

# **Feature Importance in Feature Selection**

Feature importance refers to techniques that assign scores to features based on their contribution to the performance of a predictive model. It plays a crucial role in feature selection, **which involves choosing the most relevant features for our model and discarding less useful ones**. This ultimately improves the model's accuracy, efficiency, and interpretability.

Here's how feature importance helps in feature selection:

- **Identify the most impactful features:** By understanding which features contribute most to our model's predictions, we can focus our analysis and training on those features, leading to better performance.
- **Reduce model complexity:** Removing irrelevant or redundant features, makes the model easier to train and interpret.
- **Gain insights into the data:** Feature importance scores can reveal hidden relationships and patterns within the data. This can provide valuable insights into the underlying processes or factors influencing the target variable.

There are different techniques for calculating feature importance, each with its own advantages and limitations:

- **Model-specific methods:** Some models like decision trees or linear regression have built-in methods for calculating feature importance.
- **Permutation importance:** Randomly shuffles the values of a feature and measures the decrease in model performance to assess its importance.
- **Shapley values:** Explains the marginal contribution of each feature to a prediction, accounting for interactions between features.

Choosing the right technique depends on your specific model and data. It's important to understand the assumptions and limitations of each method to interpret the results effectively.

## **Note:**

- Sometimes, important feature may get low score. So always we need to be clear with our data and its feature.

## **Forward selection**

Forward selection is a method for gradually building a model by **adding one feature at a time, aiming for the best model within a specified number of features.**

### **Forward selection works as below:**

**Empty Model Start:** Begin with a model having no features, just an intercept.

**Evaluate Individual Features:** Assess each feature's impact by calculating a performance metric

**Choose the Best Feature:** Select the feature that brings the greatest improvement in the chosen performance metric.

**Repeat the Process:** Add the chosen feature to the model and iterate, reassessing all remaining features with the updated model.

**Stop the Process:** Keep adding features and evaluating until a stopping criterion is met, like reaching a desired number of features or encountering diminishing returns.

### **Benefits of forward selection:**

**Interpretability:** Easily trace and understand how each feature contributes to the model.

**Simplicity:** Straightforward and beginner-friendly.

**Overfitting Prevention:** Gradual feature addition reduces the risk of overfitting compared to selecting all features at once.

### **Drawbacks of forward selection:**

**Suboptimal Selection:** Might miss important features that only matter in combination with others.

**Stopping Criteria Sensitivity:** Choosing the right stopping point can be subjective.

**Performance Limitations:** May be outperformed by more advanced methods in complex data scenarios.

### **Note:**

Forward selection is a good initial approach for feature selection, especially with smaller datasets or when interpretability is crucial.

```
def forward_selection(data, target, significance_level=0.05):
```

#initial\_features is a list containing the names of all features in the dataset, and best\_features will eventually store the selected features.

```
    initial_features = data.columns.tolist()
```

```
    best_features = []
```

#This initiates a while loop that continues until there are no more features left to consider.

```
    while (len(initial_features)>0):
```

```
        remaining_features = list(set(initial_features)-set(best_features))
```

```
        new_pval = pd.Series(index=remaining_features)
```

#The loop iterates through each remaining feature, fits an OLS model using the selected features plus the current new feature, and records the p-value of the new feature.

```
        for new_column in remaining_features:
```

```
            model = sm.OLS(target, sm.add_constant(data[best_features+[new_column]])).fit()
```

```
            new_pval[new_column] = model.pvalues[new_column]
```

#If the minimum p-value is below the significance level, the feature corresponding to that minimum p-value is added to best\_features. If not, the loop breaks, indicating the end of the feature selection process.

```
        min_p_value = new_pval.min()
```

```
        if(min_p_value<significance_level):
```

```
            best_features.append(new_pval.idxmin())
```

```
        else:
```

```
            break
```

```
    return best_features
```

## **Backward selection**

Backward selection takes a distinctive approach to navigating the feature landscape compared to its forward counterpart. Instead of starting small and adding features, it begins with a comprehensive model containing all features and systematically removes the least informative ones. The goal is to achieve optimal performance with the fewest features. Simply we can say like it is opposite to Forward selection approach

### **Backward selection works as below:**

**Start with a Full Model:** Begin with a model containing all available features.

**Evaluate Features Individually:** Calculate the same performance metric

**Choose the Least Helpful:** Identify the feature whose removal leads to the smallest decrease in the chosen performance metric.

**Remove the Chosen One:** Discard the least helpful feature from the model and start over from step 2, reassessing all remaining features with the updated model.

**Stop the Pruning:** Continue removing features and evaluating until a stopping criterion is met, similar to forward selection.

### **Benefits of backward selection:**

**Potentially Identifies Irrelevant Features:** Eliminates potentially irrelevant features that might not significantly contribute to the model.

**More Robust to Multicollinearity:** Can avoid selecting redundant features that are highly correlated with others.

**Potentially More Performant:** In some cases, can lead to more accurate models compared to forward selection.

### **Drawbacks of backward selection:**

**Missing Important Features:** May overlook features that only become relevant in combination with others, potentially impacting model performance.

**Less Interpretable:** The removal process can be harder to trace, making it challenging to understand individual feature contributions.

**Computationally Expensive:** Evaluating all features at each step can be computationally expensive for large datasets.

It offers a focused approach, exploring other methods ensures you harvest the most valuable insights from your data!

```

def backward_elimination(data, target,significance_level = 0.05):

    features = data.columns.tolist()

    while(len(features)>0):

        #This line adds a constant term to the feature set. In linear regression, this constant term
        #represents the intercept.

        features_with_constant = sm.add_constant(data[features])

        #This line fits an ordinary least squares (OLS) linear regression model using the current set of
        #features and calculates the p-values for each feature. The [1:] indexing removes the p-value
        #corresponding to the constant term.

        p_values = sm.OLS(target, features_with_constant).fit().pvalues[1:]

        max_p_value = p_values.max()

        #This line checks if the maximum p-value is greater than or equal to the specified
        #significance level.

        if(max_p_value >= significance_level):

            excluded_feature = p_values.idxmax()

            features.remove(excluded_feature)

        else:

            break

    return features

```