

## Module 4: Alternative approaches

## 2-step Semi-supervised Approach

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

```
x <- all_x %>% select(starts_with("COD") | starts_with("NLP"))
S <- ehr_data$main_ICDNLP
```

*# Step 1*

```
beta.step1 <- adaptive_lasso_fit(
  y = S, # surrogate
  x = x, # all X
  family = "gaussian",
  tuning = "cv"
)
```

*# Features selected*

```
names(beta.step1[abs(beta.step1) > 0])[-1]
```

```
##      [1] "COD6"      "COD8"      "COD10"     "NLP5"      "NLP7"      "NLP1"
##      [2] "NLP24"     "NLP28"     "NLP31"     "NLP33"     "NLP44"     "NLP45"
```

## 2-step Semi-supervised Approach

- i) Regress the surrogate on the features with penalized least square to get the direction of  $\beta$ .
- (ii) 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.mat
```

```
# Step 2
```

```
step2 <- glm(train_y ~ bhatx[train_data$patient_id] + S[train_data$patient_id]  
  health_count[train_data$patient_id])
```

```
beta_step2 <- coef(step2)
```

```
beta_step2
```

```
##
```

```
(Intercept)
```

```
bhatx[train_c
```

```
##
```

```
0.80767522
```

```
##
```

```
S[train_data$patient_id] health_count[train_c
```

```
##
```

```
0.14335393
```