

# Linear combiner training on random subsets of different sizes

```
library(ggplot2)
```

```
## Warning: package 'ggplot2' was built under R version 4.0.5
```

```
library(dplyr)
```

```
## Warning: package 'dplyr' was built under R version 4.0.5
```

```
##
```

```
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':
```

```
##
```

```
##      filter, lag
```

```
## The following objects are masked from 'package:base':
```

```
##
```

```
##      intersect, setdiff, setequal, union
```

```
library(scales)
```

```
## Warning: package 'scales' was built under R version 4.0.3
```

Experiment code is in the file `training_subsets_sizes_experiment.py`. Experiment on both CIFAR10 and CIFAR100 datasets. This experiment trains `WeightedLinearEnsemble` on various subsets of different sizes of data on which neural networks were trained. Four different coupling methods are used: method one and two (m1 and m2) from (Wu, Lin, and Weng 2004), Bayes covariant method (bc) from (Šuch and Barreda 2016) and a coupling method by Such, Benus and Tinajova (sbt) from (Šuch, Benuš, and Tinajová 2015). Three combining methods are used: `lda`, `logreg` and `logreg_no_interc` which employ linear discriminant analysis or logistic regression individually on each pair of classes. Goal of this experiment is to determine, for which size of the combiner training set, the ensemble achieves the best performance. The requirement of `lda`, that the predictors are normally distributed is not met by more than half of the class pairs. This may be reflected in the worsened prediction capacity of ensembles using `lda` combining method and we will consider these ensembles as less influential on the decision of the training set size.

```
metrics <- c("accuracy", "nll", "ece")
```

## CIFAR-10

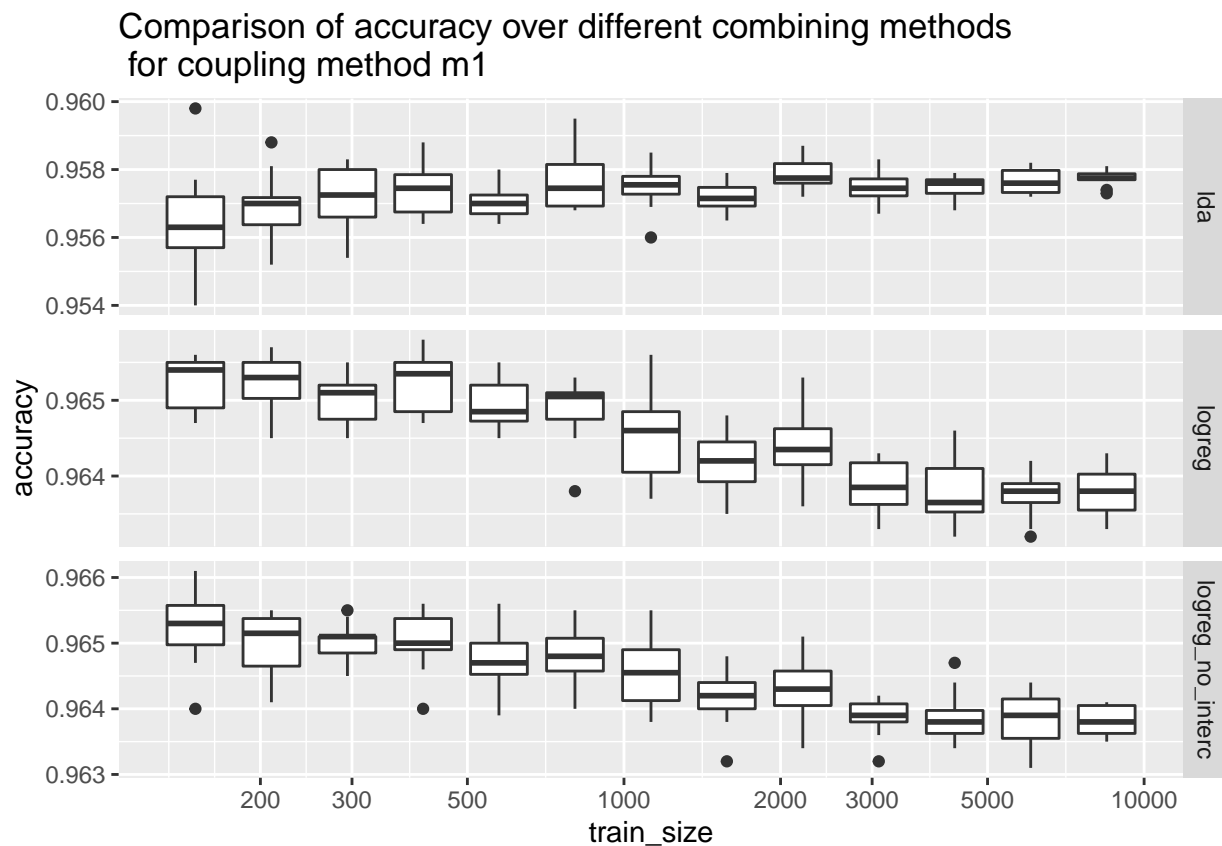
```

ens_metrics_c10 <- read.csv("../data/data_tv_5000_c10/0/exp_subsets_sizes_train_outputs/ens_metrics.csv")
nets_metrics_c10 <- read.csv("../data/data_tv_5000_c10/0/exp_subsets_sizes_train_outputs/net_metrics.csv")

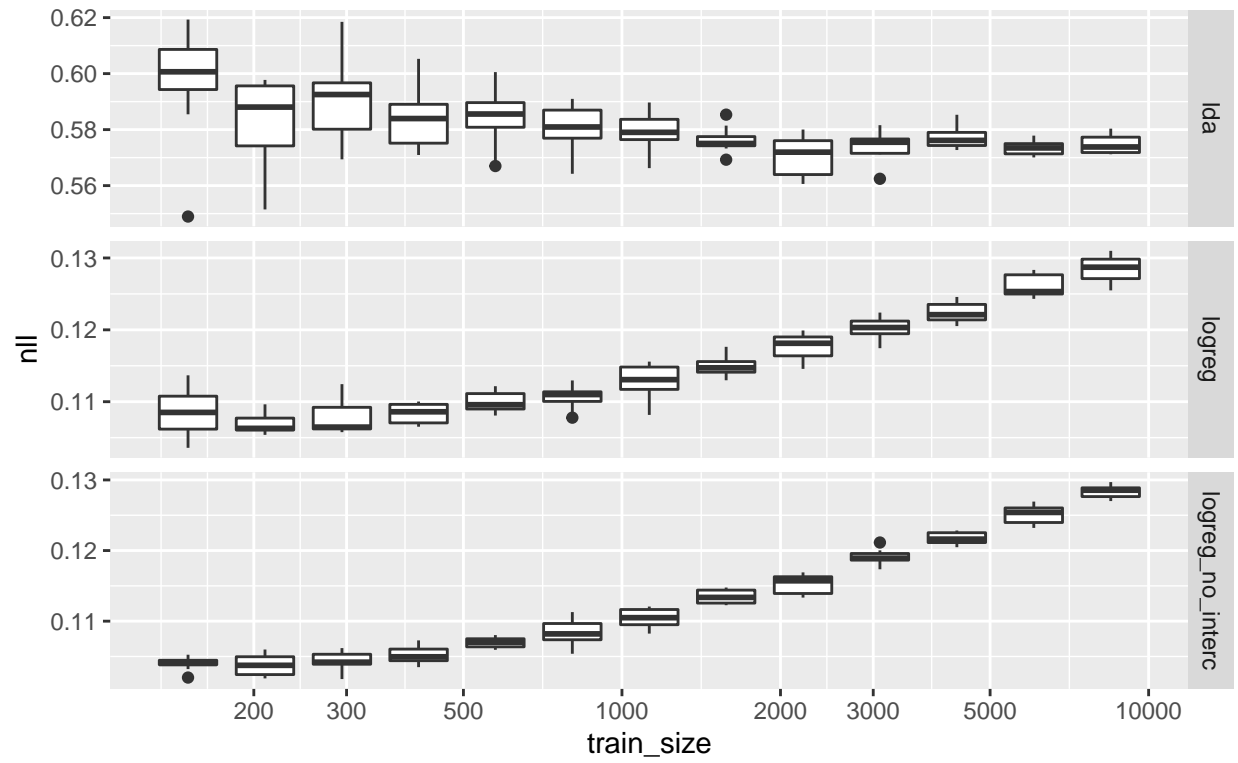
for (cp_m in unique(ens_metrics_c10$coupling_method))
{
  for (metric in metrics)
  {
    box_plt <- ens_metrics_c10 %>% filter(coupling_method==cp_m) %>%
      ggplot() +
      geom_boxplot(mapping=aes_string(x="train_size", y=metric, group="train_size")) +
      facet_grid(rows=vars(combining_method), scales="free") +
      scale_x_log10(breaks=log_breaks(n=10)) +
      ggtitle(paste0("Comparison of ", metric, " over different combining methods\n for coupling method", cp_m))

    print(box_plt)
  }
}

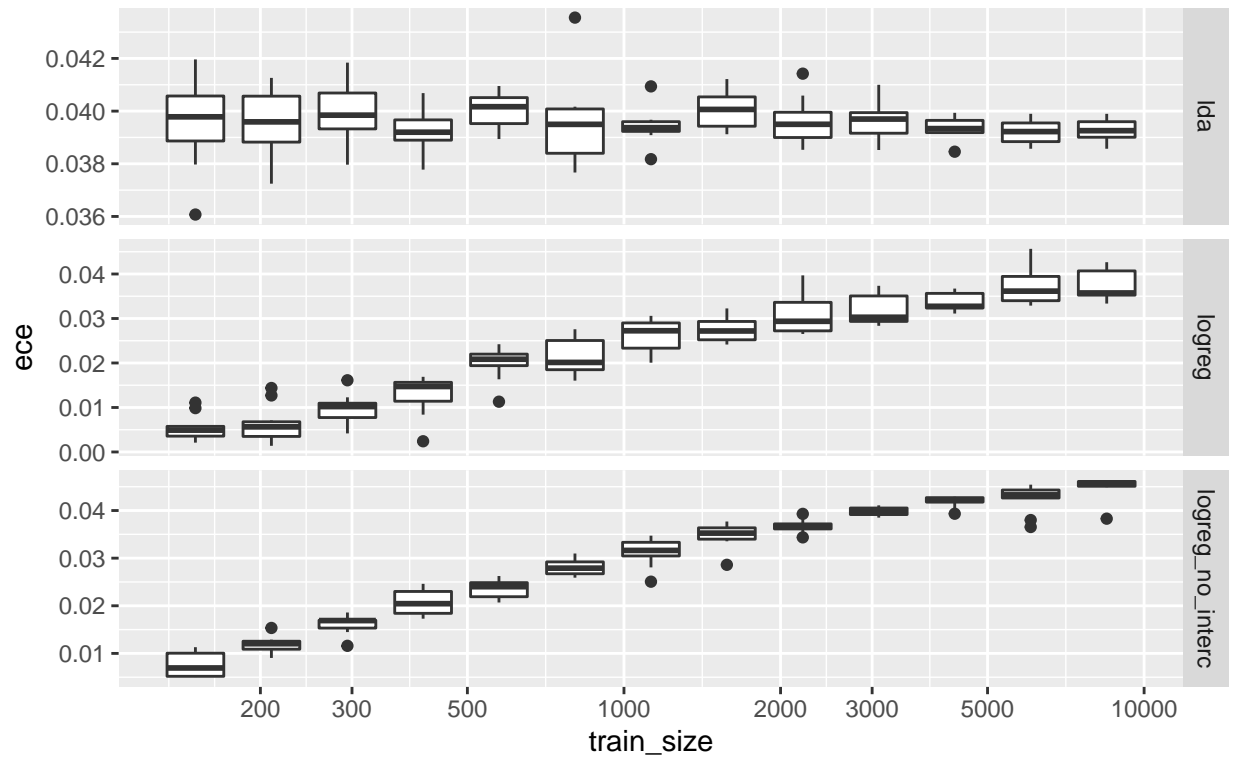
```



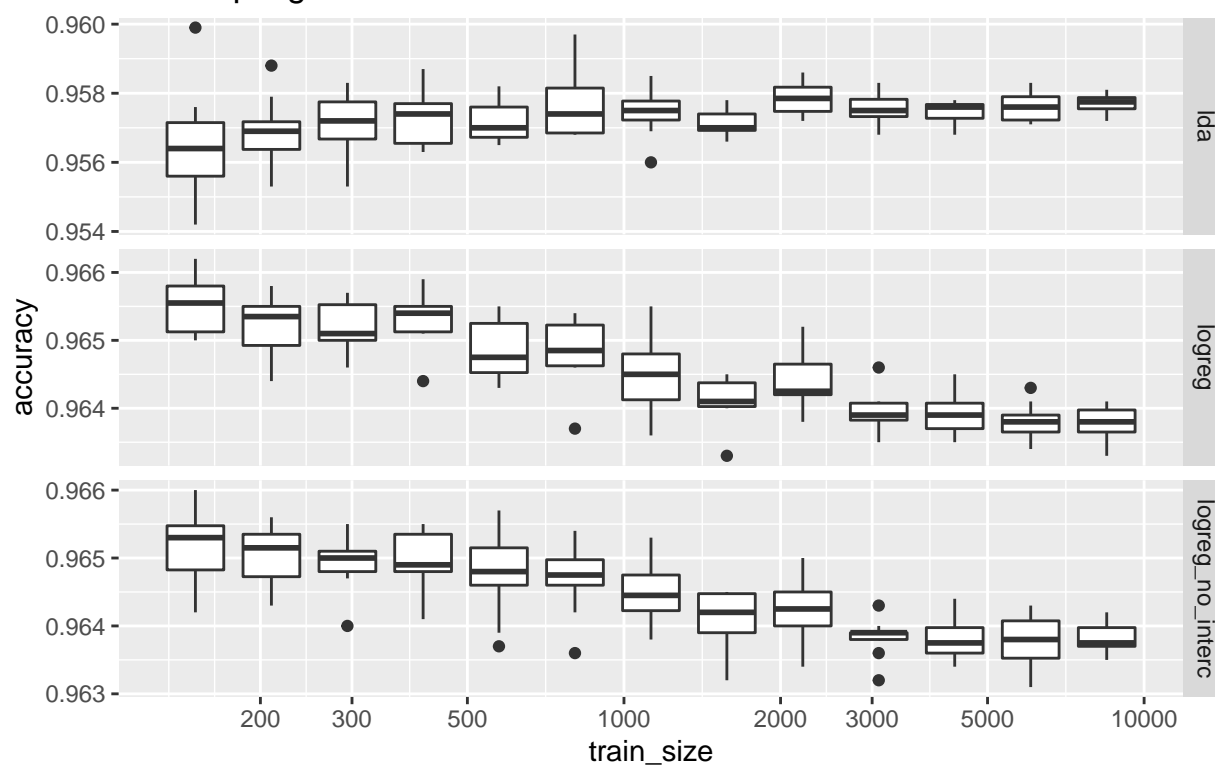
Comparison of nll over different combining methods  
for coupling method m1



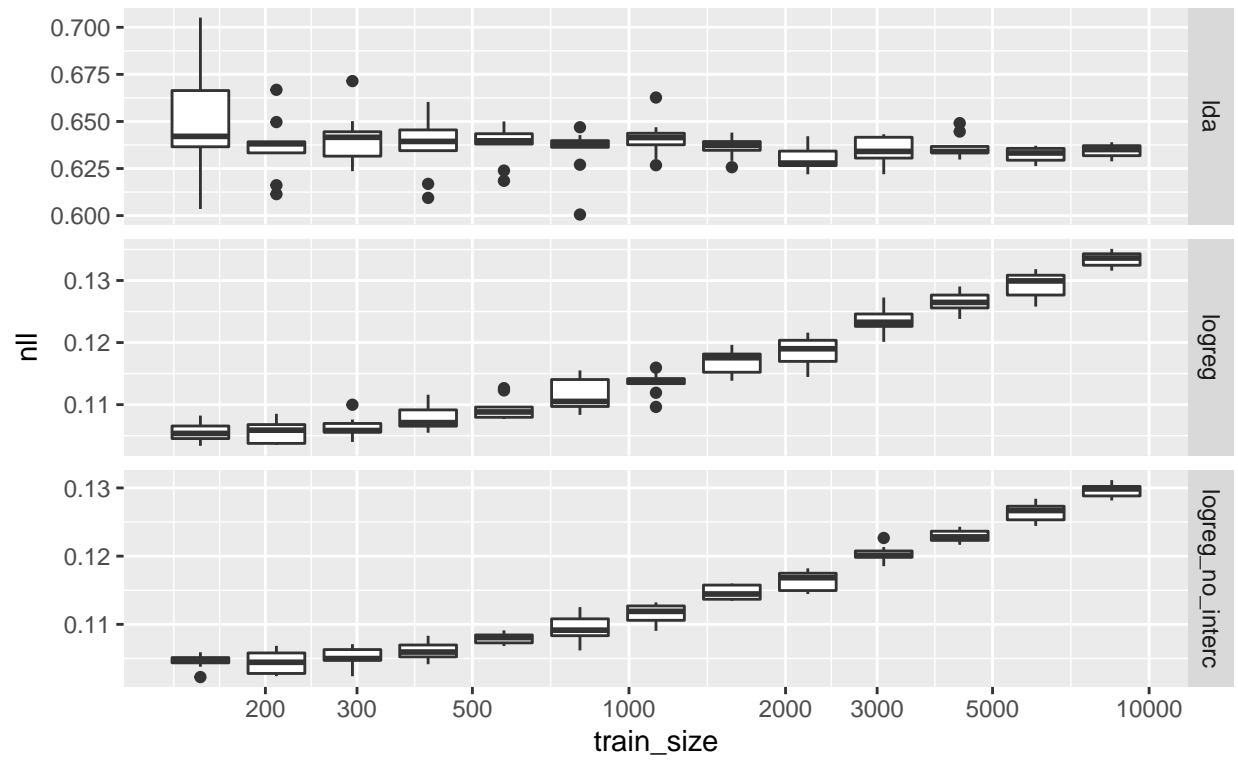
Comparison of ece over different combining methods  
for coupling method m1



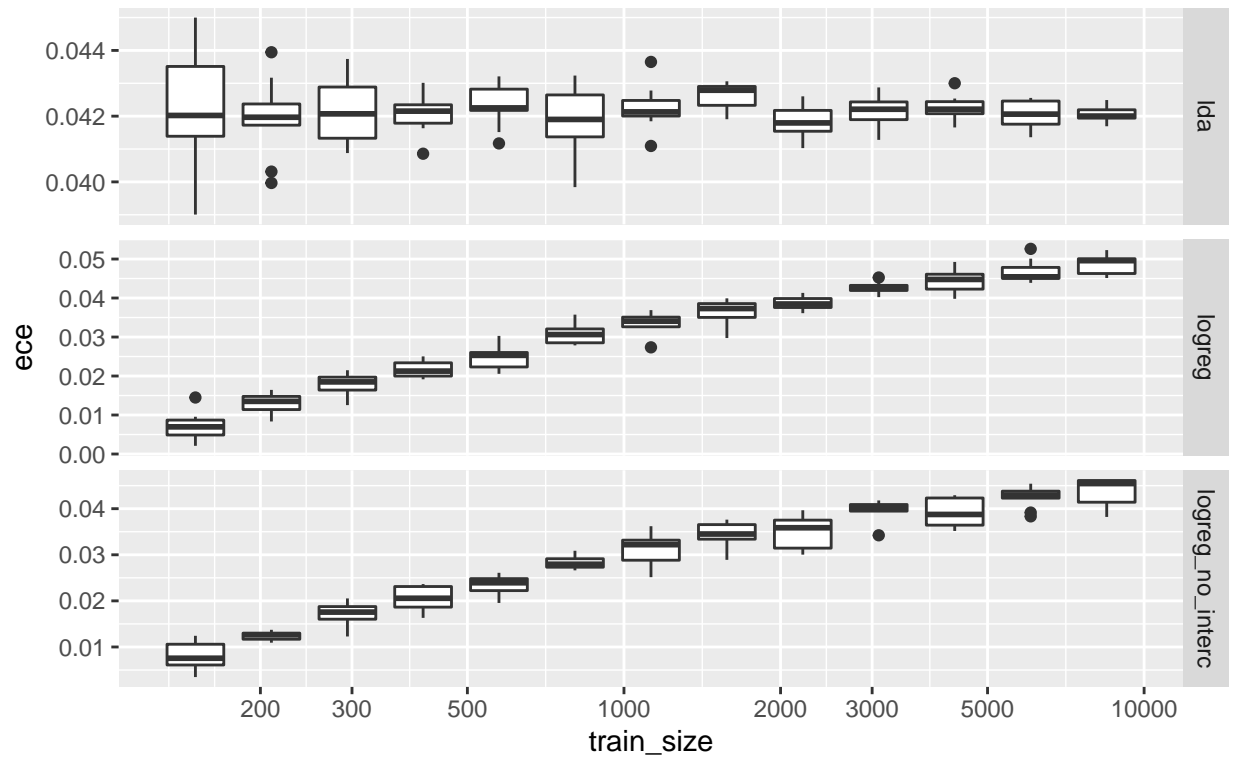
Comparison of accuracy over different combining methods  
for coupling method m2



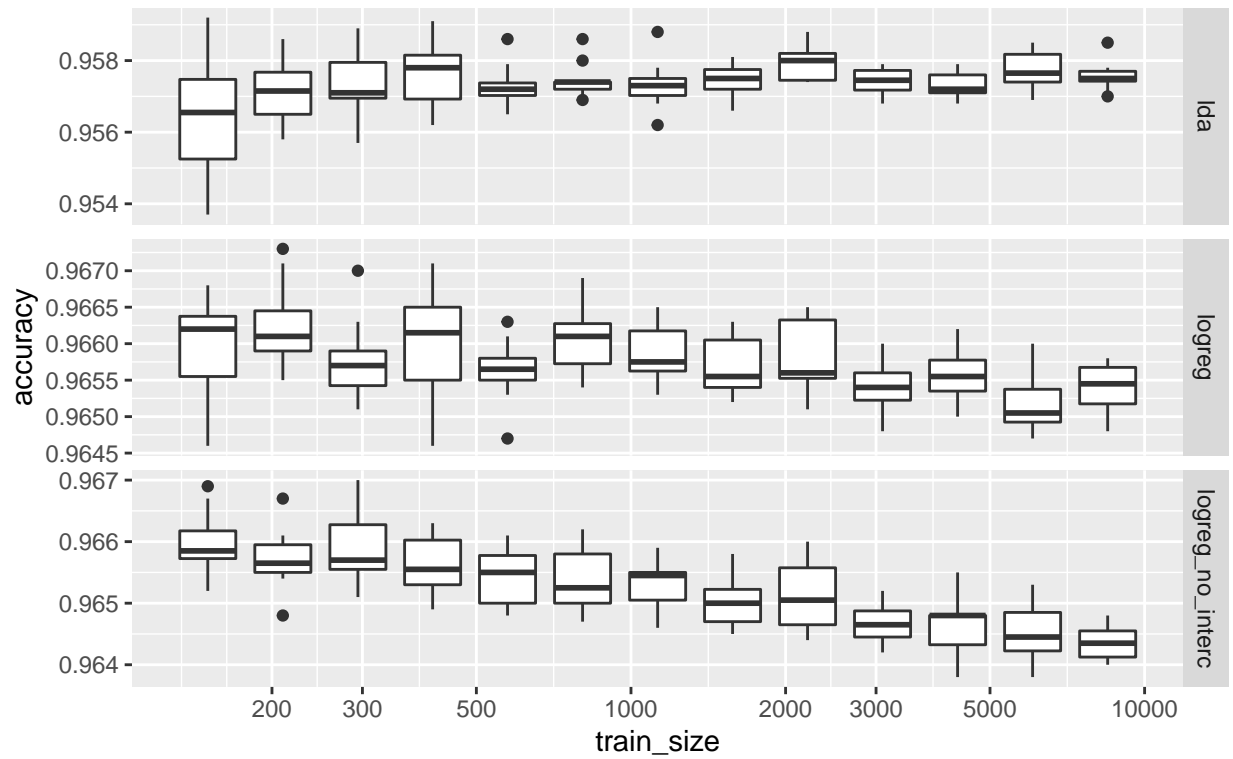
Comparison of nll over different combining methods  
for coupling method m2



Comparison of ece over different combining methods  
for coupling method m2

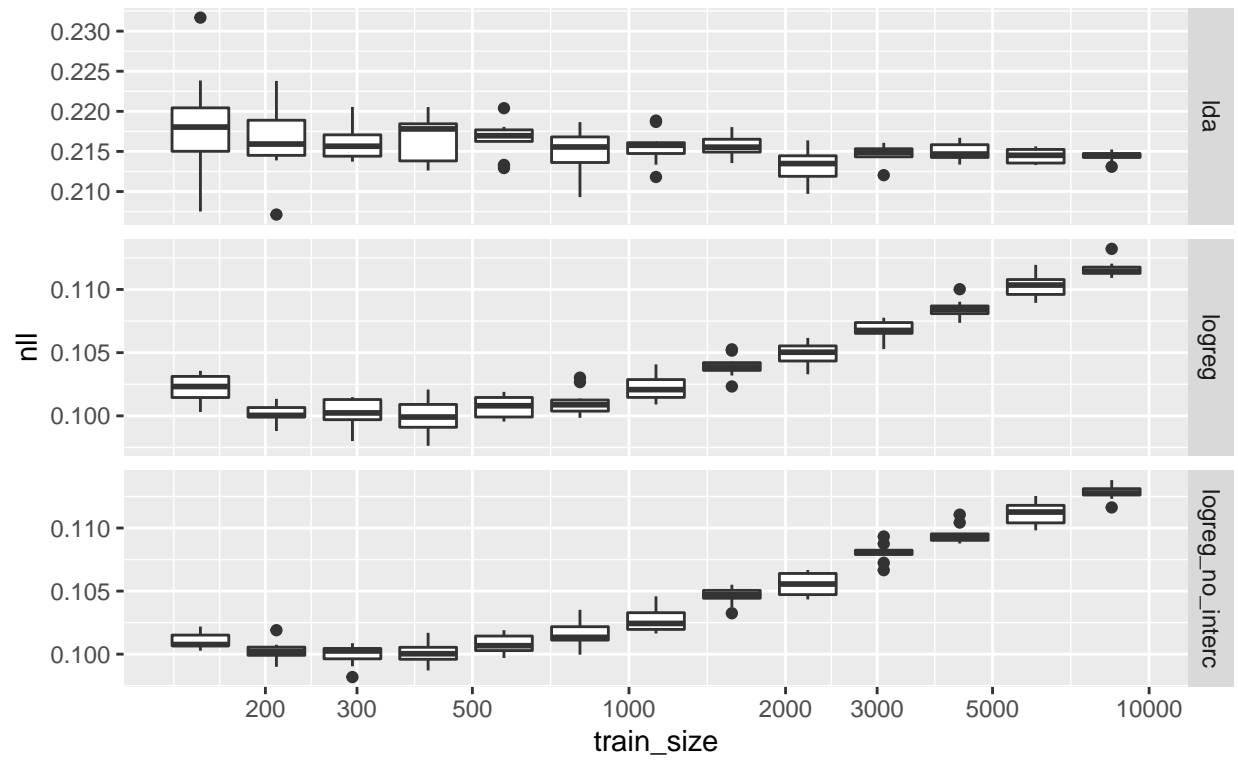


Comparison of accuracy over different combining methods  
for coupling method bc

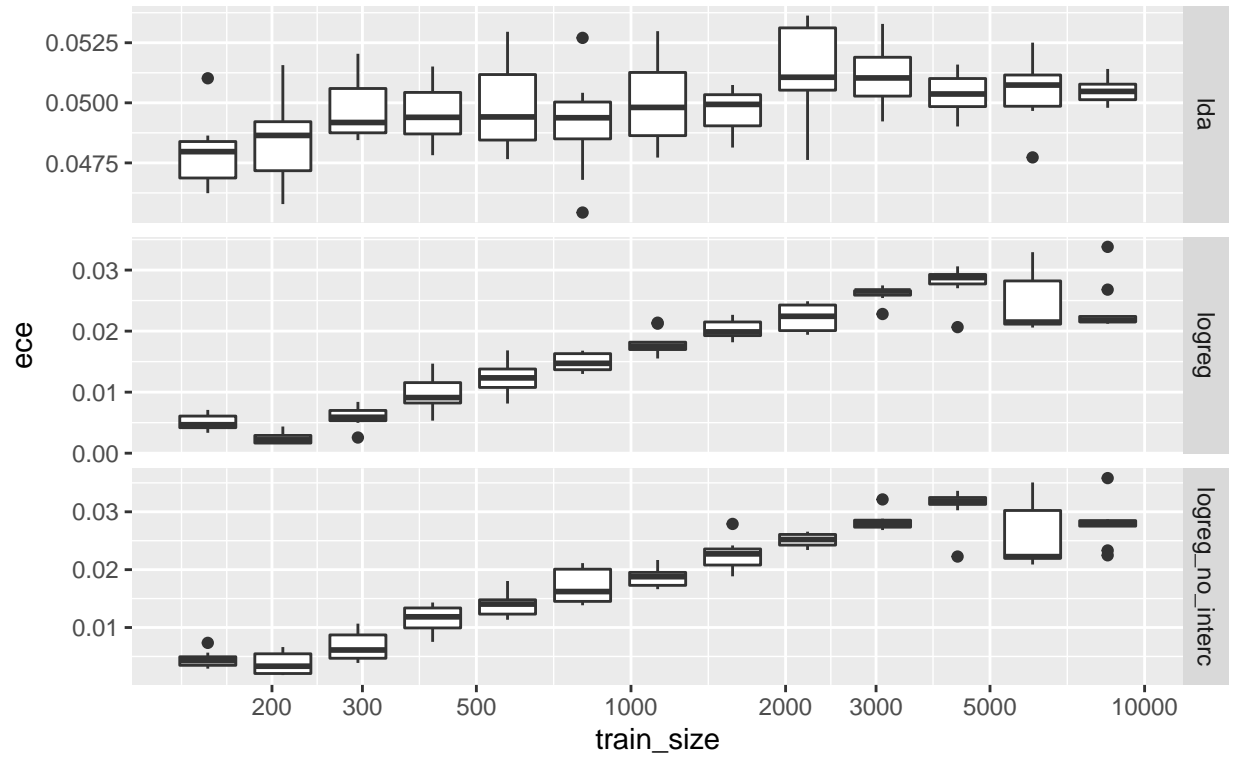




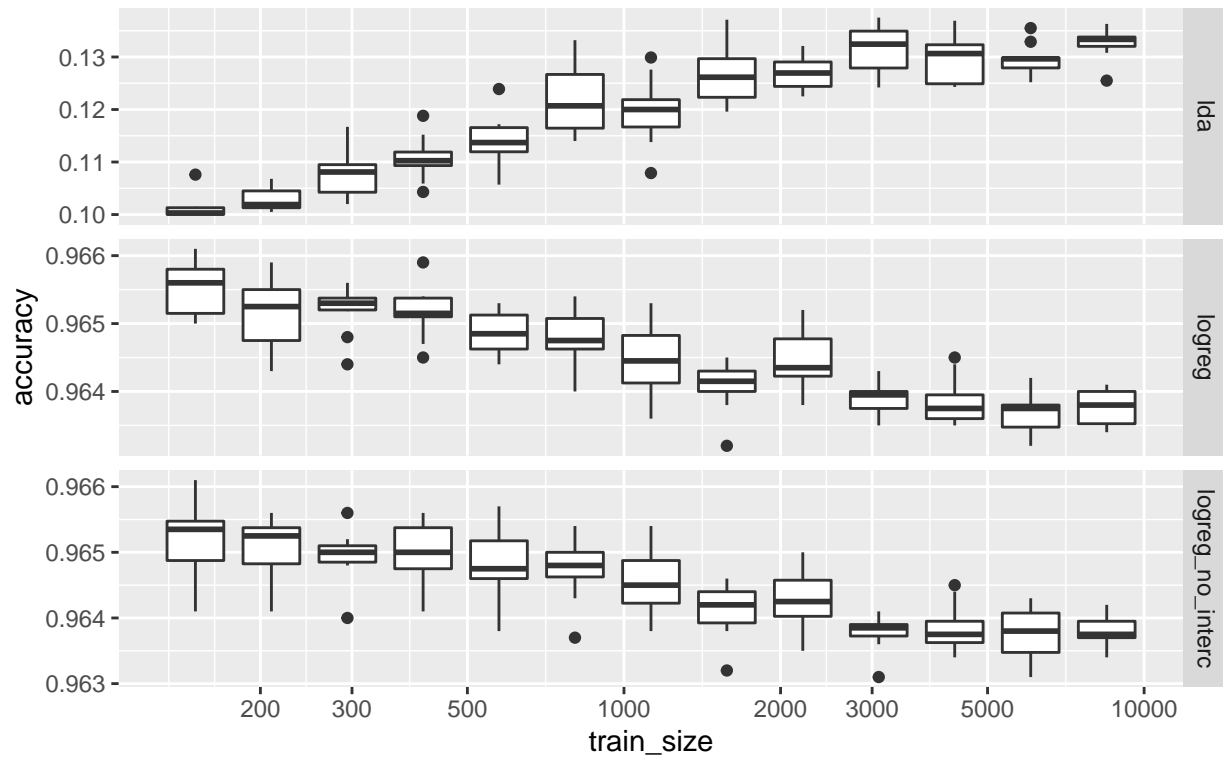
Comparison of nll over different combining methods  
for coupling method bc



Comparison of ece over different combining methods  
for coupling method bc

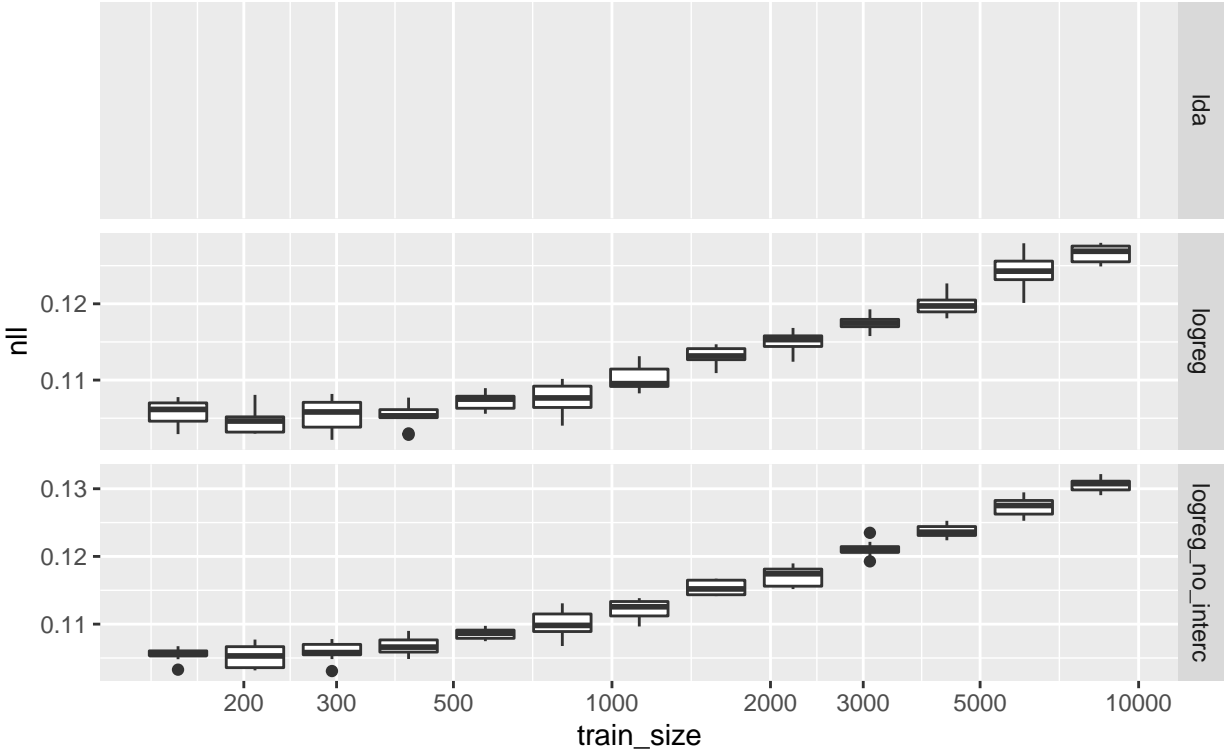


Comparison of accuracy over different combining methods  
for coupling method sbt



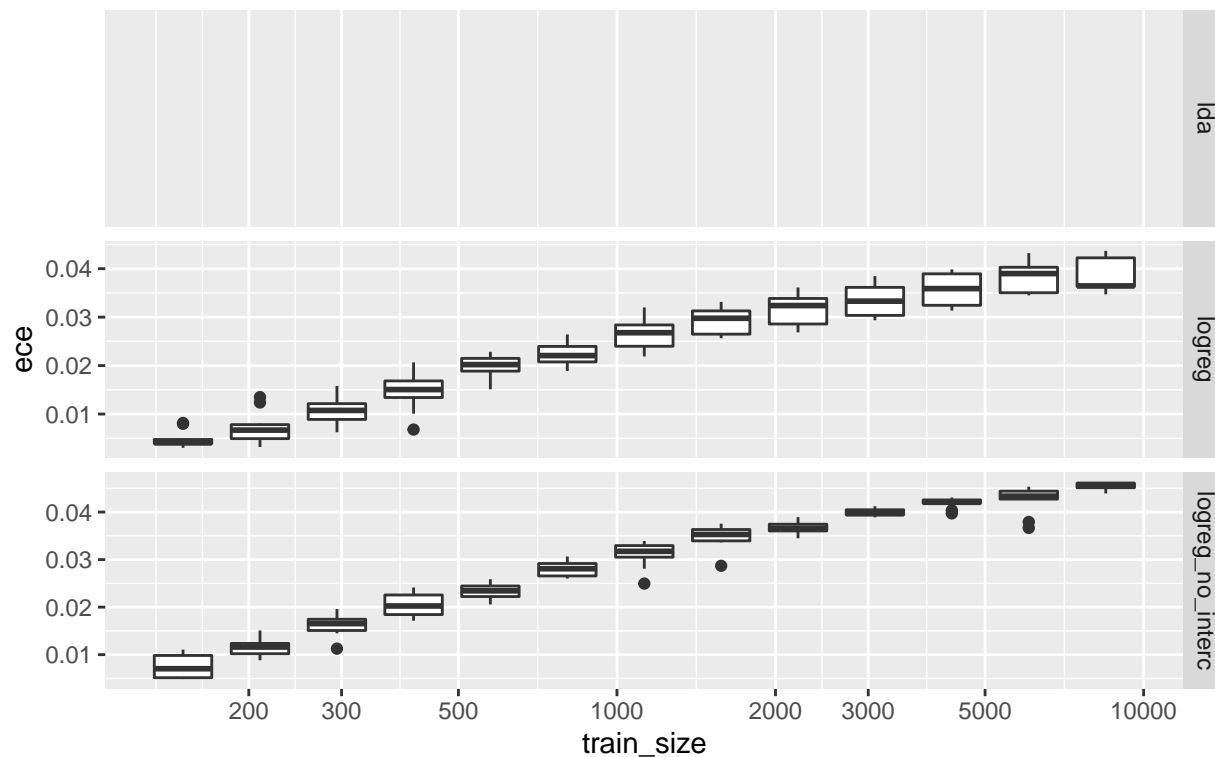
## Warning: Removed 130 rows containing non-finite values (stat\_boxplot).

Comparison of nll over different combining methods  
for coupling method sbt



## Warning: Removed 130 rows containing non-finite values (stat\_boxplot).

## Comparison of ece over different combining methods for coupling method sbt



As we can see on these plots, accuracy varies only slightly over different training set sizes. The differences among median accuracies are only about 0.05% for logistic regression combining methods and a little more for lda. The difference in accuracy of 0.05% represents only 5 samples of the test set of size 10000, therefore we don't consider these differences to provide conclusive information for the choice of training set size.

However, if we look at the metrics negative log likelihood (nll) and estimated calibration error (ece) we can see a clear increasing trend for logreg combining methods and mostly stagnating trend for combining method lda. For logreg combining methods, we can observe a stagnation or even a slight decrease in nll for several smaller training set sizes. This decrease, or stagnation changes into increase at training set size of about 500.

## CIFAR-100

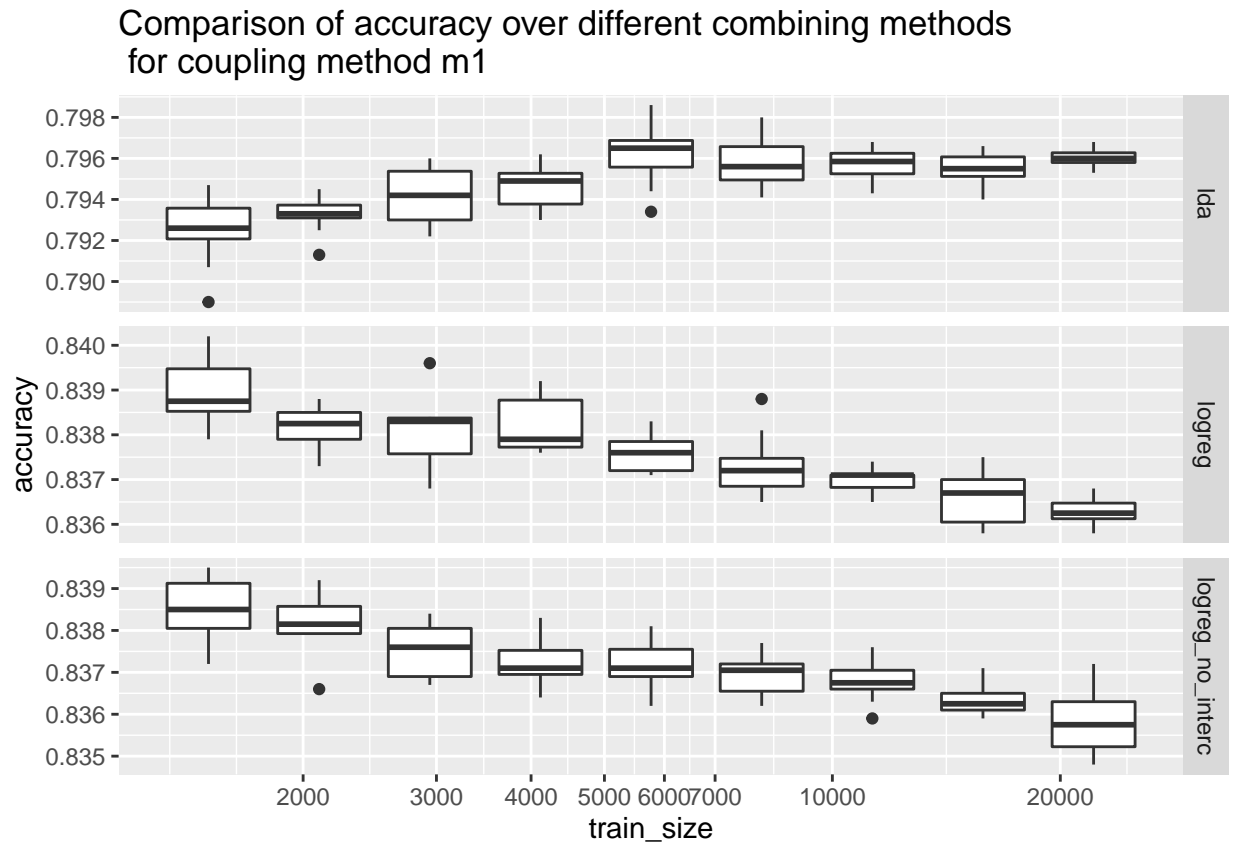
```
ens_metrics_c100 <- read.csv("../data/data_tv_5000_c100/0/exp_subsets_sizes_train_outputs/ens_metrics.csv")
nets_metrics_c100 <- read.csv("../data/data_tv_5000_c100/0/exp_subsets_sizes_train_outputs/net_metrics.csv")
```

```
for (cp_m in unique(ens_metrics_c100$coupling_method))
{
  for (metric in metrics)
  {
    box_plt <- ens_metrics_c100 %>% filter(coupling_method==cp_m) %>%
      ggplot() +
      geom_boxplot(mapping=aes_string(x="train_size", y=metric, group="train_size")) +
      facet_grid(rows=vars(combining_method), scales="free") +
```

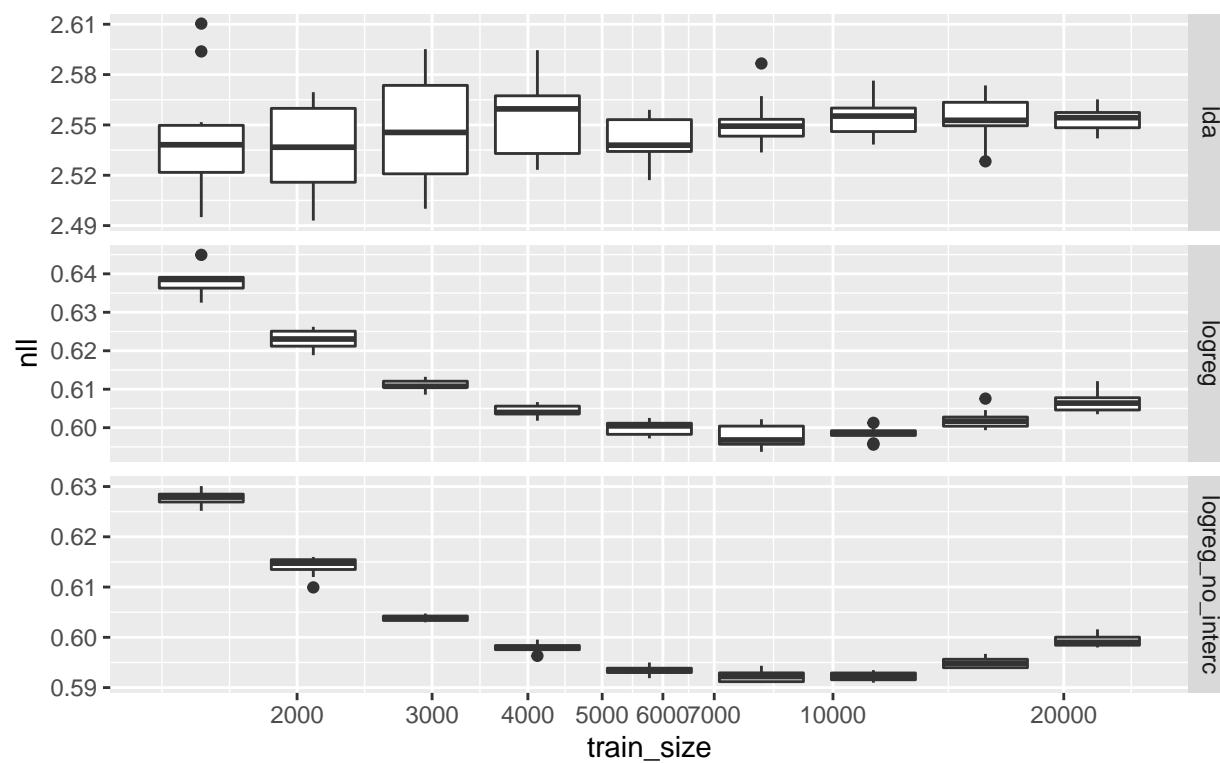
```

scale_x_log10(breaks=log_breaks(n=10)) +
  ggtitle(paste0("Comparison of ", metric, " over different combining methods\n for coupling method m1"))
print(box_plt)
}
}

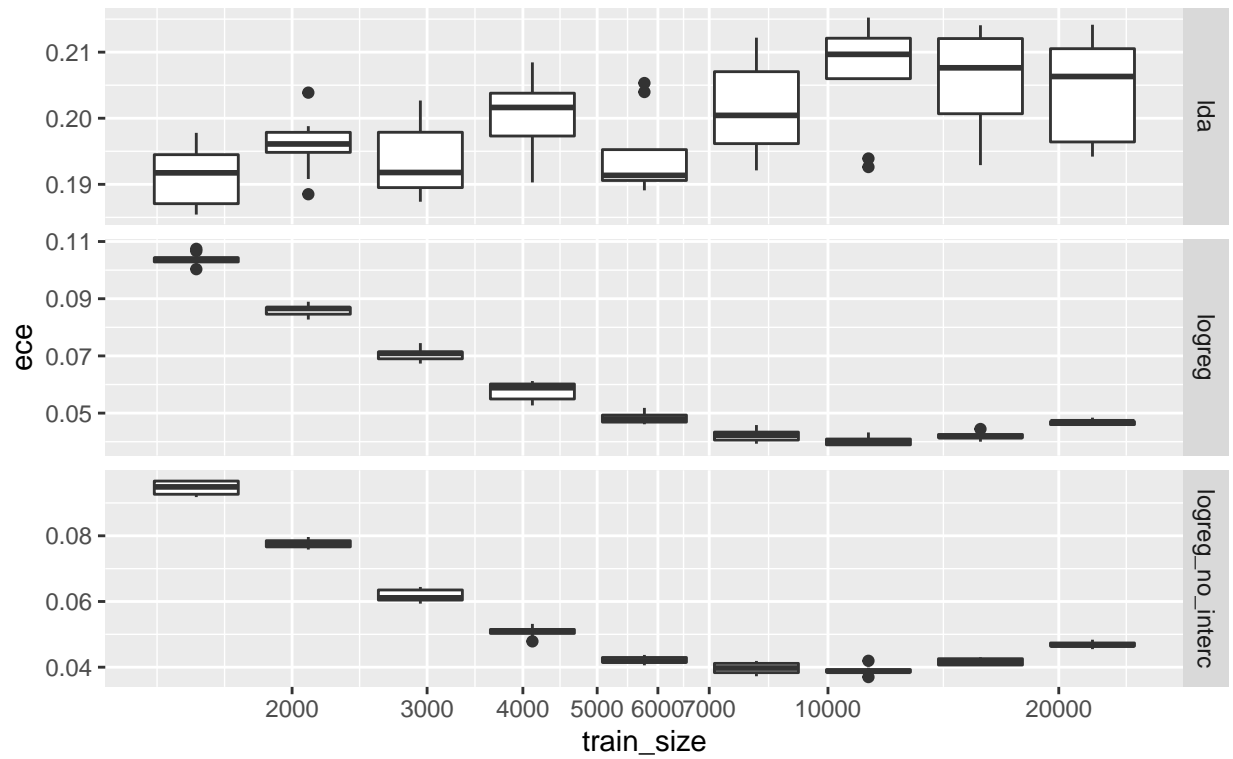
```



Comparison of nll over different combining methods  
for coupling method m1

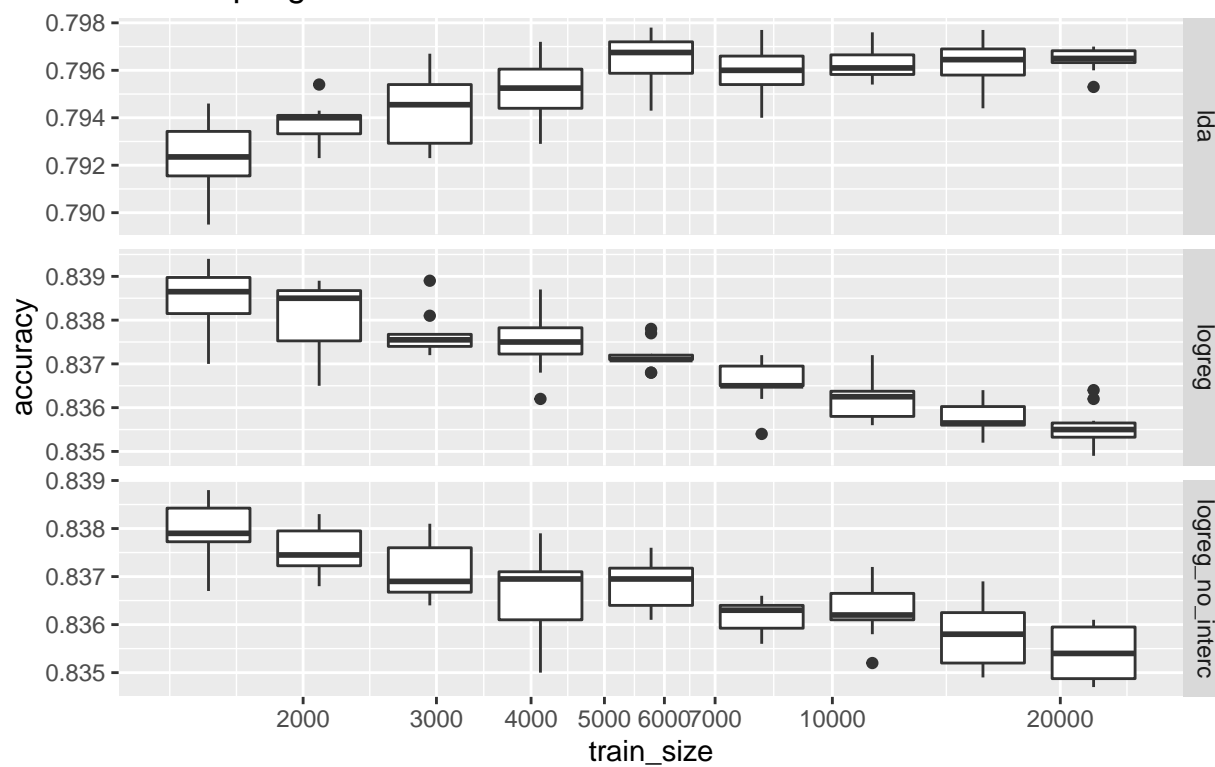


Comparison of ece over different combining methods  
for coupling method m1

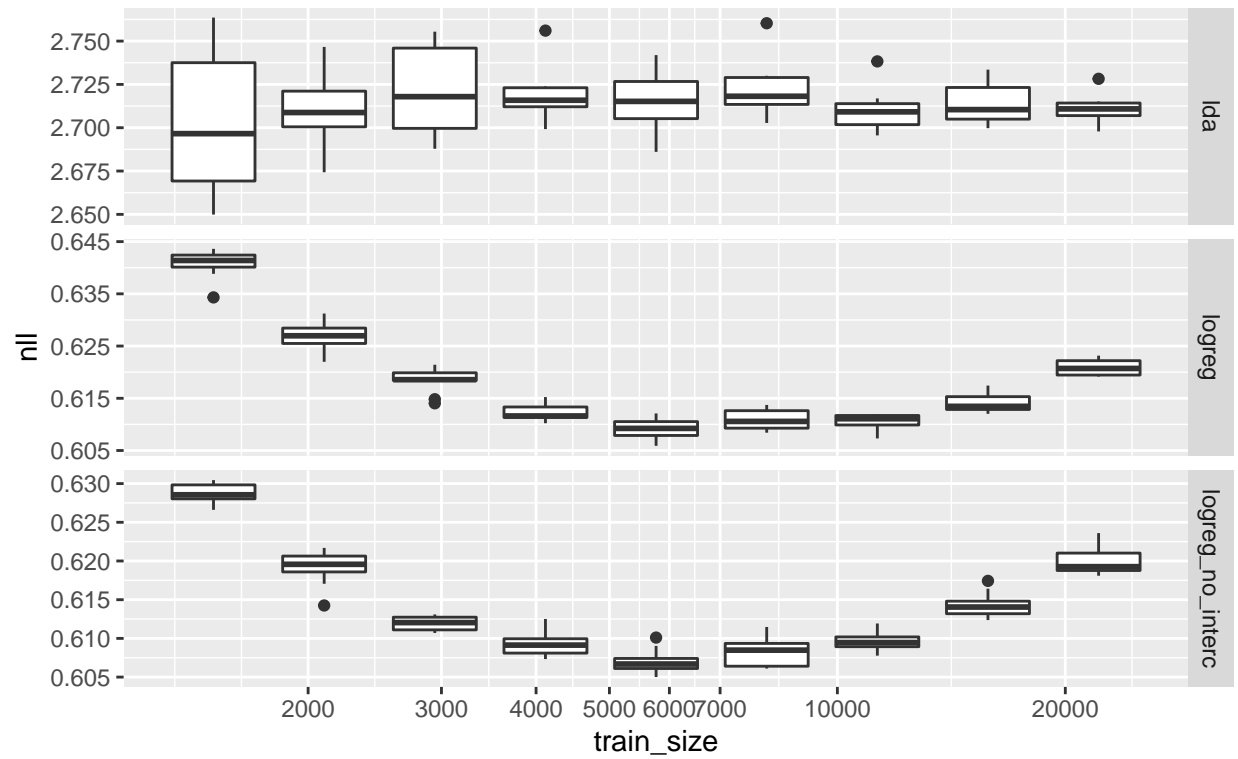




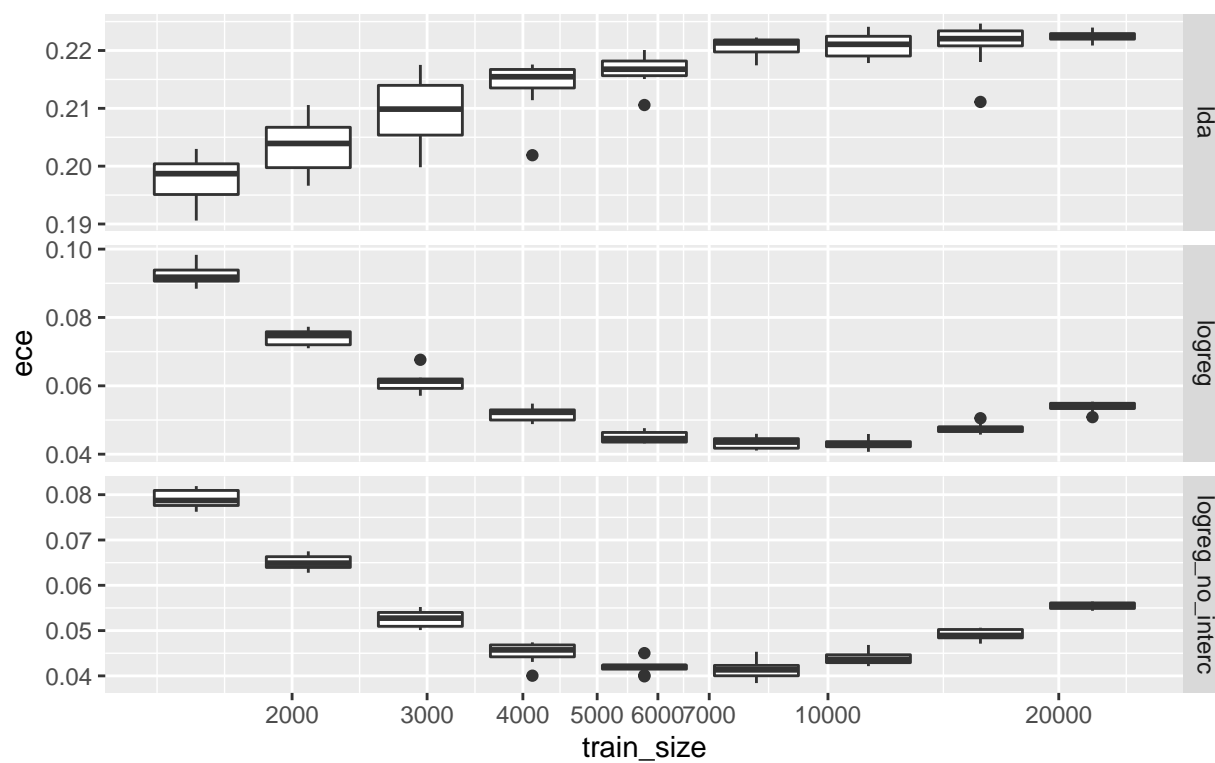
Comparison of accuracy over different combining methods  
for coupling method m2



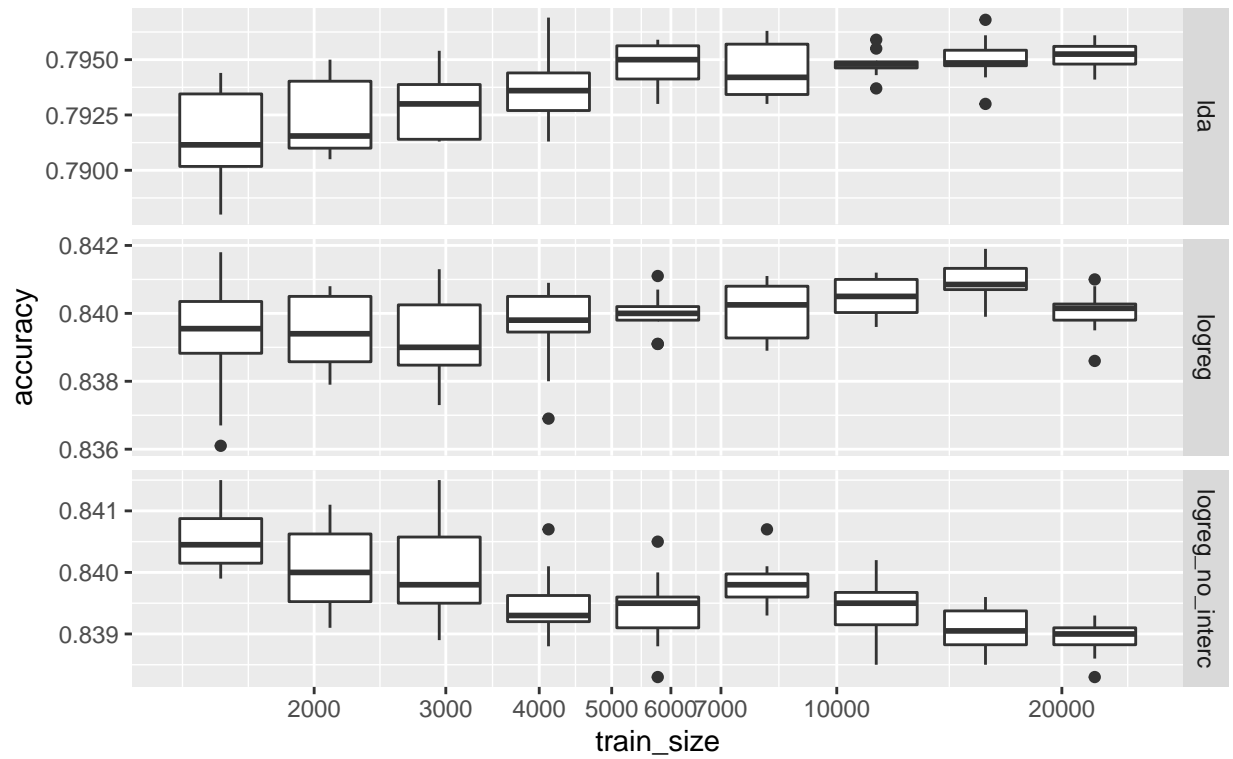
Comparison of nll over different combining methods  
for coupling method m2



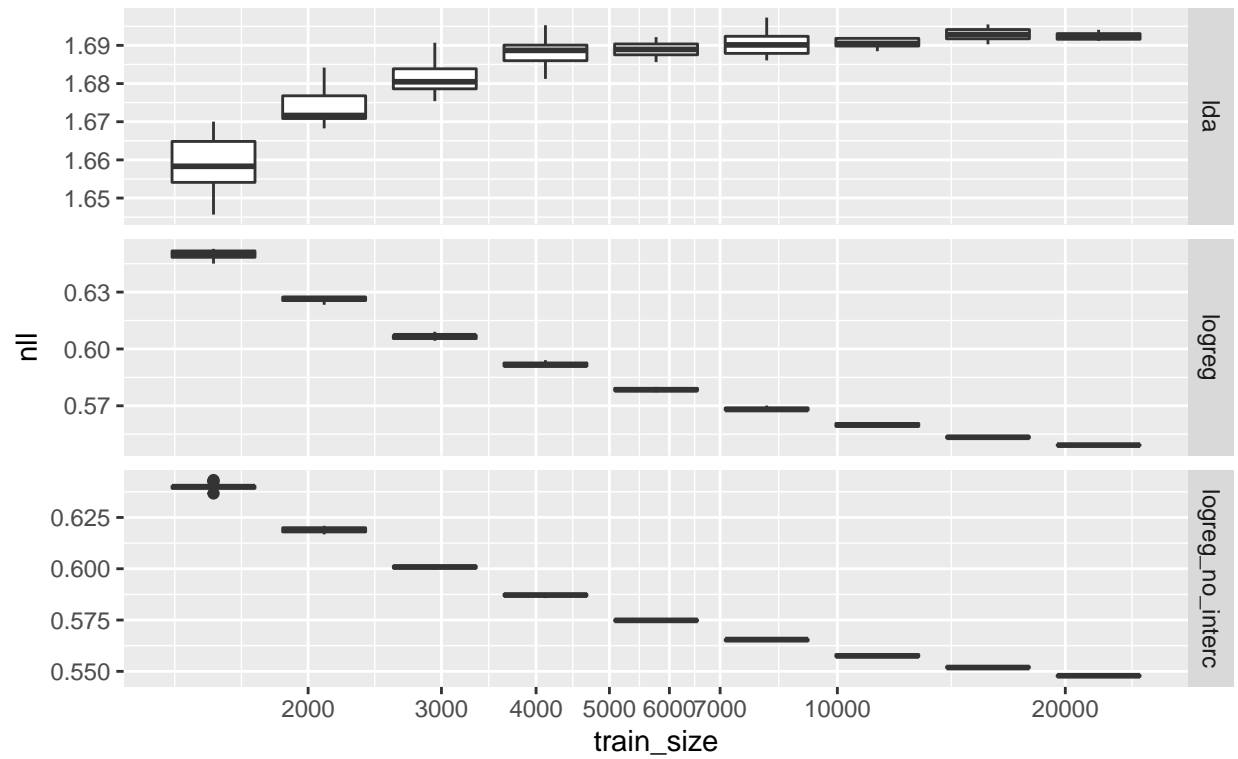
Comparison of ece over different combining methods  
for coupling method m2



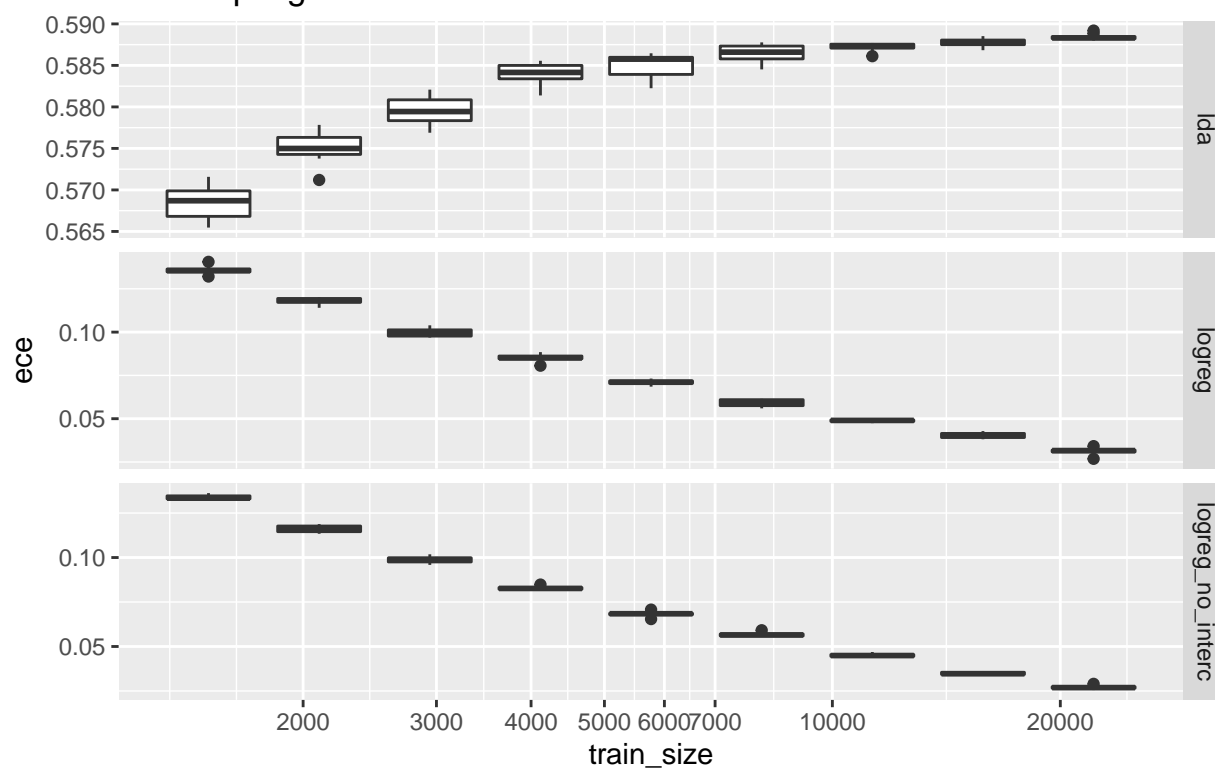
Comparison of accuracy over different combining methods  
for coupling method bc



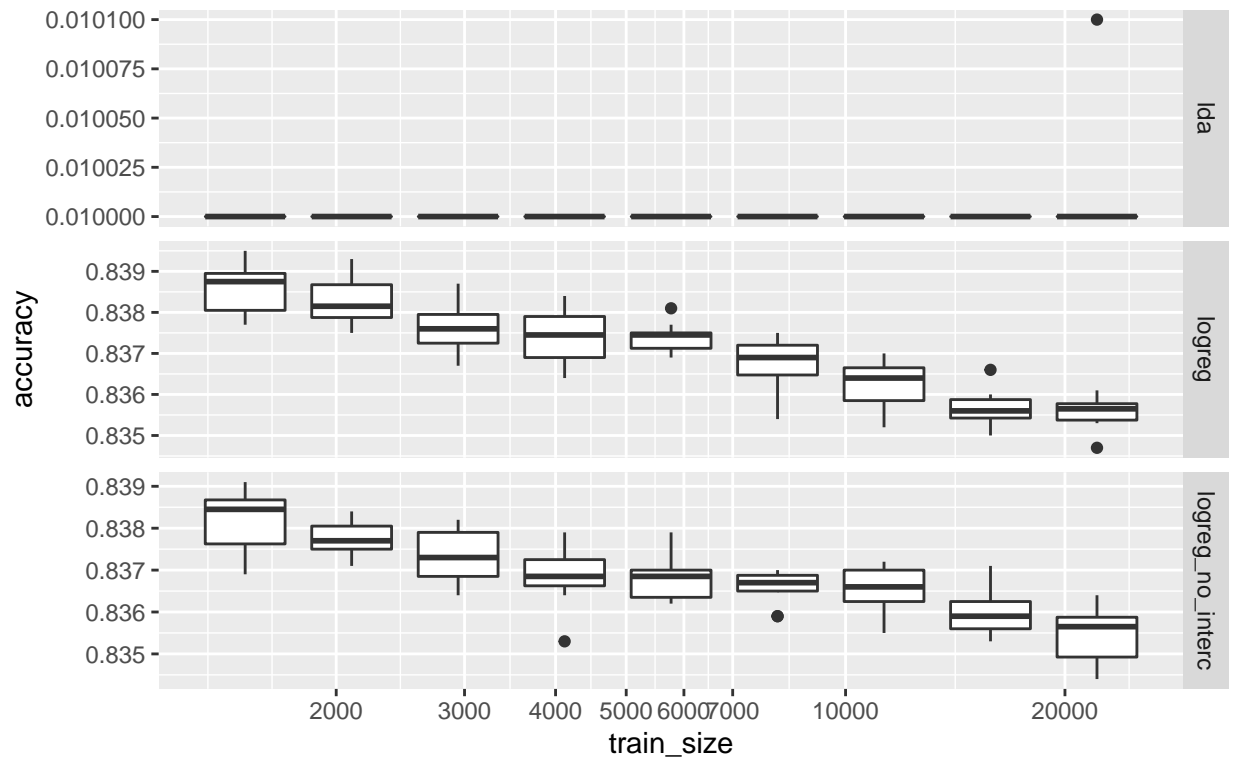
Comparison of nll over different combining methods  
for coupling method bc



Comparison of ece over different combining methods  
for coupling method bc

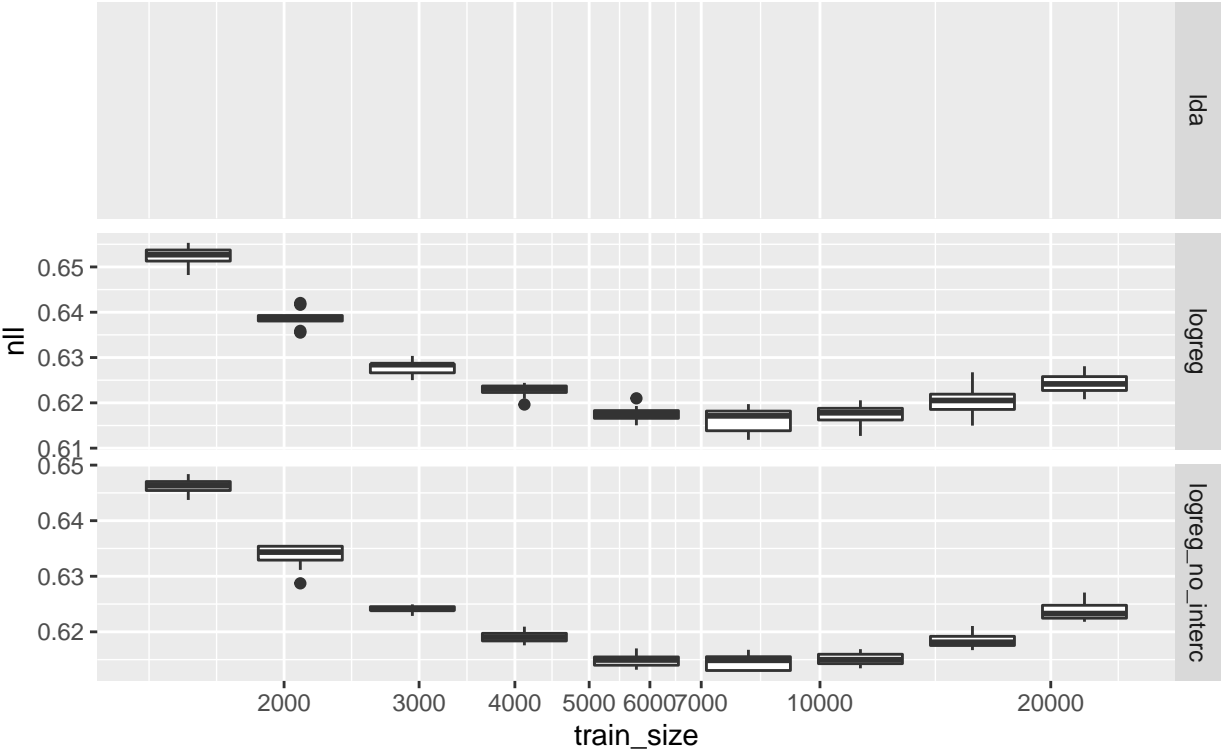


Comparison of accuracy over different combining methods  
for coupling method sbt



## Warning: Removed 90 rows containing non-finite values (stat\_boxplot).

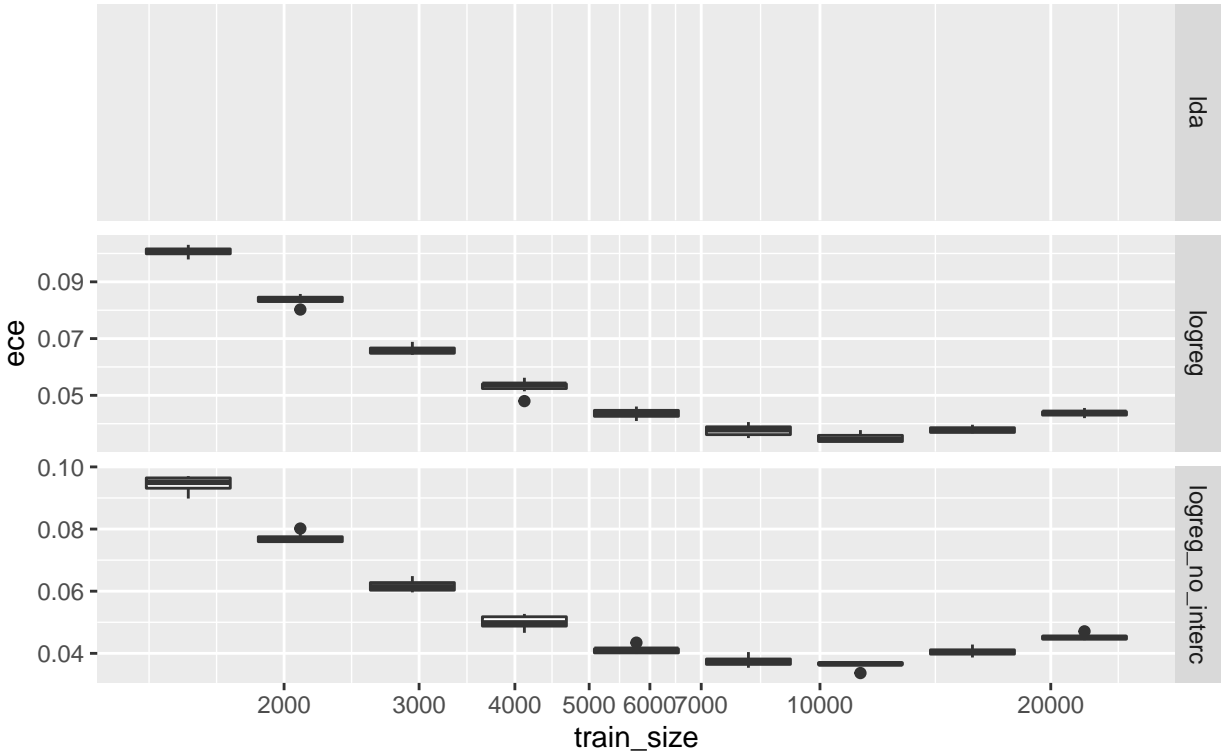
Comparison of nll over different combining methods  
for coupling method sbt



```
## Warning: Removed 90 rows containing non-finite values (stat_boxplot).
```



Comparison of ece over different combining methods  
for coupling method sbt



- Šuch, Ondrej, and Santiago Barreda. 2016. “Bayes Covariant Multi-Class Classification.” *Pattern Recognition Letters* 84: 99–106.
- Šuch, Ondrej, Štefan Benuš, and Andrea Tinajová. 2015. “A New Method to Combine Probability Estimates from Pairwise Binary Classifiers.” *Rmj* 1: 12.
- Wu, Ting-Fan, Chih-Jen Lin, and Ruby C Weng. 2004. “Probability Estimates for Multi-Class Classification by Pairwise Coupling.” *Journal of Machine Learning Research* 5 (Aug): 975–1005.