Rubén García García y Carlos Torregrosa Alcayde

# XAI3

## Index

## Introduction

This report explores the practical application of Partial Dependence Plots (PDPs) in uncovering the relationship between input features and model predictions, focusing on two domains: bike rentals and house prices. By leveraging PDPs, we dissect the influence of key variables on model outcomes, providing valuable and actionable insights for decision-making. Through concise analyses, we unravel the intricate interplay of factors driving predictions, empowering stakeholders with practical insights to optimize services and investments.

## 1.- One dimensional Partial Dependence Plot.

This report provides a comprehensive overview of the application of partial dependency plots (PDPs) to understand the impact of different features on model predictions. PDPs are robust, model-agnostic tools that demonstrate the marginal effect of one or two features on the predicted outcome of a machine learning model. The exercises involved predicting bike rentals using a random forest model and analyzing the effects of various features using PDPs. All work is meticulously version-controlled using Git and securely backed up on GitHub. **https://github.com/Carlosta1177/git-blog-demo**

To begin with, we fitted a random forest model to predict bike rentals (cnt) and used PDPs to visualize the relationships between the predicted bike counts and four features: days since 2011, temperature, humidity, and wind speed. The data preparation process involved several steps. First, the season variable was one-hot encoded to convert it into a format suitable for the model. Binary features were created for different weather conditions, specifically MISTY and RAIN, to capture the effect of varying weather on bike rentals. Additionally, temperature, humidity, and wind speed were denormalized to their original scales for more intuitive interpretation. A new feature, days_since_2011, was created to represent the days elapsed since January 1, 2011, providing a continuous time variable.

Irrelevant columns, such as instant, dteday, casual, and registered, were removed from the dataset to focus on the relevant features for prediction.
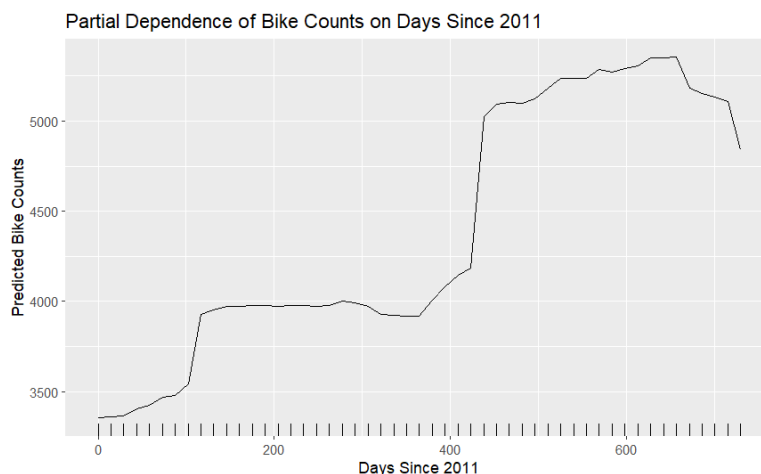
A random forest model was then trained on the prepared dataset to predict the bike rental counts. The importance of each feature was examined to ensure they contributed significantly to the model's predictions:

```
                   %IncMSE IncNodePurity
season           16.7560100     166773082
yr               19.6794392     288786685
mnth             17.3044800     101257677
holiday          -0.2549253       7280538
weekday          11.1171496      33799374
workingday        6.7933673      10979572
weathersit       16.3732829      59015870
temp             24.4014813     452975282
atemp            26.3243066     411903273
hum              30.2163698     119399981
windspeed        13.5049712      69304955
season_2          7.9149320       4968881
season_3          6.3379873      18539617
season_4         11.8441083      26625202
MISTY            14.6034254      15989124
RAIN             14.5669443      30700891
days_since_2011  38.1770273     872209997
```
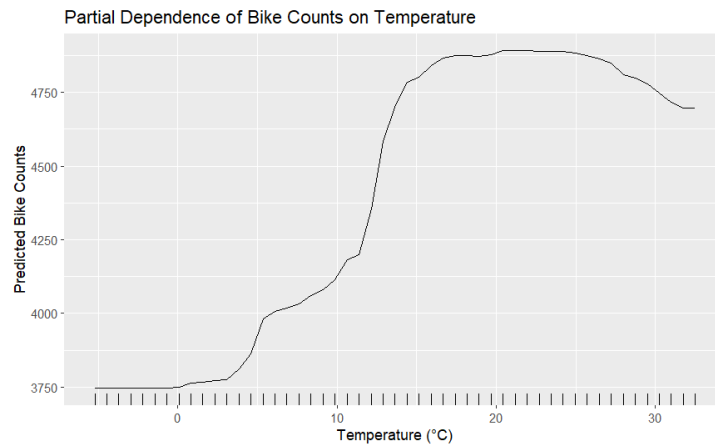
As we saw on the last lab, "days_since_2011" is the most important feature. This metrics mean that this feature is the most important for accurate predictions (%IncMSE), and plays a significant role in splitting the data effectively in the trees (IncNodePurity).

Following the model training, PDPs were generated for the features: days since 2011, temperature, humidity, and wind speed. These PDPs were visualized using ggplot2 to provide clear and insightful plots.
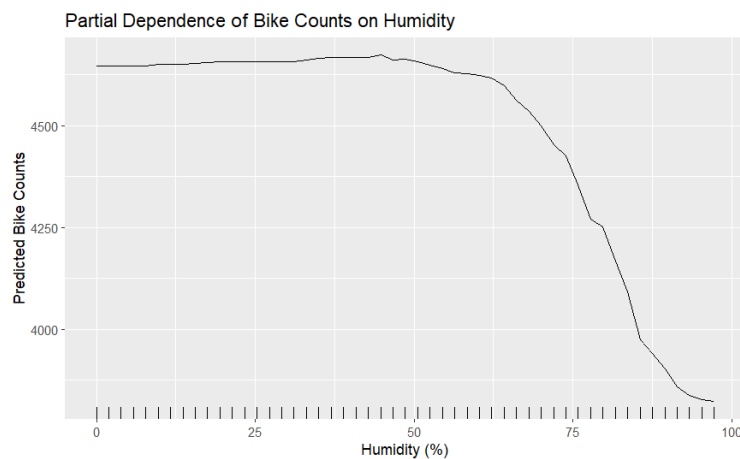
The PDP for days since 2011 shows a significant increase in bike rentals between approximately day 100 and day 150, followed by a steady increase up to around day 450. This pattern suggests a growing trend in bike usage over time, which could be due to factors such as increased availability of bikes, improved infrastructure, or a growing preference for biking.
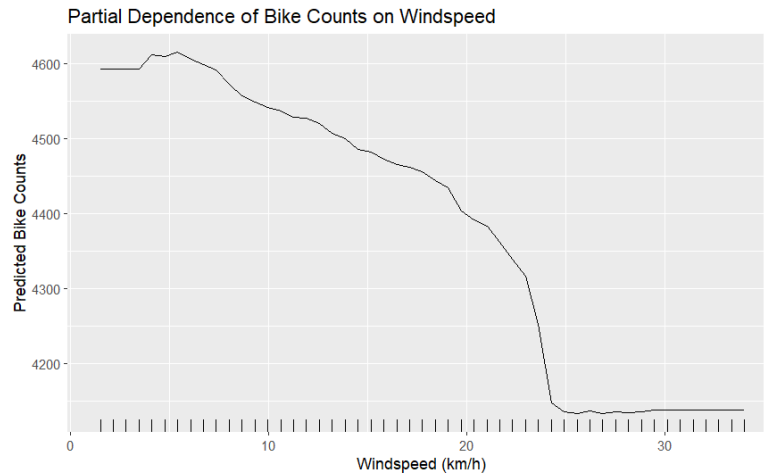
 The PDP for temperature reveals that bike rentals increase with temperature up to around 20°C, beyond which the increase in rentals slows down and eventually declines slightly after 25°C. This indicates that moderate temperatures are more favorable for biking, while extremely high temperatures might not significantly boost rentals due to potential discomfort or safety concerns.



The PDP for humidity shows an inverse relationship between humidity and bike rentals. Higher humidity levels correlate with fewer bike rentals, likely due to the discomfort associated with humid conditions. The steep decline in rentals at humidity levels above 75% indicates that very high humidity is a strong deterrent for bike usage.



Finally, the PDP for wind speed indicates a negative relationship between wind speed and bike rentals. The predicted bike counts start at around 4600 and gradually decrease as wind speed increases. At wind speeds of approximately 25 km/h, there is a significant drop in predicted bike counts to around 4100, where the line then advances horizontally, this means that from 25km/h, the number of predicted bikes will be constant and low. Higher wind speeds likely deter people from renting bikes due to the increased physical effort required to ride in windy conditions.

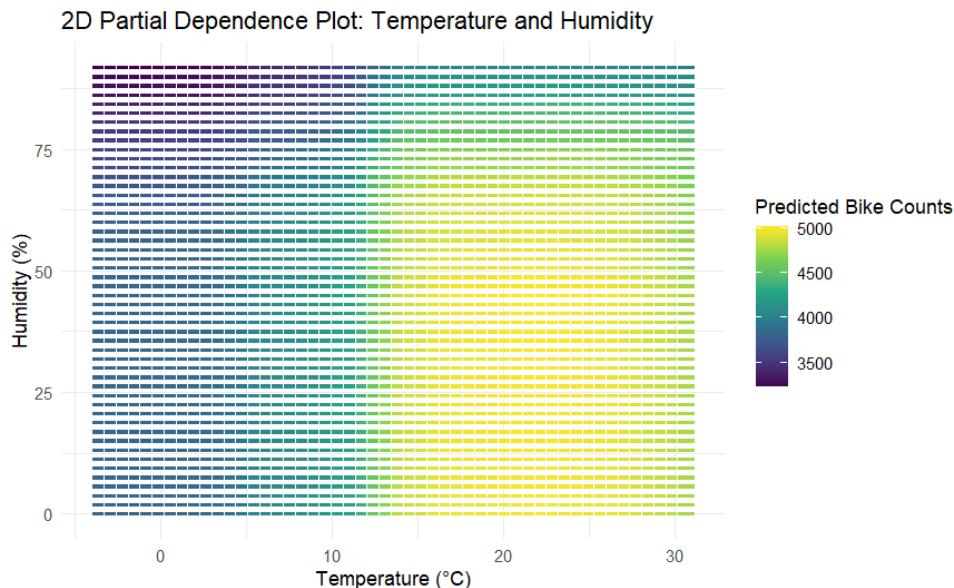Partial Dependence of Bike Counts on Windspeed

## 2.- Bidimensional Partial Dependency Plot.

This report continues the exploration of Partial Dependency Plots (PDPs) to understand the influence of different features on model predictions. In this section, we focus on generating a 2D PDP to visualize the joint effect of two features—temperature and humidity—on the predicted number of bike rentals. This visualization helps to uncover the interaction effects between these features, providing deeper insights into the model's behavior.

We used the same dataset and random forest model from the previous exercise, which involved predicting bike rental counts (cnt). The data preparation steps included one-hot encoding the season variable, creating binary features for different weather conditions (MISTY and RAIN), denormalizing temperature, humidity, and wind speed to their original scales, and creating a days_since_2011 feature. Irrelevant columns were removed to focus on the relevant features for prediction.

Due to the size of the dataset, we extracted a random sample to make the computation of the 2D PDP more efficient. Specifically, a subset of 10% of the data was sampled to generate the 2D partial dependence data. The partial function from the pdp package was used to compute the PDP for temperature and humidity with a grid resolution of 50. The resulting data was converted to a data frame for visualization.

The 2D Partial Dependency Plot shows the combined effect of temperature and humidity on the predicted number of bike rentals. The color gradient represents the predicted bike counts, with lighter colors indicating higher bike rentals and darker colors indicating lower rentals.



From the plot, we observe that the highest predicted bike counts are observed at temperatures around 20-25°C with moderate humidity levels (around 40-60%), but also with low humidity. This indicates that these conditions are most favorable for bike rentals, likely because the temperature is comfortable and the humidity is not too high, making biking enjoyable. In contrast, the predicted bike counts are significantly lower in conditions of low temperature (below 10°C) and high humidity (above 75%). These conditions are less favorable for biking, likely due to the discomfort and potential hazards associated with cold and humid weather.

At high temperatures (above 30°C), even with low humidity, the predicted bike counts start to decline. This suggests that while low humidity might be comfortable, very high temperatures can deter people from renting bikes due to the discomfort and potential health risks associated with heat. Overall, the interaction between temperature and humidity shows a clear pattern where moderate temperatures and humidity levels are most conducive to higher bike rentals. Extreme conditions, whether hot or cold, or very high humidity, tend to reduce the predicted bike counts.

The 2D PDP provides valuable insights into how temperature and humidity together influence bike rentals, highlighting the conditions that are most and least favorable for biking. Understanding these interactions helps in interpreting model predictions and can inform decision-making processes for bike rental services, such as optimal conditions for promotions or resource allocation.

# 3.- PDP to explain the price of a house.

In this exercise, a Random Forest model was applied to predict house prices using data from the `kc_house_data.csv` database. The attributes used for the prediction include the number of bedrooms, the number of bathrooms, the living area in square feet (`sqft_living`), the lot area in square feet (`sqft_lot`), the number of floors (`floors`), and the year the house was built (`yr_built`). To illustrate how these features influence the predicted house price, partial dependence plots (PDPs) were generated.
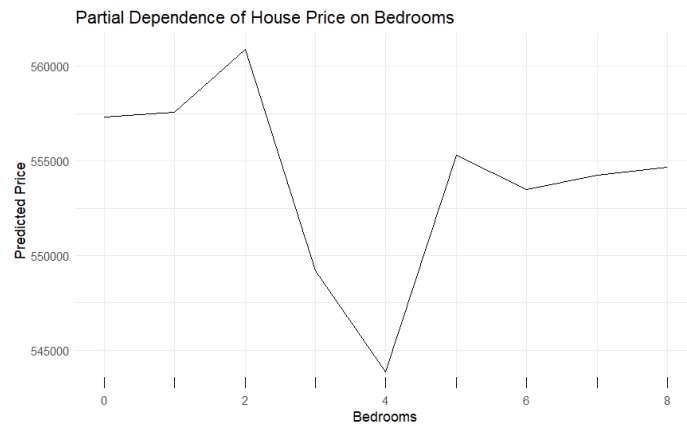
First, the dataset was loaded, and the relevant columns were selected. To facilitate processing and the generation of PDPs, a random sample of 10% of the data was taken. Subsequently, a Random Forest model was fitted to this sample to predict house prices based on the selected features.

Once the model was fitted, the importance of the features in predicting the house price was evaluated. The results of the feature importance, measured in terms of %IncMSE and IncNodePurity, show that the living area (`sqft_living`) is the most influential feature, followed by the year built (`yr_built`), the number of bathrooms, the lot area (`sqft_lot`), the number of floors, and finally, the number of bedrooms. These results indicate that features related to the size and modernity of the house have a greater impact on the predicted price.
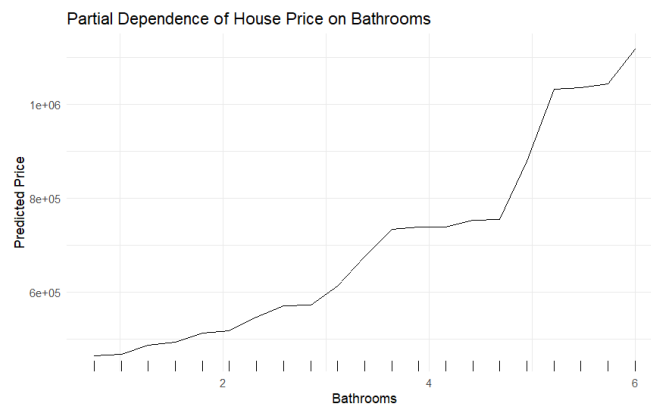
```
             %IncMSE IncNodePurity
bedrooms     7.896915  1.504339e+13
bathrooms   24.232309  5.228457e+13
sqft_living 53.371226  1.204133e+14
sqft_lot    14.421364  3.136253e+13
floors      14.917466  8.987040e+12
yr_built    37.126567  2.969178e+13
```

To visualize how these individual features affect the price prediction, PDPs were generated for each feature.
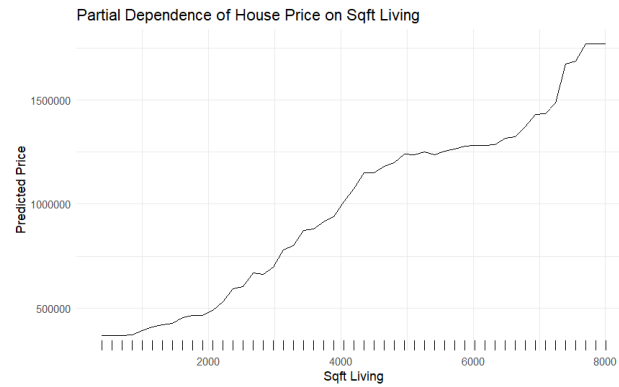
The first plot shows the partial dependence of house price on the number of bedrooms. A less direct relationship is observed, where initially, increasing the number of bedrooms from 1 to 2 significantly raises the predicted price. However, from 2 bedrooms onward, the price decreases sharply, reaching its lowest point around 4 bedrooms, after which it stabilizes and shows a slight upward trend. This fluctuation might reflect a saturation point where additional bedrooms do not proportionally increase the house value and might suggest smaller living spaces per room or less desirable features.

Partial Dependence of House Price on Bedrooms

The second plot shows the partial dependence of house price on the number of bathrooms. Here, a more consistent positive relationship is observed. As the number of bathrooms increases, so does the predicted house price. This relationship is almost linear up to about 4 bathrooms, after which the increase in predicted prices becomes more pronounced. This suggests that additional bathrooms significantly add to the house's value, likely due to increased comfort and perceived luxury.



Partial Dependence of House Price on Bathrooms

The third plot shows the partial dependence of house price on the living area (`sqft_living`). A strong positive relationship is observed, where the predicted price steadily increases as the living area increases. This behavior indicates that larger living spaces are highly valued and significantly contribute to the total house price. Beyond certain values, the increase in predicted price becomes more notable, suggesting that more spacious houses are associated with higher market values.

Partial Dependence of House Price on Sqft Living

The fourth plot shows the partial dependence of house price on the number of floors. A positive relationship is observed, where the predicted price increases as the number of floors increases from 1 to 2, and this trend continues with a more pronounced increase between 2 and 3 floors. However, beyond 3 floors, the predicted prices stabilize, suggesting that additional floors after this point do not substantially increase the house's value. This pattern may reflect market preferences for houses with more than 2 floors, possibly due to convenience, space utilization, and architectural appeal considerations.



Partial Dependence of House Price on Floors