

How to build a real-time influenza prediction model - DIY

Paul Schneider

2017-08-22

This is a fully reproducible R-tutorial on how to build a (near) real-time prediction model for influenza in Germany, based on search query data from Google, page view statistics from Wikipedia and influenza incidence data from the Robert Koch Institute. Methods are inspired by Max Kuhn's R-package caret and his book Applied Predictive Modelling. For criticism, comments and questions: schneider.paulpeter@gmail.com

1 Introduction

1.1 Digital Epidemiology

What is 'Digital Epidemiology'?

The idea behind the predictive influenza model is quite intuitive: We expect that people who get infected with influenza will tend to look up their symptoms, blablabla... Elaborate...

1.2 Milestones in Influenza prediction models

Some milestones: Google Flu Trend in 2008, Wikipedia predictions 2014, Twitter, etc...

1.3 Why Influenza?

In principle, the methods of surveillance, we use in this example, can be applied to any disease. However, some of the features of influenza make it easier to detect relevant predictors and monitor its activity. First, influenza has a strong seasonality (in most countries). Every year there is a start and an end, and during the summer months there are literally no cases of influenza. This strong annual pattern can be detected easily. In diseases, in which the incidence is more stable, there might be not enough variation to distinguish signal and noise. Second, during the flu season, between 5% and 20% of the population get sick. The online traces of rare diseases (say, achalasia) probably get lost in noise. And finally, influenza has distinct symptoms that most people will recognize as flu-symptoms, this means the symptoms are (somewhat) specific to influenza. In contrast, diseases with only subtle or very non-specific symptoms (say, fatigue) are certainly more problematic to track. Nevertheless, there is no reason to assume that other diseases can not be surveyed using similar methods. A few papers have already shown promising results in ...Ebola, TUberculosis... [ADD LINK][LINK]. We would indeed be very interested to see how predictive models perform in further diseases, like allergic rhinitis, depression or even sunburns. **If you have good data on this or know where to get it, we invite you to follow this tutorial, build your own model(s) and share your results with us and the public.**

2 Tutorial: Building ‘Influenza Nowcast’

2.1 Overview & data sources

We will build a model for predicting the influenza incidence (lab-confirmed cases) in Germany in near real-time (‘Nowcasting’) by using data from various sources. Outcome data is taken from the Robert Koch Institute, predictor data from Google Correlate, Google Trends and Wikipedia via its API and Wikishark. All analyses are performed using the statistical software R/R Studio, and plenty of its packages, particularly the caret package. Everything we use is available online for free- However, you have to have a Google account to access some of the Google Correlate data.

We will start with getting information on influenza incidence in Germany, since the availability of the outcome data determines the scope and time horizon of the model we build. We then download data on search queries and page view statistics from Google Correlate, Google Trends and Wikipedia, respectively. The data have to be re-formatted a bit to be matched with each other and the data will be pre-processed before entering into the model. Subsequently, we will use lasso regression to select relevant predictors and build a prediction model. The objective of the model will be to estimate the current influenza incidence.

We would like to point out that when we speak of ‘predictions’, we use the term in its statistical sense. We do not mean ‘forecasting’, as in forecasting future influenza cases. We just observe what people were doing on the internet until today and infer from this information how many influenza cases there are about now - We like to call that ‘Nowcasting’

If you want to transfer this approach to another setting, the only thing you need is a set of outcome data, i.e. good-quality data of the actual weekly incidence of the disease, over a sufficient time span (1 year +), in a specific country. The rest should work (more or less) automatically, unless you wish to add more data sources as predictors, which may require some more editing. Depending on the country in which you want to predict disease activity, the availability and quality of Google and Wikipedia data may differ substantially. For Wikipedia data, page view statistics can only be distinguished by language, but not by country. Language can be a good proxy in some cases (e.g. Japanese, Korean, Italian, etc.), but a pretty bad one in others (e.g. French, English, Spanish). You can look up how much of a country’s Wikipedia traffic is in a specific language, and how much Wikipedia traffic within a specific language is from a specific country. Moreover, with weekly data, you should be aware that there are divergent opinions about what is the best starting day for a week (Saturday? Sunday? Monday?).

2.1 Getting started

Before building the model, we need to install and load the required R-packages:

```
# Install and load all required packages
# This may take a while...
required_packages<-c("knitr","RCurl","ISOweek","jsonlite","ggplot2","prophet","dplyr","gtrendsR","wikidataR")

pft_packages <- function(package){
  for(i in 1:length(package)){
    if(eval(parse(text=paste("require(",package[i],")")))==0) {
      install.packages(package[i])
    }
  }
  return (eval(parse(text=paste("require(",package[i],")"))))
}

pft_packages(required_packages)

# gtrendR has to be the developer version:
# devtools::install_github("PMassicotte/gtrendsR")
# library(gtrendsR)
```

64 And we should specify a few key parameters:

```
# What is the outcome of interest? ('term' should correspond to a Wikipedia page)
term = "Influenza" # - Laboratory-confirmed cases of influenza

# For which country do we want to build the model?
country_of_interest = "DE" # Germany in ISO_3166-2 (See: https://en.wikipedia.org/wiki/ISO_3166-2)

# Which language is relevant?
language_of_interest = "de" # German in ISO_639-1 (See: https://en.wikipedia.org/wiki/List_of_ISO_639-1)

# What the relevant time horizon (i.e. the time span we have data for)
from = as.Date("2010-07-31") # Start
to = as.Date("2017-07-31") # End

# How do we split the data into training and test data?
split.at = as.Date("2016-08-01")
# --> Training data: 2010-07-31 - 2016-08-01
# --> Test Data: 2016-08-02 - 2017-07-31
```

65 2.2 Outcome data

66 We download German influenza incidence data (cases per per 10,000) from the Robert Koch Institute, from
67 August 2010 until August 2017. We uploaded a spreadsheet to github, so it is easier to access. However, you
68 can also go to Survstat and customize your data query. The data do not come in the format we need, so we
69 have to re-arrange it a bit.

```
# Load data from github repository
influenza.de = read.csv("https://raw.githubusercontent.com/projectflutrend/pft.2/master/outcome%20data/")
names(influenza.de) = c("date", "y")
head(influenza.de)
```

```
70   date      y
71 1 2009-w40 1.54 2 2009-w41 1.72 3 2009-w42 1.92 4 2009-w43 3.38 5 2009-w44 9.75 6 2009-w45 24.81
```

```
# Re-formatting 'date', and removing irrelevant data points
influenza.de$date = paste(sub("w", "W", influenza.de$date), "-1", sep="")
influenza.de$date = ISOweek2date(influenza.de$date)
influenza.de = influenza.de[influenza.de$date > as.Date("2010-07-31") & influenza.de$date <= as.Date("2017-07-31"), ]

# The RKI data does not report data during summer months, in which there are literally no influenza cases
reference = data.frame(date = seq(from=min(influenza.de$date),
                                   to=max(as.Date(influenza.de$date)),
                                   by=7))

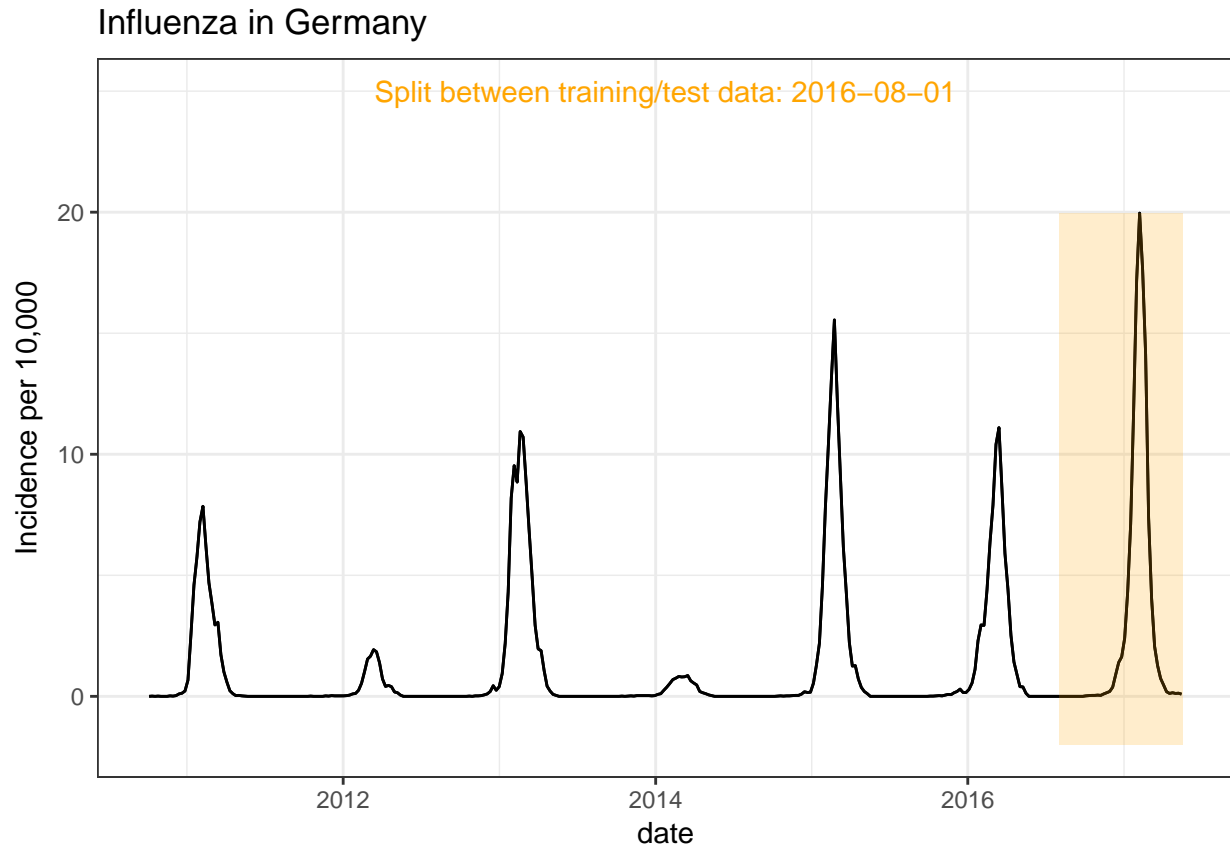
influenza.de = merge(reference, influenza.de, by="date", all=T)
influenza.de$y[is.na(influenza.de$y)] = 0

# Plotting the data
ggplot(influenza.de, aes(x=date, y=y)) +
  geom_line() +
  theme_bw() +
  geom_line() +
  annotate("rect", xmin=split.at,
          xmax=max(influenza.de$date)+2,
          ymin=min(influenza.de[,2])-2,
```

```

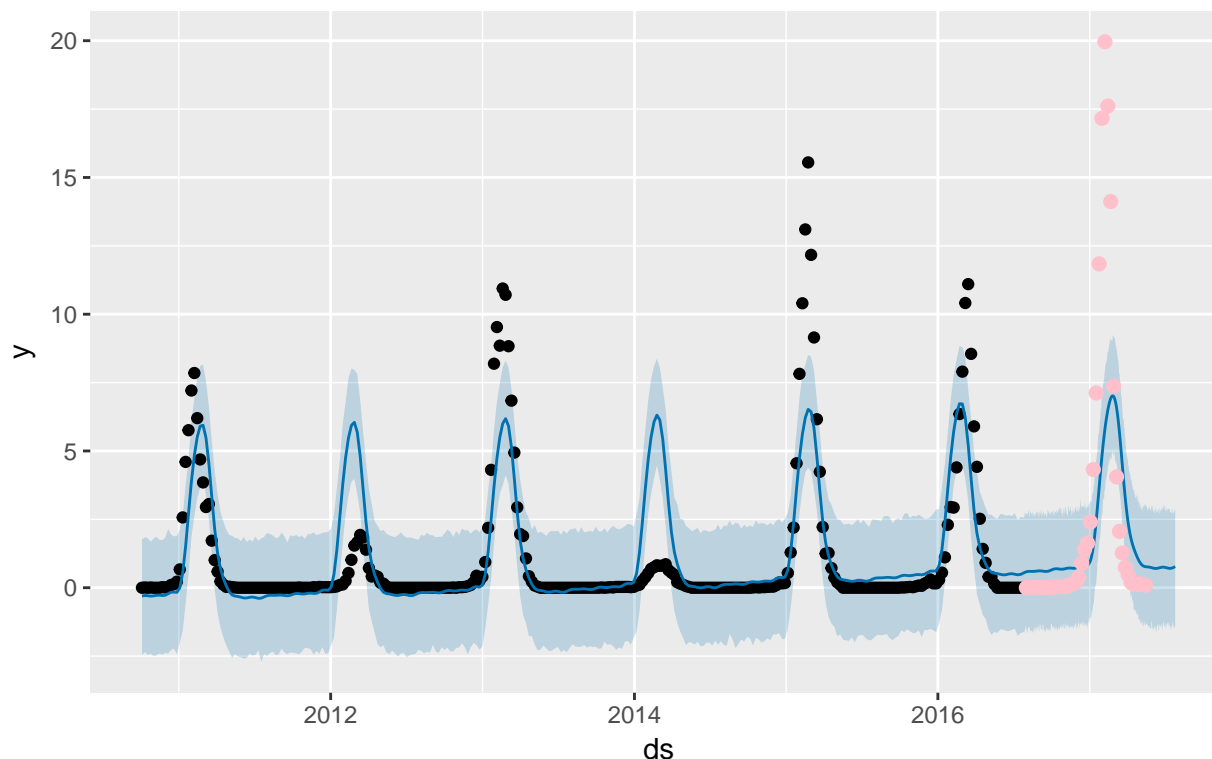
ymax=max(influenza.de$y),
fill="orange",
alpha=0.2) +
annotate("text",
x=median(influenza.de$date),
y=25,
label= paste("Split between training/test data:",split.at),
size=4,col="orange") +
ylab("Incidence per 10,000") +
ggtitle("Influenza in Germany")

```



72

73 The figure shows the influenza seasons 2010/2011 to 2016/2017, in Germany. More specifically, it shows
74 'laboratory-confirmed cases'. It should be noted that this type of outcome is probably not ideal to use as
75 a gold standard for training a digital surveillance system that is based on symptoms: If people feel like
76 having the flu, they will look up influenza on Google and Wikipedia. Flu-like symptoms that occur during
77 the summer months are, however, not caused by the influenza virus- this is why the Robert Koch Institute
78 doesn't report any data for these months and why we have to assume a flat line at zero. Our prediction
79 model has to deal with the challenge of distinguishing between people who have the real influenza and people
80 who only have influenza-like symptoms.



The prophet-forecast algorithm tries to find a pattern in the available data (2010/2011-2015/2016) and extrapolates it into the future (2016/2017). We can see that the 2016/2017 influenza season was unusually strong, compared to the last six seasons. The forecast model would have underpredicted the number of cases. Furthermore, the season started earlier than usual. So, the forecast model predicted the onset, peak and end of the season a few weeks earlier than they occurred. Later, we can use this forecast and compare it to what our Nowcast-model predicts. The metric of interest will be the Root-mean-squared-error.

2.4 Google Correlate

Google Correlate is a tool that can identify search queries that are highly correlated with any time series data you upload (Even for randomly generated numbers it will find a good match). The algorithm was also a building block of [Google Flu Trend](<https://googleblog.blogspot.co.uk/2011/05/mining-patterns-in-search-data-with.html>). See this paper for more details on how Google Trends works, and this guide on how to use it.

What you can also do is to identify queries that are correlated with another query: So, for example, find queries that were related to “Influenza” queries. But we won’t do this here.

Unfortunately, there is no Google Correlate API available for R, so you need to go to the website and upload your weekly outcome data manually. Google gives you 100 correlated search queries within the country you select (Not all countries are available) and you can download the results as a .csv file. What you need for that is a) A Google Account and b) A spreadsheet with your data in a specific format. Here, we don’t want to leak any information from the test data into the training data- So we will ask for correlated queries only for the trained data set, withholding any patterns seen in the test data set. It’s worth noting that there is a “shift series” button. It shifts your data one week in time, if you think there might be a delay between people using Google and the actual reporting of influenza cases (Only shift your data into the past, the predictors should be collected during the week the cases were reported or BEFORE the cases were reported, not afterwards).

For some reason, Google Correlate won’t provide data for the entire time span. It only gives data points until 2017-03-12 (Which is strange...). Anyway, we need to extract the keywords and download the datapoints

123 from google trends.

```
# We prepare the German Influenza data and save it in the right format
g.cor.influenza.training.upload = influenza.de[influenza.de$date<split.at,]

# Adjusting week format: making Sunday the first day of the week
g.cor.influenza.training.upload$date = g.cor.influenza.training.upload$date-1
# Saving the file in your default folder
write.table( g.cor.influenza.training.upload, col.names=FALSE,row.names = FALSE,sep=",")
```

```
124 ## 2010-10-03,0
125 ## 2010-10-10,0.01
126 ## 2010-10-17,0
127 ## 2010-10-24,0.01
128 ## 2010-10-31,0
129 ## 2010-11-07,0
130 ## 2010-11-14,0
131 ## 2010-11-21,0.02
132 ## 2010-11-28,0.01
133 ## 2010-12-05,0.03
134 ## 2010-12-12,0.1
135 ## 2010-12-19,0.13
136 ## 2010-12-26,0.22
137 ## 2011-01-02,0.67
138 ## 2011-01-09,2.57
139 ## 2011-01-16,4.6
140 ## 2011-01-23,5.76
141 ## 2011-01-30,7.21
142 ## 2011-02-06,7.85
143 ## 2011-02-13,6.2
144 ## 2011-02-20,4.69
145 ## 2011-02-27,3.85
146 ## 2011-03-06,2.95
147 ## 2011-03-13,3.06
148 ## 2011-03-20,1.72
149 ## 2011-03-27,1.01
150 ## 2011-04-03,0.59
151 ## 2011-04-10,0.23
152 ## 2011-04-17,0.11
153 ## 2011-04-24,0.03
154 ## 2011-05-01,0.04
155 ## 2011-05-08,0.02
156 ## 2011-05-15,0.01
157 ## 2011-05-22,0
158 ## 2011-05-29,0
159 ## 2011-06-05,0
160 ## 2011-06-12,0
161 ## 2011-06-19,0
162 ## 2011-06-26,0
163 ## 2011-07-03,0
164 ## 2011-07-10,0
165 ## 2011-07-17,0
166 ## 2011-07-24,0
167 ## 2011-07-31,0
168 ## 2011-08-07,0
```

169 ## 2011-08-14,0
170 ## 2011-08-21,0
171 ## 2011-08-28,0
172 ## 2011-09-04,0
173 ## 2011-09-11,0
174 ## 2011-09-18,0
175 ## 2011-09-25,0
176 ## 2011-10-02,0
177 ## 2011-10-09,0
178 ## 2011-10-16,0.01
179 ## 2011-10-23,0
180 ## 2011-10-30,0
181 ## 2011-11-06,0
182 ## 2011-11-13,0
183 ## 2011-11-20,0.02
184 ## 2011-11-27,0.01
185 ## 2011-12-04,0.03
186 ## 2011-12-11,0.02
187 ## 2011-12-18,0.02
188 ## 2011-12-25,0.02
189 ## 2012-01-01,0.02
190 ## 2012-01-08,0.03
191 ## 2012-01-15,0.05
192 ## 2012-01-22,0.1
193 ## 2012-01-29,0.12
194 ## 2012-02-05,0.25
195 ## 2012-02-12,0.54
196 ## 2012-02-19,1.02
197 ## 2012-02-26,1.54
198 ## 2012-03-04,1.66
199 ## 2012-03-11,1.93
200 ## 2012-03-18,1.84
201 ## 2012-03-25,1.39
202 ## 2012-04-01,0.72
203 ## 2012-04-08,0.41
204 ## 2012-04-15,0.45
205 ## 2012-04-22,0.4
206 ## 2012-04-29,0.18
207 ## 2012-05-06,0.15
208 ## 2012-05-13,0.05
209 ## 2012-05-20,0
210 ## 2012-05-27,0
211 ## 2012-06-03,0
212 ## 2012-06-10,0
213 ## 2012-06-17,0
214 ## 2012-06-24,0
215 ## 2012-07-01,0
216 ## 2012-07-08,0
217 ## 2012-07-15,0
218 ## 2012-07-22,0
219 ## 2012-07-29,0
220 ## 2012-08-05,0
221 ## 2012-08-12,0
222 ## 2012-08-19,0

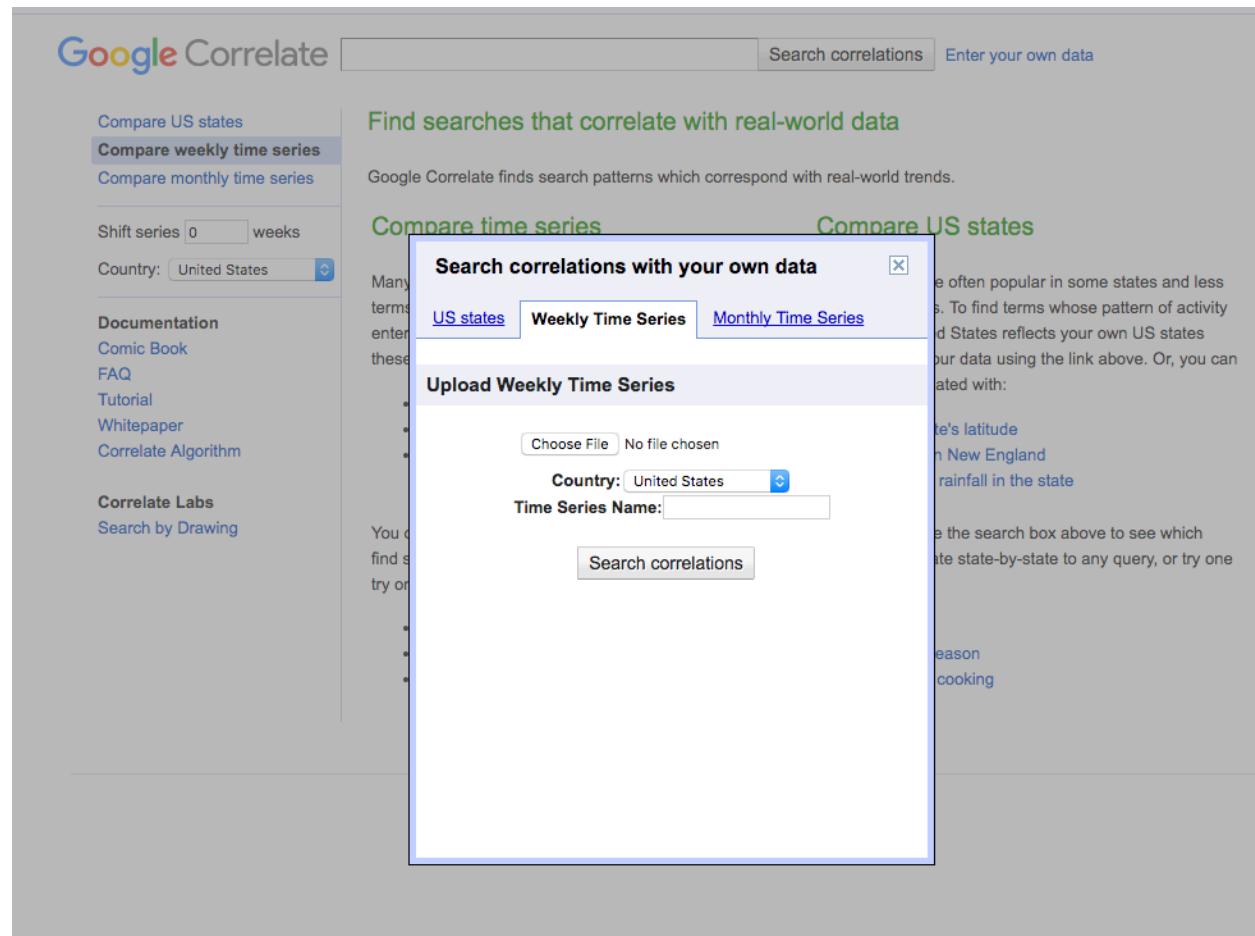
223 ## 2012-08-26,0
224 ## 2012-09-02,0
225 ## 2012-09-09,0
226 ## 2012-09-16,0
227 ## 2012-09-23,0
228 ## 2012-09-30,0
229 ## 2012-10-07,0
230 ## 2012-10-14,0
231 ## 2012-10-21,0.01
232 ## 2012-10-28,0
233 ## 2012-11-04,0.02
234 ## 2012-11-11,0.02
235 ## 2012-11-18,0.03
236 ## 2012-11-25,0.06
237 ## 2012-12-02,0.1
238 ## 2012-12-09,0.2
239 ## 2012-12-16,0.44
240 ## 2012-12-23,0.25
241 ## 2012-12-30,0.4
242 ## 2013-01-06,0.94
243 ## 2013-01-13,2.19
244 ## 2013-01-20,4.31
245 ## 2013-01-27,8.19
246 ## 2013-02-03,9.53
247 ## 2013-02-10,8.85
248 ## 2013-02-17,10.94
249 ## 2013-02-24,10.71
250 ## 2013-03-03,8.83
251 ## 2013-03-10,6.84
252 ## 2013-03-17,4.94
253 ## 2013-03-24,2.94
254 ## 2013-03-31,1.96
255 ## 2013-04-07,1.89
256 ## 2013-04-14,1.07
257 ## 2013-04-21,0.42
258 ## 2013-04-28,0.22
259 ## 2013-05-05,0.08
260 ## 2013-05-12,0.04
261 ## 2013-05-19,0
262 ## 2013-05-26,0
263 ## 2013-06-02,0
264 ## 2013-06-09,0
265 ## 2013-06-16,0
266 ## 2013-06-23,0
267 ## 2013-06-30,0
268 ## 2013-07-07,0
269 ## 2013-07-14,0
270 ## 2013-07-21,0
271 ## 2013-07-28,0
272 ## 2013-08-04,0
273 ## 2013-08-11,0
274 ## 2013-08-18,0
275 ## 2013-08-25,0
276 ## 2013-09-01,0

277 ## 2013-09-08,0
278 ## 2013-09-15,0
279 ## 2013-09-22,0
280 ## 2013-09-29,0
281 ## 2013-10-06,0
282 ## 2013-10-13,0.01
283 ## 2013-10-20,0
284 ## 2013-10-27,0.01
285 ## 2013-11-03,0.02
286 ## 2013-11-10,0.01
287 ## 2013-11-17,0.03
288 ## 2013-11-24,0.03
289 ## 2013-12-01,0.03
290 ## 2013-12-08,0.03
291 ## 2013-12-15,0.03
292 ## 2013-12-22,0.02
293 ## 2013-12-29,0.03
294 ## 2014-01-05,0.08
295 ## 2014-01-12,0.11
296 ## 2014-01-19,0.25
297 ## 2014-01-26,0.38
298 ## 2014-02-02,0.53
299 ## 2014-02-09,0.69
300 ## 2014-02-16,0.75
301 ## 2014-02-23,0.82
302 ## 2014-03-02,0.8
303 ## 2014-03-09,0.81
304 ## 2014-03-16,0.86
305 ## 2014-03-23,0.64
306 ## 2014-03-30,0.55
307 ## 2014-04-06,0.47
308 ## 2014-04-13,0.2
309 ## 2014-04-20,0.15
310 ## 2014-04-27,0.11
311 ## 2014-05-04,0.06
312 ## 2014-05-11,0.03
313 ## 2014-05-18,0
314 ## 2014-05-25,0
315 ## 2014-06-01,0
316 ## 2014-06-08,0
317 ## 2014-06-15,0
318 ## 2014-06-22,0
319 ## 2014-06-29,0
320 ## 2014-07-06,0
321 ## 2014-07-13,0
322 ## 2014-07-20,0
323 ## 2014-07-27,0
324 ## 2014-08-03,0
325 ## 2014-08-10,0
326 ## 2014-08-17,0
327 ## 2014-08-24,0
328 ## 2014-08-31,0
329 ## 2014-09-07,0
330 ## 2014-09-14,0

331 ## 2014-09-21,0
332 ## 2014-09-28,0
333 ## 2014-10-05,0
334 ## 2014-10-12,0.01
335 ## 2014-10-19,0.02
336 ## 2014-10-26,0.01
337 ## 2014-11-02,0.02
338 ## 2014-11-09,0.02
339 ## 2014-11-16,0.03
340 ## 2014-11-23,0.04
341 ## 2014-11-30,0.04
342 ## 2014-12-07,0.09
343 ## 2014-12-14,0.19
344 ## 2014-12-21,0.16
345 ## 2014-12-28,0.17
346 ## 2015-01-04,0.54
347 ## 2015-01-11,1.29
348 ## 2015-01-18,2.2
349 ## 2015-01-25,4.55
350 ## 2015-02-01,7.82
351 ## 2015-02-08,10.4
352 ## 2015-02-15,13.1
353 ## 2015-02-22,15.55
354 ## 2015-03-01,12.17
355 ## 2015-03-08,9.15
356 ## 2015-03-15,6.16
357 ## 2015-03-22,4.24
358 ## 2015-03-29,2.22
359 ## 2015-04-05,1.25
360 ## 2015-04-12,1.27
361 ## 2015-04-19,0.72
362 ## 2015-04-26,0.37
363 ## 2015-05-03,0.2
364 ## 2015-05-10,0.11
365 ## 2015-05-17,0
366 ## 2015-05-24,0
367 ## 2015-05-31,0
368 ## 2015-06-07,0
369 ## 2015-06-14,0
370 ## 2015-06-21,0
371 ## 2015-06-28,0
372 ## 2015-07-05,0
373 ## 2015-07-12,0
374 ## 2015-07-19,0
375 ## 2015-07-26,0
376 ## 2015-08-02,0
377 ## 2015-08-09,0
378 ## 2015-08-16,0
379 ## 2015-08-23,0
380 ## 2015-08-30,0
381 ## 2015-09-06,0
382 ## 2015-09-13,0
383 ## 2015-09-20,0
384 ## 2015-09-27,0

385 ## 2015-10-04,0.01
386 ## 2015-10-11,0.01
387 ## 2015-10-18,0.02
388 ## 2015-10-25,0.03
389 ## 2015-11-01,0.02
390 ## 2015-11-08,0.05
391 ## 2015-11-15,0.08
392 ## 2015-11-22,0.07
393 ## 2015-11-29,0.15
394 ## 2015-12-06,0.2
395 ## 2015-12-13,0.3
396 ## 2015-12-20,0.16
397 ## 2015-12-27,0.16
398 ## 2016-01-03,0.29
399 ## 2016-01-10,0.55
400 ## 2016-01-17,1.11
401 ## 2016-01-24,2.3
402 ## 2016-01-31,2.95
403 ## 2016-02-07,2.93
404 ## 2016-02-14,4.4
405 ## 2016-02-21,6.35
406 ## 2016-02-28,7.9
407 ## 2016-03-06,10.41
408 ## 2016-03-13,11.1
409 ## 2016-03-20,8.55
410 ## 2016-03-27,5.9
411 ## 2016-04-03,4.42
412 ## 2016-04-10,2.52
413 ## 2016-04-17,1.42
414 ## 2016-04-24,0.91
415 ## 2016-05-01,0.4
416 ## 2016-05-08,0.39
417 ## 2016-05-15,0.13
418 ## 2016-05-22,0
419 ## 2016-05-29,0
420 ## 2016-06-05,0
421 ## 2016-06-12,0
422 ## 2016-06-19,0
423 ## 2016-06-26,0
424 ## 2016-07-03,0
425 ## 2016-07-10,0
426 ## 2016-07-17,0
427 ## 2016-07-24,0

```
### --> Now, go to https://www.google.com/trends/correlate/  
### Login, upload the spreadhseet, and download results
```



```
# We have put the Google Correlate data for our data on Github:
github.url = "https://raw.githubusercontent.com/projectflutrend/pft.2/master/input%20data/correlate-g_c
# The first 10 lines are text, so we want to skip this part
g.cor.results = read.csv(skip=10,github.url)
# extracting names, except 'date' and 'y'
g.cor.keywords = names(g.cor.results)[-c(1,2)]
# R inserted "." for spaces, we have to undo this:
g.cor.keywords = gsub("\\.", " ",g.cor.keywords)

g.cor.keywords[1:20] # Showing the first 20 keywords
```

```
428 ## [1] "influenza a"           "influenza"
429 ## [3] "virusgrippe"          "j11 1"
430 ## [5] "grippe fieber"        "grippe verlauf"
431 ## [7] "influenza schnelltest" "influenza inkubationszeit"
432 ## [9] "echte grippe"         "grippe bei kindern"
433 ## [11] "grippe husten"        "grippe influenza"
```

```

434 ## [13] "ist grippe ansteckend"      "influenza grippe"
435 ## [15] "symptome grippe"           "j11 1 g"
436 ## [17] "grippe wie lange"          "symptome influenza"
437 ## [19] "wie lange dauert eine grippe" "verlauf grippe"

```

438 *If you wonder what “J11 1” might be - It’s the ICD-10 code for Influenza. Doctors print it on the letters for*
439 *sick leaves. It seems as if some patients were curious what the doctor had diagnosed them with. Other terms*
440 *refer to [5] fever, [10] Flu in Children, [11] influenza coughing or [17] Influenza how long*

441 2.5 Google trends

442 Google Trends is “a public web facility of Google Inc., based on Google Search, that shows how often a
443 particular search-term is entered relative to the total search-volume across various regions of the world, and
444 in various languages” - (from Wikipedia). It can also be used to identify search queries that are related
445 (“People who searched for this, also searched for this” - notice the difference between Google Trends and
446 Correlate!)

447 First, we are going to download the data for the 100 keywords we have identified using Google Correlate.
448 Moreover, while Google Correlate will most often provide you with good predictors, there might be additional
449 value in taking into account predictors from Google Trends as well, even if only to control for confounding or
450 media-generated over-prediction. Thus, we will download search query statistics for ‘influenza’, and for 15
451 related queries, as well as ‘influenza’-queries in Google News. If you like, you can expand your requests and
452 increase the numbers of predictors by also retrieving related keywords from related keyword, feed related
453 keywords into Google Correlate, etc.

```

# We have already set these parameters:
# term = "influenza"; country_of_interest="DE"; from="2010-07-31"; to="2017-07-31"; status= 1

# Identifying 15 related search queries
google_primer = gtrends(keyword=term,
                        geo=country_of_interest,
                        time=paste(from,to),
                        gprop ="web")

google_related = google_primer$related_queries$value[google_primer$related_queries$related_queries=="
g.trends.keywords = c(term,google_related)
print(g.trends.keywords)

```

```

454 ## [1] "Influenza"      "grippe influenza"
455 ## [3] "virus influenza" "grippe"
456 ## [5] "rki influenza"  "rki"
457 ## [7] "symptome influenza" "influenza impfung"
458 ## [9] "symptome"       "inkubationszeit"
459 ## [11] "influenza inkubationszeit" "hämophilus influenza"
460 ## [13] "influenza 2017" "influenza 2016"
461 ## [15] "influenza 2015" "fieber"
462 ## [17] "influenza ansteckend" "haemophilus influenza"
463 ## [19] "influenza impfung pferd" "influenza 2013"
464 ## [21] "influenza dauer" "influenza wie lange ansteckend"
465 ## [23] "influenza impfstoff" "influenza deutschland"
466 ## [25] "influenza viren" "was ist influenza"

```

467 **Noteworthy:** 5-year time span restriction: Google Correlate only gives you weekly data for time spans <
468 5 years. As a possible work-around, you could split your time frame into smaller chunks and download data
469 separately. But because of the way the data is generated, it does not match 100%, even after adjusting

470 the scale. Also, it is not clear whether results will be reproducible, as Google might show slightly different
471 figures in new queries.

472 The loops to download the search queries are messy and rather long. We will just load them from github,
473 where you can have a closer look at them, if you want.

474 **NOTE: FUNCTIONS DO ONLY WORK FOR THE TIME SPAN SET IN THIS EXMAPLE,**
475 **NEEDS TO BE REVISED TRAINING TESTING SPLIT NECESSARY???**

```
# Loading functions from github:
google.function <- getURL("https://raw.githubusercontent.com/projectflutrend/pft.2/master/quickfunction.R")
eval(parse(text = google.function))

# Functions loaded are:
# pft_ask_google(keyword,country_of_interest="DE",from="2010-07-31",to="2017-07-31",status= 1,prefix="g")
# To split the time span and download and merge query statistics from Google

# Rescale.gtrends(df.t1,df.t2)
# In order to rescale, we look at the overlapping time span and try to find the best mutliplicative sca

# Download query statistics for
# a) google.correlate queries:
g.cor.input = pft_ask_google(g.cor.keywords[1:3],country_of_interest="DE",from="2010-07-31",to="2017-07-31",status=1)

476 ## asking Google for statistics for influenza a - 33.3 %
477 ## asking Google for statistics for influenza - 66.7 %
478 ## asking Google for statistics for virusgrippe - 100 %

# b) trends queries:
g.trends.input = pft_ask_google(g.trends.keywords[1:3],country_of_interest="DE",from="2010-07-31",to="2017-07-31",status=1)

479 ## asking Google for statistics for Influenza - 33.3 %
480 ## asking Google for statistics for grippe influenza - 66.7 %
481 ## asking Google for statistics for virus influenza - 100 %

# c) news on influenza as a potentially relevant (negative) predictor
g.news.input = pft_ask_google("influenza",country_of_interest="DE",from="2010-07-31",to="2017-07-31",status=1)

482 ## asking Google for statistics for influenza - 100 %

google.input.data = merge(g.cor.input,g.trends.input,by="date",all=T)
google.input.data = merge(google.input.data,g.news.input,by="date",all=T)

dim(google.input.data)

483 ## [1] 366 8

head(google.input.data[,1:5])

484 ##      date g.cor.influenza.a g.cor.influenza g.cor.virusgrippe
485 ## 1 2010-W30                2.9                5                5.1
486 ## 2 2010-W31                5.0                4                5.1
487 ## 3 2010-W32                2.9                6                5.1
488 ## 4 2010-W33                3.6                6                7.2
489 ## 5 2010-W34                5.7                7                7.2
490 ## 6 2010-W35                5.7                8                4.1
491 ##      g.trends.Influenza
492 ## 1                      5
```

```

493 ## 2          4
494 ## 3          6
495 ## 4          6
496 ## 5          7
497 ## 6          8

```

498 *If you want to reproduce this example and avoid long downloading times, you can download the data [\[here\]](#)*
499 ***FIX LINK***

500 2.6 Wikipedia data

501 McIver & Brownstein made a strong case for the utility of Wikipedia page view data in influenza prediction
502 models. They showed that “Wikipedia usage accurately estimated the week of peak ILI activity 17% more
503 often than Google Flu Trends data and was often more accurate in its measure of ILI intensity”. Combining
504 Google data with Wikipedia may have further advantages.

505 We can access Wikipedia page view statistics via it’s API (for data from October 2015 until today) and
506 Wikishark (for data from January 2008 until December 2016)

- 507 • **On a side note:** On Wikipedoa, pages with similar content are linked with each others across different
508 languages. The Wikipediatrend package offers a convenient way to identify these corresponding pages:
509 For ‘Influenza’, there are about 80 similar pages. While finding the German word for influenza (it’s
510 ‘influenza’), is not too impressive, finding Wikipedia disease pages in Arabic and Japanese can be quite
511 useful. Btw, Wikipedia refers to languages in ISO 639-1 code.

```

# The 'Influenza' page in other languages
wikipediatrend::wp_linked_pages( page= "Influenza",lang="en")[1:10,2:3]

```

```

512 ##      lang      title
513 ## 1   af        Griep
514 ## 2   als       Influenza
515 ## 3   ar
516 ## 4   an        Gripe
517 ## 5   roa-rup   Aremi
518 ## 6   as
519 ## 7   ast       Gripe
520 ## 8   gn        Mba'asyparar ...
521 ## 9   ay        Jurma_usu
522 ## 10  az        Qrip

```

523 To download Wikipedia page view data, we have to turn to two sifferent sources of data. Wikipedia has
524 its own API, which allows fast and convenient open access. However, the API is only able to retrieve data
525 for > October 2015. Statistics for prior dates are stored in huge files, showing page views per hour in every
526 language. Fortunately, there is a private project, called wikishark, which we can quickly use to download
527 relevant data per week (There used to be another server with an API: stats.grok.se, but it appears to be
528 down since July 2017). For the two time periods, Wikipedia has used different metrics to count page views
529 and to detmerine individual page visitors. Thus, there might be some discrepancies. It is also important to
530 note that Wikipedia is not static. New pages are created and old pages may be changed or merged, i.e.for
531 some pages that do exist today, we won’t find any meaningful page view statistics in 2010 - So, expect (and
532 ignore) some errors when downloading the data.

533 We start with downloading page view statistics for the main influenza page. Then, we extract all links on
534 this page which refer to other wikipedia pages (On the Influenza Wikipedia page there are links to pages
535 baout Aspirin, bronchitis, Allergy etc.). In addition, we can extract statistics for all pages which link to
536 the influenza page (Most mention and link to influenza only en passant). Of you like, you can also add
537 Wikipedia pages manually, if you think people who get flu will look them up (or people who don’t have the
538 flu, but visit the infuenza page, will look them up).

539 Again, the functions to identify linked pages and to download the statistics are rather long and messy. We
540 cut it short and recommend just loading the functions from github, without going through the code. If you
541 want to have a look at it, click here.

542 **NOTE: FUNCTIONS DO ONLY WORK FOR THE TIME SPAN SET IN THIS EXMAPLE,**
543 **NEEDS TO BE REVISED**

```
# Loading functions from github:
wiki.functions <- getURL("https://raw.githubusercontent.com/projectflutrend/pft.2/master/quickfunctions.R")
eval(parse(text = wiki.functions))

# pft_wiki_lp(term = "Influenza", language_of_interest = "de", backlinked = 1 ,manual.pages=c("Halsschmerzen", "Hausmittel"))

# pft_ask_wikipedia(pages = wiki.pages, language_of_interest = "de", from = as.Date("2010-01-01"), to = as.Date("2017-07-31"), status = 1)

# Retrive potentially relevant Wikipedia pages
wiki.pages = pft_wiki_lp(term = "Influenza",
                        language_of_interest = "de",
                        backlinked = 1,
                        manual.pages=c("Halsschmerzen", "Hausmittel"))

str(wiki.pages)
```

```
544 ## chr [1:604] "Influenza" "Acetylsalicylsäure" ...

# We have identified the main influenza page, 149 pages that are links on the influenza page, and 451 pages that are links on the main influenza page.

# Now, we can download their page view statistics (Depending on the amount of relevant pages, this may take a while)
wikipedia.input.data = pft_ask_wikipedia(pages = wiki.pages[1:3],
                                        language_of_interest = "de",
                                        from = as.Date("2010-01-01"),
                                        to = as.Date("2017-07-31"),
                                        status = 1)
```

```
545 ## Downloading data for page 1 of 3 - 33.33333 %
546 ## Downloading data for page 2 of 3 - 66.66667 %
547 ## Downloading data for page 3 of 3 - 100 %
```

```
head(wikipedia.input.data[,1:4])
```

```
548 ##      date wiki.Influenza wiki.Acetylsalicylsäure
549 ## 1 2009-W53             NA                      NA
550 ## 2 2010-W01      -1.1099757                11.40327
551 ## 3 2010-W02      -0.5476942                17.00074
552 ## 4 2010-W03      -1.3283038                11.03771
553 ## 5 2010-W04      -1.3317693                12.65428
554 ## 6 2010-W05      -1.3456314                17.82682
555 ##      wiki.Adolf_Mayer_.Agronom.
556 ## 1                      NA
557 ## 2      -0.1218000
558 ## 3      4.1631075
559 ## 4     -2.8485592
560 ## 5      0.6572741
561 ## 6      1.0468112
```

562 If you want to reproduce this example and avoid long downloading times, you can download the data [here]
563 **FIX LINK**

2.4.4 Putting the data together and pre-processing

Now, that we have all the data we wanted, we merge the files into one dataframe, and split it again into predictors and outcome data, as well as training and testing/evaluation data sets.

Preprocessing: Scaling, Centering. Also, many predictors are highly correlated. Especially in the Google Correlate data set (see figure below)

```
# Combining outcome, wikipedia, google trends and google correlate
influenza.de$date = ISOweek(influenza.de$date )

# Merging by week (avoiding any Monday/Sunday or other day issues)
df.full = merge(influenza.de,google.input.data, by="date")
df.full = merge(df.full,wikipedia.input.data, by="date")

# Setting date back to a date
df.full$date = ISOweek2date(paste(df.full$date,"-1",sep="")) #
dim(df.full)
```

```
## [1] 346 12
```

```
# save(df.full, file="/users/waqr/desktop/df.full.de.2.rdata")
```

CHECK DIMS! 337 rows (=weeks),548 columns (date,outcome,546 potential predictors)

Splitting the data and pre-processing

First, we split the data into outcome/predictors, as well as training/testing data sets. From now on, we try not to leak any kind of information from the test data set to the training data set.

```
split = which(df.full$date<split.at)

df.train = df.full[split,-c(1,2)] # Predictor training data set
y.train = df.full[split,c(2)] # Outcome for training data set
date.train = df.full[split,c(1)] # Date, not a predictor but useful for plotting

df.test = df.full[-split,-c(1,2)] # Predictors for testing/evaluation data set
y.test = df.full[-split,c(2)] # Outcome for testing data set
date.test = df.full[-split,c(1)] # date for test data set
```

Second, we will scale and center all predictors - Google data comes scaled and centered already, so in this example, we will do this with Wikipedia data. This improves the performance for some of the modelling techniques and it makes it easier to compare predictors. However, we lose information about how many visits in absolute numbers a Wikipedia page had.

Third, we remove predictors that have too many (Say, >10%) missing values (See ‘wiki.Geflügel.Aufstallungsverordnung’ = translation?!) in figure below). This method we will also apply to the test data set. If we wouldn’t, we could end up with test-predictors that have too many, or only missing values. Then, we would need to re-run the whole model again without the missing test-predictor. So, we want to make sure that sufficient cases are available in both, training and test data set. When less than 10% of cases are missing, we are going to impute the missing values by “k-nearest neighbor imputation”[Add a link, description], for the training and test data set separately.

Fourth, we also want to remove predictors that have near zero variance. This means they have few unique values. Take for example the Wikipedia page on “Samuel Warren Abott”: The page was created only in 2015, so in the time before, the page had mostly 0 page views (Only rarely someone tries to access a Wikipedia page that is non-existent). In addition, this page also has some missing values. Prediction models can be adversely affected by these near zero variance predictors, and they don’t really add any value to the models.

590 You could also consider removing predictors that have a very low activity, because results could be unstable
591 if models put weight in them.

592 ***NOT RUNNING THE FOLLOWING CODE IN SHORTEND***
593 ***VERSION OF THE CODE***

```
# Removing features with >10% NAs
# in training
sum.NA.train = as.numeric(lapply(df.train,function(x){sum(is.na(x))}))
sum.NA.train = sum.NA.train > length(df.train[,1]) * 0.1
if(sum(sum.NA.train)>0){
df.train = df.train[-which(sum.NA.train)]
df.test = df.test[which(colnames(df.test) %in% colnames(df.train))]}
# and test data separately
sum.NA.test = as.numeric(lapply(df.test,function(x){sum(is.na(x))}))
sum.NA.test = sum.NA.test > length(df.test[,1]) * 0.1
if(sum(sum.NA.test)>0){
df.test = df.test[-which(sum.NA.test)]
df.train = df.train[which(colnames(df.train) %in% colnames(df.test))]}

# Removing features with near zero variance
# identify near zero-variance predictors [only in df.train!]
nearZeroVar = nearZeroVar(df.train,freqCut = 95/5 , uniqueCut = 25)
if(sum(nearZeroVar)>0){
df.train = df.train[,-nearZeroVar]
df.test = df.test[which(colnames(df.test) %in% colnames(df.train))]}

# Scaling, centering, and imputing remaining NAs
preprocess.df.train = preProcess(df.train, method=c("scale","center","knnImpute")) # why knnimpute error

df.train = predict(preprocess.df.train, newdata = df.train)
df.test = predict(preprocess.df.train,newdata = df.test)
```

594 ***NOT RUNNING THE FOLLOWING CODE IN SHORTEND***
595 ***VERSION OF THE CODE***

596 *We can also consider principal component analysis, assume methods might work better when we reduce multi-*
597 *colinerity. Especially Google Correlate data is highly correlated -> We can plot an example with just 10*
598 *predictors, but google correlate gives us 100 of these kind!*

599 2.5 Building a prediction model

600 For the purpose of this tutorial, we will only use lasso-regression to build and train our model. It's
601 a technique that automatically selects and regularizes predictors - See Tibshirani for a technical or this
602 website and video for a more practical explanation. For a real application, we can try and evaluate many
603 other methods, e.g. elastic net, multivariate adaptive regression splines, rule-based regression and other,
604 available in the caret package. The predictive accuracy of the different models can be compared in terms

605 of root-mean-squared-error, using cross-validation (see below). We can then choose whatever model is most
606 accurate, regardless of its form, complexity or the amount or types of predictors it uses. It turns out, though,
607 that lasso regression seems to be always among the top performing methods in the type of model we are
608 building. Also, compared to many other techniques, it is computationally very fast.

609 Cross-validation

610 Cross validation, rolling forward

- 611 • explain what cv is
- 612 • explain why split-cv doesn't work
- 613 • explain what rolling forward cv is
- 614 • Talk about an appropriate horizon
- 615 – ideal: cv and 1-week forward rolling evaluation plot!?
- 616 • what's a grid?

```
# control object for cross validation: rolling forward cv with fixed origin
controlObject <- trainControl(method = "timeslice",
                              initialWindow = 52,
                              horizon = 52, # 1-52 # Smaller horizon better
                              fixedWindow = FALSE,
                              allowParallel = TRUE)

# setting grids
glmnetGrid <- expand.grid(.alpha = c(0, .1, .2, .4, .6, .8, 1), .lambda = seq(.01, 1, length = 40))
# cubistGrid <- expand.grid(.committees = c(1, 5, 10, 50, 75, 100), .neighbors=c(0,1,3,5,7,9))

# parallel computing
no_cores <- detectCores() - 1
cl <- makeCluster(no_cores, type="FORK")
registerDoParallel(cl)
```

617 Model Specifications

```
# formula list
formula.list = list(glmnet.mod = list(y= y.train ,
                                       x = df.train,
                                       method = "glmnet",
                                       family = "gaussian",
                                       tuneGrid = glmnetGrid,
                                       trControl = controlObject))
```

618

```
# A function to evaluate the models
eval.function <- getURL("https://raw.githubusercontent.com/projectflutrend/pft.2/master/quickfunctions/")
eval(parse(text = eval.function))

if(1==0){
```

```

# A loop to build and evaluate the model
models.de = list(result.list = list(), eval.list = list())
for(i in 1:length(formula.list)){
  cat("Building a model:",formula.list[[i]]$method,"\n")
  tryCatch({
    models.de$result.list[[i]] = do.call("train",formula.list[[i]])
    names(models.de$result.list)[i] = names(formula.list)[i]

    models.de$eval.list[[i]] = pft_eval_model(models.de$result.list[[i]])
    names(models.de$eval.list)[i] = names(formula.list)[i]

    cat(formula.list[[i]]$method,"evaluation done! \n"),
    error=function(e) {cat("error in",formula.list[[i]]$method,"\n")})
  }

}

```

619 Evaluate the models

```

#### evaluation plots
#models.de$eval.list[[1]]$plots$pred.plot

#models.de$eval.list[[2]]$plots$pred.plot

```

620 **Concluding remarks**