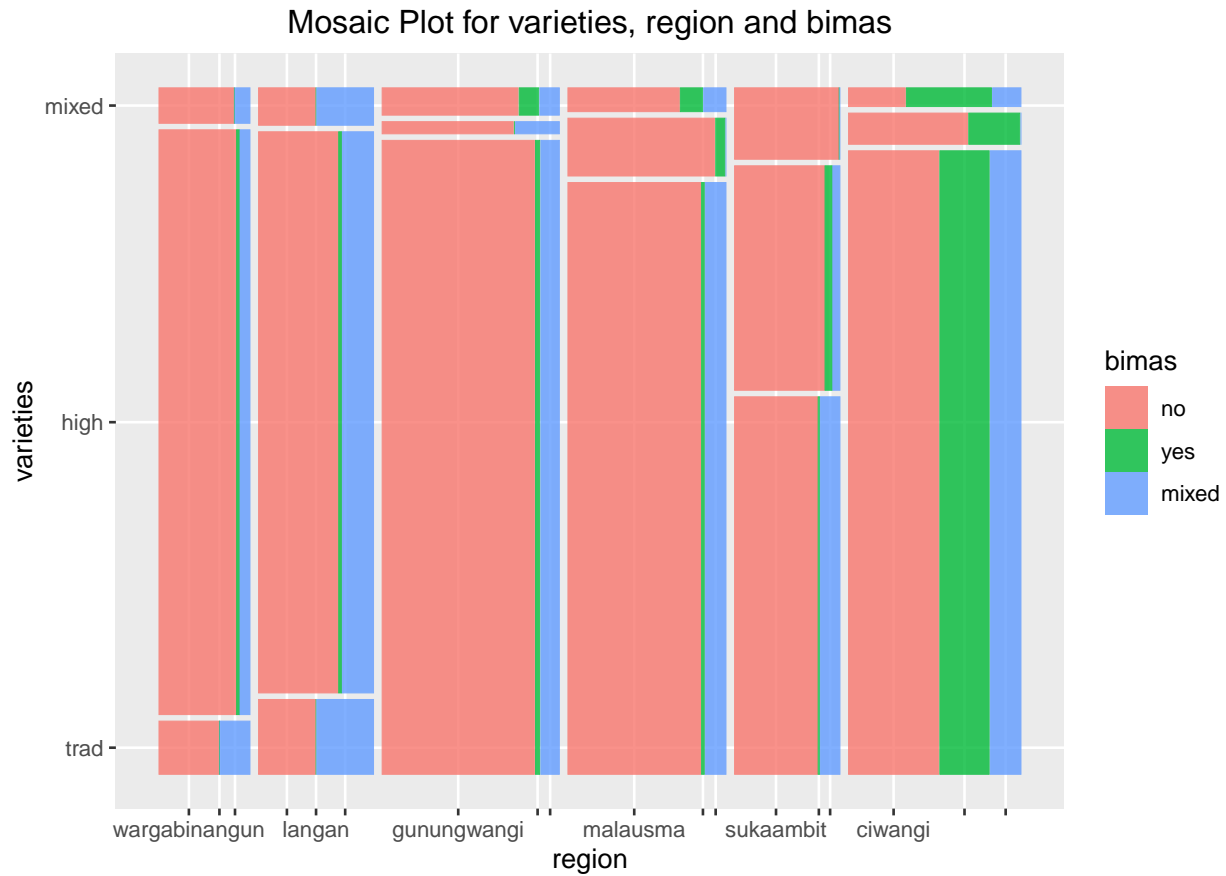
1.1 Numerical Variables

The following figure shows boxplots for the used materials and the prices paid for the materials of the respective rice farms. The boxplots for the materials show, that the distribution of all materials is right-skewed. The spread width of seed is the lowest, followed by phosphate and urea. Therefore *urea* also has the highest variance with 16166 followed by *phosphate* with 2264 and *seed* with 2048. The distribution of *urea* indicates that rice farms in Indonesia may use urea very different, caused by e.g. the bimas-status. The bimas program is a rice intensification program by the government to support local rice production by providing high-yield rice seeds as well as technical assistance. If we look at the prices for phosphate *pphosph* and urea *purea*, we can see a slight left-skewed distribution with low variance (75 for *purea* and 86 for *pphosph*). In contrast to that, the prices for seeds scatter much. The distribution of *pseed* is strongly right-skewed as well as the distribution for the rice price *price*. The price for the rice also scatters, but less than *pseed*. The two prices have a correlation of 0.67. Of course, the price of seeds affects the selling price of rice. The prices may fluctuate due to seasonal or regional factors and have an impact on each other. The distribution of labor hours is also slightly skewed to the right. Overall, the dispersion is lowest for the *famlabor*. For *hiredlabor* and *totlabor* we have a similar spread, but *totlabor* has a higher level overall. This is caused by the *hiredlabor* which is a subset of *totlabor*.



1.1 Categorical Variables

The following mosaic plot shows the distribution of of the categorical variables *varietes*, *region* and *bimas*. Overall, all regions are roughly equally represented in the data set. We can detect, that most of the farmers with the *bimas* status *yes* and *mixed* are located in the region *ciwangi*. The distribution of the different varieties is strongly dependent on the region. While the *high* varieties have the biggest share in the regions *wargabinangun* and *langan*, the *traditional* varieties are dominating the regions *gunungwangi*, *malausma* and *ciwangi*. The *mixed* varieties are only used slightly in all regions.



To test wheter the categorical variables have impact on our target variable *goutput*, one- and two-sided anovas are performed. The results of these are summarized in the following table:

formula	F-value	p-value	significant
region	22.981	< 2e-16	yes
varieties	11.764	8.94e-06	yes
bimas	14.817	4.57e-07	yes

formula	F-value	p-value	significant
region+varietes	3.847	3.96e-05	yes
region+bimas	5.651	2.94e-08	yes
varieties+bimas	0.791	0.531	no
region+varieties+bimas	0.860	0.580	no

The anova outputs show, that all of the categorical variables have a significant effect on *goutput*. The null hypothesis, that the mean of *goutput* is the same across the groups is rejected. The results of the two-sided anovas also show a significant interaction effect on *goutput*. While the interaction effect from the *region* with *varieties* and *bimas* is significant, the interaction effect of *varieties* and *bimas* and the interaction effect of all three variables is not.

1.3 Variable selection and transformation

The performance of the regression modeling is highly dependent of the variable selection and transformation. Therefore a suitable choice is very important. The variable *noutput* is a linear transformation of *goutput* as it represents *goutput* decreased by the harvesting costs. Therefore it is not used for the modeling because it would violate the multicollinearity assumption.

The variable *size* also correlates *strongly* with the target variable. This can be intuitively explained by the fact that a larger rice field naturally always produces a higher yield. Since the variables *seed*, *urea*, *phosphate* and *pesticide* are dependent on size, they are transformed into per-hectare sizes by dividing them with the respective hectare size of the farm. The *size* variable is not used for further modeling.

The variables *famlabor* and *hiredlabor* are subsets of the variable *totlabor* and are therefore transformed into the share of *totlabor* by dividing them with the amount of *totlabor*. The variable *totlabor* is after that transformed to a per-hectare size by dividing it with the *size*. The variable *wage* follows a bimodal distribution. Therefore it is transformed into a binary variable, which indicates if the respective value is over or under 100.

1.4 Model evaluation

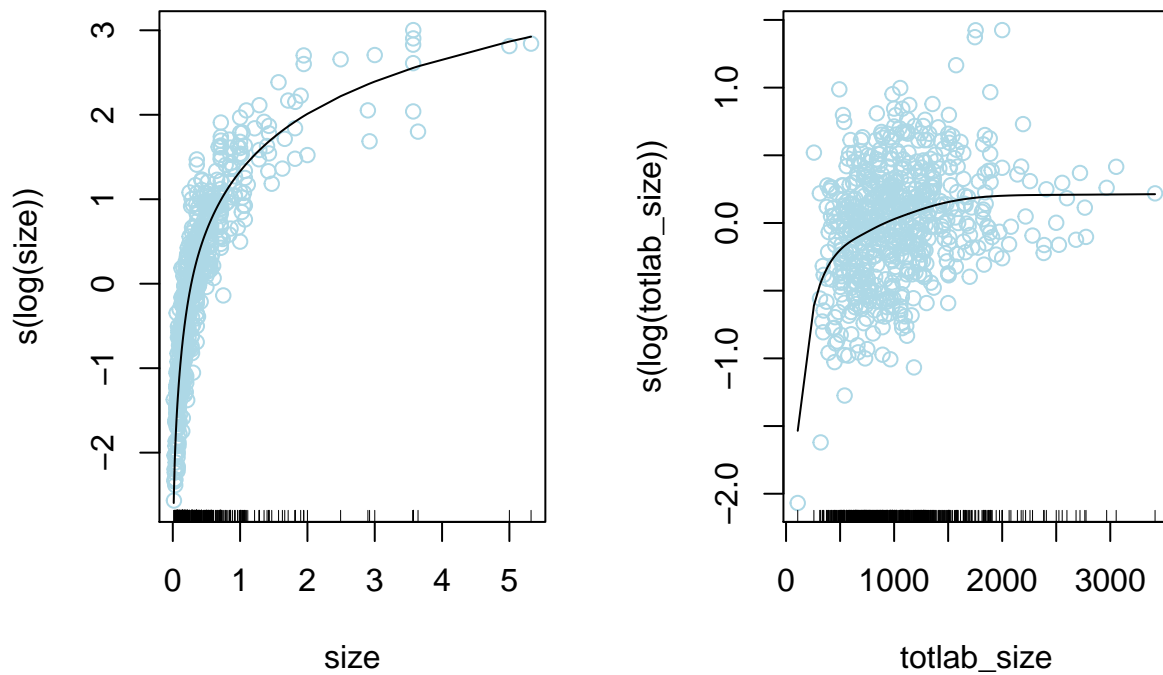
The data set will be splitted in 60-20-20 parts, where 60% of the data is used for training the model, and 20% for testing and validating respectively. In the modeling part, also cross-validation is used. To evaluate the models and compare them, different metrics will be used. The numeric metrics used are the *MSE*, which stands for the mean squared error and the *AIC*, which stands for the Akaike information criterion. Beside these metrics, also graphical analysis plots like a residual plots are used for evaluation.

3 Outliers most likely due to typos for observation 110, 947 and 1004

2 First Model

3 Second Model

GAM



Explanation: The left-hand panel indicates that holding totlab_size fixed, goutput increases with size. The right-hand panel indicates that holding size fixed, goutput increases drastically with the increased proportion of labour per size up to a proportion of 500 h / hectar and then flattens out.

GAM Comparison

table for model selection

Table 3: GAM comparison for variable selection

var	df	MSE.train	MSE.val	dev	aic	p_val_p	p_val_np	df_np
size								
s(size)	609	819023	992780	131.0	806.2	0.0000	0.0000	3
s(log(size))	609	823909	928486	111.3	706.1	0.0000	0.0070	3
labour								
s(log(totlab_size))	605	730930	774765	98.9	641.5	0.0000	0.0084	3
s(totlab_size)	605	732302	766662	100.6	651.8	0.0000	0.0001	3
urea								
s(log(urea_size))	601	765281	734252	78.2	505.5	0.0000	0.0000	3
s(urea_size)	601	769145	735393	78.3	506.3	0.0000	0.0001	3
phosphor								
s(log(phosph_size + 1))	597	682109	545771	73.9	478.6	0.0000	0.0000	3
s(phosph_size)	597	675247	550736	74.1	480.3	0.0000	0.2995	3
seed								
s(log(seed_size))	593	675922	539746	71.9	469.3	0.0001	0.0690	3
s(seed_size)	593	680228	548376	71.6	467.4	0.0002	0.0324	3
pesticide								
s(pest_size)	589	640124	504684	69.7	458.4	0.0001	0.4393	3
s(pest_size, df = 1)	592	654365	513624	70.0	455.0	0.0002	0.0019	0
price								
s(price)	585	530945	401614	65.7	430.1	0.0007	0.0000	3
family labour								
s(fam_ratio)	581	509179	436883	64.8	429.9	0.0447	0.3394	3
price info								
s(pseed)	577	483884	430733	63.9	429.0	0.0742	0.0809	3
s(pphosph)	577	470693	387321	60.5	396.0	0.0000	0.0000	3
s(pphosph)+s(purea)	573	457617	374707	59.7	395.3	0.0075	0.0903	3
wage								
s(wage)	573	417753	377066	58.4	381.5	0.0000	0.1739	3
wage_cat>100	576	443144	322736	58.9	381.7	0.0000	NA	NA
categorical variables								
bimas	574	446945	350651	57.6	371.8	0.0006	NA	NA
bimas+varieties	572	439200	329234	57.1	370.3	0.0743	NA	NA
bimas+status	572	430164	366356	57.3	372.7	0.1897	NA	NA
bimas+region	569	430469	425267	56.2	366.6	0.0013	NA	NA
final model								
s(purea)	565	412248	394646	55.2	363.2	0.0090	0.0196	3
s(purea)+varieties	563	397440	386455	54.8	362.4	0.1231	NA	NA

gam2 is better, so we will use logarithm of size

plot.gam for urea: slope at beginning of urea size => reduce df?

phosphate: this variable has lots of zeros as values. So we can not use directly log transformation because log of 0 is -inf. According to : <https://discuss.analyticsvidhya.com/t/methods-to-deal-with-zero-values-while-performing-log-transformation-of-variable/2431> log(x+1) transformation is the best way to avoid errors created by log transformation and is widely used among data scientists. So we will use this approach.

pesticide: Anova for Non-Parametric Effects p-value = 0.0001351187 => keep variable in model Anova for Parametric Effects p-value=0.433418 => reduce df of pest_size

price information

check categorical variables

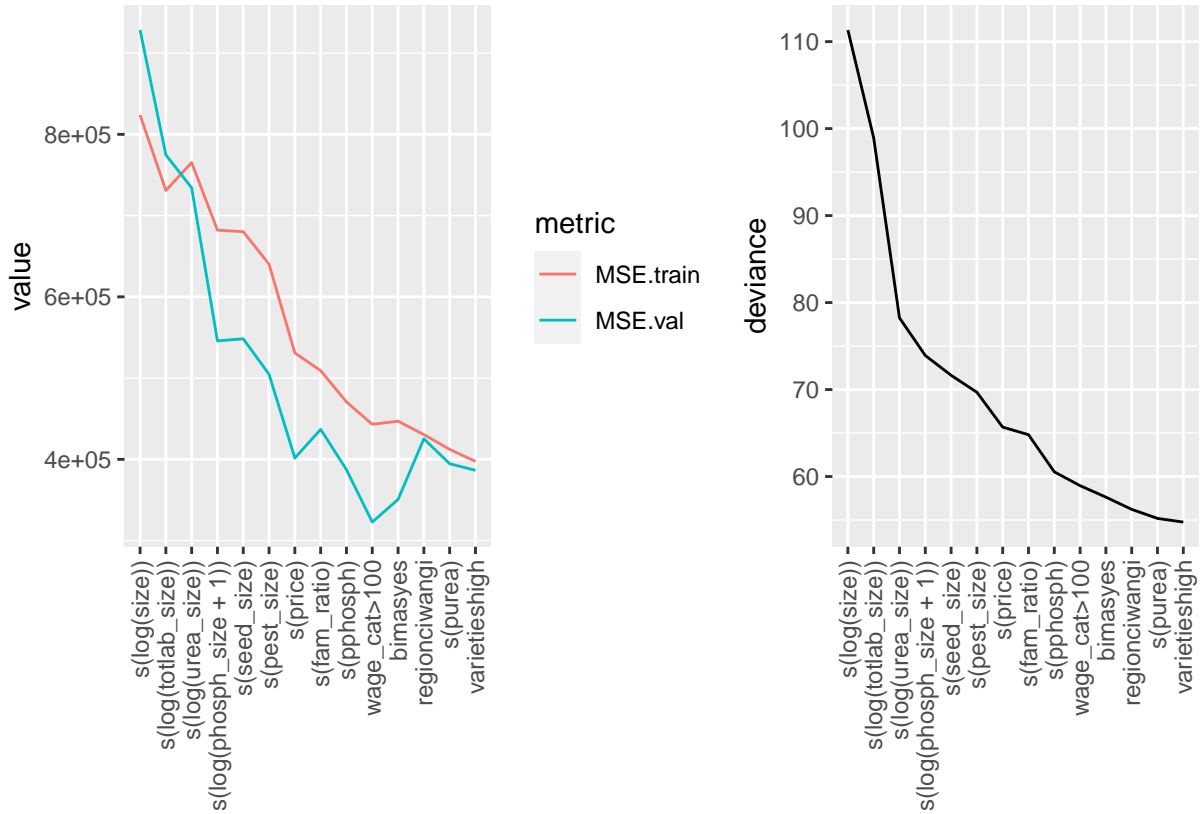
varieties: p-value 0.103104 => not enough evidence that this variable is significantly important based on a 5% significance level. BUT MSE decreases!

status: p-value 0.212109 => not important

varieties

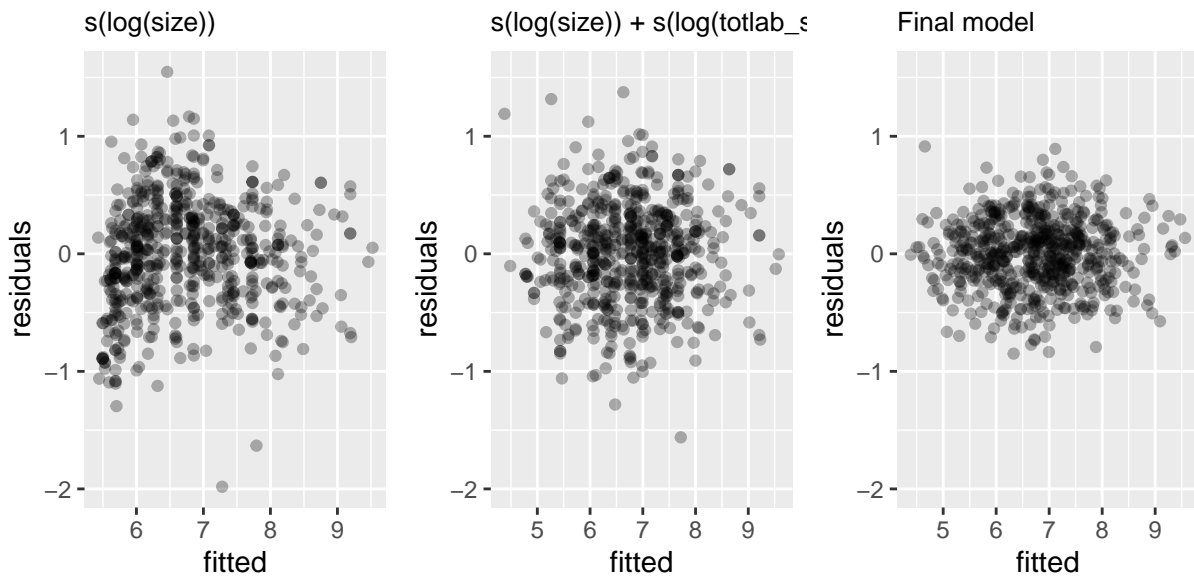
```
## Analysis of Deviance Table
##
## Model 1: log(goutput) ~ s(log(size)) + s(log(totlab_size)) + s(log(urea_size)) +
##      s(log(phosph_size + 1)) + s(seed_size) + s(pest_size) + s(price) +
##      s(fam_ratio) + s(pphosph) + wage_cat + bimas + region + s(purea)
## Model 2: log(goutput) ~ s(log(size)) + s(log(totlab_size)) + s(log(urea_size)) +
##      s(log(phosph_size + 1)) + s(seed_size) + s(pest_size) + s(price) +
##      s(fam_ratio) + s(pphosph) + wage_cat + bimas + region + s(purea) +
##      varieties
##   Resid. Df Resid. Dev Df Deviance Pr(>Chi)
## 1         565      55.188
## 2         563      54.760  2   0.42847   0.1105
```

MSE and deviance visualization

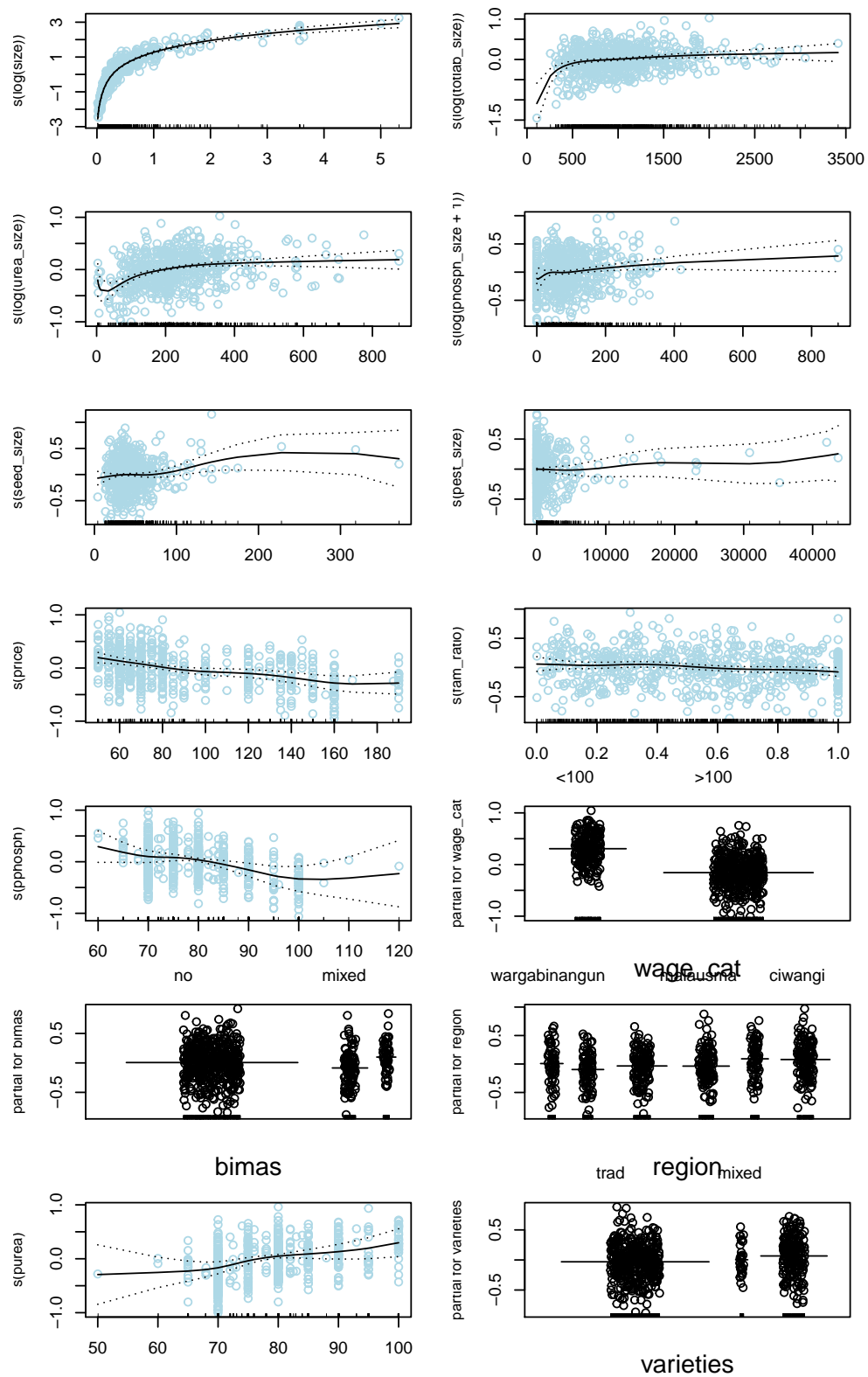


Residual Plot

Residual comparisor



Final model visualization



4 Comparision

5 Conclusion