

An Example of code written in the past - this is all one function.

Dean Gladish
CS 9537-00
Prof. Amy Greenman-Dolan

function for Different Models <- function(modname1, modname2, groupname, foldname1, foldname2, realtimecomponent, kdefiles) {

We take each of the files for the KDE model
kdefiles <- lapply Sys.glob(paste(realtimecomponent,dispath, foldname1, "/*.csv", sep = ""))

read.csv)

kdefiles <- lapply(Sys.glob(paste(realtimecomponent,dispath, foldname2, "/*.csv", sep = "")), read.csv)

First we should look at the values: #####
summary(kdefiles[[1]]\$V alue)
It seems like the KDE model has quite a larger maximum.
We also want to create four copies of each file, each of which
it will be refitted to focus on a specific week (1, 2, 3, 4)

kdeWeek1 <- kdefiles
kdeWeek2 <- kdefiles
kdeWeek3 <- kdefiles
kdeWeek4 <- kdefiles
kdeWeek1pointPredictionfort1 <- kdefiles
kdeWeek2pointPredictionfort1 <- kdefiles
kdeWeek3pointPredictionfort1 <- kdefiles
kdeWeek4pointPredictionfort1 <- kdefiles

for (i in 1:length(kdefiles)) {
We have taken each file for each week of the KDE model
and would like to look at the predictions for each of
1, 2, 3, and 4 weeks ahead.
kdeWeek1[i] <- USNationalXWeeksAhead(kdefiles[i])
kdeWeek2[i] <- USNationalXWeeksAhead(kdefiles[i])
kdeWeek3[i] <- USNationalXWeeksAhead(kdefiles[i])
kdeWeek4[i] <- USNationalXWeeksAhead(kdefiles[i])

kdeWeek1pointPredictionfort1[i] <- USNationalXWeeksAhead(kdefiles[i])
kdeWeek2pointPredictionfort1[i] <- USNationalXWeeksAhead(kdefiles[i])
kdeWeek3pointPredictionfort1[i] <- USNationalXWeeksAhead(kdefiles[i])
kdeWeek4pointPredictionfort1[i] <- USNationalXWeeksAhead(kdefiles[i])
"week one", "yes"
"week two", "yes"
"week three", "yes"
"week four", "yes"

for (i in 1:length(kdefiles)) {
We should also do this for the KDE model.
kdeWeek1[i] <- USNationalXWeeksAhead(kdefiles[i])
kdeWeek2[i] <- USNationalXWeeksAhead(kdefiles[i])
kdeWeek3[i] <- USNationalXWeeksAhead(kdefiles[i])
kdeWeek4[i] <- USNationalXWeeksAhead(kdefiles[i])
kdeWeek1pointPredictionfort1[i] <- USNationalXWeeksAhead(kdefiles[i])
kdeWeek2pointPredictionfort1[i] <- USNationalXWeeksAhead(kdefiles[i])
kdeWeek3pointPredictionfort1[i] <- USNationalXWeeksAhead(kdefiles[i])
kdeWeek4pointPredictionfort1[i] <- USNationalXWeeksAhead(kdefiles[i])
"week one", "no"
"week two", "no"
"week three", "no"
"week four", "no"

The operating problem is that all the helper functions are contained within this main function, the name "function for Different Models" is about the best descriptive that you could imagine, conveying neither its intended purpose nor the actual behavior, so everything has its unintended side effect and unwanted behavior.

The primary (functional) problem I ran into was the fact that you have to put the KDE model (this means kernel density estimation) into the second argument for modelname because otherwise you'll run into errors. Altogether the order of the

model names matters just makes this function extremely specific to this dataset. I should also talk about the structure of this function in order to better explain how it works; so, the first three arguments are there simply to encode the official names of the models and the name of the (second) group which was involved in creating them.

What the function does is a bit offhanded. It was originally designed for finding the difference, using a t-s test, between influenza - prediction model is developed by the Reach Lab and others for the CDCs FluSight initiative/contest.


```

kdeWeek2$pointPredictionforLL[[1]] <- USNationalXWeeksAhead(kdeFiles[[1]),
"week two", "yes")
kdeWeek3$pointPredictionforLL[[1]] <- USNationalXWeeksAhead(kdeFiles[[1]),
"week three", "yes")
kdeWeek4$pointPredictionforLL[[1]] <- USNationalXWeeksAhead(kdeFiles[[1]),
"week four", "yes")

```

```

dim(nrow(nrow(m1)/length(kdeWeek1))
removeStuff <- function(lis, iskdeest) {
  # Removes elements 19, 20, 21, and 22
  if (iskdeest == "iskdeest") {
    return(lis[removed(isr, c(19, 20, 21, 22))])
  }
  if (iskdeest == "smokedtest") {
    return(lis)
  }
}

```

```

# Since R works this way we can just do it all at once
}
kdeWeek1$reduced <- removeStuff(kdeWeek1, iskdeest)
kdeWeek2$reduced <- removeStuff(kdeWeek2, iskdeest)
kdeWeek3$reduced <- removeStuff(kdeWeek3, iskdeest)
kdeWeek4$reduced <- removeStuff(kdeWeek4, iskdeest)
kdeWeek1$pointPredictionforLL$reduced <- removeStuff(kdeWeek1$pointPredictionforLL, iskdeest)
kdeWeek2$pointPredictionforLL$reduced <- removeStuff(kdeWeek2$pointPredictionforLL, iskdeest)
kdeWeek3$pointPredictionforLL$reduced <- removeStuff(kdeWeek3$pointPredictionforLL, iskdeest)
kdeWeek4$pointPredictionforLL$reduced <- removeStuff(kdeWeek4$pointPredictionforLL, iskdeest)
# Here, all but the rows we want to look at are removed.

```

There is one crucial thing that we need to do as well before generating the plots. Because our files only share their respective data until week 18 of #2018 and then refer back to the last few weeks of 2017, we need to rotate the order of our files within each of our 8 files.

```

# We will just do another for-loop to replace these values and
# achieve what is functionally a rotation.
temp1 <- kdeWeek1
temp2 <- kdeWeek2
temp3 <- kdeWeek3
temp4 <- kdeWeek4
temp5 <- kdeWeek1$reduced
temp6 <- kdeWeek2$reduced
temp7 <- kdeWeek3$reduced
temp8 <- kdeWeek4$reduced

```

```

# We also need to do this for our LL values
temp9 <- kdeWeek1$pointPredictionforLL
temp10 <- kdeWeek2$pointPredictionforLL
temp11 <- kdeWeek3$pointPredictionforLL
temp12 <- kdeWeek4$pointPredictionforLL
temp13 <- kdeWeek1$pointPredictionforLL$reduced
temp14 <- kdeWeek2$pointPredictionforLL$reduced
temp15 <- kdeWeek3$pointPredictionforLL$reduced
temp16 <- kdeWeek4$pointPredictionforLL$reduced

```

```

for (i in 1:28) {
  # The modular portion (taking the remainder)
  # ensures that we do not go outside the bounds of
  # possible indices.
  kdeWeek1[[i]] <- temp1[(i+17)%28+1]
  kdeWeek2[[i]] <- temp2[(i+17)%28+1]
  kdeWeek3[[i]] <- temp3[(i+17)%28+1]
  kdeWeek4[[i]] <- temp4[(i+17)%28+1]
  kdeWeek1$reduced[[i]] <- temp5[(i+17)%28+1]
  kdeWeek2$reduced[[i]] <- temp6[(i+17)%28+1]
  kdeWeek3$reduced[[i]] <- temp7[(i+17)%28+1]
  kdeWeek4$reduced[[i]] <- temp8[(i+17)%28+1]
}

```

For example this uses a broken for loop or ref the argument passed in is a certain type of model, leading to the aforementioned error message if the KDE model is not based on the second model name (modelname2).

The commenting can't just as easily have been black form instead of using individual hashtags for each case; also, the function design bunch of unnecessary stuff which is made conditional as shown above, furthermore it includes a lot of copy-pasted code and repeated, vaguely-named operations which could just as easily have been done using more concise, dplyr/plyr-based methods.

```

kodeWeek2Reduced[[i]] <- temp2[(i+17)%28)+1])
kodeWeek3Reduced[[i]] <- temp3[(i+17)%28)+1])
kodeWeek4Reduced[[i]] <- temp4[(i+17)%28)+1])

kodeWeek1pointPredictionfor1L[[i]] <- temp1[(i+17)%28)+1])
kodeWeek2pointPredictionfor1L[[i]] <- temp2[(i+17)%28)+1])
kodeWeek3pointPredictionfor1L[[i]] <- temp3[(i+17)%28)+1])
kodeWeek4pointPredictionfor1L[[i]] <- temp4[(i+17)%28)+1])

kodeWeek1pointPredictionfor1LIRReduced[[i]] <- temp13[(i+17)%28)+1])
kodeWeek2pointPredictionfor1LIRReduced[[i]] <- temp14[(i+17)%28)+1])
kodeWeek3pointPredictionfor1LIRReduced[[i]] <- temp15[(i+17)%28)+1])
kodeWeek4pointPredictionfor1LIRReduced[[i]] <- temp16[(i+17)%28)+1])

# Now after refining the data files we need to
# look at the predictions for one week, two weeks,
# three weeks, and four weeks ahead and determine the shape of the
# Kolmogorov-Smirnov test statistics plots.

#####sum(is.na(kodeWeek1pointPredictionfor1LIRReduced[[i]]))
#####sum(is.na(kodeWeek1pointPredictionfor1LIRReduced[[i]]))

# Both sets of data files have length of 28 now.
# So there are two groups of 28 discrete distributions.

# The testStats variable has now been turned into a
# list of four vectors, each of which serves the function
# of the original testStats variable.
# This is done in order that we can generate four graphs,
# corresponding to one week, two weeks, three weeks, and
# four weeks ahead.

testStats <- vector("list", 4)

for (i in 1:4) {
  testStats[[i]] <- numeric(28)

  # In order to generate the CDF for the true K-S
  # stats, first we're going to need to generate
  # a vector corresponding to the actual CDF.
  trueTestStats <- vector("list", 4)
  for (i in 1:4) {
    trueTestStats[[i]] <- numeric(28)

    kodeWeek1CDF <- kodeWeek1Reduced
    kodeWeek2CDF <- kodeWeek2Reduced
    kodeWeek3CDF <- kodeWeek3Reduced
    kodeWeek4CDF <- kodeWeek4Reduced

    kodeWeek1CDF <- kodeWeek1Reduced
    kodeWeek2CDF <- kodeWeek2Reduced
    kodeWeek3CDF <- kodeWeek3Reduced
    kodeWeek4CDF <- kodeWeek4Reduced

    for (j in 1:28) {
      kodeWeek1CDF[[j]]$Value[j] <- sum(kodeWeek1[[j]]$Value[0:j])
      kodeWeek2CDF[[j]]$Value[j] <- sum(kodeWeek2[[j]]$Value[0:j])
      kodeWeek3CDF[[j]]$Value[j] <- sum(kodeWeek3[[j]]$Value[0:j])
      kodeWeek4CDF[[j]]$Value[j] <- sum(kodeWeek4[[j]]$Value[0:j])

      kodeWeek1CDF[[j]]$Value[j] <- sum(kodeWeek1Reduced[[j]]$Value[0:j])
      kodeWeek2CDF[[j]]$Value[j] <- sum(kodeWeek2Reduced[[j]]$Value[0:j])
      kodeWeek3CDF[[j]]$Value[j] <- sum(kodeWeek3Reduced[[j]]$Value[0:j])
      kodeWeek4CDF[[j]]$Value[j] <- sum(kodeWeek4Reduced[[j]]$Value[0:j])
    }
  }
}

```

We're basically rotating the order in
 which our files are stored within an
 array; not only does this diminish the
 reader's confidence in the accuracy of
 the code (although this does impact
 the Kolmogorov-Smirnov test accuracy),
 but it causes the entire function to
 break down and become unreadable if we
 ever want to adapt it to
 compare other predictive models, also.


```
theme(text = element_text(size = 9)) geom_hline(aes(yintercept = -0.422058), label = "0.1, yjust = -1") +
  geom_hline(aes(yintercept = 0.1186577), label = "0.1186577", label = "0.05, yjust = -0.5") + geom_hline(aes(yintercept = 0.1069283), label = "0.1, yjust = 1")
```

The test statistics seems to oscillate to some degree, and the initial spike in their values seems to indicate that a linear regression still is most likely not the best fit.

Of the true II, of course takes on a positive value that both models try to predict, but with fairly significant differences. p -Values should be extracted from these test statistics to determine the magnitude of this significance.

It makes sense that our models tend to differ more as they try to predict further into the future.

If believe we look for possible linearity we should create graphs of each individual distribution in order to determine how the models are responsible for the learner values of the first three points on the graphs.

```
# The following code will generate graphs of the probabilities
# assigned to each bin for one week ahead, two weeks ahead,
# and three weeks ahead
kde1week <- kde1week <- kde2weeks <- kde2weeks <-
kde3weeks <- kde3weeks <- kde4weeks <- kde4weeks <-
vector("list", 6)
for (i in 1:6) {
  kde1week[[i]] <- kde1week[[i]] <- kde2week[[i]] <-
  kde2week[[i]] <- kde2weeks[[i]] <- kde1weeks[[i]] <-
  kde4weeks[[i]] <- kde4weeks[[i]] <-
  numeric(15)}

```

```
0 Because kettlweek1 is a data frame of size 131
# and we basically want to explain the fact that
# the first three K_S test statistics between the
# kettle and kike models are much larger than the rest,
# we use the first for-loop to populate the kettle_week
# list with four thirteen vectors that contain all
# probability predictions for each bin for that week ahead
# We arbitrarily chose to have four vectors in each
# list because this would show the probability
# distributions that correspond to the first
# four K_S test statistics on the graph and might
# provide some insight into which model is the
# culprit in the observed increased difference.
```

Creating new vectors isn't actually necessary, but with it make the code for the actual graphs slightly smaller, if which is what we want. It will just allow us to focus on what we want.

```
for (int i = 3; i < 10; i++) {
    kode1week[i][1][0] <- kodeWeek(1)[1][5].Value[0];
    kode1week[i][1][1] <- kodeWeek(1)[1][5].Value[1];
    kode1week[i][1][2] <- kodeWeek(1)[1][5].Value[2];
    kode2week[i][1][0] <- kodeWeek(2)[1][5].Value[0];
    kode2week[i][1][1] <- kodeWeek(2)[1][5].Value[1];
    kode2week[i][1][2] <- kodeWeek(2)[1][5].Value[2];
    kode3week[i][1][0] <- kodeWeek(3)[1][5].Value[0];
    kode3week[i][1][1] <- kodeWeek(3)[1][5].Value[1];
    kode3week[i][1][2] <- kodeWeek(3)[1][5].Value[2];
    kode4week[i][1][0] <- kodeWeek(4)[1][5].Value[0];
    kode4week[i][1][1] <- kodeWeek(4)[1][5].Value[1];
    kode4week[i][1][2] <- kodeWeek(4)[1][5].Value[2];
    kode5week[i][1][0] <- kodeWeek(5)[1][5].Value[0];
    kode5week[i][1][1] <- kodeWeek(5)[1][5].Value[1];
    kode5week[i][1][2] <- kodeWeek(5)[1][5].Value[2];
}
```

```

# Now we must also define the plots
max1 <- max(kcode1$week[[1]], kcode1$week[[2]], kcode1$week[[3]], kcode1$week[[4]], kcode1$week[[5]], kcode1$week[[6]])

```

For some reason we have to define the legend not
directly using the `paste0` command.

```
legend <- paste("w", modelname1, sep = "")  
plot(w=4.5,MWeeksAhead <- g[,lineofkde1week[1]] ~ seq_along(kde1week[1])), color = -1, legend.xlab = "Week 4.5", ylab = "") %>%  
+ lineofkde1week[1] ~ seq_along(kde1week[1]), color = ~paste("w", modelname2, sep = "")) %>% g$limsfy = c(0, max(1))) %>%
```

Of course, it's well-documented with respect to the reason why I am doing these things.

This whole function is definitely overcomplicated because if you're trying to compare different models, you have to make sure that the boolean "is_kkt_test" is changed, changed not to FALSE but to the string "is not kkt test" because it's a boolean in string form — this confuses a sub-function which runs only when the boolean is true.

[illegible]


```

xlab = "Index", ylab = "Difference in Predicted ILI" %>% ggl_jmisy = c(0, max2)) %>%
gf_line(pointValues[4]) ~ seq_along(pointValues[4]), color = ~"w ahead",
xlab = "Index", ylab = "Difference in Predicted ILI", title = "Difference in Point Values" + theme(text =
element_text(size = 9))

```

```

grid.arrange(mwe11, plot, allPointValueDifferences)

```

```

grid.arrange(plotw43, AllWeeksAhead, plotw44, AllWeeksAhead, plotw45, AllWeeksAhead, plotw46, AllWeeksAhead, plotw47, AllWeeksAhead,
plotw48, AllWeeksAhead, top = paste(modelname1, " and ", modelname2, " PMFS for %s ILL All Weeks Ahead"))

```

```

grid.arrange(kStatsAllWeeksAhead,

```

```

linarityTestMaximums, actualLLAgainstMaximums, top = paste("Comparing ", groupname, "s ", modelname1,
" and ", modelname2, " Models in for Selected Weeks of 2017-2018 (Week 43, 2017 to ",

```

```

"Week 18, 2018) in Maximum Differences between the ", modelname1, " and ", modelname2, " Models' in",

```

```

"Probability Mass Functions for All ILL Bins (increment 0.1)", sep = ""))

```

```

grid.arrange(truWeekStatsAllWeeksAhead, linarityTestKStatistics, actualLLAgainstKStatistics, top = paste("Comparing ", groupname, "s ",
modelname1,

```

```

" and ", modelname2, " Models in for Selected Weeks of 2017-2018 (Week 43, 2017 to ",

```

```

"Week 18, 2018) in Kolmogorov-Smirnov Test Statistics between the ", modelname1, " and ",
modelname2, " Models' in Empirical CDF for All ILL Bins (increment 0.1)", sep = ""))

```