

# Towards A Clear Understanding Of Rural Internet: What Statistical Measures Can Be Used To Assess, Compare And Forecast Internet Speeds For Rural Canadian Communities?

A Consulting Project for Math 6627 (1/3)

Benjamin Smith

09 February 2022

## Introduction

(Quoted from the SCC website)

The Government of Canada has committed to helping 95% of Canadian households and businesses access high-speed internet at minimum speeds of 50 Mbps download and 10 Mbps upload (hereinafter referred to as the “Commitment”) by 2026, and 100% by 2030. According to the CRTC, currently 45.6% of rural community households have access to the Commitment based on what’s available to them via an Internet Service Provider (e.g. Shaw, Telus, etc.) in their region, rather than what a rural household actually realizes at home in terms of internet speeds.

For this case study, the SCC would like to understand the state of internet connectivity in both rural and underserved Canadian communities using consumer-provided data. The SCC claims that by using data directly from the consumer, it is possible to better understand connectivity in these communities as measured by the consumers in their own homes.

Specifically, the following is desired:

1. A statistical analysis of the current realized and forecasted internet speeds (upload and download) for rural and underserved communities in terms of progress towards the Commitment;
2. A comparative analysis of rural and underserved communities in terms of progress towards the Commitment; and
3. The identification of statistically reliable methods to assess and compare rural and underserved communities’s realized internet access. For this study in particular, the identification of reliable and reproducible statistical methods to understand connectivity of rural and underserved Canadian communities is critical.

## The Data

The data was made available by Ookla and Statistics Canada. One of the first things to check regarding the data is to see if any missing data is present in the dataset. The visual below shows that most of the missing data is related to population center information. Namely data on population center id, type and class (PCUID, PCTYPE, PCLASS).

```
library(tidyverse)
library(ggthemes)
library(ggspatial)
```

```
library(plotly)
library(rnaturalearth)
library(sf)
library(scales)
library(reshape2)
#dt<- readr::read_csv("./ConsultingData/ookla-canada-speed-tiles.csv")
# Accomidating for mtor
setwd("/home/ben2908")
dt<- readr::read_csv("./ookla-canada-speed-tiles.csv")

naniar::vis_miss(dt ,warn_large_data = FALSE)+
  theme(axis.text.x=element_text(angle=90))
```

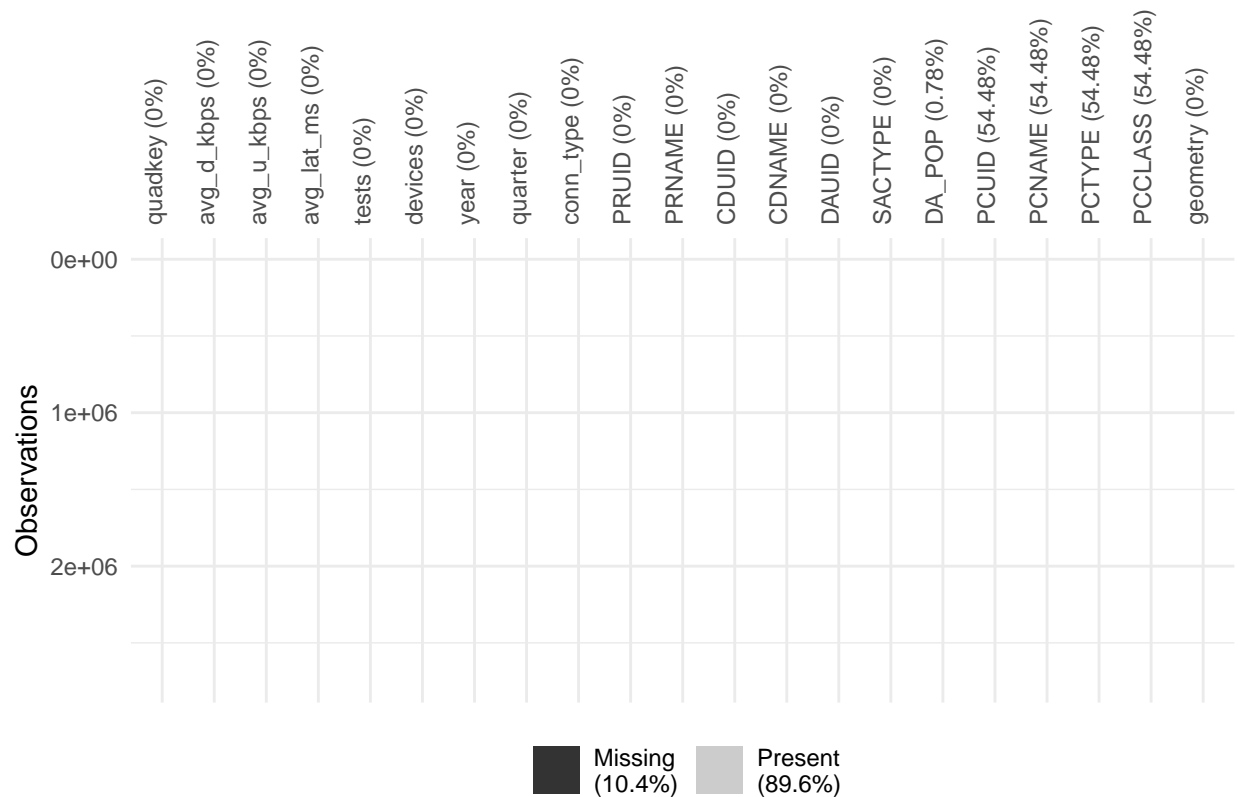


Figure 1: Missing data present in dataset provided

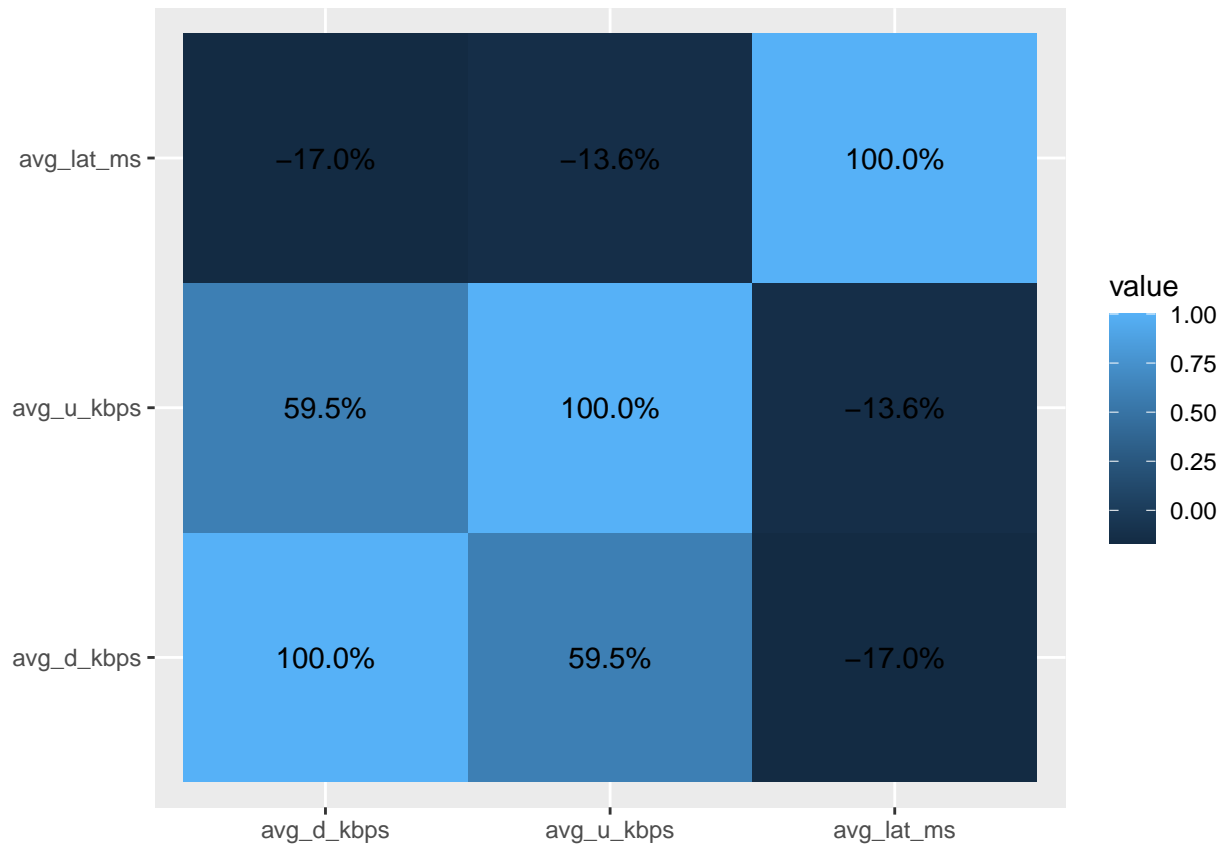
Since this information is not relevant to the analysis it does not present any challenges.

## Analysis

### Correlation Between Upload, Download and Latency times

```
dt %>%
  dplyr::select(avg_d_kbps, avg_u_kbps, avg_lat_ms) %>%
  cor() %>%
  melt() %>%
```

```
ggplot(mapping=aes(x=Var2,y=Var1,fill=value))+
  geom_tile()+
  scale_fill_gradient()+
  geom_text(mapping=aes(label=scales::percent(value)))+
  theme(axis.title = element_blank())
```



## Current Realized And Forecasted Internet Speeds (Upload And Download)

Due to the large number of observations, approaching this dataset with standard visuals was not easy to do. In lieu of this the use of directed acyclic graphs (DAGs) is employed to understand the relationship between variables available and develop the model accordingly.

```
library(ggdag)
```

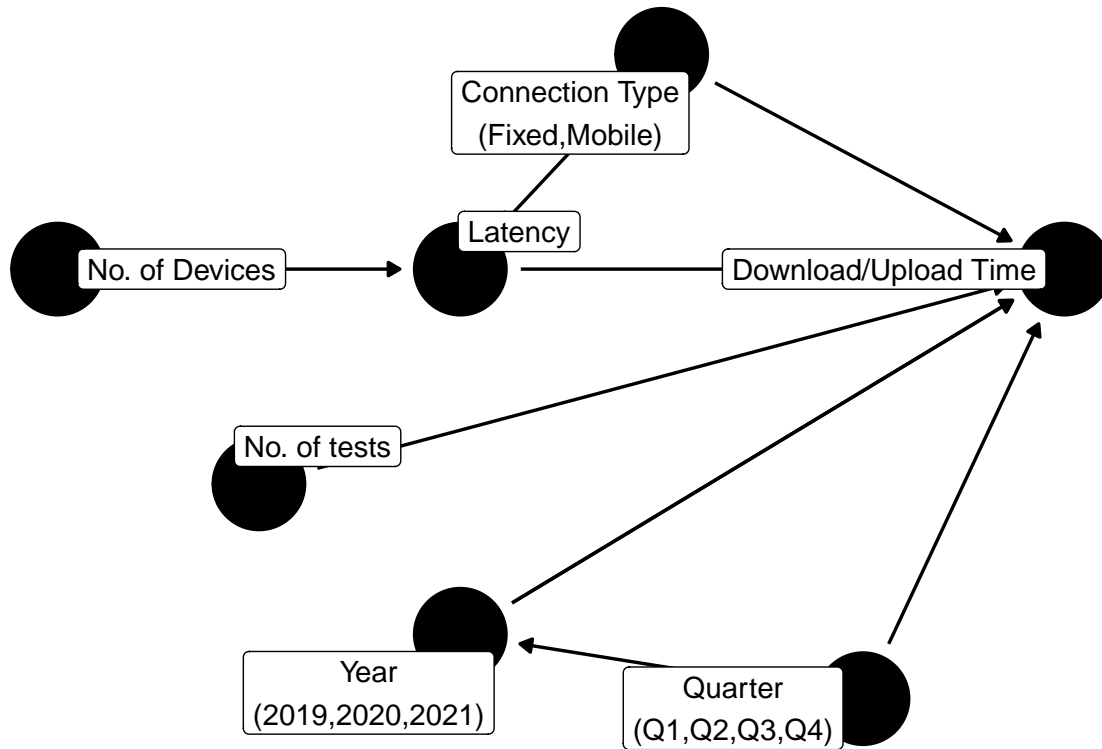
```
##
## Attaching package: 'ggdag'
## The following object is masked from 'package:stats':
##
##   filter
dagify(download_time~latency,
        download_time ~ conn_type,
        download_time ~tests,
        download_time~year,
        download_time~quarter,
        quarter~year,
```

```

latency~ conn_type,
latency~no_of_devices,
year~quarter,
labels=c("download_time"="Download/Upload Time",
         "latency" = "Latency",
         "conn_type"="Connection Type\n(Fixed,Mobile)",
         "no_of_devices"="No. of Devices",
         "tests"="No. of tests",
         "year"="Year\n(2019,2020,2021)",
         "quarter"="Quarter\n(Q1,Q2,Q3,Q4)") %>%
tidy_dagitty() %>%
mutate(xend=c(10.5,9,10.5,9,10.5,9,10.5,10.5,10.5,NA),
       yend=c(0,0,0,0,0,-1.7,0,0,0,NA),
       x=ifelse(name=="download_time",10.5,
                ifelse(name=="latency",9,
                      ifelse(name=="conn_type",9.5,
                            ifelse(name=="tests",8.5,
                                  ifelse(name=="year",9,
                                        ifelse(name=="quarter",10,
                                              8)))))),
       y=ifelse(name=="download_time",0,
                ifelse(name=="latency",0,
                      ifelse(name=="conn_type",1,
                            ifelse(name=="tests",-1,
                                  ifelse(name=="year",-1.7,
                                        ifelse(name=="quarter",-2,
                                              0)))))),
       effectType=ifelse(name %in% c("download_time",
                                     "latency",
                                     "conn_type",
                                     "no_of_devices",
                                     "tests",
                                     "year",
                                     "quarter"),"Fixed","Random")) %>%
ggdag(text=FALSE,use_labels = "label")+
ggtitle("DAG of relationship between Fixed Effects and Download/Upload Time")+
theme_dag()

```

## DAG of relationship between Fixed Effects and Download/Upload Time



It is from this DAG that the following mixed model is constructed.

$$Y = X\beta + Zb$$

Where  $X$  is the design matrix for the fixed effects,  $Z$  is the design matrix of the random effects and  $\beta$  and  $b$  are the fixed and random effects vectors. The fixed and random effects are:

Fixed Effects	Random Effect
No. of devices	Province
Connection Type	
No. of Tests	
Year	
Quarter	
Year*Quarter	

After running this model in SAS note the following:

```
X_download<- model.matrix(avg_d_kbps~ devices + conn_type+ tests + year*quarter + PRNAME,data=dt %>%
  dplyr::select(avg_d_kbps,PRNAME, devices, conn_type, tests, year,quarter))
```

```
library(lme4)
```

```
## Loading required package: Matrix
```

```
##
```

```
## Attaching package: 'Matrix'
```

```

## The following objects are masked from 'package:tidyr':
##
##   expand, pack, unpack
fit_download_time <- lmer(avg_d_kbps~ devices + conn_type+ tests + year*quarter + (1|PRNAME), data= data)

## Warning: Some predictor variables are on very different scales: consider
## rescaling
summary(fit_download_time)

## Linear mixed model fit by REML ['lmerMod']
## Formula: avg_d_kbps ~ devices + conn_type + tests + year * quarter + (1 |
##   PRNAME)
##   Data: data.frame(dt, stringsAsFactors = T)
##
## REML criterion at convergence: 70171806
##
## Scaled residuals:
##      Min       1Q   Median       3Q      Max
## -13.2528  -0.6578  -0.2950   0.4008  11.2018
##
## Random effects:
##   Groups   Name                Variance Std.Dev.
##   PRNAME   (Intercept)  6.664e+08  25814
##   Residual                        6.975e+09  83515
## Number of obs: 2751464, groups:  PRNAME, 13
##
## Fixed effects:
##              Estimate Std. Error t value
## (Intercept)  -2.622e+07  2.632e+05 -99.603
## devices       2.252e+03  1.015e+01  221.957
## conn_typemobile 2.324e+03  1.273e+02  18.248
## tests         2.173e+01  2.234e+00   9.727
## year          1.301e+04  1.303e+02  99.870
## quarterQ2     -1.183e+07  3.686e+05 -32.095
## quarterQ3     -1.160e+07  3.581e+05 -32.382
## quarterQ4     -2.035e+07  3.606e+05 -56.438
## year:quarterQ2  5.858e+03  1.825e+02  32.102
## year:quarterQ3  5.742e+03  1.773e+02  32.391
## year:quarterQ4  1.008e+04  1.785e+02  56.479
##
## Correlation of Fixed Effects:
##              (Intr) devices cnn_ty tests  year  qrtrQ2 qrtrQ3 qrtrQ4 yr:qQ2
## devices      -0.012
## conn_typmb1   0.011  0.065
## tests         0.023 -0.756  0.053
## year          -1.000  0.012 -0.011 -0.023
## quarterQ2     -0.713  0.009  0.000 -0.013  0.714
## quarterQ3     -0.735  0.003 -0.012 -0.016  0.735  0.524
## quarterQ4     -0.730 -0.004 -0.004 -0.015  0.730  0.521  0.536
## year:qrtrQ2    0.713 -0.009  0.000  0.013 -0.714 -1.000 -0.524 -0.521
## year:qrtrQ3    0.735 -0.003  0.012  0.016 -0.735 -0.524 -1.000 -0.536  0.524
## year:qrtrQ4    0.730  0.004  0.004  0.015 -0.730 -0.521 -0.536 -1.000  0.521
##
##              yr:qQ3

```

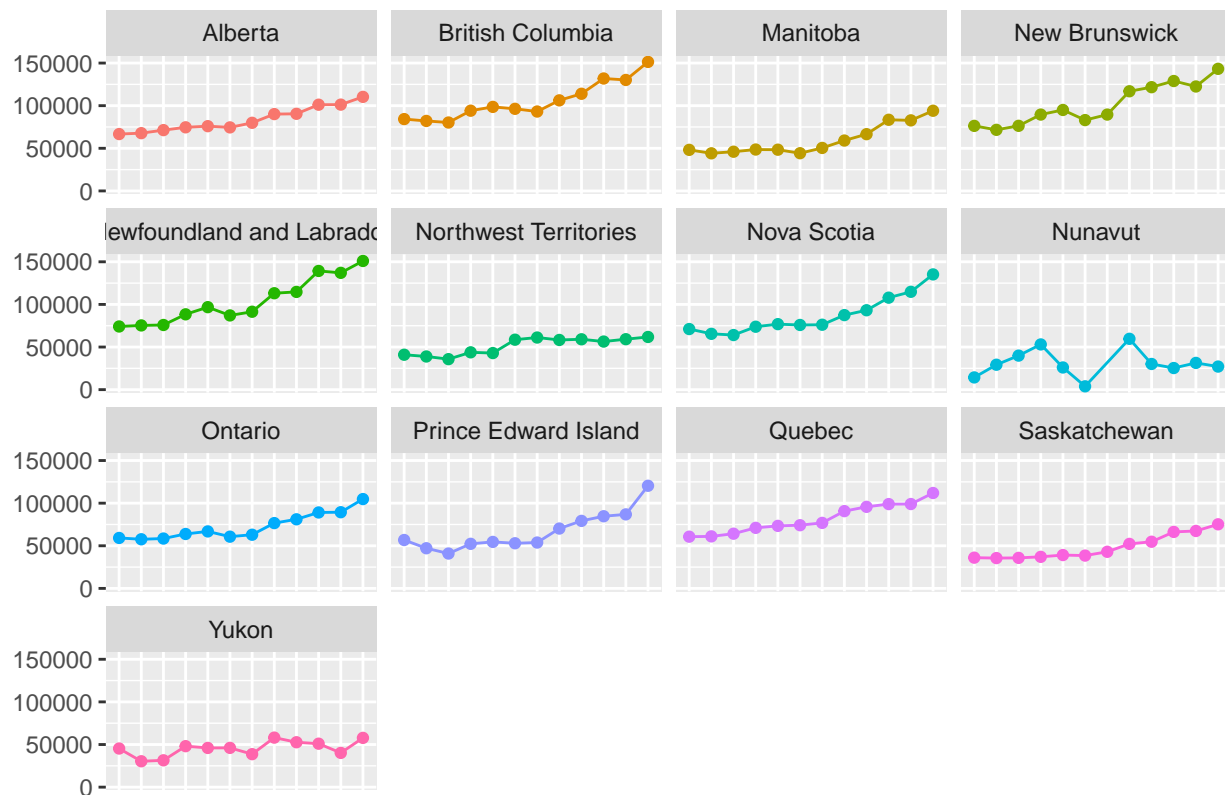
```
## devices
## conn_ttypmb1
## tests
## year
## quarterQ2
## quarterQ3
## quarterQ4
## year:qrtrQ2
## year:qrtrQ3
## year:qrtrQ4 0.536
## fit warnings:
## Some predictor variables are on very different scales: consider rescaling
```

## Rural and Underserved communities in terms of progress towards the Commitment

```
dt %>%
  dplyr::select(avg_d_kbps, PRNAME, devices, conn_type, tests, year, quarter) %>%
  group_by(PRNAME, year, quarter) %>%
  summarise(avg_download_province = mean(avg_d_kbps)) %>%
  mutate(year_quarter=paste(year, quarter),
         Province= ifelse(grepl("\\/", PRNAME), PRNAME %>% str_extract('.*(=?= \\/)'), PRNAME)) %>%
  ggplot(mapping=aes(x= year_quarter, y=avg_download_province, color=Province, group=Province))+
  geom_point()+
  geom_line()+
  facet_wrap(~Province)+
  ggtitle("Download Speed")+
  theme(legend.position = "none",
        axis.ticks.x = element_blank(),
        axis.text.x = element_blank(),
        axis.title = element_blank())
```

## `summarise()` has grouped output by 'PRNAME', 'year'. You can override using the `.groups` argument.

## Download Speed

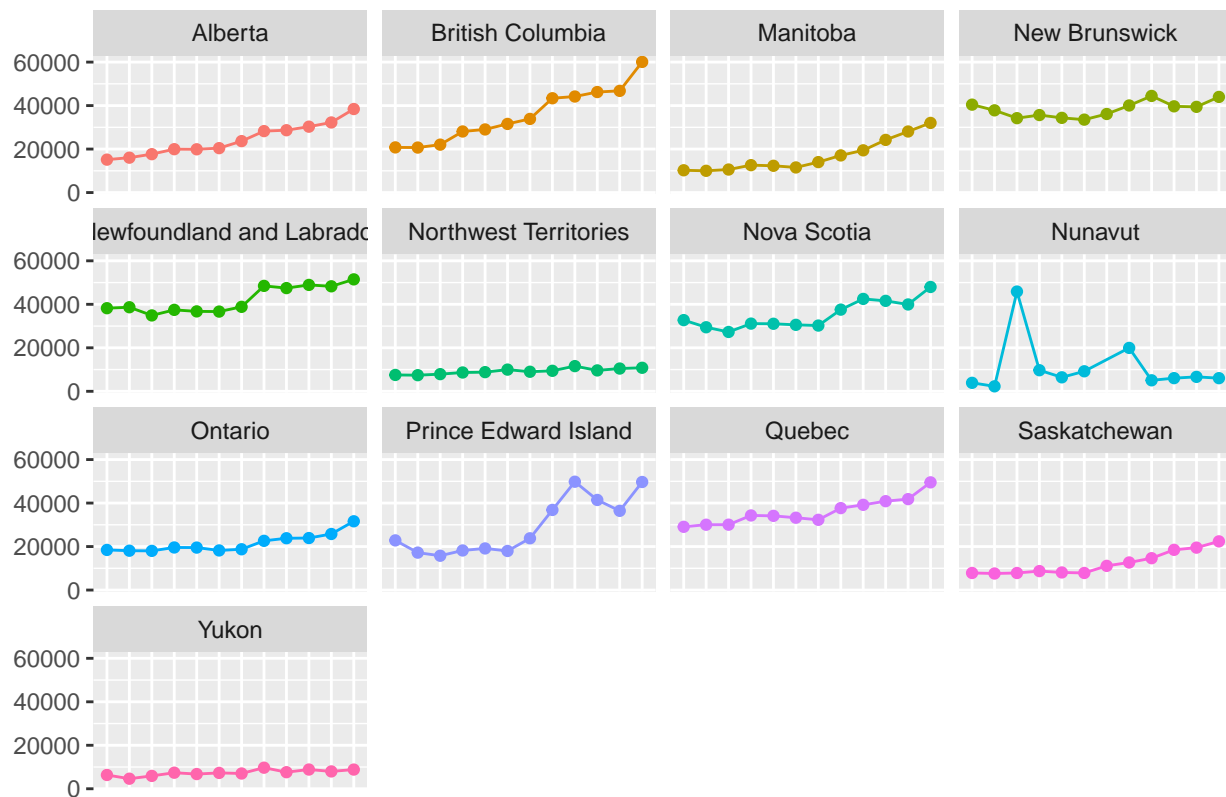


```
dt %>%
  dplyr::select(avg_u_kbps, PRNAME, devices, conn_type, tests, year, quarter) %>%
  group_by(PRNAME, year, quarter) %>%
  summarise(avg_upload_province = mean(avg_u_kbps)) %>%
  mutate(year_quarter = paste(year, quarter),
         Province = ifelse(grepl("\\/", PRNAME), PRNAME %>% str_extract('.*(=? \\/)'), PRNAME)) %>%
  ggplot(mapping = aes(x = year_quarter, y = avg_upload_province, color = Province, group = Province)) +
  geom_point() +
  geom_line() +
  facet_wrap(~Province) +
  ggtitle("Upload Speed") +
  theme(legend.position = "none",
        axis.ticks.x = element_blank(),
        axis.text.x = element_blank(),
        axis.title = element_blank())
```

## `summarise()` has grouped output by 'PRNAME', 'year'. You can override using the `.groups` argument.



## Upload Speed



The identification of statistically reliable methods to assess and compare rural and underserved communities's realized internet access.

- Test test test.

## Maps

```
# Convert dataset to sf object
dt <- dt %>% st_as_sf(wkt = "geometry")

world_map <- ne_countries(scale = "large", returnclass = 'sf')
canada_map <- world_map %>% filter(name == "Canada")

# dt %>%
#   ggplot()+
#   geom_sf()

# Loading images
setwd("/home/ben2908")

df <- read_sf('lookla-canada-speed-tiles.shp')

df2 <- read_sf('lpr_000b16a_e.shp')

# ggplot()+
#   geom_sf(data=df2)+
```

```
# geom_sf(data=df,mapping=aes(color=avg_d_kbps))+  
# ggtitle("Average Download Times in Canada")+  
# facet_wrap(year~quarter)
```

```
# ggplot()+  
# geom_sf(data=df2)+  
# geom_sf(data=df,mapping=aes(color=avg_u_kbps))+  
# ggtitle("Average Upload Times in Canada")+  
# facet_wrap(year~quarter)
```

```
# ggplot()+  
# geom_sf(data=df2)+  
# geom_sf(data=df,mapping=aes(color=avg_lat_ms))+  
# ggtitle("Average Latency Times in Canada")+  
# facet_wrap(year~quarter)
```

## References

## Code Appendix

### SAS Code

```
/*Update File Path Accordingly*/  
FILENAME REFFILE '.../ookla-canada-speed-tiles.csv';  
PROC IMPORT DATAFILE=REFFILE  
    DBMS=CSV  
    OUT=DT;  
    GETNAMES=YES;  
RUN;  
  
PROC CONTENTS DATA=DT;  
RUN;  
  
/*Download Time*/  
proc mixed data=DT method=reml covtest;  
class PRNAME quarter year conn_type;  
model avg_d_kbps= devices conn_type tests year quarter quarter*year;  
random PRNAME/s;  
run;  
  
/*Upload Time*/  
proc mixed data=DT method=reml covtest;  
class PRNAME quarter year conn_type;  
model avg_u_kbps= devices conn_type tests year quarter quarter*year;  
random PRNAME/s;  
run;
```