

# WB XAI-2 HW 1

Agata Kaczmarek  
24 03 2021

## Cel

Dla dostępnego zbioru danych dotyczącego telefonów (parametry dla wybranych modeli oraz ich ceny) przeanalizować w jaki sposób poszczególne parametry wpływają na przewidywaną cenę urządzenia oraz spróbować wyjaśnić. Model został stworzony przy wykorzystaniu lasu losowego z pakietu ranger.

Przyjrzyjmy się naszym danym, czy zawierają one jakieś braki? Jeśli tak, to je usuniemy.

```
which(is.na(phones))

## [1]  992 1011 1058 1097 1130 1198 1229 1234 1391 1400 1406 1413 1420 1425 1432
## [16] 1439 1441 1442 1444 1448 1450 1452 1456 1457 1458 1460 1461 1463 1467 1468
## [31] 1469 1472 1473 1476 1484 1485 1511 1535 1540 1544 1568 1571 1578 1587 1601
## [46] 1612 1619 1630 1639 1643 1645 1646 1648 1654

#usuwany brak
phones2<- phones[complete.cases(phones),]
```

## Predykcja dla wybranej obserwacji

Dla naszych danych tworzę model, następnie dla wybranej obserwacji ze zbioru danych wyliczam predykcję modelu.

```
#tworzymy model
model <- ranger::ranger(price~., data=phones2, seed=123)
model

## Ranger result
##
## Call:
## ranger::ranger(price ~ ., data = phones2, seed = 123)
##
## Type:                               Regression
## Number of trees:                     500
## Sample size:                         368
## Number of independent variables:     10
## Mtry:                                3
## Target node size:                    5
## Variable importance mode:            none
## Splitrule:                           variance
## OOB prediction error (MSE):          310268.6
## R squared (OOB):                     0.8800473

#predykcja dla pierwszego
predict(model, phones2[1,])$predictions

## [1] 2699.072

#dane podane
phones[1,]$price

## [1] 1999
```

Model przewiduje cenę danego urządzenia na podstawie dostępnych mu parametrów.

## Wyliczanie dekompozycji

Dla wcześniej wybranej obserwacji przy pomocy Break Down zobaczymy dekompozycję - powie nam ona w jaki sposób "myślał" model wyliczając przewidywaną cenę. Czyli które własności uznał za pozytywnie wpływające na cenę (podwyższające ją) a które negatywnie (obniżające ją).

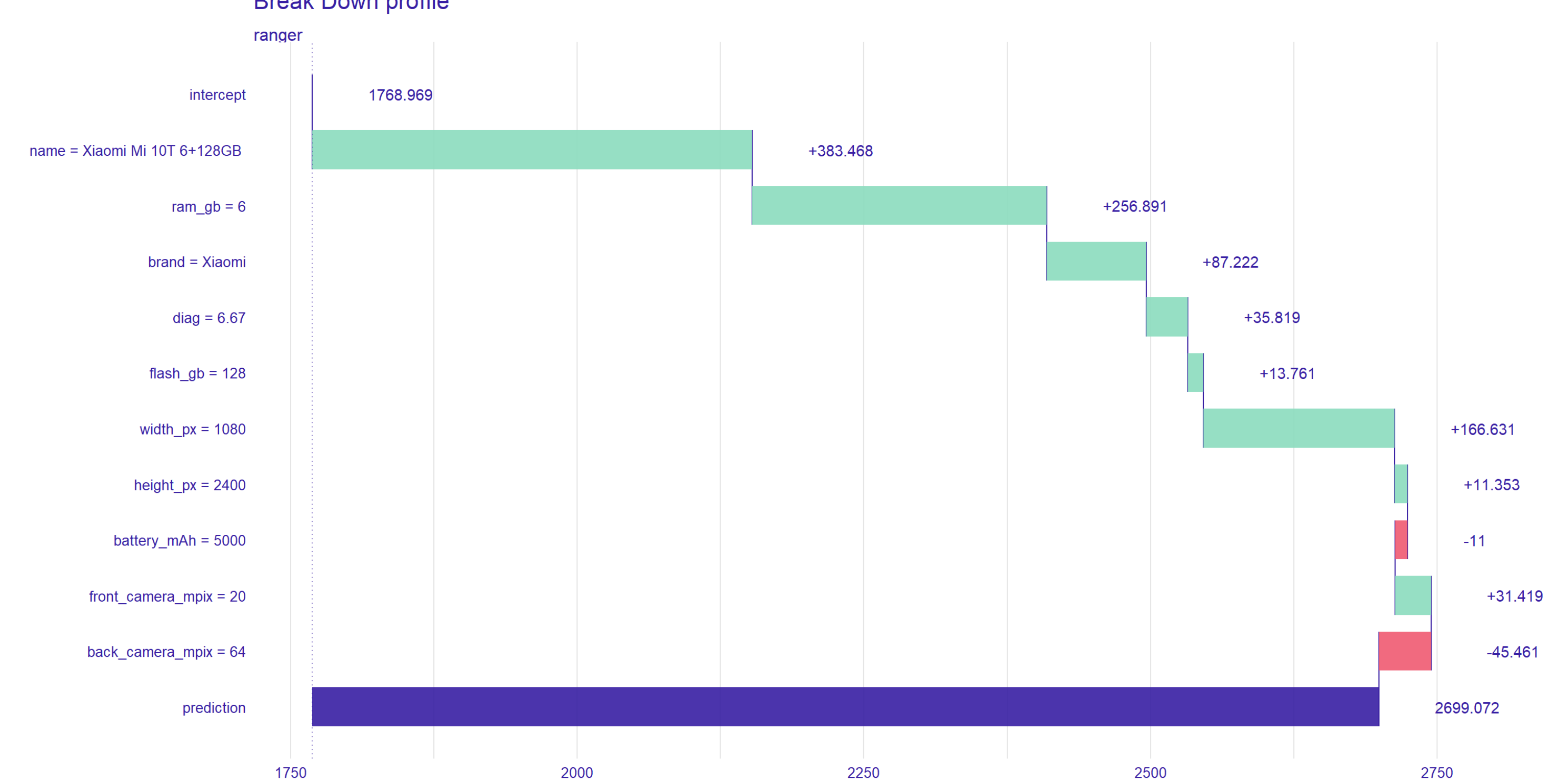
```
#explainer
explainer <- explain(model = model,
  data = phones2[,,-11],
  y = phones2$price,
  colorize = FALSE)

## Preparation of a new explainer is initiated
## -> model label      : ranger ( default )
## -> data              : 368 rows 10 cols
## -> target variable   : 368 values
## -> predict function  : yhat.ranger will be used ( default )
## -> predicted values  : No value for predict function target column. ( default )
## -> model_info        : package ranger , ver. 0.12.1 , task regression ( default )
## -> predicted values  : numerical, min = 242.4232 , mean = 1768.969 , max = 6886.142
## -> residual function : difference between y and yhat ( default )
## -> residuals        : numerical, min = -645.9523 , mean = 6.01745 , max = 1912.858
## A new explainer has been created!

explainer$predict_function

## function (X_model, newdata, ...)
## UseMethod("yhat")
## <bytecode: 0x0000000015a8f588>
## <environment: namespace:DALEX>

#Break Down
phones2_bd_1<-predict_parts(explainer,
  new_observation = phones[1,], type="break_down")
plot(phones2_bd_1)
```

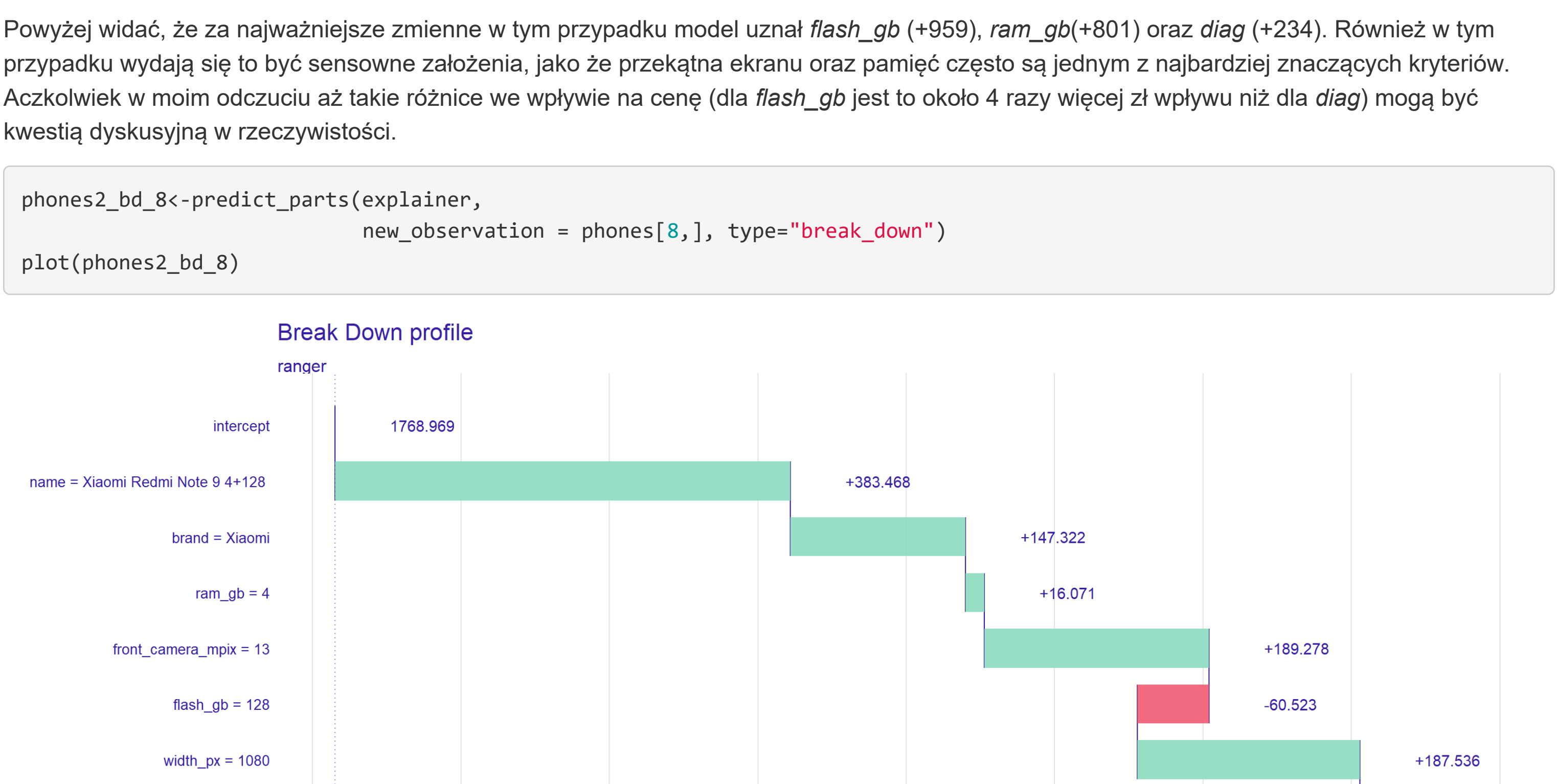


*prediction* odnacza przewidywaną cenę dla telefonu o wyżej pokazanych właściwościach. Jak widać największy pozytywny wpływ na cenę miały *name* (+383), *width\_px* (+166) oraz *ram\_gb* (+256). Ma to sens ponieważ wiemy, że często ludzie są skłonni zapłacić więcej za polecany im np przez znajomych telefon. Co więcej model uznał, że znaczny wpływ na cenę może mieć szerokość ekranu, ale co ciekawe, w tym przypadku, nie wysokość. W rzeczywistości szerokość i wysokość ekranu mogą być brane pod uwagę z podobnym priorytetem, podobnie jak przekątna ekranu, która dla modelu nie miała aż takiego wielkiego wpływu. Trzecią właściwością mającą największe znaczenie wg modelu była ilość pamięci RAM w GB, co też wydaje się być "rzeczywistym" założeniem.

## Najważniejsze zmienne

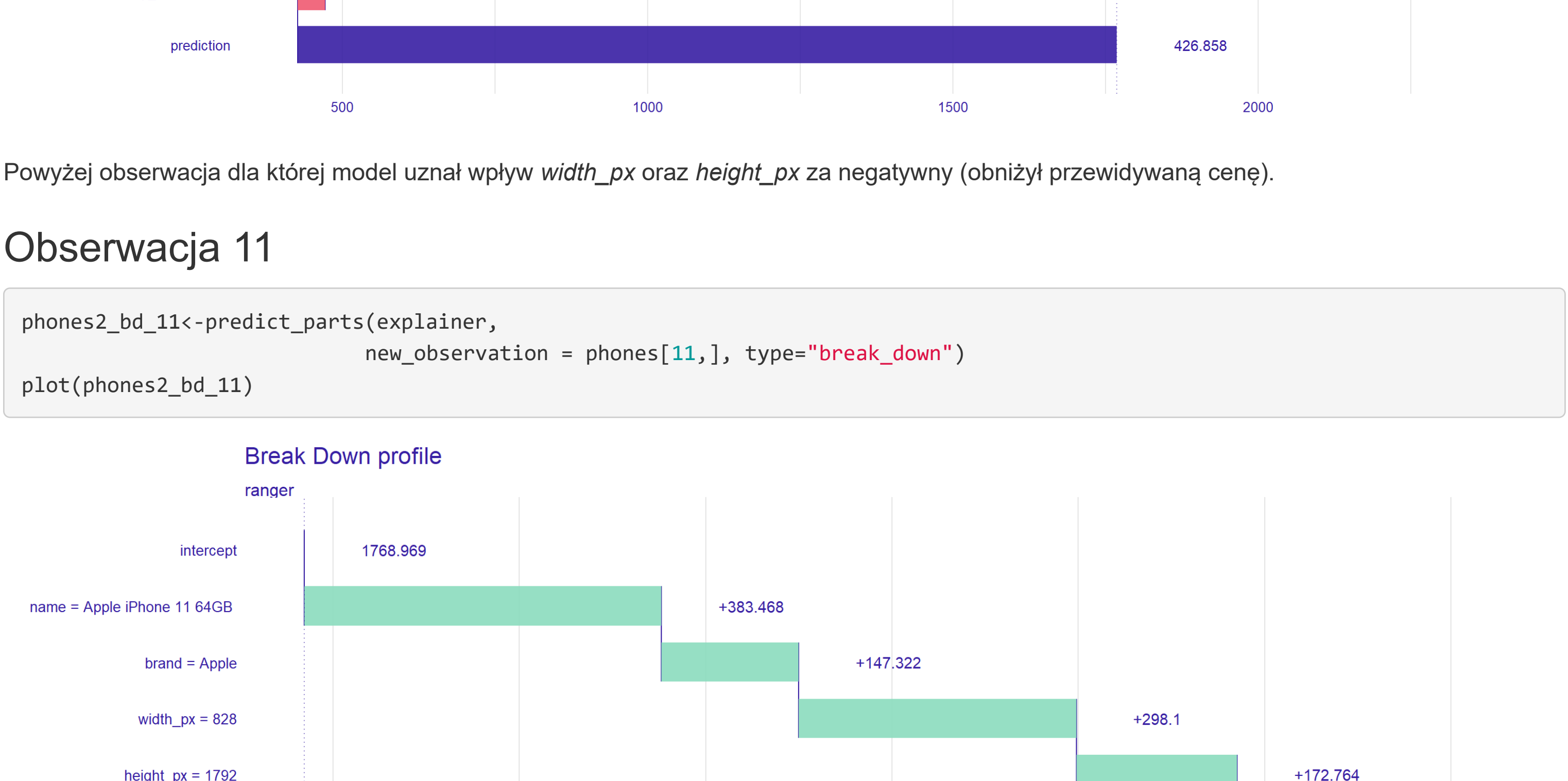
Wybór dwóch obserwacji, które mają najważniejsze inne zmienne.

```
phones2_bd_3<-predict_parts(explainer,
  new_observation = phones[3,], type="break_down")
plot(phones2_bd_3)
```



Powyżej widać, że za najważniejsze zmienne w tym przypadku model uznał *flash\_gb* (+959), *ram\_gb* (+801) oraz *diag* (+234). Również w tym przypadku wydaje się to być sensowne założenie, jako że przekątna ekranu oraz pamięć często są jednym z najbardziej znaczących kryteriów. Aczkolwiek w moim odczuciu aż takie różnice we wpływie na cenę (dla *flash\_gb* jest to około 4 razy więcej niż wpływu niż dla *diag*) mogą być kwestią dyskusyjną w rzeczywistości.

```
phones2_bd_8<-predict_parts(explainer,
  new_observation = phones[8,], type="break_down")
plot(phones2_bd_8)
```



Tutaj natomiast za najważniejsze model uznał *name* (+383), *front\_camera\_mpix* (+189) oraz *width\_px* (+187). Te zmienne, chociaż inne niż we wcześniejszym przykładzie, również wyglądają na sensowne założenia.

## Obserwacje mające inne efekty dla tych samych zmiennych

### Obserwacja 31

```
phones2_bd_31<-predict_parts(explainer,
  new_observation = phones[31,], type="break_down")
plot(phones2_bd_31)
```



Powyżej obserwacja dla której model uznał wpływ *width\_px* oraz *height\_px* za negatywny (obniżył przewidywaną cenę).

### Obserwacja 11

```
phones2_bd_11<-predict_parts(explainer,
  new_observation = phones[11,], type="break_down")
plot(phones2_bd_11)
```



Zaś tutaj model uznał wpływ *width\_px* oraz *height\_px* za pozytywny (podwyższył przewidywaną cenę).

Powyżej możemy zobaczyć, że wartości dla zmiennych *width\_px* oraz *height\_px* nie odbiegają od siebie w znaczny sposób. W obserwacji 31 szerokość ekranu to 720, a wysokość 1560, podczas gdy w obserwacji 11 te wartości odpowiednio to 828 oraz 1792. A jednak wpływ tych wartości na przewidywaną cenę ostateczna jest diametralnie inny. Co więcej, sama cena modelu bardzo różni się w obu przypadkach - dla obserwacji 31 jest to 443, natomiast 11 - 2859.

W tej sytuacji możemy powiedzieć, że rozważane tu dwie zmienne nie są jedynymi, które mają wpływ na końcowy wynik (cenę). Podczas tworzenia przewidywanych cen, model ma do dyspozycji 10 różnych zmiennych. A to znaczy, że wartości pozostałych 8 zmiennych będą również miały wpływ na cenę końcową, a także na istotność wymienionych tu zmiennych.

## Wnioski

Model stworzony na potrzeby tego zadania przewidywał ceny telefonów na podstawie 10 różnych ich cech. Jak można zauważyć po przytoczonych tu przykładach, nie możemy powiedzieć, że zawsze któraś konkretna cecha będzie miała największy wpływ na cenę, a któraś najmniejszy. Wpływ poszczególnych zmiennych i wartości jest bardziej skomplikowany do znalezienia, dla różnych modeli (np. po zmienienu wartości *seed* dla modelu) potrafi się on znacznie różnić. Oprócz tego, w naszych przykładach było widać, że dla poszczególnych predykcji różne cechy miały inny wpływ na końcową cenę. O ile wnioski na podstawie wyjaśnień naszego modelu w powyższych przykładach wydają się być "sensowne" i "rzeczywiste", o tyle założenie z góry że jedna konkretna zmienna ma największy wpływ na cenę dla całego zbioru danych, okazuje się błędne. Żeby wyciągnąć globalne wnioski, potrzebowalibyśmy użyć innych narzędzi.