

As we can see, all eight features' importance picked out by permutation algorithm show a statistical significant, which suggests again, permutation importance method reflects a very decent true feature importance.

## Wine quality Data (Classification)

Lastly, let's test on a classification problem.

```
In [378]: wine_df = pd.read_csv('winequality-white.csv', sep=';')
wine_df['quality'] = wine_df['quality'].apply(lambda x: 1 if x > 5 else 0)

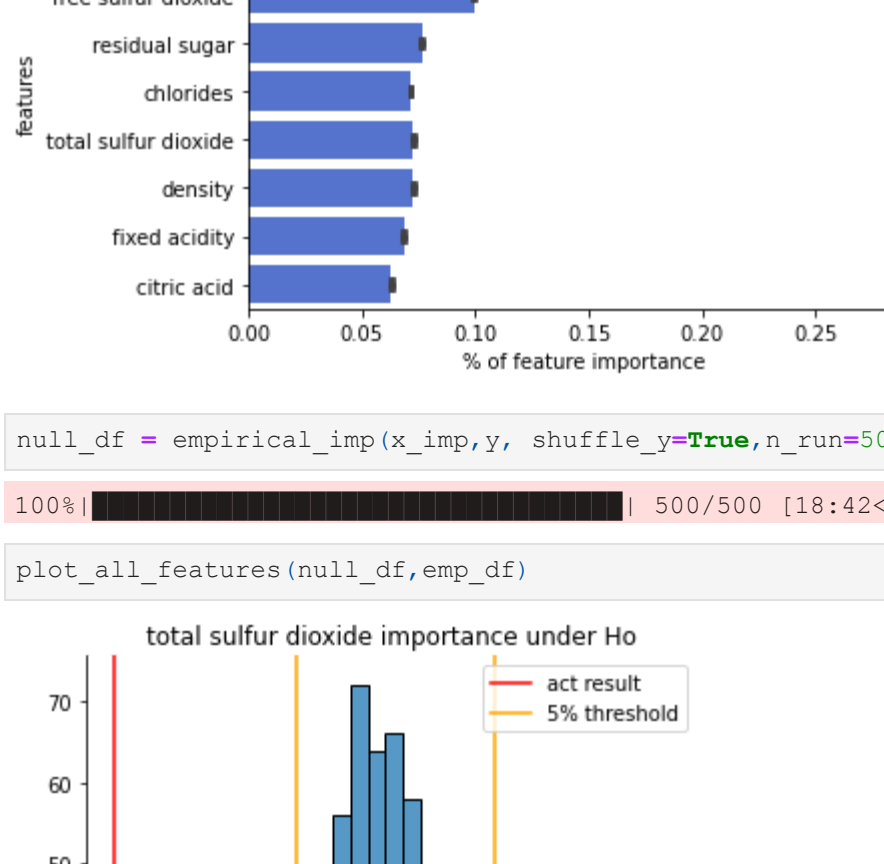
In [381]: x, y = wine_df.iloc[:, :-1], wine_df['quality']

In [387]: x.shape
Out[387]: (4898, 11)
```

## check featimp

```
In [396]: noise = pd.Series(np.random.normal(0, 1, x.shape[0]), name='noise')
X_noise = x.join(noise)
X_noise_scal = pd.DataFrame(StandardScaler().fit_transform(X_noise), columns=list(X_noise))

In [398]: rf = RandomForestClassifier(oob_score=True)
permu_result = permutation_importances(rf, X_noise_scal, y)
visualizing(permu_result)
```

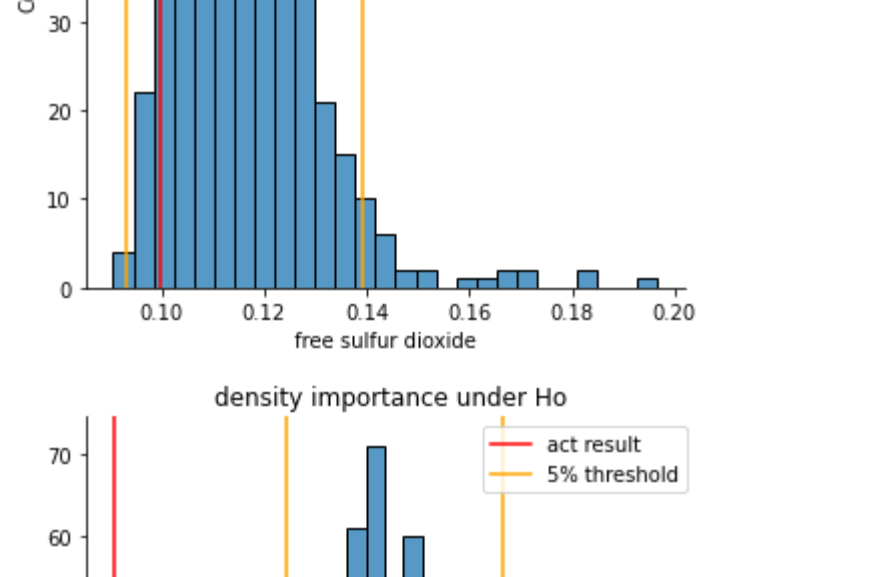


## Automatic feature selection

```
In [401]: features, xs, ys = auto_feat_select(x, y, rf)
print(f"Selected important features are {features}")

Selected important features are ['alcohol', 'volatile acidity', 'free sulfur dioxide', 'density', 'chlorides', 'total sulfur dioxide', 'residual sugar', 'citric acid', 'fixed acidity']

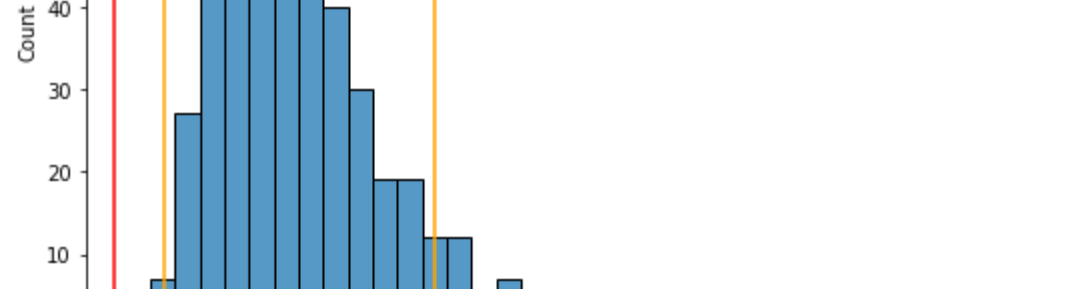
In [402]: sns.pointplot(x = [str(i) for i in xs], y = ys)
plt.xlabel('number of top important features')
plt.ylabel('log_loss')
plt.title('Log loss v.s. number of top important features');
```



## empirical p-values for feature importances

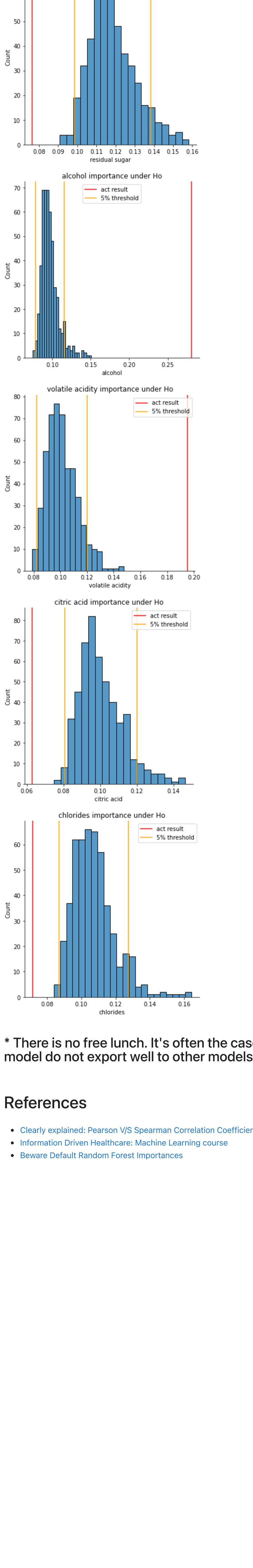
```
In [403]: # take out the important features marked by permutation importance
# (above noise column)
x_imp = StandardScaler().fit_transform(x[features])
x_imp = pd.DataFrame(x_imp, columns=list(x[features]))

In [405]: emp_df = empirical_imp(x_imp, y, n_run=100)
sns.barplot(data=emp_df, orient='h', color='royalblue', cspsize=2)
plt.xlabel('% of feature importance')
plt.ylabel('features')
plt.title('feature important with 95% CI');
```



```
In [407]: null_df = empirical_imp(x_imp, y, shuffle_y=True, n_run=500)

In [408]: plot_all_features(null_df, emp_df)
```



\* There is no free lunch. It's often the case that features determined for one model do not export well to other models.

## References

- Clearly explained: Pearson V/S Spearman Correlation Coefficient
- Information Driven Healthcare: Machine Learning course
- Beware Default Random Forest Importances