

Modélisation Prédictive Rapport

Valentin Gözl, Laura Fuentes

February 28, 2023

Contents

Introduction	2
Choix du type de modèle	2
Linear Models	3
Random Forest	4
Modèles additifs généralisés	5
ARIMA et Kalman Filter	7
Pipeline basée sur le modèle qgam	9
Aggrégation d’experts	12
Conclusion	15

Introduction

Nous avons un jeu de données regroupant différentes variables en rapport avec la consommation énergétique française pendant la période de 2012 à 2021. Notre but est de construire un modèle qui permet de prédire la consommation française en énergie pendant la période du Covid.

Le premier réflexe est de télécharger l'ensemble des packages et diviser le set train en deux pour pouvoir tester nos modèles avant de les soumettre. Nous avons ainsi choisi la période de 2012 - 2019 comme train et 2019-(15/04/2020) comme test.

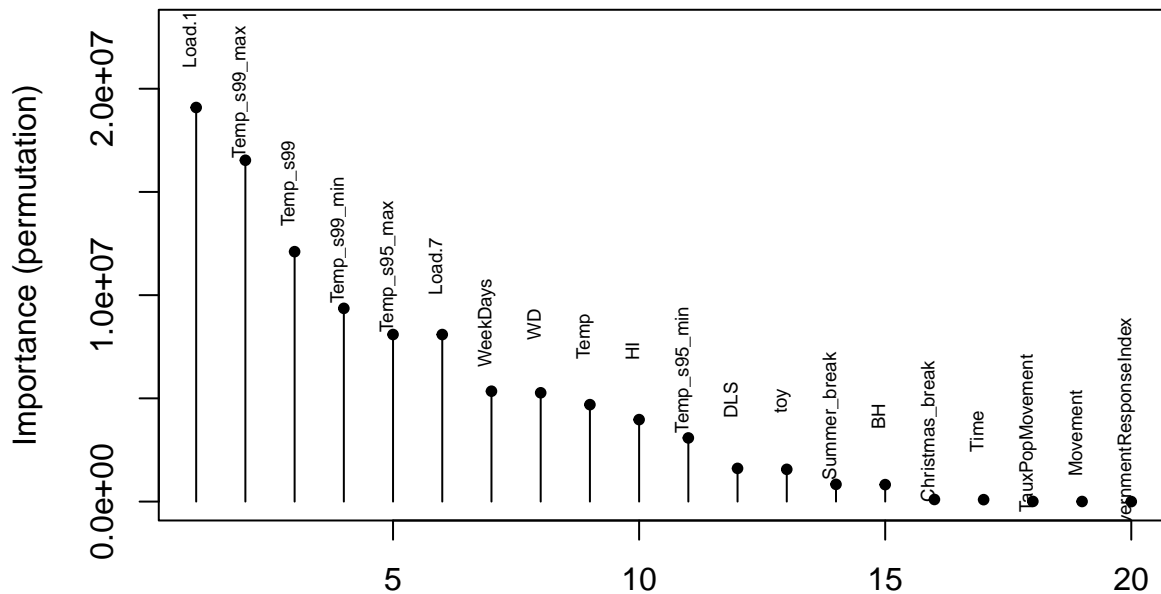
Choix du type de modèle

Pour commencer, nous avons créé tout d'abord, la variable **WeekDays2**. Il s'agit d'une version modifiée de la variable WeekDays qui distingue les jours laboraux, samedis et dimanches. Nous avons d'autre part récupéré des données des différents mouvements sociaux et le pourcentage de population mobilisée. Pendant les jours avec des manifestations, **Movement** est 1 et 0 sinon. Les données viennent de SNCF. Ensuite nous avons créé une variable mesurant la température ressentie. Le problème de cette dernière variable concernait les nombreuses valeurs NA's, ainsi que la représentativité au niveau national des stations météorologiques constituant les données.

Pour comprendre quelles variables sont plus significatives, et argumenter le choix, nous allons effectuer une random forest, et regarder l'importance des variables.

```
equation <- "Load~ Time + toy + Temp + Load.1 + Load.7 + Temp_s99 + WeekDays + BH + Temp_s95_max +
  Temp_s99_max + Summer_break + Christmas_break +
  Temp_s95_min +Temp_s99_min + DLS + GovernmentResponseIndex + TauxPopMovement + Movement + HI + WD"
rf <- ranger(equation, data=Data0, importance = 'permutation')

#####importance plot
imp <- rf$variable.importance
imp <- sort(imp)
o <- order(imp, decreasing=T)
nom <- names(imp)
plot(c(1:length(imp)), imp[o], type='h', ylim = c(0, max(imp) + max(imp)/5), xlab='', ylab='Importance
K <- length(imp)
text(tail(c(1:length(imp)), K), tail(imp[o]+max(imp)/8), K), labels= tail(nom[o], K), pos=3, srt=90, adj
points(c(1:length(imp)), imp[o], pch=20)
```



Nous pouvons ainsi bien remarquer que les variables à plus forte importance sont: Load.1, Load.7, les variables relatives à la température, WeekDays, WD, BH, toy, Summer_break, DLS and Christmas_break. Nous avons également vérifié que la variable **Movement** (=taux de la population qui a manifesté) n'était pas explicative à l'aide de ANOVA.

```
formula1 <- "Load ~ Temp + WeekDays + WD + BH+ toy + Summer_break + DLS + Christmas_break"
formula2 <- "Load ~ Temp + WeekDays + WD + BH+ toy + Summer_break + DLS + Christmas_break + Movement"
small_lm <- lm(formula1%>%as.formula, data=Data0)
large_lm <- lm(formula2%>%as.formula, data=Data0)
anova.res <- anova(small_lm, large_lm)
summary(anova.res)
```

##	Res.Df	RSS	Df	Sum of Sq
##	Min. :2931	Min. :3.954e+10	Min. :1	Min. :4987977
##	1st Qu.:2931	1st Qu.:3.954e+10	1st Qu.:1	1st Qu.:4987977
##	Median :2932	Median :3.954e+10	Median :1	Median :4987977
##	Mean :2932	Mean :3.954e+10	Mean :1	Mean :4987977
##	3rd Qu.:2932	3rd Qu.:3.954e+10	3rd Qu.:1	3rd Qu.:4987977
##	Max. :2932	Max. :3.954e+10	Max. :1	Max. :4987977
##			NA's :1	NA's :1
##	F	Pr(>F)		
##	Min. :0.3698	Min. :0.5432		
##	1st Qu.:0.3698	1st Qu.:0.5432		
##	Median :0.3698	Median :0.5432		
##	Mean :0.3698	Mean :0.5432		
##	3rd Qu.:0.3698	3rd Qu.:0.5432		
##	Max. :0.3698	Max. :0.5432		
##	NA's :1	NA's :1		

Linear Models

Simple linear model

Pour comprendre et appréhender le cadre d'étude nous commencerons par effectuer un modèle simple. C'est-à-dire un modèle linéaire avec les co-variables choisies précédemment. Nous avons considéré que la consommation de la veille changeait en fonction du jour de la semaine. C'est pour cela que nous avons décidé

de créer une fonction de la consommation de la veille en fonction de chaque catégorie de **WeekDays**.

```
formula <- "Load ~ Load.1 + WeekDays + Load.7 + Temp + Temp_s99_max + Temp_s99_min + WeekDays + WD + BH"
slm <- lm(formula %>% as.formula, data = Data0)
pred1 = predict(slm, newdata = Data1)
```

```
## Warning in predict.lm(slm, newdata = Data1): prediction from a rank-deficient
## fit may be misleading
```

```
rmse(Data1$Load, pred1)
```

```
## [1] 2153
```

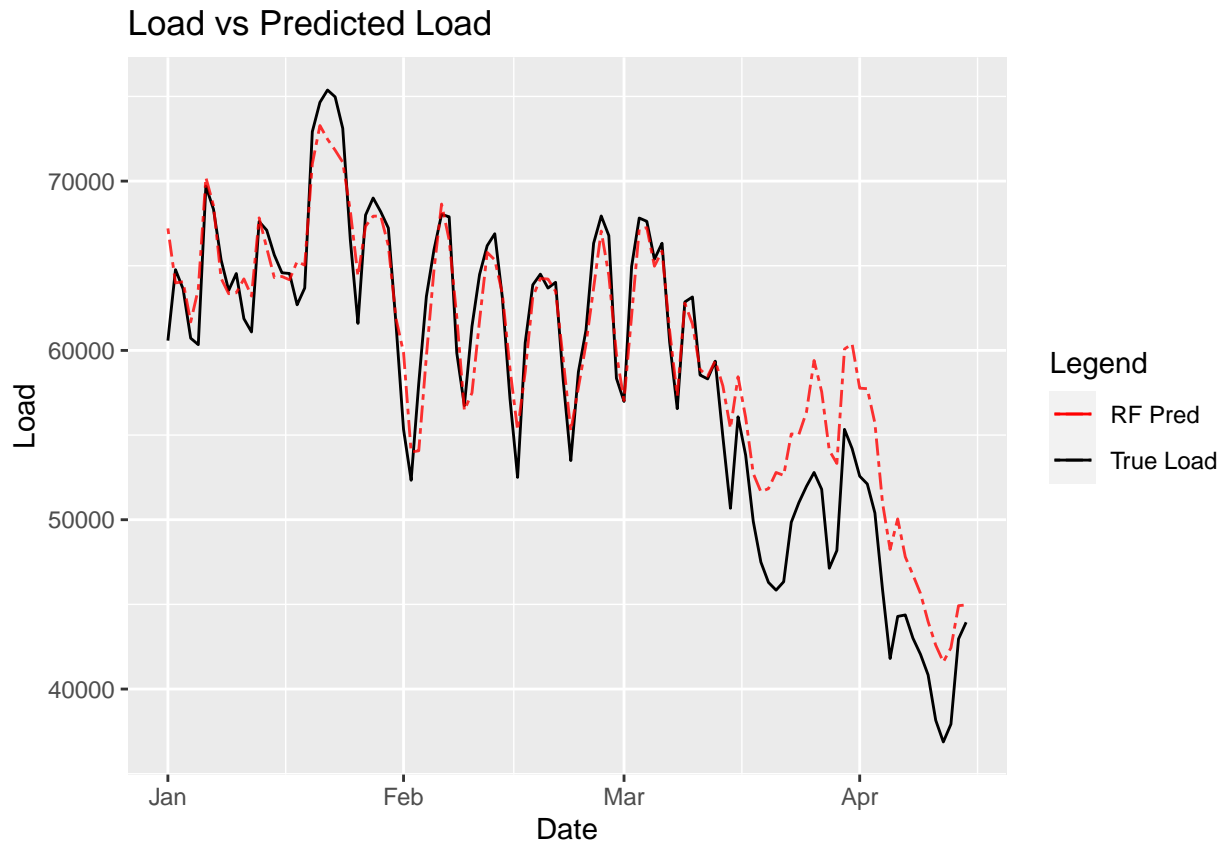
Polynomial transformations

Nous avons mis en place des transformations polynomiales sur le modèle linéaire comme une première approche de complexification du modèle. Pour cela, nous avons ajouté les co-variables relatives à la température et Load.7 au carré.

Random Forest

```
formule <- "Load ~ Month + Temp_s95_min + Temp_s95_max + HI + TauxPopMovement + Time + toy + Temp + Load"

rf <- ranger::ranger(formule, data = Data0, importance = 'permutation')
rf.forecast <- predict(rf, data = Data1)$predictions
plot_data <- tibble(x = Data1$Date, y = Data1$Load, y1 = rf.forecast)
p <- ggplot(plot_data, aes(x)) +
  geom_line(aes(y = y, color = "True Load")) +
  geom_line(aes(y = y1, color = "RF Pred", linetype = "twodash", alpha = 0.8)) +
  labs(x = "Date", y = "Load", title = "Load vs Predicted Load") +
  scale_color_manual(name = "Legend", values = c("True Load" = "black", "RF Pred" = "red"))
p
```



```
rmse(Data1$Load, rf.forecast)
```

```
## [1] 3104
```

Modèles additifs généralisés

Choix de la partie linéaire et spline

Dans la suite, nous avons considéré de mettre en place des Modèles Additifs généralisés. Pour cela, nous avons d'abord distingué les variables à mettre dans la partie linéaire du modèle, puis dans la partie spline. Nous avons intégré les variables qualitatives ainsi que la consommation de la veille en fonction du jour de la semaine dans la partie linéaire. On a ajouté dans la partie spline les variables ayant une notion de temporalité comme la consommation de la semaine ou les températures. Nous avons également regroupé dans une même spline des variables ayant une relation logique, comme c'était le cas avec la température et le temps. Nous avons également créé une fonction spline pour chaque jour de la semaine pour la variable `toy` pour ne pas négliger l'effet des jours de la semaine sur la consommation annuelle.

```
equation <- "Load ~ Load.1:as.factor(WeekDays) + HI + BH + Christmas_break + Summer_break + DLS + s(Temp
```

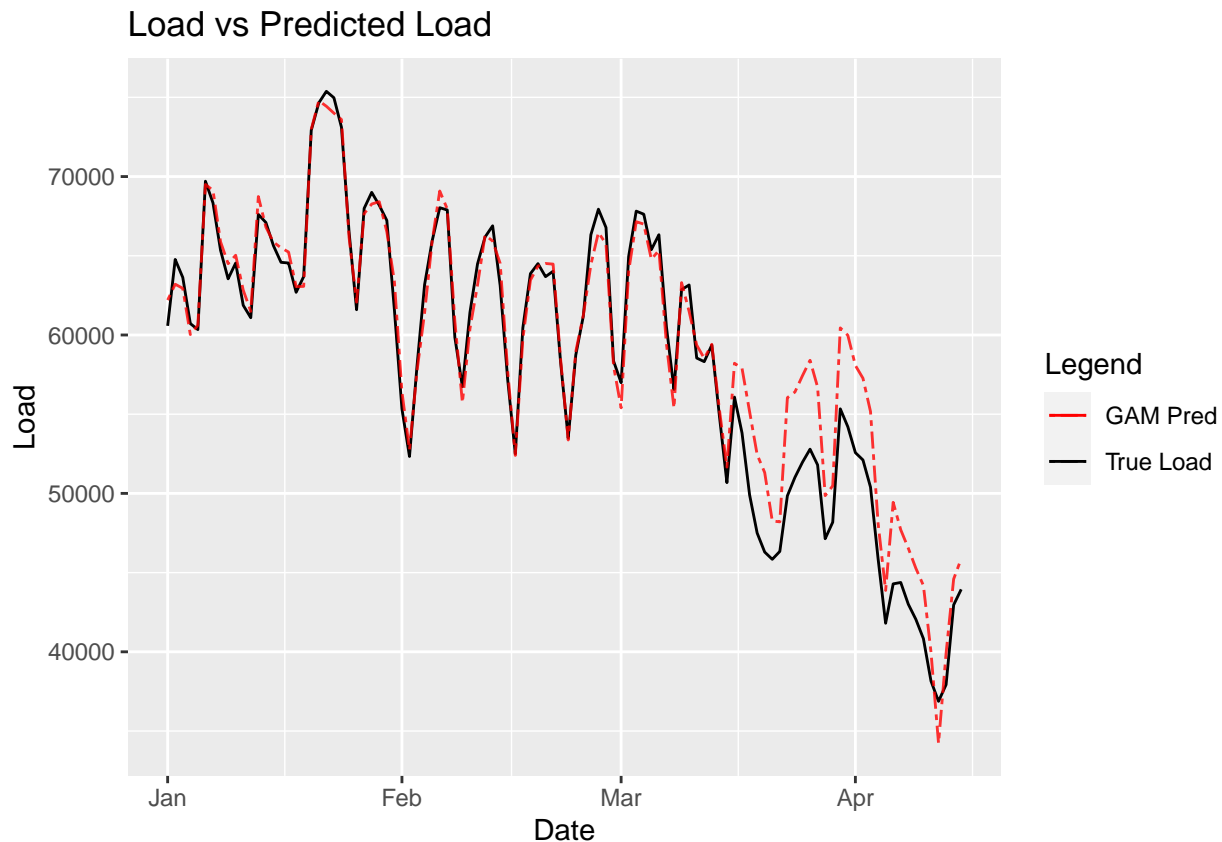
```
gam<-gam(equation%>%as.formula, data=Data0)
```

```
gam.forecast <- predict(gam, newdata=Data1)
```

```
plot_data <- tibble(x=Data1$Date, y=Data1$Load, y1=gam.forecast)
```

```
p <- ggplot(plot_data, aes(x)) +
  geom_line(aes(y=y, color="True Load")) +
  geom_line(aes(y=y1, color="GAM Pred", linetype="twodash", alpha=0.8)) +
  labs(x="Date", y="Load", title="Load vs Predicted Load") +
  scale_color_manual(name="Legend", values=c("True Load"="black", "GAM Pred"="red"))
```

p



```
rmse(Data1$Load, gam.forecast)
```

```
## [1] 2316
```

Pour améliorer le rendement du modèle, nous avons tenté de comprendre l'origine des erreurs à partir des courbes de consommation. Nous avons constaté que les erreurs commencent à s'accroître au niveau du mois de mars 2020, juste au niveau du début de la période covid. Ceci s'explique du fait que la variable **GouvernementResponseIndex** comprend des valeurs nulles pendant des années, et celles-ci explosent dans une courte période d'un mois, laissant peu de temps d'entraînement sur la pandémie. Pour simuler le comportement de la population pendant le confinement avec les données que l'on avait déjà, nous avons pensé aux samedis. En effet, nous avons émis l'hypothèse qu'un jour de confinement était comparable en termes de consommation à un jour de weekend comme un samedi. Dans cet esprit, nous avons créé la variable WD, qui modifie le jour de la semaine à samedi s'il y a confinement (la `GouvernementResponseIndex >= 70`), et maintient le jour de la semaine inchangé sinon. !! -> mettre online-learning

Nous avons également utilisé la fonction `gam.check` pour améliorer le rendement du modèle `gam`. Celle-ci nous a permis d'ajuster la dimension des bases des splines. Nous avons ainsi incrémenter les valeurs de `k` quand la `p`-value était très petite. Pour la variable `toy`, utiliser un `k` très grand faisait tourner le modèle trop long, nous avons donc pris `k=30` même si la `p`-value était encore petite. Enfin, nous avons également vérifié que les résidus étaient bien gaussiens à chaque fois à partir de l'histogramme issu du plot. !!! COMENTER ?!!

GAM et régression quantile: `qgam`

Dans la suite on utilisera le package `qgam`, et en particulier la fonction `qgam`. Celle-ci ajuste un modèle additif ainsi qu'une régression quantile sur un unique quantile. On utilise ici la même équation qu'auparavant, il suffit juste d'ajuster la variable "qu", correspondant au quantile. Après plusieurs essais, nous avons remarqué qu'on obtient de meilleurs résultats avec "qu" autour de 0.4. En effet, en fixant le quantile à 0.4, on change la

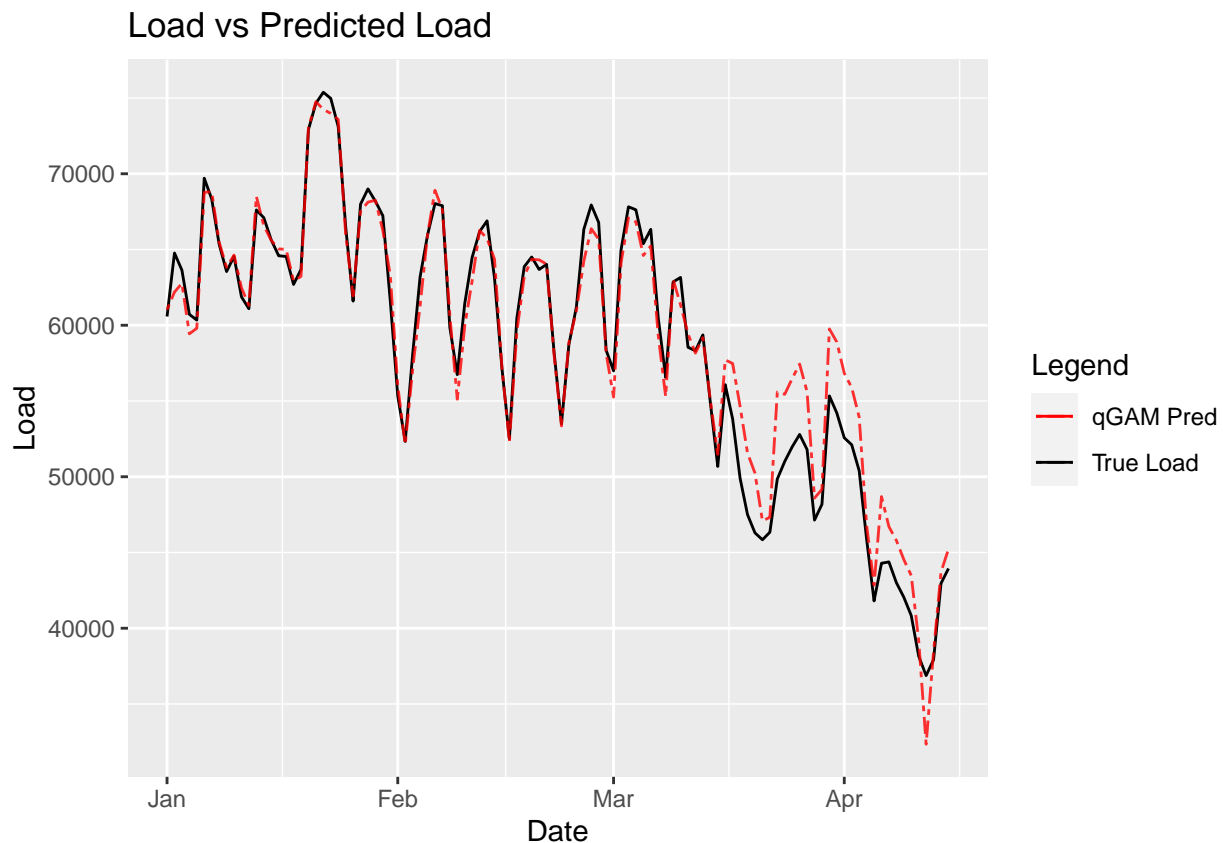
fonction de perte. On introduit ainsi un biais, qui permet de s'ajuster mieux aux données lors de la période du covid.

```
gam9<-qgam(equation%>%as.formula, data=Data0, qu=0.4)

## Estimating learning rate. Each dot corresponds to a loss evaluation.
## qu = 0.4.....done

gam9.forecast <- predict(gam9, newdata=Data1)

plot_data <- tibble(x=Data1$Date, y=Data1$Load, y1=gam9.forecast)
p <- ggplot(plot_data, aes(x)) +
  geom_line(aes(y=y, color="True Load")) +
  geom_line(aes(y=y1, color="qGAM Pred", linetype="twodash", alpha=0.8)) +
  labs(x="Date", y="Load", title="Load vs Predicted Load") +
  scale_color_manual(name="Legend", values=c("True Load"="black", "qGAM Pred"="red"))
p
```



```
rmse(Data1$Load, gam9.forecast)
```

```
## [1] 1955
```

ARIMA et Kalman Filter

ARIMA

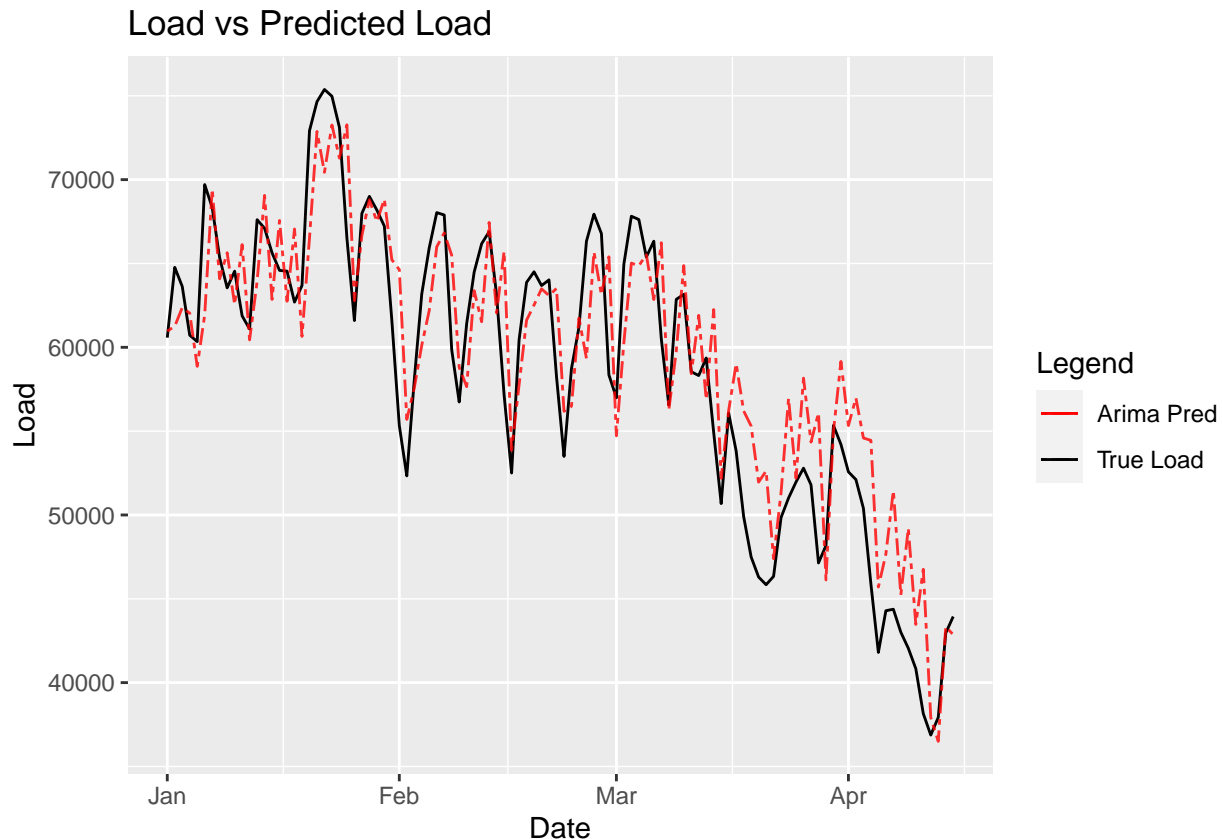
ARIMA (Autoregressive Integrated Moving Average) est un modèle de série chronologique qui utilise les valeurs et les erreurs passées pour prévoir les valeurs futures d'une série. Le modèle peut être ajusté pour capturer les tendances et la saisonnalité, et implique de différencier les données pour les rendre stationnaires. L'ARIMA est largement utilisé pour faire des prévisions sur la base de tendances historiques.

```

arima.fit <- forecast::Arima(gam9.forecast, order = c(1,1,2), seasonal = c(0,0,2))
arima.predict <- fitted(arima.fit)

plot_data <- tibble(x=Data1$Date, y=Data1$Load, y1=arima.predict)
p <- ggplot(plot_data, aes(x)) +
  geom_line(aes(y=y, color="True Load")) +
  geom_line(aes(y=y1, color="Arima Pred", linetype="twodash", alpha=0.8)) +
  labs(x="Date", y="Load", title="Load vs Predicted Load") +
  scale_color_manual(name="Legend", values=c("True Load"="black", "Arima Pred"="red"))
p

```

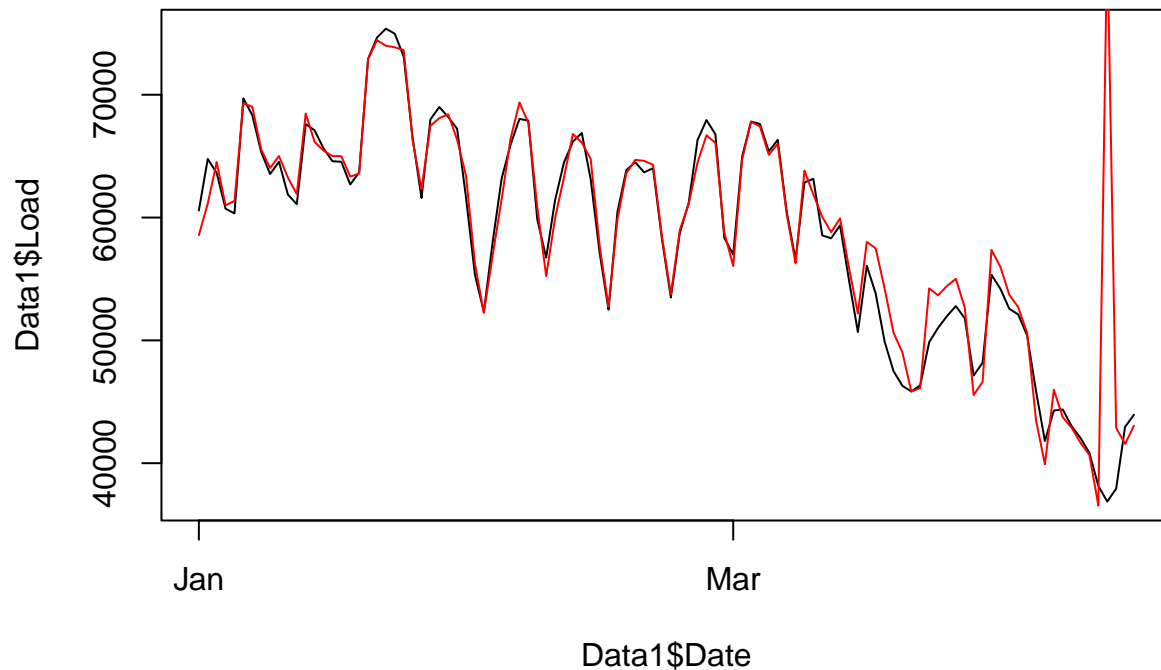


```
rmse(arima.predict, Data1$Load)
```

```
## [1] 3969
```

Filtre Kalman

Une fois le modèle qgam amélioré le modèle qgam, nous avons décidé d'implémenter le online-learning. Cette méthode nous permet d'aborder le problème de données d'entraînement insuffisantes pour la période covid. En effet, le filtre Kalman dynamique met à jour l'estimation des poids séquentiellement, au fur et mesure des prédictions ... ajouter ici ce que dynamique fait en particulier. Pour simplifier la situation covid, nous avons créé la variable GRI_factor, étant une variable catégorielle à trois facteurs (none: 0-50, medium: 50-62, high:62-70). Ceci nous permet faire une mise au point de la situation covid avec du online-learning. Malheureusement, les résultats étaient meilleurs sans cette variable. C'est pour cela que nous avons décidé de continuer avec le modèle qgam précédent pour le online learning.



Pipeline basée sur le modèle qgam

Étant arrivés au bout des améliorations de qgam, et tout en ayant tenté arima et le filtre kalman, nous avons considéré d'autres modèles vus en cours pour comparer les performances. On a ainsi décidé de garder notre équation sur la qgam et de l'implémenter ensuite sur d'autres modèles. On a ainsi étudié les résidus. Ceci va ainsi nous permettre de ???? \ Nous avons ainsi décidé de tester les forêts aléatoires sur les résidus de qgam. Après avoir appliqué l'effet des forêts aléatoires sur le modèle, nous avons amélioré davantage la performance à l'aide de Arima.

```
Nblock<-10
borne_block<-seq(1, nrow(Data0), length=Nblock+1)%>%floor
block_list<-list()
l<-length(borne_block)
for(i in c(2:(l-1)))
{
  block_list[[i-1]] <- c(borne_block[i-1]:(borne_block[i]-1))
}
block_list[[l-1]]<-c(borne_block[l-1]:(borne_block[l]))

blockRMSE<-function(equation, block)
{
  g<- gam(as.formula(equation), data=Data0[-block,])
  forecast<-predict(g, newdata=Data0[block,])
  return(forecast)
}

Block_forecast<-lapply(block_list, blockRMSE, equation=equation)%>%unlist
Block_residuals <- Data0$Load-Block_forecast

####estimation of GAM, GAM effects
g <- gam9
g.forecast <- predict(g, newdata=Data1)
```

```

terms0 <- predict(g, newdata=Data0, type='terms')
terms1 <- predict(g, newdata=Data1, type='terms')
colnames(terms0) <- paste0("gterms_", c(1:ncol(terms0)))
colnames(terms1) <- paste0("gterms_", c(1:ncol(terms1)))

Data0_rf <- data.frame(Data0, terms0)
residualsCV <- Block_residuals

Data0_rf$residuals <- residualsCV
Data0_rf$res.48 <- c(residualsCV[1], residualsCV[1:(length(residualsCV)-1)])
Data0_rf$res.336 <- c(residualsCV[1:7], residualsCV[1:(length(residualsCV)-7)])

Data1_rf <- data.frame(Data1, terms1)
residuals <- Data1_rf$Load - gam9.forecast
Data1_rf$residuals <- residuals
Data1_rf$res.48 <- c(residuals[1], residuals[1:(length(residuals)-1)])
Data1_rf$res.336 <- c(residuals[1:7], residuals[1:(length(residuals)-7)])

cov <- "Time + toy + Temp + Load.1 + Load.7 + Temp_s99 + WeekDays + BH + Temp_s95_max + Temp_s99_max + 
Temp_s95_min +Temp_s99_min + DLS + GovernmentResponseIndex + res.48 + res.336 +"
gterm <-paste0("gterms_", c(1:ncol(terms0)))
gterm <- paste0(gterm, collapse='+')
cov <- paste0(cov, gterm, collapse = '+')
formule_rf <- paste0("residuals", "~", cov)
rf_gam<- ranger::ranger(formule_rf, data = Data0_rf, importance = 'permutation')
rf_gam.forecast <- predict(rf_gam, data = Data1_rf)$predictions+ g.forecast

rmse(y=Data1$Load, ychap=rf_gam.forecast)

## [1] 1877

rf_gam$variable.importance%>%sort

```

## GovernmentResponseIndex	Summer_break	gterms_4
## 0.0000	940.5953	1402.5147
## gterms_11	gterms_15	gterms_17
## 2553.0828	2978.8074	4069.7756
## gterms_16	gterms_13	gterms_12
## 11211.0152	11591.9690	13924.5274
## Christmas_break	gterms_3	gterms_5
## 14881.1763	15969.2436	20063.5504
## DLS	gterms_8	res.336
## 27339.6767	36087.2203	38081.5193
## WeekDays	gterms_7	gterms_14
## 38698.9074	48744.7089	58378.9778
## gterms_2	BH	gterms_1
## 73699.5371	76106.2958	107780.4395
## Temp_s95_min	Temp	Temp_s99_min
## 111815.0582	115775.3755	123911.6507
## Temp_s95_max	Temp_s99_max	Temp_s99
## 127467.2865	127988.1977	128811.4388
## gterms_9	Load.7	gterms_18
## 130055.3509	131333.5721	148409.6593
## res.48	gterms_10	Load.1

```
##          150414.4203          161142.3806          170269.0572
##          Time          toy          gterms_6
##          173681.9470          178510.7395          189330.8968
```

```
Block_residuals.ts <- ts(Block_residuals, frequency=7)
fit.arima.res <- auto.arima(Block_residuals.ts,max.p=3,max.q=4, max.P=2, max.Q=2, trace=T,ic="aic", method="ML")
```

```
##
## Fitting models using approximations to speed things up...
```

```
## ARIMA(2,1,2)(1,0,1)[7] with drift          : 48455.59
## ARIMA(0,1,0) with drift          : 50047.92
## ARIMA(1,1,0)(1,0,0)[7] with drift          : 49275.05
## ARIMA(0,1,1)(0,0,1)[7] with drift          : 48676.12
## ARIMA(0,1,0) : 50045.94
## ARIMA(2,1,2)(0,0,1)[7] with drift          : 48618.76
## ARIMA(2,1,2)(1,0,0)[7] with drift          : 48466.6
## ARIMA(2,1,2)(2,0,1)[7] with drift          : Inf
## ARIMA(2,1,2)(1,0,2)[7] with drift          : Inf
## ARIMA(2,1,2) with drift          : 48640.09
## ARIMA(2,1,2)(0,0,2)[7] with drift          : 48619.8
## ARIMA(2,1,2)(2,0,0)[7] with drift          : Inf
## ARIMA(2,1,2)(2,0,2)[7] with drift          : Inf
## ARIMA(1,1,2)(1,0,1)[7] with drift          : 48501.34
## ARIMA(2,1,1)(1,0,1)[7] with drift          : 48457.88
## ARIMA(3,1,2)(1,0,1)[7] with drift          : 48461.34
## ARIMA(2,1,3)(1,0,1)[7] with drift          : Inf
## ARIMA(1,1,1)(1,0,1)[7] with drift          : 48499.64
## ARIMA(1,1,3)(1,0,1)[7] with drift          : Inf
## ARIMA(3,1,1)(1,0,1)[7] with drift          : 48467.54
## ARIMA(3,1,3)(1,0,1)[7] with drift          : 48457.95
## ARIMA(2,1,2)(1,0,1)[7] : 48453.68
## ARIMA(2,1,2)(0,0,1)[7] : 48617.07
## ARIMA(2,1,2)(1,0,0)[7] : 48464.6
## ARIMA(2,1,2)(2,0,1)[7] : Inf
## ARIMA(2,1,2)(1,0,2)[7] : Inf
## ARIMA(2,1,2) : 48638.38
## ARIMA(2,1,2)(0,0,2)[7] : 48618.12
## ARIMA(2,1,2)(2,0,0)[7] : Inf
## ARIMA(2,1,2)(2,0,2)[7] : Inf
## ARIMA(1,1,2)(1,0,1)[7] : 48499.43
## ARIMA(2,1,1)(1,0,1)[7] : 48456.01
## ARIMA(3,1,2)(1,0,1)[7] : 48459.34
## ARIMA(2,1,3)(1,0,1)[7] : Inf
## ARIMA(1,1,1)(1,0,1)[7] : 48497.72
## ARIMA(1,1,3)(1,0,1)[7] : 48476.44
## ARIMA(3,1,1)(1,0,1)[7] : 48465.56
## ARIMA(3,1,3)(1,0,1)[7] : 48454.25
```

```
##
## Now re-fitting the best model(s) without approximations...
```

```
## ARIMA(2,1,2)(1,0,1)[7] : 40113.13
```

```
##
## Best model: ARIMA(2,1,2)(1,0,1)[7]
```

```
#Best model: ARIMA(2,1,2)(2,0,0)[7]
#saveRDS(fit.arima.res, "../Results/tif.arima.res.RDS")
ts_res_forecast <- ts(c(Block_residuals.ts, Data1$Load-gam9.forecast), frequency= 7)
refit <- Arima(ts_res_forecast, model=fit.arima.res)
prevARIMA.res <- tail(refit$fitted, nrow(Data1))
gam9.arima.forecast <- gam9.forecast + prevARIMA.res
```

Aggrégation d'experts

Comme dernière méthode, nous avons décidé de mettre en place un agrégation d'experts pour extraire une combinaison de prédicteurs qui puissent améliorer davantage la performance du modèle. Pour cela, nous avons regroupé les différents prédicteurs dans une variable experts. Dans cette agrégation d'experts, nous avons utilisé les différents modèles gam, qgam (avec et sans arima), une forêt aléatoire comprenant toutes les variables, un filtre kalman, la pipeline... Le modèle obtenu par combinaison des différents prédicteurs obtient une performance bien meilleure que celle obtenue auparavant. Comme on peut le voir dans la figure, le modèle BOA permet d'ajuster les coefficients des modèles de façon variable sur chaque partie.

```
experts <- cbind(gam9.forecast, pred2.rf.forecast, gam.forecast, gam9.kalman.Dyn, arima.predict, as.num)
nom_exp <- c("qgam", "polynomial_lm", "rf", "gam", "kalman", "arima", "pipeline")
colnames(experts) <- nom_exp
```

```
rmse_exp <- apply(experts, 2, rmse, y=Data1$Load)
sort(rmse_exp)
```

```
##      pipeline      qgam polynomial_lm      gam      rf
##      1505      1955      2000      2316      3104
##      arima      kalman
##      3969      4700
```

```
cumsum_exp <- apply(Data1$Load-experts, 2, cumsum)
```

```
#Correction du biais
```

```
expertsM2000 <- experts-2000
expertsP2000 <- experts+2000
experts <- cbind(experts, expertsM2000, expertsP2000)
colnames(experts) <-c(nom_exp, paste0(nom_exp, "M"), paste0(nom_exp, "P"))
cumsum_exp <- apply(Data1$Load-experts, 2, cumsum)
```

```
par(mfrow=c(1,1))
agg <- mixture(Y = Data1$Load, experts = experts, loss.gradient=TRUE)
summary(agg)
```

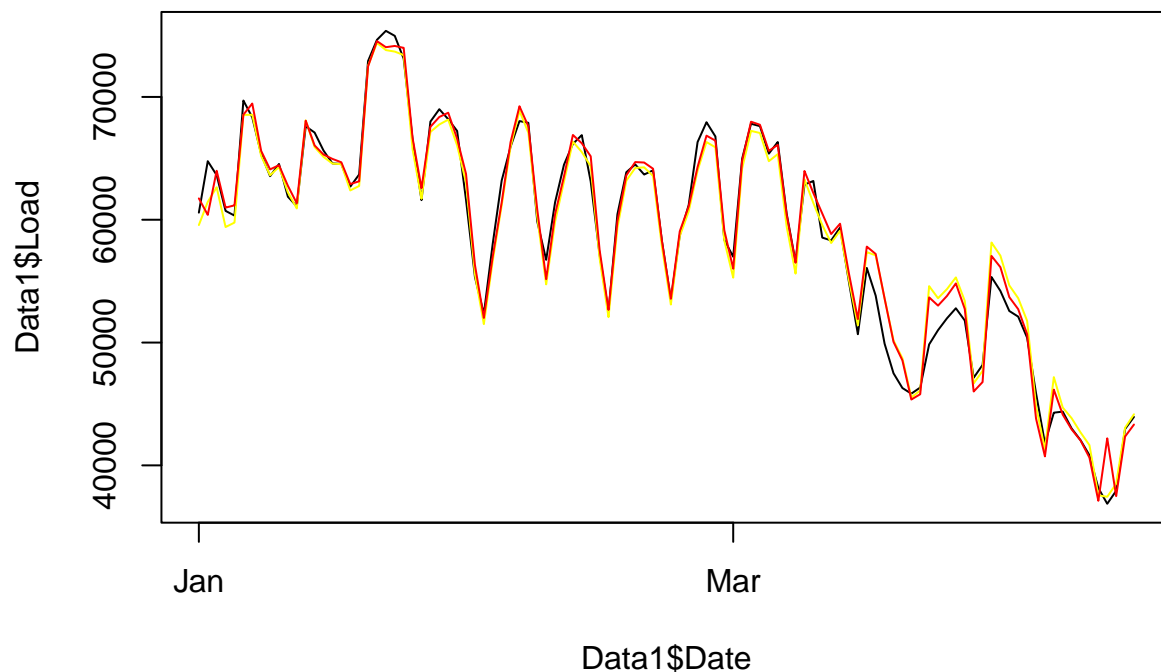
```
## Aggregation rule: MLpol
## Loss function: squareloss
## Gradient trick: TRUE
## Coefficients:
##      qgam polynomial_lm rf      gam kalman arima pipeline qgamM polynomial_lmM
## 0.0587      0.0283 0 0.0261      0 0.0127      0.101 0.146      0.118
##      rfM gamM kalmanM arimaM pipelineM qgamP polynomial_lmP rfP gamP kalmanP
## 0.0912 0.117      0 0.104      0.186      0      0      0      0      0
## arimaP pipelineP
##      0      0.0112
##
## Number of experts: 21
## Number of observations: 106
```

```
## Dimension of the data: 1
##
##          rmse  mape
## MLpol    1340 0.0174
## Uniform  2000 0.0289

####PLOT EXPERT AGGREGATION
#plot(agg)
or <- oracle(Y=Data1$Load, experts);or

## Call:
## oracle.default(Y = Data1$Load, experts = experts)
##
## Coefficients:
##      qgam polynomial_lm      rf      gam kalman      arima pipeline      qgamM
## 3.03e-21      3.03e-21 3.03e-21 3.03e-21 0.0139 3.03e-21      0.295 3.03e-21
## polynomial_lmM      rfM      gamM kalmanM      arimaM pipelineM      qgamP
##      3.03e-21 3.03e-21 3.03e-21 0.0631 3.03e-21      0.338 3.03e-21
## polynomial_lmP      rfP      gamP kalmanP      arimaP pipelineP
##      3.03e-21 3.03e-21 3.03e-21 0.0466 3.03e-21      0.243
##
##
##          rmse  mape
## Best expert oracle:    1510 0.0193
## Uniform combination:    2000 0.0289
## Best convex oracle:    1360 0.0180

par(mfrow=c(1,1))
plot(Data1$Date, Data1$Load, type='l')
lines(Data1$Date, or$prediction, type='l', col='yellow')
lines(Data1$Date, agg$prediction, type='l', col='red')
```



```
rmse(agg$prediction, y=Data1$Load)
```

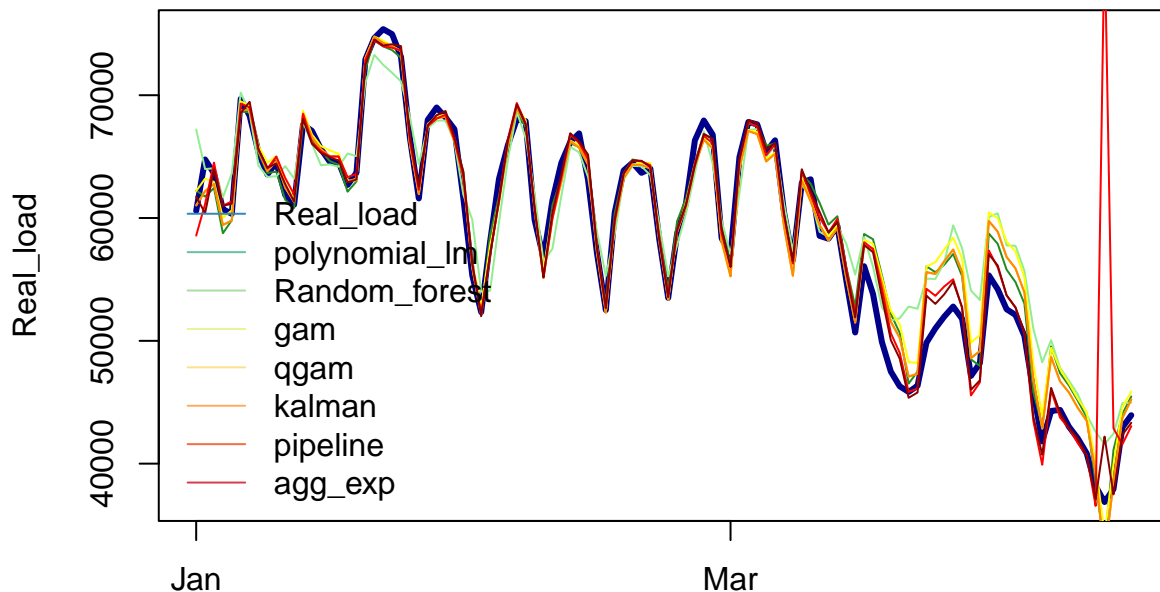
```
## [1] 1340
```

La figure finale avec presque tous nos modèles sur les données train (sel_a):

```
models = c(Data1$Load, pred2, rf.forecast, gam.forecast, gam9.forecast, gam9.kalman.Dyn, gam9.arima.for
K <- ncol(models)
col <- rev(RColorBrewer::brewer.pal(n = max(min(K,8),6),name = "Spectral"))[1:min(K,8)]
nom_mod <- c("Real_load", "polynomial_lm", "Random_forest", "gam", "qgam", "kalman", "pipeline", "agg_exp")

plot(Data1$Date, Data1$Load, type='l', ylab='Real_load', col="darkblue", main="Prédictions", lwd=3 )
lines(Data1$Date, pred2, type='l', ylab = "polynomial_lm", col='forestgreen')
#Random forest
lines(Data1$Date, rf.forecast, type='l', ylab='rf', col= "lightgreen")
#gam
lines(Data1$Date, gam.forecast, type='l', ylab='gam', col="yellow")
#qgam
lines(Data1$Date, gam9.forecast, type='l', ylab='qgam', col="orange")
#kalman
lines(Data1$Date, gam9.forecast, type='l', ylab='qgam', col="darkorange")
#Pipeline
lines(Data1$Date, gam9.kalman.Dyn, type='l', ylab='qgam', col="red")
#Aggregation d'experts
lines(Data1$Date, agg$prediction, type='l', ylab='agg_exp', col='darkred')
legend("bottomleft", col=col, legend=nom_mod, lty=1, bty='n')
```

Prédictions



Nous pouvons voir sur la figure, que le modèle le plus performant est effectivement issu de la combinaison des autres modèles. !!!! COMMENTER !!!

%%% A Rajouter les critères utilisés pour des bons modèles.

Conclusion

Dans ce travail, nous avons testé plusieurs modèles différents pour prédire la charge électrique de la France. Si nos modèles les plus simples donnaient déjà des résultats exploitables, les plus complexes étaient encore plus prometteurs. Le plus grand défi a été le fait que la période de test a été fixée pendant la pandémie de covid, rendant la prévision plus difficile. Notre meilleur modèle est un modèle additif généralisé avec régression quantile. L'utilisation de techniques telles que l'agrégation d'experts et l'apprentissage en ligne a encore amélioré notre précision. Un autre défis, était le fait qu'au niveau des jeux de données train issus de sel_a, on n'avait pas encore du covid. Ceci créait un biais entre les scores prévus et ceux sur Kaggle. En effet, en entraînant le jeu de données train tout entier, on avait un mois de covid, ce qui améliorait nettement nos résultats.