# PS10 Finley

Caden Finley

April 2023

## 1 Workflow

library(tidyverse) library(tidymodels) library(magrittr) library(modelsummary) library(rpart) library(e1071) library(kknn) library(nnet) library(kernlab) library(readr)

set.seed(100)

income ¡- $read_c sv$("$http://archive.ics.uci.edu/ml/machine-learning-databases/adult/adult.data", col_n ames = FALSE) names(income) < -c$("age","workclass","fnlwgt","edu

Clean up the data Drop unnecessary columns income Make sure continuous variables are formatted as numeric income Make sure discrete variables are formatted as factors income Combine levels of factor variables that currently have too many levels income Advanced = c("Masters","Doctorate","Prof-school"), Bachelors = c("Bachelors"), SomeCollege = c("Some-college","Assoc-acdm","Assoc-voc"), HSgrad = c("HS-grad","12th"), HSdrop = c("11th","9th","7th-8th","1st-4th","10th","5th-6th","Preschool") ), $marital.status = fct_c ollapse(marital.status, Married = c("Married-civ-spouse","Married-spouse-absent","Married-AF-spouse"), Divorced = c("Divorced","Separated"), Widowed = c("Widowed"), NeverMarried = c("Never-married")), race = fct_c ollapse(race, White = c("White"), Black = c("Black"), Asian = c("Asian-Pac-Islander"), Other = c("Other","Amer-Indian-Eskimo")), workclass = fct_c ollapse(workclass, Private = c("Private"), SelfEmp = c("Self-emp-not-inc","Self-emp-inc"), Gov = c("Federal-gov","Local-gov","State-gov"), Other = c("Without-pay","Never-worked","?")), occupation = fct_c ollapse(occupation, BlueCollar = c("?","Craft-repair","Farming-fishing","Handlers-cleaners","Machine-op-inspct","Transport-moving"), WhiteCollar = c("Adm-clerical","Exec-managerial","Prof-specialty","Sales","Tech-support"), Services = c("Armed-Forces","Other-service","Priv-house-serv","Protective-serv")))$

tidymodels time! $income_s plit < -initial_s plit(income, prop = 0.8) income_t rain < -training(income_s plit) income_t est < -testing(income_s plit)$

logistic regression print('Starting LOGIT') set up the task and the engine $tune_l ogit_s pec < -logistic_r eg(penalty = tune(), tuning parameter mixture = 11 = lasso, 0 = ridge) set_e ngine("glmnet") set_m ode("classification")$

define a grid over which to try different values of the regularization parameter lambda $lambda_g rid < -grid_r egular(penalty(), levels = 50)$

3-fold cross-validation $rec_f olds < -vfold_c v(income_t rain, v = 3)$

Workflow $rec_wf < -workflow()add_model(tune_logit_spec)add_formula(high.earner\ education +$
$marital.status + race + workclass + occupation + relationship + sex + age +$
$capital.gain + capital.loss + hours)$

Tuning results $rec_res < -rec_wftune_grid(resamples = rec_folds, grid =$
$lambda_grid)$

what is the best value of lambda? $top_acc < -show_best(rec_res, metric =$
$"accuracy")best_acc < -select_best(rec_res, metric = "accuracy")final_logit_lasso <$
$-finalize_workflow(rec_wf, best_acc)$

print('*********** LOGISTIC REGRESSION **************') $logit_test <$
$-last_fit(final_logit_lasso, income_split)collect_metrics()$

$logit_testtop_acc$

combine results into a nice tibble (for later use) $logit_ans < -top_acclogit_ansmutate(alg =$
$"logit")$

tree model print('Starting TREE') set up the task and the engine $tune_tree_spec <$
$-decision_tree(min_n = tune(), tuningparametertree_depth = tune(), tuningparametercost_complexity =$
$tune(), tuningparameter)set_engine("rpart")set_mode("classification")$

define a set over which to try different values of tuning parameters $tree_grid <$
$-grid_latin_hypercube(min_n(), tree_depth(), cost_complexity(), size = 20)$

3-fold cross-validation $tree_folds < -vfold_cv(income_train, v = 3)$

Workflow $tree_wf < -workflow()add_model(tune_tree_spec)add_formula(high.earner\ education +$
$marital.status + race + workclass + occupation + relationship + sex + age +$
$capital.gain + capital.loss + hours)$

Tuning results $tree_res < -tree_wftune_grid(resamples = tree_folds, grid =$
$tree_grid)$

what are the best values of the tuning parameters? $top_acc_tree < -show_best(tree_res, metric =$
$"accuracy")best_acc_tree < -select_best(tree_res, metric = "accuracy")final_tree <$
$-finalize_workflow(tree_wf, best_acc_tree)$

print('*********** DECISION TREE **************') $tree_test < -last_fit(final_tree, income_split)collect_$

$tree_testtop_acc_tree$

combine results into a nice tibble (for later use) $tree_ans < -top_acc_treetree_ansmutate(alg =$
$"tree")$

k-nearest neighbors print('Starting KNN') set up the task and the engine
$tune_knn_spec < -nearest_neighbor(weight_func = tune(), tuningparameterneighbors =$
$tune()tuningparameter)set_engine("kknn")set_mode("classification")$

define a set over which to try different values of tuning parameters $knn_grid <$
$-grid_latin_hypercube(weight_func(), neighbors(), size = 20)$

3-fold cross-validation $knn_folds < -vfold_cv(income_train, v = 3)$

Workflow $knn_wf < -workflow()add_model(tune_knn_spec)add_formula(high.earner\ education +$
$marital.status + race + workclass + occupation + relationship + sex + age +$
$capital.gain + capital.loss + hours)$

Tuning results $knn_res < -knn_wftune_grid(resamples = knn_folds, grid =$
$knn_grid)$

what are the best values of the tuning parameters? $top_acc_knn < -show_best(knn_res, metric =$
$"accuracy")best_acc_knn < -select_best(knn_res, metric = "accuracy")final_knn <$
$-finalize_workflow(knn_wf, best_acc_knn)$

print('*********** K-NEAREST NEIGHBORS *************') $\text{knn}_t est < -last_f it(final_k nn, income_s plit) collect_m etrics()$

$\text{knn}_t esttop_a cc_k nn$

combine results into a nice tibble (for later use) $\text{knn}_a ns < -top_a cc_k nn$

SVM

print('Starting SVM') set up the task and the engine $tune_s vm_s pec < -svm_p oly(degree = tune(), tuning parameter : degree of polynomial kernel scale = tune(), tuning parameter : scaling factor for kernel cost = tune(), tuning parameter : cost of violation epsilon = 0.1 epsilon for epsilon-insensitive loss function) set_e ngine("LiblineaR") set_m ode("classification")$

define a grid over which to try different values of tuning parameters $svm_g rid < -grid_t une(degree = seq(1, 5, by = 1), try degrees from 1 to 5 scale = seq(0.1, 1, by = 0.1), try scaling factors from 0.1 to 1 cost = 10^s eq(-3, 3, by = 1) try costs from 0.001 to 1000)$

Workflow $svm_w f < -workflow() add_m odel(tune_s vm_s pec) add_f ormula(high.earner\ education + marital.status + race + workclass + occupation + relationship + sex + age + capital.gain + capital.loss + hours)$

Tuning results $svm_r es < -svm_w f tune_g rid(resamples = rec_f olds, grid = svm_g rid)$

what is the best combination of tuning parameters? $top_a cc_s vm < -show_b est(svm_r es, metric = "accuracy") best_a cc_s vm < -select_b est(svm_r es, metric = "accuracy") final_s vm < -finalize_w orkflow(svm_w f, best_a cc_s vm) print('***********SUPPORT VECTOR MACHINE (SVM) ************') svm_t est < -last_f it(final_s vm, income_s plit) collect_m etrics()$

$svm_t esttop_a cc_s vm$

combine results into a nice tibble (for later use) $svm_a ns < -top_a cc_s vm svm_a ns mutate(alg = "svm")$

combine answers $all_a ns < -bind_r ows(logit_a ns, tree_a ns, nnet_a ns, knn_a ns, svm_a ns) data summary_d f(all_a ns$

## 2 Explanation

I apologize for the messy code and lack of answers, I was not sure how to interpret some of what was required in the problem set. I had issues with the time constraint this time so I could not get it presented like I wanted to. Alongside that, getting this code to translate over in LaTex was giving me trouble, so again I apologize for the sloppy work.