

Exploring Tools for Interpretable Machine Learning

Dr. Juan Ordúz

PyData Global 2021

Outline

Introduction

Data Set ([6])

Models Fit ([9])

Model Explainability ([6], [4])

- Model Specific

 - Beta Coefficients and Weight Effects

 - Tree ensembles

- Model Agnostic

 - PDP and ICE Plots

 - Permutation Importance

 - SHAP

References

Introduction

Aim and Scope of the Talk

What? In this talk we want to test various ways of getting a better understanding on how machine learning (ML) models generate predictions and how features interact with each other. This is in general not straight forward and key components are

- ▶ Domain knowledge on the problem.
- ▶ Understanding on the input data.
- ▶ Understanding the logic behind the ML algorithms.

How? We are going to work out a concrete example.

Introduction

Aim and Scope of the Talk

What? In this talk we want to test various ways of getting a better understanding on how machine learning (ML) models generate predictions and how features interact with each other. This is in general not straight forward and key components are

- ▶ Domain knowledge on the problem.
- ▶ Understanding on the input data.
- ▶ Understanding the logic behind the ML algorithms.

How? We are going to work out a concrete example.

References

This talk is based on this blog post ([9]), which itself is based on these two amazing references:

- ▶ Interpretable Machine Learning, A Guide for Making Black Box Models Explainable by Christoph Molnar ([6])
- ▶ Interpretable Machine Learning with Python by Serg Masís ([4])

Interpretable ML \neq Causality

See [7] for a Brief History of Interpretable ML and comments on Causality

Note that the methods discussed in this notebook are not related with causality.

References

- ▶ Be Careful When Interpreting Predictive Models in Search of Causal Insights by Scott Lundberg (one of the core developers of SHAP)
- ▶ Statistical Rethinking, A Bayesian Course with Examples in R and Stan by Richard McElreath ([5])
- ▶ Causal Inference: The Mixtape by Scott Cunningham ([2])

Interpretable ML \neq Causality

See [7] for a Brief History of Interpretable ML and comments on Causality

Note that the methods discussed in this notebook are not related with causality.

References

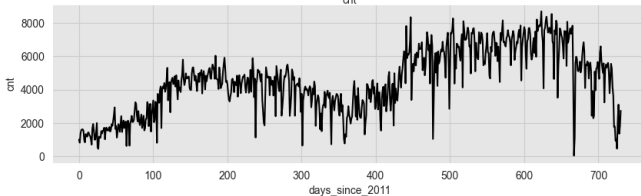
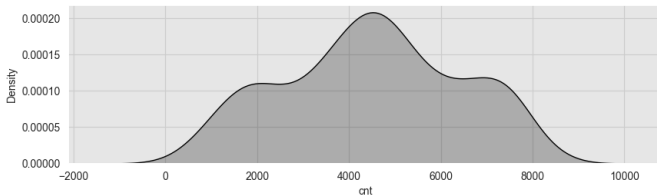
- ▶ Be Careful When Interpreting Predictive Models in Search of Causal Insights by Scott Lundberg (one of the core developers of SHAP)
- ▶ Statistical Rethinking, A Bayesian Course with Examples in R and Stan by Richard McElreath ([5])
- ▶ Causal Inference: The Mixtape by Scott Cunningham ([2])

Remark: The article General Pitfalls of Model-Agnostic Interpretation Methods for Machine Learning Models ([8]) is highly recommended to understand the challenges, limitations and recommendations for some of the model-agnostic methods discussed in this talk.

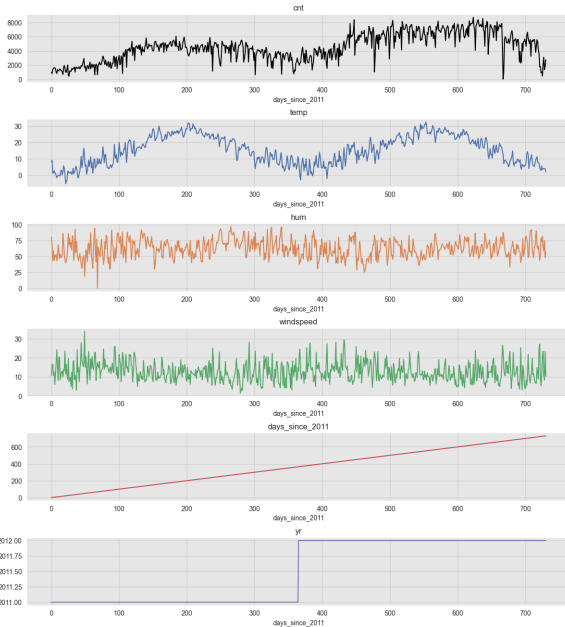
Target Variable - cnt: Daily Bike Rents

	season	yr	mnth	holiday	weekday	workingday	weathersit	temp	hum	windspeed	cnt	days_since_2011
0	SPRING	2011	JAN	NO HOLIDAY	SAT	NO WORKING DAY	MISTY	8.175849	80.5833	10.749882	985	0
1	SPRING	2011	JAN	NO HOLIDAY	SUN	NO WORKING DAY	MISTY	9.083466	69.6087	16.652113	801	1
2	SPRING	2011	JAN	NO HOLIDAY	MON	WORKING DAY	GOOD	1.229108	43.7273	16.636703	1349	2
3	SPRING	2011	JAN	NO HOLIDAY	TUE	WORKING DAY	GOOD	1.400000	59.0435	10.739832	1562	3
4	SPRING	2011	JAN	NO HOLIDAY	WED	WORKING DAY	GOOD	2.666979	43.6957	12.522300	1600	4

cnt: Target Variable

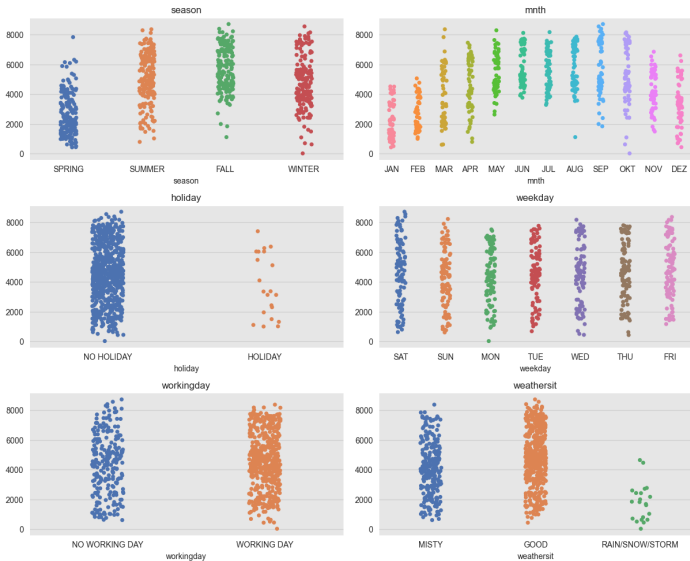


Continuous Regressors

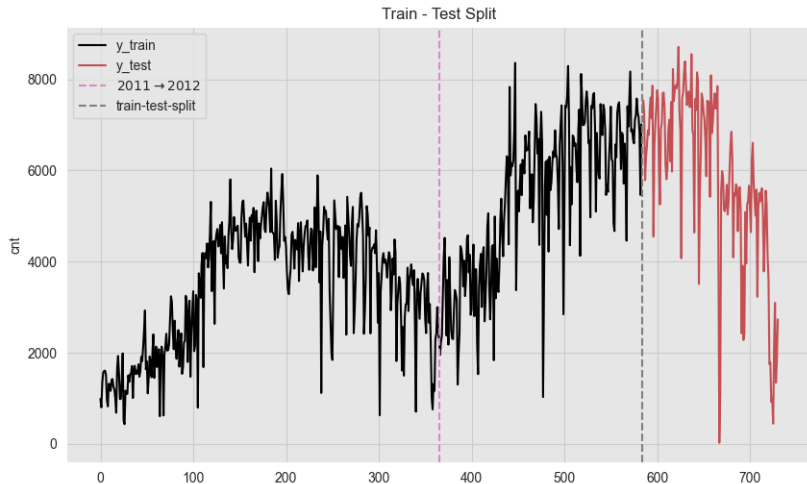


Categorical Regressors

cnt distribution over categorical_features



Train-Test Split



Models

Two model flavours

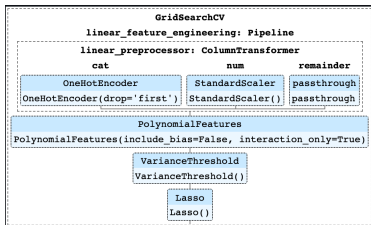


Figure 1: Lasso model with second order polynomial interactions ([10]).

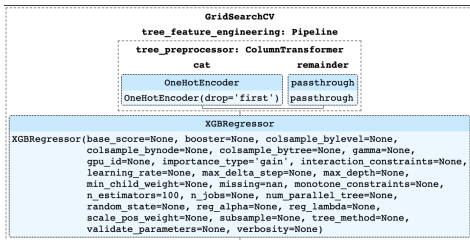
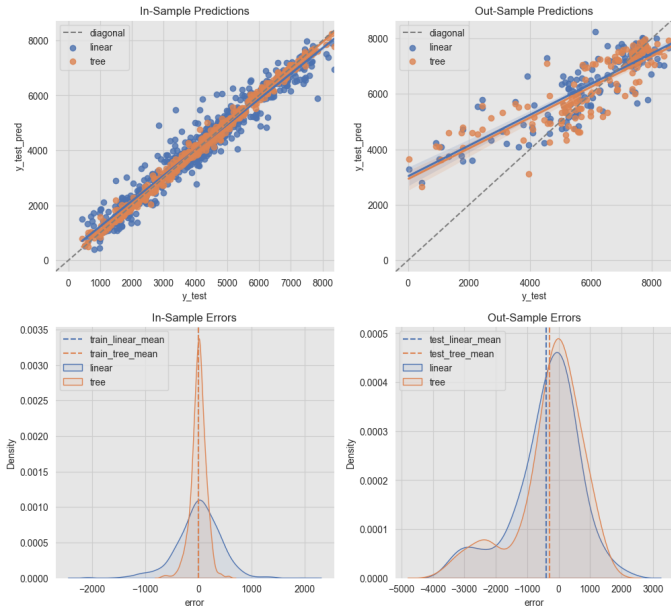
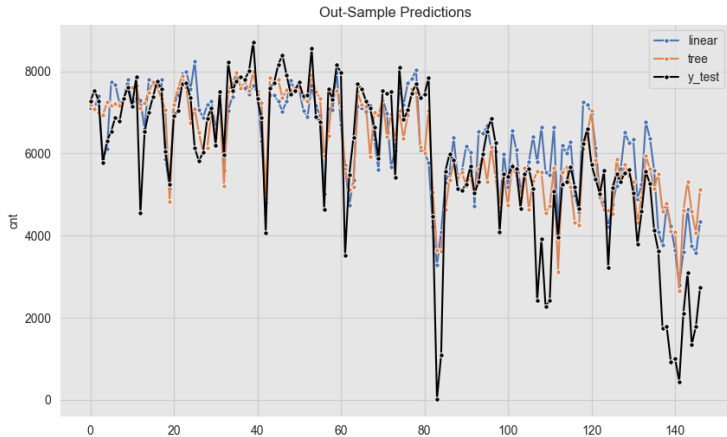


Figure 2: XGBoost regression model ([1]).

Out of sample performance - Errors Distribution



Out of sample performance - Predictions



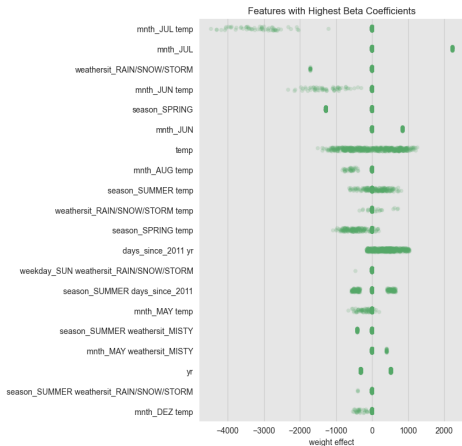
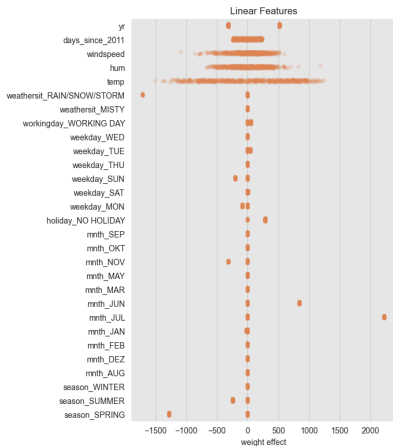
β coefficients

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \varepsilon, \quad \text{where } \varepsilon \sim N(0, \sigma^2)$$

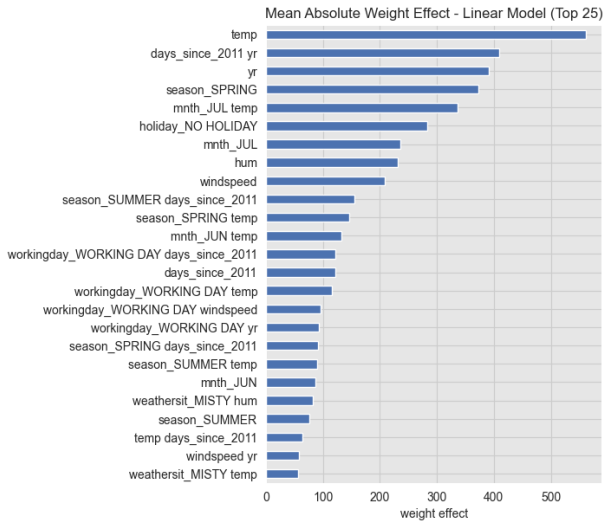
	linear_features	coef_	abs_coef_
0	mnth_JUL temp	-2305.096894	2305.096894
1	mnth_JUL	2227.672335	2227.672335
2	weathersit_RAIN/SNOW/STORM	-1710.469071	1710.469071
3	mnth_JUN temp	-1299.644413	1299.644413
4	season_SPRING	-1279.629779	1279.629779
5	mnth_JUN	845.229031	845.229031
6	temp	646.609622	646.609622
7	mnth_AUG temp	-523.011653	523.011653
8	season_SUMMER temp	489.319256	489.319256
9	weathersit_RAIN/SNOW/STORM temp	-482.660271	482.660271
10	season_SPRING temp	465.512410	465.512410
11	days_since_2011 yr	465.079169	465.079169
12	weekday_SUN weathersit_RAIN/SNOW/STORM	-462.286059	462.286059
13	season_SUMMER days_since_2011	454.137278	454.137278
14	mnth_MAY temp	-445.268148	445.268148
15	season_SUMMER weathersit_MISTY	-408.809531	408.809531
16	mnth_MAY weathersit_MISTY	404.790954	404.790954
17	yr	403.199142	403.199142
18	season_SUMMER weathersit_RAIN/SNOW/STORM	-394.157306	394.157306
19	mnth_DEZ temp	363.222114	363.222114

Weight Effects $\beta_i x_i$

Effect Weight Distribution



Weight Effects Importance $w_i = \frac{1}{n} \sum_{i=1}^n |\beta_i x_i|$



Weight Effects: Temperature (z-transform)

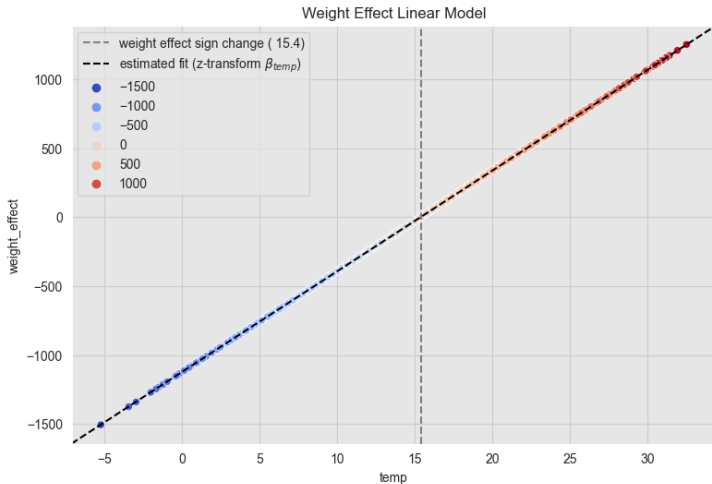


Figure 3: This plot just shows the effect of the linear term *temp* and not the interactions.

Weight Effects: Interactions

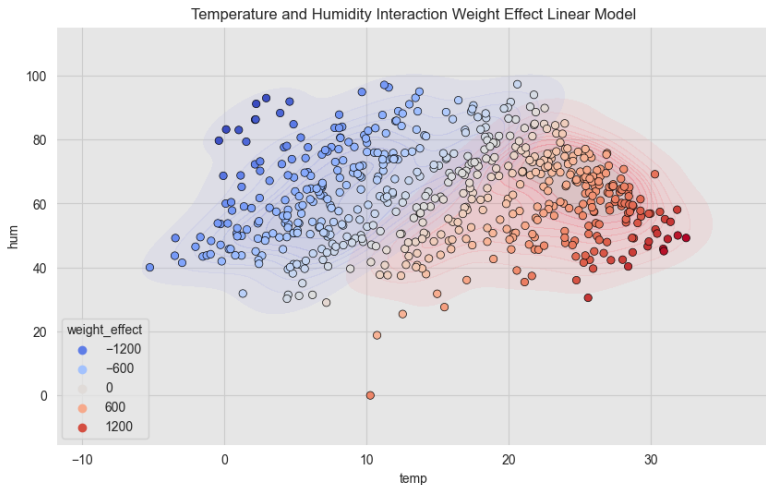


Figure 4: We can visualize the interaction between *temp* and *hum* by computing the total weight effect as

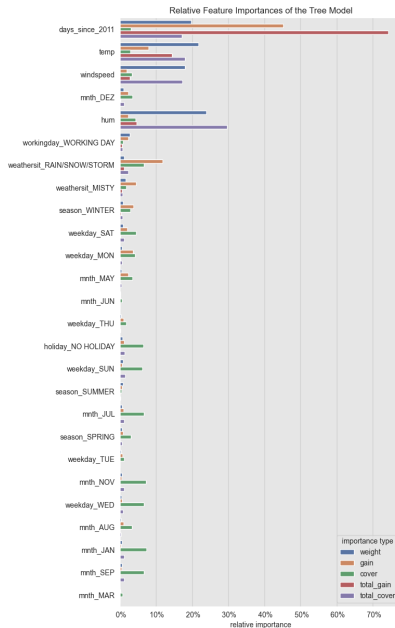
$$\beta_{temp}X_{temp} + \beta_{hum}X_{hum} + \beta_{temp \times hum}X_{temp}X_{hum}.$$

Explaining Individual Predictions



Figure 5: Weight effects of the linear model for observation 284. Left: All weight effects. Right: Weight effects of the linear terms.

Feature Importance Metrics: XGBoost



Partial Dependence Plot (PDP) & Individual Conditional Expectation (ICE) ([6, Section 5.1])

- ▶ Let x_S be the features for which the partial dependence function should be plotted and x_C be other features used in the machine learning model. One can estimate the *dependence function* as

$$\hat{f}_{x_S}(x_S) = \frac{1}{n} \sum_{i=1}^n \hat{f}(x_S, x_C^i)$$

where \hat{f} is the model prediction function, x_C^i are actual feature values (not in S) and n is the number points.

- ▶ Similar to a PDP, an individual conditional expectation (ICE) plot shows one line per instance. That is, for each instance in $\{(x_S^i, x_C^i)\}_{i=1}^n$, we plot \hat{f}_S as a function of x_S^i while leaving x_C^i fixed.

PDP & ICE Examples (1D)

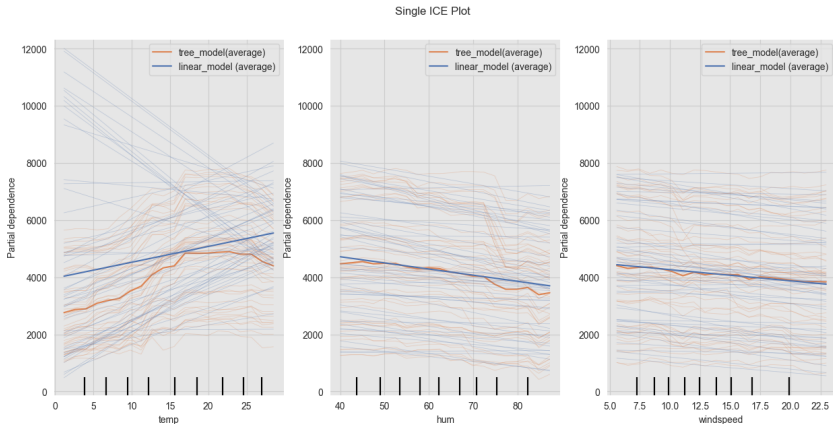
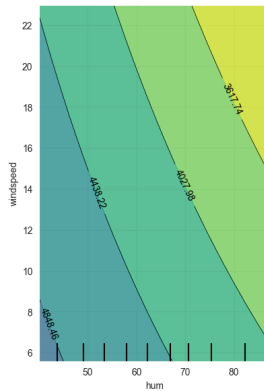
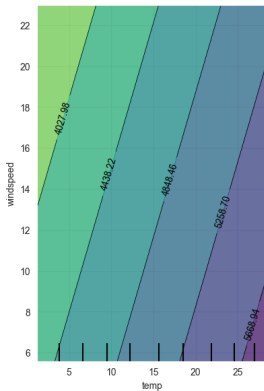
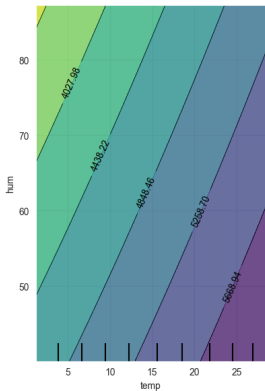


Figure 6: PDP & ICE plots for some numerical variables for the linear and XGBoost models.

PDP & ICE Examples (2D)

Pair ICE Plot - Linear Model



	linear_features	coef_	abs_coef_
6	temp	646.609622	646.609622
29	hum	-281.701209	281.701209
32	windspeed	-264.190697	264.190697
114	hum windspeed	-31.993914	31.993914
130	temp hum	15.127133	15.127133
155	temp windspeed	-0.000000	0.000000

Permutation Importance

Measures the increase in the prediction error of the model after we permuted the feature's values, which breaks the relationship between the feature and the true outcome ([6, Section 5.6]).

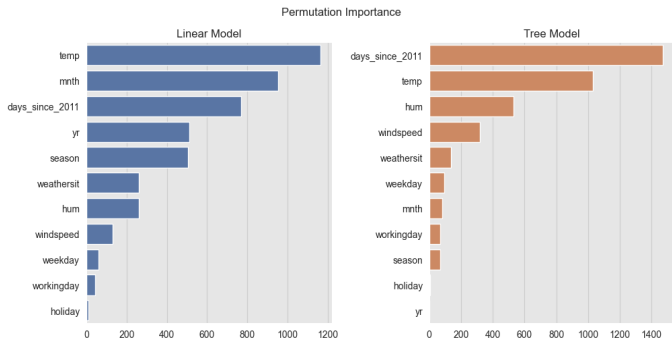


Figure 7: The permutation importance for these two models have *days_since_2011* and *temp* on their top 3 ranking, which partially explain the trend and seasonality components respectively (see [6, Figure 5.32]).

SAHP Values

Definition

For each data instance x

- ▶ Sample coalitions $z'_k \in \{0, 1\}^M$, where M , is the maximum coalition size.
- ▶ Get prediction for each z'_k . For features not in the coalition we replace their values with random samples from the dataset (background data).
- ▶ Compute the weight for each z'_k , with the SHAP kernel,

$$\pi_x(z') = \frac{(M-1)}{\binom{M}{|z'|} |z'| (M - |z'|)}$$

- ▶ Fit weighted linear model.
- ▶ Return Shapley values, i.e. the coefficients from the linear model.

SHAP Values

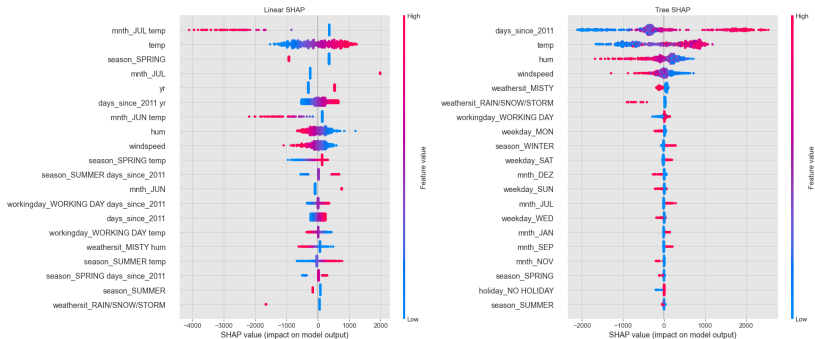
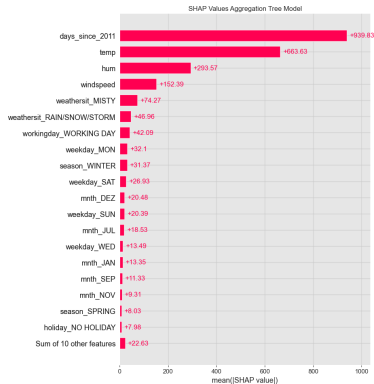
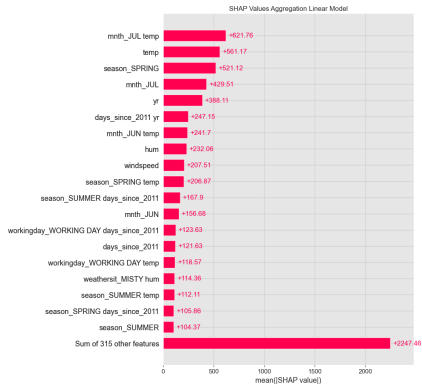


Figure 8: SHAP values per data instance. The x position of the dot is determined by the SHAP value of that feature, and dots "pile up" along each feature row to show density. Color is used to display the original value of a feature ([3]).

Mean Abs SHAP Values



SHAP Values: Temperature

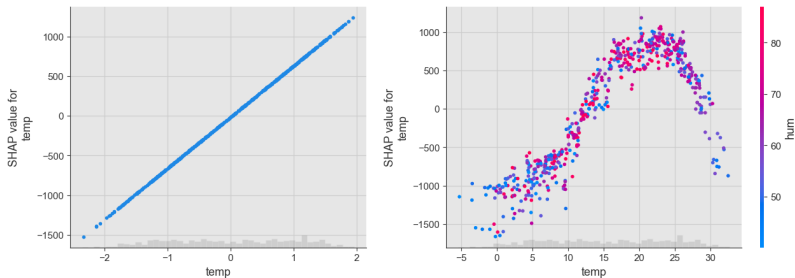


Figure 9: This figure shows the SHAP values as a function of temperature. Compare with Figure 6

SHAP Values: Observation 284

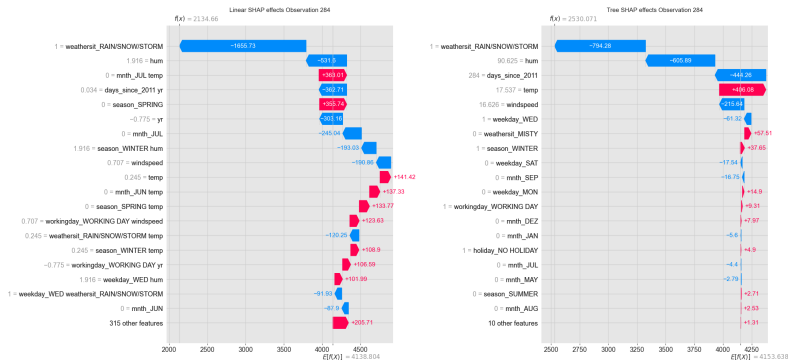


Figure 10: This *waterfall plot* shows how the SHAP values of each feature move the model output from our prior expectation under the background data distribution, to the final model prediction given the evidence of all the features ([3]). Compare with Figure 5.

References I

- [1] Tianqi Chen and Carlos Guestrin.
XGBoost: A scalable tree boosting system.
In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16, pages 785–794, New York, NY, USA, 2016. ACM.
- [2] Scott Cunningham.
Causal Inference: The Mixtape.
Yale University Press, 2021.
- [3] Scott M Lundberg and Su-In Lee.
A unified approach to interpreting model predictions.
In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 4765–4774. Curran Associates, Inc., 2017.
- [4] Serg Masís.
Interpretable Machine Learning with Python.
Packrat, 2021.
<https://github.com/PacktPublishing/Interpretable-Machine-Learning-with-Python>.

References II

- [5] Richard McElreath.
Statistical Rethinking: A Bayesian Course with Examples in R and Stan, 2nd Edition.
CRC Press, 2 edition, 2020.
- [6] Christoph Molnar.
Interpretable Machine Learning.
2019.
<https://christophm.github.io/interpretable-ml-book/>.
- [7] Christoph Molnar, Giuseppe Casalicchio, and Bernd Bischl.
Interpretable machine learning – a brief history, state-of-the-art and challenges,
2020.
- [8] Christoph Molnar, Gunnar Klippel, Julia Herbringer, Timo Freiesleben, Susanne Dandl, Christian A. Scholbeck, Giuseppe Casalicchio, Moritz Grosse-Wentrup, and Bernd Bischl.
General pitfalls of model-agnostic interpretation methods for machine learning models, 2021.
- [9] Juan Orduz.
Exploring tools for interpretable machine learning.
https://juanitorduz.github.io/interpretable_ml/, Jul 2021.

References III





- [10] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay.

Scikit-learn: Machine learning in Python.

Journal of Machine Learning Research, 12:2825–2830, 2011.

Thank You!

Contact

- ▶  <https://juanitorduz.github.io>
- ▶  github.com/juanitorduz
- ▶  [juanitorduz](https://twitter.com/juanitorduz)
- ▶  juanitorduz@gmail.com

