



"Confidential and property of H2O.ai. All rights reserved"

Target Mean Encoding

what?

Replace categorical variables with the mean of the response

Categorical variables increase the number of features (dimensionality) and can cause us to overfit

why?

H2O Target Mean Encoding

```
h2o.target_encode_apply(data, x, y, target_encode_map, holdout_type,  
                        fold_column = NULL, blended_avg = TRUE,  
                        noise_level = NULL, seed = -1)
```



```
1 def mean_target_encoding(data, x, y, fold_column):  
2     grouped_data = data[[x, fold_column, y]].group_by([x, fold_column])  
3     grouped_data.sum(na = "ignore").count(na = "ignore")  
4     df = grouped_data.get_frame().as_data_frame()  
5     df_list = []  
6     nfold = int(data[fold_column].max()) + 1  
7     for j in range(0, nfold):  
8         te_x = "te_{}".format(x)  
9         sum_y = "sum_{}".format(y)  
10        oof = df.loc[df[fold_column] != j, [x, sum_y, "nrow"]]  
11        stats = oof.groupby([x]).sum()  
12        stats[x] = stats.index  
13        stats[fold_column] = j  
14        stats[te_x] = stats[sum_y] / stats["nrow"]  
15        df_list.append(stats[[x, fold_column, te_x]])  
16    return h2o.H2OFrame(pd.concat(df_list))
```



Target Mean Encoding

What?

Replace categorical variables with the mean of the response

Why?

Categorical variables increase the number of features (dummy encoding) and can cause us to overfit