

## Module 4: Alternative approaches

## 2-step Semi-supervised Approach

1. Regress the surrogate on the features with penalized least square to get the direction of  $\beta$ .

```
# COD + NLP + HU.
x <- ehr_data_transformed %>% select(starts_with("health")
  starts_with("COD") | starts_with("NLP"))
S <- ehr_data_transformed$main_NLP

# Step 1.
beta_step1 <- adaptive_lasso_fit(
  y = S, # surrogate
  x = x, # all X
  family = "gaussian",
  tuning = "cv"
)
```

## 2-step Semi-supervised Approach

1. Regress the surrogate on the features with penalized least squares to get the direction of  $\beta$ .
2. Regress the outcome on the linear predictor to get the intercept and multiplier for the  $\beta$ .

```
# Linear predictor without intercept.
bhatx <- linear_model_predict(beta = beta_step1, x = as.matrix(x))

# Step 2.
step2 <- glm(
  train_y ~ bhatx[train_data$patient_id] + S[train_data$patient_id],
  family = "binomial"
)
beta_step2 <- coef(step2)
beta_step2
```

```
##                (Intercept) bhatx[train_data$patient_id]
##                -0.3504483                1.2620270
##      S[train_data$patient_id]
##                0.4012988
```

```
# Recover beta.
beta <- beta_step2[2] * beta_step1
```

# Compare selected features

```
# LASSO.
```

```
names(beta_lasso[!beta_lasso == 0])[-1]
```

```
## [1] "NLP93" "NLP104" "NLP304"
```

```
## [4] "main_NLP" "healthcare_utilization"
```

```
# ALASSO.
```

```
names(beta_lasso[!beta_lasso == 0])[-1]
```

```
## [1] "NLP304" "main_NLP" "healthcare_utilization"
```

# Compare selected features

```
# PheCAP.  
feature_selected
```

```
## Feature(s) selected by surrogate-assisted feature extraction (SAFE)  
## [1] "main_ICD" "main_NLP" "NLP56" "NLP93" "NLP274" "NLP306"
```

```
# Two Step.  
names(beta[!beta == 0])[-1]
```

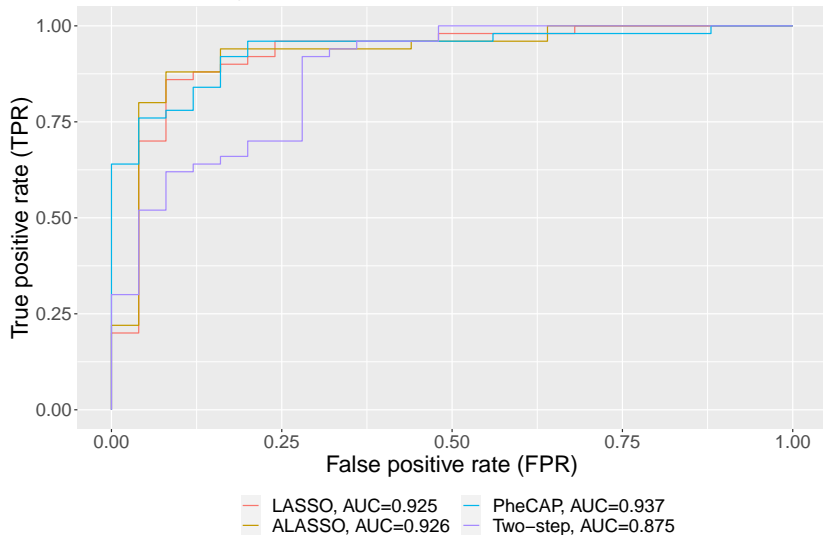
```
## [1] "COD3" "COD6" "COD8" "COD10" "NLP5" "NLP7" "NLP9" "NLP10"  
## [9] "NLP18" "NLP21" "NLP24" "NLP26" "NLP28" "NLP29" "NLP31" "NLP33"  
## [17] "NLP40" "NLP44" "NLP50" "NLP53" "NLP56" "NLP59" "NLP60" "NLP62"  
## [25] "NLP68" "NLP69" "NLP73" "NLP74" "NLP76" "NLP81" "NLP92" "NLP93"  
## [33] "NLP95" "NLP96" "NLP98" "NLP103" "NLP104" "NLP116" "NLP120" "NLP127"  
## [41] "NLP140" "NLP146" "NLP150" "NLP160" "NLP161" "NLP172" "NLP176" "NLP178"  
## [49] "NLP179" "NLP183" "NLP189" "NLP190" "NLP192" "NLP195" "NLP199" "NLP200"  
## [57] "NLP202" "NLP203" "NLP206" "NLP207" "NLP212" "NLP218" "NLP220" "NLP225"  
## [65] "NLP231" "NLP237" "NLP243" "NLP246" "NLP250" "NLP266" "NLP274" "NLP281"  
## [73] "NLP287" "NLP288" "NLP291" "NLP294" "NLP295" "NLP298" "NLP299" "NLP300"  
## [81] "NLP301" "NLP302" "NLP304" "NLP306" "NLP309" "NLP318" "NLP321" "NLP326"  
## [89] "NLP334" "NLP338" "NLP339" "NLP342" "NLP343" "NLP347" "NLP349" "NLP350"  
## [97] "NLP351" "NLP357" "NLP359" "NLP361" "NLP362" "NLP365" "NLP369" "NLP380"  
## [105] "NLP387" "NLP395" "NLP396" "NLP403" "NLP405" "NLP407" "NLP417" "NLP431"  
## [113] "NLP434" "NLP435" "NLP436" "NLP437" "NLP440" "NLP446" "NLP447" "NLP451"  
## [121] "NLP452" "NLP456" "NLP457" "NLP463" "NLP465" "NLP468" "NLP473" "NLP475"  
## [129] "NLP482" "NLP483" "NLP484" "NLP486" "NLP490" "NLP495" "NLP500" "NLP507"  
## [137] "NLP510" "NLP516" "NLP517" "NLP523" "NLP529" "NLP533" "NLP534" "NLP536"  
## [145] "NLP539" "NLP541" "NLP544" "NLP547" "NLP548" "NLP549" "NLP554" "NLP560"  
## [153] "NLP561" "NLP562" "NLP564" "NLP574"
```

# ROC

```
mu <- beta_step2[1] +  
  as.numeric(as.matrix(x[test_data$patient_id, ]  
  %*% beta[-1]) +  
  as.numeric(beta_step2[3] %*% S[test_data$patient_id])  
  
# Expit.  
y_hat_twostep <- plogis(mu)  
  
roc_twostep <- roc(test_y, y_hat_twostep)
```

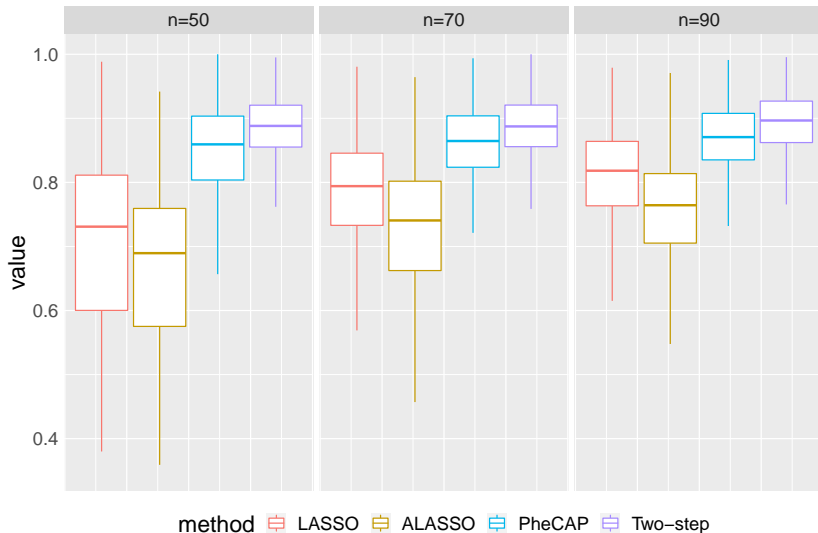
# ROC

The operating receiver characteristic (ROC) curve



# Model Evaluation

Area under the ROC curve (AUC) from 600 simulations





# MAP

```
# Use un-transformed data; MAP requires sparse matrix.
# Create sparse matrix for surrogates.
data_fit <- sparsify(
  PheCAP::ehr_data %>%
    select(main_ICD, main_NLP) %>%
    rename(ICD = main_ICD) %>% data.table()
)

# Create sparse matrix for HU.
note <- Matrix(
  PheCAP::ehr_data$healthcare_utilization,
  ncol = 1, sparse = TRUE
)
model_map <- MAP(mat = data_fit, note = note, full.output = TRUE)
```

```
## #####
## MAP only considers patients who have note count data and
##       at least one nonmissing variable!
## ####
## Here is a summary of the input data:
## Total number of patients: 10000
##   ICD main_NLP note   Freq
## 1 YES      YES   YES 10000
## ####
```

```
y_hat_map <- model_map$scores[data$validation_set]
roc_map <- roc(test_y, y_hat_map)
```

# ROC

The operating receiver characteristic (ROC) curve

