# Interpretable Machine Learning

# **Permutation Feature Importance (PFI)**





Figure: Bike Sharing Dataset

#### Learning goals

- Understand how PFI is computed
- Understanding strengths and weaknesses
- Testing Importance

#### PERMUTATION FEATURE IMPORTANCE (PFI) Breiman (2001)

**Idea:** "Destroy" feat. of interest  $x_j$  by perturbing it s.t. it becomes uninformative, e.g., randomly permute obs. in  $x_j$  (marginal distribution  $\mathbb{P}(x_j)$  stays the same). PFI for features  $x_s$  using test data  $\mathcal{D}$ :

- Measure the error without permuting feat. and with permuted feat. values  $\tilde{x}_S$
- Repeat permuting the feat. (e.g., *m* times) and avg. the difference of both errors:

$$\widehat{PFI}_{S} = \frac{1}{m} \sum_{k=1}^{m} \mathcal{R}_{emp}(\hat{f}, \widetilde{\mathcal{D}}_{(k)}^{S}) - \mathcal{R}_{emp}(\hat{f}, \mathcal{D}),$$
  
where  $\mathcal{R}_{emp}(\hat{f}, \mathcal{D}) = \frac{1}{n} \sum_{(x,y) \in \mathcal{D}} L(\hat{f}(x), y)$ 



#### PERMUTATION FEATURE IMPORTANCE (PFI) Breiman (2001)

**Idea:** "Destroy" feat. of interest  $x_j$  by perturbing it s.t. it becomes uninformative, e.g., randomly permute obs. in  $x_j$  (marginal distribution  $\mathbb{P}(x_j)$  stays the same). PFI for features  $x_s$  using test data  $\mathcal{D}$ :

- Measure the error without permuting feat. and with permuted feat. values  $\tilde{x}_S$
- Repeat permuting the feat. (e.g., *m* times) and avg. the difference of both errors:

$$\begin{split} \widehat{\textit{PFI}}_{\mathcal{S}} &= \frac{1}{m} \sum_{k=1}^{m} \mathcal{R}_{\text{emp}}(\hat{f}, \tilde{\mathcal{D}}_{(k)}^{\mathcal{S}}) - \mathcal{R}_{\text{emp}}(\hat{f}, \mathcal{D}), \\ \text{where } \mathcal{R}_{\text{emp}}(\hat{f}, \mathcal{D}) &= \frac{1}{n} \sum_{(x,y) \in \mathcal{D}} L(\hat{f}(x), y) \end{split}$$

The data  $\mathcal{D}$  where  $x_S$  is replaced with  $\tilde{x}^S$  is denoted as  $\tilde{\mathcal{D}}^S$ . Example of permuting feature  $x_S$  with  $S = \{1\}$  and m = 6:

${\cal D}$				$ ilde{\mathcal{D}}^{\mathcal{S}}_{(1)}$			$ ilde{\mathcal{D}}^{\mathcal{S}}_{(2)}$			$ ilde{\mathcal{D}}^{\mathcal{S}}_{(3)}$			$ ilde{\mathcal{D}}^{m{S}}_{(4)}$			$ ilde{\mathcal{D}}^{\mathcal{S}}_{(5)}$			$ ilde{\mathcal{D}}^{S}_{(6)}$		
$\mathbf{x}_1$	<b>x</b> <sub>2</sub>	<b>X</b> 3	<u>→</u>	$\mathbf{x}_{S}$	<b>X</b> <sub>2</sub>	<b>X</b> 3	$\mathbf{x}_{S}$	<b>X</b> 2	<b>X</b> 3	$\mathbf{x}_{S}$	<b>X</b> <sub>2</sub>	<b>X</b> 3	$\mathbf{X}_S$	<b>X</b> <sub>2</sub>	<b>X</b> 3	$\mathbf{X}_S$	<b>X</b> 2	<b>X</b> 3	$\mathbf{x}_{S}$	<b>X</b> 2	$\mathbf{X}_3$
1	4	7		1	4	7	2	4	7	2	4	7	1	4	7	3	4	7	3	4	7
2	5	8		2	5	8	1	5	8	3	5	8	3	5	8	1	5	8	2	5	8
3	6	9		3	6	9	3	6	9	1	6	9	2	6	9	2	6	9	1	6	9

Note: The *S* in  $x_S$  refers to a **S**ubset of features for which we are interested in their effect on the prediction. Here: We calculate the feature importance for one feature at a time |S| = 1.









- **1. Perturbation:** Sample feature values from the distribution of  $x_S$  ( $P(X_S)$ ).
  - $\Rightarrow$  Randomly permute feature  $x_S$

 $\Rightarrow$  Replace original feature with permuted feature  $\tilde{x}_S$  and create data  $\tilde{\mathcal{D}}^S$  containing  $\tilde{x}_S$ 

 $\mathcal{D}$ 





- **1.** Perturbation: Sample feature values from the distribution of  $x_S$  ( $P(X_S)$ ).
  - $\Rightarrow$  Randomly permute feature  $x_S$

 $\Rightarrow$  Replace original feature with permuted feature  $\tilde{x}_S$  and create data  $\tilde{\mathcal{D}}^S$  containing  $\tilde{x}_S$ 

2. Prediction: Make predictions for both data, i.e.,  ${\cal D}$  and  $\tilde{{\cal D}}^{{\cal S}}$ 



 $\mathcal{D}$ 



#### 3. Aggregation:

• Compute the loss for each observation in both data sets





#### 3. Aggregation:

- Compute the loss for each observation in both data sets
- Take the difference of both losses  $\Delta L$  for each observation

$$\mathcal{R}_{emp}(\hat{f}, \tilde{\mathcal{D}}^{\mathcal{S}}_{(k)}) - \mathcal{R}_{emp}(\hat{f}, \mathcal{D})$$





#### 3. Aggregation:

- Compute the loss for each observation in both data sets
- Take the difference of both losses  $\Delta L$  for each observation
- Average this change in loss across all observations Note: This is equivalent to computing  $\mathcal{R}_{emp}$  on both data sets and taking the difference





#### 3. Aggregation:

- Compute the loss for each observation in both data sets
- Take the difference of both losses  $\Delta L$  for each observation
- Average this change in loss across all observations
- Repeat perturbation and average over multiple repetitions

#### **EXAMPLE: BIKE SHARING DATASET**





#### Interpretation:

- Year (yr) and Temperature (temp) are most important features
- Destroying information about yr by permuting it increases mean absolute error of model by 816
- 5% and 95% quantile of repetitions due multiple permutations are shown as error bars

• Interpretation: PFI is the increase of model error when feature's information is destroyed



- Interpretation: PFI is the increase of model error when feature's information is destroyed
- Results can be unreliable due to random permutations
   ⇒ Solution: Average results over multiple repetitions



- Interpretation: PFI is the increase of model error when feature's information is destroyed
- Results can be unreliable due to random permutations
   ⇒ Solution: Average results over multiple repetitions
- Permuting features despite correlation with other features can lead to unrealistic combinations of feature values (since under dependence P(x<sub>j</sub>, x<sub>-j</sub>) ≠ P(x<sub>j</sub>)P(x<sub>-j</sub>)) → Extrapolation issue



- Interpretation: PFI is the increase of model error when feature's information is destroyed
- Results can be unreliable due to random permutations
   ⇒ Solution: Average results over multiple repetitions
- Permuting features despite correlation with other features can lead to unrealistic combinations of feature values (since under dependence P(x<sub>j</sub>, x<sub>-j</sub>) ≠ ℙ(x<sub>j</sub>)ℙ(x<sub>-j</sub>)) → Extrapolation issue
- PFI automatically includes importance of interaction effects with other features
   ⇒ Permutation also destroys information of interactions where permuted feature is involved
  - $\Rightarrow$  Importance of all interactions with the permuted feature are contained in PFI score



- Interpretation: PFI is the increase of model error when feature's information is destroyed
- Results can be unreliable due to random permutations
   ⇒ Solution: Average results over multiple repetitions
- Permuting features despite correlation with other features can lead to unrealistic combinations of feature values (since under dependence P(x<sub>j</sub>, x<sub>-j</sub>) ≠ P(x<sub>j</sub>)P(x<sub>-j</sub>)) → Extrapolation issue
- PFI automatically includes importance of interaction effects with other features
   ⇒ Permutation also destroys information of interactions where permuted feature is involved
  - $\Rightarrow$  Importance of all interactions with the permuted feature are contained in PFI score
- Interpretation of PFI depends on whether training or test data is used



#### **COMMENTS ON PFI - EXTRAPOLATION**

**Example:** Let  $y = x_3 + \epsilon_y$  with  $\epsilon_y \sim N(0, 0.1)$  where  $x_1 := \epsilon_1, x_2 := x_1 + \epsilon_2$  are highly correlated ( $\epsilon_1 \sim N(0, 1), \epsilon_2 \sim N(0, 0.01)$ ) and  $x_3 := \epsilon_3, x_4 := \epsilon_4$ , with  $\epsilon_3, \epsilon_4 \sim N(0, 1)$ . All noise terms are independent. Fitting a LM yields  $\hat{f}(\mathbf{x}) \approx 0.3x_1 - 0.3x_2 + x_3$ .



#### **COMMENTS ON PFI - EXTRAPOLATION**

**Example:** Let  $y = x_3 + \epsilon_y$  with  $\epsilon_y \sim N(0, 0.1)$  where  $x_1 := \epsilon_1, x_2 := x_1 + \epsilon_2$  are highly correlated ( $\epsilon_1 \sim N(0, 1), \epsilon_2 \sim N(0, 0.01)$ ) and  $x_3 := \epsilon_3, x_4 := \epsilon_4$ , with  $\epsilon_3, \epsilon_4 \sim N(0, 1)$ . All noise terms are independent. Fitting a LM yields  $\hat{f}(\mathbf{x}) \approx 0.3x_1 - 0.3x_2 + x_3$ .



Hexbin plot of  $x_1$ ,  $x_2$  before permuting  $x_1$  (left), after permuting  $x_1$  (center), and PFI scores (right)



### **COMMENTS ON PFI - EXTRAPOLATION**

**Example:** Let  $y = x_3 + \epsilon_y$  with  $\epsilon_y \sim N(0, 0.1)$  where  $x_1 := \epsilon_1, x_2 := x_1 + \epsilon_2$  are highly correlated ( $\epsilon_1 \sim N(0, 1), \epsilon_2 \sim N(0, 0.01)$ ) and  $x_3 := \epsilon_3, x_4 := \epsilon_4$ , with  $\epsilon_3, \epsilon_4 \sim N(0, 1)$ . All noise terms are independent. Fitting a LM yields  $\hat{f}(\mathbf{x}) \approx 0.3x_1 - 0.3x_2 + x_3$ .



Hexbin plot of  $x_1, x_2$  before permuting  $x_1$  (left), after permuting  $x_1$  (center), and PFI scores (right)  $\Rightarrow x_1$  and  $x_2$  should be irrelevant for the prediction  $\hat{f}(\mathbf{x})$  for

 $\{\mathbf{x} : \mathbb{P}(\mathbf{x}) > 0\}$  as  $0.3x_1 - 0.3x_2 \approx 0$  $\Rightarrow$  PFI evaluates model on unrealistic obs. outside  $\mathbb{P}(\mathbf{x}) \rightsquigarrow x_1, x_2$  are considered relevant (PFI > 0)



#### **COMMENTS ON PFI - INTERACTIONS**

**Example:** Let  $x_1, \ldots, x_4$  be independently and uniformly sampled from  $\{-1, 1\}$  and

 $y := x_1 x_2 + x_3 + \epsilon_Y$  with  $\epsilon_Y \sim N(0, 1)$ 

Fitting a LM yields  $\hat{f}(x) \approx x_1 x_2 + x_3$ .



#### **COMMENTS ON PFI - INTERACTIONS**

**Example:** Let  $x_1, \ldots, x_4$  be independently and uniformly sampled from  $\{-1, 1\}$  and

$$y := x_1 x_2 + x_3 + \epsilon_Y$$
 with  $\epsilon_Y \sim N(0, 1)$ 

Fitting a LM yields  $\hat{f}(x) \approx x_1 x_2 + x_3$ .

Although  $x_3$  alone contributes as much to the prediction as  $x_1$  and  $x_2$  jointly, all three are considered equally relevant.

 $\Rightarrow$  PFI does not fairly attribute the performance to the individual features.





#### **COMMENTS ON PFI - TEST VS. TRAINING DATA**

**Example:**  $x_1, \ldots, x_{20}, y$  are independently sampled from  $\mathcal{U}(-10, 10)$ . An xgboost model with default hyperparameters is fit on a small training set of 50 observations. The model overfits heavily.





**Figure:** While PFI on test data considers all features to be irrelevant, PFI on train data exposes the features on which the model overfitted.

#### **COMMENTS ON PFI - TEST VS. TRAINING DATA**

**Example:**  $x_1, \ldots, x_{20}, y$  are independently sampled from  $\mathcal{U}(-10, 10)$ . An xgboost model with default hyperparameters is fit on a small training set of 50 observations. The model overfits heavily.





**Figure:** While PFI on test data considers all features to be irrelevant, PFI on train data exposes the features on which the model overfitted.

Why? PFI can only be nonzero if the permutation breaks a dependence in the data. Spurious correlations help the model perform well on train data but are not present in the test data.

 $\Rightarrow$  If you are interested in which features help the model to generalize, apply PFI on test data.

#### **IMPLICATIONS OF PFI**

Can we get insight into whether the ...

- feature  $x_j$  is causal for the prediction?
  - $PFI_j \neq 0 \Rightarrow$  model relies on  $x_j$
  - As the training vs. test data example demonstrates, the converse does not hold



#### **IMPLICATIONS OF PFI**

Can we get insight into whether the ...

- feature  $x_j$  is causal for the prediction?
  - $PFI_j \neq 0 \Rightarrow$  model relies on  $x_j$
  - As the training vs. test data example demonstrates, the converse does not hold
- 2 feature  $x_j$  contains prediction-relevant information?
  - *PFI<sub>j</sub>* ≠ 0 ⇒ x<sub>j</sub> is dependent of y or it's covariates x<sub>-j</sub> or both (due to extrapolation)
  - *x<sub>j</sub>* is not exploited by model (regardless of whether it is useful for *y* or not)
     ⇒ *PFI<sub>j</sub>* = 0



### **IMPLICATIONS OF PFI**

Can we get insight into whether the ...

- feature  $x_i$  is causal for the prediction?
  - $PFI_j \neq 0 \Rightarrow$  model relies on  $x_j$
  - As the training vs. test data example demonstrates, the converse does not hold
- Ifeature x<sub>j</sub> contains prediction-relevant information?
  - *PFI<sub>j</sub>* ≠ 0 ⇒ x<sub>j</sub> is dependent of y or it's covariates x<sub>-j</sub> or both (due to extrapolation)
  - *x<sub>j</sub>* is not exploited by model (regardless of whether it is useful for *y* or not)
     ⇒ *PFI<sub>j</sub>* = 0
- model requires access to  $x_j$  to achieve it's prediction performance?
  - As the extrapolation example demonstrates, such insight is not possible



• PIMP was originally introduced for random forest's built-in permutation feature importance



- PIMP was originally introduced for random forest's built-in permutation feature importance
- PIMP investigates whether the PFI score significantly differs from 0
   ⇒ Useful because PFI can be non-zero due to stochasticity



- PIMP was originally introduced for random forest's built-in permutation feature importance
- PIMP investigates whether the PFI score **significantly** differs from 0 ⇒ Useful because PFI can be non-zero due to stochasticity
- PIMP tests the *H*<sub>0</sub>-hypothesis: Feature is independent of the target *y* (unimportant)



- PIMP was originally introduced for random forest's built-in permutation feature importance
- PIMP investigates whether the PFI score **significantly** differs from 0 ⇒ Useful because PFI can be non-zero due to stochasticity
- PIMP tests the *H*<sub>0</sub>-hypothesis: Feature is independent of the target *y* (unimportant)
- Sampling under *H*<sub>0</sub>: Permute target *y*, retrain model, compute PFI scores (repeat)
  - $\Rightarrow$  Permuting *y* breaks relationship to all features
  - $\Rightarrow$  By computing PFI scores again, we obtain distribution of PFI scores under  $H_0$



- PIMP was originally introduced for random forest's built-in permutation feature importance
- PIMP investigates whether the PFI score **significantly** differs from 0 ⇒ Useful because PFI can be non-zero due to stochasticity
- PIMP tests the *H*<sub>0</sub>-hypothesis: Feature is independent of the target *y* (unimportant)
- Sampling under *H*<sub>0</sub>: Permute target *y*, retrain model, compute PFI scores (repeat)
  - $\Rightarrow$  Permuting *y* breaks relationship to all features
  - $\Rightarrow$  By computing PFI scores again, we obtain distribution of PFI scores under  $H_0$
- Compute p-value the tail probability under *H*<sub>0</sub> and use it as a new importance measure



### **TESTING IMPORTANCE (PIMP)**

PIMP algorithm:

- For  $m \in \{1, \ldots, n_{repetitions}\}$ :
  - Permute response vector y
  - Retrain model with data **X** and permuted y
  - Compute feature importance  $PFI_i^m$  for each feature *j* (under  $H_0$ )



### **TESTING IMPORTANCE (PIMP)**

PIMP algorithm:

- For  $m \in \{1, \ldots, n_{repetitions}\}$ :
  - Permute response vector y
  - Retrain model with data **X** and permuted y
  - Compute feature importance  $PFI_i^m$  for each feature j (under  $H_0$ )
- 2 Train model with  $\mathbf{X}$  and unpermuted y



### **TESTING IMPORTANCE (PIMP)**

PIMP algorithm:

- For  $m \in \{1, \ldots, n_{repetitions}\}$ :
  - Permute response vector y
  - Retrain model with data X and permuted y
  - Compute feature importance  $PFI_i^m$  for each feature *j* (under  $H_0$ )
- 2 Train model with X and unpermuted y
- For each feature  $j \in \{1, \ldots, p\}$ :
  - Fit probability distribution of the feature importance values *PFI<sub>j</sub><sup>m</sup>*, *m* ∈ {1,..., *n<sub>repetitions</sub>*} (choice between Gaussian, lognormal, gamma or non-parametric)
  - Compute feature importance *PFI<sub>j</sub>* for the model without permutation of *y* (under *H*<sub>1</sub>)
  - Retrieve the p-value of PFI<sub>j</sub> based on the fitted distribution



#### PIMP FOR EXTRAPOLATION EXAMPLE

**Recall:**  $y = x_3 + \epsilon_y$  with  $\epsilon_y \sim N(0, 0.1)$ ,  $x_1$ ,  $x_2$  highly correlated but independent of y,  $x_4$  is independent of y and all other variables. Fitting a LM yields  $\hat{f}(\mathbf{x}) \approx 0.3x_1 - 0.3x_2 + x_3$ .





- Histograms:  $H_0$  distribution of PFI scores after permuting y (1000 repetitions)
- Red: PFI score estimated on unpermuted *y* (under *H*<sub>1</sub>) → compare against *H*<sub>0</sub> distribution
- Results: Although PFI for x<sub>1</sub> and x<sub>2</sub> is nonzero (red), PIMP considers them not significantly relevant (p-value > 0.05)

▶ Romano et al. (2010)

• When should we reject the *H*<sub>0</sub>-hypothesis for a feature?



#### ▶ Romano et al. (2010)

- When should we reject the *H*<sub>0</sub>-hypothesis for a feature?
- The larger the number of features, the more tests need to be performed by PIMP  $\rightsquigarrow$  **Multiple testing problem**: If multiplicity of tests is not taken into account, the probability that some of the true  $H_0$ -hypothesis is rejected (type-I error) by chance may be large



#### ▶ Romano et al. (2010)

- When should we reject the *H*<sub>0</sub>-hypothesis for a feature?
- The larger the number of features, the more tests need to be performed by PIMP  $\rightsquigarrow$  **Multiple testing problem**: If multiplicity of tests is not taken into account, the probability that some of the true  $H_0$ -hypothesis is rejected (type-I error) by chance may be large
- Accounting for multiplicity of individual tests can be achieved by controlling an appropriate error rate, e.g., the **family-wise error rate** (FWE: probability of at least one type-I error)



#### ▶ Romano et al. (2010)

- When should we reject the *H*<sub>0</sub>-hypothesis for a feature?
- The larger the number of features, the more tests need to be performed by PIMP  $\rightsquigarrow$  **Multiple testing problem**: If multiplicity of tests is not taken into account, the probability that some of the true  $H_0$ -hypothesis is rejected (type-I error) by chance may be large
- Accounting for multiplicity of individual tests can be achieved by controlling an appropriate error rate, e.g., the **family-wise error rate** (FWE: probability of at least one type-I error)
- One classical method to control the FWE is the **Bonferroni correction** which rejects a null hypothesis if its p-value is smaller than  $\alpha/m$  with *m* as the number of performed parallel tests

