

Projet NF26

Données : France de 2005 à 2014

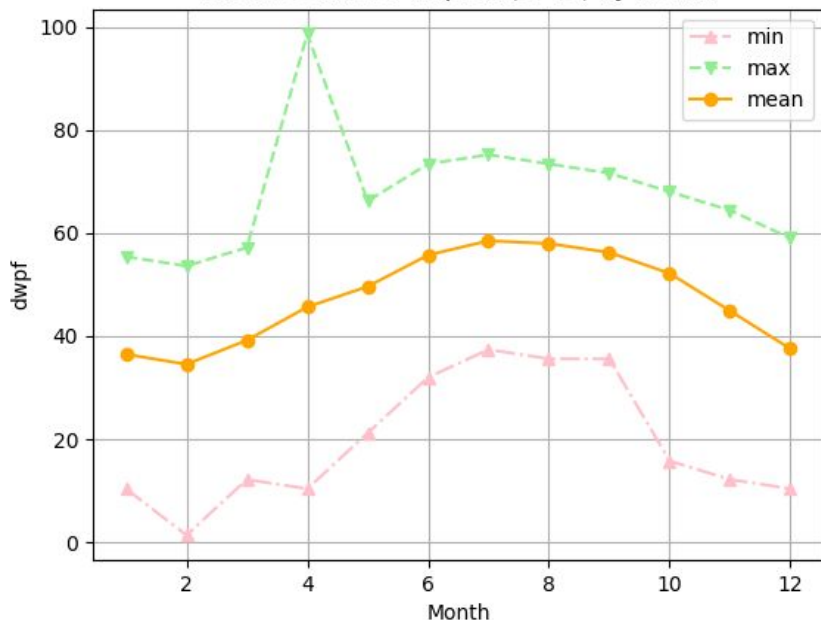
Hugo DURET
Haojie LU

I. - Résultats du projet : Q1

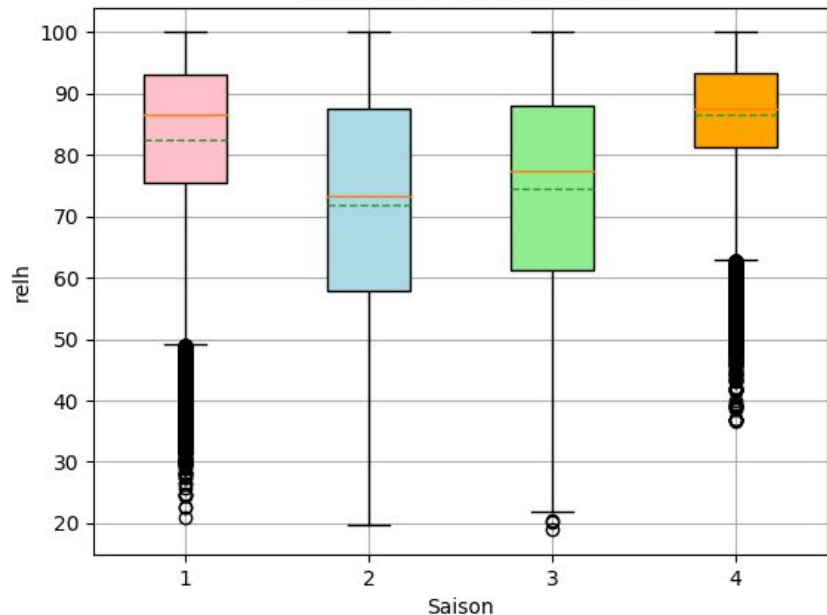
→ `linechartparmois(point, indicateur)` ⇒ `linechartparmois((3,43), 'dwpf')`

→ `boxplot(point, indicateur)` ⇒ `boxplot('LFQQ', 'relh')`

The line chart of dwpf at (3, 43) by month

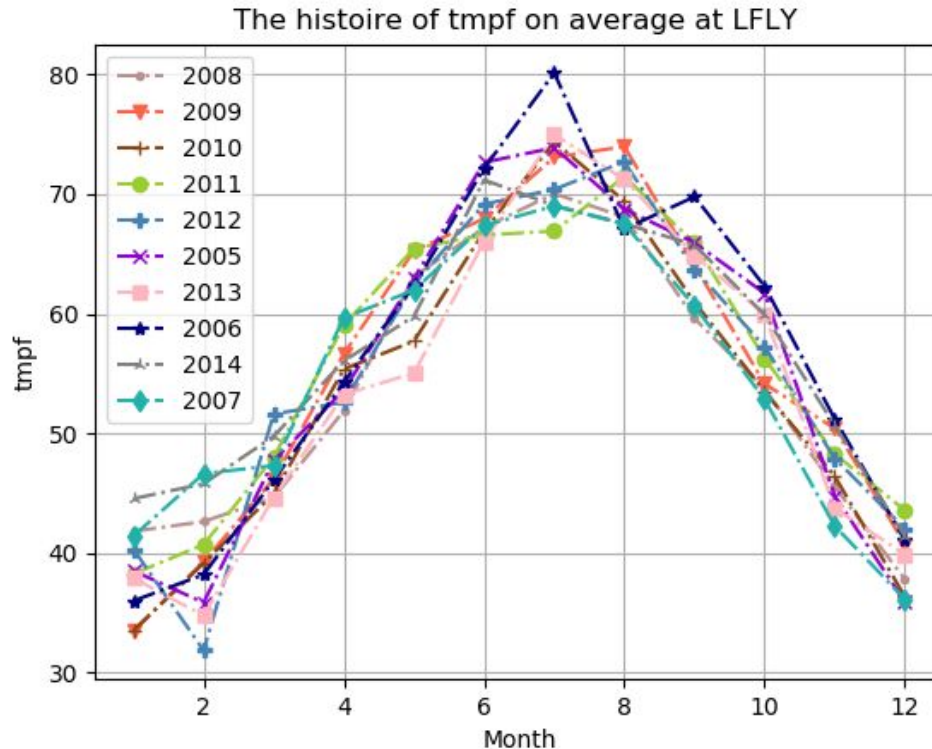


The boxplot of relh at LFQQ



I. - Résultats du projet : Q1

→ *linechart_histoire(point, indicateur)* ⇒ *linechart_histoire('LFLY', 'tmpf')*

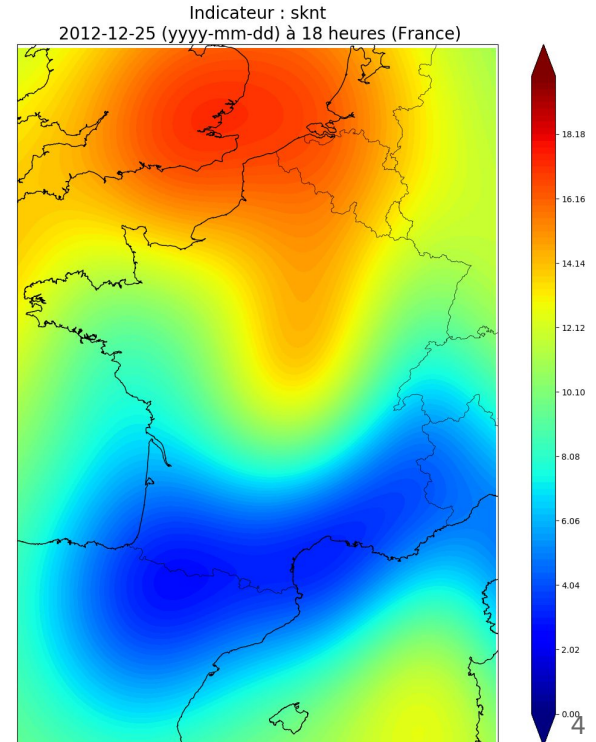
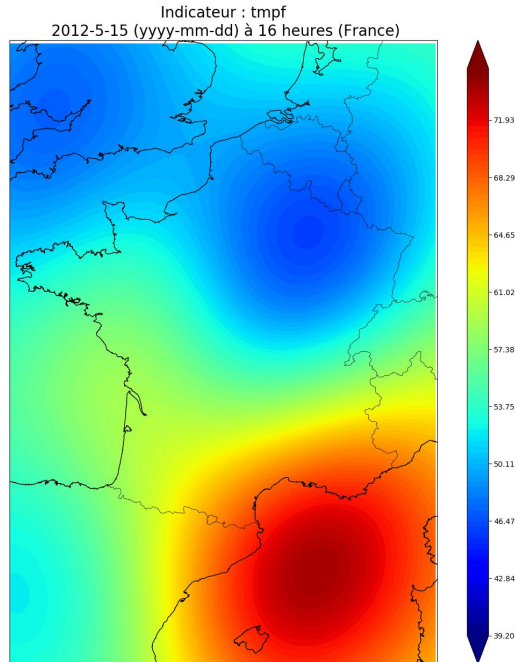


I. - Résultats du projet : Q2

→ `get_plot_per_hour_and_indicator(year, month, day, hour, indicator)`

⇒ `get_plot_per_hour_and_indicator(2012, 5, 15, 16, 'tmpf')`

⇒ `get_plot_per_hour_and_indicator(2012, 12, 25, 18, 'sknt')`

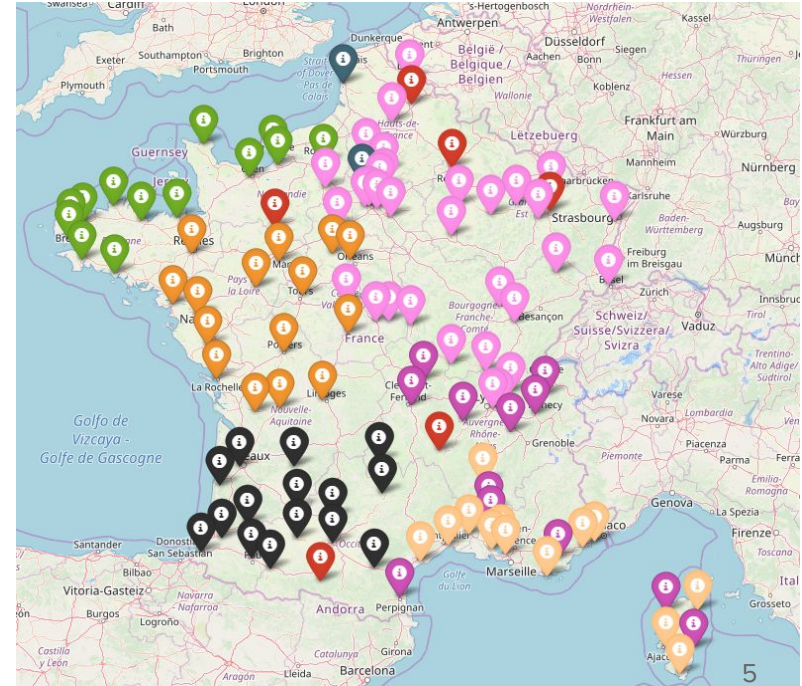
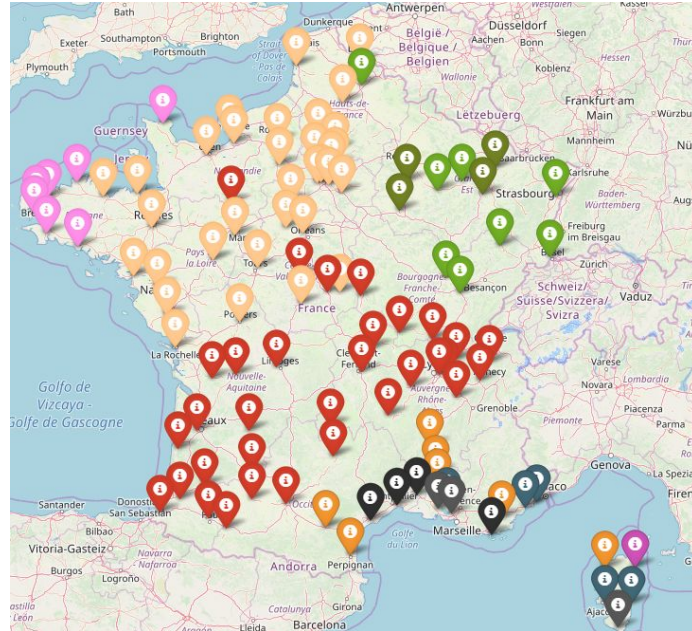


I. - Résultats du projet : Q3

→ *mapofclusterisation*(start_date, end_date, nombre_de_groupe(10 par défaut))

⇒ *mapofclusterisation*('2012-06-09', '2012-08-09')

⇒ *mapofclusterisation*('2012-03-09', '2012-08-22',8)



II. - Choix de stockage

Q1

Clé de partitionnement	Clés de tri
station	year, month, day, hour, minute

Q2

Clé de partitionnement	Clés de tri
year, month, day, hour	minute, station

Q3

Clé de partitionnement	Clés de tri
timestamp_day	hour, minute, station

III- Traitement : Q1

Line Charts par mois :

$(\text{station}, \text{année}, \text{mois}, \text{jour}, \text{ind}) \Rightarrow (\text{mois}, (1, \text{ind}, \text{ind}, \text{ind}))$

$r : (a, b) \Rightarrow (a[0]+b[0], \min(a[1], b[1]), \max(a[2], b[2]), a[3]+b[3])$

$(\text{mois}, (s_0, s_1, s_2, s_3)) \Rightarrow (\text{mois}, (s_1, s_2, s_3/s_0))$

Boxplots pour la saisonnalité :

$(\text{station}, \text{année}, \text{mois}, \text{jour}, \text{ind}) \Rightarrow (\text{saison}, \text{ind}) \Rightarrow (\text{saison}, [\text{itérable}])$

On crée les boxplots à partir du mapping obtenu précédemment.

III- Traitement : Q1

Line Charts par mois pour toute la période :

$(\text{station}, \text{année}, \text{mois}, \text{jour}, \text{ind}) \Rightarrow (\text{année}, \text{mois}, (1, \text{ind}))$

$r : (a, b) \Rightarrow (a+b)$

$((\text{année}, \text{mois}), (s_0, s_1)) \Rightarrow ((\text{année}, \text{mois}), (s_1/s_0))$

$((\text{année}, \text{mois}), s_0) \Rightarrow (\text{année}, (\text{mois}, s_0)) \Rightarrow (\text{année}, [\text{itérable}])$

III- Traitement : Q2

Création d'une heatmap pour un indicateur :

Données des stations pour une heure

Krigeage

Heatmap

III- Traitement : Q3

Traitement de données

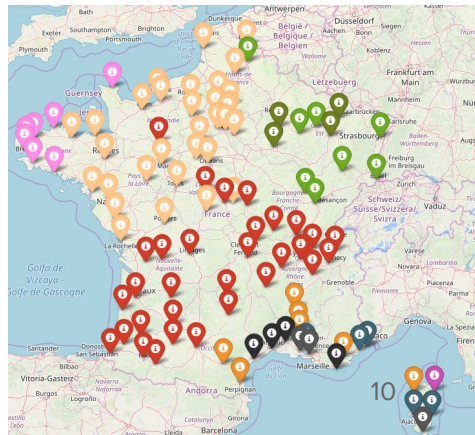
- Principe : numérique, non nul => (tmpf, dwpf, relh, drct, sknt, alti, vsby, skyl1, feel)
- Moyenne

Clusterisation : Kmeans

- Initialiser au hasard K centres des catégories.
- Mapping : le centre i le plus proche, la distance d'Euler, une pair clé-valeur ($i, (1, \text{enregistrement})$).
- Réduction et Mapping : mettre à jour le centre de chaque catégorie avec le moyen de tous les enregistrements actuels de cette catégorie.

Représentation

- Folium
- (latitude, longitude)
- HTML



Conclusion

Gros volume de données

Problèmes dans les données (valeurs nulles, stations sans données avant 2010...)

Choix de stockage spécifique à une question

Cassandra pour le stockage en Orienté Colonnes

Spark pour des calculs en streaming

MERCI