

# Analyse des données de pluies au Japon

Mathias Ballot

Mars 2024

# Introduction

Le Japon est régulièrement touché par des désastres naturels tel que des typhons, des averses violentes ou des inondations. Le taux de précipitations par ans est deux fois plus élevé que la moyenne mondiale, et tombe en majorité durant la saison des pluies entre mai et août. Face a de tels évènements, le développement d'outils de mesure et de prédiction des précipitations urge.

Le Japon étant une île montagneuse l'eau des précipitations a tendance à rapidement ruisseler vers la mer. Des barrages permettent de contrôler le flux des précipitations, mais sont souvent débordées par les averses violentes devenues de plus en plus fréquentes [4]. Des plans d'actions prévoient le relâchement anticipée de l'eau stockée dans les barrages afin d'éviter le débordements lors de l'averse. Les barrages régulent à eux seules près de 50% de l'eau utiliser dans les foyers japonais. L'énergie hydroélectrique qui représente 8% de la production électrique japonaise reste l'énergie renouvelable la plus fiable dans un tel pays. Afin d'optimiser la production électrique du barrage l'eau relâchée doit être remplacée par celle provenant des averse. Pour cela des estimations précises de la quantité de pluies tombant dans les alentours du barrage est nécessaire. Les prévisions doivent étre estimées 2 à 3 heures en avances, ce qui correspond à l'afflux de l'eau des précipitations dans les réservoirs [5].

Pour répondre à ces besoins le Japon développe depuis les années 1950 des radars d'observations météorologique de plus en plus performants [1]. Ces radars qui n'étaient capable que de détecter des zones de pluies peuvent maintenant estimer l'intensité et la forme des gouttes de pluie. En 1964, un radar a été construit au point culminant du Japon : le mont Fuji. Il était à l'époque le radar le plus haut du monde.

Une autre manière de mesurer les précipitations consiste à disséminer des jauges de mesures sur le sol Japonais. Ces jauges permettent de mesurer avec précision la quantité de pluie en un point.

Les résultats suivants sont une analyse des données de pluies au Japon. Les modèles de machine développés à partir de ces données ne sont pas présenter dans ce papier, mais le code est accessible depuis le gitHub.

De tel prédictions servent à réguler les quantités d'eau stockée dans les barrages, à la préventions de désastres naturels, l'approvisionnement d'eau dans les foyers, l'agriculture...

# Chapter 1

## Analyse des données

### 1.1 Japan Meteorological Agency

La "Japan Meteorological Agency" décompose l'analyse des précipitations en 2 catégories [3].

- Estimation météorologique : Consiste à déterminer la quantité et l'intensité des précipitations passées ou présentes. Ces données sont collectées par des jauges et des radars.
- Prévision météorologique : prédictions des futures précipitations à partir d'un modèle météorologique. Ce modèle peut être basé sur des équations mathématique, des règles thermodynamiques, la théorie des fluides...

#### 1.1.1 Estimation météorologique

Dans cette partie, nous nous intéressons aux estimations météorologiques. Pour estimer la pluie toutes les 30 minutes avec une résolution spatiale de 1 km, la "Japan Meteorological Agency" utilise deux instruments de mesures :

- Les jauges de pluie : Conteneur physique collectant une certaine quantité de pluie. Les jauges sont des instruments de mesures fiables et précis. Une jauge ne peut mesurer la pluie qu'en un seul emplacement, il est donc nécessaire d'en déployer partout afin de mesurer les précipitations sur le sol japonais. Plus de 10.000 jauge sont disséminées à intervalle régulier de 7 km à travers le territoire japonais, collectant des données toutes les 10 minutes ou 1 heures.
- Les radars de pluie : Radar émettant des impulsions courtes de manière périodique. Les ondes émises sont partiellement réfléchies par les gouttes de pluies, donnant des informations sur l'intensité, la position et les déplacements des précipitations. Un radar peut couvrir à lui seul une surface de  $500\text{km}^2$ . 46 radars mesurent les précipitations au Japon, générant une

cartographie des pluies avec une résolution spatiale de 1 km et à intervalle de 5 minutes.

Pour réaliser les estimations météorologique, les données radars sont sommées puis calibrer grâce aux mesures de jauges. L'atténuation des ondes radar et le caractère uniforme des précipitations sont utilisés pour affiner les résultats. Dans la figure 1.1 la figure centrale est obtenu en utilisant les mesures radars (à gauche) et les mesures de jauges (à droite). On remarque que les mesures de jauges sont ponctuels, contrairement aux mesures radars.

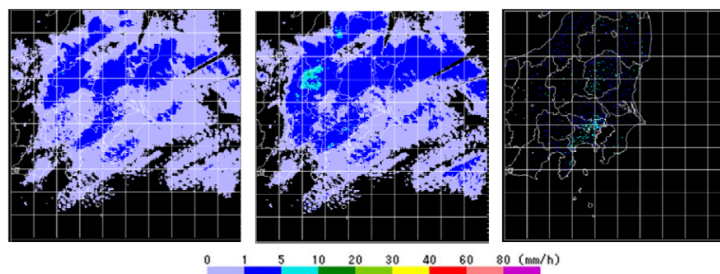


Figure 1.1: Gauche : mesures radar; Droite : mesures jauge;  
Centre : Estimations combinées

### 1.1.2 Prévisions météorologique

La prévision météorologique sur court terme (6 heures) se base sur 2 éléments :

- Extrapolation des estimations : Les estimations mesurées par les radars et jauge sont généralisée aux 6 heures suivantes. L'intensité et la direction du vent est pris en compte afin de prédire les mouvements des précipitations sur les prochaines heures. La topologie montagneuse du Japon causant les nuages à se condenser est aussi prise en compte.
- Modèle méso-échelle : Le terme méso-échelle est une échelle de grandeur représentant les phénomènes météorologique entre 2km et 2.000km (échelle planétaire). Un tel modèle se base sur des équations de masse, de dynamisme, d'énergie interne et toute ce qui touche la mécanique des fluides.

Les extrapolations sont précises sur les 3 premières heures de prévisions. Le modèle méso-échelle vient compléter les extrapolation sur des heures de prévision plus lointaines. Dans la figure 1.2 les extrapolations révèlent une zone intense de pluie au nord-est du Japon (encadré en rouge), alors que le modèle méso-échelle montre de forte précipitations au coeur du Japon (encadré en bleue). Le mélange des deux prévisions se rapproche des quantités de précipitation réellement obtenus.

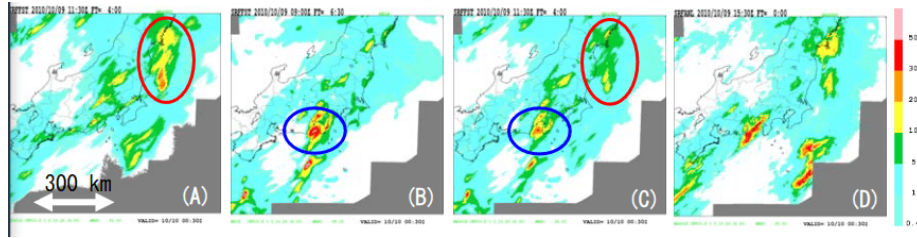


Figure 1.2: A: Extrapolation; B: Modèle méso-échelle; C : Combinaisons de A et B; D: Valeur réellement mesurée

### 1.1.3 Exemple d'application

Les prévisions météorologique sont utilisées pour mesurer l'indice d'humidité dans le sol. Des risques de glissement de terrains sont fréquent sur les sols ayant un fort indice. L'eau des précipitations restent stocke dans le sol pendant un certain temps avant de s'écouler, ces désastres peuvent donc être causé par des pluies ayant eu lieu quelques jours plus tôt. La figure 1.3 montre que le taux d'humidité dans le sol permet de prévoir un tel désastre.

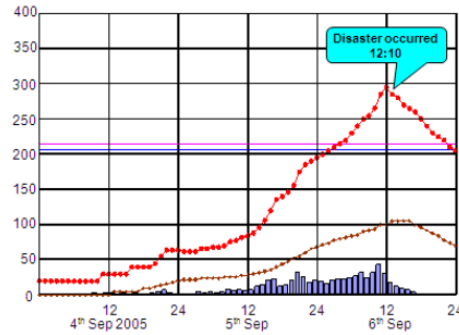


Figure 1.3: Rouge : indice d'eau dans le sol; Brown: précipitations cumulées sur les dernières 24 heures; barres: précipitations sur 1 heure

## 1.2 Régression linéaire

À partir des données radar et jauge, l'objectif est de développer un modèle de machine learning capable d'estimer la pluie. Il ne s'agit pas ici de prévoir la pluie, mais simplement d'établir un modèle capable déterminer les quantités de précipitations à un instant  $t$ , à partir de mesures radar passé, présente et future.

### 1.2.1 Données

Nous disposons de deux types de données collectées au coeur de l'île principale du Japon sur les 3 premiers mois de 2010:

- Radar (fig 1.4a) : Les données radars sont prises à intervalle temporel de 10 minutes. Les radars quadrillent la surface avec une dimension de  $204 * 160 = 21624$ . Chaque mesure représente la quantité de précipitations sur les 10 dernières minutes.
- Jauge (fig 1.4b) : Les données de jauge sont récupérées toutes les heures. Un total de 371 jauge couvrent cette surface. Les mesures représentent la quantité moyenne de pluie tombée toutes les 10 minutes sur la dernière heure. Pour chaque heure nous avons donc une seule valeur.

Cette zone est intéressante car elle contient les Alpes Japonaises, une chaîne de montagne traversant le coeur du Japon. Nous pouvons voir que des précipitations de quelques millimètres ont lieu dans ces montagnes (encadré en rouge). Vers Kanazawa (encadré en vert) le vent froid provenant de Sibérie est humidifié par la mer du Japon (située du côté de la Chine et la Corée). Ce vent est ensuite bloqué par les montagnes générant des nuages de précipitations.

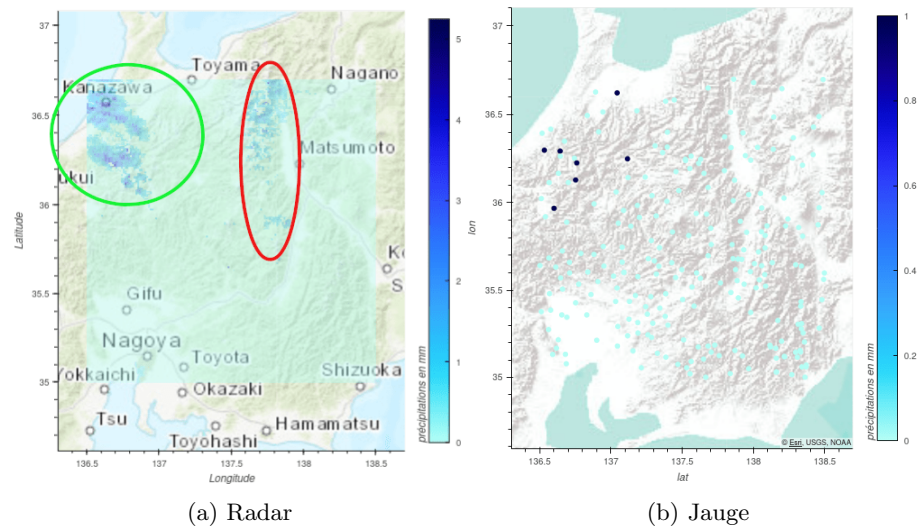


Figure 1.4: Mesure des précipitations le 2 janvier 2010 au Japon

### Analyse des données

La première étape en machine learning consiste à analyser ses données. Sur le graphe 1.5 on peut comparer les mesures du radar et de la jauge en un emplacement. Les mesures sont toujours exprimées en millimètres mesure sur

une certaine période. 1 millimètres de pluie correspond à 1 litre par mètre carré. On remarque que les mesure ne correspondent pas parfaitement. Alors que la jauge mesure des faible précipitations durant la journée du 13 février, le radar ne semble rien mesurer. La jauge est une mesure plus fiable que le radar à cet emplacement, on peut donc supposer que des interférences on empêchée le radar de mesurée la pluie à cet emplacement.

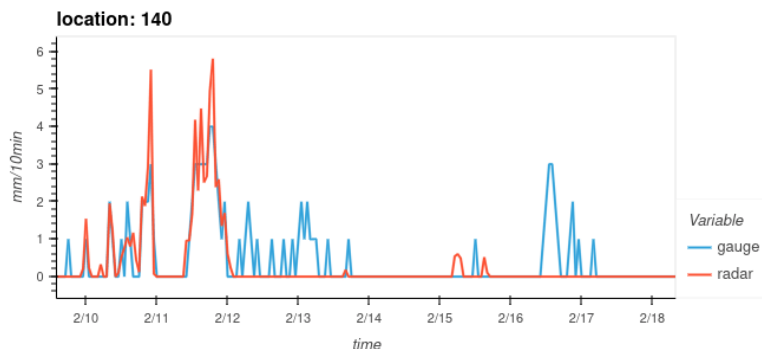


Figure 1.5: Quantité de pluie mesurée à un emplacement contenant une jauge

Pour aller plus loin nous cherchons à observer les mesure de radar potentiellement fausse. Par exemple si le radar mesure une forte précipitations à un emplacement, mais que les 8 emplacements voisins l'entourant ne mesure aucune pluie alors on peut estimer que la forte précipitation mesurée est fausse. Pour trouver ces mesures aberrantes nous calculons à chaque emplacements la plus petite différence de pluie entre cet emplacement et ces 8 voisins l'entourant. Si cette valeur est très élevé (ou très petite) on peut supposer que l'erreur provient du radar. Sur la figure 1.6 par exemple, la mesure de 25mm de pluie ne fait pas de sens.

L'une des caractéristiques de nos données comme nous pouvons le voir ici, est que la distribution des quantités de précipitation par tranche de 10 minutes n'a rien de gaussien. En effet 90 pourcent des données auront pour valeur 0 à 1 millimètres de pluie, et seuls quelques extrêmes se démarqueront avec des précipitations de 20 millimètres. Le traitement de valeurs extrêmes est généralement plus compliqué dans notre cas. En effet même si 20 millimètres de pluie est une valeur normale, elle peut paraître extrême comparé à la quantité énorme de valeurs proche de 0.

Les données de jauge non-plus ne sont pas parfaite, de nombreuses jauge ne semblent pas fonctionner car elles ne mesurent aucune pluie durant 3 mois. Certaines jauge nécessitent une intervention humaine pour récupérer les données de précipitations, d'autre sont automatiquement mesurée par outil d'enregistrement. Les erreur de mesure peuvent provenir de ces deux facteurs.

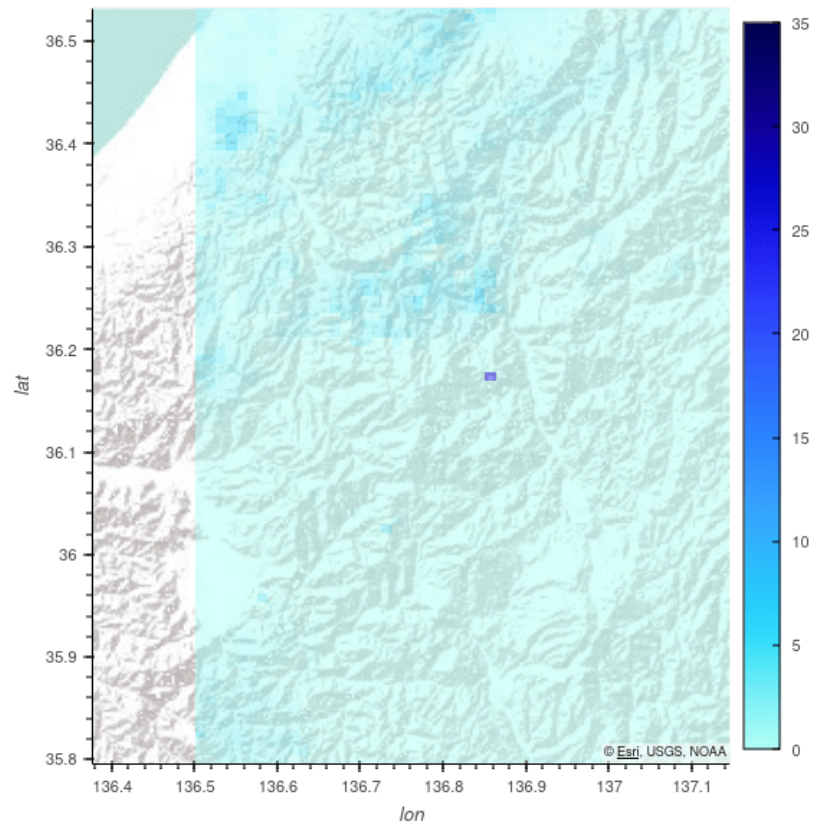


Figure 1.6: Valeurs aberrantes

### Traitement des données

Les données de jauge n'ayant aucune valeur sont simplement évincées, ces jauges ne seront donc pas utiliser dans notre modèle de machine learning. Les valeurs aberrantes sont réduites aux valeurs de leur plus proche voisin. On considère une valeur aberrante lorsque la mesure effectuée en un point est 20 mm plus élevé que ses 8 plus proches voisins.

Pour définir un modèle de machine learning il faut définir les données entrant dans ce modèle et celle sortant. Dans tout les modèles que nous testerons, les données entrantes seront les valeurs mesurée par le radar. Ces données sont normaliser afin d'être toute compris entre 0 et 1. Cela permet de mettre toute les donnée sur une même échelle. Les données rentrant seront du même format que les données entrante.

Il faut ensuite définir les exemples annotés servant de base d'apprentissage pour notre modèle. Les exemples annotés seront les valeurs mesurée par les jauges toute les heures.



Nous allons créer différent modèle de régression linéaire ayant pour objectif d'estimer la pluie tombée au niveau des jauges a partir des données radars. Nous ne prendrons les données que d'une jauge au départ afin de faciliter le problème, puis les modèles linéaires suivant utiliserons les données de toute les jauges. Les données en entrées sont donc les mesures radars effectuée au niveau de la jauge toute les 10 minutes. L'objectif est d'obtenir une estimation heure par heure proche des valeurs de jauge.

### 1.2.2 Régression lineaire sur une jauge

Avant même de créer notre premier modèle de régression linéaire, nous allons construire un modèle simple qui ne demande aucune connaissance et machine learning. Il s'agit simplement de calculer la moyenne des données radar pour chaque heure. Il nous servira de base de comparaison pour évaluer les modèles linéaires.

#### Définition du modèle linéaire

Nous créons ensuite un modèle de régression linéaire  $f$ .

$$f(X) = a * x_0 + b * x_1 + c * x_2 + d * x_3 + e * x_4 + f * x_5 \quad (1.1)$$

$$f(X) = \Theta X \quad (1.2)$$

Où :

$$\begin{aligned} x_i &= \text{valeur du radar a la minute 0, 10, 20... aussi appelé features} \\ a, b, c, d, e, f &= \text{paramètres de notre modèle} \\ X^T &= (x_0 \ x_1 \ x_2 \ x_3 \ x_4 \ x_5) \\ \Theta &= (a \ b \ c \ d \ e \ f) \end{aligned}$$

Puis une fonction coût  $J$  permettant d'évaluer les prédictions de notre modèle linéaire. La fonction que nous utilisons est l'erreur quadratique moyenne calculée par :

$$J(\Theta) = \frac{1}{n} \sum_{i=1}^n (y_{i\_jauge} - f(X_i))^2$$

Où :

$$\begin{aligned} n &= \text{le nombre d'échantillons dans notre jeu d'entraînement.} \\ f(X_i) &= \text{correspond à la valeur obtenu par notre modèle sur le } i\text{-ème échantillon.} \\ y_{i\_jauge} &= \text{la valeur mesurée par la jauge, c'est le label associé a } X_i. \end{aligned}$$

La fonction coût ne dépend que des paramètres ( $\Theta$ ) de notre modèle. Plus cette fonction est petite pour un échantillon  $X_i$ , plus les prédictions de notre modèle se rapproche des valeurs idéales de jauge. Nous cherchons donc à obtenir la meilleur combinaison de paramètre  $\Theta$  tel que la fonction coût soit le plus petit possible sur le jeu d'entraînement de notre modèle.

La méthode que nous utilisons pour trouver les meilleurs paramètres  $\Theta$  est un algorithme d'apprentissage supervisé nommée descente du gradient [2]. On appelle apprentissage supervisé un algorithme apprenant à associer un input  $x$  à un output  $y$ , en s'entraînant sur un jeu de données composé d'input labélisé. Le modèle est entraîné sur ces input labélisé durant un certain nombre d'itération. À chaque itération du modèle sur notre jeu d'entraînement nous calculons  $f(X)$  et nous comparons notre résultat avec le label  $y_{i\_jauge}$  en calculant  $J(\Theta)$ . À chaque itération les paramètres  $\Theta$  du modèle sont légèrement modifiés afin de minimiser la fonction coût lors de la prochaine itération. Cette étape est répétée pendant un certain nombre d'époques, jusqu'à ce que la fonction coût ait atteint son minimum.

### Évaluation du modèle sur la jauge d'entraînement

Le modèle de régression linéaire est entraîné sur le jeu de données d'entraînement pendant 100 000 époques. Dans la figure 1.7 on peut observer l'évolution décroissante de la fonction coût, jusqu'à atteindre une limite. Lors de l'entraînement il est important de calculer la valeur de la fonction coût sur les données d'entraînements et sur les données test. Cela permet de vérifier que notre modèle n'est pas en sur-apprentissage. En effet si la fonction coût des données de test commence à croître alors que celle des données d'entraînement suit une allure régulière, cela signifie que le modèle s'est excessivement entraîné à représenter le jeu d'entraînement et n'est plus capable de prédire d'autres données.

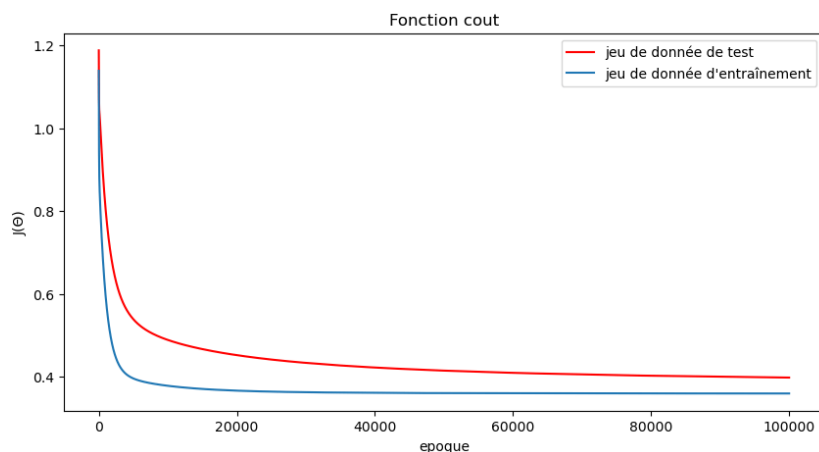


Figure 1.7: Évolution de la fonction coût

La figure 1.8 nous montre les résultats obtenus par notre modèle (trait plein bleu) en comparaison du modèle simpliste calculant la moyenne (tiret orange). La courbe rouge est la valeur idéale de jauge. On remarque que notre modèle linéaire n'est pas plus efficace que la valeur moyenne et semble faire les mêmes erreurs de prédiction. En appliquant notre modèle sur toutes les jauges et en

comparant sa précision à la moyenne des radars , on se rend compte que le modèle est 15% moins précis que la moyenne des radars. Sur la figure 1.8 on peut observer certain défaut de notre modèle. Le modèle de régression linéaire n'arrive pas a représenter la valeur minimum de 0 et semble avoir une valeur constante aux alentours de 0.1mm de pluie. De plus, les fort pics de pluies sont sous-évaluer par la régression linéaire. Ces problèmes proviennent du fait que notre modèle est une fonction linéaire tentant de représenter un problème de plus haut degré polynomiale.

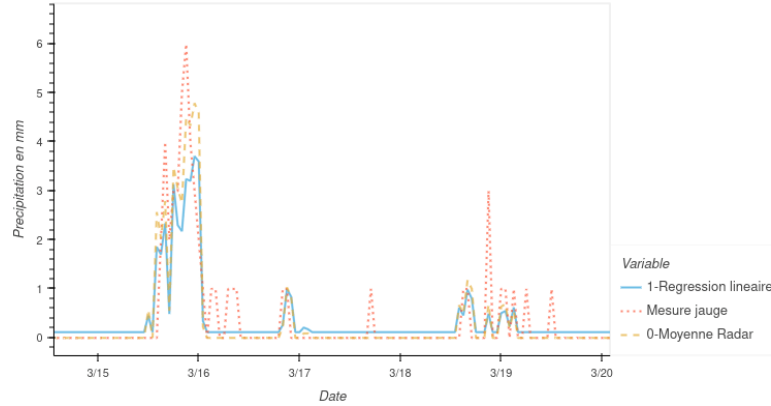


Figure 1.8: Mesure moyenne des radars comparé à la régression linéaire

### Évaluation du modèle sur toute les jauges

Le modèle à été entraînée sur une jauge. Le but du modèle est d'être généralisé à n'importe quelle mesure de jauge et de radar. Nous évaluons donc les performances de notre modèle sur des jauge qui lui sont totalement inconnus. Si les données d'entraînements correspondent à un zone géographique plane, lorsque des données provenant d'une zone montagneuse sont inséré il y a de forte chance que les résultats soient mauvais. La figure ci dessous compare le modèle lorsqu'il est entraînée dans une zone montagneuse au bord de la mer 1.9b, et un entraînement sur une zone retranchée dans les terres 1.9a. Les points en vert représentent les jauges ou le modèle fait une meilleur prédiction que une simple moyenne. Les points en bleu et rouge représentent respectivement une prédiction à peu près égale à la moyenne, ou une meilleur prédiction de la part de la moyenne. On peut voir que en fonction de la position des données d'entraînements les prédictions changent drastiquement. Lorsque le modèle est entraîné dans les terres, il n'arrive pas à prédire les précipitations dans le Sud-Est, zone montagneuse au bord de la cote. Au contraire lorsque le modèle est entraînée dans le Sud-Est, quelques bonne prédictions sont réalisées dans cette zone montagneuse. Cependant le reste des prédictions sont fausse. On dit que les covariables sont décalées : la distribution du jeu de donnée d'entraînement

ne correspond pas au reste des données de test. La zone au Sud-Est a tendance à avoir un temps calme avec de fort pic de pluie, ce qui se démarque des autres climats. Notre modèle ne performe donc pas bien car la distribution des données d'entraînement et de test sont différentes. Pour que notre programme se généralise, nous allons entraîner le modèle sur toute les données de jauge.

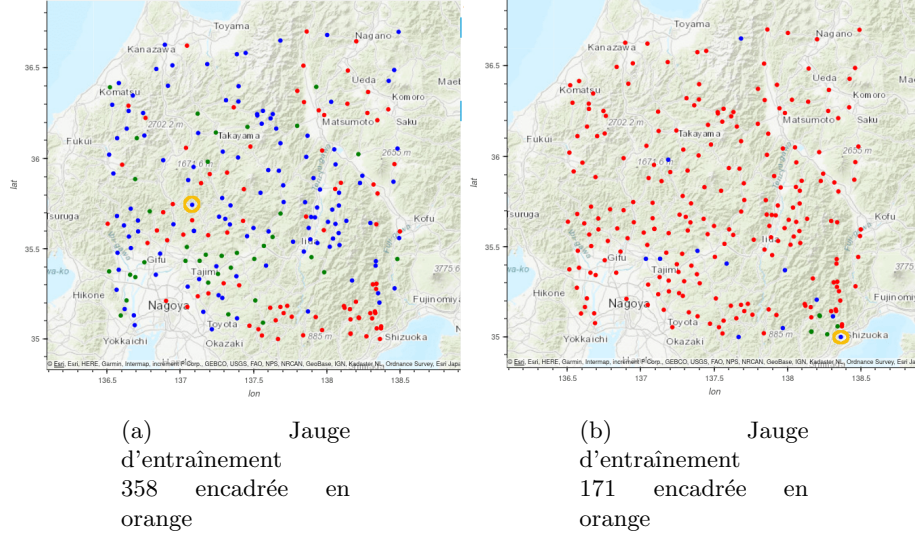


Figure 1.9: Performance du modèle. rouge : la moyenne performe mieux que le modèle. Vert : le modèle performe mieux. Bleu : moyenne et modèle ont à peu près les mêmes performances.

### 1.2.3 Entraînement sur toute les jauges

Jusqu'à maintenant le modèle de régression linéaire ne s'entraînait que sur les données d'une jauge. Nous allons maintenant l'entraîner sur tout les emplacements de jauges, ce qui permet au modèle de s'exercer sur des données plus variées et représentative de la réalité. En entraînant notre modèle de cette manière sa précision dépasse celle de la valeur moyenne de peu. La figure 1.10 nous montre les performances du modèle par rapport à la moyenne. Contrairement aux résultats de la section précédente le modèle arrive à égaler ou dépasser la moyenne sur la plupart des points d'évaluations. On remarque cependant que la zone située au Sud-Est n'obtient toujours pas de bon résultats avec le modèle. En comparant les mesures dans cette zone avec les autres jauge, on voit que les pics de pluie sont plus puissants. En regardant la quantité de pluie tombée sur les 3 mois d'observation, on observe un surplus de 20% dans la zone sud-est.

En analysant la fonction linéaire que l'on a entraînée, on remarque aussi que les coefficients  $\Theta$  semblent tendre vers des valeurs proche de  $1/6$  1.11. Le modèle linéaire calcul donc la moyenne des valeurs radars pour chaque heures,

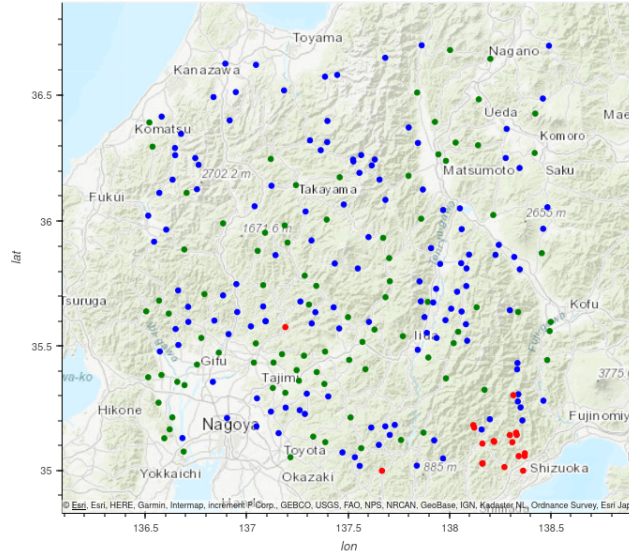


Figure 1.10: Performance du modèle entraîné sur toutes les jauges. rouge : la moyenne performe mieux que le modèle. Vert : le modèle performe mieux. Bleu : moyenne et modèle ont à peu près les mêmes performances.

cela confirme que la fonction moyenne est la fonction linéaire la plus efficace pour ce problème. En réalité les valeurs  $\Theta$  ne tendent pas exactement vers  $1/6$ , mais vers une valeur légèrement plus faible. De plus le modèle semble donner plus d'importance au coefficient représentant la valeurs du radar à 50 minutes. Notre modèle est donc une valeur moyenne atténué donnant plus d'importance à un coefficient.

La régression linéaire est un premier pas vers des modèles plus compliquées du machine learning. Grâce à ce modèle nous pouvons avoir une première estimations des précipitations sur le territoire japonais. Un tel modèle est très simple mais nous a permis de démontrer l'importance de certaine caractéristique permettant de prédire les précipitations, comme la prise en compte des dimensions temporels lorsque l'on estime une précipitation. Cependant ce modèle n'est pas assez complexe pour estimer une système non-linéaire.

### 1.3 Importance des caractéristiques temporelles pour estimer les précipitations

Les données utilisées pour le modèle linéaire sont les mesures radars sur l'heure que l'on estime. Ces mesures sont nommées caractéristiques ("features" en anglais). Pour améliorer notre modèle l'ajout de caractéristique est nécessaire.

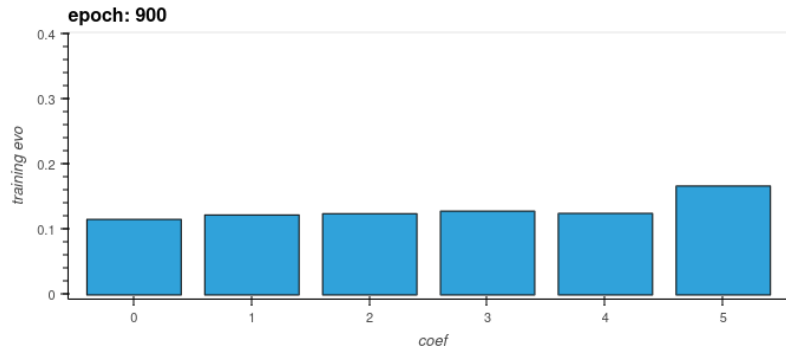


Figure 1.11: Coefficients finaux obtenus par la régression linéaire. Les coefficients représentent les mesures du radar toute les 10 minutes durant l’heure estimée

Par exemple l’ajout d’informations spatiales ou de données atmosphériques permettent d’améliorer les performance d’un modèle. Dans cette section nous étudions l’importance de caractéristique temporels.

### 1.3.1 Modèle linéaire avec une fenêtre temporelle

Pour estimer la pluie à une certaine heure, les modèles précédents utilisent les 6 données radars mesurées sur cette heure. Dans le modèle suivant l’objectif est d’étendre les données temporel que nous utilisons. Pour une estimation a une certaine heure  $H$  nous utilisons donc les données radars liées à l’heure  $H$ , mais aussi  $H+1$  et  $H-1$ . Dans ce cas nous avons une fenêtre temporelle de taille 3.

Une fois de plus, la précision de notre modèle s’améliore légèrement. Sur la figure 1.12 on peut observer l’évolution de la fonction coût en fonction de la fenêtre temporelle choisis. Il est intéressant de remarquer que lorsque la fenêtre temporelle devient trop large, les performances du modèle diminue. Le surplus de données empêche les plus importantes de se démarquer. En effet toute les caractéristiques temporelles ne se valent pas. Par exemple les valeurs de pluie ayant eu lieu en  $H+100$  minutes sont beaucoup moins importante que celle ayant lieu en  $H$ . La figure 1.13 montre les coefficients obtenus par la régression linéaire. Les coefficients sont ordonnés de  $-7H$  à  $+7H$  par tranche de 10 minutes. Par exemple, le coefficient  $-7$  représente 70 minute avant la mesure de la jauge en  $H$ . En s’entraînant le modèle a de lui même atténuer les coefficients représentant des heures éloignés de l’heure  $H$ . On voit sur cette figure que les coefficients les plus important sont situés entre  $-10$  et  $10$  (soit  $H-100$  minutes et  $H+100$  minutes), ce qui permet au mesures temporel adjacentes à notre heure  $H$  d’impacter fortement le résultat de la régression.

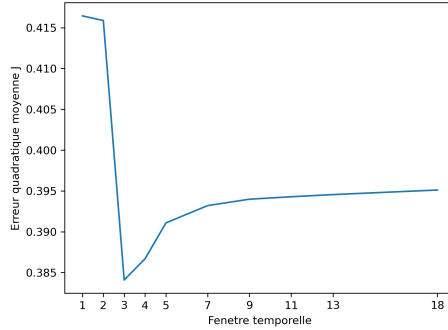


Figure 1.12: Évolution de la fonction coût en fonction de la taille de la fenêtre temporelle

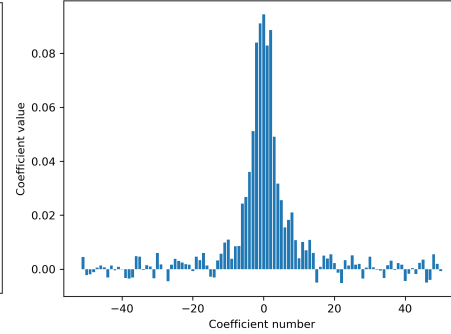


Figure 1.13: Coefficient du modèle pour une fenêtre de taille  $-7H, +7H$

### 1.3.2 Analyse des caractéristique avant l'entraînement

Précédemment nous avons analysé l'importance de nos caractéristiques après avoir entraîné notre modèle. Dans le cas où l'entraînement de notre modèle est long, obtenir des informations sur nos caractéristiques en amont de l'entraînement est nécessaire. Pour cela il existe différentes méthodes d'analyse des caractéristiques tel que la corrélation par rapport à la cible, et la PCA (Principal Component Analysis).

#### Corrélation avec la cible

Le coefficient de corrélation est une mesure quantifiant la relation linéaire entre 2 variables. Dans notre cas ce coefficient est utilisé pour calculer la corrélation entre une caractéristique et la cible de notre modèle. Pour chaque mesure du radar au temps  $H$ ,  $H+10$  minutes,  $H-10$  minutes etc.. nous nous observons la corrélation avec les mesures de la jauge au temps  $H$ . La figure 1.14 montre les coefficients de corrélations obtenus. Un coefficient de corrélation est compris entre -1 et 1. Plus la valeur est proche de 1 (ou -1), plus les variables sont positivement (ou négativement) corrélées, c'est à dire qu'elles évoluent de la même manière. On peut voir que les caractéristiques temporelles au temps  $H$  à  $H+6$  ont les coefficients de corrélations les plus élevés. Plus l'on s'éloigne des caractéristiques au temps  $H$ , moins les variables sont linéairement corrélées. Pour réduire le nombre de caractéristiques que l'on choisit il suffit de définir un seuil de corrélation. Toute caractéristique ayant un coefficient plus faible que le seuil sont rejetées.

#### PCA

Principal Component Analysis est une méthode permettant de condenser l'ensemble des caractéristiques envoyées en input du modèle. Dans notre cas l'objectif est de

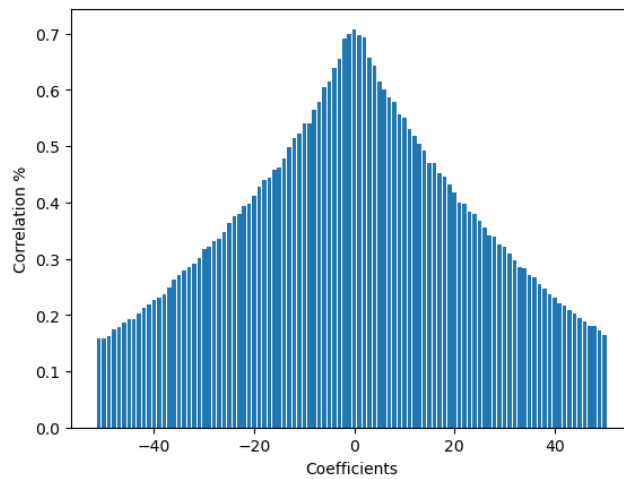


Figure 1.14: Coefficients de corrélations entre les caractéristiques et la cible du modèle

réduire les caractéristiques choisies en un ensemble de 6 nouvelles caractéristiques qui seront une combinaison linéaire des précédentes. On peut voir la PCA comme une réorganisation de nos caractéristiques afin de les ordonner en fonction de la quantité d'information qu'elles fournissent. Contrairement à l'exemple précédent ou l'

- caractéristiques temporelles des précipitations (coef, pca, corrélation)

Lors de l'étude PCA : consisterait à sélectionner les caractéristiques les plus importantes de notre jeu de données, et utiliser ces caractéristiques en entrée de la régression

- limitation modèle linéaire Ajout de plus de caractéristique (les données voisines, un plus large temporel ...) impossible pour le modèle linéaire

- Limitation dans les features choisies pour le modèle On n'utilise pas les informations spatiales autour de notre point. C'est pour ça qu'on voit que notre modèle linéaire ne fait que copier les données radars : là où le radar mesure beaucoup de pluie, le modèle linéaire aussi. (Mettre une image où le radar et la prédiction monte, alors que la jauge ne mesure pas de pluie)



# Bibliography

- [1] *Japan Meteorological Agency*. URL: <https://www.jma.go.jp/jma/en/copyright.html>.
- [2] Machine Learnia. *DESCENTE DE GRADIENT (GRADIENT DESCENT) - ML#4*. 2019. URL: [https://www.youtube.com/watch?v=rcl\\_YRyoLIY](https://www.youtube.com/watch?v=rcl_YRyoLIY).
- [3] Kazuhiko Nagata. “Quantitative precipitation estimation and quantitative precipitation forecasting by the Japan Meteorological Agency”. In: *RSMC Tokyo–Typhoon Center Technical Review* 13 (2011), pp. 37–50.
- [4] Ryota Nakamura and Yukihiro Shimatani. “Extreme-flood control operation of dams in Japan”. In: *Journal of Hydrology: Regional Studies* 35 (2021), p. 100821. ISSN: 2214-5818. DOI: <https://doi.org/10.1016/j.ejrh.2021.100821>. URL: <https://www.sciencedirect.com/science/article/pii/S2214581821000501>.
- [5] Satoru Oishi, Toshihiko Tahara, and Mariko Ogawa. “Study on optimization of the operation of dams using ensemble prediction and a distributed rainfall-runoff model”. In: *HIC*. 2018, pp. 1584–1588.