

forecast_timeseries_egjohnson

Elizabeth

March 28, 2019

1. Read in data

```
# read in data as character
tcdat.raw <- fread("app_data.csv", colClasses=c("character"))
tcdat.raw %>% tail
```

begin_time	average_response_time	throughput
1509234840	121.66768262254809	130255
1509234900	116.72438778317895	130836
1509234960	115.62114492094678	116567
1509235020	114.93987520315277	118753
1509235080	123.2707866886065	111213
1509235140	161.7974806350771	129874

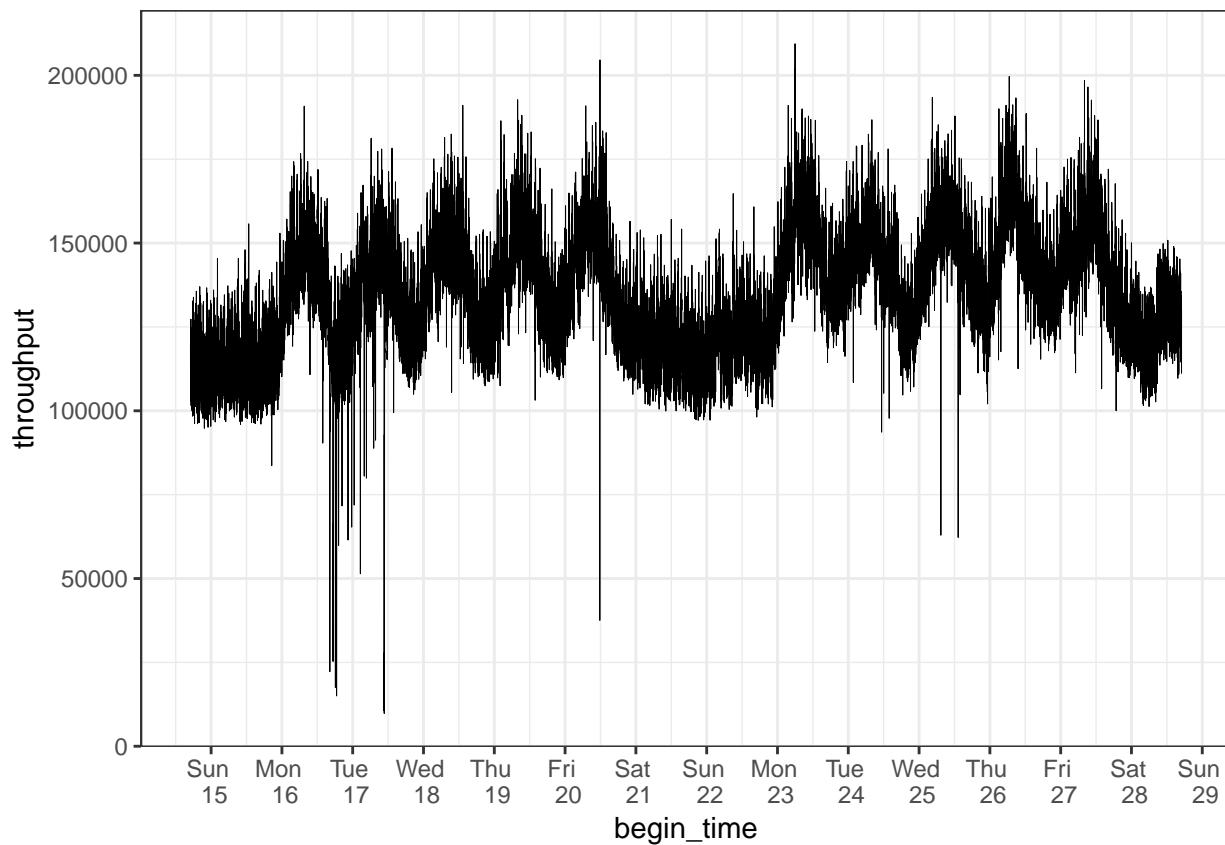
```
#change UNIX timestamp to actual date
#seconds since Jan 01 1970. (UTC)
#convert characters to numeric
library(tidyr)
tcdat.raw %>% dplyr::mutate(begin_time=anytime(as.integer(begin_time))) %>% dplyr::mutate_if(is.character,
tcdat %>% head
```

begin_time	average_response_time	throughput
2017-10-14 17:00:00	203.2001	119586
2017-10-14 17:01:00	144.1432	112493
2017-10-14 17:02:00	127.0975	108595
2017-10-14 17:03:00	134.4365	104563
2017-10-14 17:04:00	127.9597	119842
2017-10-14 17:05:00	138.5804	111901

2. Inspect Data

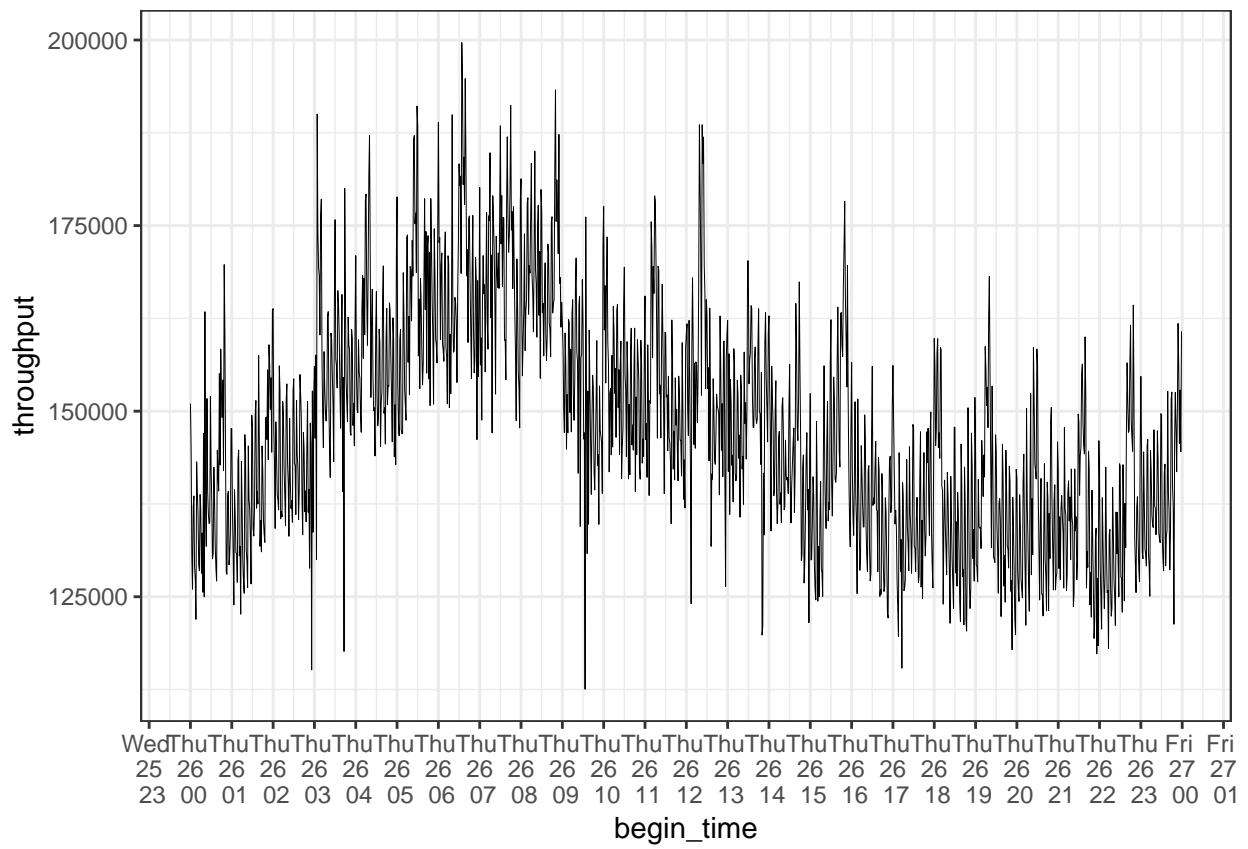
Throughput peaks in the middle of the day - but only on weekdays

```
ggplot(tcdat, aes(x = begin_time, y = throughput)) +
  geom_line(size = 0.1) +
  scale_x_datetime(breaks = date_breaks("1 day"), labels=date_format("%a \n %d", tz=Sys.timezone()))+
  theme_bw()
```



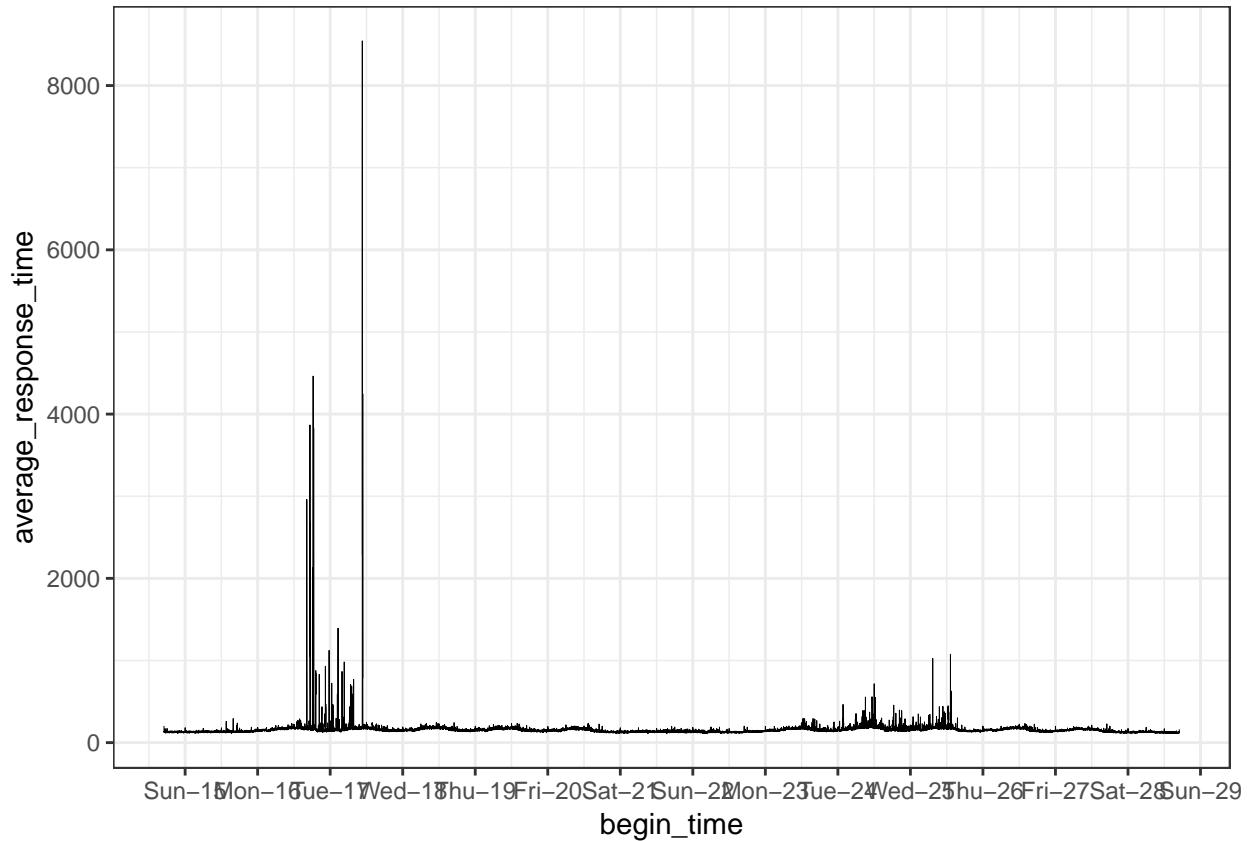
On a weekday throughput typically peaks prior to noon then tapers off.

```
tcdat %>% dplyr::filter(day(begin_time) == 26) -> tc.plot  
  
ggplot(tc.plot, aes(x = begin_time, y = throughput)) +  
  geom_line(size = 0.1) +  
  scale_x_datetime(breaks = date_breaks("1 hour"), labels = time_format("%a\n%d\n%H", tz = Sys.timezone()))  
  theme_bw()
```



response time spikes on some weekdays

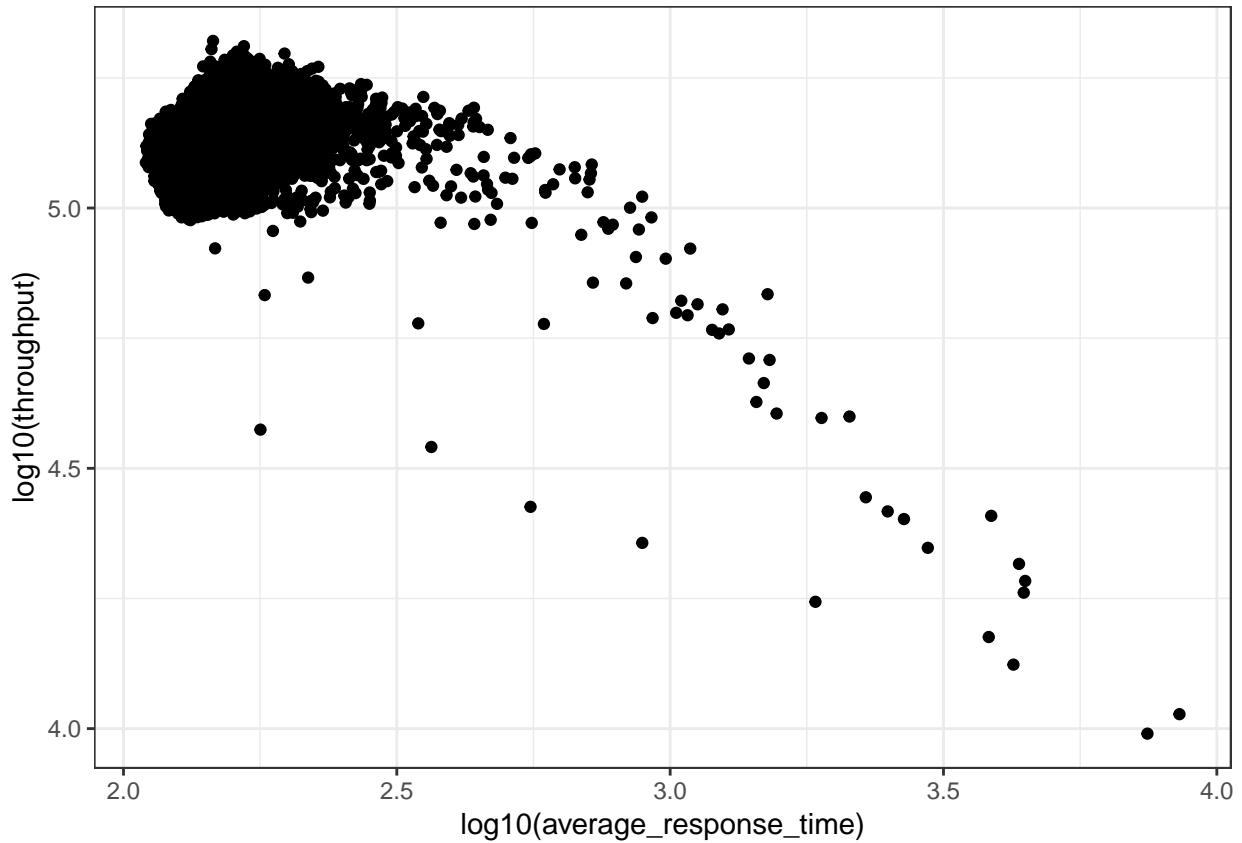
```
ggplot(tcdat, aes(x = begin_time, y = average_response_time)) +
  geom_line(size = 0.1) +
  scale_x_datetime(breaks = date_breaks("1 day"), labels=date_format("%a-%d", tz=Sys.timezone()))+
  theme_bw()
```



3. Characterize relationship/processes linking two data sets

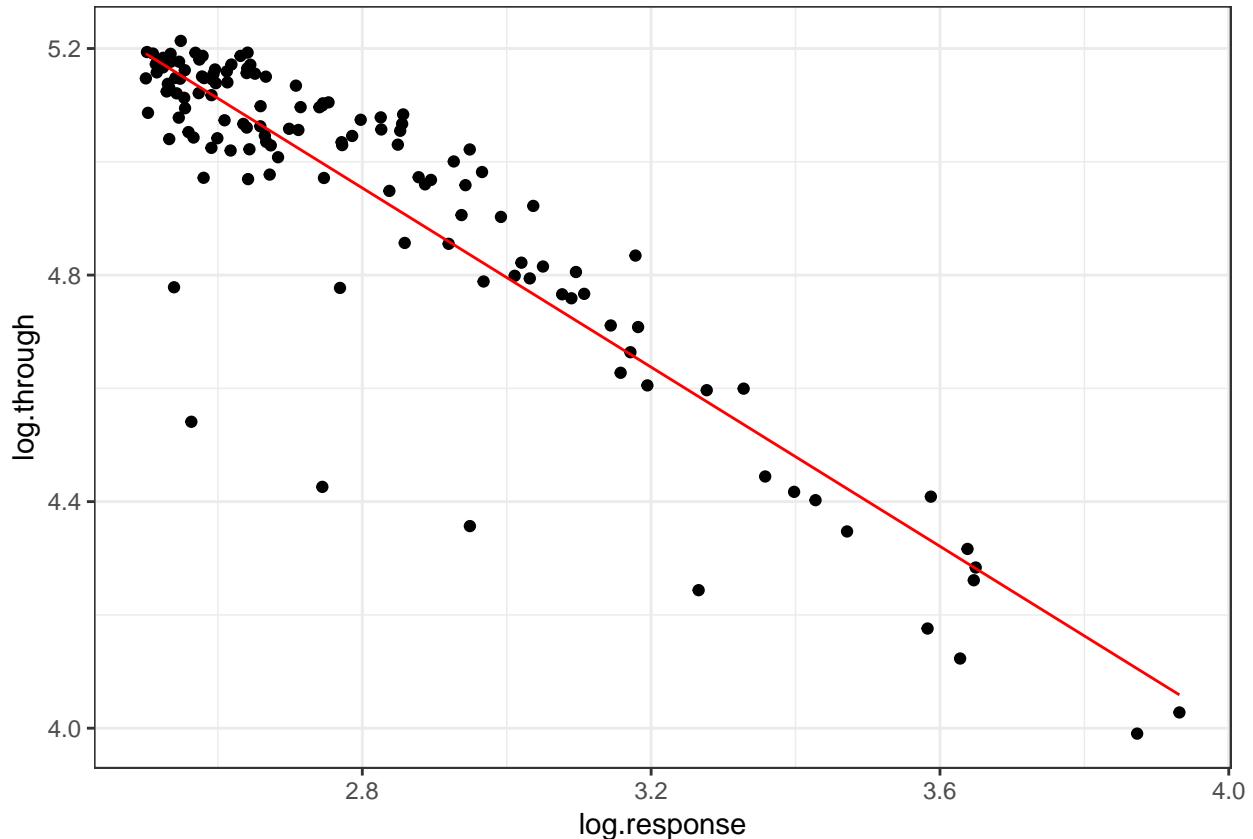
Graph of throughput and response time (log-log)

```
#plot model of both data sets vs eachother
ggplot(tcdat, aes(y=log10(throughput), x =log10(average_response_time))) + geom_point()+
  geom_point(size =0.5) +
  theme_bw()
```



Process 1: extremely long response times bottleneck throughput (log-log)

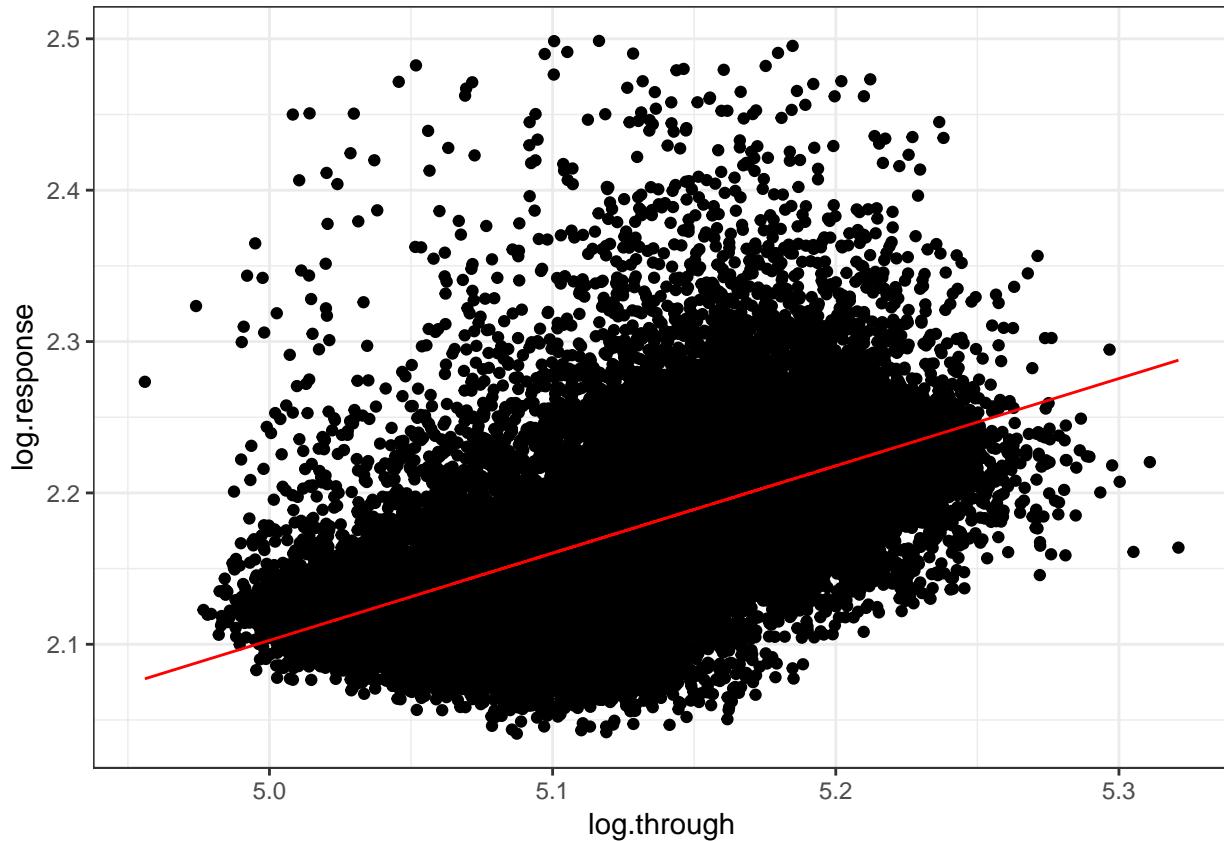
```
#take log of both throughput and response time
tcdat %>% mutate(log.response=log10(average_response_time),log.through=log10(throughput))
#create subset of data that has response time effect
tcdat.res<-tcdat %>% filter(log.response>2.5)
# estimate impact response has on throughput
res.impact.throughput.mod <-lm(log.through ~ log.response,data=tcdat.res)
#plot model of responses impact on throughput
tcdat.res$predict.throughput<-predict(res.impact.throughput.mod,log.response=tcdat.res$log.response,data=tcdat.res)
ggplot(tcdat.res, aes(y=log.through, x =log.response)) + geom_point() +
  geom_point(size =0.5) +
  geom_line(color='red',aes(y=tcdat.res$predict.throughput,x=tcdat.res$log.response))+ theme_bw()
```



Process 2: Increased throughput can slow down response time (log-log)

```
#create subset of data around throughput impact of throughput on response time
# (the more traffic the slower the response time)
tcdat.thr<-tcdat %>% filter(log.response<2.5 & log.through > 4.95)
# estimate impact response has on throughput
response.throughput.mod <-lm(log.response ~ log.through,data=tcdat.thr)
#plot model of responses impact on throughput
tcdat.thr$predict.response<-predict(response.throughput.mod,log.through=tcdat.thr$log.through,data=tcdat.thr)
response.resid<-tcdat.thr$average_response_time - 10^(predict(response.throughput.mod)) #calculate residuals
```

```
ggplot(tcdat.thr, aes(x=log.through, y =log.response)) + geom_point()+
  geom_point(size =0.5) +
  geom_line(color='red',aes(y=tcdat.thr$predict.response,x=tcdat.thr$log.through))+theme_bw()
```



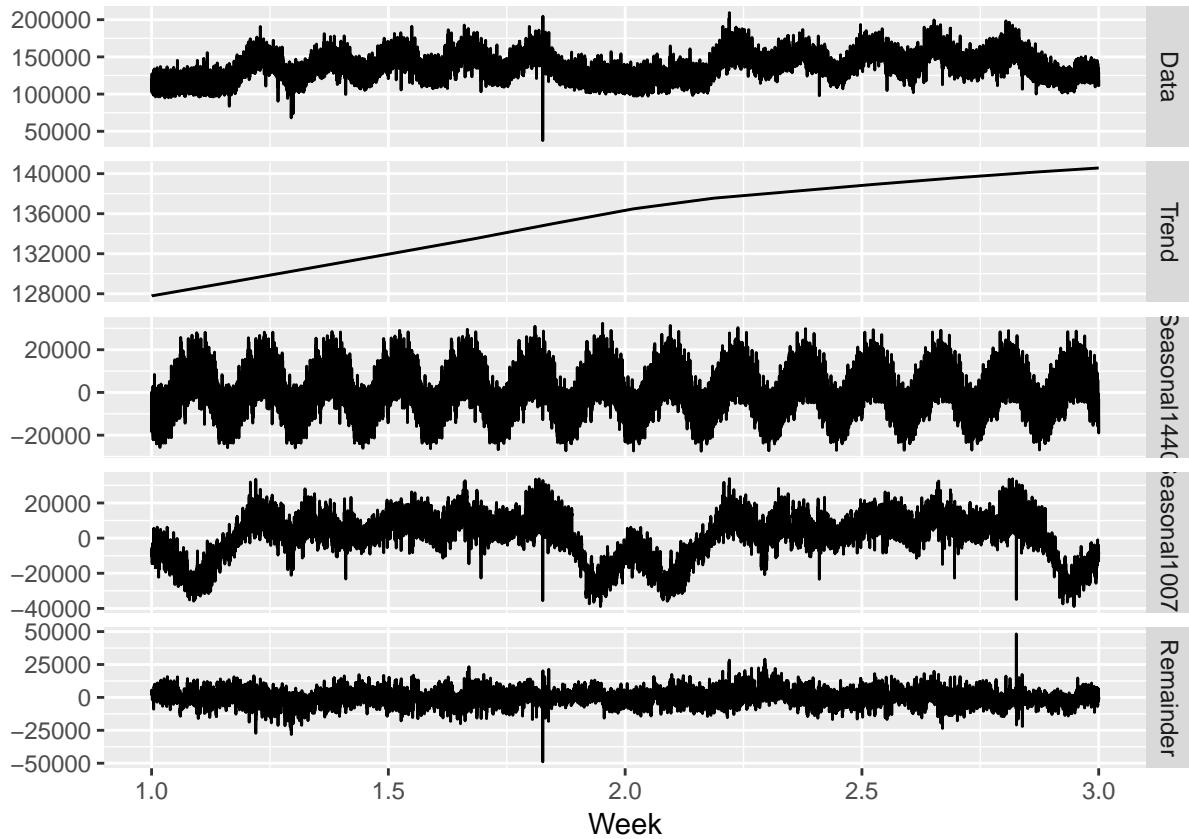
4. Decompose timeseries and account for any remaining structured noise

account for structured noise in forecasting model by subtracting out throughput spikes caused by extremely slow response times

```
# #1. put an NA in throughput outlier/spikes
tcdat %>% mutate(throughput=case_when(log.response > 2.5 ~ NA_real_ ,TRUE ~ throughput))
# #2. smooth out these by interpolation
tcdat$throughput<-na.interpolation(tcdat$throughput, option = "linear")
```

Decompose Timeseries into a weekly period, a daily period, trend and random noise

```
tcdat.ts <- forecast::msts(tcdat$throughput,seasonal.periods = c(24*60,24*60*7-1))
tcdat.ts %>% mstl(.) ->tcdat.mstl
tcdat.mstl %>% autoplot() + xlab("Week")
```



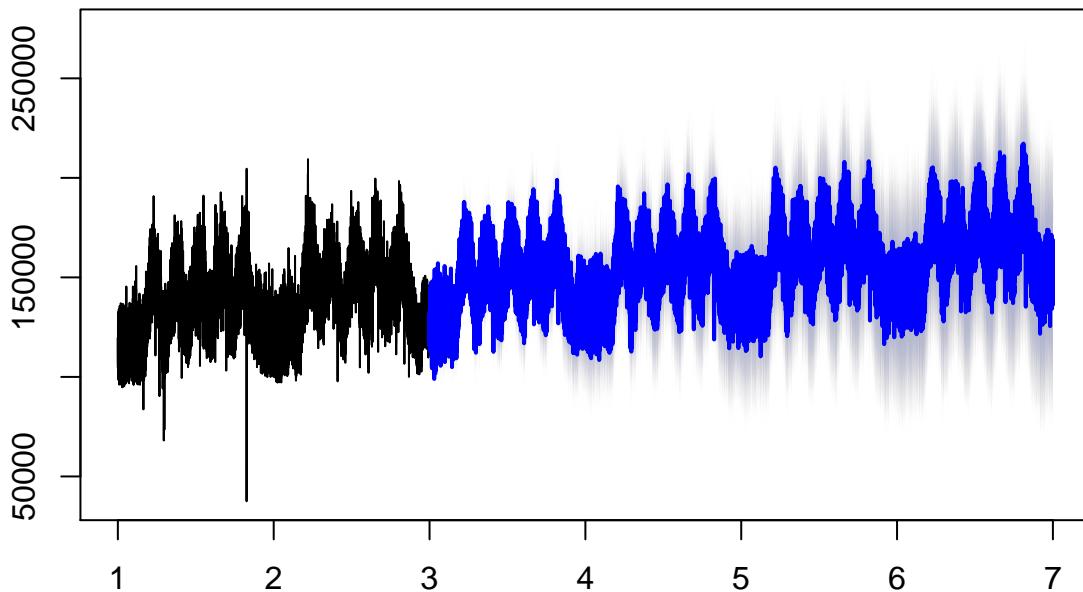
5. Forecast timeseries

Initial forecast of timeseries using STL decomposition and ARIMA for trend and noise component

```
# STL: local regression for decomposition of cyclic components
# non seasonal arima model for trend and noise component
# 2 autoregressive terms (p)
# 1 difference needed for stationarity (d)
# 1 lagged forecast errors terms (q)

#create model that describes 4 weeks of data
tcdat.ts %>%stlf(method='arima',level=c(80,95),h =24*60*7*4) ->ts.model
#use model to forecast out 4 weeks
the.forecast<-forecast(ts.model,robust=TRUE)
plot(the.forecast)
```

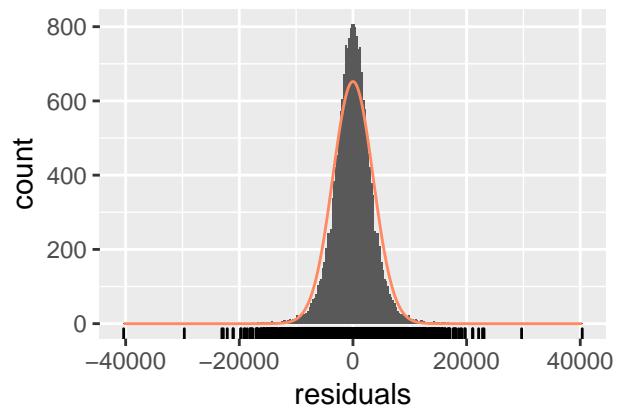
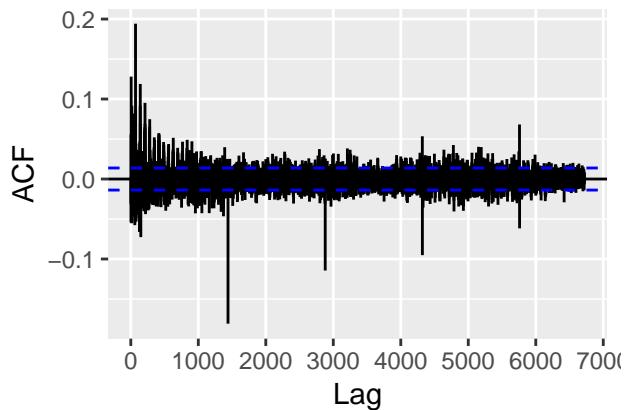
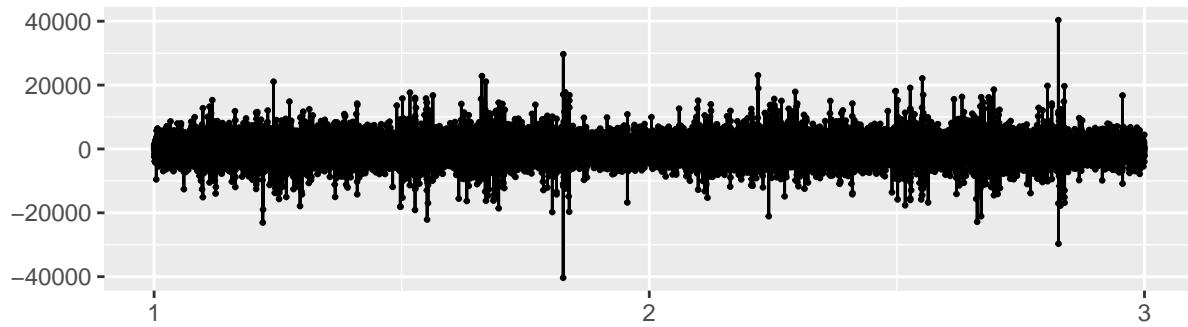
Forecasts from STL + ARIMA(2,1,1) with drift



```
checkresiduals(ts.model)
```

```
## Warning in checkresiduals(ts.model): The fitted degrees of freedom is based
## on the model used for the seasonally adjusted data.
```

Residuals from STL + ARIMA(2,1,1) with drift



```
##
```

```

## Ljung-Box test
##
## data: Residuals from STL + ARIMA(2,1,1) with drift
## Q* = 21904, df = 4028, p-value < 2.2e-16
##
## Model df: 4. Total lags used: 4032

```

create forecast dataframe

```

tcdat %>% mutate(begin_time2=begin_time+minutes(1)) %>%
filter(max(begin_time2)==begin_time2) %>% select(begin_time2) %>%
.$begin_time %>% as.character() -> start.new.time.vec
the.forecast <- data.frame(the.forecast)
forecast.date<-seq(from=ymd_hms(start.new.time.vec),length.out=dim(the.forecast)[1],by="1 min")
the.forecast %>% data.frame() %>% cbind(forecast.date) -> the.forecast

```

6. Adjust for changes in user behavior (Daylight Savings Time)

adjust for end of daylight savings time (app usage should be shifted)

```

# Sunday, November 5, 2017, 2:00:00 am clocks were turned backward 1 hour to 1:00am (US Canada)
# UTC does not change with a change of seasons
# but usage pattern of app will shift according to local time change

# for example: post time change, the 8am UTC usage will look like previous 7am UTC usage during daylight
# during daylight savings people were getting to work an hour earlier in UTC time

# ID daylight savings observations and adjust throughput predictions (fall back)
the.forecast %>% mutate(dst=case_when(forecast.date<ymd_hms("2017-11-05 01:00:00")~"yes",TRUE ~ "no"))
the.forecast %>% filter(dst=="no" & month(forecast.date)==11 & day(forecast.date)==5 & hour(forecast.date)==0)
the.forecast %>% filter(dst=="no") %>% mutate(forecast.date=ymd_hms(forecast.date,tz = "UTC")+60*60->time)
the.forecast.adj<-rbind(the.forecast %>% filter(dst=="yes"),fall.back.hr,the.forecast.dst.shift)
the.forecast.adj %>% arrange(forecast.date) %>% head

```

Point.Forecast	Lo.80	Hi.80	Lo.95	Hi.95	forecast.date	dst
132488.8	127951.3	137026.2	125549.3	139428.2	2017-10-28 17:00:00	yes
122422.9	117457.7	127388.2	114829.2	130016.6	2017-10-28 17:01:00	yes
125373.1	119887.3	130858.9	116983.3	133762.9	2017-10-28 17:02:00	yes
115980.8	110241.6	121719.9	107203.5	124758.0	2017-10-28 17:03:00	yes
127887.7	121954.7	133820.8	118813.9	136961.5	2017-10-28 17:04:00	yes
125441.0	119382.8	131499.1	116175.8	134706.1	2017-10-28 17:05:00	yes

plot adjusted forecast

```

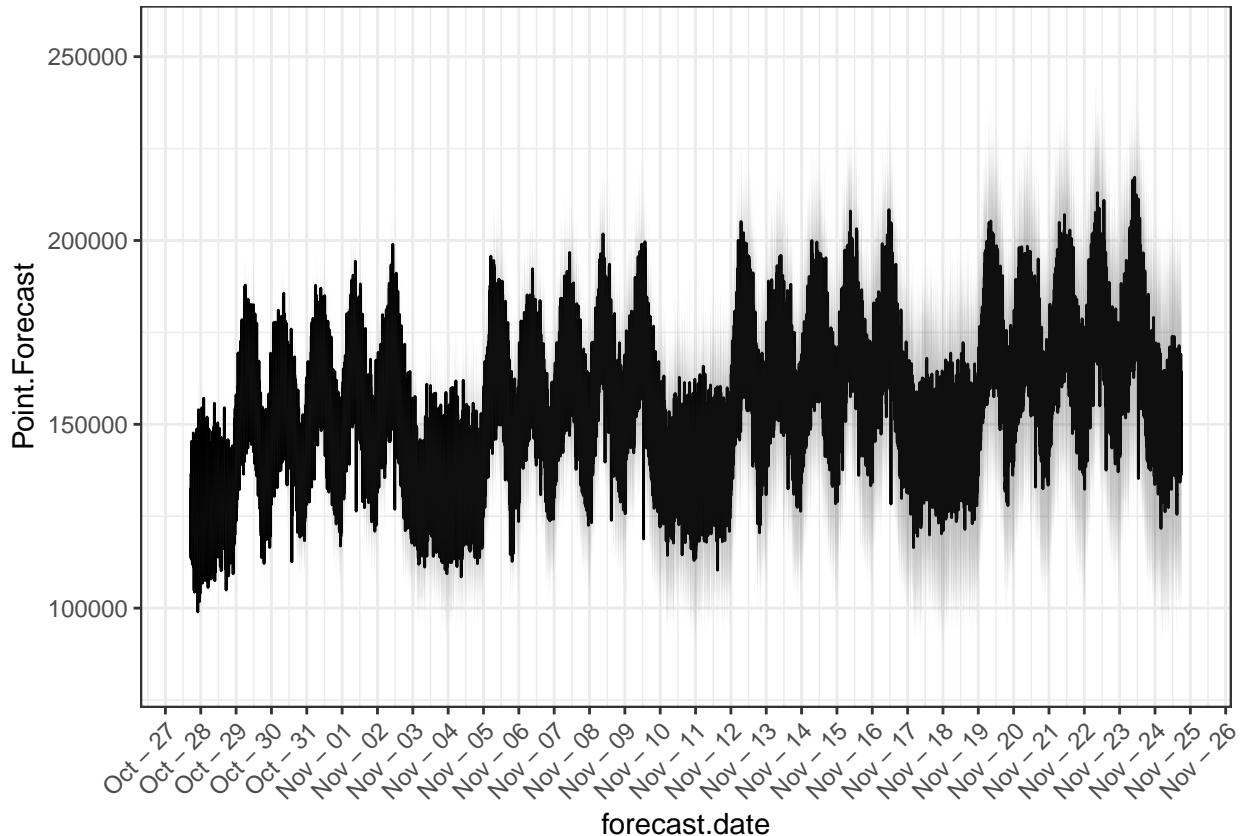
#the.forecast.adj %>% mutate(wkd=weekdays.POSIXt(forecast.date)) %>% filter(!wkd %in% #c("Saturday","Sunday"))
the.forecast.adj %>% mutate(wkd=weekdays.POSIXt(forecast.date))->the.forecast.adj.plot

```

```

p1 <- ggplot(the.forecast.adj.plot, aes(forecast.date, Point.Forecast))+
  geom_line() + geom_ribbon(data=the.forecast.adj, aes(ymin=Lo.80,ymax=Hi.80),alpha=0.3) +
  scale_x_datetime(breaks = date_breaks("1 day"), labels=date_format("%b - %d", tz=Sys.timezone())) +
  theme_bw() +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
p1

```



7. Construct A probable average_response_time timecourse for the projection

account for process 2 (throughput slowing down response time)

```

#use forecasted throughput timeline to forecast response time with fitted function describing process 1
the.forecast.adj$resp_p2_component_log<-predict(response.throughput.mod,data.frame(log.through=log10(the
forecast.adj %>%
  mutate(resp_p2_component=10^resp_p2_component_log) %>%
  select(-resp_p2_component_log) %>%
  mutate(resp_p2_component=resp_p2_component+sample(response.resid,size=length(resp_p2_component),replac
e=T)

```

response time timecourse: account for process 1 (response time decreasing throughput)

```
#randomize response times while keeping the day specific structure to create probable trajectory
set.seed(181) #121

# function to randomize slowdown events while preserving day/week structure
reshuffle.spikes<-function(shift.sec,type,the.seed){set.seed(the.seed)
  #pull of only large spikes in log response for the first week and second week
  the.week<-switch(type,week1=tcdat %>% filter(log.response>2.5 & day(begin_time) <21),
                     week2=tcdat %>% filter(log.response> 2.5 & day(begin_time) > 21))
  #split off a week into days and resample the large response times on each day
  wk.shift<-switch(type,1,0)
  the.week %>% mutate(wkd=weekdays(begin_time)) %>% group_by(wkd) %>%
    mutate(average_response_time=sample(average_response_time,size=length(average_response_time),replace=TRUE))
    mutate(begin_time=begin_time+shift.sec+60*60*24*7*(wk.shift))%>% ungroup

  week.no<-sample(c(2,2,1,1),4,replace=T)
  seed.no<-sample(c(1:50),4)
  # create 4 week forecast of probable events
  response.spikes.randomized<-rbind(reshuffle.spikes(60*60*24*7*1,week.no[1],seed.no[1]),
                                      reshuffle.spikes(60*60*24*7*2,week.no[2],seed.no[2]),
                                      reshuffle.spikes(60*60*24*7*3,week.no[3],seed.no[3]),
                                      reshuffle.spikes(60*60*24*7*4,week.no[4],seed.no[4])) %>% select(begin_time,average_response_time,seed.no)

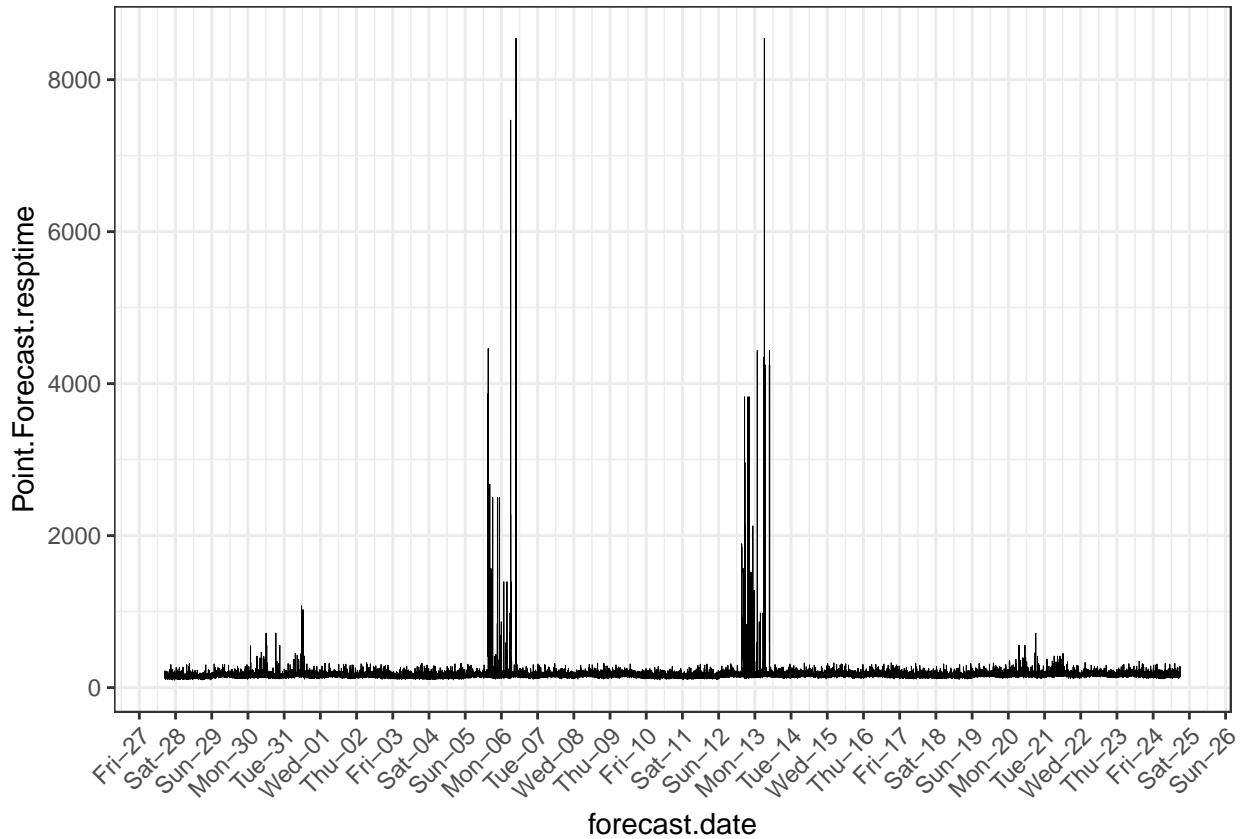
  # 2. merge strucutured events into the full forecasted data frame
  merge(response.spikes.randomized %>% mutate(forecast.date=as.character(forecast.date)),the.forecast.adj)
  mutate(forecast.date=ymd_hms(forecast.date)) -> the.forecast.adj
```

combine both process 1 and 2 to get simulated timecourse for response time

```
the.forecast.adj %<>%
  mutate(Point.Forecast.resptime=case_when(is.na(average_response_time_spike) ~ resp_p2_component, TRUE
                                             ~ average_response_time_spike))

ggplot(the.forecast.adj,aes(forecast.date,Point.Forecast.resptime))+geom_line(size=0.1)+  

  scale_x_datetime(breaks = date_breaks("1 day"), labels=date_format("%a-%d",tz=Sys.timezone()))+
  theme_bw()+
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
```

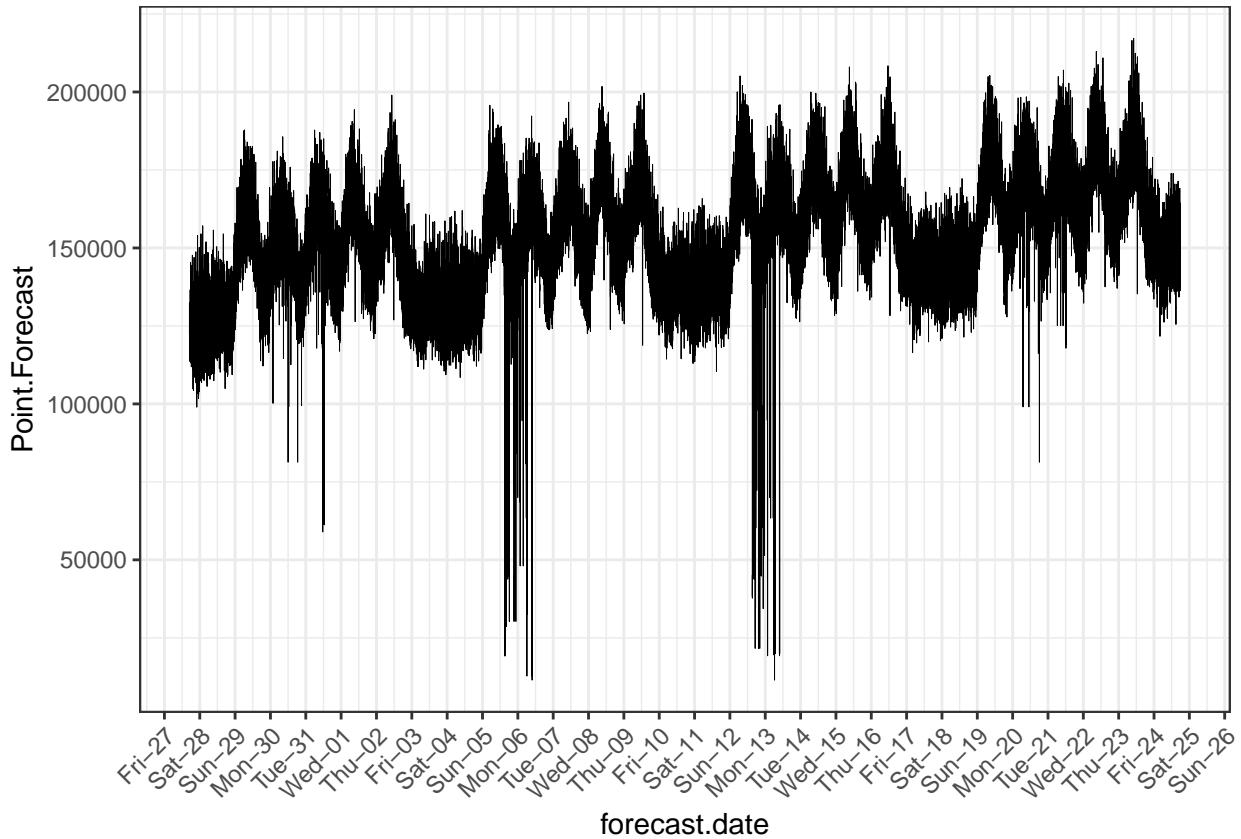


8. Account for the impact of the simulated response time slowdown events in the forecasted projection of throughput:

```
# write function to predict throughput drops from response time spikes
predict.throughput<-function(spike.dat){
  log10.through<-coefficients(res.impact.throughput.mod)[2]*log10(spike.dat)+coefficients(res.impact.throughput.mod)[1]
  10^log10.through}

# predict throughput drops
the.forecast.adj$throughput.impact<-predict.throughput(the.forecast.adj$average_response_time_spike)
# account for impact of response slowdowns events on throughputs
the.forecast.adj %>% mutate(Point.Forecast=case_when(is.na(throughput.impact) ~ Point.Forecast, TRUE ~ Point.Forecast))

ggplot(the.forecast.adj,aes(forecast.date,Point.Forecast))+geom_line(size =0.09) +
  scale_x_datetime(breaks = date_breaks("1 day"), labels=date_format("%a-%d",tz=Sys.timezone()))+
  theme_bw()+
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
```



```
# 9 create 4 week forecasted data frame in same format as original
```

```
dim(the.forecast.adj)[1] -> num.minutes.in.forecast
tcdat.raw %>% tail %>% arrange(desc(begin_time)) %>% .[1,1] %>% as.integer -> last.unix.timestamp
the.seq=seq(from=last.unix.timestamp+60,length.out=num.minutes.in.forecast,by=60)
the.forecast.adj %>% mutate(begin_time=the.seq)%>% select(begin_time,Point.Forecast.resptime,Point.Forecast.throughput) %>% rename(average_response_time=Point.Forecast.resptime,throughput=Point.Forecast) -> forecast.dat.final

#check that forecast outputs make sense
library(broom)
dim(forecast.dat.final)[1] -> num.minutes.in.forecast
print(paste0("number of days ",floor(num.minutes.in.forecast/(60*24))))
```

```
## [1] "number of days 28"
r.dat<-rbind(tcdat$average_response_time %>% summary %>% broom::tidy(.),forecast.dat.final$average_response_time %>% cbind(data=c("raw","forecast"),.) %>% cbind(variable=c("average_response_time"),.))
t.dat<-rbind(tcdat$throughput %>% summary %>% broom::tidy(.),forecast.dat.final$throughput %>% summary %>% data.frame %>% cbind(data=c("raw","forecast"),.) %>% cbind(variable=c("throughput"),.))

rbind(r.dat,t.dat)
```

variable	data	minimum	q1	median	mean	q3	maximum
average_response_time	raw	109.8967	133.4094	147.5379	157.2613	166.7865	8543.577
average_response_time	forecast	99.0734	147.4522	160.6072	168.7383	175.7324	8543.577
throughput	raw	37529.0000	121771.2500	134795.5000	135497.7106	148345.0000	209418.000
throughput	forecast	11450.5676	140107.3930	152890.2221	153085.5199	166161.6858	217195.381

```
#write data
forecast.dat.final %>% write.csv(., "forecast_data_app_4wks_johnson.csv")
```