

XGBoost

Contents

1	Upload-Rate Prediction	1
1.1	Reading the Data	1
1.2	Create the Prediction Tasks for Each Provider	3
1.3	Create Data Splitting Strategies for Testing and Validation	4
1.4	Create the Prediction Pipeline	5
1.5	Parameter Tuning	6
1.6	Create Learners with Tuned Hyperparameters	7
1.7	Validation Results	7
1.8	Feature Importance	11
2	Download-Rate Prediction	14
2.1	Reading the Data	14
2.2	Create the Prediction Tasks for Each Provider	16
2.3	Parameter Tuning	16
2.4	Create Learners with Tuned Hyperparameters	17
2.5	Validation Results	17
2.6	Feature Importance	21

```
library(tidyverse)
```

```
library(ggplot2)
```

```
library(mlr3)
```

```
library(mlr3learners)
```

```
library(mlr3pipelines)
```

```
library(mlr3tuning)
```

```
library(mlr3filters)
```

```
library(paradox)
```

```
future::plan("multiprocess")
```

```
## Warning: Strategy 'multiprocess' is deprecated in future (>= 1.20.0). Instead,  
## explicitly specify either 'multisession' or 'multicore'. In the current R  
## session, 'multiprocess' equals 'multisession'.
```

1 Upload-Rate Prediction

1.1 Reading the Data

```
data_dir = "../datasets/"
```

```
results_dir = "../prediction_results/"
```

```
dataset_ul = read_csv(  
  str_c(data_dir, "dataset_ul.csv"),  
  col_types = cols(  
    drive_id = col_integer(),
```



```
## $ rsrq_db      <dbl> -5, -6, -5, -6, -6, -10, -8, -11, -11, -10, -9, -1...
## $ rssnr_db     <dbl> 22, 11, 29, 13, 16, 13, 7, 0, 8, 2, 24, 10, 22, 15...
## $ cqi          <dbl> 10, 13, 15, 12, 9, 15, 10, 9, 9, 7, 10, 9, 12, 15,...
## $ ta          <dbl> 7, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 3, 3,...
## $ enodeb       <fct> 103068, 114809, 114809, 114809, 114809, 114809, 11...
## $ f_mhz        <dbl> 1720, 1720, 1720, 1720, 1720, 1720, 1720, 1720, 17...
## $ payload_mb   <dbl> 4.0, 2.0, 4.0, 9.0, 8.0, 6.0, 5.0, 4.0, 3.0, 2.0, ...
## $ throughput_mbits <dbl> 24.52, 14.86, 16.27, 12.68, 14.59, 13.13, 16.37, 1...
```

```
dataset_ul_vodafone = filter(dataset_ul, provider=="vodafone")
glimpse(dataset_ul_vodafone)
```

```
## Rows: 1,828
## Columns: 15
## $ row_id_original <int> 4341, 4342, 4343, 4344, 4345, 4346, 4347, 4348, 43...
## $ drive_id       <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,...
## $ timestamp      <dtm> 2018-12-10 09:09:03, 2018-12-10 09:09:21, 2018-12...
## $ scenario       <fct> campus, campus, campus, campus, campus, campus, ca...
## $ provider       <fct> vodafone, vodafone, vodafone, vodafone, vodafone, ...
## $ velocity_mps   <dbl> 11.70, 8.22, 8.00, 10.30, 12.28, 0.00, 0.00, 0.00, ...
## $ rsrp_dbm       <dbl> -121, -108, -111, -106, -110, -94, -95, -92, -98, ...
## $ rsrq_db        <dbl> -15, -9, -13, -8, -9, -7, -7, -8, -6, -10, -7, -8, ...
## $ rssnr_db       <dbl> -8, 2, 6, 5, 9, 23, 23, 24, 14, 1, 14, 12, 14, 7, ...
## $ cqi            <dbl> 4, 2, 6, 11, 10, 15, 12, 15, 12, 6, 15, 10, 11, 7, ...
## $ ta             <dbl> 63, 21, 16, 16, 16, 16, 16, 16, 16, 16, 16, 16, 16...
## $ enodeb         <fct> 51044, 52316, 50026, 50026, 50026, 50026, 50026, 5...
## $ f_mhz          <dbl> 1770, 1770, 1770, 1770, 1770, 1770, 1770, 1770, 17...
## $ payload_mb     <dbl> 6.0, 10.0, 0.1, 2.0, 6.0, 0.1, 0.1, 0.5, 7.0, 0.1, ...
## $ throughput_mbits <dbl> 1.29, 3.18, 0.05, 2.93, 8.79, 5.16, 4.73, 10.13, 1...
```

1.2 Create the Prediction Tasks for Each Provider

```
make_task = function(dataset, task_id) {
  task = TaskRegr$new(
    id = task_id,
    backend = dataset %>% select(-drive_id, -timestamp, -provider, -scenario),
    target = "throughput_mbits"
  )

  task$col_roles$name = "row_id_original"
  task$col_roles$feature = setdiff(task$col_roles$feature, "row_id_original")

  return(task)
}
```

```
task_ul_o2 = make_task(dataset_ul_o2, "task_ul_o2")
task_ul_o2
```

```
## <TaskRegr:task_ul_o2> (2039 x 10)
## * Target: throughput_mbits
## * Properties: -
## * Features (9):
##   - dbl (8): cqi, f_mhz, payload_mb, rsrp_dbm, rsrq_db, rssnr_db, ta,
##     velocity_mps
##   - fct (1): enodeb
```

```

task_ul_tmobile = make_task(dataset_ul_tmobile, "task_ul_tmobile")
task_ul_tmobile

## <TaskRegr:task_ul_tmobile> (2301 x 10)
## * Target: throughput_mbits
## * Properties: -
## * Features (9):
##   - dbl (8): cqi, f_mhz, payload_mb, rsrp_dbm, rsrq_db, rssnr_db, ta,
##     velocity_mps
##   - fct (1): enodeb

task_ul_vodafone = make_task(dataset_ul_vodafone, "task_ul_vodafone")
task_ul_vodafone

## <TaskRegr:task_ul_vodafone> (1828 x 10)
## * Target: throughput_mbits
## * Properties: -
## * Features (9):
##   - dbl (8): cqi, f_mhz, payload_mb, rsrp_dbm, rsrq_db, rssnr_db, ta,
##     velocity_mps
##   - fct (1): enodeb

```

1.3 Create Data Splitting Strategies for Testing and Validation

The outer resampling is used for the train/validation split.

```

get_row_ids_by_drive_ids = function(task, dataset, drive_ids) {
  result = (tibble(task$row_names) %>%
    inner_join(dataset, by=c("row_name"="row_id_original")) %>%
    filter(drive_id %in% drive_ids))$row_id
  return(result)
}

make_outer_resampling = function(task, dataset, drive_ids_train, drive_ids_test) {
  row_ids_train = get_row_ids_by_drive_ids(task, dataset, drive_ids_train)
  row_ids_test = get_row_ids_by_drive_ids(task, dataset, drive_ids_test)

  result = rsmp("custom")
  result$instantiate(task, train_sets=list(row_ids_train), test_sets=list(row_ids_test))

  return(result)
}

```

The inner resampling is used for the parameter tuning on the training set.

```

make_inner_resampling = function(task, dataset, last_drive_id) {
  train_sets = list()
  test_sets = list()

  for (cur_last_drive_id_train in 2:(last_drive_id-1)) {
    drive_ids_train = 1:cur_last_drive_id_train
    drive_ids_test = cur_last_drive_id_train + 1

    row_ids_train = get_row_ids_by_drive_ids(task, dataset, drive_ids_train)
    row_ids_test = get_row_ids_by_drive_ids(task, dataset, drive_ids_test)
  }
}

```

```

train_sets[[length(train_sets)+1]] = row_ids_train
test_sets[[length(test_sets)+1]] = row_ids_test
}

result = rsmp("custom")
result$instantiate(task, train_sets=train_sets, test_sets=test_sets)

return(result)
}

```

1.4 Create the Prediction Pipeline

```

make_learner = function(nrounds=100, eta=NULL, gamma=NULL, lambda=NULL) {
  factor_encoding = po(
    "encode",
    method = "one-hot",
    affect_columns = selector_type("factor")
  )
  xgboost = lrn("regr.xgboost")

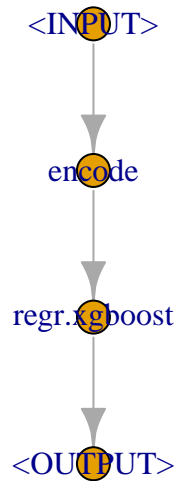
  if (!is.null(nrounds)) {
    xgboost$param_set$values = mlr3misc::insert_named(
      xgboost$param_set$values,
      list(nrounds=nrounds)
    )
  }
  if (!is.null(eta)) {
    xgboost$param_set$values = mlr3misc::insert_named(
      xgboost$param_set$values,
      list(eta=eta)
    )
  }
  if (!is.null(gamma)) {
    xgboost$param_set$values = mlr3misc::insert_named(
      xgboost$param_set$values,
      list(gamma=gamma)
    )
  }
  if (!is.null(lambda)) {
    xgboost$param_set$values = mlr3misc::insert_named(
      xgboost$param_set$values,
      list(lambda=lambda)
    )
  }

  pipe = factor_encoding %>% PipeOpLearner$new(xgboost)
  learner = GraphLearner$new(pipe)
  return(learner)
}

```

Here we can see the prediction pipeline:

```
make_learner()$graph$plot()
```



1.5 Parameter Tuning

```
parameter_space = ParamSet$new(list(  
  ParamInt$new("regr.xgboost.nrounds", lower=100, upper=1000),  
  ParamDbl$new("regr.xgboost.eta", lower=0.01, upper=1),  
  ParamDbl$new("regr.xgboost.gamma", lower=0, upper=10),  
  ParamDbl$new("regr.xgboost.lambda", lower=0, upper=10)  
))  
  
get_tuning_result = function(task, dataset, grid_resolution, n_evals) {  
  tuning_instance = TuningInstanceSingleCrit$new(  
    task = task,  
    learner = make_learner(),  
    resampling = make_inner_resampling(task, dataset, last_drive_id=7),  
    measure = msr("regr.mae"),  
    terminator = trm("evals", n_evals=n_evals),  
    search_space = parameter_space$clone(deep = TRUE),  
    store_benchmark_result = TRUE,  
    check_values = TRUE  
  )  
  
  tuner = tnr("grid_search", resolution = grid_resolution)  
  tuner$optimize(tuning_instance)  
  
  return(tuning_instance)  
}
```

```

tuning_result_ul_o2 = get_tuning_result(task_ul_o2, dataset_ul, grid_resolution = 20, n_evals = 10)

tuning_result_ul_tmobile = get_tuning_result(task_ul_tmobile, dataset_ul, grid_resolution = 20, n_evals = 10)

tuning_result_ul_vodafone = get_tuning_result(task_ul_vodafone, dataset_ul, grid_resolution = 20, n_evals = 10)

tuning_result_ul = bind_rows(
  tibble(tuning_result_ul_o2$result) %>% mutate(provider="o2"),
  tibble(tuning_result_ul_tmobile$result) %>% mutate(provider="tmobile"),
  tibble(tuning_result_ul_vodafone$result) %>% mutate(provider="vodafone"),
) %>% select("provider", "regr.xgboost.nrounds", "regr.xgboost.eta", "regr.xgboost.gamma", "regr.xgboost.lambda")

knitr::kable(tuning_result_ul)

```

provider	regr.xgboost.nrounds	regr.xgboost.eta	regr.xgboost.gamma	regr.xgboost.lambda
o2	621	0.1142105	3.157895	3.684210
tmobile	242	0.3747368	6.842105	6.842105
vodafone	858	0.1142105	6.842105	8.947368

1.6 Create Learners with Tuned Hyperparameters

```

learner_ul_o2 = make_learner(
  nrounds = tuning_result_ul_o2$result$regr.xgboost.nrounds,
  eta = tuning_result_ul_o2$result$regr.xgboost.eta,
  gamma = tuning_result_ul_o2$result$regr.xgboost.gamma,
  lambda = tuning_result_ul_o2$result$regr.xgboost.lambda
)

learner_ul_tmobile = make_learner(
  nrounds = tuning_result_ul_tmobile$result$regr.xgboost.nrounds,
  eta = tuning_result_ul_tmobile$result$regr.xgboost.eta,
  gamma = tuning_result_ul_tmobile$result$regr.xgboost.gamma,
  lambda = tuning_result_ul_tmobile$result$regr.xgboost.lambda
)

learner_ul_vodafone = make_learner(
  nrounds = tuning_result_ul_vodafone$result$regr.xgboost.nrounds,
  eta = tuning_result_ul_vodafone$result$regr.xgboost.eta,
  gamma = tuning_result_ul_vodafone$result$regr.xgboost.gamma,
  lambda = tuning_result_ul_vodafone$result$regr.xgboost.lambda
)

```

1.7 Validation Results

```

resampling_result_ul_o2 = resample(
  task = task_ul_o2,
  learner = learner_ul_o2,
  resampling = make_outer_resampling(task_ul_o2, dataset_ul, drive_ids_train=1:7, drive_ids_test=8:10),
  store_models = TRUE
)

```

```
resampling_result_ul_tmobile = resample(
  task = task_ul_tmobile,
  learner = learner_ul_tmobile,
  resampling = make_outer_resampling(task_ul_tmobile, dataset_ul, drive_ids_train=1:7, drive_ids_test=8),
  store_models = TRUE
)

resampling_result_ul_vodafone = resample(
  task = task_ul_vodafone,
  learner = learner_ul_vodafone,
  resampling = make_outer_resampling(task_ul_vodafone, dataset_ul, drive_ids_train=1:7, drive_ids_test=8),
  store_models = TRUE
)

predictions_ul_o2 = as.data.table(resampling_result_ul_o2$prediction())
glimpse(tibble(predictions_ul_o2))

## Rows: 615
## Columns: 3
## $ row_id    <int> 137, 138, 139, 140, 141, 142, 143, 144, 145, 146, 147, 148...
## $ truth     <dbl> 4.47, 2.59, 2.26, 1.09, 0.77, 0.19, 0.26, 0.65, 1.45, 1.12...
## $ response  <dbl> 3.1943078, 2.4845176, 3.0184193, 2.8502736, 2.4734516, 1.3...

predictions_ul_tmobile = as.data.table(resampling_result_ul_tmobile$prediction())
predictions_ul_vodafone = as.data.table(resampling_result_ul_vodafone$prediction())

validation_results_ul = bind_rows(
  tibble(predictions_ul_o2) %>%
    inner_join(tibble(task_ul_o2$row_names), by="row_id") %>%
    inner_join(dataset_ul, by=c("row_name"="row_id_original")),
  tibble(predictions_ul_tmobile) %>%
    inner_join(tibble(task_ul_tmobile$row_names), by="row_id") %>%
    inner_join(dataset_ul, by=c("row_name"="row_id_original")),
  tibble(predictions_ul_vodafone) %>%
    inner_join(tibble(task_ul_vodafone$row_names), by="row_id") %>%
    inner_join(dataset_ul, by=c("row_name"="row_id_original"))
)
glimpse(validation_results_ul)

## Rows: 1,840
## Columns: 18
## $ row_id    <int> 137, 138, 139, 140, 141, 142, 143, 144, 145, 146, ...
## $ truth     <dbl> 4.47, 2.59, 2.26, 1.09, 0.77, 0.19, 0.26, 0.65, 1....
## $ response  <dbl> 3.1943078, 2.4845176, 3.0184193, 2.8502736, 2.4734...
## $ row_name  <int> 137, 138, 139, 140, 141, 142, 143, 144, 145, 146, ...
## $ drive_id  <int> 8, 8, 8, 8, 8, 8, 8, 8, 8, 8, 8, 8, 8, 8, 8, 9, 9,...
## $ timestamp <dtm> 2018-12-11 09:04:11, 2018-12-11 09:04:22, 2018-12...
## $ scenario  <fct> campus, campus, campus, campus, campus, campus, ca...
## $ provider  <fct> o2, o2, o2, o2, o2, o2, o2, o2, o2, o2, o2, o2, o2...
## $ velocity_mps <dbl> 0.00, 6.11, 9.39, 8.45, 11.68, 0.00, 0.00, 0.00, 4...
## $ rsrp_dbm  <dbl> -89, -92, -94, -98, -102, -100, -101, -101, -100, ...
## $ rsrq_db   <dbl> -9, -12, -14, -15, -16, -17, -16, -16, -17, -14, -...
## $ rssnr_db  <dbl> 13, 3, -1, -3, -5, -7, -6, -5, -8, 1, -7, -1, -2, ...
## $ cqi       <dbl> 11, 5, 5, 4, 2, 3, 4, 4, 4, 6, 3, 5, 5, 2, 4, 6, 3...
## $ ta       <dbl> 7, 7, 7, 7, 7, 12, 12, 12, 12, 12, 12, 12, 12, 12,...
```



```
## $ enodeb      <fct> 52410, 52410, 52410, 52410, 52410, 52900, 52900, 5...
## $ f_mhz       <dbl> 880, 880, 880, 880, 880, 880, 880, 880, 880, 880, ...
## $ payload_mb  <dbl> 0.1, 0.5, 3.0, 9.0, 7.0, 3.0, 2.0, 2.0, 6.0, 3.0, ...
## $ throughput_mbits <dbl> 4.47, 2.59, 2.26, 1.09, 0.77, 0.19, 0.26, 0.65, 1....
```

```
all(validation_results_ul$truth == validation_results_ul$throughput_mbits)
```

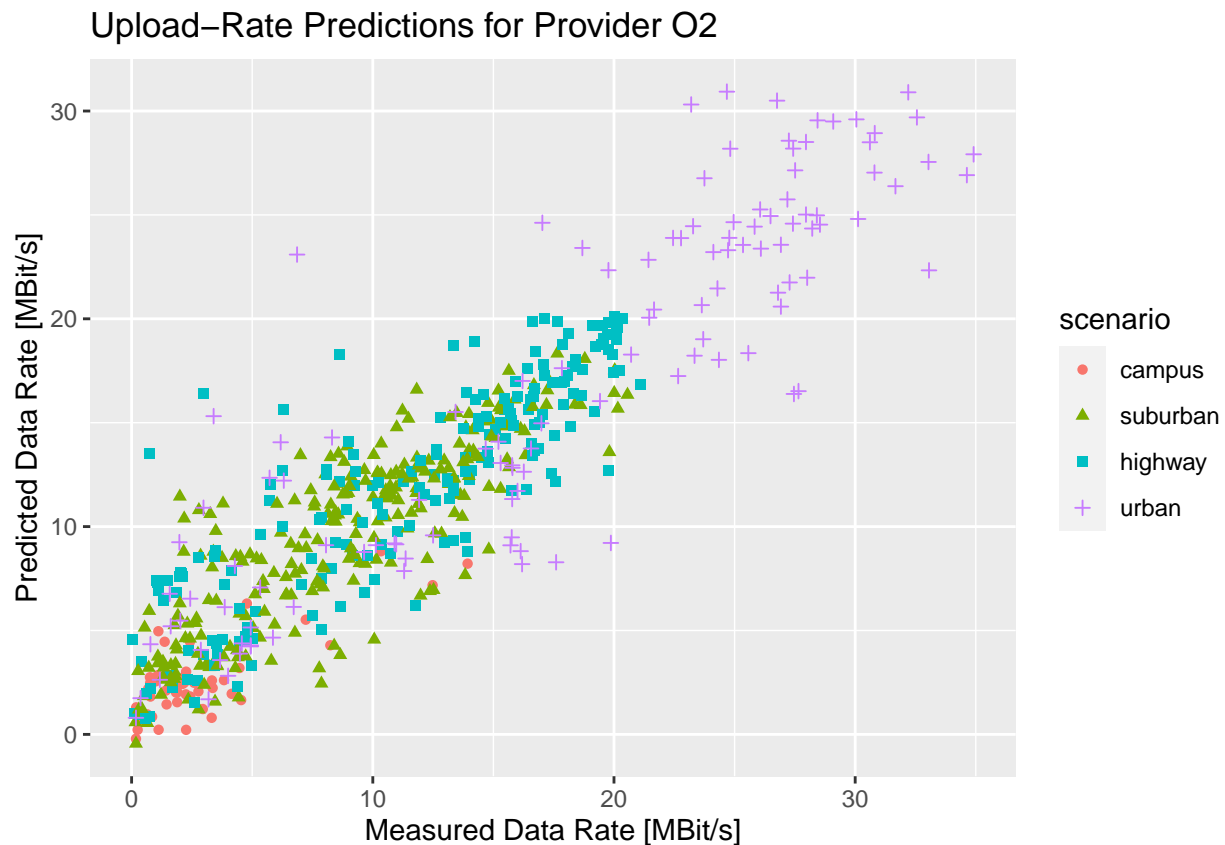
```
## [1] TRUE
```

```
validation_results_ul = validation_results_ul %>%
  rename(prediction_xgboost=response) %>%
  select(-truth, -row_id, -row_name)

write_csv(validation_results_ul, str_c(results_dir, "predictions_xgboost_ul.csv"))
```

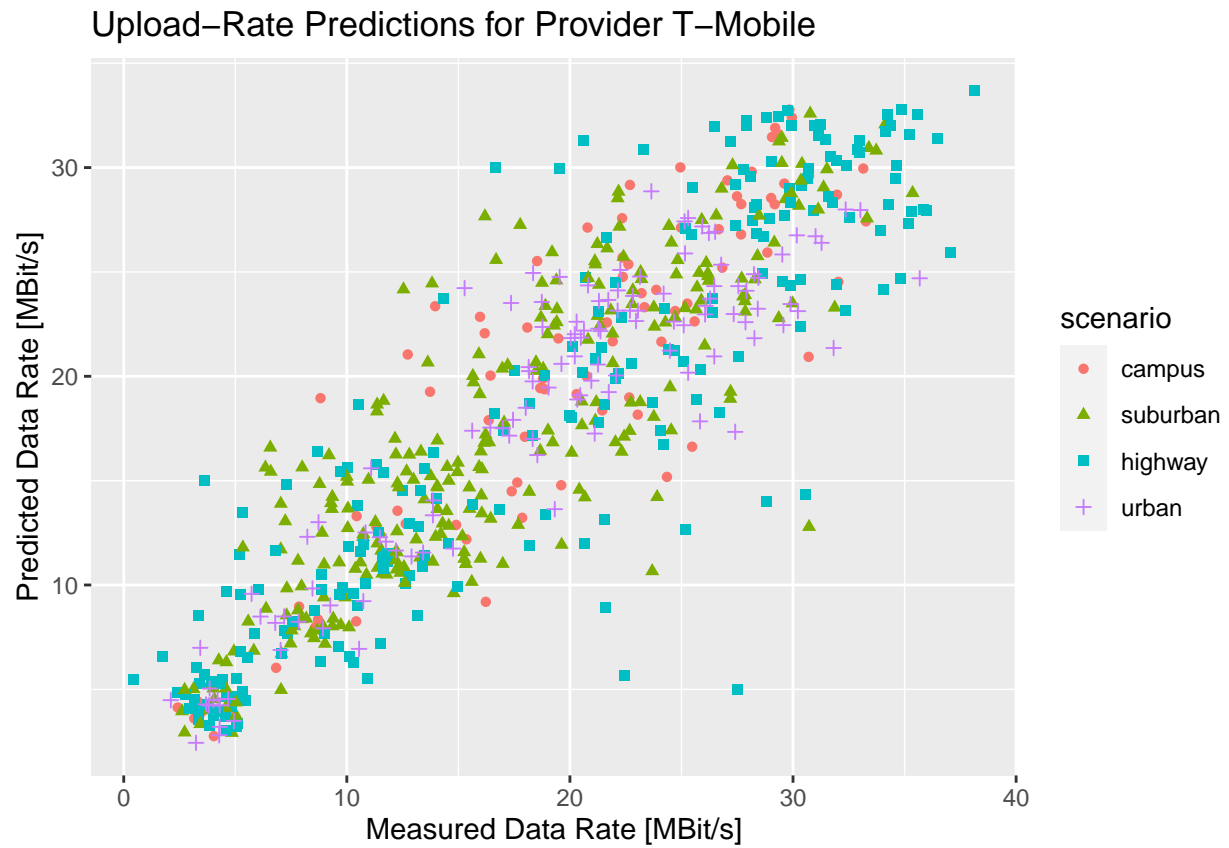
1.7.1 Scatter Plots

```
ggplot(filter(validation_results_ul, provider=="o2"), aes(x=throughput_mbits, y=prediction_xgboost)) +
  geom_point(aes(color=scenario, shape=scenario)) +
  xlab("Measured Data Rate [MBit/s]") +
  ylab("Predicted Data Rate [MBit/s]") +
  ggtitle("Upload-Rate Predictions for Provider O2")
```



```
ggplot(filter(validation_results_ul, provider=="tmobile"), aes(x=throughput_mbits, y=prediction_xgboost)) +
  geom_point(aes(color=scenario, shape=scenario)) +
  xlab("Measured Data Rate [MBit/s]") +
  ylab("Predicted Data Rate [MBit/s]") +
```

```
ggtitle("Upload-Rate Predictions for Provider T-Mobile")
```



```
ggplot(filter(validation_results_ul, provider=="vodafone"), aes(x=throughput_mbits, y=prediction_xgboost)) +
  geom_point(aes(color=scenario, shape=scenario)) +
  xlab("Measured Data Rate [MBit/s]") +
  ylab("Predicted Data Rate [MBit/s]") +
  ggtitle("Upload-Rate Predictions for Provider Vodafone")
```

Upload–Rate Predictions for Provider Vodafone



1.8 Feature Importance

1.8.1 XGBoost Gain

```
importance_ul_o2 = tibble(xgboost::xgb.importance(
  feature_names = resampling_result_ul_o2$learners[[1]]$model$regr.xgboost$model$feature_names,
  model = resampling_result_ul_o2$learners[[1]]$model$regr.xgboost$model
)) %>% mutate(provider = "o2")
```

[20:10:30] WARNING: amalgamation/./src/objective/regression_obj.cu:174: reg:linear is now deprecated

```
importance_ul_tmobile = tibble(xgboost::xgb.importance(
  feature_names = resampling_result_ul_tmobile$learners[[1]]$model$regr.xgboost$model$feature_names,
  model = resampling_result_ul_tmobile$learners[[1]]$model$regr.xgboost$model
)) %>% mutate(provider = "tmobile")
```

[20:10:30] WARNING: amalgamation/./src/objective/regression_obj.cu:174: reg:linear is now deprecated

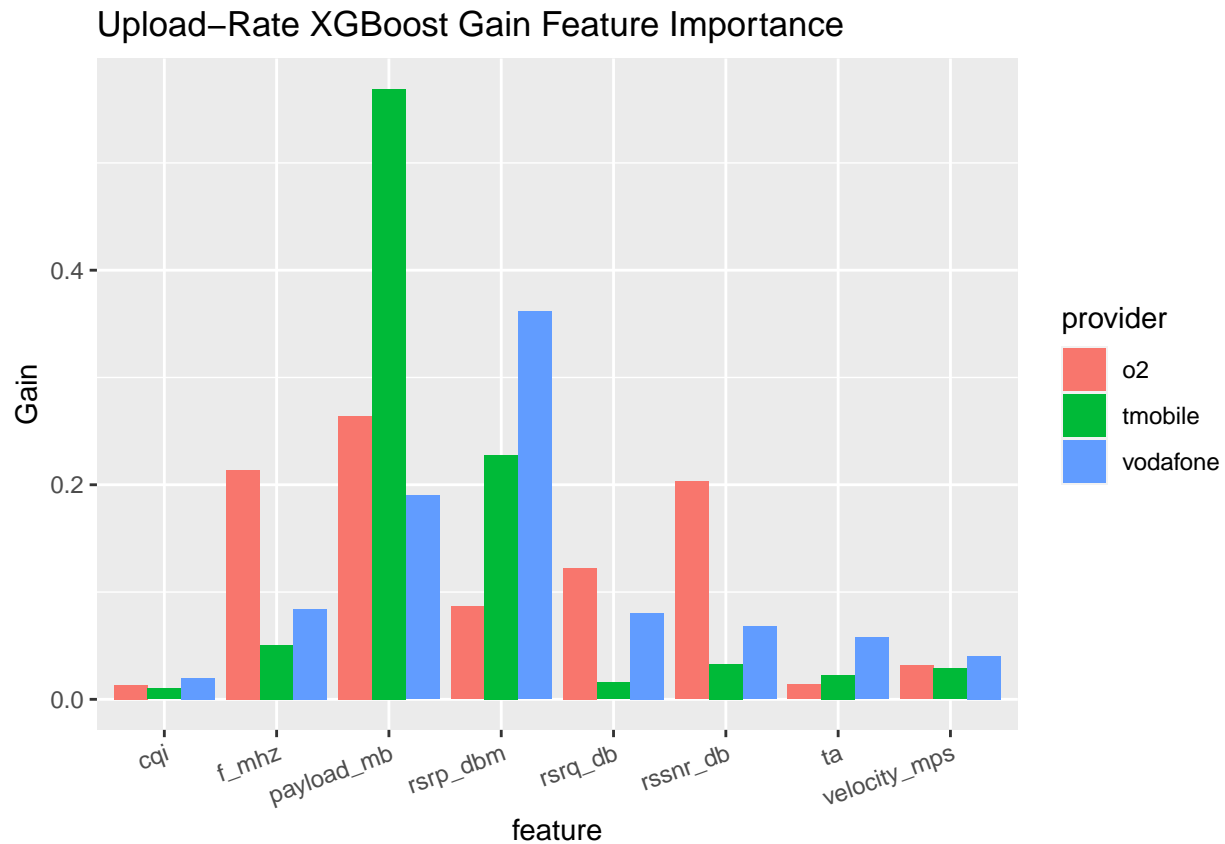
```
importance_ul_vodafone = tibble(xgboost::xgb.importance(
  feature_names = resampling_result_ul_vodafone$learners[[1]]$model$regr.xgboost$model$feature_names,
  model = resampling_result_ul_vodafone$learners[[1]]$model$regr.xgboost$model
)) %>% mutate(provider = "vodafone")
```

[20:10:30] WARNING: amalgamation/./src/objective/regression_obj.cu:174: reg:linear is now deprecated

```
importance_ul = bind_rows(
  importance_ul_o2,
  importance_ul_tmobile,
```

```
importance_ul_vodafone
) %>% filter(!str_starts(Feature, "enodeb"))

ggplot(importance_ul) +
  geom_bar(position = "dodge", aes(x = Feature, y = Gain, fill = provider), stat="identity") +
  xlab("feature") +
  ylab("Gain") +
  scale_x_discrete(guide = guide_axis(angle = 20)) +
  ggtitle("Upload-Rate XGBoost Gain Feature Importance")
```



1.8.2 Permutation

```
uninstantiate_resampling = function(resampling) {
  new_resampling = new.env()
  class(new_resampling) = class(resampling)
  for (val in ls(resampling, all.names = TRUE)) {
    if (val != "is_instantiated") {
      assign(val, get(val, envir=resampling), envir = new_resampling)
    }
  }
  new_resampling$is_instantiated = FALSE

  return(new_resampling)
}
```

```

num_permutation_sims_ul = 1

filter_permutation_o2_ul = flt("permutation",
  learner = learner_ul_o2$clone(),
  resampling = uninstantiate_resampling(
    make_outer_resampling(task_ul_o2, dataset_ul, drive_ids_train=1:7, drive_ids_test=8:10)
  ),
  measure = msr("regr.mae"),
  standardize = TRUE,
  nmc=num_permutation_sims_ul
)
filter_permutation_o2_ul$calculate(task_ul_o2)
permutation_ul_o2 = tibble(as.data.table(filter_permutation_o2_ul)) %>% mutate(provider="o2")

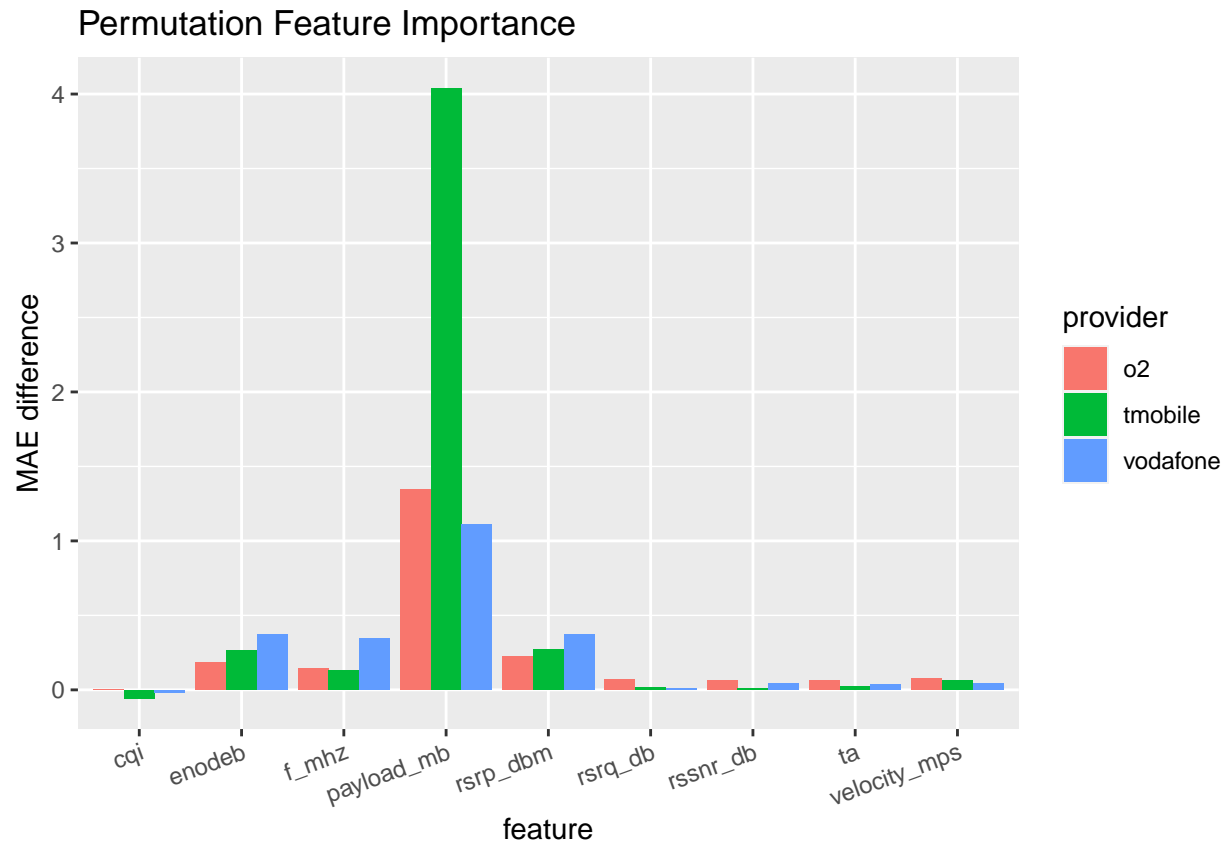
filter_permutation_tmobile_ul = flt("permutation",
  learner = learner_ul_tmobile$clone(),
  resampling = uninstantiate_resampling(
    make_outer_resampling(task_ul_tmobile, dataset_ul, drive_ids_train=1:7, drive_ids_test=8:10)
  ),
  measure = msr("regr.mae"),
  standardize = TRUE,
  nmc=num_permutation_sims_ul
)
filter_permutation_tmobile_ul$calculate(task_ul_tmobile)
permutation_ul_tmobile = tibble(as.data.table(filter_permutation_tmobile_ul)) %>% mutate(provider="tmob")

filter_permutation_vodafone_ul = flt("permutation",
  learner = learner_ul_vodafone$clone(),
  resampling = uninstantiate_resampling(
    make_outer_resampling(task_ul_vodafone, dataset_ul, drive_ids_train=1:7, drive_ids_test=8:10)
  ),
  measure = msr("regr.mae"),
  standardize = TRUE,
  nmc=num_permutation_sims_ul
)
filter_permutation_vodafone_ul$calculate(task_ul_vodafone)
permutation_ul_vodafone = tibble(as.data.table(filter_permutation_vodafone_ul)) %>% mutate(provider="vodafone")

permutation_ul = bind_rows(
  permutation_ul_o2,
  permutation_ul_tmobile,
  permutation_ul_vodafone
)

ggplot(permutation_ul) +
  geom_bar(position = "dodge", aes(x = feature, y = score, fill = provider), stat="identity") +
  xlab("feature") +
  ylab("MAE difference") +
  scale_x_discrete(guide = guide_axis(angle = 20)) +
  ggtitle("Permutation Feature Importance")

```



2 Download-Rate Prediction

2.1 Reading the Data

```
dataset_dl = read_csv(
  str_c(data_dir, "dataset_dl.csv"),
  col_types = cols(
    drive_id = col_integer(),
    scenario = col_factor(),
    provider = col_factor(),
    ci = col_factor(),
    enodeb = col_factor()
  )
) %>% select(
  drive_id,
  timestamp,
  scenario,
  provider,
  velocity_mps,
  rsrp_dbm,
  rsrq_db,
  rssnr_db,
  cqi,
  ta,
  enodeb,
```



```
## $ scenario      <fct> campus, campus, campus, campus, campus, campus, ca...
## $ provider      <fct> vodafone, vodafone, vodafone, vodafone, vodafone, ...
## $ velocity_mps  <dbl> 11.70, 8.22, 8.00, 10.60, 10.30, 12.28, 11.45, 0.0...
## $ rsrp_dbm      <dbl> -121, -108, -111, -113, -106, -110, -93, -94, -95,...
## $ rsrq_db       <dbl> -15, -9, -13, -11, -8, -9, -5, -7, -7, -8, -6, -6,...
## $ rssnr_db      <dbl> -8, 2, 6, 1, 5, 9, 21, 23, 23, 24, 14, 23, 13, 1, ...
## $ cqi           <dbl> 4, 2, 6, 6, 11, 10, 14, 15, 12, 15, 12, 14, 15, 6,...
## $ ta            <dbl> 63, 21, 16, 16, 16, 16, 16, 16, 16, 16, 16, 16, 16...
## $ enodeb        <fct> 51044, 52316, 50026, 50026, 50026, 50026, 50026, 5...
## $ f_mhz         <dbl> 1865, 1865, 1865, 1865, 1865, 1865, 1865, 1865, 18...
## $ payload_mb    <dbl> 0.1, 5.0, 1.0, 3.0, 8.0, 4.0, 0.5, 5.0, 6.0, 0.5, ...
## $ throughput_mbits <dbl> 3.54, 18.57, 5.22, 3.97, 11.68, 35.91, 25.32, 62.7...
```

2.2 Create the Prediction Tasks for Each Provider

```
task_dl_o2 = make_task(dataset_dl_o2, "task_dl_o2")
task_dl_o2
```

```
## <TaskRegr:task_dl_o2> (2033 x 10)
## * Target: throughput_mbits
## * Properties: -
## * Features (9):
##   - dbl (8): cqi, f_mhz, payload_mb, rsrp_dbm, rsrq_db, rssnr_db, ta,
##     velocity_mps
##   - fct (1): enodeb
```

```
task_dl_tmobile = make_task(dataset_dl_tmobile, "task_dl_tmobile")
task_dl_tmobile
```

```
## <TaskRegr:task_dl_tmobile> (2300 x 10)
## * Target: throughput_mbits
## * Properties: -
## * Features (9):
##   - dbl (8): cqi, f_mhz, payload_mb, rsrp_dbm, rsrq_db, rssnr_db, ta,
##     velocity_mps
##   - fct (1): enodeb
```

```
task_dl_vodafone = make_task(dataset_dl_vodafone, "task_dl_vodafone")
task_dl_vodafone
```

```
## <TaskRegr:task_dl_vodafone> (2170 x 10)
## * Target: throughput_mbits
## * Properties: -
## * Features (9):
##   - dbl (8): cqi, f_mhz, payload_mb, rsrp_dbm, rsrq_db, rssnr_db, ta,
##     velocity_mps
##   - fct (1): enodeb
```

2.3 Parameter Tuning

```
tuning_result_dl_o2 = get_tuning_result(task_dl_o2, dataset_dl, grid_resolution = 20, n_evals = 10)

tuning_result_dl_tmobile = get_tuning_result(task_dl_tmobile, dataset_dl, grid_resolution = 20, n_evals = 10)

tuning_result_dl_vodafone = get_tuning_result(task_dl_vodafone, dataset_dl, grid_resolution = 20, n_evals = 10)
```



```
tuning_result_dl = bind_rows(
  tibble(tuning_result_dl_o2$result) %>% mutate(provider="o2"),
  tibble(tuning_result_dl_tmobile$result) %>% mutate(provider="tmobile"),
  tibble(tuning_result_dl_vodafone$result) %>% mutate(provider="vodafone"),
) %>% select("provider", "regr.xgboost.nrounds", "regr.xgboost.eta", "regr.xgboost.gamma", "regr.xgboost.lambda")

knitr::kable(tuning_result_dl)
```

provider	regr.xgboost.nrounds	regr.xgboost.eta	regr.xgboost.gamma	regr.xgboost.lambda
o2	384	0.0100000	8.947368	2.631579
tmobile	858	0.0621053	7.894737	8.947368
vodafone	716	0.3747368	7.368421	6.842105

2.4 Create Learners with Tuned Hyperparameters

```
learner_dl_o2 = make_learner(
  nrounds = tuning_result_dl_o2$result$regr.xgboost.nrounds,
  eta = tuning_result_dl_o2$result$regr.xgboost.eta,
  gamma = tuning_result_dl_o2$result$regr.xgboost.gamma,
  lambda = tuning_result_dl_o2$result$regr.xgboost.lambda
)

learner_dl_tmobile = make_learner(
  nrounds = tuning_result_dl_tmobile$result$regr.xgboost.nrounds,
  eta = tuning_result_dl_tmobile$result$regr.xgboost.eta,
  gamma = tuning_result_dl_tmobile$result$regr.xgboost.gamma,
  lambda = tuning_result_dl_tmobile$result$regr.xgboost.lambda
)

learner_dl_vodafone = make_learner(
  nrounds = tuning_result_dl_vodafone$result$regr.xgboost.nrounds,
  eta = tuning_result_dl_vodafone$result$regr.xgboost.eta,
  gamma = tuning_result_dl_vodafone$result$regr.xgboost.gamma,
  lambda = tuning_result_dl_vodafone$result$regr.xgboost.lambda
)
```

2.5 Validation Results

```
resampling_result_dl_o2 = resample(
  task = task_dl_o2,
  learner = learner_dl_o2,
  resampling = make_outer_resampling(task_dl_o2, dataset_dl, drive_ids_train=1:7, drive_ids_test=8:10),
  store_models = TRUE
)

resampling_result_dl_tmobile = resample(
  task = task_dl_tmobile,
  learner = learner_dl_tmobile,
  resampling = make_outer_resampling(task_dl_tmobile, dataset_dl, drive_ids_train=1:7, drive_ids_test=8:10),
  store_models = TRUE
)
```

```

resampling_result_dl_vodafone = resample(
  task = task_dl_vodafone,
  learner = learner_dl_vodafone,
  resampling = make_outer_resampling(task_dl_vodafone, dataset_dl, drive_ids_train=1:7, drive_ids_test=8:10),
  store_models = TRUE
)

predictions_dl_o2 = as.data.table(resampling_result_dl_o2$prediction())
glimpse(tibble(predictions_dl_o2))

## Rows: 609
## Columns: 3
## $ row_id    <int> 146, 147, 148, 149, 150, 151, 152, 153, 154, 155, 156, 157...
## $ truth     <dbl> 3.72, 2.56, 0.69, 0.36, 0.55, 0.36, 0.37, 0.28, 0.12, 1.36...
## $ response  <dbl> 9.0443287, 1.6223296, 4.5180016, 0.6935747, 1.0766277, 0.6...

predictions_dl_tmobile = as.data.table(resampling_result_dl_tmobile$prediction())
predictions_dl_vodafone = as.data.table(resampling_result_dl_vodafone$prediction())

validation_results_dl = bind_rows(
  tibble(predictions_dl_o2) %>%
    inner_join(tibble(task_dl_o2$row_names), by="row_id") %>%
    inner_join(dataset_dl, by=c("row_name"="row_id_original")),
  tibble(predictions_dl_tmobile) %>%
    inner_join(tibble(task_dl_tmobile$row_names), by="row_id") %>%
    inner_join(dataset_dl, by=c("row_name"="row_id_original")),
  tibble(predictions_dl_vodafone) %>%
    inner_join(tibble(task_dl_vodafone$row_names), by="row_id") %>%
    inner_join(dataset_dl, by=c("row_name"="row_id_original"))
)
glimpse(validation_results_dl)

## Rows: 1,923
## Columns: 18
## $ row_id      <int> 146, 147, 148, 149, 150, 151, 152, 153, 154, 155, ...
## $ truth       <dbl> 3.72, 2.56, 0.69, 0.36, 0.55, 0.36, 0.37, 0.28, 0....
## $ response    <dbl> 9.0443287, 1.6223296, 4.5180016, 0.6935747, 1.0766...
## $ row_name    <int> 146, 147, 148, 149, 150, 151, 152, 153, 154, 155, ...
## $ drive_id    <int> 8, 8, 8, 8, 8, 8, 8, 8, 9, 9, 9, 9, 9, 9, 9, 9,...
## $ timestamp   <dtm> 2018-12-11 09:04:12, 2018-12-11 09:04:23, 2018-12...
## $ scenario    <fct> campus, campus, campus, campus, campus, campus, ca...
## $ provider    <fct> o2, o2, o2, o2, o2, o2, o2, o2, o2, o2, o2, o2...
## $ velocity_mps <dbl> 0.00, 7.20, 8.74, 9.01, 9.00, 11.68, 0.00, 0.00, 1...
## $ rsrp_dbm    <dbl> -89, -92, -88, -99, -105, -102, -100, -101, -100, ...
## $ rsrq_db     <dbl> -9, -12, -12, -15, -16, -16, -17, -16, -14, -16, -...
## $ rssnr_db    <dbl> 13, 3, 7, -4, -3, -5, -7, -5, 1, -4, -7, -6, 0, -2...
## $ cqi         <dbl> 11, 5, 6, 4, 3, 2, 3, 4, 6, 5, 3, 4, 5, 6, 5, 6, 4...
## $ ta         <dbl> 7, 7, 7, 7, 7, 7, 12, 12, 12, 7, 12, 12, 12, 12, 1...
## $ enodeb     <fct> 52410, 52410, 52410, 52410, 52410, 52410, 52900, 5...
## $ f_mhz      <dbl> 850, 850, 850, 850, 850, 850, 850, 850, 850, 850, ...
## $ payload_mb  <dbl> 6.0, 2.0, 7.0, 2.0, 8.0, 2.0, 8.0, 6.0, 0.5, 1.0, ...
## $ throughput_mbits <dbl> 3.72, 2.56, 0.69, 0.36, 0.55, 0.36, 0.37, 0.28, 0....

```

```
all(validation_results_dl$truth == validation_results_dl$throughput_mbits)

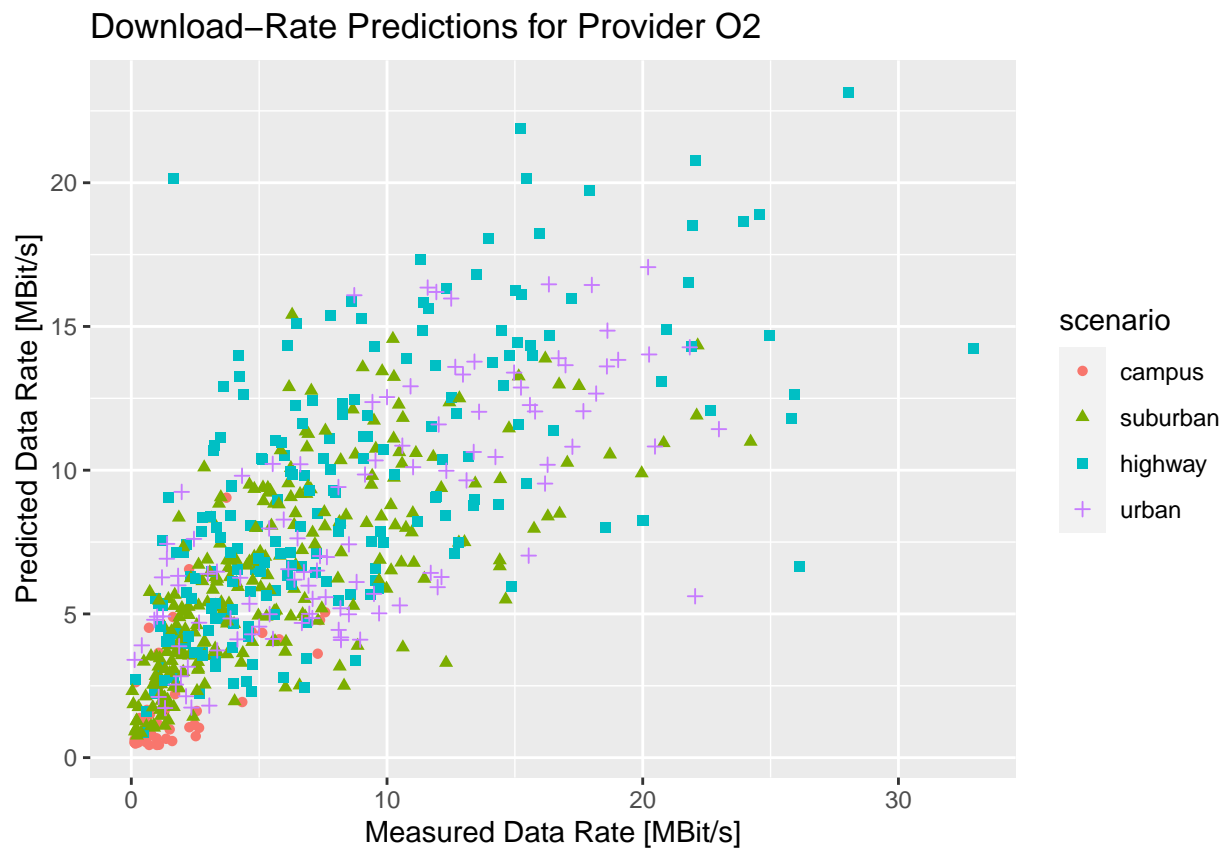
## [1] TRUE

validation_results_dl = validation_results_dl %>%
  rename(prediction_xgboost=response) %>%
  select(-truth, -row_id, -row_name)

write_csv(validation_results_dl, str_c(results_dir, "predictions_xgboost_dl.csv"))
```

2.5.1 Scatter Plots

```
ggplot(filter(validation_results_dl, provider=="o2"), aes(x=throughput_mbits, y=prediction_xgboost)) +
  geom_point(aes(color=scenario, shape=scenario)) +
  xlab("Measured Data Rate [MBit/s]") +
  ylab("Predicted Data Rate [MBit/s]") +
  ggtitle("Download-Rate Predictions for Provider O2")
```



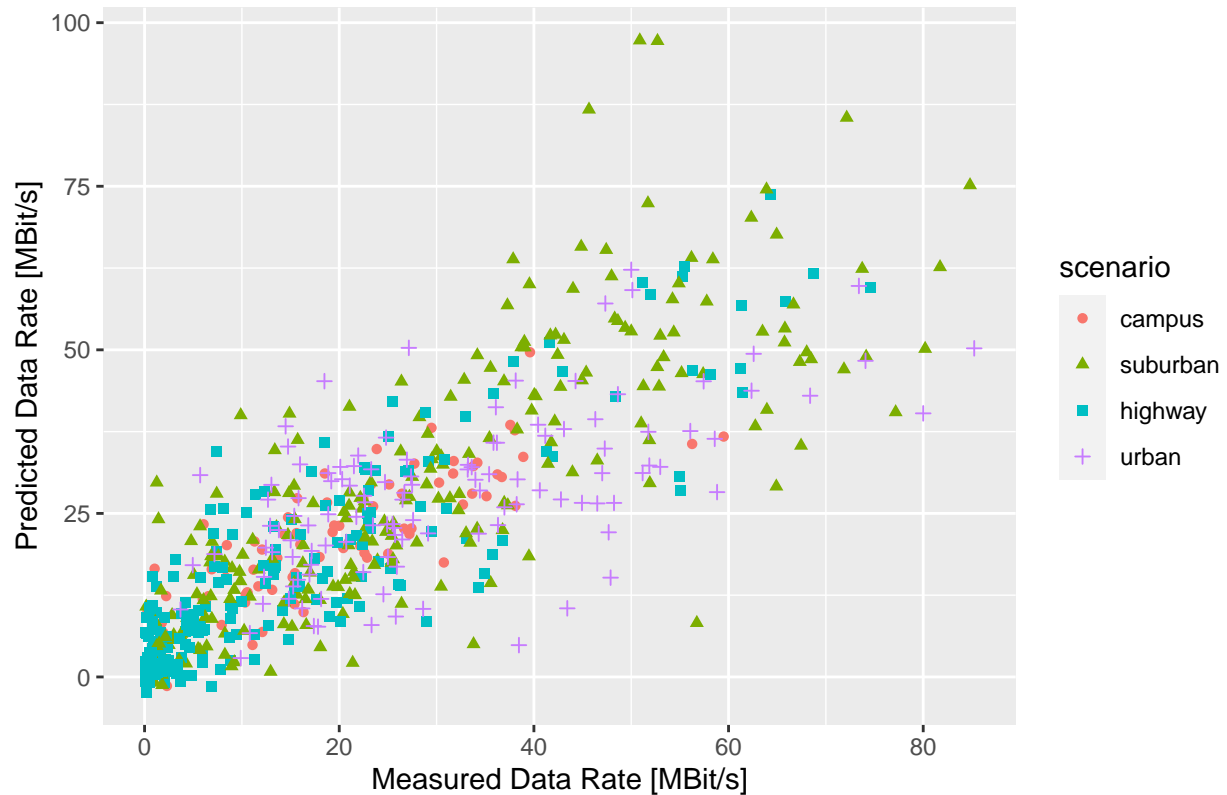
```
ggplot(filter(validation_results_dl, provider=="tmobile"), aes(x=throughput_mbits, y=prediction_xgboost)) +
  geom_point(aes(color=scenario, shape=scenario)) +
  xlab("Measured Data Rate [MBit/s]") +
  ylab("Predicted Data Rate [MBit/s]") +
  ggtitle("Download-Rate Predictions for Provider T-Mobile")
```

Download-Rate Predictions for Provider T-Mobile



```
ggplot(filter(validation_results_dl, provider=="vodafone"), aes(x=throughput_mbits, y=prediction_xgboost)) +
  geom_point(aes(color=scenario, shape=scenario)) +
  xlab("Measured Data Rate [MBit/s]") +
  ylab("Predicted Data Rate [MBit/s]") +
  ggtitle("Download-Rate Predictions for Provider Vodafone")
```

Download–Rate Predictions for Provider Vodafone



2.6 Feature Importance

2.6.1 XGBoost Gain

```
importance_dl_o2 = tibble(xgboost::xgb.importance(
  feature_names = resampling_result_dl_o2$learners[[1]]$model$regr.xgboost$model$feature_names,
  model = resampling_result_dl_o2$learners[[1]]$model$regr.xgboost$model
)) %>% mutate(provider = "o2")
```

[20:13:49] WARNING: amalgamation/./src/objective/regression_obj.cu:174: reg:linear is now deprecated

```
importance_dl_tmobile = tibble(xgboost::xgb.importance(
  feature_names = resampling_result_dl_tmobile$learners[[1]]$model$regr.xgboost$model$feature_names,
  model = resampling_result_dl_tmobile$learners[[1]]$model$regr.xgboost$model
)) %>% mutate(provider = "tmobile")
```

[20:13:50] WARNING: amalgamation/./src/objective/regression_obj.cu:174: reg:linear is now deprecated

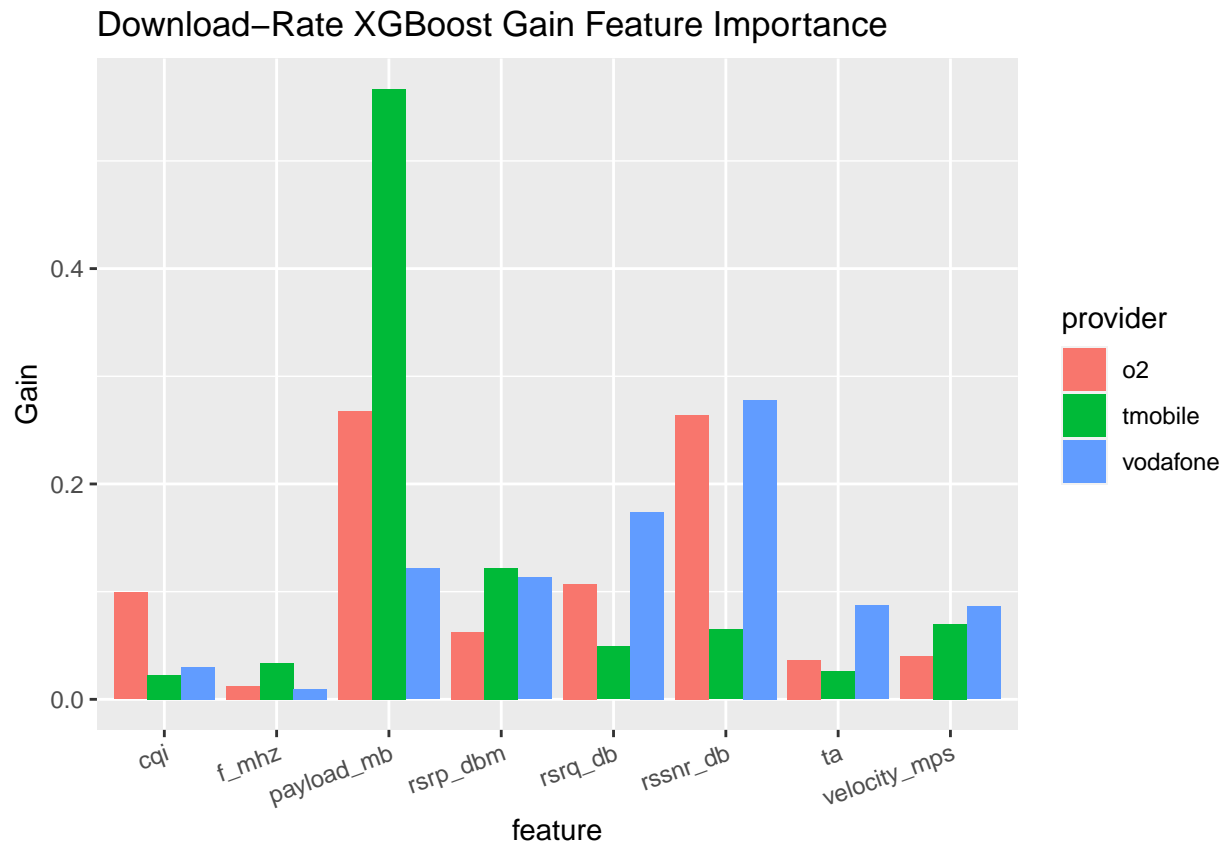
```
importance_dl_vodafone = tibble(xgboost::xgb.importance(
  feature_names = resampling_result_dl_vodafone$learners[[1]]$model$regr.xgboost$model$feature_names,
  model = resampling_result_dl_vodafone$learners[[1]]$model$regr.xgboost$model
)) %>% mutate(provider = "vodafone")
```

[20:13:50] WARNING: amalgamation/./src/objective/regression_obj.cu:174: reg:linear is now deprecated

```
importance_dl = bind_rows(
  importance_dl_o2,
  importance_dl_tmobile,
```

```
importance_dl_vodafone
) %>% filter(!str_starts(Feature, "enodeb"))

ggplot(importance_dl) +
  geom_bar(position = "dodge", aes(x = Feature, y = Gain, fill = provider), stat="identity") +
  xlab("feature") +
  ylab("Gain") +
  scale_x_discrete(guide = guide_axis(angle = 20)) +
  ggtitle("Download-Rate XGBoost Gain Feature Importance")
```



2.6.2 Permutation

```
num_permutation_sims_dl = 1

filter_permutation_o2_dl = flt("permutation",
  learner = learner_dl_o2$clone(),
  resampling = uninstantiate_resampling(
    make_outer_resampling(task_dl_o2, dataset_dl, drive_ids_train=1:7, drive_ids_test=8:10)
  ),
  measure = msr("regr.mae"),
  standardize = TRUE,
  nmc = num_permutation_sims_dl
)

filter_permutation_o2_dl$calculate(task_dl_o2)
permutation_dl_o2 = tibble(as.data.table(filter_permutation_o2_dl)) %>% mutate(provider="o2")
```

```

filter_permutation_tmobile_dl = flt("permutation",
  learner = learner_dl_tmobile$clone(),
  resampling = uninstantiate_resampling(
    make_outer_resampling(task_dl_tmobile, dataset_dl, drive_ids_train=1:7, drive_ids_test=8:10)
  ),
  measure = msr("regr.mae"),
  standardize = TRUE,
  nmc = num_permutation_sims_dl
)
filter_permutation_tmobile_dl$calculate(task_dl_tmobile)
permutation_dl_tmobile = tibble(as.data.table(filter_permutation_tmobile_dl)) %>% mutate(provider="tmob")

filter_permutation_vodafone_dl = flt("permutation",
  learner = learner_dl_vodafone$clone(),
  resampling = uninstantiate_resampling(
    make_outer_resampling(task_dl_vodafone, dataset_dl, drive_ids_train=1:7, drive_ids_test=8:10)
  ),
  measure = msr("regr.mae"),
  standardize = TRUE,
  nmc = num_permutation_sims_dl
)
filter_permutation_vodafone_dl$calculate(task_dl_vodafone)
permutation_dl_vodafone = tibble(as.data.table(filter_permutation_vodafone_dl)) %>% mutate(provider="vodafone")

permutation_dl = bind_rows(
  permutation_dl_o2,
  permutation_dl_tmobile,
  permutation_dl_vodafone
)

ggplot(permutation_dl) +
  geom_bar(position = "dodge", aes(x = feature, y = score, fill = provider), stat="identity") +
  xlab("feature") +
  ylab("MAE difference") +
  scale_x_discrete(guide = guide_axis(angle = 20)) +
  ggtitle("Permutation Feature Importance")

```

