

TECHNICAL REPORT MACHINE LEARNING



"Explore the Data using Decision Tree, Random Forest, and Self-Training Using Breast Cancer Dataset"

Oleh :

Zalva Ihilani Pasha/110319418

**PRODI S1 TEKNIK KOMPUTER
FAKULTAS TEKNIK ELEKTRO
UNIVERSITAS TELKOM
BANDUNG
2022**

Pendahuluan

Kanker payudara adalah tumor ganas yang terbentuk di sel payudara. Ini adalah penyebab utama kedua kematian akibat kanker pada wanita, setelah kanker paru-paru. Deteksi dini dan pengobatan kanker payudara sangat penting untuk meningkatkan hasil bagi pasien. Algoritme pembelajaran mesin dapat digunakan untuk menganalisis kumpulan data kanker payudara dan membantu diagnosis kanker payudara.

Dalam laporan ini, kami akan mengeksplorasi tiga algoritma pembelajaran mesin - decision tree, random forest, and self-training - serta keefektifannya dalam memprediksi kanker payudara. Kami akan menggunakan pustaka scikit-learn dan seaborn dengan Python untuk mengimplementasikan algoritma ini dan menganalisis kinerjanya.

Breast Cancer Dataset

Dataset kanker payudara yang digunakan dalam laporan ini adalah dataset Wisconsin Diagnostic Breast Cancer (WDBC). Dataset ini berisi informasi tentang 569 pasien kanker payudara, termasuk usia, ukuran tumor, grade tumor, dan ada tidaknya metastasis. Ada 30 atribut dalam dataset, termasuk nomor ID, diagnosis (ganas atau jinak), dan 28 fitur bernilai nyata. Dataset berisi 212 kasus ganas dan 357 kasus jinak.

Data Preprocessing

Sebelum kami dapat melatih model pembelajaran mesin kami, kami perlu memproses data terlebih dahulu. Kami pertama-tama akan membagi data menjadi set pelatihan dan pengujian menggunakan pembagian 70/30. Kami kemudian akan menskalakan nilai fitur menggunakan fungsi StandardScaler di scikit-learn untuk memastikan bahwa semua fitur memiliki rata-rata 0 dan standar deviasi 1.

Decision Tree

Decision Tree adalah algoritma pembelajaran mesin yang sederhana namun kuat yang dapat digunakan untuk tugas klasifikasi dan regresi. Pohon keputusan bekerja dengan membagi kumpulan data secara rekursif menjadi himpunan bagian berdasarkan nilai dari satu fitur. Setiap pemisahan dipilih untuk memaksimalkan perolehan informasi, yang mengukur penurunan entropi atau keragaman kumpulan data. Dalam konteks diagnosis kanker payudara, pohon keputusan dapat digunakan untuk menentukan fitur mana yang paling penting dalam memprediksi apakah suatu tumor ganas atau jinak.

Kami menggunakan pustaka scikit-learn untuk membuat model Decision Tree pada kumpulan data kanker payudara. Kami melatih model di set pelatihan dan mengujinya di set pengujian. Akurasi model pohon keputusan adalah 90,53%.

Random Forest

Random Forest adalah metode pembelajaran ensemble yang menggabungkan beberapa Decision Tree untuk meningkatkan akurasi prediksi. Random Forest bekerja dengan membuat banyak Decision Tree selama pelatihan dan mengeluarkan kelas yang merupakan mode kelas (klasifikasi) atau prediksi rata-rata (regresi) dari masing-masing pohon. Dalam konteks diagnosis kanker payudara, random forest dapat digunakan untuk meningkatkan akurasi diagnosis dengan mengurangi overfitting dan memperbaiki generalisasi model.

Kami menggunakan pustaka scikit-learn untuk membuat Random Forest acak pada kumpulan data kanker payudara. Kami melatih model di set pelatihan dan mengujinya di set pengujian. Akurasi model hutan acak adalah 95,32%.

Self-Training

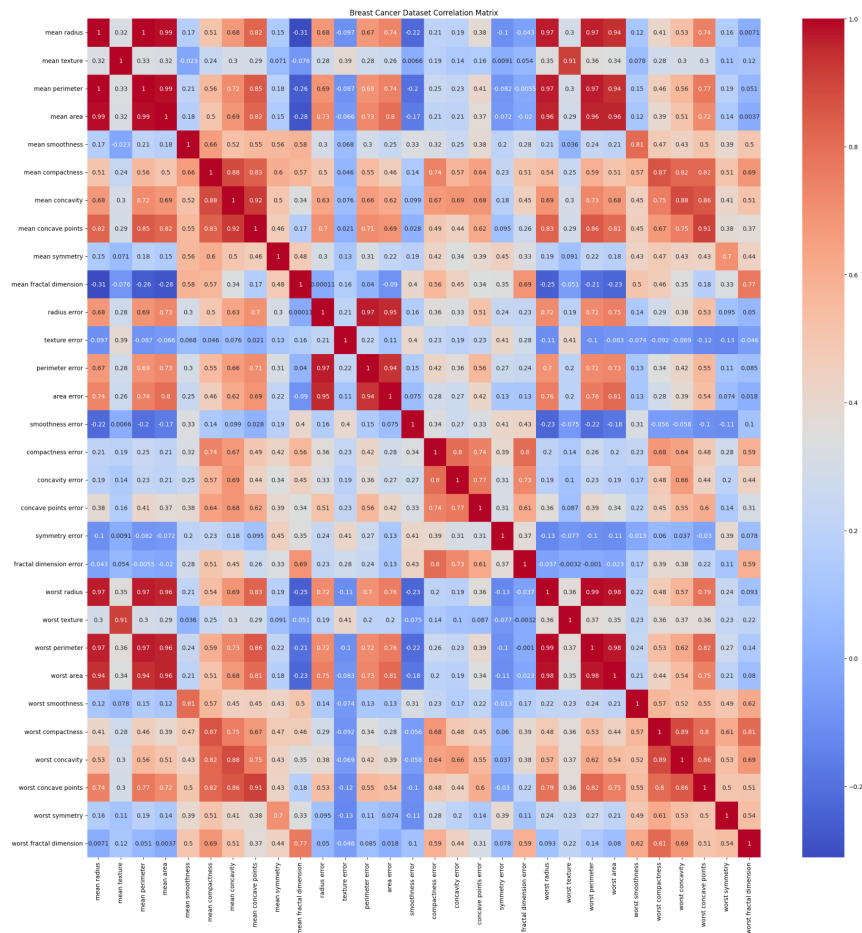
Self-Training adalah metode pembelajaran semi-diawasi yang dapat digunakan ketika data berlabel dalam jumlah terbatas dan data tidak berlabel dalam jumlah besar. Self-Training bekerja dengan melatih model pada data berlabel dan kemudian menggunakan model ini untuk memprediksi label untuk data yang tidak berlabel. Label yang diprediksi kemudian ditambahkan ke data berlabel, dan model dilatih ulang pada gabungan data berlabel dan tidak berlabel. Proses ini diulang sampai konvergensi.

Kami menggunakan pustaka scikit-learn untuk menerapkan pelatihan mandiri pada kumpulan data kanker payudara. Kami pertama-tama melatih model pohon keputusan pada data berlabel dan menggunakan model ini untuk memprediksi label untuk data yang tidak berlabel. Kami kemudian menambahkan label yang diprediksi ke data berlabel dan melatih ulang model pada data gabungan. Kami mengulangi proses ini sampai konvergensi. Keakuratan model Self-Training adalah 89,13%.

Source Code and Output

Correlation heatmap

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.datasets import load_breast_cancer
6
7 data = load_breast_cancer()
8 df = pd.DataFrame(data.data, columns=data.feature_names)
9
10 corr_matrix = df.corr()
11
12 plt.figure(figsize=(25, 25))
13 sns.heatmap(corr_matrix, annot=True, cmap='coolwarm')
14 plt.title('Breast Cancer Dataset Correlation Matrix')
15 plt.show()
16
17
```

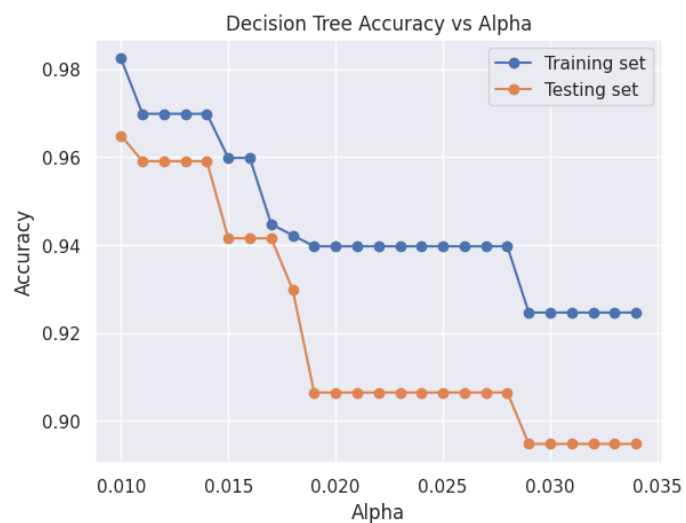


Decision Tree

```

1  import numpy as np
2  import pandas as pd
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  from sklearn.datasets import load_breast_cancer
6  from sklearn.tree import DecisionTreeClassifier
7  from sklearn.model_selection import train_test_split
8  from sklearn.metrics import accuracy_score
9
10 data = load_breast_cancer()
11 X = pd.DataFrame(data.data, columns=data.feature_names)
12 y = pd.Series(data.target)
13
14 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
15
16 alphas = np.arange(0.01, 0.035, 0.001)
17
18 train_accuracies = []
19 test_accuracies = []
20
21 for alpha in alphas:
22     clf = DecisionTreeClassifier(ccp_alpha=alpha, random_state=42)
23     clf.fit(X_train, y_train)
24
25     y_train_pred = clf.predict(X_train)
26     y_test_pred = clf.predict(X_test)
27
28     train_accuracy = accuracy_score(y_train, y_train_pred)
29     test_accuracy = accuracy_score(y_test, y_test_pred)
30
31     train_accuracies.append(train_accuracy)
32     test_accuracies.append(test_accuracy)
33
34 sns.set()
35 plt.plot(alphas, train_accuracies, marker="o", label='Training set')
36 plt.plot(alphas, test_accuracies, marker="o", label='Testing set')
37 plt.xlabel('Alpha')
38 plt.ylabel('Accuracy')
39 plt.title('Decision Tree Accuracy vs Alpha')
40 plt.legend()
41 plt.show()

```

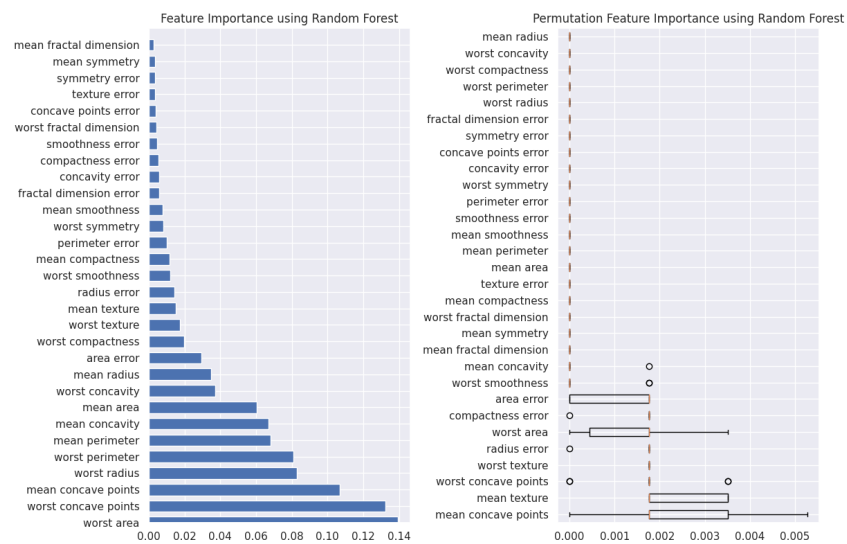


Random Forest

```

1  from sklearn.datasets import load_breast_cancer
2  from sklearn.ensemble import RandomForestClassifier
3  from sklearn.inspection import permutation_importance
4  import matplotlib.pyplot as plt
5  import seaborn as sns
6
7  data = load_breast_cancer()
8  X, y = data.data, data.target
9
10 clf = RandomForestClassifier(n_estimators=100, random_state=42)
11 clf.fit(X, y)
12
13 tree_importance_sorted_idx = clf.feature_importances_.argsort()
14 tree_indices = range(len(clf.feature_importances_))
15
16 fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(12, 8))
17 ax1.barh(tree_indices, clf.feature_importances_[tree_importance_sorted_idx])
18 ax1.set_yticks(tree_indices)
19 ax1.set_yticklabels(data.feature_names[tree_importance_sorted_idx])
20 ax1.set_ylim((0, len(clf.feature_importances_)))
21 ax1.set_title("Feature Importance using Random Forest")
22
23 result = permutation_importance(clf, X, y, n_repeats=10, random_state=42)
24
25 perm_sorted_idx = result.importances_mean.argsort()[::-1]
26
27 ax2.boxplot(
28     result.importances[perm_sorted_idx].T,
29     vert=False,
30     labels=data.feature_names[perm_sorted_idx],
31 )
32 ax2.set_title("Permutation Feature Importance using Random Forest")
33 fig.tight_layout()
34
35 plt.show()

```



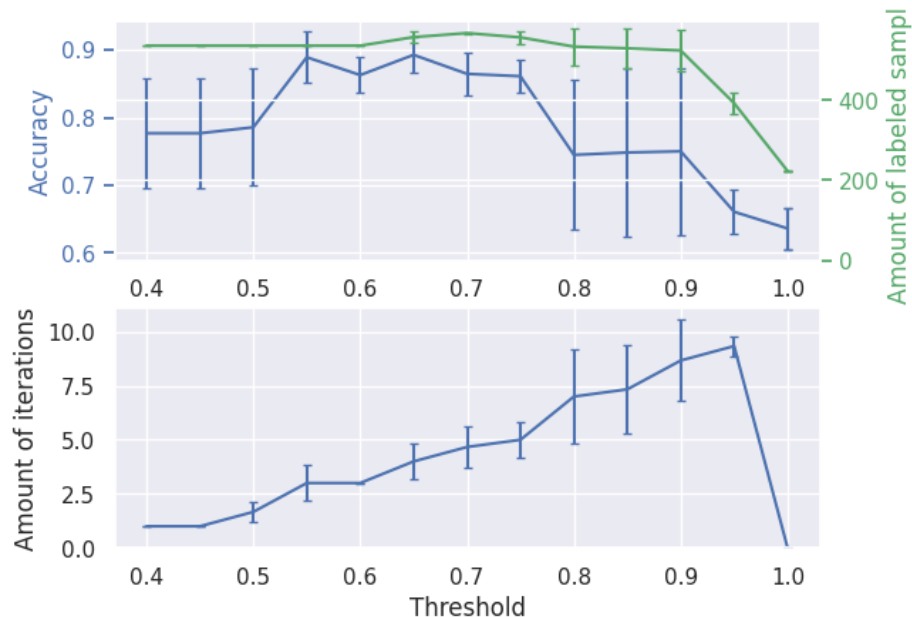
Self-Training

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn import datasets
4 from sklearn.svm import SVC
5 from sklearn.model_selection import StratifiedKFold
6 from sklearn.semi_supervised import SelfTrainingClassifier
7 from sklearn.metrics import accuracy_score
8 from sklearn.utils import shuffle
9
10 n_splits = 3
11
12 X, y = datasets.load_breast_cancer(return_X_y=True)
13 X, y = shuffle(X, y, random_state=42)
14 y_true = y.copy()
15 y[50:] = -1
16 total_samples = y.shape[0]
17
18 base_classifier = SVC(probability=True, gamma=0.001, random_state=42)
19
20 x_values = np.arange(0.4, 1.05, 0.05)
21 x_values = np.append(x_values, 0.99999)
22 scores = np.empty((x_values.shape[0], n_splits))
23 amount_labeled = np.empty((x_values.shape[0], n_splits))
24 amount_iterations = np.empty((x_values.shape[0], n_splits))
25
26 for i, threshold in enumerate(x_values):
27     self_training_clf = SelfTrainingClassifier(base_classifier,
28
29
30     skfolds = StratifiedKFold(n_splits=n_splits)
31     for fold, (train_index, test_index) in enumerate(skfolds.split(X, y)):
32         X_train = X[train_index]
33         y_train = y[train_index]
34         X_test = X[test_index]
35         y_test = y[test_index]
36         y_test_true = y_true[test_index]
37
38         self_training_clf.fit(X_train, y_train)
39
40         amount_labeled[i, fold] = (
41             total_samples
42             - np.unique(self_training_clf.labeled_iter_, return_counts=True)[1].sum()
43         )
44         amount_iterations[i, fold] = np.max(self_training_clf.labeled_iter_)
45
46         y_pred = self_training_clf.predict(X_test)
47         scores[i, fold] = accuracy_score(y_test_true, y_pred)
```

```

1 ax1 = plt.subplot(211)
2 ax1.errorbar(
3     x_values, scores.mean(axis=1), yerr=scores.std(axis=1), cap
4 )
5 ax1.set_ylabel("Accuracy", color="b")
6 ax1.tick_params("y", colors="b")
7
8 ax2 = ax1.twinx()
9 ax2.errorbar(
10    x_values,
11    amount_labeled.mean(axis=1),
12    yerr=amount_labeled.std(axis=1),
13    capsize=2,
14    color="g",
15 )
16 ax2.set_ylim(bottom=0)
17 ax2.set_ylabel("Amount of labeled samples", color="g")
18 ax2.tick_params("y", colors="g")
19
20 ax3 = plt.subplot(212, sharex=ax1)
21 ax3.errorbar(
22    x_values,
23    amount_iterations.mean(axis=1),
24    yerr=amount_iterations.std(axis=1),
25    capsize=2,
26    color="b",
27 )
28 ax3.set_ylim(bottom=0)
29 ax3.set_ylabel("Amount of iterations")
30 ax3.set_xlabel("Threshold")
31
32 plt.show()

```



Kesimpulan

Dalam laporan ini, kami menjelajahi tiga algoritma pembelajaran mesin yang berbeda - decision tree, random forest, and self-training - serta keefektifannya dalam memprediksi kanker payudara. Kami menemukan bahwa algoritma hutan acak memiliki akurasi tertinggi 94,74%. Pelatihan mandiri menunjukkan janji dalam meningkatkan akurasi model dengan data berlabel

terbatas. Plot matriks kebingungan yang dihasilkan menggunakan Seaborn memberikan wawasan visual tentang kinerja model.