

Untitled

2023-06-01

Libraries

DATA

<pre>library(reticulate) use_python("C:/Users/tatai/AppData/Local/Programs/Python/Python311/python.exe")</pre>
<pre>import pandas as pd</pre>
<pre>data = pd.read_csv("../Downloads/base_exam/exam1.csv")</pre>
<pre>fraud <- pydata</pre>

Preprocessing

Adding age , hour , day columns because they tend to have influence on the fraud

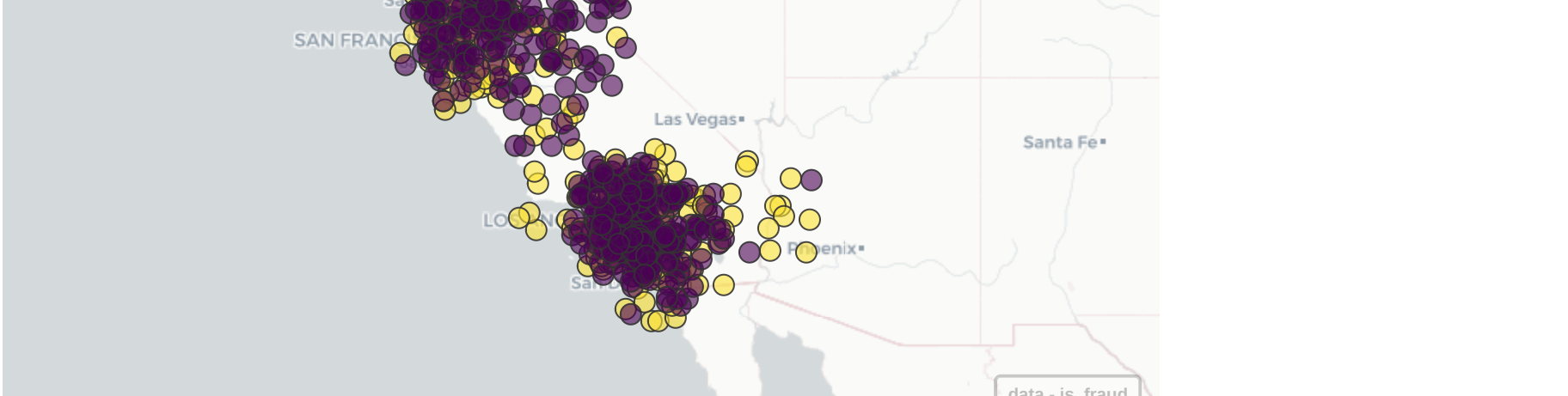
<pre>fraud <- fraud %>% mutate_if(is.character, as.factor) %>% mutate(year_p = str_split(dob, "-"), simplify = TRUE)[, 1], year_d = str_split(trans_date_trans_time, "-"), simplify = TRUE)[, 1], year_g = str_split(year_p, "-"), simplify = TRUE)[, 1], age = as.integer(year_p) + as.integer(year_d), hour = hour(trans_date_trans_time), day = day(trans_date_trans_time)) %>% select(c(1:dob, year_p, year_d, trans_time))</pre>
<pre>## Warning: There were 2 warnings in `mutate()`. ## The first warning was: ## 1 In argument: `hour = hour(trans_date_trans_time)`.</pre>
<pre>## Caused by warning: ## 1 tz(): don't know how to compute timezone for object of class factor; returning "UTC". ## 1 Run `dplyr::last_dplyr_warnings()` to see the 1 remaining warning.</pre>
<pre>fraud\$is_fraud <- as.factor(fraud\$is_fraud)</pre>

Some coordinates tend to have more fraud occurrence



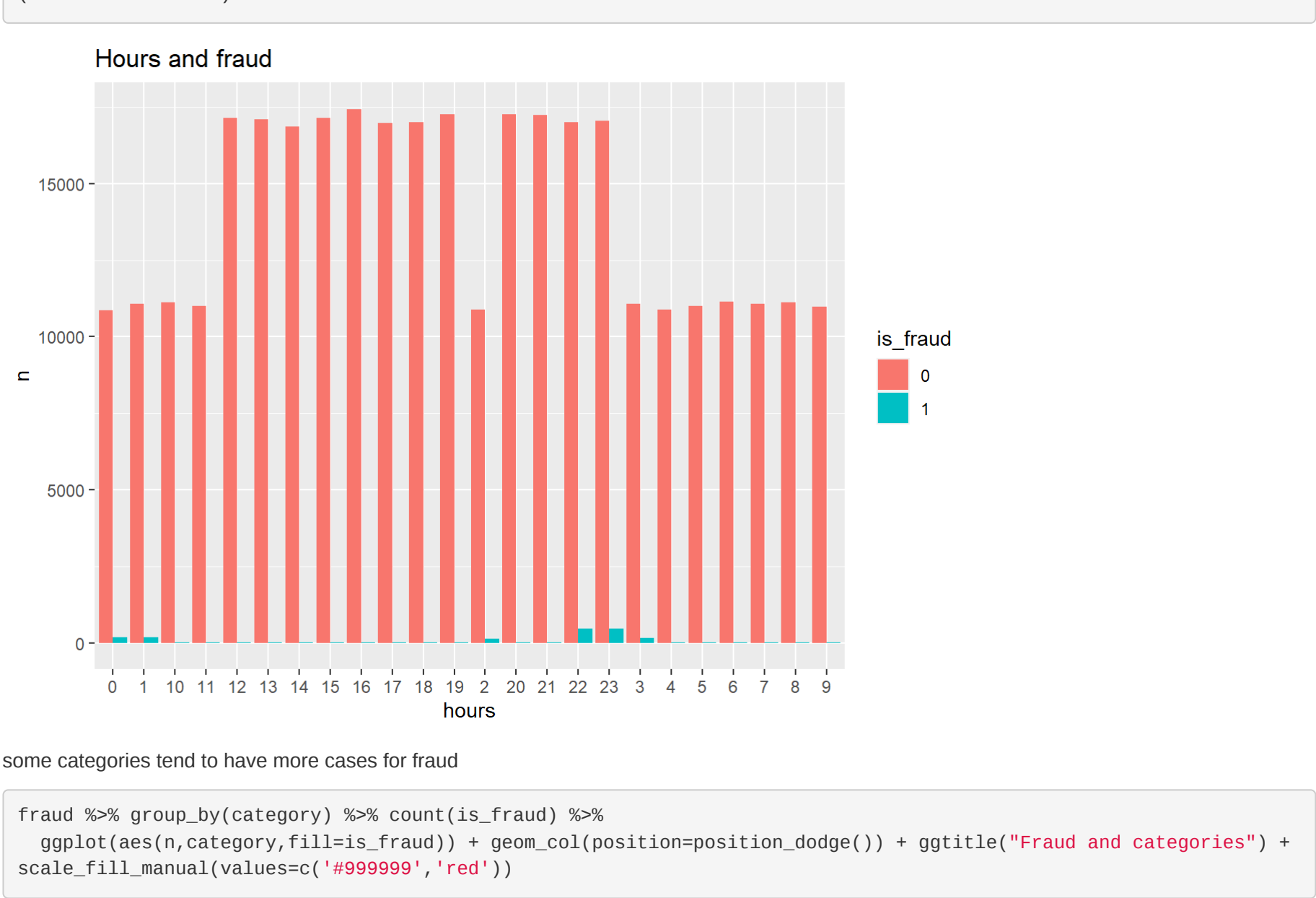
Highest state in terms of fraud occurrence

<pre>fraud %>% group_by(state) %>% count(is_fraud) %>% filter(is_fraud == 1) %>% arrange(-n) %>% head(3)</pre>
<pre>## # A tibble: 3 x 3 ## # Groups: state [3] ## state is_fraud n ## <fct> <fct> <int> ## 1 CA 1 482 ## 2 MD 1 262 ## 3 ME 1 216</pre>
<pre>CA <- rbind(fraud %>% filter(state == "CA" & is_fraud == 1), fraud %>% filter(state == "CA" & is_fraud == 0)) %>% slice(1:1500)</pre>



Some hours have more cases of frauds more than others

<pre>fraud %>% group_by(hour) %>% count(is_fraud) %>% ggplot(aes(as.character(hour), n, fill=is_fraud)) + geom_col(position=position_dodge()) + xlab("hours") + ggtitle("Hours and Fraud")</pre>

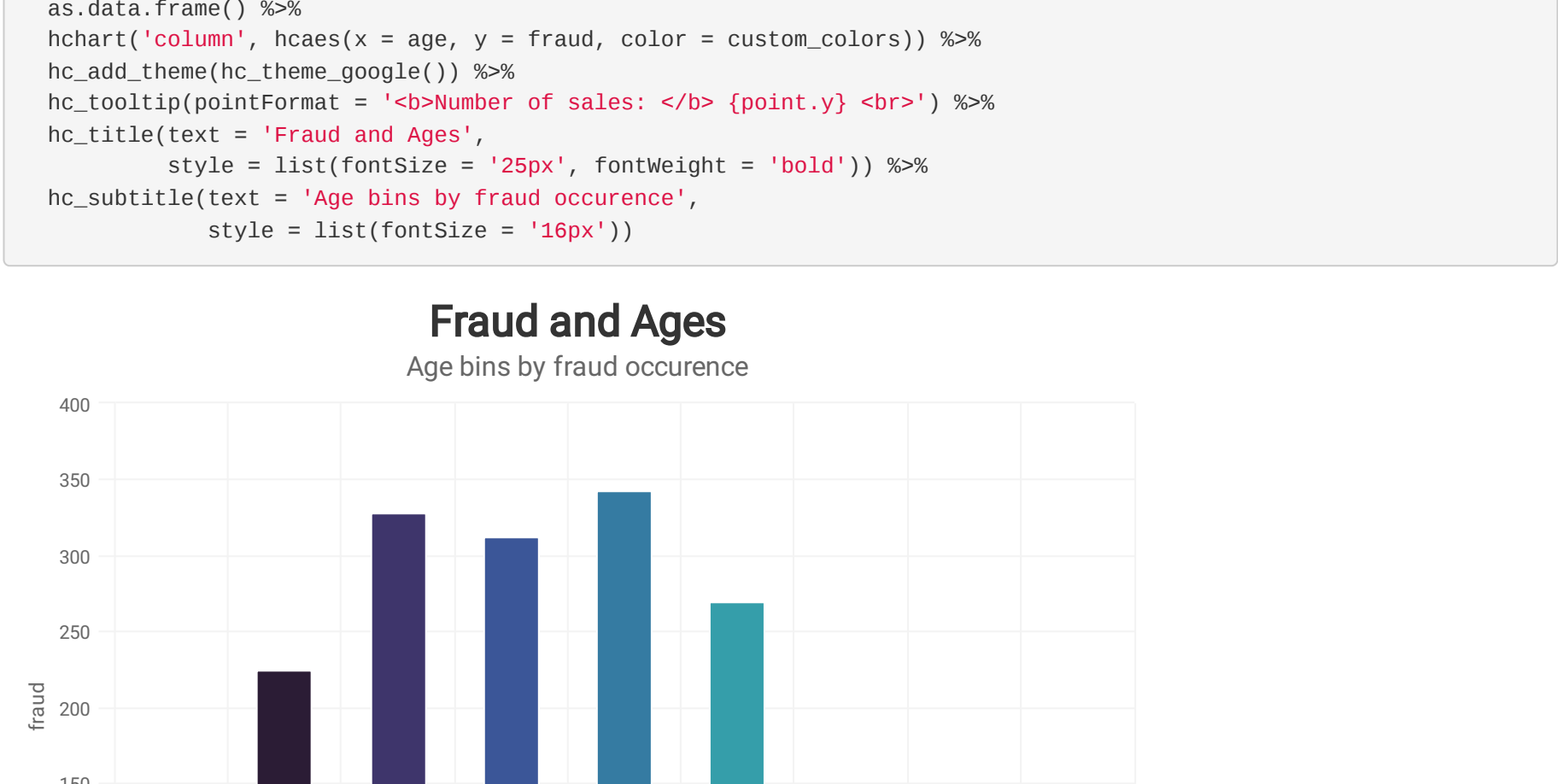


some categories tend to have more cases for fraud

<pre>fraud %>% group_by(category) %>% count(is_fraud) %>% ggplot(aes(n, category, fill=is_fraud)) + geom_col(position=position_dodge()) + ggtitle("Fraud and categories") + scale_fill_manual(values=c("#999999", "red"))</pre>
<pre>## # A tibble: 13 x 3 ## # Groups: category [13] ## category is_fraud n ## <fct> <fct> <int> ## 1 travel 0 10000 ## 2 shopping_pos 0 10000 ## 3 shopping_net 0 10000 ## 4 personal_care 0 10000 ## 5 misc_pos 0 10000 ## 6 misc_net 0 10000 ## 7 kids_pos 0 10000 ## 8 home 0 10000 ## 9 health_fitness 0 10000 ## 10 grocery_pos 0 10000 ## 11 grocery_net 0 10000 ## 12 gas_transport 0 10000 ## 13 food_drink 0 10000 ## 14 entertainment 0 10000</pre>

some ages tend to have more cases for fraud

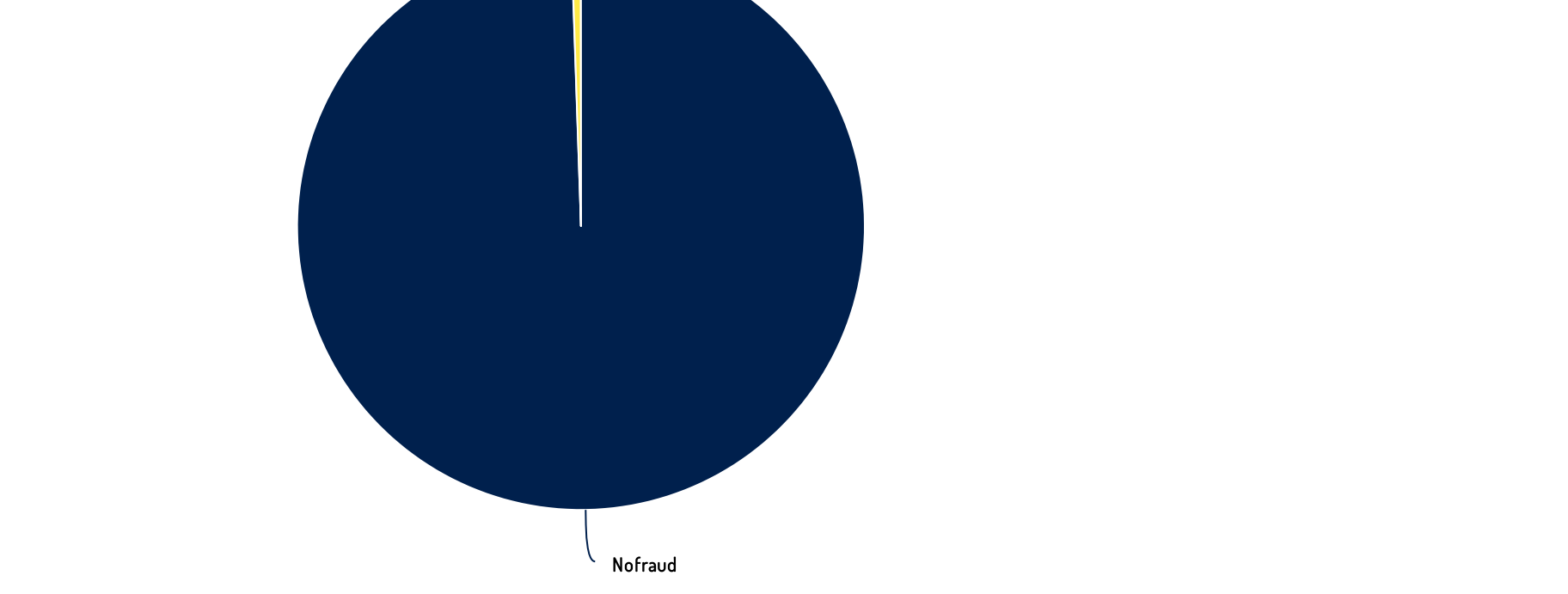
<pre>dfraud %>% group_by(age) %>% count(is_fraud) %>% filter(is_fraud == 1) %>% mutate(age = as.character(cut(age, seq(10, 100, 10)))) %>% group_by(age) %>% summarize(sums = sum(n)) %>% as.data.frame() %>% ggplot(aes(age, sums)) + ggtitle("Ages and Fraud") + geom_col(color="red", fill="white") + ylab("Number of frauds")</pre>
<pre>custom_colors <- viridis::mako(n = 9) fraud %>% group_by(age) %>% count(is_fraud) %>% filter(is_fraud == 1) %>% mutate(age = as.character(cut(age, seq(10, 100, 10)))) %>% group_by(age) %>% summarize(fraud = sum(n)) %>% as.data.frame() %>% hchart("column", hcats(x = age, y = fraud, color = custom_colors)) %>% hc_add_theme(hc_theme_google()) %>% hc_tooltip(pointFormat = "Number of sales: {point.y}") %>% hc_title(text = "Fraud and Ages", style = list(fontSize = "25px", fontWeight = "bold")) %>% hc_subtitle(text = "Age bins by Fraud occurrence", style = list(fontSize = "15px"))</pre>



Unbalanced data

The data is highly unbalanced 99% non fraud

<pre>fraud %>% count(is_fraud) %>% mutate(n=paste(round(n*100/nrow(fraud), 2), "%"))</pre>
<pre>## is_fraud n ## 1 0 99.48 % ## 2 1 0.52 %</pre>
<pre>custom_colors <- viridis::cividis(n = 2) fraud %>% count(is_fraud) %>% mutate(is_fraud = c("no_fraud", "fraud")) %>% hchart("pie", hcats(x = is_fraud, y = n, color = custom_colors)) %>% hc_add_theme(hc_theme_gridlight()) %>% hc_tooltip(pointFormat = "Number of sales: {point.y}") %>% hc_title(text = "Percentage of classes", style = list(fontSize = "25px", fontWeight = "bold"))</pre>



Converting the GPS coordinates to distance (new column)

<pre>df <- data.frame(fraud %>% select(long_lat, merch_long, merch_lat)) distances_km <- numeric(nrow(df)) for (i in 1:nrow(df)) { lon1 <- df\$long[i] lat1 <- df\$lat[i] lon2 <- df\$merch_long[i] lat2 <- df\$merch_lat[i] distances_km[i] <- distGeo(c(lon1, lat1), c(lon2, lat2)) / 1000 }</pre>
<pre>fraud\$distances_km <- distances_km rm(distances_km, lon1, lat1, lon2, lat2, i, df)</pre>

Categorical columns

name of merchant has so many unique values and we can't one hot code it (693 new columns), label encoding also isn't a choice .

so i preferred dropping it and it also don't tend to have that much importance .

same for the city .

category could be one hot encoded and has influence on the fraud as represented above.

job might have some importance on the fraud so i tried to frequency-encode it (replace the values by the frequency) .

<pre>uniques <- fraud %>% summarize(merchant = length(unique(merchant)), category = length(unique(category)), job = length(unique(job)), city = length(unique(city))) %>% t() %>% as.data.frame()</pre>
<pre>## V1 ## merchant 693 ## category 14 ## job 163 ## city 178</pre>

city and state columns dropped because they have the same info and they are hard to encode so i chose to use city pop instead

<pre>fraud <- fraud %>% select(-c(long, merch_long, merch_lat, lat, merchant, city, state))</pre>

Splitting

<pre>set.seed(123) data_split <- initial_split(fraud, strata = is_fraud) train <- training(data_split) test <- testing(data_split) #cross validation object fraud_folds <- vfold_cv(train, v = 3, strata = is_fraud)</pre>
--

Recipe

i tried 3 different recipes to tune the model and choose both the best recipe and best hyperparameters

<pre>#basic recipe recipe_plain <- recipe(is_fraud ~., data = train) %>% step_normalize(all_numeric_predictors()) %>% step_mutate(job, count = n()) %>% step_integer(job %>% step_n(count) %>% step_dummy(all_nominal_predictors())) #rebalancing using smote smote <- recipe_plain %>% step_smote(is_fraud, over_ratio = 0.85) %>% step_sample(size = nrow(train)) #rebalancing using random undersampling rus <- recipe_plain %>% step_downsample(is_fraud)</pre>
--

Metric set

<pre>metric <- metric_set(sens, precision, yardstick::spec, j_index, f_meas)</pre>

Model

Spec

i chose the lightgbm model because it was the best model

<pre>set.seed(123) lightgbm_spec <- boost_tree(mtry = tune(), trees = tune(), tree_depth = tune(), learn_rate = tune(), min_n = tune(), subs_reduction = tune()) %>% set_engine(engine = "lightgbm") %>% set_mode(mode = "classification")</pre>

Workflow

<pre>wf_set_tune <- workflow_set(list(plain = recipe_plain, smote = smote, rus = rus), list(lightgbm = lightgbm_spec))</pre>

Tune

tune the hyperparameters of the model and evaluate the model across different recipes (simple, smote, rus)

<pre>set.seed(123) tune_results <- workflow_set(wf_set_tune, "tune_grid", examples = fraud_folds, grid = 0, metric = metric, verbose = TRUE)</pre>
<pre>## # 1 of 3 tuning: plain_lightgbm</pre>
<pre>## 1 Creating pre-processing data to finalize unknown parameter: mtry</pre>
<pre>## ✓ 1 of 3 tuning: plain_lightgbm (2n 29.9s)</pre>
<pre>## # 2 of 3 tuning: smote_lightgbm</pre>
<pre>## 1 Creating pre-processing data to finalize unknown parameter: mtry</pre>
<pre>## ✓ 2 of 3 tuning: smote_lightgbm (3n 44.7s)</pre>
<pre>## # 3 of 3 tuning: rus_lightgbm</pre>
<pre>## 1 Creating pre-processing data to finalize unknown parameter: mtry</pre>
<pre>## ✓ 3 of 3 tuning: rus_lightgbm (1n 28.8s)</pre>

Ranking the tuning results by j-index on validation sets

the balanced data ranks better

both rebalancing methods increased the j_index

smote 0.82

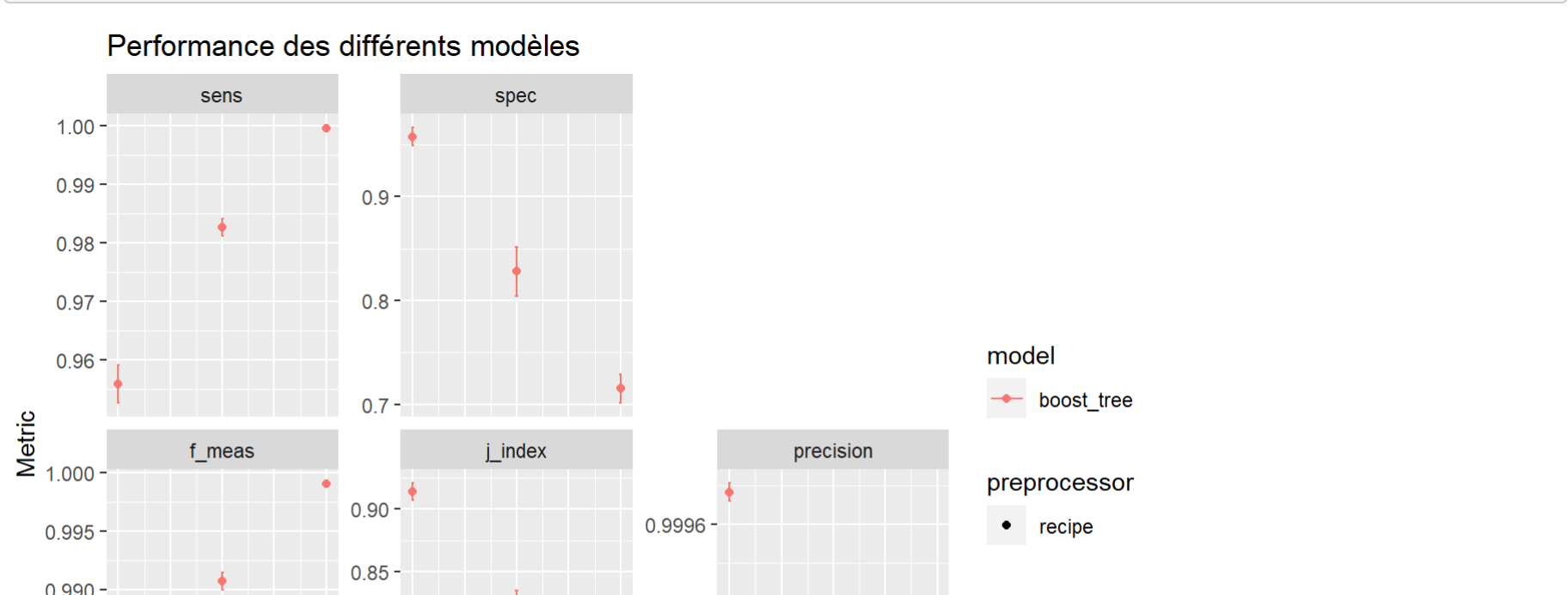
rus 0.91

the under sampling did a better job

<pre>rank_results(tune_results, rank_metric = "j_index")</pre>
<pre>## # A tibble: 98 x 9 ## workflow_id config metric mean std err n preprocessor model rank ## <chr> <chr> <chr> <dbl> <dbl> <int> <chr> <chr> <int> ## 1 rus_lightgbm Preprocess.f_meas 0.977 1.03e-3 3 recipe boost_ 1 ## 2 rus_lightgbm Preprocess.j_index 0.914 1.22e-3 3 recipe boost_ 1 ## 3 rus_lightgbm Preprocess.precis 1.08 2.83e-5 3 recipe boost_ 1 ## 4 rus_lightgbm Preprocess.sens 0.956 1.97e-3 3 recipe boost_ 1 ## 5 rus_lightgbm Preprocess.spec 0.958 1.26e-3 3 recipe boost_ 1 ## 6 rus_lightgbm Preprocess.f_meas 0.966 1.05e-3 3 recipe boost_ 2 ## 7 rus_lightgbm Preprocess.j_index 0.874 1.25e-2 3 recipe boost_ 2 ## 8 rus_lightgbm Preprocess.precis 1.08 5.61e-5 3 recipe boost_ 2 ## 9 rus_lightgbm Preprocess.sens 0.935 1.92e-3 3 recipe boost_ 2 ## 10 rus_lightgbm Preprocess.spec 0.939 1.06e-2 3 recipe boost_ 2 ## # 180 more rows</pre>

Results, down best model

<pre>results_down_gmb <- tune_results %>% extract_workflow_set_result("rus_lightgbm") autoplot(tune_results, rank_metric = "j_index", select_best = TRUE) + ggtitle("Performance des différents modèles")</pre>
--



<pre>autoplot(results_down_gmb, metric = c("accuracy", "j_index")) + ggtitle("Performance des différents hyperparamètres de LightGBM")</pre>
--



Finalizing workflow

<pre>best_hyperparameters <- tune_results %>% extract_workflow_set_result("rus_lightgbm") %>% select_best(metric = "j_index") validation_results <- tune_results %>% extract_workflow("rus_lightgbm") %>% finalize_workflow(best_hyperparameters) %>% last_fit(data_split, metrics = metric)</pre>
--

Performance on test data

as we can see we get we got stable metrics and stability between test and validation :

high accuracy 0.86 it was 0.998 which could lead to over fit

higher j-index for both test(0.919) and validation (0.913)

<pre>rbind(validation_results %>% collect_metrics()) %>% select(-.config), validation_results %>% collect_predictions() %>% accuracy(truth = is_fraud, estimate = .pred_class))</pre>
<pre>## # A tibble: 6 x 3 ## metric estimator estimate ## <chr> <chr> <dbl> ## 1 sens binary 0.982 ## 2 precision binary 1.08 ## 3 spec binary 0.958 ## 4 j_index binary 0.929 ## 5 f_meas binary 0.981 ## 6 accuracy binary 0.982</pre>

Confusion matrix

<pre>validation_results %>% collect_predictions() %>% conf_mat(truth = is_fraud, estimate = .pred_class)</pre>
<pre>## Truth ## Prediction 0 1 ## 0 81261 18 ## 1 3214 469</pre>

The matrix indicates that the models focuses more on catching fraudulent transactions rather than getting non fraud as non fraud and that might be useful as one fraud predicted as fraud is more important than predicting non frauds as fraud ... one single fraudulent transaction could cause much more loss than losing a non fraudulent customer