

Module 2: Supervised learning

Split data into train and test

```
dim(train_x)
```

```
## [1] 106 587
```

```
length(train_y)
```

```
## [1] 106
```

```
dim(test_x)
```

```
## [1] 75 587
```

```
length(test_y)
```

```
## [1] 75
```

LASSO logistic regression

```
# Choose best lambda using CV
beta.lasso <- lasso_fit(
  x = log(train_x + 1), y = train_y,
  tuning = "cv", family = "binomial"
)
```

```
# Features Selected
names(beta.lasso[abs(beta.lasso) > 0])[-1]
```

```
## [1] "COD2"      "COD10"     "NLP1"      "NLP17"     "NLP56"     "NLP82"
## [7] "NLP93"     "NLP104"    "NLP118"    "NLP130"    "NLP144"    "NLP164"
## [13] "NLP172"    "NLP193"    "NLP199"    "NLP222"    "NLP231"    "NLP265"
## [19] "NLP274"    "NLP280"    "NLP297"    "NLP299"    "NLP346"    "NLP362"
## [25] "NLP375"    "NLP382"    "NLP396"    "NLP401"    "NLP409"    "NLP435"
## [31] "NLP451"    "NLP462"    "NLP488"    "NLP533"    "NLP536"    "NLP552"
## [37] "NLP568"    "main_NLP"
```

ALASSO logistic regression

```
# Fit Adaptive LASSO
```

```
beta.lasso <- adaptive_lasso_fit(  
  x = log(train_x + 1), y = train_y,  
  tuning = "cv", family = "binomial"  
)
```

```
# ALASSO features selected - we show the features, please also print out the standardized coefficients so v
```

```
names(beta.lasso[abs(beta.lasso) > 0])[-1]
```

```
## [1] "NLP56"      "NLP93"      "NLP104"     "NLP118"     "NLP222"     "NLP231"  
## [7] "NLP265"     "NLP280"     "NLP297"     "NLP299"     "NLP409"     "NLP536"  
## [13] "main_NLP"
```

```
# LASSO features selected
```

```
names(beta.lasso[abs(beta.lasso) > 0])[-1]
```

```
## [1] "COD2"      "COD10"     "NLP1"      "NLP17"     "NLP56"     "NLP82"  
## [7] "NLP93"     "NLP104"    "NLP118"    "NLP130"    "NLP144"    "NLP164"  
## [13] "NLP172"    "NLP193"    "NLP199"    "NLP222"    "NLP231"    "NLP265"  
## [19] "NLP274"    "NLP280"    "NLP297"    "NLP299"    "NLP346"    "NLP362"  
## [25] "NLP375"    "NLP382"    "NLP396"    "NLP401"    "NLP409"    "NLP435"  
## [31] "NLP451"    "NLP462"    "NLP488"    "NLP533"    "NLP536"    "NLP552"  
## [37] "NLP568"    "main_NLP"
```

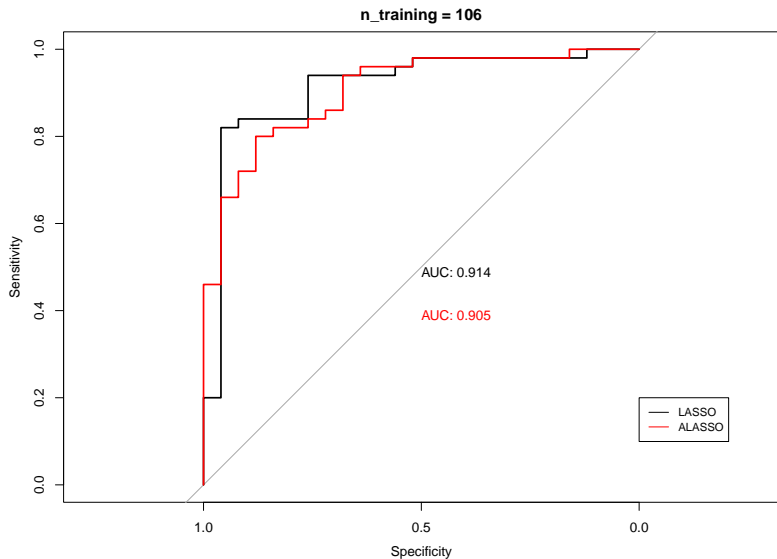
Get model predictions + ROC curve

```
# Prediction on testing set (LASSO)  
y_hat.lasso <- linear_model_predict(  
  beta = beta.lasso, x = log(test_x + 1),  
  probability = TRUE  
)
```

```
# Prediction on testing set (ALASSO)  
y_hat.alasso <- linear_model_predict(  
  beta = beta.alasso, x = log(test_x + 1),  
  probability = TRUE  
)
```

```
roc.lasso <- roc(test_y, y_hat.lasso)  
roc.alasso <- roc(test_y, y_hat.alasso)  
# as expected alasso selects less features
```

LASSO vs. ALASSO



LASSO vs. ALASSO at $\text{FPR} = 0.10$

```
roc_full.lasso <- get_roc(y_true = test_y, y_score = y_hat.lasso) %>% data.frame()
get_roc_parameter(0.1, roc_full.lasso)
```

```
##      cutoff pos.rate FPR  TPR      PPV      NPV      F1
## 1 0.8573308 0.5933333 0.1 0.84 0.9438202 0.7377049 0.8888889
```

```
roc_full.lasso <- get_roc(y_true = test_y, y_score = y_hat.lasso) %>% data.frame()
get_roc_parameter(0.1, roc_full.lasso)
```

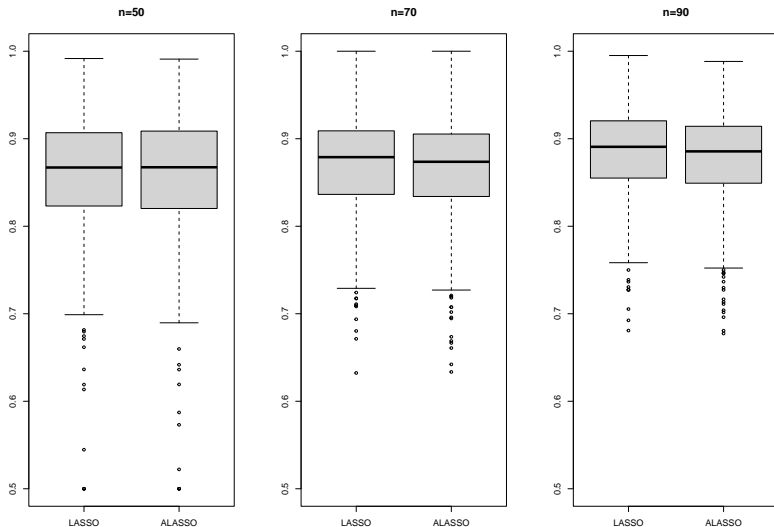
```
##      cutoff pos.rate FPR  TPR      PPV      NPV      F1
## 1 0.9228109 0.5133333 0.1 0.72 0.9350649 0.6164384 0.8135593
```

LASSO vs. ALASSO with different training set size

- ▶ Randomly sample training size = 50, 70, 90
- ▶ Use the remaining data as the test set
- ▶ Repeat 600 times

```
auc_supervised <- validate_supervised(  
  dat = labeled_data, nsim = 600,  
  n.train = c(50, 70, 90)  
)
```


LASSO vs. ALASSO with different training set size



Random Forest and SVM

```
# Random forest
model_rf <- rfsrc(y ~ ., data = data.frame(y = train_y, x = train_x))
y_hat.rf <- predict(model_rf,
                    newdata = data.frame(x = test_x))$predicted
roc.rf <- roc(test_y, y_hat.rf)
# Use the Phecap functions to compute lasso, alssso etc and be consistnet
# clean the helper functions to only have what you need
# same for packages
# Make sure the slides match the markdown - please update
```

```
# SVM
model_svm <- SVMmaj::svmmaj(X = train_x, y = train_y)
y_hat.svm <- predict(model_svm, test_x)
roc.svm <- roc(test_y, y_hat.svm)
```

ROC curves

