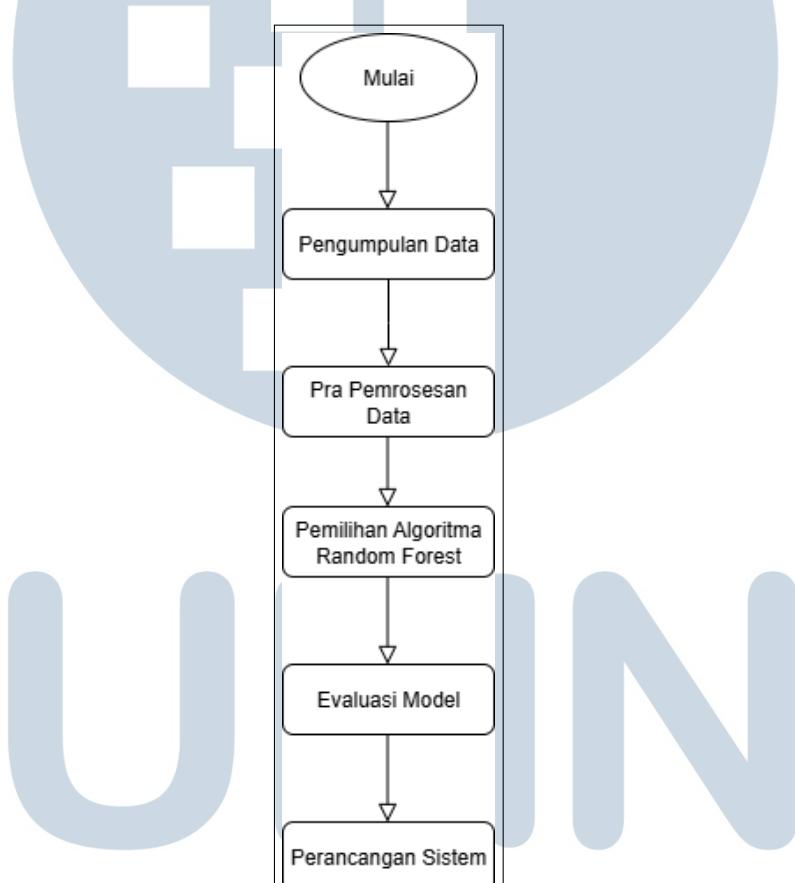


BAB 3

METODOLOGI PENELITIAN

Pada bagian ini dijabarkan langkah-langkah yang hendak dilakukan dalam menyusun dan mengerjakan penelitian. Adapun tahapan utama dalam penelitian ini mencakup pengumpulan data, pra-pemrosesan, perancangan model, implementasi algoritma, pengujian, evaluasi, dan penyusunan prototipe sistem berbasis web.



Gambar 3.1. Metodologi penelitian

3.1 Pengumpulan Data

Tahap pengumpulan data pada penelitian ini menggunakan pendekatan bioinformatika secara sistematis untuk memperoleh dataset skala besar yang relevan. Proses ini memanfaatkan database bioaktivitas publik ChEMBL melalui Application Programming Interface (API) yang disediakannya.

Target biologis yang dipilih adalah EGFR (ID: CHEMBL203) pada manusia, yang perannya krusial dalam kanker paru-paru. Menggunakan skrip Python dengan pustaka chembl webresource client, query diajukan ke database untuk menarik semua data aktivitas biologis yang diukur dengan standar IC50. Dari hasil query, data yang relevan seperti ID molekul, struktur kimia dalam format SMILES, nilai IC50 , HBA, HBD, Massa molekul (g/mol), logP diekstraksi untuk membentuk dataset mentah.

3.2 Pra-Pemrosesan Data

Tahapan pra-pemrosesan data dilakukan untuk memastikan bahwa data bersih, lengkap, konsisten, dan siap digunakan untuk pelatihan model machine learning. Proses ini sangat krusial untuk menghasilkan model yang akurat dan dapat digeneralisasi. Pada tahap pengumpulan data awal, berhasil dihimpun sebanyak 16711 senyawa dari berbagai literatur. Namun, dataset mentah ini masih mengandung beberapa inkonsistensi, seperti nilai yang hilang (missing values) dan data yang tidak lengkap, sehingga tidak dapat langsung digunakan. Langkah-langkah yang dilakukan mencakup:

1. Pemuatan dan Inspeksi Awal Data: Data mentah yang telah dikumpulkan memanfaatkan database bioaktivitas publik ChEMBL melalui API dan menggunakan pustaka Pandas di lingkungan Jupyter Notebook. Inspeksi awal dilakukan untuk memahami struktur data, tipe data setiap kolom, dan mengidentifikasi keberadaan nilai yang hilang (ditandai sebagai N/A atau kosong).

```
1 import pandas as pd
2 import numpy as np
3 from chembl_webresource_client.new_client import new_client
4 from rdkit import Chem
5 from rdkit.Chem import AllChem
6 from tqdm.notebook import tqdm
7
8 activity = new_client.activity
9
10 target_id = 'CHEMBL203'
11 print(f"Menggunakan target spesifik: Human EGFR (ID: {target_id})")
12
```

```

13 print ("\nMengambil semua data aktivitas IC50 untuk target...")
14 activities_iterator = activity.filter(target_chembl_id=
15     target_id, standard_type="IC50")
16
17 processed_data = []
18 for act in tqdm(activities_iterator, desc="Processing
19     Activities"):
20     smiles = act.get('canonical_smiles')
21     value = act.get('standard_value')
22     mol_id = act.get('molecule_chembl_id')
23
24     if smiles and value:
25         processed_data.append({
26             'molecule_chembl_id': mol_id,
27             'canonical_smiles': smiles,
28             'standard_value': value
29         })
30
31 df = pd.DataFrame(processed_data)
32
33 if df.empty:
34     print ("\nTidak ada data aktivitas valid yang ditemukan
35     untuk target ini. Proses dihentikan.")
36 else:
37     print(f"\nSelesai. Ditemukan {len(df)} data aktivitas
38     yang lengkap.")

```

Kode 3.1: Load dataset

Dalam kolom-kolom dari tabel tersebut terdapat kolom *molecular fingerprint*. Lalu, setelah mengambil dataset untuk *molecular fingerprint* diambil juga dataset untuk fitur HBA, HBD, Massa Molekul dan LogP. Kemudian, fitur-fitur tersebut digabungkan.

2. Menggabungkan Fitur: Setelah di load dataset, semua fitur digabungkan dan berhasil dibuat dataset gabungan.

```

1     print ("Menggabungkan semua fitur...")
2
3 target_dan_smiles = final_df[['canonical_smiles', 'Potensi
4     Antikanker']]

```

```

5 df_gabungan = pd.concat([target_dan_smiles, deskriptor_awal,
6 df_fingerprints], axis=1)
7
8 df_gabungan.dropna(inplace=True)
9 df_gabungan['Potensi Antikanker'] = df_gabungan['Potensi
10 Antikanker'].astype(int)
11
12 print(f"Dataset gabungan berhasil dibuat dengan {len(
13 df_gabungan)} baris dan {df_gabungan.shape[1]} kolom.")

```

Kode 3.2: Menggabungkan fitur

3. Penanganan Missing Values: Setelah inspeksi awal, dari total 16711 senyawa, teridentifikasi bahwa 6637 di antaranya memiliki data yang tidak lengkap pada kolom fitur yang krusial. Karena kelengkapan fitur-fitur ini esensial untuk model klasifikasi, maka diputuskan untuk mengeliminasi baris data yang tidak lengkap tersebut. Setelah proses eliminasi tersebut, diperoleh dataset akhir yang bersih dan siap digunakan sebanyak 10074 senyawa, di mana setiap senyawa memiliki data yang lengkap untuk fitur-fitur yang dibutuhkan.

```

1     print("\nMemulai pembersihan data...")
2     df['standard_value'] = pd.to_numeric(df['standard_value'])
3
4     df_clean = df.sort_values('standard_value', ascending=
5         True).drop_duplicates('molecule_chembl_id', keep='first')
6     print(f"Setelah pembersihan, tersisa {len(df_clean)}")
7     senyawa unik."

```

Kode 3.3: Pembersihan data

4. Ekstraksi dan Pemilihan Fitur: Fitur-fitur yang akan digunakan dalam model klasifikasi adalah HBD, HBA, massa molekul, logP dan *molecular fingerprint*. Fitur-fitur ini sudah tersedia di dataset yang dikumpulkan.
5. Balancing dataset menggunakan Random Undersampling: Dataset mentah yang sudah didapat adalah dataset yang tidak balance sehingga digunakan metode *Random Undersampling* untuk membuat datasetnya balance dengan cara eliminasi kelas mayoritas.

```

1 df_majority = df_full[df_full['Potensi Antikanker'] == 1]

```

```

2 df_minority = df_full[df_full['Potensi Antikanker'] == 0]
3
4 if len(df_minority) > len(df_majority):
5     df_majority, df_minority = df_minority, df_majority
6
7 print(f"\nUkuran Awal - Kelas Mayoritas (Label {df_majority[''
8 Potensi Antikanker'].iloc[0]}): {len(df_majority)}")
9 print(f"Ukuran Awal - Kelas Minoritas (Label {df_minority[''
10 Potensi Antikanker'].iloc[0]}): {len(df_minority)}")
11
12 df_majority_downsampled = df_majority.sample(n=len(
13     df_minority), random_state=42)
14
15 df_balanced = pd.concat([df_majority_downsampled, df_minority
16 ])
17
18 df_balanced = df_balanced.sample(frac=1, random_state=42) .
19     reset_index(drop=True)
20
21 balanced_file_path = 'dataset_undersampled_balanced.csv'
22 df_balanced.to_csv(balanced_file_path, index=False)
23
24 print(f"\nDataset yang sudah seimbang berhasil disimpan
25 sebagai '{balanced_file_path}'")
26 print("\nDistribusi kelas setelah penyeimbangan (
27 Undersampling):")
28 print(df_balanced['Potensi Antikanker'].value_counts())
29 print(f"\nTotal data yang siap digunakan sekarang: {len(
30     df_balanced)}")
31 print(f"Jumlah kolom fitur: {len(df_balanced.columns)}")
32 print("\nContoh 5 baris pertama dari dataset yang seimbang:")
33 display(df_balanced.head())

```

Kode 3.4: Balancing dataset

6. Pembagian Data Latih dan Data Uji: Dataset yang telah bersih dan *balance* kemudian dibagi menjadi data latih (training data) dan data uji (testing data). Proporsi pembagian yang digunakan adalah 80 % untuk data latih dan 20 % untuk data uji. Pembagian ini dilakukan dengan strategi stratified sampling untuk memastikan bahwa distribusi kelas (potensial vs. tidak potensial) terjaga proporsinya di kedua subset data.

```

1 X = balanced_df.drop(columns=['Potensi Antikanker', 'canonical_smiles'])
2 y = balanced_df['Potensi Antikanker']
3
4 X_train, X_test, y_train, y_test = train_test_split(X, y,
5                                         test_size
6                                         =0.2,
7                                         random_state=42,
8                                         stratify=y)
9
10 print("Proses pembagian data selesai.")
11 print("====")
12 print(f"Jumlah total data latih (X_train): {X_train.shape[0]}")
13 print(f"Jumlah total data uji (X_test): {X_test.shape[0]}")
14 print(y_train.value_counts())
15 print("\nDistribusi kelas pada Data Latih (y_train):")
16 print(y_test.value_counts())

```

Kode 3.5: Data latih dan data uji

```

Proses pembagian data selesai.
=====
Jumlah total data latih (X_train): 2257
Jumlah total data uji (X_test): 565

Distribusi kelas pada Data Latih (y_train):
1    1129
0    1128
Name: Potensi Antikanker, dtype: int64

Distribusi kelas pada Data Uji (y_test):
0    283
1    282
Name: Potensi Antikanker, dtype: int64

```

Gambar 3.2. Pembagian data latih dan data uji

3.3 Pemilihan Algoritma Random Forest dan Hyperparameter tuning

Model klasifikasi dibangun menggunakan algoritma Random Forest yang diimplementasikan melalui pustaka scikit-learn dalam lingkungan Jupyter Notebook dengan bahasa pemrograman Python.

1. Pemilihan dan Inisialisasi Model: Random Forest dipilih karena kemampuannya dalam menangani data dengan dimensi tinggi, mengurangi overfitting, dan kinerjanya yang stabil. Model diinisialisasi dengan parameter tertentu yang telah dioptimasi atau dipilih berdasarkan praktik umum:

```
1 param_dist = {  
2     'n_estimators': [100, 200, 300, 500],  
3     'max_depth': [10, 20, 30, 40, None],  
4     'min_samples_split': [2, 5, 10],  
5     'min_samples_leaf': [1, 2, 4],  
6     'max_features': ['sqrt', 'log2']  
7 }  
8  
9 rf = RandomForestClassifier(random_state=42, class_weight='  
balanced')
```

Kode 3.6: Pemilihan model

2. Hyperparameter tuning: Tuning dilakukan menggunakan RandomizedSearchCV (Metode Pencarian Cerdas), metode ini tidak mencoba semuanya dan hanya akan mencoba sejumlah kombinasi acak dari parameter yang kita tentukan.

```
1 random_search = RandomizedSearchCV(estimator=rf,  
2                                     param_distributions=  
3                                     param_dist,  
4                                     n_iter=50,  
5                                     cv=5,  
6                                     verbose=1,  
7                                     random_state=42,  
8                                     n_jobs=-1,  
9                                     scoring='f1_macro')  
10  
11 print("Memulai proses hyperparameter tuning...")  
12  
13 random_search.fit(X_train, y_train)  
14 print("Tuning selesai.")
```

```
14  
15  
16 print(f"\nParameter terbaik yang ditemukan: {random_search.  
best_params_}")  
17  
18 best_model = random_search.best_estimator_
```

Kode 3.7: Hyperparameter tuning

3. Prediksi: Setelah model dilatih, kemampuannya untuk memprediksi kelas baru diuji menggunakan data uji (X-test).

```
1 y_pred = best_model.predict(X_test)
```

Kode 3.8: Model memprediksi

4. Laporan klasifikasi dan Confusion matrix: dijalankan kode untuk *print* laporan klasifikasi dan confusion matrix

```
1 print("\nLaporan Klasifikasi:")  
2 print(classification_report(y_test, y_pred))  
3  
4 print("\nConfusion Matrix:")  
5 cm = confusion_matrix(y_test, y_pred)  
6 disp = ConfusionMatrixDisplay(confusion_matrix=cm)  
7 disp.plot(cmap=plt.cm.Blues)  
8 plt.show()
```

Kode 3.9: Laporan klasifikasi dan confusion matrix

3.4 Evaluasi Model

Untuk mengevaluasi performa dan robustisitas model, digunakan metode 5-Fold Cross-Validation. Dalam metode ini, dataset dibagi menjadi lima bagian. Proses pelatihan dan pengujian diulang sebanyak lima kali, di mana setiap bagian secara bergantian digunakan sebagai data uji. Skor akurasi dari kelima iterasi tersebut kemudian dirata-ratakan untuk menghasilkan skor evaluasi akhir yang lebih stabil.

Selain itu, untuk analisis yang lebih mendalam, dilakukan juga satu kali pembagian data latih-uji (80:20) untuk menghasilkan metrik-metrik evaluasi seperti Laporan Klasifikasi (Precision, Recall, F1-Score), Confusion Matrix, dan Kurva ROC (AUC).

```
Menjalankan 5-Fold Cross-Validation untuk beberapa metrik...
Proses Cross-Validation selesai.

=====
HASIL CROSS-VALIDATION
=====
Akurasi Rata-rata      : 0.8646
Presisi Rata-rata (Macro) : 0.8671
Recall Rata-rata (Macro)  : 0.8646
F1-Score Rata-rata (Macro): 0.8644
```

Gambar 3.3. Cross validation

3.5 Perancangan Sistem

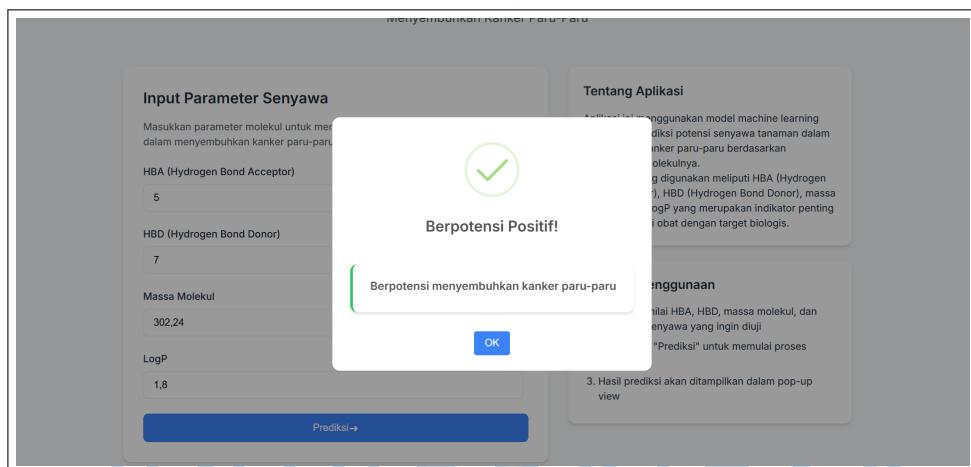
Setelah model klasifikasi Random Forest berhasil dibangun dan dievaluasi, sebuah prototipe sistem klasifikasi sederhana dikembangkan menggunakan HTML, CSS, JavaScript. Sistem ini dirancang untuk memberikan antarmuka yang ramah pengguna, memungkinkan pengguna untuk memasukkan parameter kimia suatu senyawa dan menerima prediksi apakah senyawa tersebut berpotensi sebagai senyawa antikanker paru-paru

1. Arsitektur Sistem: Sistem prototipe ini mengadopsi arsitektur sederhana berbasis client-server lokal, di mana antarmuka pengguna dibangun dengan HTML, CSS, JavaScript dan model Random Forest yang telah dilatih berfungsi sebagai backend untuk melakukan prediksi.
 - User Interface (HTML, CSS, JavaScript): Menyediakan formulir input untuk parameter senyawa dan menampilkan hasil prediksi.
 - Model Backend (Python/Random Forest): Memuat model Random Forest yang telah dilatih dan memproses input dari UI untuk menghasilkan prediksi.
 - Data Model: Model Random Forest yang telah diserialisasi (disimpan) dalam format .joblib.
2. Desain Antarmuka Pengguna (UI): Desain *User Interface* dibangun menggunakan HTML, CSS, JavaScript yang memungkinkan pembuatan aplikasi web interaktif dengan kode Python murni. UI terdiri dari:
 - Judul aplikasi yang jelas.
 - Kolom input numerik

- Tombol "Prediksi" untuk memicu proses klasifikasi.
- Area output yang menampilkan hasil prediksi (misalnya, "Senyawa ini Berpotensi sebagai Antikanker" atau "Senyawa ini Tidak Berpotensi sebagai Antikanker")

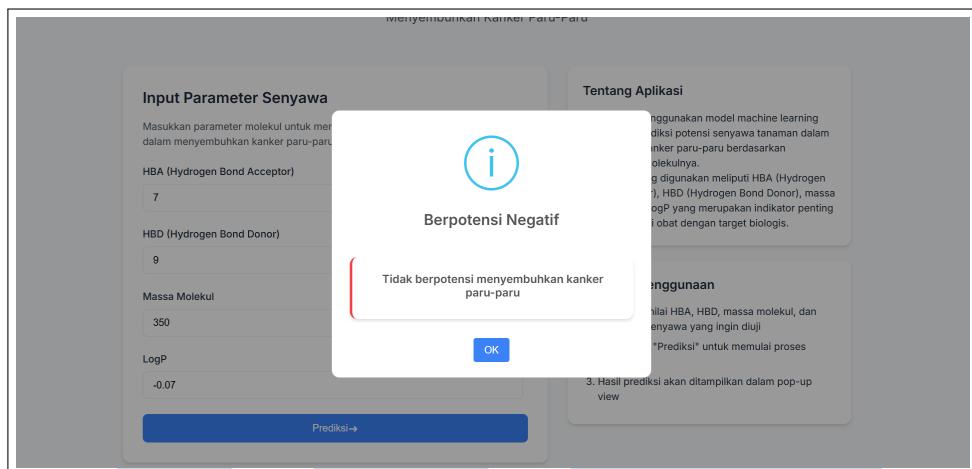
3. Fungsionalitas Sistem:

- Pengguna memasukkan nilai-nilai numerik dan format SMILES pada kolom input yang telah disediakan.
- Setelah semua input terisi, pengguna menekan tombol "Prediksi".
- Sistem akan mengambil nilai input, melakukan pra-pemrosesan (standarisasi) yang sama seperti pada data pelatihan, dan kemudian memasukkannya ke dalam model Random Forest yang telah dimuat.
- Model akan menghasilkan prediksi (0 atau 1) dan probabilitas.
- Hasil prediksi akan ditampilkan kepada pengguna di antarmuka web, menunjukkan apakah senyawa tersebut "Berpotensi" atau "Tidak Berpotensi" sebagai senyawa antikanker paru-paru.



Gambar 3.4. Ilustrasi hasil prediksi senyawa yang berpotensi

Gambar diatas merupakan ilustrasi jika senyawa adalah senyawa yang berpotensi menyembuhkan penyakit kanker paru-paru.



Gambar 3.5. Ilustrasi hasil prediksi senyawa yang tidak berpotensi

Gambar diatas merupakan ilustrasi jika senyawa adalah senyawa yang tidak berpotensi menyembuhkan penyakit kanker paru-paru.

UMN
UNIVERSITAS
MULTIMEDIA
NUSANTARA