



# **Uczenie maszynowe w R.**

## **Zaglądamy do czarnej skrzynki**

---

Bałtycki Festiwal Nauki

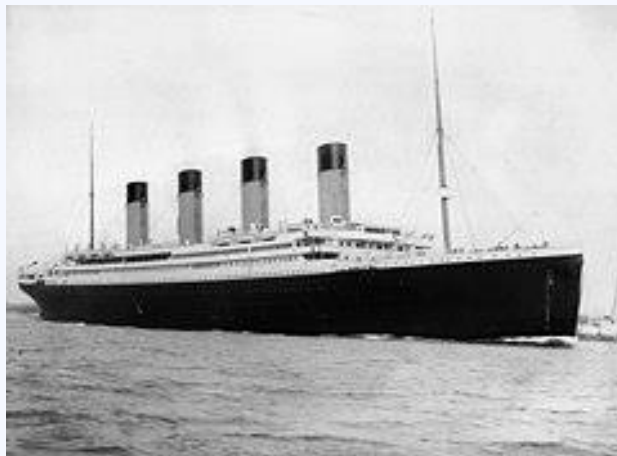




# Uczenie maszynowe

Dziedzina sztucznej inteligencji zajmująca się tworzeniem algorytmów, które potrafią uczyć się ze zbioru danych. W przypadku niektórych algorytmów, szczególnie tych bardziej skomplikowanych, trudno jest zrozumieć, jak dokładnie działają, stąd określenie **“czarna skrzynka”**.

# RMS Titanic



Był to największy i najbardziej luksusowy statek pasażerski swoich czasów, wyruszył w dziewiczy rejs z **Southampton** do **Nowego Jorku** 10 kwietnia 1912 roku. Niestety, 14 kwietnia statek zderzył się z górą lodową na północnym Atlantyku, co spowodowało poważne uszkodzenia i zatonięcie statku kilka godzin później. Tragedia ta pochłonęła życie **ponad 1500 osób** i stała się jednym z najbardziej tragicznych wydarzeń w historii morskiej.

# Tabela *titanic* to dane przygotowane przez prof. P. Biecką opisujące pasażerów Titanika.



```
titanic<-archivist::aread("pbiecek/models/27e5c")
```



```
View(titanic)
```

	gender	age	class	embarked	country	fare	sibsp	parch	survived
1	male	42.0000000	3rd	Southampton	United States	7.1100	0	0	no
2	male	13.0000000	3rd	Southampton	United States	20.0500	0	2	no
3	male	16.0000000	3rd	Southampton	United States	20.0500	1	1	no
4	female	39.0000000	3rd	Southampton	England	20.0500	1	1	yes
5	female	16.0000000	3rd	Southampton	Norway	7.1300	0	0	yes
6	male	25.0000000	3rd	Southampton	United States	7.1300	0	0	yes
7	male	30.0000000	2nd	Cherbourg	France	24.0000	1	0	no
8	female	28.0000000	2nd	Cherbourg	France	24.0000	1	0	yes
9	male	27.0000000	3rd	Cherbourg	Lebanon	18.1509	0	0	yes
10	male	20.0000000	3rd	Southampton	Finland	7.1806	0	0	yes
11	male	30.0000000	3rd	Southampton	Sweden	7.0500	0	0	no
12	male	27.0000000	3rd	Southampton	England	8.0100	0	0	no
13	female	40.0000000	3rd	Southampton	Sweden	9.0906	1	0	no
14	male	0.8333333	3rd	Southampton	England	9.0700	0	1	yes
15	female	18.0000000	3rd	Southampton	England	9.0700	0	1	yes

# Modele, których użyjemy

## Model regresji liniowej

przewidywanie prawdopodobieństwa wystąpienia zdarzenia oparte na zależności między zmienną zależną a jedną lub więcej zmiennymi niezależnymi

## Model lasów losowych

tworzy i łączy wiele drzew decyzyjnych, aby zwiększyć dokładność i stabilność przewidywań dla zadań klasyfikacji i regresji

## Model *gradient boosting*

stopniowe dodawanie nowych, prostych modeli (często drzew decyzyjnych) do zestawu już istniejących, aby coraz lepiej przewidywać wyniki, poprawiając błędy krok po kroku

## Model maszyny wektorów nośnych

znajduje linię lub płaszczyznę, która najlepiej oddziela różne grupy danych, pomagając w podejmowaniu decyzji o ich klasyfikacji



#model regresji logistycznej:

```
titanic_lmr <- rms::lrm(survived == "yes" ~ gender + rms::rccs(age) + class +  
  sibsp + parch + fare + embarked, titanic)
```

#model lasów losowych

```
set.seed(1313)
```

```
titanic_rf <- randomForest::randomForest(survived ~ class + gender + age +  
  sibsp + parch + fare + embarked, data = titanic, na.action=na.roughfix)
```

#model gradient boosting

```
set.seed(1313)
```

```
titanic_gbm <- gbm::gbm(survived == "yes" ~ class + gender + age +  
  sibsp + parch + fare + embarked, data = titanic,  
  n.trees = 15000, distribution = "bernoulli")
```

#model maszyny wektorów nosnych:

```
set.seed(1313)
```

```
titanic_svm <- e1071::svm(survived == "yes" ~ class + gender + age +  
  sibsp + parch + fare + embarked, data = titanic,  
  type = "C-classification", probability = TRUE)
```



# Porównywanie siły predykcyjnej modeli

```
model_performance(titanic_lmr_exp)$measures$auc  
model_performance(titanic_rf_exp)$measures$auc  
model_performance(titanic_gbm_exp)$measures$auc  
model_performance(titanic_svm_exp)$measures$auc
```

Model regresji logistycznej:	Model lasów losowych:	Model gradient boosting:	Model maszyny wektorów nośnych:
0,8174447	0,8636533	<b>0,8666712</b>	<b>0,8129198</b>

# Przykładowi pasażerowie

## Rose

*class: 1st  
gender: female  
age: 17  
sibsp: 0  
parch: 1  
fare: 850\$  
embarked:  
Southampton*

## Jack

*class: 3rd  
gender: male  
age: 20  
sibsp: 0  
parch: 0  
fare: ok 10\$  
embarked:  
Southampton*



### #przewidywania modeli dla pasazerow

```
{print(predict(titanic_lmr_exp, jack))
  print(predict(titanic_rf_exp, jack))
  print(predict(titanic_gbm_exp, jack))
  print(predict(titanic_svm_exp, jack))
}

{print(predict(titanic_lmr_exp, rose))
  print(predict(titanic_rf_exp, rose))
  print(predict(titanic_gbm_exp, rose))
  print(predict(titanic_svm_exp, rose))}
```

# Przewidywania modeli

## Rose



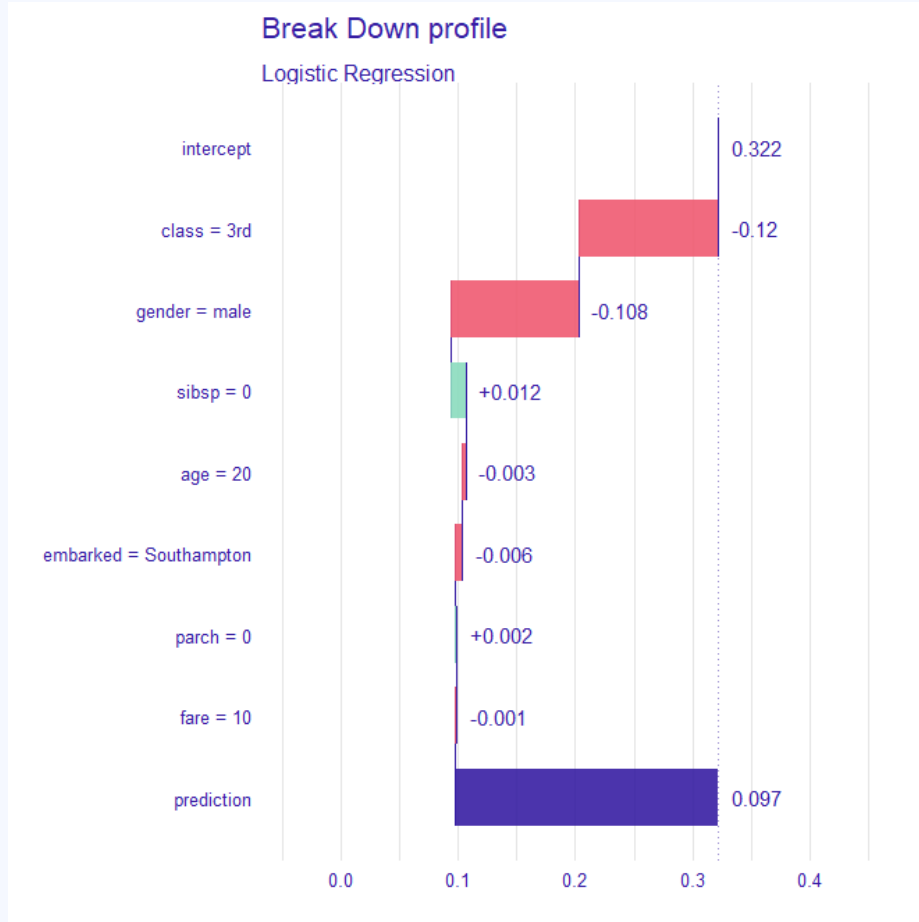
- Model regresji logistycznej: 98,04%
- Model lasów losowych: 96,8%
- Model gradient boosting: 96,899%
- Model maszyny wektorów nośnych: 43,57%

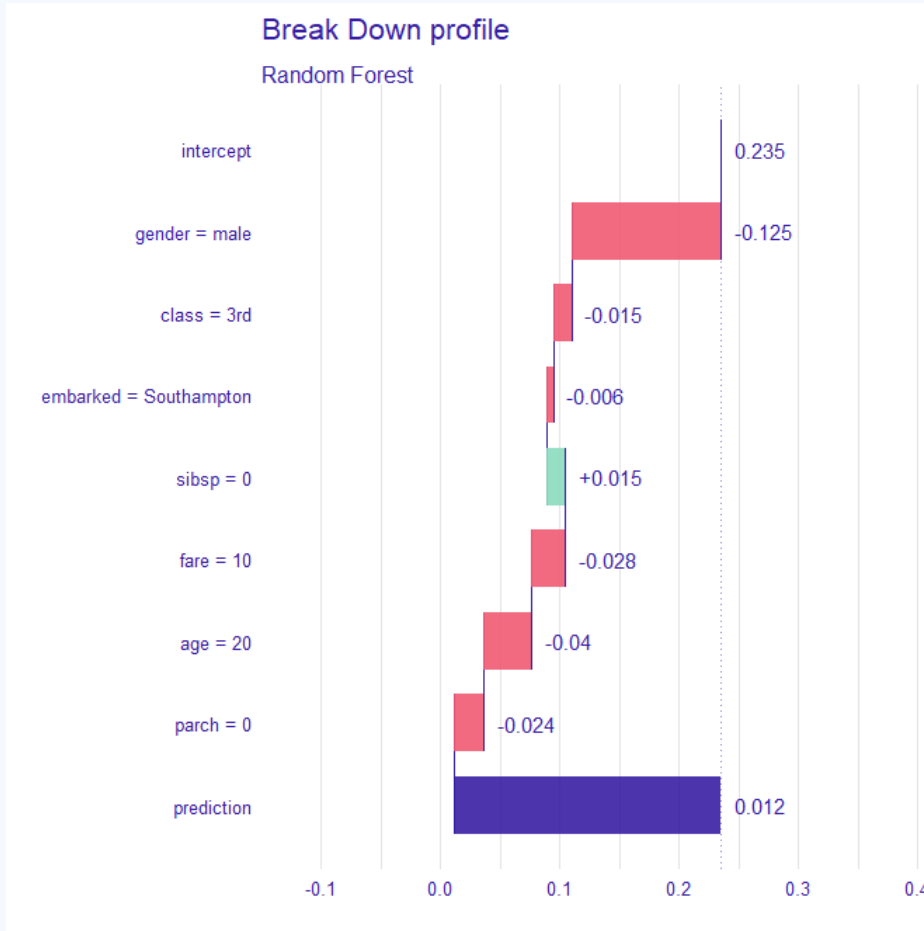
## Jack

- Model regresji logistycznej: 9,75%
- Model lasów losowych: 1,2%
- Model gradient boosting: 2,37%
- Model maszyny wektorów nośnych: 1,72%

# Co wpływa na wynik w danym modelu?

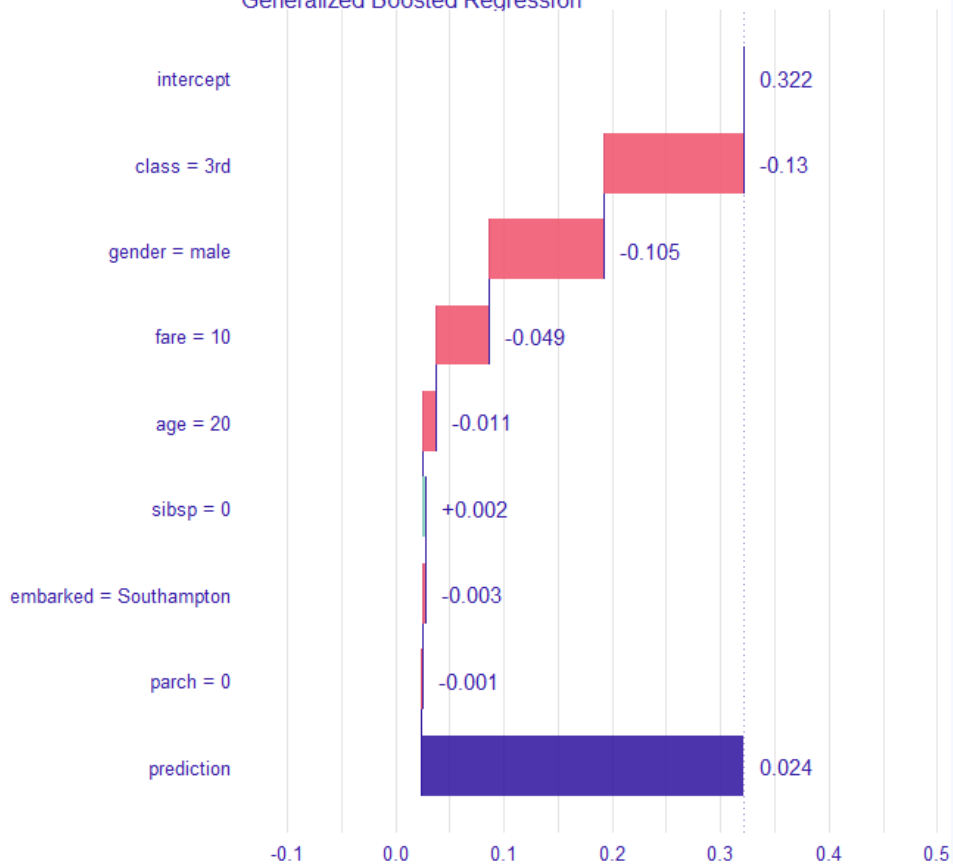






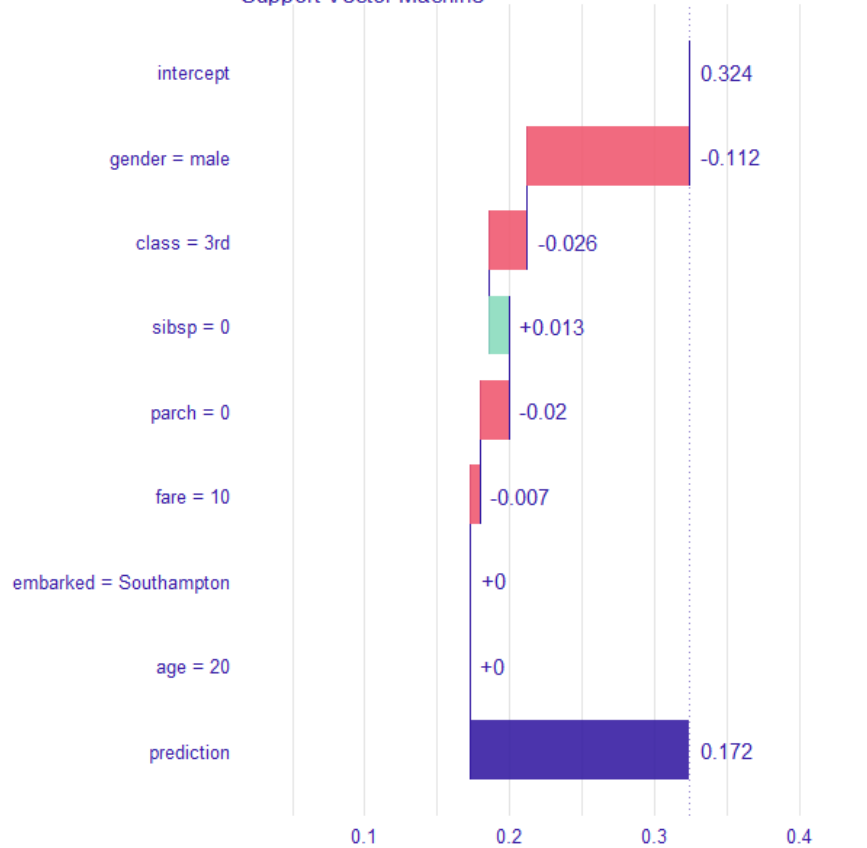
## Break Down profile

Generalized Boosted Regression



## Break Down profile

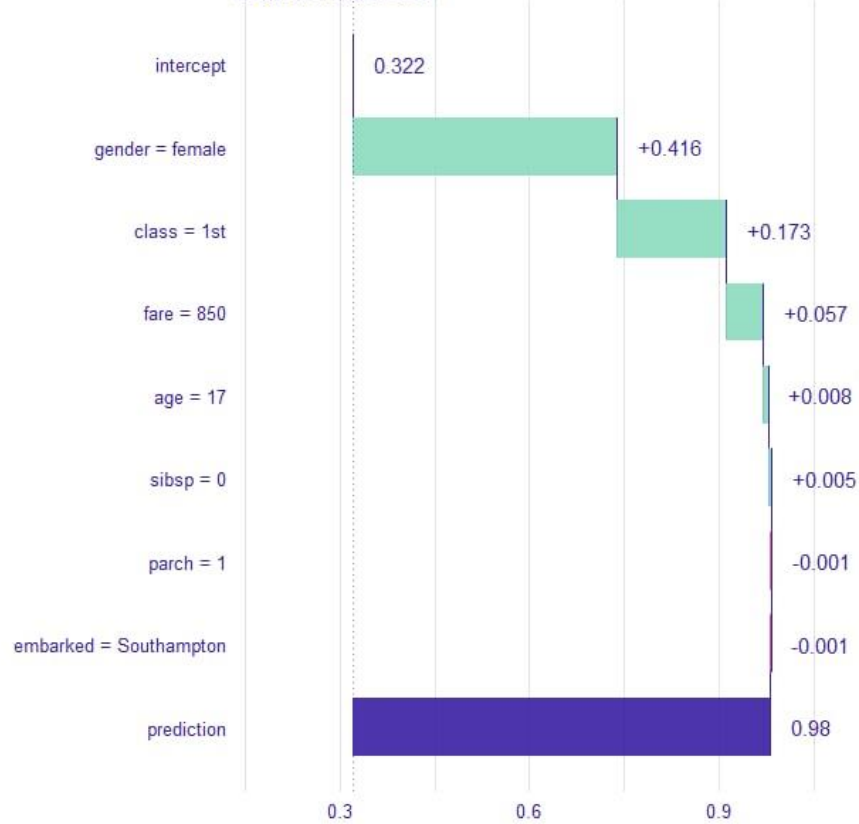
Support Vector Machine





## Break Down profile

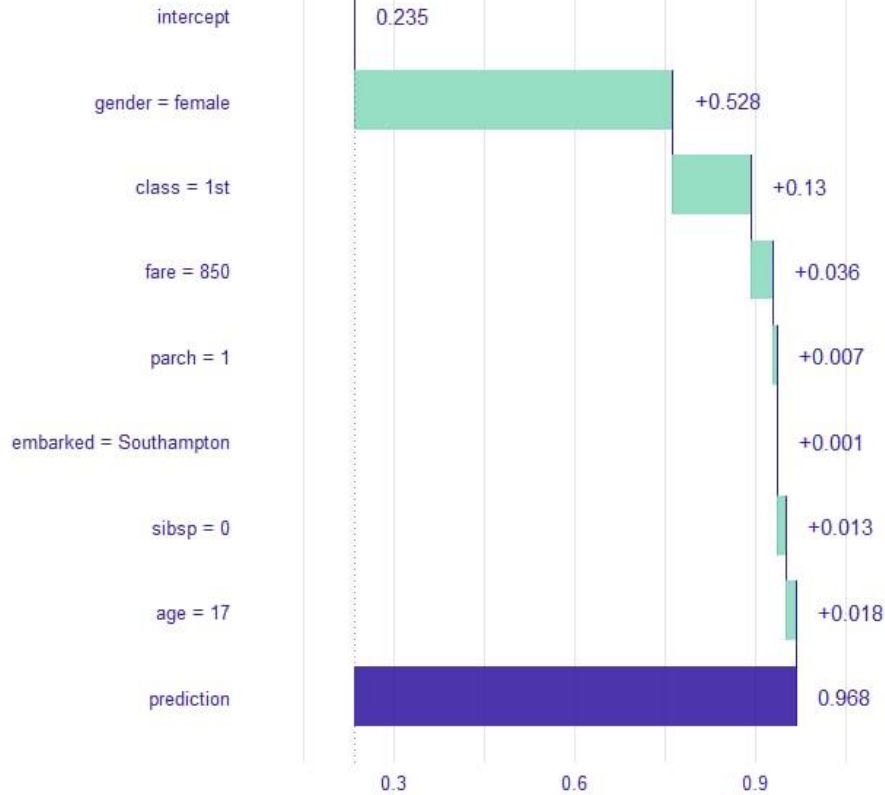
Logistic Regression





## Break Down profile

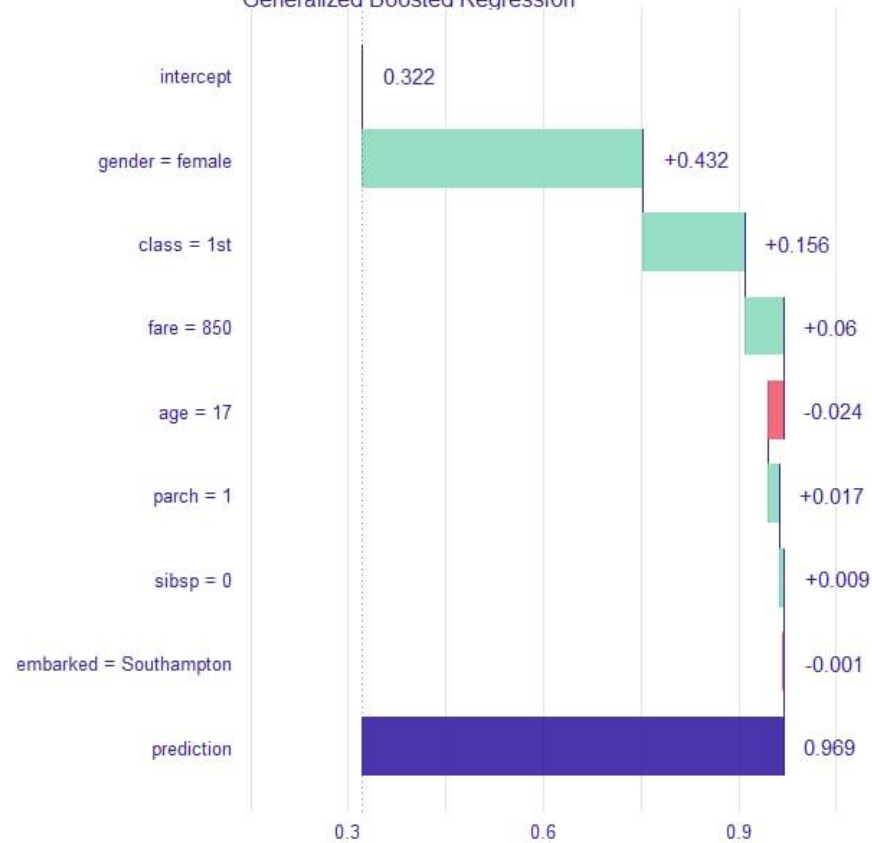
Random Forest





## Break Down profile

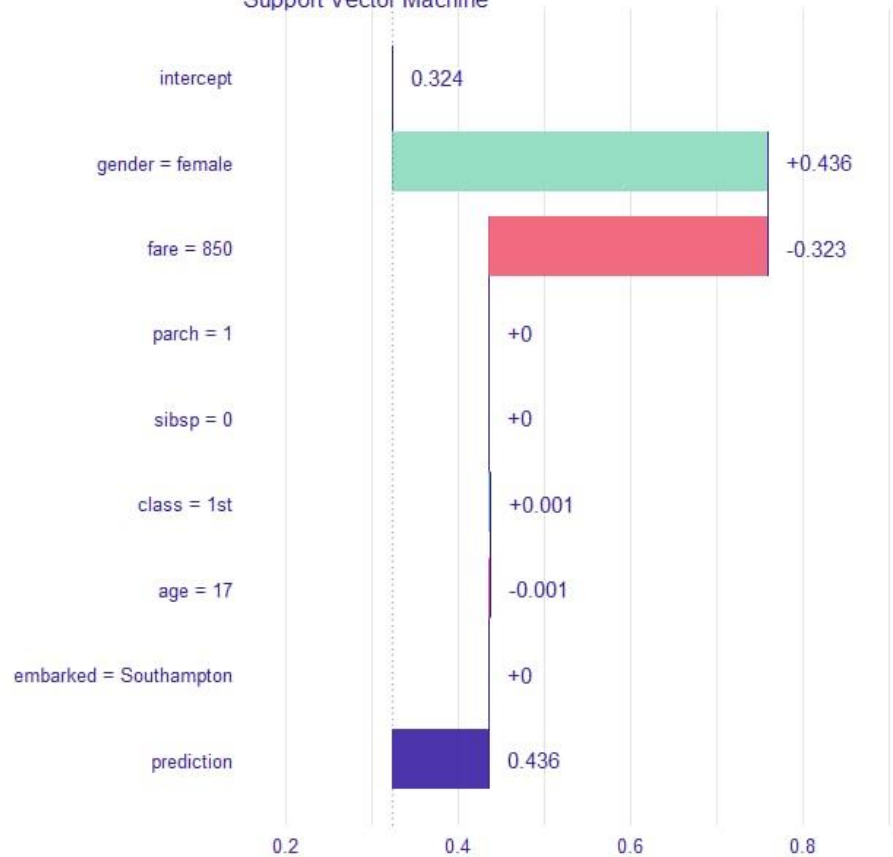
Generalized Boosted Regression



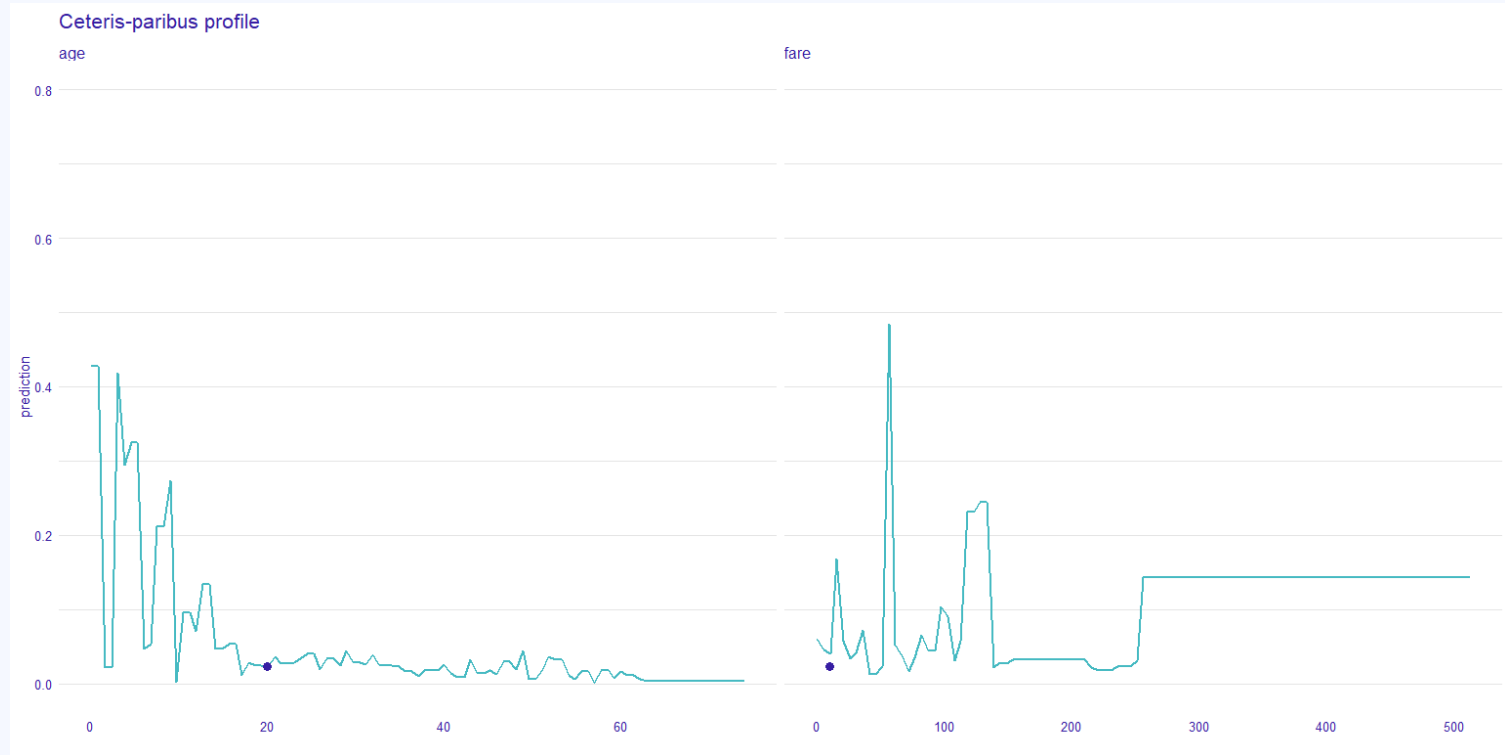


## Break Down profile

Support Vector Machine

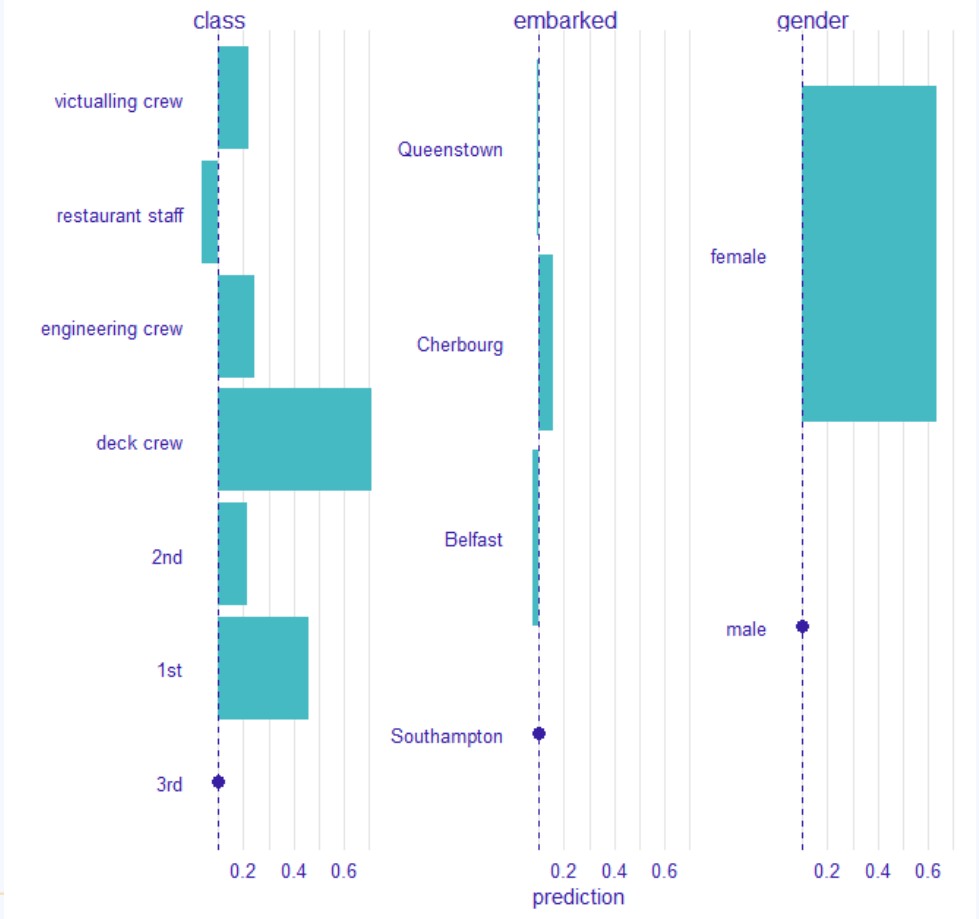


```
cp_titanic_gbm <- predict_profile(explainer = titanic_gbm_exp,  
                                new_observation = jack)  
  
plot(cp_titanic_gbm, variables = c("age", "fare")) +  
  ggtitle("Ceteris-paribus profile", "") + ylim(0, 0.8)
```

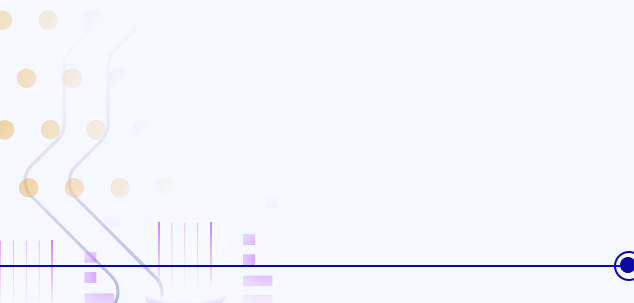
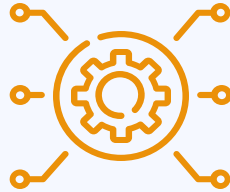


```
cp_titanic_lmr <- predict_profile(explainer = titanic_lmr_exp,  
                                new_observation = jack)  
  
plot(cp_titanic_lmr, variables = c("class", "embarked", "gender"),  
     variable_type = "categorical", categorical_type = "bars") +  
  ggtitle("Ceteris-paribus profile", "")
```

### Ceteris-paribus profile

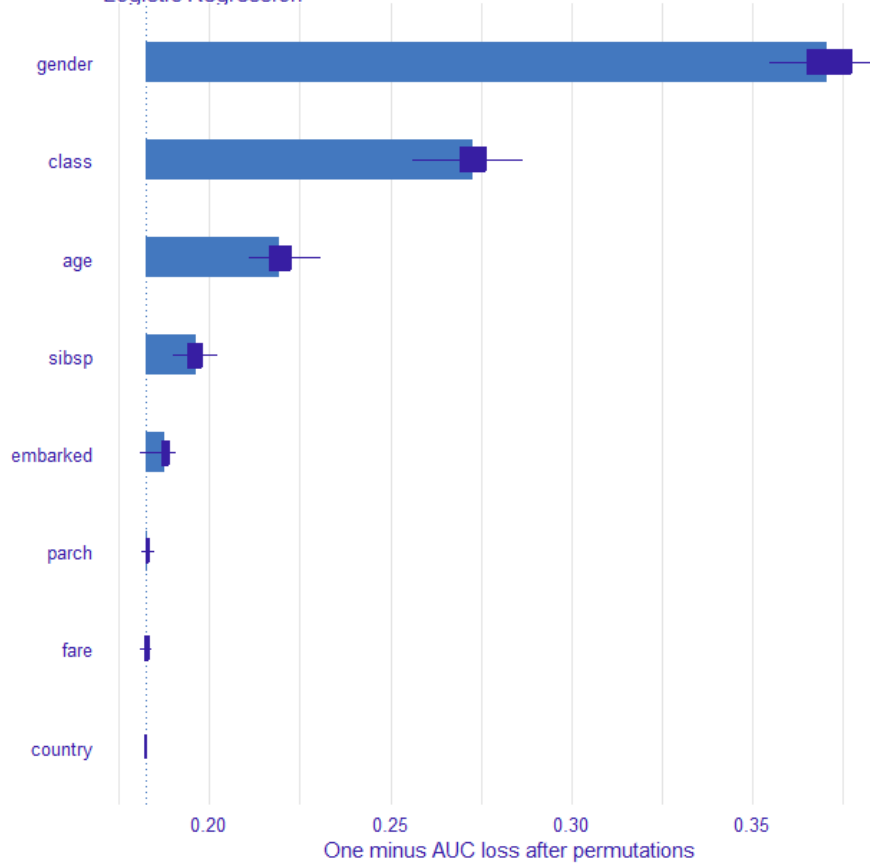


**Od jakich zmiennych  
zależy predykcja?**



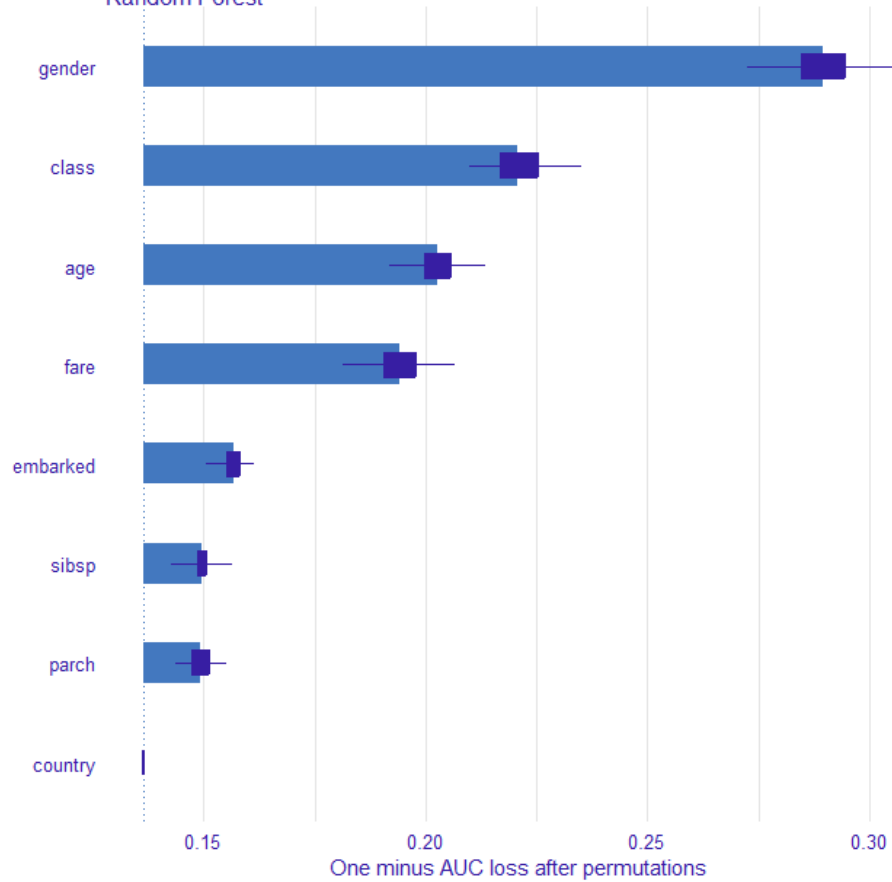
## Feature Importance

created for the Logistic Regression model  
Logistic Regression



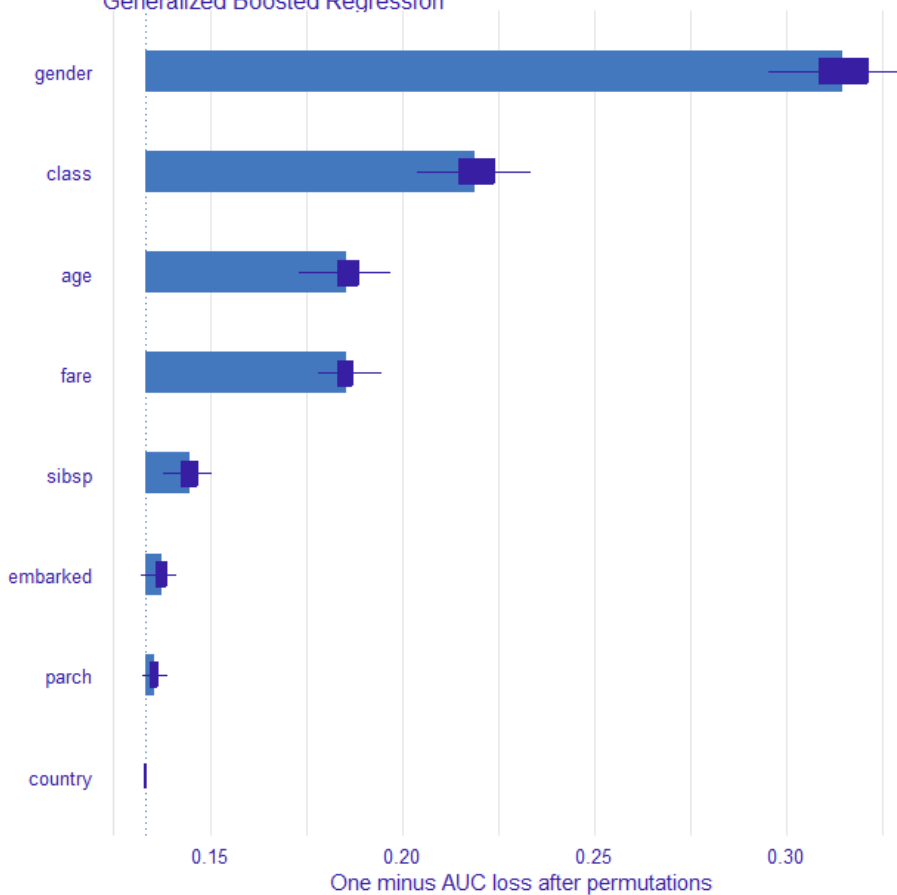
## Feature Importance

created for the Random Forest model  
Random Forest



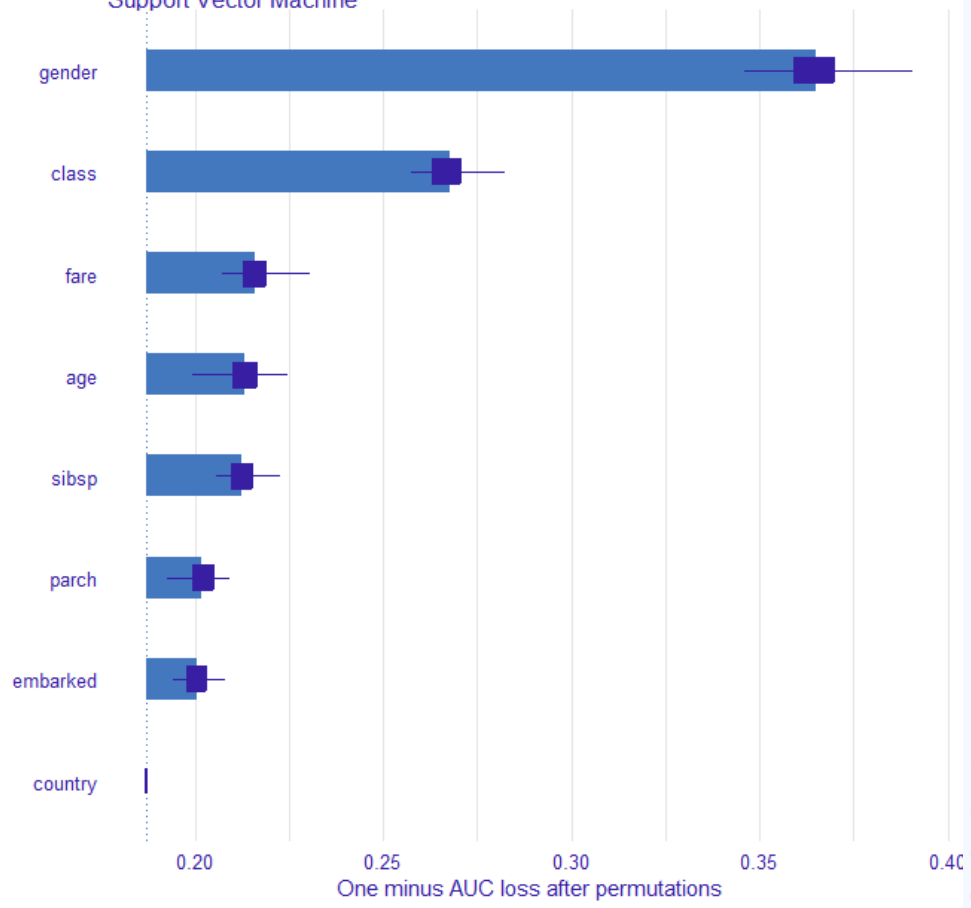
## Feature Importance

created for the Generalized Boosted Regression model  
Generalized Boosted Regression

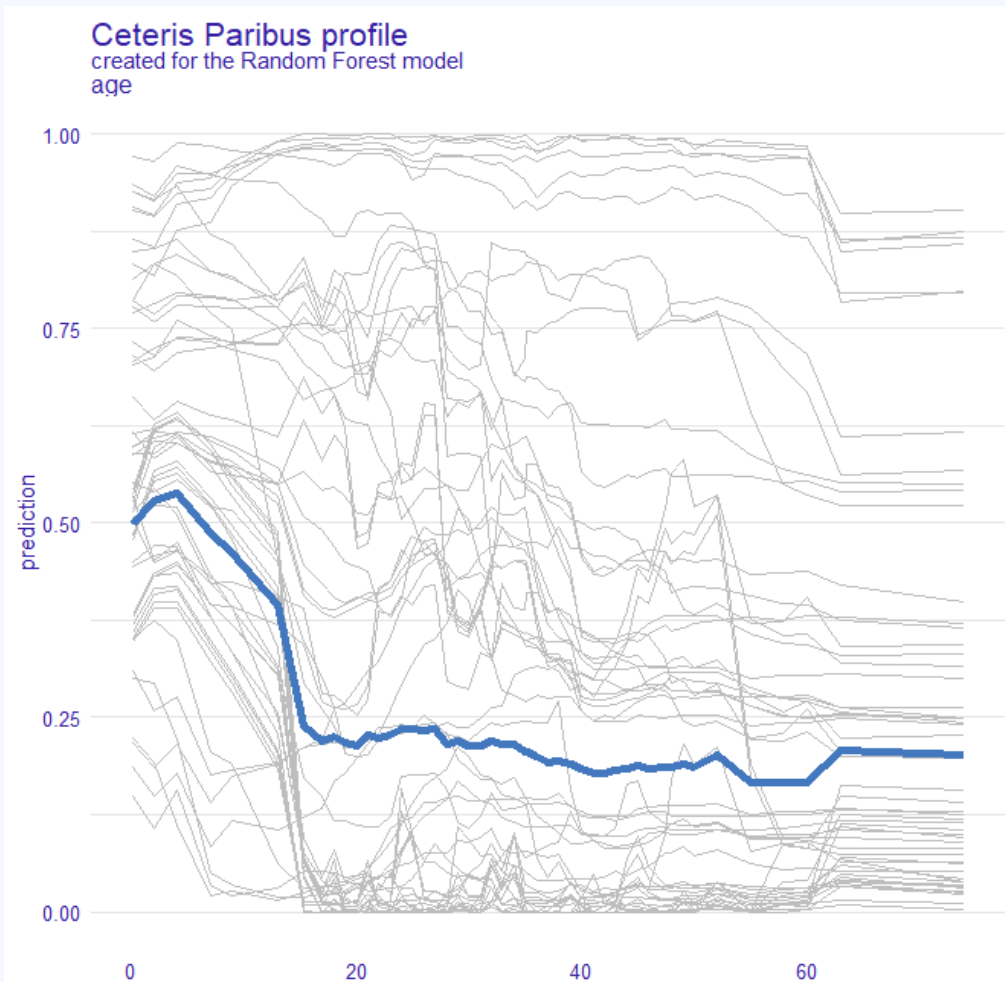


### Feature Importance

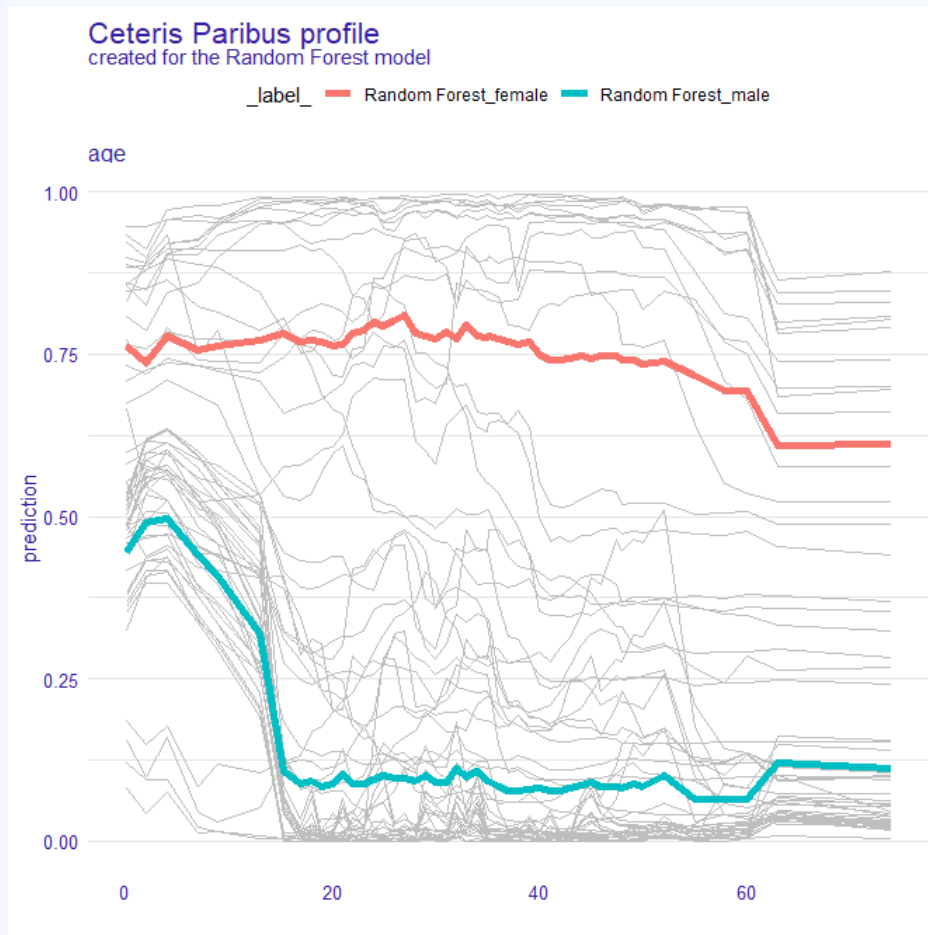
created for the Support Vector Machine model  
Support Vector Machine



# Jak wiek wpływa na przeżywalność?



# Jak płeć wpływa na przeżywalność?



# Teraz pora na Ciebie 😊

```
98 #twoja kolej
99 #twoje_imie <- data.frame(
100 #class = factor(, levels = c("1st", "2nd", "3rd", "deck crew", "engineering crew", "restaurant staff", "victualling crew")),
101 #gender = factor(, levels = c("female", "male")),
102 #age = ,
103 #sibsp = ,
104 #parch = ,
105 #fare = ,
106 #embarked = factor(, levels = c("Belfast", "Cherbourg", "Queenstown", "Southampton"))
107 #)
```