



UNIVERSITAS INDONESIA

**PENGEMBANGAN WEB CRAWLER UNTUK PENDETEKSIAN FRAUD PADA
SITUS MERCHANT PT NUSA SATU INTI ARTHA**

SKRIPSI

**BAHY HELMI HARTOYO PUTRA
1606918124**

**FAKULTAS ILMU KOMPUTER
PROGRAM STUDI SISTEM INFORMASI
DEPOK
BULAN TAHUN**



UNIVERSITAS INDONESIA

**PENGEMBANGAN WEB CRAWLER UNTUK PENDETEKSIAN FRAUD PADA
SITUS MERCHANT PT NUSA SATU INTI ARTHA**

SKRIPSI

**Diajukan sebagai salah satu syarat untuk memperoleh gelar
S.Kom**

**BAHY HELMI HARTOYO PUTRA
1606918124**

**FAKULTAS ILMU KOMPUTER
PROGRAM STUDI SISTEM INFORMASI
DEPOK
BULAN TAHUN**

HALAMAN PERSETUJUAN

Judul : Pengembangan Web Crawler untuk Pendeteksian Fraud pada Situs Merchant PT Nusa Satu Inti Artha
Nama : Bahy Helmi Hartoyo Putra
NPM : 1606918124

Laporan Skripsi ini telah diperiksa dan disetujui.

Tanggal Bulan Tahun

Adila Alfa Krisnadhi S.Kom., M.Sc

Pembimbing Skripsi

Pembimbing Kedua Anda

Pembimbing Skripsi

HALAMAN PERNYATAAN ORISINALITAS

**Skripsi ini adalah hasil karya saya sendiri,
dan semua sumber baik yang dikutip maupun dirujuk
telah saya nyatakan dengan benar.**

Nama : Bahy Helmi Hartoyo Putra
NPM : 1606918124
Tanda Tangan :

Tanggal : Tanggal Bulan Tahun

HALAMAN PENGESAHAN

Skripsi ini diajukan oleh :

Nama : Bahy Helmi Hartoyo Putra

NPM : 1606918124

Program Studi : Sistem Informasi

Judul Skripsi : Pengembangan Web Crawler untuk Pendeteksian
Fraud pada Situs Merchant PT Nusa Satu Inti Artha

Telah berhasil dipertahankan di hadapan Dewan Penguji dan diterima sebagai bagian persyaratan yang diperlukan untuk memperoleh gelar S.Kom pada Program Studi Sistem Informasi, Fakultas Ilmu Komputer, Universitas Indonesia.

DEWAN PENGUJI

Pembimbing 1 : Adila Alfa Krisnadhi S.Kom., M.Sc ()

Pembimbing 2 : Pembimbing Kedua Anda ()

Penguji 1 : Penguji Pertama Anda ()

Penguji 2 : Penguji Kedua Anda ()

Ditetapkan di : Depok

Tanggal : Tanggal Bulan Tahun

HALAMAN PERNYATAAN PERSETUJUAN PUBLIKASI TUGAS AKHIR UNTUK KEPENTINGAN AKADEMIS

Sebagai sivitas akademik Universitas Indonesia, saya yang bertanda tangan di bawah ini:

Nama : Bahy Helmi Hartoyo Putra
NPM : 1606918124
Program Studi : Sistem Informasi
Fakultas : Ilmu Komputer
Jenis Karya : Skripsi

demikian pengembangan ilmu pengetahuan, menyetujui untuk memberikan kepada Universitas Indonesia **Hak Bebas Royalti Noneksklusif (Non-exclusive Royalty Free Right)** atas karya ilmiah saya yang berjudul:

Pengembangan Web Crawler untuk Pendeteksian Fraud pada Situs Merchant PT Nusa
Satu Inti Artha

berserta perangkat yang ada (jika diperlukan). Dengan Hak Bebas Royalti Noneksklusif ini Universitas Indonesia berhak menyimpan, mengalihmedia/formatkan, mengelola dalam bentuk pangkalan data (database), merawat, dan memublikasikan tugas akhir saya selama tetap mencantumkan nama saya sebagai penulis/pencipta dan sebagai pemilik Hak Cipta.

Demikian pernyataan ini saya buat dengan sebenarnya.

Dibuat di : Depok
Pada tanggal : Tanggal Bulan Tahun
Yang menyatakan

(Bahy Helmi Hartoyo Putra)

KATA PENGANTAR

Pendahuluan. Ucapan Terima Kasih:

1. Pembimbing.
2. Dosen.
3. Instansi.
4. Orang tua.
5. Sahabat.
6. Teman.

Penulis menyadari bahwa laporan Skripsi ini masih jauh dari sempurna. Oleh karena itu, apabila terdapat kesalahan atau kekurangan dalam laporan ini, Penulis memohon agar kritik dan saran bisa disampaikan langsung melalui *e-mail* emailanda@mail.id.

Depok, Tanggal Bulan Tahun

Bahy Helmi Hartoyo Putra

ABSTRAK

Nama : Bahy Helmi Hartoyo Putra
Program Studi : Sistem Informasi
Judul : Pengembangan Web Crawler untuk Pendeteksian Fraud
pada Situs Merchant PT Nusa Satu Inti Artha

Isi abstrak.

Kata kunci:

Keyword satu, kata kunci dua

ABSTRACT

Name : Bahy Helmi Hartoyo Putra
Program : Sistem Informasi
Title : Development of Web Crawler for Fraud Detection on PT Nusa
Satu Inti Artha Merchant Site

Abstract content.

Key words:

Keyword one, keyword two

DAFTAR ISI

HALAMAN JUDUL	i
LEMBAR PERSETUJUAN	ii
LEMBAR PERNYATAAN ORISINALITAS	iii
LEMBAR PENGESAHAN	iv
LEMBAR PERSETUJUAN PUBLIKASI ILMIAH	v
KATA PENGANTAR	vi
ABSTRAK	vii
Daftar Isi	ix
Daftar Gambar	xii
Daftar Tabel	xiii
Daftar Kode Program	xiv
1 PENDAHULUAN	1
1.1 Latar Belakang	1
1.2 Permasalahan	3
1.2.1 Definisi Permasalahan	3
1.2.2 Batasan Permasalahan	3
1.3 Tujuan Penelitian	3
1.4 Tahapan Penelitian	4
1.5 Sistematika Penulisan	4
2 STUDI LITERATUR	5
2.1 <i>Financial Technology</i>	5
2.1.1 <i>Payment Financial Technology</i>	5
2.1.2 <i>Know Your Business</i>	5
2.2 <i>Web Crawler</i>	6
2.2.1 <i>Crawl Frontier</i>	6
2.2.2 <i>Web Driver</i>	6
2.2.3 <i>HTTP Programming</i>	6
2.3 <i>Web Scraping</i>	7
2.3.1 <i>HTML Parsing</i>	7
2.3.2 <i>Regular Expression</i>	8

2.4	<i>Classifier</i>	10
2.4.1	Decision Tree Classifier	10
2.4.2	Random Forest Classifier	12
2.4.3	Extreme Gradient Boost Classifier	13
2.4.4	Bernoulli Naive Bayes Classifier	14
2.5	Validasi model	15
2.5.1	<i>Nested Cross Validation</i>	16
2.5.2	Metrik Evaluasi	17
3	DESAIN DAN IMPLEMENTASI	18
3.1	Arsitektur Web Crawler	18
3.2	Proses Pelatihan Model	19
3.3	Modul <i>Helper</i>	20
3.3.1	Inisiasi & Atur Ulang Browser	20
3.3.2	<i>URL Format Handler</i>	21
3.3.3	<i>Hyperlink Gatherer</i>	22
3.3.3.1	<i>Static Hyperlink Gatherer</i>	22
3.3.3.2	<i>Dynamic Hyperlink Gatherer</i>	24
3.3.4	<i>Paragraph Extractor</i>	25
3.3.4.1	Static Paragraph Extractor	26
3.3.4.2	Dynamic Paragraph Extractor	27
3.4	Modul Ekstraksi Fitur	30
3.4.1	<i>Broken Link Checker</i>	30
3.4.2	<i>TnC Checker</i>	32
3.4.3	<i>About Us Checker</i>	34
3.4.4	<i>Contact Info Checker</i>	35
3.5	Model Scorer	37
4	EKSPERIMEN DAN ANALISIS	39
4.1	Alur Kerja Eksperimen	39
4.2	Evaluasi Kinerja Model	42
4.2.1	<i>Decision Tree Classifier</i>	43
4.2.1.1	Parameter Terbaik Hasil Grid Search pada Inner Loop	44
4.2.1.2	Hasil Test pada Outer Loop dan Test Data	44
4.2.2	<i>Random Forest Classifier</i>	45
4.2.2.1	Parameter Terbaik Hasil Grid Search pada Inner Loop	45
4.2.2.2	Hasil Test pada Outer Loop dan Test Data	46
4.2.3	<i>Extreme Gradient Boost Classifier</i>	46
4.2.3.1	Parameter Terbaik Hasil Grid Search pada Inner Loop	47
4.2.3.2	Hasil Test pada Outer Loop dan Test Data	47
4.2.4	<i>Bernoulli Naive Bayes Classifier</i>	48
4.2.4.1	Parameter Terbaik Hasil Grid Search pada Inner Loop	48
4.2.4.2	Hasil Test pada Outer Loop dan Test Data	49
4.2.5	Kinerja Model Terbaik	50
4.2.6	Kombinasi Hyperparameter Terbaik	50

5	PENUTUP	54
5.1	Kesimpulan	54
5.2	Saran	55
	Daftar Referensi	56
	LAMPIRAN	1
	Lampiran 1: Judul Lampiran 1	2

DAFTAR GAMBAR

2.1	Struktur <i>tree</i> pada Decision Tree Classifier, diambil dari DataCamp Community	10
2.2	Cara menghitung <i>gini impurity</i> pada Decision Tree Classifier	11
2.3	<i>Training time</i> dan performa XGBoost dibandingkan model lain	13
3.1	Arsitektur Web Crawler	18
3.2	Proses Training	19
4.1	Skema <i>nested cross-validation</i>	43
4.2	Rata-rata skor AUC yang didapatkan setiap model pada <i>outer loop</i>	50

DAFTAR TABEL

2.1	Contoh <i>response code</i> yang dikembalikan <i>HTTP response message</i>	8
2.2	Contoh penggunaan syntax <i>regex</i>	9
3.1	Konversi masukan URL menjadi keluaran yang sesuai	22
4.1	Struktur tabel yang dikirimkan oleh tim EDU	39
4.2	Fitur yang diekstrak dari proses <i>feature engineering</i>	41
4.3	Kombinasi parameter terbaik dari <i>inner loop</i> model Decision Tree Classifier	44
4.4	Percobaan kombinasi parameter terbaik ke <i>outer loop & test data</i> model Decision Tree Classifier	44
4.5	Kombinasi parameter terbaik dari <i>inner loop</i> model Random Forest Clas- sifier	45
4.6	Percobaan kombinasi parameter terbaik ke <i>outer loop & test data</i> model Random Forest Classifier	46
4.7	Kombinasi parameter terbaik dari <i>inner loop</i> model XGBoost Classifier .	47
4.8	Percobaan kombinasi parameter terbaik ke <i>outer loop & test data</i> model XGBoost Classifier	47
4.9	Kombinasi parameter terbaik dari <i>inner loop</i> model Bernoulli Naive Bayes Classifier	48
4.10	Percobaan kombinasi parameter terbaik ke <i>outer loop & test data</i> model Bernoulli Naive Bayes Classifier	49
4.11	Ringkasan metrik terbaik keempat model	51
4.12	Ringkasan <i>confusion matrix</i> terbaik keempat model	51

DAFTAR KODE PROGRAM

3.1	Fungsi inisiasi dan atur ulang <i>browser</i>	20
3.2	Fungsi <i>hyperlink gathering</i>	22
3.3	Fungsi <i>hyperlink gathering dynamic</i>	24
3.4	Fungsi <i>paragraph extractor</i>	26
3.5	Fungsi <i>paragraph extractor dynamic</i>	28
3.6	Fungsi <i>broken link checker</i>	30
3.7	Fungsi <i>refund policy matcher</i>	33
3.8	Fungsi <i>about us check</i>	34
3.9	Fungsi <i>telephone matcher</i>	36
3.10	Fungsi <i>email matcher</i>	36
3.11	Fungsi <i>model scorer</i>	37

BAB 1

PENDAHULUAN

Pada bab ini, akan dijelaskan tentang latar belakang dan permasalahan yang diselesaikan pada penelitian ini.

1.1 Latar Belakang

Financial technology atau yang biasa disebut *fintech* merupakan industri yang tengah berkembang pesat di Indonesia. Kemajuan yang signifikan ini ditunjukkan oleh dua sektor dominan yang merajai pasar *fintech* Indonesia, yaitu sektor pembayaran dan pendanaan. Berdasarkan data yang dihimpun oleh Asosiasi Fintech Indonesia, laju pertumbuhan nilai transaksi meningkat sebanyak 15% dari Juni 2018 ke Juni 2019 pada sektor pembayaran. Sedangkan pada sektor pendanaan, laju pertumbuhan meningkat secara drastis sebanyak 97.6% dari Desember 2018 ke Juni 2019 [1].

Berkembangnya layanan *fintech* mendorong perusahaan-perusahaan baru bermunculan dan ikut serta dalam menyelenggarakan layanan *fintech*. Pada Desember 2019, tercatat sebanyak 61 perusahaan telah tergabung pada sektor pembayaran dan 164 perusahaan pada sektor pendanaan [1]. Meskipun sektor pembayaran bukan merupakan sektor *fintech* dengan pemain dan pertumbuhan terbesar saat ini, namun potensi untuk berkembang pada sektor ini sangat memungkinkan untuk terjadi.

Sebuah riset yang dikeluarkan oleh Metra Digital Innovation Venture dan Mandiri Capital memprediksi bahwa pada tahun 2020 angka *Gross Transaction Volume* (GTV) industri pembayaran *mobile* di Indonesia akan mencapai US\$30 miliar. Hal ini menunjukkan kenaikan yang cukup signifikan sebesar 159% dari tahun 2016. Selain itu, berkembangnya industri pembayaran *mobile* ini juga diprediksi akan berkontribusi terhadap total *Gross Domestic Product* (GDP) Indonesia sebanyak 3% [2].

Dari sekian banyak perusahaan *fintech* yang bergerak pada sektor pembayaran, salah satu yang cukup lama berada di sektor ini adalah PT Nusa Satu Inti Artha, atau lebih dikenal dengan DOKU, DOKU berdiri pada tahun 2007 dan sejak saat itu fokus perusahaan adalah mengembangkan dan menyediakan layanan pembayaran secara *online* kepada pelaku bisnis dan masyarakat. Saat ini, DOKU telah terkoneksi dengan lebih dari 20 bank dan institusi finansial di Indonesia. Layanan pembayaran DOKU meliputi penyediaan metode pembayaran bagi *merchant* dan pelanggannya (*payment gateway*) dan layanan transfer melalui berbagai *channel* seperti kartu kredit, transfer bank, *e-wallet* dan

juga pembayaran *offline* di beberapa toko-toko pilihan.

Saat ini DOKU telah digunakan oleh lebih dari 100.000 *merchant* dalam kedua layanannya, yaitu *payment gateway* dan juga *transfer service*. Untuk dapat bergabung menjadi *merchant* yang menggunakan layanan DOKU, pemilik *merchant* perlu melalui serangkaian proses registrasi. Proses registrasi tersebut dinamakan *Know Your Business* (KYB). Dalam mengoptimalkan proses registrasi ini, DOKU membentuk sebuah tim khusus untuk mengurus proses *end-to-end* registrasi, yang dinamakan dengan tim EDU atau *Early Detection Unit*. Tugas utama dari tim ini adalah mengevaluasi *merchant* yang melakukan proses registrasi dan memastikan bahwa *merchant* tersebut telah memenuhi syarat & ketentuan yang ditentukan oleh DOKU. Jika *merchant* dianggap telah sesuai, maka *merhcant* dapat melanjutkan ke proses integrasi sistem pembayaran, namun apabila tidak sesuai maka *merchant* akan diberikan peringatan dan diminta keterangan lebih lanjut.

Salah satu hal penting dalam proses pemeriksaan ini adalah kelengkapan komponen situs *merchant* yang bersangkutan. Ada beberapa syarat yang ditetapkan oleh tim EDU untuk sebuah situs agar dapat dianggap sebagai situs yang valid. Diantaranya adalah pemeriksaan ketersediaan tautan utama (*contact us, about us, terms & condition*), jumlah tautan tidak aktif, ketersediaan informasi kontak, dan juga adanya *policy* terkait pengembalian barang dan uang. Keseluruhan proses ini masih dilakukan secara manual oleh anggota tim EDU. Menurut keterangan pihak DOKU¹, *rate* pemeriksaan sebuah situs yang berada pada *waiting list* registrasi berkisar antara 35 - 50 situs/hari, atau sekitar 5 situs/jam². Selain itu, untuk mencapai *rate* tersebut dibutuhkan 2 orang tim EDU yang ikut serta dalam pengecekan. *Rate* pengecekan ini dianggap masih terlalu lambat dan masih melibatkan terlalu banyak sumber daya manusia. Di satu sisi diketahui juga bahwa rata-rata *merchant* yang mendaftar mencapai 100 *merchant* per hari. Oleh karena itu, pada tahun 2020 tim DOKU mencoba untuk melakukan otomasi terhadap proses ini, dengan harapan *rate* pemeriksaan harian dapat meningkat dan beban pekerjaan yang dimiliki tim EDU dapat terbantu.

Dengan adanya *web crawler*, rata-rata pengecekan satu buah situs hanya memakan waktu sekitar 15 detik³. Dengan *rate* 15 detik/situs, maka dalam satu jam *web crawler* dapat melakukan pengecekan sebanyak 240 situs. Sehingga, jika dikonversikan ke *rate* harian, dengan menggunakan *web crawler*, tim DOKU dapat mengecek sebanyak 1620 situs. Angka ini jauh meningkat dibandingkan pengecekan manual. Sehingga, efisiensi

¹Diwakili oleh Pak Reza Farasdak Abdat, Product Owner DOKU

²Dengan perhitungan 9 jam kerja dan dengan asumsi bahwa seluruh jam kerja staf yang melakukan pengecekan didedikasikan untuk ini

³Diambil dari rata-rata percobaan yang dilakukan bersama tim DOKU terhadap 100 situs *merchant* yang disampel secara *random* dengan kecepatan internet 20 MBps

rate pengecekan dapat ditingkatkan sebesar 4700%⁴.

Pengecekan yang dilakukan oleh *web crawler* meliputi pengecekan kelengkapan komponen situs dan prediksi kategori situs (*fraud* atau tidak *fraud*). Akurasi model *classifier* yang dimiliki oleh *web crawler* juga sudah cukup baik dengan angka 0.864 dan AUC sebesar 0.953.

1.2 Permasalahan

Sebutkan permasalahan penelitian Anda dari latar belakang tersebut.

1.2.1 Definisi Permasalahan

Berikut ini adalah rumusan permasalahan dari penelitian yang dilakukan:

- Bagaimana cara membuat sebuah *web crawler* untuk melakukan ekstraksi data pada situs *merchant*?
- Fitur apa saja yang dapat digunakan untuk mendeteksi situs *merchant* yang *fraud*?
- Bagaimana performa model yang dibuat dalam mendeteksi situs *merchant* yang *fraud*?

1.2.2 Batasan Permasalahan

Berikut ini adalah asumsi yang digunakan sebagai batasan penelitian ini:

- Situs *merchant* yang digunakan dalam penelitian adalah situs yang terdaftar dalam proses registrasi *online merchant* PT Nusa Satu Inti Artha.
- Output dari klasifikasi adalah dua *class* yaitu *fraud* dan tidak *fraud* beserta probabilitasnya.

1.3 Tujuan Penelitian

Berikut ini adalah tujuan penelitian yang dilakukan:

- Membuat sebuah *web crawler* yang bisa melakukan ekstraksi data pada sebuah situs *merchant*.
- Mendapatkan fitur-fitur yang bisa digunakan untuk mendeteksi sebuah situs *merchant* yang *fraud*.

⁴Asumsi *best case*, dengan *rate* 15 detik/situs

- Melihat performa model yang dibuat dalam mendeteksi situs *merchant* yang *fraud*.

1.4 Tahapan Penelitian

Tahapan-tahapan yang dilakukan penulis dalam melakukan penelitian adalah sebagai berikut:

1. Studi literatur

Pada tahap ini, dipelajari teori-teori yang terkait dengan penelitian ini untuk mendapatkan konsep dasar yang dibutuhkan dalam mencapai tujuan penelitian.

2. Desain dan implementasi

Pada tahap ini, dilakukan desain dari aplikasi *web crawler* yang akan dibuat dan pemanfaatan *machine learning* yang digunakan untuk melakukan klasifikasi hasil akhir. Hasil prediksi akan dianalisa dan ditarik sebuah kesimpulan dari keseluruhan penelitian ini.

1.5 Sistematika Penulisan

Sistematika penulisan laporan adalah sebagai berikut:

- Bab 1 PENDAHULUAN

Bab ini mencakup latar belakang, cakupan penelitian, dan pendefinisian masalah.

- Bab 2 STUDI LITERATUR

Bab ini mencakup pemaparan terminologi dan teori yang terkait dengan penelitian berdasarkan hasil tinjauan pustaka yang telah dilakukan.

- Bab 3 DESAIN DAN IMPLEMENTASI

Bab ini menjelaskan desain dari aplikasi *web crawler* yang dibuat oleh penulis dan proses implementasinya pada perusahaan.

- Pecah bab 3 menjadi dua bab (desain dan implementasi)

- Bab 4 EKSPERIMEN DAN ANALISIS

Bab ini menjelaskan hasil eksperimen yang dilakukan untuk menguji kebenaran implementasi dari *web crawler* yang dibuat.

- Bab 5 PENUTUP

Bab ini mencakup kesimpulan akhir penelitian dan saran terkait pengembangan selanjutnya.

BAB 2

STUDI LITERATUR

Untuk memulai penelitian, dibutuhkan kerangka berpikir yang sesuai untuk permasalahan yang ingin dipecahkan. Untuk membentuk kerangka berpikir yang sesuai, perlu dikaitkan dengan hasil studi literatur yang telah dilakukan. Oleh karena itu, pada bab ini, akan dijelaskan hasil studi literatur yang telah dilakukan yang telah dikaitkan dengan kerangka kerja untuk penelitian ini.

2.1 *Financial Technology*

Financial technology atau biasa disingkat *fintech* adalah sebuah industri finansial yang mengaplikasikan teknologi untuk kemajuan kegiatan finansialnya [3]. Dalam definisi lain, *fintech* juga dapat disebut sebagai sebuah aplikasi, proses, produk, atau model bisnis baru dalam industri finansial, yang terdiri dari satu atau lebih layanan keuangan dan dilaksanakan melalui proses *end-to-end* melalui internet [4]. Perusahaan yang bergerak pada bidang *fintech* dapat terdiri dari berbagai macam perusahaan, mulai dari *startup*, perusahaan jasa keuangan, ataupun perusahaan teknologi yang menyediakan solusi bagi perusahaan jasa keuangan yang telah ada.

2.1.1 *Payment Financial Technology*

Industri *fintech* sangat luas dengan berbagai macam layanan yang ditawarkan. Sehingga *fintech* dapat dikategorikan menjadi beberapa segmen, salah satunya adalah *payment* [5]. Industri *fintech payment* atau sebagian menyebutnya *digital payment* berfokus pada pelayanan dan penyediaan sistem pembayaran berbasis digital. Contoh dari layanan ini adalah layanan dompet digital, *payment gateway*, dan *transfer services*.

2.1.2 *Know Your Business*

Know Your Business merupakan suatu rangkaian proses yang bertujuan untuk melakukan verifikasi terhadap suatu bisnis [6]. KYB biasa diterapkan pada sebuah perusahaan *fintech* dalam melakukan *screening* terhadap calon mitra kerjanya, baik itu instansi finansial ataupun *merchant*. Hal ini merupakan sebuah protokol standar pada jenis bisnis B2B. KYB dilakukan untuk menghindari terjadinya kerjasama dengan calon mitra yang memiliki tendensi kriminal seperti penipuan identitas (*identity fraud*), bisnis palsu (*fake*

merchant fraud), pencucian uang, pendanaan teroris, perusahaan tempurung, dan lain-lainnya.

2.2 *Web Crawler*

Web crawler atau dalam terminologi singkat disebut *crawler* adalah sebuah bot yang memanfaatkan internet untuk melakukan penelusuran terhadap *World Wide Web* secara sistematis [7]. Pemanfaatan *web crawler* dapat meliputi berbagai kebutuhan, seperti *web indexing*, otomasi *maintenance* situs, dan pengekstrakan data situs (*web scraping*).

2.2.1 *Crawl Frontier*

Crawl Frontier merupakan sebuah bagian dari sistem *crawling* yang mengatur *logic*, urutan, dan ketetapan yang akan diterapkan oleh sebuah *web crawler* dalam mengunjungi sebuah situs. Hal pertama yang akan dilakukan oleh sebuah *crawl frontier* adalah mengumpulkan URL yang akan dituju. Kumpulan dari URL ini dinamakan *seeds* [8]. Seeds selanjutnya dapat dimanfaatkan kedepannya sesuai dengan *logic* yang telah ditentukan oleh *frontier*.

2.2.2 *Web Driver*

Web Driver adalah sebuah *interface* dan protokol yang dapat digunakan untuk mengatur behavior sebuah *web browser*. WebDriver biasa dimanfaatkan untuk melakukan *automatic testing* dengan menuliskan serangkaian instruksi menggunakan sebuah bahasa pemrograman. Bahasa pemrograman yang didukung meliputi Java, .Net, PHP, Python, Perl, dan Ruby. Instruksi yang telah ditulis akan dieksekusi dan akan menjalankan *web browser* yang dipilih sesuai dengan *logic* yang dituliskan. Beberapa *web browser* yang dapat dikontrol melalui WebDriver adalah Google Chrome, Mozilla Firefox, Internet Explorer, Safari, Opera, dan GhostDriver [9].

2.2.3 *HTTP Programming*

HyperText Transfer Protocol (HTTP) adalah sebuah *application-level protocol* yang digunakan untuk sistem informasi terdistribusi dan kolaboratif [10]. Salah satu kegunaan HTTP yang sering digunakan adalah HTTP *request*. Cara kerja sebuah HTTP *request* adalah dengan mengirimkan sebuah *request* ke server yang dituju, dimana data tersebut disimpan, lalu server mengembalikan respon berupa *status code* dan konten (data) yang

diminta. Siklus seperti ini disebut sebagai *request-response cycle*¹.

Untuk dapat melakukan *request*, pengguna harus menuliskan sebuah program komputer yang dapat memungkinkan proses komunikasi antara komputer pengguna dan server dapat terjadi. Hal ini disebut sebagai *HTTP Programming*. Saat ini ada banyak *library* yang dapat dimanfaatkan untuk melakukan hal ini, salah satunya *library requests*² yang dibangun di atas bahasa pemrograman Python.

2.3 Web Scraping

Web scraping merupakan sebuah terminologi yang digunakan untuk menggambarkan proses esktraksi data dari sebuah situs [11]. *Web scraping* dapat digunakan untuk beberapa kebutuhan seperti *information extraction*, *content analysis*, *content ranking*, *price engine*, dan banyak hal lainnya. Proses *scraping* dapat dilakukan secara manual oleh pengguna atau disematkan secara otomatis di dalam sebuah *web crawler*.

Web scraping yang dilakukan secara manual oleh pengguna biasanya dilakukan untuk keperluan tidak berulang (*single use*) dan memang spesifik terhadap suatu situs. Sedangkan, untuk melakukan ekstraksi data secara berulang dan dalam skala yang besar, maka proses *scraping* harus dijalankan menggunakan bantuan *bot* atau biasa disebut *web crawler*.

Web crawler berfungsi sebagai sebuah *bot* yang akan membentuk sebuah jaring-jaring, jalan, atau peta, untuk mengunjungi berbagai macam-macam halaman yang ada di situs tersebut [12]. Setiap *crawler* mengunjungi sebuah halaman, maka proses *scraping* akan dijalankan. Informasi ataupun data yang ada dalam halaman tersebut akan diekstrak dan disimpan atau langsung diproses sesuai dengan kebutuhan pengguna. Halaman web biasanya dibangun di atas sebuah *text-based mark-up language* seperti HTML dan XHTML. Karena hal ini, tidak semua informasi yang diekstrak biasanya dapat langsung dimengerti dan diinterpretasikan oleh pengguna. Sehingga perlu ada proses lanjutan lainnya yang dapat menyempurnakan proses *scraping* yaitu dengan melakukan *HTML parsing*.

2.3.1 HTML Parsing

Objek yang diterima dalam sebuah *HTTP request* merupakan sebuah entitas *HTTP response message*. Di dalam *HTTP response message*, terdapat beberapa komponen yang termuat, diantaranya *status code*, *header*, dan *body* [10]. *Status code* merepresentasikan status balasan dari *request* yang dikirimkan oleh pengguna. *Status code* yang lazim ditemui dalam melakukan proses *crawling* dapat dilihat pada Tabel 2.1

¹ Menurut situs pembelajaran daring *w3schools*, https://www.w3schools.com/whatis/whatis_http.asp

² <https://pypi.org/project/requests/>

Status Code	Deskripsi
0	Blocked by Robots.txt/DNS Lookup Failed/- Connection Timeout/Connection Refused/No Response
200	OK
301	Moved Permanently
302	Moved Temporarily
400	Bad Request
403	Forbidden
404	Page Not Found
410	Removed
429	Too Many Requests
500	Internal Server Error

Tabel 2.1: Contoh *response code* yang dikembalikan *HTTP response message*

Terkait konten ataupun informasi utama dari halaman yang dituju, akan dikembalikan di dalam *body*. Untuk dapat melakukan ekstraksi informasi yang ada pada *body* dari *HTML response* maka perlu dilakukan teknik *HTML parsing*. Salah satu *library* yang dapat dimanfaatkan dalam melakukan *parsing* adalah BeautifulSoup³. BeautifulSoup memungkinkan untuk mengubah dokumen HTML yang kompleks menjadi objek Python yang berbentuk *tree* dan dapat dengan mudah dinavigasikan oleh pengguna dengan cara mengakses tag-tag yang ada pada *tree* HTML tersebut.

2.3.2 Regular Expression

Regular expression atau biasa disebut *regex* merupakan sebuah konstruksi teks yang disusun untuk mencocokkan sebuah pola tertentu. Regex bisa digunakan untuk melakukan operasi *matching* (pencocokan), *locating* (pelokasian), dan *manipulating* (manipulasi), pada sebuah *string*. Pola (*pattern*) yang dapat dibuat menggunakan regex sangat banyak, beberapa contoh yang paling sering digunakan tertulis dalam Tabel ?? di bawah ini.

³<https://www.crummy.com/software/BeautifulSoup/bs4/doc/>

Regex	Contoh Syntax	Contoh String	Penjelasan
<code>\d</code>	<code>urutan_\d</code>	<code>urutan_5</code>	Mencocokkan satu digit angka 0-9
<code>[abc]</code>	<code>karakter_[abc]</code>	<code>karakter_a</code>	Mencocokkan karakter "a" atau "b" atau "c" sebanyak 1 kali
<code>[abc]+</code>	<code>karakter_[abc]+</code>	<code>karakter_aaa</code>	Mencocokkan karakter "a" atau "b" atau "c" sebanyak 1 kali sampai tidak terbatas
<code>[^abc]</code>	<code>karakter_[^abc]</code>	<code>karakter_s</code>	Mencocokkan karakter selain "a" atau "b" atau "c" sebanyak 1 kali
<code>\s</code>	<code>satu\sdua</code>	<code>satu dua</code>	Mencocokkan karakter <i>whitespace</i> (<i>spasi, tab, newline</i>)
<code>[a-zA-z]</code>	<code>karakter_[a-zA-z]</code>	<code>karakter_Y</code>	Mencocokkan karakter dari <i>range</i> "a" sampai "z" atau "A" sampai "Z" sebanyak 1 kali
<code>(word1 word2)</code>	<code>ini (word1 word2)</code>	<code>ini word1</code>	Mencocokkan kata "word1" atau "word2" sebanyak 1 kali

Tabel 2.2: Contoh penggunaan syntax *regex*

Pada Tabel 2.2, yang ditampilkan hanya merupakan beberapa *regex* yang biasa dan sering digunakan dan termuat dalam *common regex usage* di beberapa situs. Penggunaan *syntax regex* dalam bahasa pemrograman Python sendiri dapat dilakukan menggunakan *built-in library* yang dimiliki, yang dinamakan *re*⁴. Ada beberapa metode yang dapat dilakukan menggunakan *re*, diantaranya sebagai berikut.

- **re.match():** Menentukan apakah *regex* yang dibuat dapat menemukan *pattern* yang cocok pada awal string yang diberikan.
- **re.search():** Melakukan *scanning* di keseluruhan string, mencari apakah ada *pattern* yang sesuai dengan *regex* yang ditentukan.

⁴<https://docs.python.org/3/howto/regex.html>

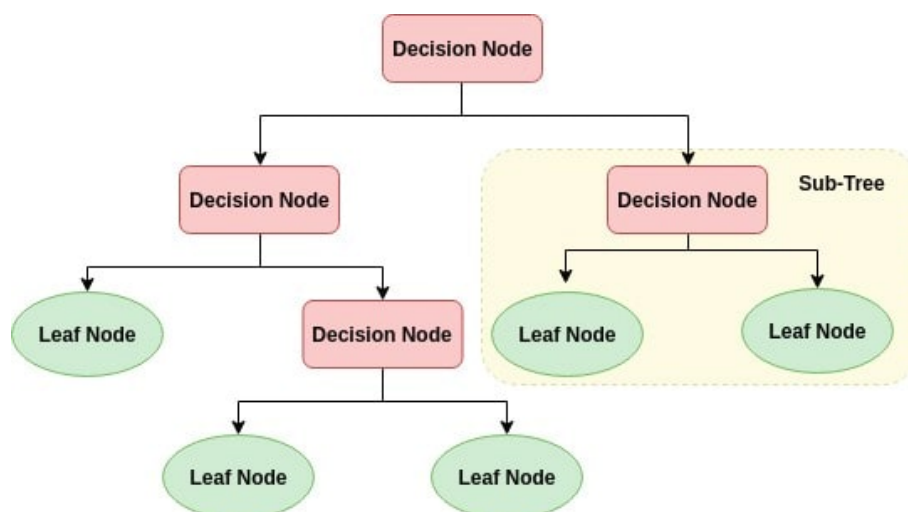
- **re.findall()**: Mencari semua substring dalam sebuah string dimana *pattern* regex dapat ditemukan. Fungsi ini mengembalikan semua substring dalam bentuk list.

2.4 Classifier

Classifier merupakan sebuah algoritma yang memanfaatkan atau mempelajari input (*training data*) untuk melakukan prediksi terhadap label kategori (*class*) dari sebuah *data point* yang diberikan (*testing data*) [13]. *Classifier* banyak dimanfaatkan untuk banyak kasus di dunia nyata, seperti klasifikasi *spam email*, klasifikasi transaksi *fraud*, dan banyak hal lainnya yang merupakan implementasi dari *pattern recognition*⁵.

2.4.1 Decision Tree Classifier

Decision Tree Classifier merupakan salah satu model prediktif yang biasa digunakan dalam *data mining*, *machine learning*, dan statistik. *Decision tree* menerapkan algoritma *non-parametric unsupervised learning* dan dapat digunakan untuk persoalan regresi maupun klasifikasi. *Classifier* jenis ini akan membangun sebuah *tree* yang akan melakukan partisi secara rekursif, menciptakan sebuah *tree* baru (*sub-tree*), sampai homogenitas dari setiap node terminimalisir. Strukturnya yang mudah dimengerti dan menyerupai *flowchart* memudahkan penggunaannya untuk mengerti dan menginterpretasikan hasil klasifikasinya. Struktur dari *decision tree* dapat dilihat lebih jelas pada Gambar 2.1 dibawah ini.



Gambar 2.1: Struktur *tree* pada Decision Tree Classifier, diambil dari DataCamp Community

⁵Lebih lengkapnya tentang *pattern recognition* dapat diakses pada tautan berikut <https://deepai.org/machine-learning-glossary-and-terms/pattern-recognition>

Pada sebuah *decision tree*, *node* yang berada paling atas disebut dengan *root node*. *Root node* merupakan partisi pertama dari keseluruhan data yang ada. Berikutnya, untuk setiap *node* yang berada di sisi dalam (*internal node*), merepresentasikan atribut (*feature*) yang digunakan pada model. Lalu, setiap percabangan (*branch*) merepresentasikan *decision rule* yang dibuat dan setiap *node* terluar (*leaf node*) merepresentasikan hasil akhir (*outcome*) dari partisi data.

$$\text{Gini impurity} = \sum_{i=1}^k p_i(1 - p_i) \text{cross-entropy} = - \sum_{i=1}^k p_i \log(p_i)$$

Gambar 2.2: Cara menghitung *gini impurity* pada Decision Tree Classifier

Dalam kasus *outcome categorical*, salah satu cara yang paling umum dalam melakukan *splitting* adalah dengan menghitung *gini impurity*. Pada persamaan di Gambar 2.2, k mewakili jumlah *class* (kategori) yang ada pada *dataset*. Sedangkan p_i merepresentasikan proporsi *case* (data) yang merupakan bagian dari *class* i . Cara untuk melakukan *splitting* adalah dengan menghitung *gini impurity* pada *initial node* dan pada kedua *child node* yang dihasilkan. Jika hasil *weighted sum gini impurity* dari kedua *child node* lebih kecil daripada *gini impurity* yang ada pada *initial node*, maka *split* dianggap menghasilkan *improvement* untuk lebih mempartisi data dan *split* akan dilakukan [14].

Untuk menggunakan Decision Tree Classifier dalam Python, *scikit-learn*⁶ telah menyediakan sebuah *library* yang dapat digunakan secara komprehensif untuk melakukan *training*, *tuning*, dan *predicting*. Ada beberapa parameter penting yang perlu dipahami dalam menggunakan *library* ini, diantaranya adalah:

- **criterion** (“gini”, “entropy”, default=“gini”)

Menentukan standar hitung yang digunakan untuk melakukan *split*. Untuk *split* berdasarkan *gini impurity*, dapat menggunakan “gini”. Sedangkan untuk menghitung berdasarkan *information gain*, dapat menggunakan “entropy”.

- **max_depth** (“int”, default=None)

Menentukan kedalaman maksimum dari *tree* yang dibuat. Jika dibiarkan None (*default*), maka *tree* akan terus berekspansi sampai semua *leaf* sudah murni atau jumlah *dataset* pada *leaf* sudah lebih sedikit dari jumlah yang ditentukan pada *min_samples_split*

- **min_samples_split** (int, float, default=2)

Jumlah minimum data pada sebuah *leaf* yang dibutuhkan untuk sebuah *leaf* agar

⁶<https://scikit-learn.org/>

dapat melakukan *split*. Jika yang dimasukkan adalah *integer*, maka angka tersebut akan menjadi angka minimum data yang dimaksud. Sedangkan apabila nilai *float* yang dimasukkan, maka akan ditranslasikan menjadi presentase minimum sampel.

- **min_samples_leaf** (int, float, default=1)

Jumlah minimum data yang harus ada pada sebuah *leaf*. *Split* hanya dapat dilakukan apabila pada jumlah data pada *child node* yang dihasilkan berisi minimum sejumlah/di atas jumlah yang ditetapkan.

2.4.2 Random Forest Classifier

Random Forest Classifier adalah sebuah *tree-based* model yang dikembangkan dari algoritma dari Decision Tree Classifier. Faktor yang membedakan diantara keduanya adalah jumlah *tree* yang dibangun saat proses *fitting*. Jika pada Decision Tree Classifier *tree* yang dibangun hanya berjumlah satu, pada Random Forest Classifier, seperti namanya, *forest*, algoritma ini akan membangun lebih dari satu *tree*. Pada metode ini, proses *training* dilakukan menggunakan metode *bagging*. Ide umum dari metode *bagging* adalah dengan menggabungkan berbagai macam *learner* untuk menciptakan hasil keseluruhan yang lebih baik. Dengan adanya diversitas *tree* dan fitur-fitur yang dipilihnya, secara umum, performa yang dihasilkan model akan menjadi lebih baik.

Model ini disebut *random* dikarenakan ada sebuah perilaku acak yang terjadi dalam membangun kumpulan *treenya*. Alih-alih memilih fitur terpenting dari seluruh fitur yang ada dalam melakukan *splitting* pada sebuah *node*, *random forest* akan menentukan terlebih dahulu subset fitur yang akan digunakan secara *random* baru melakukan pemilihan fitur dengan kualitas split terbaik dari subset tersebut [15].

Sebuah *tree* akan dibangun sesuai parameter yang ditentukan layaknya yang terjadi pada Decision Tree Classifier. Jadi pada dasarnya, parameter yang dapat digunakan untuk kedua model ini dapat terbilang sama. Namun, model ini tidak akan berhenti pada satu *tree*, melainkan akan kembali membangun *tree* berikutnya dengan skema pemilihan fitur acak, sampai jumlah *tree* yang diinginkan tercapai. Saat semua *tree* sudah dibangun, maka nantinya proses *predicting* akan melalui proses *voting*. Proses *voting* adalah proses pemilihan suara terbanyak dari keseluruhan hasil prediksi *tree* yang ada. Misalkan, pada sebuah proses *fitting* Random Forest Classifier menghasilkan 10 *tree* dan masing-masing *tree* telah memiliki hasil prediksinya tersendiri terhadap sebuah *test data*. Jika 8 dari 10 *tree* memprediksi bahwa *test data* tersebut milik sebuah *class* 1 dan 2 sisanya memprediksi *class* 0, maka model akan mengembalikan hasil prediksi dengan *vote* terbanyak, yaitu *test data* milik *class* 1.

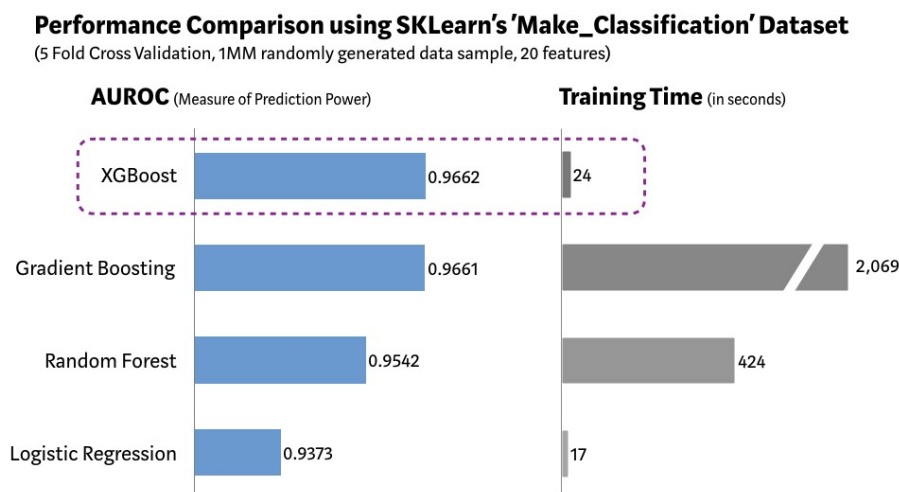
Pemilihan dataset dalam proses *fitting* Random Forest Classifier juga dapat dilakukan

dengan teknik *bootstrap*. Bootstrap merupakan sebuah metode *sampling* statistika, yaitu dengan cara mengambil sebuah proporsi sampel pada dataset secara iteratif dan acak dengan menerapkan *replacement* [16]. Salah satu keunggulan dari metode ini adalah bisa dihasilkannya lebih banyak data dalam melakukan estimasi tanpa harus mengumpulkan lebih banyak *training data*.

2.4.3 Extreme Gradient Boost Classifier

Extreme Gradient Boost atau biasa dikenal dengan XGBoost adalah sebuah algoritma yang baru-baru ini dikembangkan dan mendominasi pada berbagai bidang *applied machine learning* dan di beberapa kompetisi-kompetisi daring [17]. XGBoost sendiri merupakan sebuah implementasi dari algoritma Gradient Boosted Decision Tree yang telah dimodifikasi sedemikian rupa sehingga memiliki performa yang lebih baik dalam perihal kecepatan dan performa. Model ini pertama kali dikenalkan oleh seorang *researcher* yang memiliki spesialisasi dalam bidang *machine learning*, Tianqi Chen⁷.

XGboost merupakan sebuah metode *ensemble learning* [18]. Terkadang, bergantung pada satu model *machine learning* tidak cukup. *Ensemble learning* menawarkan sebuah solusi yang sistematis untuk menggabungkan kemampuan prediksi dari berbagai macam *learner*. Hasil dari penggabungan ini adalah sebuah model tunggal yang memberikan keluaran yang diaggregasikan dari beberapa model (*learner*). Dengan melakukan aggregasi dan penggabungan ini, maka hasil prediksi yang dibuat akan menjadi lebih baik dan memiliki variansi yang rendah.



Gambar 2.3: Training time dan performa XGBoost dibandingkan model lain

XGBoost memiliki dukungan dari komunitas Data Scientist yang kuat. Saat ini *open*

⁷<https://tqchen.com/>

source proyek XGBoost telah memiliki lebih dari 350 kontributor dan 3600 *commit* pada *repository* GitHub⁸. Konsekuensi dari pengembangan dan kontribusi yang luas dari algoritma ini adalah:

- Dapat digunakan dalam berbagai macam permasalahan: Dapat digunakan pada permasalahan regresi, klasifikasi, dan *ranking*.
- Portabilitas: Dapat dijalankan dengan lancar di berbagai macam OS (Windows, Linux, Mac OS) dan di berbagai macam *cloud server* seperti AWS, Azure, dan ekosistem lain seperti Spark, Flink.
- Bahasa pemrograman: Dapat dijalankan di berbagai macam bahasa pemrograman diantaranya C++, Python, R, Java, Scala, dan Julia
- Kecepatan dan Ketepatan: Dapat dilihat pada Gambar 2.3 bahwa menggunakan *dataset* yang sama, XGBoost dapat menampilkan hasil yang lebih baik dengan nilai 0.96 hanya dengan waktu training selama 24 detik

2.4.4 Bernoulli Naive Bayes Classifier

Metode-metode Naive Bayes adalah sekumpulan metode *supervised learning* yang didasari pada penerapan teorema Bayes. Teorema Bayes menerapkan asumsi *naive*, yaitu adanya asumsi independensi yang kondisional diantara satu fitur dengan fitur lainnya diketahui nilai dari *class variable*. Teorema Bayes dapat dinyatakan sebagai berikut:

$$P(y | x_1, \dots, x_n) = \frac{P(y)P(x_1, \dots, x_n | y)}{P(x_1, \dots, x_n)}$$

Pada teorema Bayes, y merupakan *class label* yang diketahui dan x_1 sampai dengan x_n merupakan *feature vector* yang dependen. Teorema ini berlaku pada keseluruhan Naive Bayes Classifiers yang ada, umumnya yang membedakan antara satu jenis Naive Bayes Classifier dan lainnya adalah asumsi dari masing-masing jenis terhadap distribusi dari $P(x_i | y)$ [19].

Walaupun Naive Bayes dikatakan sebagai algoritma yang sangat sederhana, namun algoritma ini telah banyak terbukti cukup baik dalam menghadapi berbagai macam permasalahan nyata di dunia, seperti pada kasus *document classification* dan *spam filtering* [20]. Jumlah *training data* yang dibutuhkan algoritma ini tidak perlu banyak untuk dapat melakukan prediksi yang baik. Selain itu, algoritma ini juga cepat dalam melakukan proses *training* dan *predicting* dibandingkan algoritma rumit lainnya.

⁸<https://github.com/dmlc/xgboost/>

Spesifik terhadap salah satu jenis dari Naive Bayes, yaitu Bernoulli Naive Bayes algoritma ini merupakan algoritma yang menerapkan cara Naive Bayes melakukan *training* dan *classification* terhadap data-data yang terdistribusikan secara Bernoulli multivariat. Maksudnya dari Bernoulli multivariat disini adalah saat ada berbagai macam fitur-fitur yang dijadikan input dan setiap fitur yang dimiliki tersebut bersifat *binary-valued* (Bernoulli, *boolean*). Jika fitur-fitur tersebut tidak bersifat biner, maka binerisasi akan dilakukan terhadap nilai-nilai dari fitur-fitur yang tidak bersifat biner tersebut.

$$P(x_i | y) = P(i | y)x_i + (1 - P(i | y))(1 - x_i)$$

Decision rule untuk Bernoulli Naive Bayes ditampilkan pada persamaan di atas. Dimana pada Bernoulli Naive Bayes, kemunculan dari fitur i tidak diacuhkan seperti pada Multinomial Naive Bayes. Pada algoritma Bernoulli, ketidakmunculan fitur i diberikan penalti namun tetap dimasukan kedalam perhitungan.

Pengimplementasian algoritma Bernoulli Naive Bayes dapat menggunakan model BernoulliNB⁹ pada *scikit-learn*. Ada beberapa parameter penting yang perlu dipahami dalam menggunakan *library* ini, diantaranya adalah:

- **alpha** (float, default="1.0")
Additive (Laplace/Lidstone) *smoothing* parameter pada *decision rule* (0 untuk tidak ada *smoothing* sama sekali)
- **binarize** (float, None, default=0.0)
Threshold yang dapat ditentukan untuk melakukan binerisasi dari *sample features* yang bukan merupakan *binary-valued*. Jika None, maka diasumsikan fitur yang diberikan sudah merupakan fitur biner.
- **fit_prior** (boolean, default=True)
Menentukan apakah model perlu mempelajari probabilitas *class prior* atau tidak. Jika False, maka *prior* uniform untuk dua *class* akan digunakan.

2.5 Validasi model

Untuk mengukur kualitas dari sebuah model yang telah dibuat, maka sebuah validasi harus dilakukan. Validasi model merupakan sebuah kunci untuk dapat mengembangkan model yang telah dibuat secara iteratif. Dalam melakukan validasi, ada beberapa teknik yang dapat diadopsi, diantaranya adalah dengan melakukan *nested cross validation* dan melihat metrik evaluasi suatu model.

⁹https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.BernoulliNB.html

2.5.1 *Nested Cross Validation*

Nested Cross validation merupakan salah satu teknik validasi model yang dikembangkan dari teknik *cross validation* pada umumnya. *Cross validation* sendiri merupakan sebuah teknik validasi yang digunakan untuk mengestimasi kualitas model secara tidak bias. Hal ini dikarenakan *cross validation* dapat menghindari *overfitting* dan *selection bias* pada model [21]. Metode ini bekerja dengan membagi dataset menjadi k bagian. Dari k bagian yang ada, akan diambil 1 bagian sebagai *test set* dan sisanya ($k-1$ bagian) lainnya akan menjadi *train set*. Pembagian ini akan dilakukan *looping* sebanyak k kali, sehingga semua bagian data berkesempatan menjadi *test set*.

Sedangkan pada *nested cross validation*, proses k -fold *cross validation* yang biasa dilakukan hanya menjadi salah satu bagian dari metode ini. K -fold pada metode ini akan dilakukan pada dua *loop*, yaitu pada *outer loop* dan *inner loop*. Tujuan dari pembagian ini adalah agar *test data* tetap terjaga *pure* pada *fold* di *outer loop* karena CV dilakukan pada *inner loop* di *training set* *outer loop* tersebut [22]. Algoritma dari *nested cross validation* secara lengkap dapat dijabarkan sebagai berikut:

1. Bagi keseluruhan dataset menjadi K *folds* secara random
2. Untuk setiap fold $k = 1, 2, \dots, K$: merupakan *outer loop* digunakan untuk evaluasi model dengan hyperparameter yang terpilih
 - (a) Jadikan *fold k* sebagai *test*
 - (b) Jadikan seluruh data selain yang ada di *fold k* sebagai *trainval*
 - (c) Secara *random*, lakukan *split* pada *trainval* menjadi L *folds*
 - (d) Untuk setiap fold $l = 1, 2, \dots, L$: merupakan *inner loop* digunakan untuk *hyperparameter tuning*
 - i. Jadikan *fold l* sebagai *val*
 - ii. Jadikan seluruh data selain yang ada pada *test* dan *val* sebagai *train*
 - iii. Lakukan *training* menggunakan setiap kombinasi *hyperparameter* pada *train* lalu evaluasi pada *val*. Catat *performance metric* dari *hyperparameter* tersebut
 - (e) Untuk setiap kombinasi *hyperparameter*, lakukan kalkulasi rata-rata *metric score* pada L *folds* yang ada lalu pilih kombinasi yang terbaik
 - (f) Train sebuah model menggunakan *hyperparameter* terbaik pada *trainval*. Evaluasi performanya pada *test* lalu simpan *score* untuk *fold k*
3. Lakukan kalkulasi rata-rata skor yang ada pada seluruh K *folds*

2.5.2 Metrik Evaluasi

Untuk mengukur kualitas prediksi dari model, sebuah metrik harus ditentukan sebagai acuan. Ada berbagai macam metrik yang dapat dijadikan sebagai acuan penilaian untuk sebuah model *classifier*, diantaranya adalah Confusion Matrix, Akurasi, AUC, Precision, Recall, dan F1 Score.

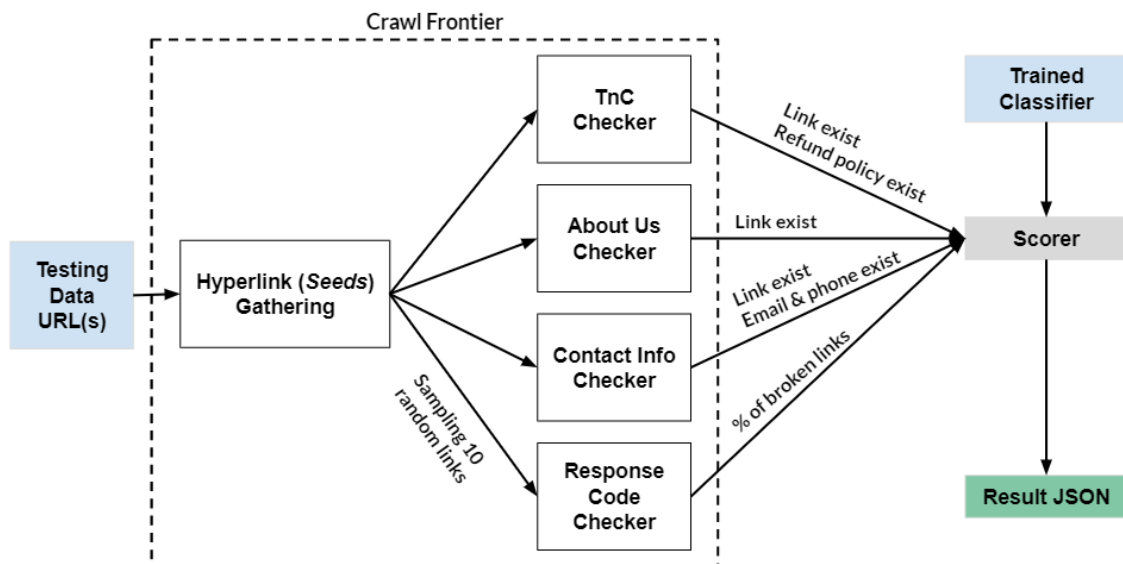
- **Confusion matrix:** Sebuah matriks untuk menampilkan visualisasi dari performa model dalam melakukan klasifikasi. Setiap kolom pada *confusion matrix* merepresentasikan jumlah data yang terprediksi pada *class* yang bersangkutan dan setiap baris merepresentasikan jumlah data aktual pada *class* yang bersangkutan atau dapat pula terjadi sebaliknya [23]. Adapun 4 klasifikasi elemen yang ada pada *confusion matrix* dapat didefinisikan sebagai berikut:
 - True Positive (TP): Jumlah data positif yang diprediksi positif.
 - False Positive (FP): Jumlah data negatif yang diprediksi positif.
 - True Negative (TN): Jumlah data negatif yang diprediksi negatif.
 - False Negative (FN): Jumlah data positif yang diprediksi negatif.
- **AUC:** Probabilitas sebuah model bisa mengurutkan probabilitas data positif lebih tinggi dibandingkan data negatif. Dengan $AUC = 1$, jika dibuat sebuah garis bilangan horizontal dari 0 - 1, maka semua probabilitas prediksi data positif akan berada di sebelah kanan data negatif.
- **Akurasi:** Rasio prediksi benar (*class* positif dan *class* negatif) dibanding keseluruhan data. Dapat dihitung dengan rumus $(TP+TN) / (TP+FP+FN+TN)$.
- **Precision:** Berapa proporsi benar dari prediksi *class* positif yang model lakukan. Precision dapat dihitung dengan rumus $TP / (TP+FP)$.
- **Recall:** Berapa proporsi benar dari keseluruhan data positif yang ada yang dapat diprediksi oleh model. Recall dapat dihitung dengan rumus $TP / (TP+FN)$.
- **F1 Score:** Perbandingan rata-rata Precision dan Recall yang dibobotkan. Dapat dihitung dengan rumus $2 * (Recall * Precision) / (Recall + Precision)$.

BAB 3

DESAIN DAN IMPLEMENTASI

Pada bab ini akan dijelaskan mengenai desain dan implementasi *web crawler* untuk melakukan pendeteksian *fraud*. Desain dan implementasi mencakup arsitektur *web crawler*, desain proses pelatihan model, implementasi modul *helper*, implementasi modul ekstraksi fitur, dan *model scorer* pada *crawler*.

3.1 Arsitektur Web Crawler



Gambar 3.1: Arsitektur Web Crawler

Web crawler dibangun di atas bahasa pemrograman Python menggunakan sebuah *framework microservice* yang ringkas dan mudah diimplementasikan bernama Flask¹. *Web crawler* dapat diakses melalui sebuah *endpoint*² yang menerima sebuah *request* GET berupa URL atau sekumpulan URL situs *merchant* yang telah dikompilasi ke dalam satu *file* CSV (*comma-seperated values*). URL tersebut akan melalui serangkaian *logic* (*process*) yang diatur di dalam sebuah *crawl frontier*. Proses tersebut meliputi pengumpulan *hyperlink* yang terdapat pada halaman muka situs *merchant* (*hyperlink gathering*) dan pengecekan kelengkapan atribut situs, secara bergantian, dengan mengunjungi tautan-tautan pada situs tersebut yang telah terkumpul dari proses *hyperlink gathering*.

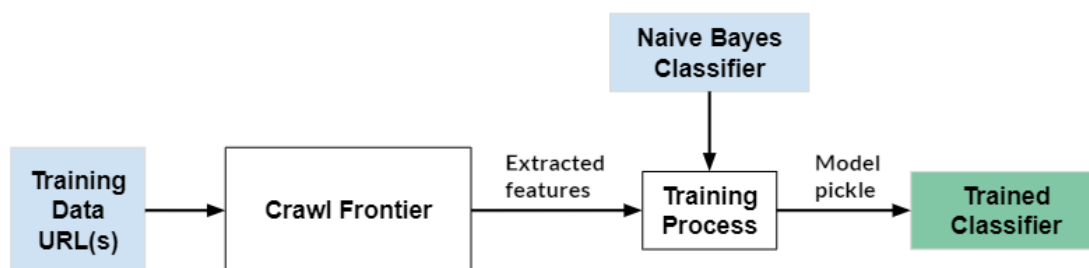
¹<https://flask.palletsprojects.com/en/1.1.x/>

²Demo untuk *endpoint* ini dapat diakses pada <http://demodemotest.com>

Setelah menjalani proses pengecekan yang terdapat pada *crawl frontier*, keluaran dari setiap proses tersebut akan dikumpulkan menjadi serangkaian fitur. Fitur kemudian dijadikan sebagai *input* bagi *classifier* pada proses *scoring*. Hasil dari *scoring* dan fitur-fitur yang telah diekstrak oleh *crawler* akan dijadikan satu di dalam sebuah JSON dan dikembalikan sebagai *response* dari penembakan *endpoint*.

3.2 Proses Pelatihan Model

Pada Gambar 3.1 terdapat sebuah komponen yang dinamakan *trained classifier*. Proses yang terjadi untuk menghasilkan *trained classifier* tersebut digambarkan pada Gambar 3.2 berikut.



Gambar 3.2: Proses Training

Proses *training* dimulai dengan mengakses sebuah *endpoint*³ yang menerima *input* berupa URL atau sekumpulan URL. URL tersebut merupakan URL yang ingin dijadikan sebagai *training data*. URL kemudian akan diproses oleh *crawl frontier* dan menghasilkan fitur-fitur yang telah diekstrak. Fitur-fitur tersebut kemudian akan menjadi *input* bagi *classifier* yang dipilih, yaitu Naive Bayes Classifier. Naive Bayes Classifier yang digunakan pada proses ini adalah Bernoulli Naive Bayes. Hal ini didasarkan pada kemampuannya yang baik dalam mengolah input fitur yang bersifat *binary* untuk melakukan prediksi.

Setelah proses *fitting* model dilakukan dan hasil evaluasi performa model dirasa sudah tepat, maka model yang sudah melalui proses *training* tersebut akan disimpan ke dalam bentuk *pickle*. Pickle memungkinkan untuk menyimpan sebuah *python object* menjadi sebuah *binary file* [24]. *Binary file* tersebut nantinya dapat kembali diakses dan digunakan kembali sebagai *python object* walaupun dijalankan dalam *session* dan waktu yang berbeda.

³Demo untuk *endpoint* ini dapat diakses pada <http://demodemotest.com>

3.3 Modul *Helper*

Dalam menjalankan *web crawler*, terdapat beberapa fungsi yang memang digunakan secara berulang-ulang di keseluruhan implementasi. Hal ini dikarenakan fungsi ini merupakan fungsi-fungsi yang bersifat umum dan memang didesain untuk memudahkan modul-modul spesifik untuk dapat berjalan. Fungsi-fungsi ini dinamakan fungsi *helper* dan diantaranya adalah *reset browser*, *url format handler*, *hyperlink gatherer* dan *paragraph extractor*.

3.3.1 Inisiasi & Atur Ulang Browser

Dalam mengunjungi halaman-halaman yang telah ditentukan, ada dua cara umum yang dapat dilakukan oleh *crawler*, yaitu dengan mengirimkan *request* dan melakukan *parsing response* atau mengunjungi halaman dengan *WebDriver* dan melakukan *parsing* terhadap konten di halaman tersebut. Untuk cara kedua, yaitu menggunakan *WebDriver*, maka untuk setiap sesi *crawling* perlu dilakukan inisiasi dari *WebDriver* itu sendiri. *Selenium WebDriver* digunakan untuk bernavigas dan *browser* yang digunakan pada *web crawler* ini adalah Google Chrome (*ChromeDriver*⁴).

```

1  ## Initiate ChromeDriver (WebDriver) for the 1st time
2  chrome_options = webdriver.ChromeOptions()
3  chrome_options.add_argument('--ignore-certificate-errors')
4  driver = webdriver.Chrome(chrome_options=chrome_options)
5  hyperlinks_dynamic = False
6  dynamic_links = []
7  dynamic_texts = []
8
9  ## Reset ChromeDriver function
10 def reset_browser():
11     ''' Browser will be reset if an exception occurs when running di WebDriver
12         and will be reset after one crawling job finished '''
13
14     global driver
15     global hyperlinks_dynamic
16     global dynamic_links
17     global dynamic_texts
18
19     ## Quit driver
20     driver.close()
21     driver.quit()
22
23     chrome_options = webdriver.ChromeOptions()
24     chrome_options.add_argument('--ignore-certificate-errors')
25
26     driver = webdriver.Chrome(chrome_options=chrome_options)
27     hyperlinks_dynamic = False

```

⁴<https://chromedriver.chromium.org/getting-started>

```

28     dynamic_links = []
29     dynamic_texts = []

```

Kode 3.1: Fungsi inisiasi dan atur ulang *browser*

Pada 7 baris pertama di luar fungsi *reset_browser* yang termuat pada Kode 3.1, ditampilkan bagaimana *ChromeDriver* diinisiasi untuk pertama kali pada awal sesi. *ChromeDriver* dapat menerima argumen-argumen yang dapat dimasukkan ke dalam parameter *chrome_options*. Parameter yang penulis gunakan adalah `"--ignore-certificate-errors"`. Parameter ini berfungsi untuk mengabaikan *SSL Certificate Error* yang kerap terjadi pada situs-situs yang tidak memiliki kelengkapan sertifikasi yang baik. Jika tidak diabaikan, maka situs-situs tersebut tidak dapat diakses dan diekstrak kontennya. Terdapat 3 variabel penting yang diinisiasi pada bagian ini juga, yaitu *hyperlinks_dynamic* (*boolean*), *dynamic_links* (*list*), dan *dynamic_texts* (*list*). Kegunaan *hyperlinks_dynamic* adalah sebagai *flag* apakah *ChromeDriver* sudah pernah dijalankan untuk mengekstrak konten ataukah belum. Jika sudah, maka *flag* akan berubah menjadi `True` sehingga apabila *ChromeDriver* dijalankan untuk kedua kalinya pada sesi yang sama, proses *gathering hyperlinks* dan *paragraph extractor* yang dilakukan dengan cara dinamis tidak perlu dijalankan kembali. Hal ini dilakukan dengan tujuan efisiensi waktu dan *bandwidth* karena konten pada *running* pertama sudah berhasil diekstrak dan disimpan sehingga *crawler* hanya perlu menggunakan ulang konten tersebut.

Fungsi *reset_browser* sendiri memiliki tujuan untuk melakukan atur ulang (*resetting*) terhadap *browser* yang sedang berjalan. Yang terjadi adalah *browser* akan menutup dirinya sendiri dan melakukan inisiasi ulang seperti yang terjadi pada awal sesi. Hal seperti ini diperlukan saat terjadi *TimeoutExceptionError* pada *ChromeDriver* ataupun saat sesi *crawling* untuk satu situs telah berakhir. Mengatur dan menginisiasi ulang *ChromeDriver* pada dua kondisi tersebut memastikan *browser* sudah *clean* dan bisa digunakan kembali untuk mengeksekusi situs ataupun perintah yang berjalan berikutnya.

3.3.2 URL Format Handler

Fungsi ini merupakan sebuah fungsi singkat yang melakukan *handling* terhadap *input* URL dari pengguna yang tidak sesuai dengan *library-library* yang digunakan. *Library request* ataupun *ChromeDriver* yang digunakan pada *web crawler* memiliki format tertentu yang dapat diterima sebagai masukan. Salah satu contohnya adalah URL yang diterima hanya URL yang memiliki skema, yaitu yang diawali dengan `http://` ataupun `https://`. Sedangkan, pengguna terkadang malas ataupun lupa untuk memasukkan skema tersebut ke dalam URL yang diinput. Hal ini akan menyebabkan *error* dan *library* tidak akan dapat berjalan. Adapun variasi kasus masukan yang ditangani pada fungsi ini dapat dilihat

lebih lengkap pada Tabel 3.1 berikut.

No	Masukan	Keluaran	Deskripsi
1	" https://facebook.com "	"https://facebook.com"	Menghilangkan <i>whitespace</i> di belakang dan depan URL
2	"http://facebook.com"	"https://facebook.com"	Mengubah seluruh URL http:// menjadi https://
3	"https://facebook.com/"	"https://facebook.com"	Menghilangkan <i>backslash</i> di belakang URL

Tabel 3.1: Konversi masukan URL menjadi keluaran yang sesuai

3.3.3 *Hyperlink Gatherer*

Pada proses *hyperlink gathering*, yang terjadi adalah *crawler* akan melakukan ekstraksi pada seluruh *hyperlink* yang ada di halaman antarmuka situs tersebut. Pada implementasinya, terdapat dua cara untuk melakukan ekstraksi *hyperlink*, bergantung pada jenis halaman utama yang dimiliki oleh situs yang dituju. Untuk situs dengan halaman statis, maka *crawler* akan mengirimkan *HTTP request* dan melakukan *parsing* terhadap *response body* yang didapatkan untuk mengekstraksi *tag link* (<a href>) yang ada pada *response* tersebut. Sedangkan untuk melakukan ekstraksi pada halaman muka yang dinamis⁵, *crawler* akan menggunakan *WebDriver* untuk mengekstraksi *tag link* yang terdapat pada halaman dinamis tersebut.

3.3.3.1 *Static Hyperlink Gatherer*

Dalam implementasi kali ini, penulis memanfaatkan sebuah *library* yang memungkinkan *http programming* menjadi lebih mudah di Python, yaitu *requests-html*⁶. Pada dasarnya *requests-html* merupakan pengembangan dari *library requests* yang dikeluarkan oleh Python, hanya saja terdapat beberapa fungsi olahan lanjutan yang diimplementasikan oleh *library ini* sehingga *code* dapat ditulis dan dijalankan dengan lebih efisien.

```
1 def get_hyperlinks(url):
```

⁵Halaman dinamis didefinisikan sebagai halaman yang memerlukan *javascript* untuk melakukan *rendering* terhadap keseluruhan kontennya

⁶<https://github.com/psf/requests-html>

```

2     """Return all absolute hyperlinks within the home url"""
3
4     base_url = url_format_handler(url)
5     session = HTMLSession()
6     r = session.get(base_url, headers = {'User-Agent':
7         np.random.choice(user_agent_list)}, timeout=30)
8     res = list(r.html.absolute_links)
9
10    ## If r-html anchor failed, concate manually
11    res_final = [""]
12    for url in res:
13        if not str(url).startswith("http"):
14            res_final.append(str(base_url + url))
15        else:
16            res_final.append(str(url))
17
18    ## Check if domain expired/redirects
19    domain = base_url.split("//")[1]
20    if any(domain in url for url in res_final):
21        pass
22    else:
23        res_final = [""]
24
25    ## If website does not return hyperlink
26    if len(res) == 0 or len(res_final) == 0:
27        res = [""]
28    else:
29        res = res_final
30
31    return res

```

Kode 3.2: Fungsi *hyperlink gathering*

Fungsi *get_hyperlinks* yang terdapat pada Kode 3.2 menerima input berupa *string* URL dari halaman muka sebuah situs *merchant*. Link tersebut kemudian akan dimasukkan kedalam sebuah *helper function* bernama *url_format_handler* untuk melakukan normalisasi URL agar sesuai dengan format yang bisa diterima library *requests-html*. Kemudian sebuah *session* akan dibuat dengan situs tersebut. Pada *crawler* ini, dalam setiap pembuatan *session*, *user-agent* yang digunakan pada *header* akan selalu diubah secara *random*. Hal ini dilakukan untuk menghindari *traffic* yang berasal dari *crawler* dianggap sebagai suatu serangan DDoS⁷.

Setelah *session* terbuat dan *response* didapatkan, maka *hyperlink* dapat diekstrak dari *entity response* dengan memanfaatkan *parser* dari *requests-html* dengan memanggil *method html.aboslute_links*. *Hyperlink* yang telah diekstrak akan dikumpulkan dalam bentuk *list*.

Apabila *library* gagal mengembalikan *hyperlink* dengan *anchor*,⁸ maka pengga-

⁷Distributed Denial-of-Service terjadi saat seorang peretas membanjiri jaringan lalu lintas data sebuah server sehingga mengganggu operasi sistem tersebut

⁸URL yang dikembalikan tidak mengandung *base URL*, seperti misalnya hanya */contact-us* bukan nama-

bugan akan dilakukan secara manual. Apabila situs tidak mengirimkan *response*, maka fungsi ini akan mengembalikan sebuah *list* kosong. List yang mengandung *hyperlink* tersebut akan dimanfaatkan dalam beberapa proses lainnya.

3.3.3.2 *Dynamic Hyperlink Gatherer*

Untuk situs-situs yang dianggap memiliki konten dinamis, maka fungsi ini akan dipanggil. Fungsi ini merupakan fungsi yang akan dijalankan pada tahap awal pengumpulan *hyperlink* apabila menggunakan fungsi *get_hyperlink* tidak dapat mengembalikan *link* apapun dari proses ekstraksi statisnya. Fungsi ini juga akan dipanggil pada modul-modul *checker* apabila *crawler* tidak dapat menemukan *link-link* untuk 3 fitur utama situs (*tnc*, *contact us*, *about us*). Hal ini dikarenakan pada fungsi *get_hyperlinks_dynamic* ini, *WebDriver* akan mengecek juga pada *text* dari sebuah link tidak hanya pada URL nya, sehingga apabila *link* tidak mengandung kata kunci dari halaman yang dituju, maka URL tetap dapat terdeteksi. Contohnya adalah apabila *hyperlink* dari halaman *Contact Us* "https://situssaya.com" bukan "http://situssaya.com/contact-us" atau sejenisnya, melainkan merupakan URL yang *generic* seperti "http://situssaya.com/content=01".

```

1 def get_hyperlinks_dynamic(url):
2     global hyperlinks_dynamic
3     global dynamic_links
4     global dynamic_texts
5
6     try:
7         ## Do not gather (again) if has been gathered before
8         if hyperlinks_dynamic == True:
9             print("- Using previously gathered hyperlinks")
10            links, texts = dynamic_links, dynamic_texts
11        else:
12            print("- Gathering hyperlinks dynamically")
13            driver.set_page_load_timeout(60)
14            driver.get(url)
15            elems = driver.find_elements_by_xpath("//a[@href]")
16            links = []
17            texts = []
18            for elem in elems:
19                links.append(elem.get_attribute("href"))
20                texts.append(elem.text)
21
22            ## Set to true, collect hyperlinks
23            hyperlinks_dynamic = True
24            dynamic_links = links
25            dynamic_texts = texts
26        except Exception as e:
27            reset_browser()
28            print(e)

```

situs.com/contact-us.

```

29         links = []
30         texts = []
31
32     return links, texts

```

Kode 3.3: Fungsi *hyperlink gathering dynamic*

Fungsi *get_hyperlinks_dynamic* pada Kode 3.3 dimulai dengan memeriksa apakah fungsi ini sudah pernah dijalankan sebelumnya. Variabel *hyperlinks_dynamic* menandakan hal ini, apabila nilainya *True* maka fungsi ini sudah pernah dijalankan sehingga *dynamic_links* dan *dynamic_texts* tidak perlu diekstrak kembali dan cukup langsung dikembalikan. Sedangkan apabila nilainya *False*, maka *ChromeDriver* akan dijalankan untuk mengakses halaman muka yang dituju. *Browser* ini memiliki batasan *timeout* 60 detik⁹ untuk memuat seluruh konten yang ada pada halaman yang dituju. Jika *loading* belum selesai dalam rentang waktu tersebut, maka *TimeoutException* akan terjadi dan *browser* akan diinisiasi ulang tanpa ada konten yang terekstrak. Sedangkan apabila dalam 60 detik halaman berhasil dimuat, maka ekstraksi elemen akan dilakukan dengan melakukan pencarian berdasarkan *xpath*. *ChromeDriver* akan mencari elemen dengan *tag link* (<a href>) dan mengekstrak *attribute* (URL yang dituju pada *property href*) serta mengekstrak teks yang ditampilkan oleh *tag* tersebut. Setelah diekstrak, maka akan disimpan ke dalam variabel *links* dan *texts* dan akan dikembalikan di akhir fungsi.

3.3.4 Paragraph Extractor

Setiap *crawler* butuh untuk melakukan ekstraksi informasi, maka fungsi *paragraph extractor* akan. Tujuan utama fungsi ini adalah untuk melakukan proses *scraping* pada sebuah halaman. Ekstraksi informasi yang dilakukan ini nantinya akan dibutuhkan oleh 3 modul *checker* pada *crawler*, yaitu pada *TnC Checker*, *About Us Checker*, dan *Contact Info Checker*. Fungsi ini membantu untuk melakukan ekstraksi semua konten teks yang ada pada halaman-halaman yang bersangkutan. Ada dua jenis fungsi ekstraksi yang diimplementasi pada *crawler*, yaitu *paragraph_extractor* dan *paragraph_extractor_dynamic* dengan masing-masing tujuan yang sama namun cara yang berbeda. Fungsi *paragraph_extractor* digunakan sebagai esktraktor konten statis yang ada pada sebuah halaman dengan mengirim *HTTP request* dan melakukan *parsing* terhadap *body response* yang dikembalikan. Sedangkan sebaliknya, *paragraph_extractor_dynamic* berfungsi sebagai ekstraktor konten dinamis dengan memanfaatkan *ChromeDriver* seperti pada fungsi *get_hyperlinks_dynamic*.

⁹Rentang waktu ini dipilih berdasarkan hasil diskusi dengan tim DOKU dan dianggap sebagai waktu tunggu yang ideal

3.3.4.1 Static Paragraph Extractor

Static paragraph extractor fungsi *default* yang akan dijalankan saat ingin melakukan ekstraksi informasi. Fungsi ini memanfaatkan *library request* untuk mengirimkan sebuah *HTTP request* ke halaman yang dituju dan menampung *response content*nya. Elemen-elemen yang diperlukan untuk diekstrak adalah semua jenis teks yang ada pada halaman tersebut. Hal ini dilakukan dengan memanfaatkan *library BeautifulSoup4* untuk melakukan *parsing* dari *response content* yang didapatkan.

```

1 def paragraph_extractor(url):
2     ''' Extract all paragraph (texts) found in a given URL '''
3
4     try:
5         ## Send request
6         page = requests.get(url, headers = {'User-Agent':
7             np.random.choice(user_agent_list)}, timeout=30, verify=False)
8
9         ## Parse content
10        soup = bs(page.content, 'html.parser')
11        all_ps = soup.find_all("p") + soup.find_all("em") + soup.find_all("li") +
12            soup.find_all("address")\
13            + soup.find_all("h1") + soup.find_all("h2") + soup.find_all("h3") +
14            soup.find_all("h4") + soup.find_all("h5")\
15            + soup.find_all("h6") + soup.find_all("a") + soup.find_all("strong")
16
17        ## CloudFare email getter
18        if re.search('data-cfemail', str(soup)) is not None:
19            email_code = re.search('data-cfemail="(.*?)"', str(soup)).group(1)
20            email = str(decode_email(email_code))
21        else:
22            email = ""
23
24        ## Metadata getter
25        meta_property = soup.find("meta", property="og:description")
26        meta_name = soup.find("meta", {"name": "description"})
27
28        if meta_property is not None:
29            meta_property = soup.find("meta", property="og:description")["content"]
30        else:
31            meta_property = ""
32        if meta_name is not None:
33            meta_name = soup.find("meta", {"name": "description"})["content"]
34        else:
35            meta_name = ""
36
37        ## Div with class:address getter
38        div_address = str(soup.find("div", {"class": 'address'}))
39        if div_address is None:
40            div_address = ""
41
42        list_p = []
43        for p in all_ps:
44            list_p.append(unidecode.unidecode(p.getText()) + "\n")

```



```

43     ## Join all paragraphs element
44     paragraf = "".join(list_p)
45     paragraf += meta_property + meta_name + email + div_address
46     paragraf = paragraf.lower()
47
48     except Exception as e:
49         print(e)
50         paragraf = ""
51
52     return paragraf

```

Kode 3.4: Fungsi *paragraph extractor*

Pada Kode 3.4, yang akan dilakukan pertama kali adalah mengirimkan *request* ke URL terkait dengan beberapa parameter. Parameter tersebut adalah *user agent* pada *headers*, *timeout*, dan *verify*. *User agent* yang digunakan pada setiap *request* akan dirotasi seperti yang telah dijelaskan pada Subbab 3.3.3.1. Sedangkan *timeout* yang ditentukan untuk setiap *request* adalah 30 detik, jika melebihi waktu tersebut maka *ConnectionError* akan dimunculkan. Parameter *verify* digunakan untuk melakukan *bypass* terhadap situs-situs yang memiliki sertifikasi tidak baik agar tetap bisa dilakukan *crawling*. Hal ini sama seperti opsi “--ignore-certificate-error” yang ada pada *ChromeDriver*.

Memasuki bagian *parsing*, *beautifulsoup* akan melakukan *parsing* dokumen HTML yang diterima dari proses *request* yang dikirimkan. Adapun untuk mencari seluruh teks yang ada pada halaman tersebut dapat dilakukan dengan metode *find_all* yang dimiliki oleh *library BeautifulSoup*. Semua teks dengan *tag-tag* yang dirinci pada kode akan diekstrak dan disimpan dalam bentuk *string* pada variabel *all_ps*. Selain teks umum pada sebuah halaman, ada beberapa elemen lain yang perlu diekstrak karena dianggap dapat menyimpan informasi yang juga penting ataupun tidak terkekstrak pada tahap pertama. Elemen pertama adalah data *email* terenkripsi jika situs tersebut menggunakan layanan *CloudFlare*¹⁰, yang tentunya perlu didekripsi menggunakan fungsi *email.decode*. Elemen lainnya adalah *metadata* yang mengandung *description* serta *tag <div>* yang memiliki *class address*. Hal ini penting untuk diekstrak karena beberapa informasi tersebut dibutuhkan bagi modul-modul *checker* yang mengimplementasikan fungsi ini. Fungsi ini mengembalikan variabel *paragraf* yang merupakan satu string gabungan dari keseluruhan teks yang ada pada halaman tersebut.

3.3.4.2 Dynamic Paragraph Extractor

Dalam melakukan ekstraksi, tidak selamanya fungsi *paragraph extractor* biasa berhasil melakukan ekstraksi konten yang diinginkan. Hal ini dapat terjadi karena dua hal, antara

¹⁰Sebuah perusahaan *internet security* yang menyediakan berbagai macam jasa keamanan untuk situs. Dikutip dari <https://www.cloudflare.com>

memang halaman yang dituju tidak dapat diakses, ataupun konten di halaman tersebut belum lengkap. Merujuk pada alasan kedua, ketidaklengkapan konten dapat dikarenakan halaman tersebut sebenarnya bersifat dinamis. Halaman yang bersifat dinamis tidak akan langsung melakukan *rendering* terhadap semua elemen yang sebenarnya dimiliki. Halaman seperti ini biasanya membutuhkan waktu tunggu beberapa detik agar *script* yang ada pada halamannya dapat berjalan untuk melakukan *rendering* konten halaman secara utuh. Sayangnya, *delayed content* tersebut tidak dapat didapatkan dengan mengirim *request* biasa menggunakan *library requests*. Oleh karena itu dibuat fungsi *paragraph extractor dynamic* yang memanfaatkan *ChromeDriver*, sebuah *browser* yang dapat diotomasi, untuk melakukan *rendering* dan mendapatkan konten-konten dinamis tersebut.

```

1 def paragraf_extractor_dynamic(url):
2     ''' Extract all paragraph (texts) found in a given dynamic page URL '''
3
4     try:
5         ## Send request
6         print("-- Extracting paragraphs dynamically")
7         driver.set_page_load_timeout(60)
8         driver.get(url)
9
10        ## Parse content
11        soup = bs(driver.page_source, 'html.parser')
12
13        all_ps = driver.find_elements_by_tag_name("p") +
14        driver.find_elements_by_tag_name(
15            "em") + driver.find_elements_by_tag_name("li") +
16        driver.find_elements_by_tag_name("address") \
17            + driver.find_elements_by_tag_name("h1") +
18        driver.find_elements_by_tag_name(
19            "h2") + driver.find_elements_by_tag_name("h3") +
20        driver.find_elements_by_tag_name(
21            "h4") + driver.find_elements_by_tag_name("h5") \
22            + driver.find_elements_by_tag_name("h6") +
23        driver.find_elements_by_tag_name(
24            "a") + driver.find_elements_by_tag_name("strong") +
25        driver.find_elements_by_tag_name("span")
26
27        texts = []
28        for p in all_ps:
29            texts.append(p.text + "\n")
30
31        ## CloudFare email getter
32        if re.search('data-cfemail', str(soup)) is not None:
33            email_code = re.search('data-cfemail="(.*?)"', str(soup)).group(1)
34            email = str(decode_email(email_code))
35        else:
36            email = ""
37
38        ## Metadata getter
39        meta_property = soup.find("meta", property="og:description")
40        meta_name = soup.find("meta", {"name": "description"})

```

```

35
36     if meta_property is not None:
37         meta_property = soup.find("meta", property="og:description")["content"]
38     else:
39         meta_property = ""
40     if meta_name is not None:
41         meta_name = soup.find("meta", {"name": "description"})["content"]
42     else:
43         meta_name = ""
44
45     ## Div with class:address getter
46     div_address = str(soup.find("div", {"class": 'address'}))
47     if div_address is None:
48         div_address = ""
49
50     paragraf = "".join(texts)
51     paragraf += meta_property + meta_name + email + div_address
52     paragraf = paragraf.lower()
53 except Exception as e:
54     reset_browser()
55     print(e)
56     paragraf = ""
57
58 return paragraf

```

Kode 3.5: Fungsi *paragraph extractor dynamic*

Fungsi *paragraph_extractor_dynamic* yang ada pada Kode 3.5 sebenarnya sangat mirip dengan fungsi *paragraph_extractor* yang ada pada Kode 3.4. Yang membedakan antara keduanya hanya pada proses mengirimkan *request* dan mendapatkan *response*. Pada fungsi ini, *library requests* tidak digunakan, melainkan menggunakan *ChromeDriver* yang telah diinisiasi pada awal sesi. Alasan dibalik penggunaan ini adalah, *Selenium WebDriver* yang diimplementasikan menggunakan *ChromeDriver* ini memiliki kemampuan untuk melakukan ekstraksi terhadap keseluruhan konten, yaitu yang statis maupun dinamis. *ChromeDriver* akan menunggu selama waktu *timeout* yang telah diset untuk semua konten dapat terekstrak. Apabila sampai ke waktu yang telah ditentukan halaman yang dituju masih tidak memberikan respon ataupun belum menyelesaikan proses *loading*-nya, maka *TimeoutException* akan dikeluarkan.

Ekstraksi konten pada sekumpulan *tag* yang melambangkan teks dilakukan dengan *method find_elements_by_tag_name()*. *Method* ini akan mengembalikan sekumpulan *WebElement*¹¹ yang nantinya dapat diekstrak teks di dalamnya dengan memanggil *.text*. Sedangkan untuk ekstraksi konten-konten pelengkap lainnya tidak ada perbedaan dengan kode pada Kode 3.4 karena pada dasarnya *CloudFare email*, *metadata*, dan *tag div address* memang sudah merupakan elemen statis. Fungsi ini juga mengembalikan variabel *paragraf* yang merupakan satu string gabungan dari keseluruhan teks yang ada pada ha-

¹¹ Cara *WebDriver* merepresentasikan sebuah HTML *element*

laman tersebut.

3.4 Modul Ekstraksi Fitur

3.4.1 *Broken Link Checker*

Bagian ini memiliki fungsi untuk melakukan pengecekan terhadap jumlah *broken link* yang terdapat pada situs *merchant*. Fungsi ini memanfaatkan keluaran dari pengumpulan *hyperlink* yang dijalankan oleh Kode 3.2 ataupun Kode 3.3. *Library* yang digunakan pada fungsi ini adalah *grequests*¹², yang merupakan pengembangan dari *library requests* dan *gevent*¹³. *Library* ini digunakan untuk melakukan *HTTP request* secara asinkronus. *Request* secara asinkronus diperlukan untuk mempercepat proses pengecekan pada lebih dari satu URL.

```

1 def broken_link_score(df, hyperlinks):
2     """ Return score of broken link in a website
3         1 for >= 50 %
4         0 for < 50 %
5     """
6
7     print("Checking broken link...")
8
9     ## Avoid extraction of external links
10    avoid = pd.Series(hyperlinks).str.contains("wa.me") |
11    pd.Series(hyperlinks).str.contains("youtube") | \
12    pd.Series(hyperlinks).str.contains("linkedin") |
13    pd.Series(hyperlinks).str.contains("facebook") | \
14    pd.Series(hyperlinks).str.contains("cloudflare") |
15    pd.Series(hyperlinks).str.contains("twitter") | \
16    pd.Series(hyperlinks).str.contains("github") |
17    pd.Series(hyperlinks).str.contains("instagram") | \
18    pd.Series(hyperlinks).str.contains("tokopedia") |
19    pd.Series(hyperlinks).str.contains("bukalapak") | \
20    pd.Series(hyperlinks).str.match("tel") |
21    pd.Series(hyperlinks).str.contains("gitlab") | \
22    pd.Series(hyperlinks).str.contains("Tel") |
23    pd.Series(hyperlinks).str.contains("jobstreet") | \
24    pd.Series(hyperlinks).str.contains("download") |
25    pd.Series(hyperlinks).str.contains("google") | \
26    pd.Series(hyperlinks).str.contains("javaScript") |
27    pd.Series(hyperlinks).str.contains("_blank")
28    hyperlinks = list(pd.Series(hyperlinks)[~avoid].values)
29
30    ## If length of hyperlinks collected > 10, do sampling
31    if len(hyperlinks) > 10:
32        hyperlinks = sample(hyperlinks, 10)
33

```

¹²<https://pypi.org/project/grequests/>

¹³<https://pypi.org/project/gevent/>

```

25     ## Send request
26     rs = (grequests.get(x, \
27         headers = {}, \
28         timeout=20, verify=False) for x in hyperlinks)
29     rs_res = grequests.map(rs, size = 3)
30
31     links = {}
32     i = 0
33     if len(hyperlinks) == 1 and hyperlinks[0] == "":
34         links = "No hyperlinks gathered"
35     else:
36         ## Get all response code from sampled links
37         for response in rs_res:
38             try:
39                 links[response.request.url] = str(response)
40             except Exception as e:
41                 print(e)
42                 ## No response/timeout return this
43                 links[hyperlinks[i]] = 'No Response/Timeout'
44                 i += 1
45
46     status_not_ok = np.count_nonzero(np.array(rs_res, dtype=str) != '<Response
47     [200]>')
48     status_length = len(rs_res)
49
50     try:
51         ## Do scoring
52         score = status_not_ok/status_length*100
53     except Exception as e:
54         print(e)
55         score = 100
56
57     res_df = pd.DataFrame({"merchant_name": df['merchant_name'].values[0],
58         "broken_link_score": 0 if score < 50 else 1, \
59         "links_response": str(links)}, index=[0])
60
61     print("Broken links checked.\n")
62
63     return res_df

```

Kode 3.6: Fungsi *broken link checker*

Pada fungsi ini, yang dilakukan adalah melakukan pengecekan *response code* satu per satu terhadap kumpulan URL (*hyperlinks*) yang telah diekstrak pada proses sebelumnya. Jika *list hyperlink* mengandung lebih dari 10 *hyperlink*, maka akan dilakukan *sampling* dengan cara *random sample* sebanyak 10. Namun, jika *list* mengandung kurang dari 10 *hyperlink*, maka *sampling* tidak perlu dilakukan dan pengecekan dilakukan ke semua URL yang ada. Sebelum *request* dikirimkan ke URL yang ada pada list, URL akan terlebih dahulu disaring sehingga *link-link* eksternal yang ada pada situs tersebut tidak masuk ke dalam perhitungan fungsi¹⁴.

¹⁴Link yang mengarah ke pranala luar situs bukan merupakan bagian dari *scope* pengecekan tim DOKU

List yang telah disaring dan berisi kumpulan URL tersebut akan diolah oleh *library requests* dengan mengirimkan *request* secara asinkronus terhadap semua URL tersebut. Untuk setiap *request* akan diberikan *timeout* selama 20 detik dan dalam satu detiknya *requests* dapat mengirimkan sampai dengan 3 *request* secara bersamaan. Setelah *request* dikirimkan, *response code* dari setiap URL akan disimpan dalam sebuah *dictionary*. *Key* dari *dictionary* tersebut adalah URL dan *valuenya* adalah *response code* dari URL yang bersangkutan. Setelah semua *response code* didapatkan, maka dapat dihitung presentase *hyperlink* yang berstatus OK¹⁵ dengan membaginya dengan keseluruhan jumlah *hyperlink* yang diolah.

Pada akhirnya fungsi akan mengembalikan *dataframe* yang berisikan nama *merchant* (*merchant_name*), presentase *broken link* (*broken_link_score*) yang telah dibinerisasi¹⁶, serta *dictionary* yang mengandung *response code* dari keseluruhan URL yang diolah (*links_response*).

3.4.2 TnC Checker

Bagian ini memiliki fungsi untuk melakukan pengecekan informasi syarat dan ketentuan pada sebuah situs *merchant*. Informasi yang akan dicek adalah ada/tidaknya sebuah URL yang mengarah pada halaman syarat & ketentuan serta ada/tidaknya kata kunci yang menyebutkan tentang pengembalian barang atau uang (*refund policy*) di dalam halaman tersebut. Fungsi yang digunakan pada *tnc checker* ini adalah fungsi *tnc_score*¹⁷ untuk menjalankan keseluruhan algoritma *tnc checker* mulai dari mencari halaman syarat & ketentuan sampai mencari konten *refund policy* di dalamnya. Untuk melakukan pencarian konten *refund policy* di setiap kandidat halaman, maka ada dua fungsi *helper* lain yang dapat digunakan, yaitu *paragraph_extractor* dan fungsi *refund_policy_matcher*.

Fungsi *tnc score* menerima masukan berupa kumpulan *hyperlink* yang telah dikumpulkan pada proses sebelumnya. Fungsi ini akan melakukan pengecekan pada seluruh isi dari *list hyperlink* yang diberikan dan memeriksa apakah pada URL tersebut mengandung kata-kata yang telah ditentukan pada *keyword_tnc*. Apabila ada, maka URL-URL tersebut akan dikumpulkan ke dalam satu *list* baru dan ditandai sebagai kandidat halaman syarat & ketentuan.

Untuk mengekstraksi konten dari kandidat-kandidat halaman tersebut, maka fungsi *paragraph extractor* akan dipanggil untuk setiap halaman yang ada. Fungsi *paragraph extractor* akan mengembalikan sebuah *string* yang merupakan kumpulan teks dari halaman

¹⁵Didefinisikan dengan URL yang mengembalikan *response code* 200 OK

¹⁶Binerisasi menggunakan *threshold* 50%. Apabila jumlah *link* yang mati berada di atas atau sama dengan 50% maka akan ditransformasikan menjadi 1, sementara dibawah itu menjadi 0

¹⁷Dapat dilihat pada Lampiran

tersebut. Untuk setiap *string* yang diekstrak dari sebuah halaman, maka akan diperiksa eksistensi *refund policy*nya dengan menggunakan fungsi *refund_policy_matcher*. Cara kerja *refund_policy_matcher* dapat dilihat pada Kode 3.7.

Jika berhasil menemukan konten *refund policy* pada suatu halaman maka fungsi ini akan langsung mengembalikan hasilnya, yaitu sebuah *flag* yang menandakan bahwa pada situs *merchant* ini terdapat halaman syarat & ketentuan (*link_tnc_exist*: 1) dan di dalamnya terdapat kata kunci yang berhubungan dengan *refund policy* (*tnc_refund_policy_exist*: 1). Namun, apabila tidak ditemukan salah satu diantara dua komponen tersebut, *tnc_score* akan melanjutkan proses pencarian ke halaman muka dari situs tersebut, dan mencari kata kunci di halaman itu. Jika masih tidak ditemukan, maka proses pencarian akan diulang kembali namun menggunakan fungsi *paragraph_extractor_dynamic* dan *get_hyperlinks_dynamic*, dikarenakan ada kemungkinan situs tersebut sebenarnya memiliki konten dinamis yang tidak bisa diekstrak menggunakan metode statis.

```

1 def refund_policy_matcher(paragraf):
2     keyword_refund = ['refunds', 'refund', 'refund policy', 'return', 'returns',
3                       'return policy', 'pengembalian', 'pengembalian dana', 'mengembalikan dana',
4                       'dikembalikan', 'pengembalian', 'mengembalikan', 'retur', 'tukar', 'penukaran']
5
6     ## Tokenizing & cleaning paragraf into words
7     try:
8         lang = 'indonesian' if tb(paragraf).detect_language() == 'id' else
9         'english'
10        tokenizer = RegexpTokenizer(r'\w+')
11        words = tokenizer.tokenize(paragraf)
12
13        ## Removing stop words
14        stop_words = nltk.corpus.stopwords.words(lang)
15
16        cleaned_words = []
17        for word in words:
18            if word not in stop_words:
19                cleaned_words.append(word)
20
21    except Exception as e:
22        print(e)
23        cleaned_words = ""
24
25    mask = pd.Series(cleaned_words).str.contains("|".join(keyword_refund))
26
27    count_refund = np.count_nonzero(mask)
28    if count_refund > 0:
29        return 1
30    else:
31        return 0

```

Kode 3.7: Fungsi *refund policy matcher*

Fungsi *refund_policy_matcher* pada Kode 3.7 berfungsi untuk mencari kata kunci dari

masukan *string* yang diterima. Kata kunci yang dicari merupakan sekumpulan kata-kata yang berhubungan dengan proses pengembalian uang dan barang. Fungsi tersebut melakukan penyocokan kata kunci dengan sekumpulan *string* yang telah ditokenisasi pada halaman yang dicari. Tokenisasi memudahkan proses pencocokan kata karena *string* yang tadinya berbentuk paragraf telah dipecah menjadi kumpulan-kumpulan kata. Untuk melakukan *removal* terhadap kata-kata yang tidak berkaitan, pembersihan *stop words*¹⁸ juga dilakukan. Setelah itu, akan dihitung ada berapa banyak *keyword* yang ditemukan dalam kumpulan kata-kata yang telah diolah. Jika terdapat lebih dari 1, maka fungsi akan melakukan *return* 1, yang mana dianggap sebuah *keyword* terkait *refund policy* telah berhasil ditemukan.

3.4.3 About Us Checker

Bagian ini memiliki fungsi pengecekan yang paling sederhana. Pada bagian ini, *crawler* hanya akan melakukan pengecekan apakah situs *merchant* memiliki URL yang mengarah pada sebuah halaman yang berisi penjelasan tentang situs tersebut (*about us*, tentang kami, dan sebagainya) atau tidak. Proses pencarian halaman informasi ini dilakukan melalui dijalankan dengan melakukan pencocokan kata kunci dengan setiap URL yang ada pada kumpulan *hyperlink*.

```

1 def about_us_check(df, hyperlinks):
2     """Return boolean of about us link existance"""
3
4     print("Checking about us...")
5     keyword_about = ["about", "tentang"]
6     avoid = ["conditioner", "termurah", "termahal", "expired", "expire", "renewal"]
7
8     about_mask =
9         pd.Series(hyperlinks).str.lower().str.contains('|'.join(keyword_about)) \
10             & ~pd.Series(hyperlinks).str.lower().str.contains('|'.join(avoid))
11
12     about_count = 1 if np.count_nonzero(np.array(hyperlinks)[about_mask]) >= 1
13     else 0
14
15     if about_count == 1:
16         print("About Us Link: %s" % np.array(hyperlinks)[about_mask][0])
17
18     ## Check for inexact link
19     try:
20         if about_count < 1:
21             base_url = str(df['website'].values[0])
22             links, texts = get_hyperlinks_dynamic(url_format_handler(base_url))
23             if len(links) > 0:
24                 for i in range(len(links)):
25                     ## About Us Link Finder

```

¹⁸*Stop words* adalah sebuah kumpulan kata yang lazim digunakan di berbagai macam bahasa. *Stop words* biasanya dihilangkan agar pemrosesan dapat berfokus pada kata yang penting saja.


```

23         if
24             pd.Series(str(texts[i])).str.lower().str.contains('|'.join(keyword_about)).any():
25                 print("About Us Link: %s" % links[i])
26                 about_count = 1
27     except Exception as e:
28         print(e)
29         pass
30
31     res_df = pd.DataFrame({"merchant_name": df['merchant_name'].values[0],
32                           "link_about_us_exist": int(about_count)},
33                           index=[0])
34
35     print("About us checked.\n")
36
37     return res_df

```

Kode 3.8: Fungsi *about us check*

Pada Kode 3.8, fungsi *about_us_check* akan melakukan pengecekan pada setiap URL yang ada di masukan (*list hyperlinks*). Untuk setiap URL, dicek apakah terdapat kata kunci yang termuat dalam *list keyword_about*. Jika dapat ditemukan maka *flag about_count* akan diubah *valuenya* menjadi 1. Apabila tidak ditemukan, maka akan dicari *hyperlinks* baru menggunakan metode pencarian dinamis.

Pencarian dinamis dilakukan dengan cara mengumpulkan ulang seluruh *hyperlink* yang ada di halaman muka dengan fungsi *get_hyperlinks_dynamic*. Apabila ternyata hasil ekstraksi *hyperlink* mengeluarkan hasil yang berbeda dan *link about us* dapat ditemukan, maka dapat disimpulkan *tag-tag* link yang ada pada halaman utama dirender dengan metode dinamis ataupun memiliki *link* yang tidak mengandung *keyword* di URL nya seperti contoh yang sudah dijelaskan pada Subbab 3.3.4.2. Dengan menerapkan pencarian *hyperlink* secara dinamis, ekstraksi informasi menjadi lebih akurat dikarenakan adanya *double check* dan tidak hanya bergantung pada pencarian statis.

3.4.4 Contact Info Checker

Bagian ini berfungsi untuk melakukan pengecekan terhadap informasi kontak yang terdapat pada situs *merchant*. Informasi yang dicek adalah ada/tidaknya URL yang mengarah ke halaman kontak dan juga ada/tidaknya *email* serta nomor telepon pada situs tersebut. Pencarian *email* dan nomor telepon dilakukan pada dua *search spaces*, yaitu pada halaman kontak dan pada halaman muka situs itu sendiri. Dalam bagian ini, terdapat dua algoritma *regular expression* yang diimplementasikan, yaitu *telephone_matcher* dan juga *email_matcher* untuk mencari nomor telepon dan email pada paragraf yang diekstrak.

Untuk mencari URL yang mengarah ke halaman kontak, fungsi *contact_us_score*¹⁹ akan menjalankan algoritma yang sama seperti yang dijalankan oleh *tnc_score* pada Sub-

¹⁹Dapat dilihat pada Lampiran

printing terhadap email tersebut untuk keperluan *debugging* dan fungsi akan melakukan *return flag* 1.

3.5 Model Scorer

Modul ini memiliki fungsi untuk melakukan *scoring* terhadap situs *merchant* yang telah diekstrak fitur-fiturnya. Scoring dilakukan dengan menggunakan *trained classifier* yang telah diimplementasikan dan mengeluarkan *output* berupa *prediction class* dan *prediction probability* kelas itu berasal dari *class* berlabel 1.

```

1 def calculate_score(features):
2     """Return fraud prediction score of a website, 0-100 (Good - Bad)"""
3
4     ## Post API URL, change the URL to the host site URL
5     url = 'http://127.0.0.1:5000/api/v1/model'
6
7     ## Process Test Data
8     df = pd.DataFrame(features, index=[0])
9     columns = ['broken_link_score', 'link_contact_us_exist', 'cu_email_exist', \
10    'cu_phone_number_exist', 'link_about_us_exist', 'link_tnc_exist', \
11    'tnc_refund_policy_exist']
12     test_df = df[columns]
13     data = test_df.to_json()
14
15     ## Post to Model API
16     headers = {'content-type': 'application/json', 'Accept-Charset': 'UTF-8'}
17     req = requests.post(url, data=data, headers=headers).text
18     score = float(req)
19
20     df['prediction_prob'] = score
21     df['prediction_class'] = 1 if score >= 0.5 else 0
22     res = df
23
24     return res

```

Kode 3.11: Fungsi *model scorer*

Setelah semua fitur diekstrak pada proses-proses sebelumnya, maka fungsi *calculate_score* seperti yang dapat dilihat pada Kode 3.11 akan menjalankan tugasnya untuk melakukan *scoring* berdasarkan kumpulan fitur tersebut. Masukan yang diterima merupakan sebuah *dictionary* yang berisi pengembalian hasil dari modul-modul yang telah dijalankan. Hanya 7 kolom yang akan diambil dari keseluruhan kolom yang ada, yaitu fitur-fitur yang memang menjadi input bagi model *machine learning* dan akan dibahas lebih detil pada Bab 4.

Fitur-fitur tersebut akan ditranslasikan ke dalam bentuk JSON dengan format {"fitur1": value, "fitur2": value, ..., "fiturN": value}. JSON ini akan ditembak ke *end-*

*point model scoring*²⁰ dengan *POST request*, yang mana *endpoint* tersebut akan mengembalikan *prediction probability class 1* dari fitur-fitur tersebut. Jika *probability* berada di atas 0.5 maka akan dianggap sebagai *class 1*, jika tidak maka 0. Threshold dapat diatur secara manual pada fungsi ini apabila ada *preferensi* tersendiri setelah dilakukan evaluasi model lagi pada waktu yang akan datang.

²⁰Kode *endpoint* untuk *model scoring* dapat dilihat pada Lampiran

BAB 4

EKSPERIMEN DAN ANALISIS

4.1 Alur Kerja Eksperimen

Sebelum *web crawler* siap untuk masuk ke tahap produksi, maka perlu ada evaluasi yang mendalam terhadap model yang ditanamkan pada *web crawler* tersebut. Oleh karena itu, sebuah alur kerja eksperimen disusun untuk mengevaluasi performa model secara komprehensif. Adapun langkah-langkah yang dilakukan pada bagian ini meliputi proses pengolahan data awal dan evaluasi hyperparameter model menggunakan *nested cross-validation*.

Pada penelitian ini data yang digunakan untuk membuat model yang akan diimplementasikan di *crawler* merupakan situs-situs yang terdaftar sebagai *online merchant* DOKU. Eksperimen dilakukan menggunakan metode *supervised learning* dimana model akan menerima input untuk dipetakan menjadi output dari sebuah dataset yang telah dilabel. Dataset didapatkan dari tim *Early Detection Unit* (EDU) DOKU yang mereka kumpulkan dan verifikasi secara manual. Data awal diberikan dalam bentuk CSV dengan total 292 baris dan 3 kolom. Kolom terdiri dari nama *merchant* (*merchant_name*), alamat situs (*website*), dan juga status dari situs tersebut (*class*). Struktur data awal yang berbentuk file CSV tersebut digambarkan pada Tabel 4.1.

merchant_name	website	label
Nyetak.ID	www.nyetak.id	APPROVED
YoyoMats Indonesia	www.yoyomats.shopify.com	REJECTED

Tabel 4.1: Struktur tabel yang dikirimkan oleh tim EDU

Data awal yang diberikan oleh tim EDU bukan merupakan data yang bisa langsung dijadikan input (*training data*) untuk model *machine learning*. Data ini harus melalui proses ekstraksi fitur yang telah dijelaskan di Gambar 3.2 pada Bab 3. Proses ekstraksi fitur, akan menambah kolom pada Tabel 4.1 sebanyak 8 kolom tambahan, yang merupakan fitur hasil ekstraksi *crawler*. Data yang dihasilkan memiliki dimensi 292 baris dan 14 kolom. Namun, untuk proses evaluasi dan *training*, tidak semua fitur digunakan sebagai input bagi model *machine learning*, beberapa fitur hanya merupakan fitur informatif yang diminta oleh pihak DOKU untuk ditampilkan pada JSON *response crawler*. Keseluruhan fitur yang diekstrak beserta penjelasannya dapat dilihat pada Tabel 4.2 berikut.

No	Fitur	Deskripsi	Catatan
1.	merchant_name	Nama <i>merchant</i> yang bersangkutan. Nilai: String	-
2.	website	Alamat situs <i>merchant</i> yang bersangkutan. Nilai: String	-
3.	links_response	Sebuah <i>dictionary</i> yang menampilkan <i>response code</i> dari setiap <i>hyperlink</i> yang disampel Nilai: link (key): <i>response code</i> (value)	-
4.	label	Label kategori dari <i>merchant</i> yang bersangkutan. Nilai: 0(APPROVED)/1(REJECTED)	-
5.	broken_link_score	Binerisasi persentase tautan yang mengembalikan <i>response code</i> 200 dari semua <i>hyperlink</i> yang disampel. Nilai: 0(0-50%)/1(51-100%)	Pengambilan sampel dilakukan dengan <i>random sampling</i> (n = 10). <i>Hyperlink</i> didefinisikan sebagai semua tautan yang berada di beranda situs web.
6.	link_contact_us_exist	Menandakan adanya tautan halaman kontak di situs <i>merchant</i> . Nilai: 0/1	-

7.	cu_email_exist	Menandakan apakah adanya alamat email tercantum di situs <i>merchant</i> . Nilai: 0/1	Ruang pencarian meliputi halaman beranda dan halaman kontak.
8.	cu_phone_number_exist	Menandakan apakah adanya nomor telepon tercantum di situs <i>merchant</i> . Nilai: 0/1	Ruang pencarian meliputi halaman beranda dan halaman kontak.
9.	link_about_us_exist	Menandakan adanya tautan <i>About Us</i> (atau yang sejenisnya) di situs <i>merchant</i> Nilai: 0/1	-
10.	link_tnc_exist	Menandakan adanya tautan halaman syarat & ketentuan di situs <i>merchant</i> Nilai: 0/1	-
11.	tnc_refund_policy_exist	Menandakan adanya ketentuan pengembalian barang dan/atau uang di situs <i>merchant</i> Nilai: 0/1	Ruang pencarian meliputi halaman beranda dan halaman syarat & ketentuan.

Tabel 4.2: Fitur yang diekstrak dari proses *feature engineering*

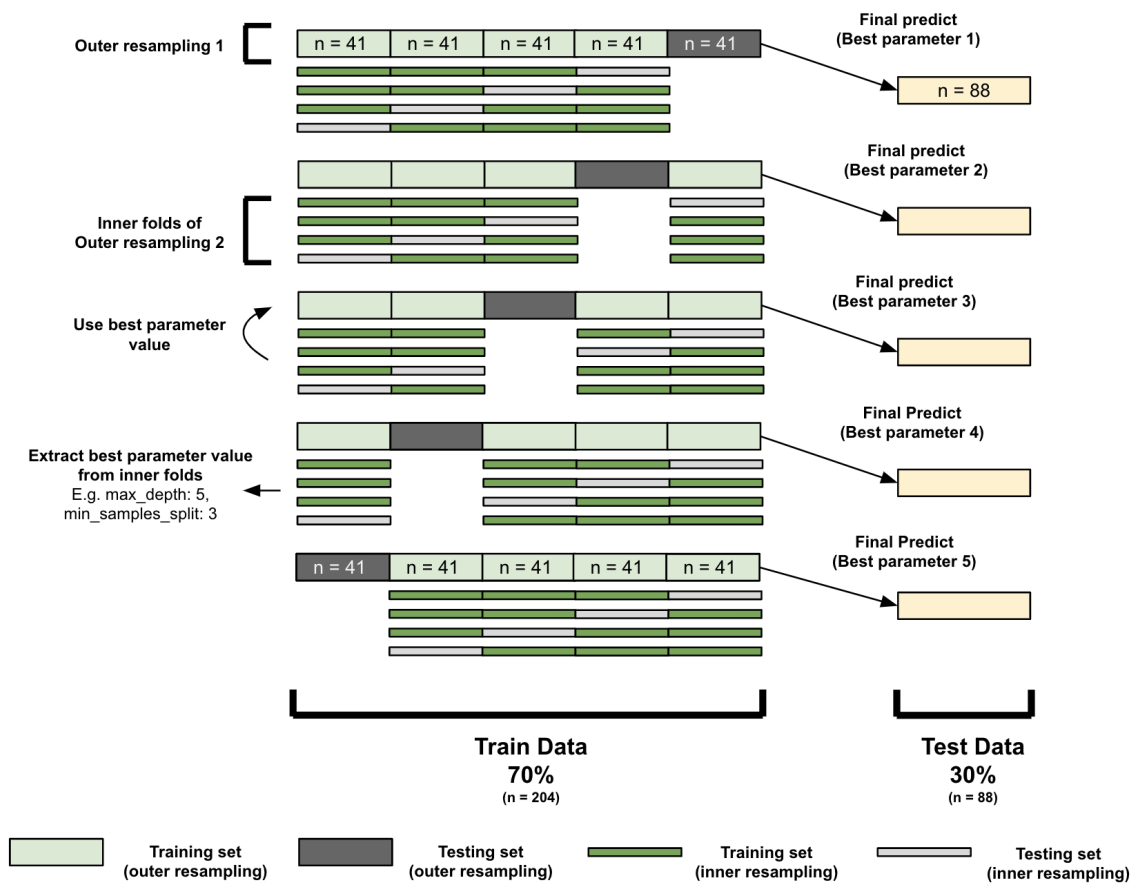
Dari Tabel 4.2, hanya akan ada 8 kolom yang diproses. Kolom-kolom tersebut adalah satu kolom kategori (*label*) dan 7 kolom fitur dari proses *feature engineering* (*broken_link_score*, *link_tnc_exist*, *tnc_refund_policy_exist*, *link_about_us_exist*, *link_contact_us_exist*, *cu_email_exist*, *cu_telephone_number_exist*). Delapan kolom ini akan digunakan sebagai masukan (*training data*) untuk model-model *machine learning* yang akan dievaluasi. Sedangkan ketiga kolom lainnya (*merchant_name*, *website*, *links_response*) hanya akan ditampilkan pada JSON *response crawler* sebagai informasi pelengkap.

4.2 Evaluasi Kinerja Model

Pada tahapan evaluasi kinerja model, ada empat jenis model *classifier* yang digunakan untuk bereksperimen. Diantaranya adalah *decision tree classifier*, *bernoulli naive bayes classifier*, *extreme gradient boost classifier*, dan *random forest classifier*. Training dan evaluasi setiap model dilakukan pada sebuah *notebook* dengan menggunakan bantuan library *scikit-learn*. Metode yang digunakan untuk memilih model terbaik beserta *hyperparameter*nya adalah *nested cross-validation*.

Pada setiap model, akan dipilih beberapa kandidat parameter beserta nilai-nilainya. Pemilihan parameter didasarkan pada pemahaman penulis akan pengaruh parameter yang bersangkutan pada proses prediksi yang akan dilakukan oleh model. Proses evaluasi ini menggunakan 292 baris data yang terdiri atas 215 data dengan label 0 dan 72 data dengan label 1. Pada langkah awal, data ini akan dibagi menjadi dua yaitu 70% *train data* dan 30% *test data*. *Test data* ini tidak akan digunakan sampai langkah terakhir evaluasi, sedangkan untuk *train data*, merujuk pada metode *nested cross-validation*, akan distribusikan dengan metode *stratified* ke dalam 5 fold pada *outer loop* dan pada setiap *train data* di setiap *fold* (*resampling*) pada *outer loop*, data akan dibagi lagi menjadi 4 fold (*resampling*) di *inner loop*. Setelah mendapatkan rata-rata skor terbaik untuk kombinasi parameter x_1 , x_2 , ..., x_n untuk setiap *inner loop* dengan menggunakan metode *hyperparameter tuning* Grid Search (*GridSearchCV*¹), maka kombinasi parameter tersebut akan difit ke *training data* pada *outer loop* yang bersangkutan (bukan keseluruhan *training data*) dan disimpan hasilnya. Untuk memastikan stabilitas model, maka 5 kandidat kombinasi parameter terbaik dari 5 fold yang ada di *outer loop* akan difit ke keseluruhan *training data* dan dites ke *testing data* yang belum pernah digunakan. Dari situ bisa didapatkan kesimpulan skor terbaik beserta parameter terbaiknya. Metrik yang akan digunakan untuk mengevaluasi skor adalah *Area Under the Curve* (AUC). Ilustrasi terkait mekanisme ini dapat dilihat lebih lengkap pada Gambar 4.1.

¹https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html

Gambar 4.1: Skema *nested cross-validation*

4.2.1 Decision Tree Classifier

Pada percobaan ini, *library* yang akan digunakan adalah *DecisionTreeClassifier*² yang dimiliki oleh *scikit-learn*. Penggunaannya cukup langsung mengimpor model tersebut ke dalam *notebook*.

Parameter yang akan dipilih sebagai kandidat di model ini adalah *max_depth*, *min_samples_split*, dan *min_samples_leaf*. Ketiga parameter ini dipilih karena dianggap dapat memberikan variasi kombinasi kedalaman *tree* yang dibuat. Kedalaman *tree* yang dibuat akan berpengaruh pada *overfitting/underfitting* yang terjadi saat model melakukan *fitting* terhadap *training data*. Kumpulan nilai parameter yang dicoba pada proses ini adalah sebagai berikut:

- **max_depth:** [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15]
- **min_samples_split:** [2, 5, 10, 15, 20, 25, 30]
- **min_samples_leaf:** [1, 3, 5, 10, 15]

²<https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html>

4.2.1.1 Parameter Terbaik Hasil Grid Search pada Inner Loop

outer_resampling	max_depth	min_samples_split	min_samples_leaf	inner_avg_score
1	3.0	25.0	3.0	0.912866
2	4.0	20.0	1.0	0.944893
3	3.0	25.0	5.0	0.891810
4	4.0	20.0	5.0	0.883346
5	4.0	15.0	5.0	0.885606

Tabel 4.3: Kombinasi parameter terbaik dari *inner loop* model Decision Tree Classifier

Pada Tabel 4.3, dapat dilihat bahwa kombinasi parameter terbaik dikeluarkan oleh *inner loop* yang terdapat pada *outer resampling* (*fold*) 2. Kombinasi parameter terbaik yang dimiliki adalah *max_depth* dengan nilai 4, *min_samples_split* dengan nilai 20, dan *min_samples_leaf* dengan nilai 1. Skor AUC rata-rata yang berhasil dihasilkan oleh kombinasi parameter ini di *inner loop* pada *outer resampling* 2 sebesar **0.944893**. Sedangkan, untuk rata-rata dari keseluruhan hasil *inner loop* pada seluruh *outer resampling* adalah **0.9037042**.

4.2.1.2 Hasil Test pada Outer Loop dan Test Data

outer_resampling	outer_test_score	test_data_score
1	0.910606	0.937672
2	0.780303	0.943871
3	0.933333	0.937672
4	0.934848	0.943871
5	0.970219	0.928375

Tabel 4.4: Percobaan kombinasi parameter terbaik ke *outer loop* & *test data* model Decision Tree Classifier

Hasil yang ditampilkan pada Tabel 4.4 merupakan hasil prediksi menggunakan parameter terbaik yang telah didapatkan dari *inner loop* CV yang dilakukan di Tabel 4.3. Prediksi dilakukan ke *testing set* yang ada pada setiap *outer resampling* dan juga ke 25% *test data* yang sudah dipisahkan sejak awal proses. Hasil prediksi terbaik terhadap *testing set* pada *outer resampling* ada di kombinasi parameter 4 dengan skor AUC sebesar **0.934848**. Sedangkan untuk hasil prediksi ke *test data*, hasil tertinggi ada pada 2 kombinasi parameter, yaitu kombinasi 2 dan 4 dengan skor AUC sebesar **0.943871**. Rata-rata skor pada kolom *outer_test_score* adalah **0.9058618** dan rata-rata skor pada kolom *test_data_score*

adalah **0.9382922**.

4.2.2 *Random Forest Classifier*

Pada percobaan ini, *library* yang akan digunakan adalah *RandomForestClassifier*³ yang dimiliki oleh *scikit-learn*. Penggunaannya cukup langsung mengimpor model tersebut ke dalam *notebook*.

Parameter yang akan dipilih sebagai kandidat di model ini sebenarnya hampir sama dengan Decision Tree. Namun ada dua parameter tambahan yang merupakan ciri khas sebuah *ensemble* model, yaitu *n_estimators* dan *bootstrap*. Parameter-parameter yang dipilih ini menentukan seberapa banyak *learner* yang akan dibangun dan seperti apa setiap *learner* tersebut akan dibangun. Kumpulan nilai parameter yang dicoba pada proses ini adalah sebagai berikut:

- **n_estimators**: [10, 25, 50, 75, 100]
- **max_depth**: [3, 5, 7, 10]
- **bootstrap**: [True, False]
- **min_samples_split**: [3, 5, 7]

4.2.2.1 Parameter Terbaik Hasil Grid Search pada Inner Loop

outer _resampling	n_estimators	max_depth	bootstrap	min_samples _split	inner_avg _score
1	50.0	7.0	True	7.0	0.938858
2	10.0	7.0	True	3.0	0.954415
3	50.0	3.0	False	7.0	0.935697
4	10.0	3.0	False	3.0	0.922688
5	10.0	3.0	False	3.0	0.917045

Tabel 4.5: Kombinasi parameter terbaik dari *inner loop* model Random Forest Classifier

Pada Tabel 4.5, dapat dilihat bahwa kombinasi parameter terbaik dikeluarkan oleh *inner loop* yang terdapat pada *outer resampling* (*fold*) 2. Kombinasi parameter terbaik yang dimiliki adalah *n_estimators* dengan nilai 10, *max_depth* dengan nilai 7, *bootstrap* dengan nilai True, dan *min_samples_split* dengan nilai 3. Skor AUC rata-rata yang berhasil dihasilkan oleh kombinasi parameter ini di *inner loop* pada *outer resampling* 2 sebesar **0.954415**. Sedangkan, untuk rata-rata dari keseluruhan hasil *inner loop* pada seluruh

³<https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>

outer resampling adalah **0.9337406**.

4.2.2.2 Hasil Test pada Outer Loop dan Test Data

outer_resampling	outer_test_score	test_data_score
1	0.915152	0.942837
2	0.824242	0.931474
3	0.918182	0.943526
4	0.966667	0.953168
5	0.959248	0.953168

Tabel 4.6: Percobaan kombinasi parameter terbaik ke *outer loop* & *test data* model Random Forest Classifier

Hasil yang ditampilkan pada Tabel 4.6 merupakan hasil prediksi menggunakan parameter terbaik yang telah didapatkan dari *inner loop* CV yang dilakukan di Tabel 4.5. Prediksi dilakukan ke *testing set* yang ada pada setiap *outer resampling* dan juga ke 25% *test data* yang sudah dipisahkan sejak awal proses. Hasil prediksi terbaik terhadap *testing set* pada *outer resampling* ada di kombinasi parameter 4 dengan skor AUC sebesar **0.966667**. Sedangkan untuk hasil prediksi ke *test data*, hasil tertinggi ada pada 2 kombinasi parameter, yaitu kombinasi 4 dan 5 dengan skor AUC sebesar **0.953168**. Rata-rata skor pada kolom *outer_test_score* adalah **0.9166982** dan rata-rata skor pada kolom *test_data_score* adalah **0.9448346**.

4.2.3 *Extreme Gradient Boost Classifier*

Pada percobaan ini, *library* yang akan digunakan adalah *XGBoost Classifier*⁴ untuk Python. Penggunaan dari *library* ini adalah dengan melakukan instalasi via *package manager* Python (*pip*) lalu dapat langsung diimport ke dalam *notebook*.

Karena XGBoost juga merupakan sebuah *tree-based mode*, maka parameter yang akan dipilih sebagai kandidat di model ini juga sebenarnya hampir sama dengan Decision Tree. Namun ada satu *parameter* yang merupakan ciri khas dari algoritma *boosting*, yaitu *learning_rate*. Parameter ini membantu untuk meregularisasikan model. *Learning_rate* mereduksi pengaruh dari setiap individu *learner (tree)* sehingga membuka ruang bagi *tree* berikutnya untuk dapat melakukan improvisasi terhadap model [25]. Kumpulan nilai parameter yang dicoba pada proses ini adalah sebagai berikut:

- **learning_rate:** [0.05, 0.10, 0.15, 0.20, 0.25, 0.30]

⁴https://xgboost.readthedocs.io/en/latest/python/python_api.html

- **max_depth**: [3, 4, 5, 6, 8, 10, 12, 15]
- **min_child_weight**: [1, 3, 5, 7]

4.2.3.1 Parameter Terbaik Hasil Grid Search pada Inner Loop

outer_resampling	learning_rate	max_depth	min_child_weight	inner_avg_score
1	0.10	3.0	5.0	0.924282
2	0.05	3.0	5.0	0.936703
3	0.05	3.0	5.0	0.908294
4	0.25	3.0	5.0	0.916014
5	0.20	3.0	5.0	0.911364

Tabel 4.7: Kombinasi parameter terbaik dari *inner loop* model XGBoost Classifier

Pada Tabel 4.7, dapat dilihat bahwa kombinasi parameter terbaik dikeluarkan oleh *inner loop* yang terdapat pada *outer resampling (fold)* 2. Kombinasi parameter terbaik yang dimiliki adalah *learning_rate* dengan nilai 0.05, *max_depth* dengan nilai 3, dan *min_child_weight* dengan nilai 5. Skor AUC rata-rata yang berhasil dihasilkan oleh kombinasi parameter ini di *inner loop* pada *outer resampling* 2 sebesar **0.936703**. Sedangkan, untuk rata-rata dari keseluruhan hasil *inner loop* pada seluruh *outer resampling* adalah **0.9193314**.

4.2.3.2 Hasil Test pada Outer Loop dan Test Data

outer_resampling	outer_test_score	test_data_score
1	0.927273	0.952479
2	0.868182	0.953168
3	0.921212	0.953168
4	0.895455	0.945937
5	0.949843	0.945937

Tabel 4.8: Percobaan kombinasi parameter terbaik ke *outer loop* & *test data* model XGBoost Classifier

Hasil yang ditampilkan pada Tabel 4.8 merupakan hasil prediksi menggunakan parameter terbaik yang telah didapatkan dari *inner loop* CV yang dilakukan di Tabel 4.7. Prediksi dilakukan ke *testing set* yang ada pada setiap *outer resampling* dan juga ke 25% *test data* yang sudah dipisahkan sejak awal proses. Hasil prediksi terbaik terhadap *testing set* pada *outer resampling* ada di kombinasi parameter 5 dengan skor AUC sebesar **0.949843**.

Sedangkan untuk hasil prediksi ke *test data*, hasil tertinggi ada pada 2 kombinasi parameter, yaitu kombinasi 2 dan 3 dengan skor AUC sebesar **0.953168**. Rata-rata skor pada kolom *outer_test_score* adalah **0.912393** dan rata-rata skor pada kolom *test_data_score* adalah **0.9501378**.

4.2.4 Bernoulli Naive Bayes Classifier

Pada percobaan ini, *library* yang akan digunakan adalah *BernoulliNB*⁵ untuk Python. Penggunaan dari *library* dapat langsung diimport ke dalam *notebook*.

Pada dasarnya semua algoritma Naive Bayes hampir tidak memiliki *hyperparameter* untuk dilakukan *tuning*. Hal ini dikarenakan adanya asumsi independen antara satu fitur dengan fitur lainnya dan *decision rule* yang dimilikinya. Tanpa menentukan *hyperparameter* sebenarnya Bernoulli Naive Bayes sudah akan cukup mengeneralisir hasil prediksinya. Namun ada beberapa *setup* parameter yang dapat dilakukan di awal yang dapat sedikit banyak mengubah hasil model. Parameter-parameter dan nilai-nilainya adalah sebagai berikut:

- **alpha**: [0.5, 1.0]
- **fit_prior**: [True, False]

4.2.4.1 Parameter Terbaik Hasil Grid Search pada Inner Loop

outer_resampling	alpha	fit_prior	inner_avg_score
1	1.0	True	0.923498
2	0.5	True	0.946604
3	1.0	True	0.931831
4	0.5	True	0.927207
5	0.5	True	0.922348

Tabel 4.9: Kombinasi parameter terbaik dari *inner loop* model Bernoulli Naive Bayes Classifier

Pada Tabel 4.9, dapat dilihat bahwa kombinasi parameter terbaik dikeluarkan oleh *inner loop* yang terdapat pada *outer resampling (fold)* 2. Kombinasi parameter terbaik yang dimiliki adalah *alpha* dengan nilai 0.5 dan *fit_prior* dengan nilai True. Skor AUC rata-rata yang berhasil dihasilkan oleh kombinasi parameter ini di *inner loop* pada *outer resampling* 2 sebesar **0.946604**. Sedangkan, untuk rata-rata dari keseluruhan hasil *inner loop* pada seluruh *outer resampling* adalah **0.9302976**.

⁵https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.BernoulliNB.html

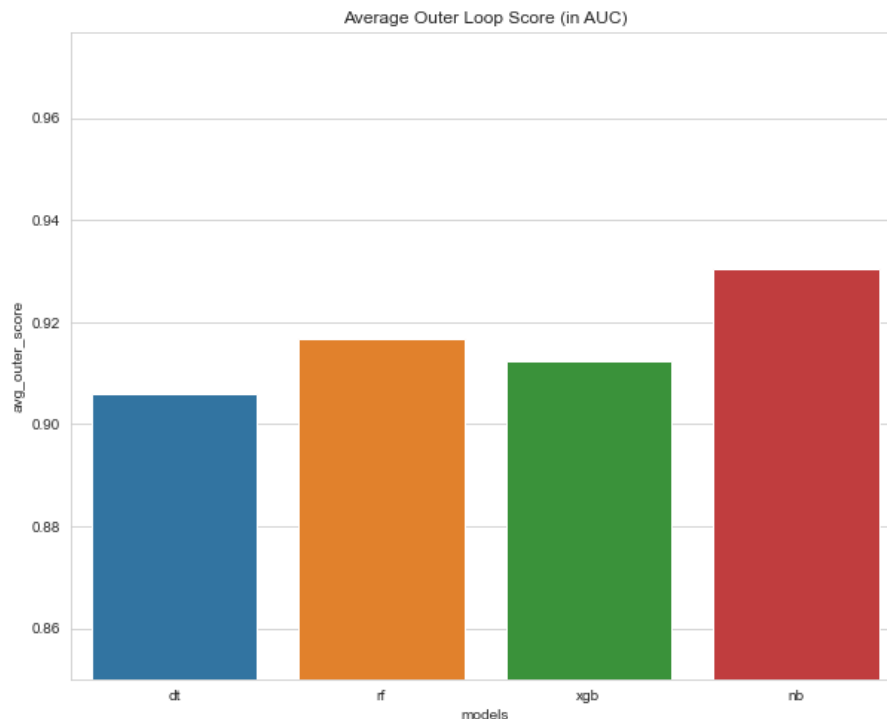
4.2.4.2 Hasil Test pada Outer Loop dan Test Data

outer_resampling	outer_test_score	test_data_score
1	0.963636	0.953168
2	0.877273	0.953168
3	0.930303	0.953168
4	0.924242	0.953168
5	0.956113	0.953168

Tabel 4.10: Percobaan kombinasi parameter terbaik ke *outer loop* & *test data* model Bernoulli Naive Bayes Classifier

Hasil yang ditampilkan pada Tabel 4.10 merupakan hasil prediksi menggunakan parameter terbaik yang telah didapatkan dari *inner loop* CV yang dilakukan di Tabel 4.9. Prediksi dilakukan ke *testing set* yang ada pada setiap *outer resampling* dan juga ke *25% test data* yang sudah dipisahkan sejak awal proses. Hasil prediksi terbaik terhadap *testing set* pada *outer resampling* ada di kombinasi parameter 1 dengan skor AUC sebesar **0.963636**. Sedangkan untuk hasil prediksi ke *test data*, hasil tertinggi ada pada semua kombinasi parameter dengan skor AUC sebesar **0.953168**. Rata-rata skor pada kolom *outer_test_score* adalah **0.9303134** dan rata-rata skor pada kolom *test_data_score* adalah **0.953168**.

4.2.5 Kinerja Model Terbaik



Gambar 4.2: Rata-rata skor AUC yang didapatkan setiap model pada *outer loop*

Untuk dapat melihat performa model terbaik berdasarkan hasil generalisasi yang dicapai menggunakan *nested cross validation*, maka kita hanya perlu melihat hasil rata-rata dari *outer_test_score* setiap model. Pada Gambar 4.2, dapat dilihat bahwa semua model memiliki rata-rata AUC di atas 0.9. Hasil ini menunjukkan bahwa keseluruhan model sebenarnya memiliki performa yang stabil dan baik. Semua model telah berhasil mencapai probabilitas di atas 90% untuk dapat mengurutkan keluaran prediksi secara tepat.

Decision Tree Classifier merupakan model terendah dengan skor 0.9058618, disusul oleh XGBoost Classifier dengan 0.912393, lalu Random Forest Classifier dengan 0.9166982, dan skor tertinggi ada pada Bernoulli Naive Bayes Classifier dengan skor 0.9303134.

4.2.6 Kombinasi Hyperparameter Terbaik

Untuk mendapatkan *hyperparameter* terbaik dari 5 kombinasi yang dihasilkan setiap model, maka penulis melakukan *testing* menggunakan setiap kombinasi tersebut ke sebuah *test data* yang sudah dipisahkan dari awal. Data dilakukan *training* pada keseluruhan *train data* (70% dari *dataset* yang dimiliki penulis) menggunakan setiap kombinasi parameter dan dilakukan *testing* ke *test data* serta dicatat skor AUC terbaik yang dihasilkan.

Dari 4 model yang diuji, dapat dilihat bahwa nilai AUC terbaik yang diambil dari kombinasi parameter terbaik keempat model berada pada rentang yang cukup tinggi saat memprediksi *test data*. Skor paling rendah ada pada model Decision Tree Classifier yaitu 0.943871 dan skor paling tinggi dipegang oleh 3 model lainnya secara bersamaan yaitu Random Forest Classifier, XGBoost Classifier, dan Bernoulli Naive Bayes Classifier dengan skor 0.953168. Kemampuan keempat model dalam melakukan prediksi jika diukur dengan AUC, menunjukkan bahwa keempat model tersebut memiliki kemampuan yang baik dalam mengurutkan probabilitas prediksi data positif untuk berada lebih tinggi dibandingkan probabilitas prediksi data negatif. Oleh karena itu, penulis merasa perlu melihat lebih dalam kepada metrik-metrik lainnya, seperti melihat Akurasi, Presisi, Recall, dan F1 dan Confusion Matrix dari setiap model untuk dapat menarik kesimpulan model mana yang merupakan model terbaik terhadap *test data* dengan kombinasi parameter yang dimilikinya.

model	auc	accuracy	balanced _accuracy	precision	recall	f1_score
dt	0.944	0.920	0.871	0.895	0.773	0.829
rf	0.953	0.909	0.879	0.818	0.818	0.818
xgb	0.953	0.898	0.871	0.783	0.818	0.800
nb	0.953	0.864	0.864	0.679	0.864	0.760

Tabel 4.11: Ringkasan metrik terbaik keempat model

model	true_positive	true_negative	false_positive	false_negative
dt	17	64	2	5
rf	18	62	4	4
xgb	18	61	5	4
nb	19	57	8	3

Tabel 4.12: Ringkasan *confusion matrix* terbaik keempat model

Melihat ringkasan yang ditampilkan pada Tabel 4.11 dan Tabel 4.12, ada dua model yang dapat disorot, yaitu Random Forest Classifier dan Bernoulli Naive Bayes Classifier. Model pertama adalah Random Forest Classifier (kode model **rf**). Model ini memiliki nilai yang stabil dan seimbang di berbagai macam metrik. Tidak ada satupun metrik yang jatuh di bawah angka 0.8. Mengasumsikan *cost* untuk *false positive* dan *false negative* sama, model ini memiliki performa paling bagus jika dilihat dari akurasi yang telah diseimbangkan⁶. *Reliability* model dalam memprediksi label positif yang ditampilkan

⁶Balanced accuracy = (TPR+TNR) / 2

oleh *precision* juga telah menunjukkan angka yang cukup baik di 0.818. Kemampuan model untuk menangkap keseluruhan label positif juga sudah baik dengan nilai *recall* sebesar 0.818. Model ini juga memiliki jumlah kesalahan prediksi yang cukup sedikit, yaitu hanya 4 *false positive* dan 4 *false negative*. Random Forest Classifier dengan kombinasi parameter 4⁷, dapat dikatakan sebagai model yang terbaik diantara keempat model yang ada dalam melakukan prediksi terhadap *test data*.

Selain model *rf*, ada juga model lain yang dapat disorot yaitu Bernoulli Naive Bayes Classifier (kode model **nb**). Pada fase awal implementasi, hasil prediksi yang dikeluarkan oleh *web crawler* nantinya akan kembali dievaluasi oleh *end-user*, yaitu tim EDU DOKU. Namun, tim EDU hanya akan melakukan evaluasi manual terhadap situs-situs *merchant* yang diprediksi sebagai *merchant fraud* oleh *web crawler*. Hal ini dilakukan karena tim EDU DOKU memiliki penilaian *cost* yang berbeda terhadap *false negative* dan *false positive* yang dikeluarkan oleh model. Mereka menganggap bahwa adanya cukup banyak *false positive* masih lebih sedikit secara *cost* dibandingkan dengan adanya *merchant fraud* yang diloloskan, walaupun sedikit, sebagai *merchant* yang tidak *fraud* (*false negative*)⁸. Dalam kata lain, *cost false negative* akan menjadi lebih tinggi dibandingkan *cost false positive*. Hal ini lumrah terjadi dalam penerapan *machine learning* untuk melakukan *decision making* di dunia bisnis. Oleh karena itu, untuk menyesuaikan preferensi tersebut, sebuah model yang memiliki *recall* tertinggi dapat dipilih, meskipun memiliki nilai *precision* yang cukup rendah, namun tetap *acceptable*. Model tersebut adalah model Naive Bayes Classifier dengan kombinasi parameter 1⁹. Dengan model ini, angka *recall* dapat mencapai 0.864 dan hanya akan memiliki 3 situs *merchant fraud* yang salah diklasifikasikan. Model ini juga mengembalikan angka *true positive* terbanyak dengan jumlah 19 prediksi positif.

Dari hasil evaluasi di atas, dapat disimpulkan bahwa terdapat dua kandidat model yang dapat dipilih untuk diimplementasikan di *web crawler*. Sejauh ini, model yang diterapkan pada *web crawler* adalah Bernoulli Naive Bayes Classifier. Hal ini dipilih berdasarkan beberapa pertimbangan yang didiskusikan dalam beberapa pertemuan bersama tim Engineering dan tim Product dari pihak DOKU. Diantaranya, model Bernoulli Naive Bayes lebih intuitif dan sederhana untuk dipahami cara bekerjanya. Selain itu, evaluasi model menggunakan 5 jenis kombinasi parameter pada *inner loop*, *outer loop*, dan *test data* juga menunjukkan hasil yang paling konsisten dibandingkan model lainnya. Rata-rata skor AUC pada *inner loop* mencapai angka 0.9302976, sedangkan pada *outer loop* sebesar 0.9303134, dan pada *test data* sebesar 0.953168. Untuk skor individualnya, Bernoulli Naive Bayes Classifier juga tidak memiliki hasil prediksi dengan skor AUC di bawah

⁷Merujuk pada Tabel 4.6

⁸Berdasarkan informasi yang disampaikan oleh Product Owner DOKU, Pak Reza Farasdak

⁹Merujuk pada Tabel 4.10

0.877273. Dengan kombinasi parameter yang cukup sederhana, model *nb* dapat menunjukkan performa yang baik dan telah berhasil mencapai tujuan awal dari pengembangan *web crawler* untuk melakukan klasifikasi situs *merchant* yang *fraud* secara tepat.

BAB 5

PENUTUP

5.1 Kesimpulan

Penelitian ini telah berhasil menciptakan sebuah aplikasi *web crawler* yang dapat digunakan untuk melakukan ekstraksi fitur pada sebuah situs *merchant* dan melakukan *scoring* terhadap situs tersebut menggunakan algoritma *machine learning*. Aplikasi dibangun dengan memanfaatkan sebuah *micro web framework* berbasis Python yang dinamakan Flask. Aplikasi ini dibuat untuk mempersingkat proses verifikasi pendaftaran *merchant* DOKU. Terdapat empat tahapan yang akan dilewati dalam penggunaan *web crawler*. Dimulai dari memasukkan URL situs *merchant*, proses ekstraksi *hyperlink*, proses ekstraksi fitur situs *merchant*, dan proses prediksi kategori (*scoring*) situs *merchant* menggunakan model *machine learning*.

Pada penelitian ini, fitur-fitur yang dapat digunakan untuk melakukan prediksi kategori sebuah situs *merchant* adalah sekumpulan 7 fitur biner. Pada dasarnya fitur-fitur ini merupakan sebuah *completeness checker*. Fitur-fitur ini disusun dan ditentukan bersama dengan tim EDU DOKU dengan mereferensikan kepada SOP pengecekan situs *merchant* yang selama ini mereka miliki. Fitur-fitur tersebut, diantaranya adalah 'broken_link_score', 'link_contact_us_exist', 'cu_email_exist', 'cu_phone_number_exist', 'link_about_us_exist', 'link_tnc_exist', 'tnc_refund_policy_exist' yang mana penjelasan lebih lengkap setiap fiturnya telah dibahas pada Bab 4.

Model *machine learning* yang digunakan pada penelitian ini adalah Bernoulli Naive Bayes Classifier. Model ini dipilih setelah melalui tahap evaluasi kinerja model menggunakan 292 data yang telah dilabel untuk kedua kategori. Model ini memiliki nilai AUC terbaik sebesar 0.953 dan akurasi sebesar 0.864. Untuk setiap situs *merchant*, rata-rata waktu yang diperlukan *web crawler* untuk melakukan prediksi adalah 20 detik. Dengan kecepatan prediksi tersebut dan ketepatan prediksi yang telah terbukti baik, adanya *web crawler* ini dapat meningkatkan efisiensi proses pengecekan situs yang dilakukan oleh tim EDU. Dalam 1 jam, *web crawler* dapat memeriksa 240 situs atau 235 situs lebih banyak dibandingkan pengecekan manual yang hanya sebanyak 5 situs. Adanya *web crawler* dapat meningkatkan efisiensi *rate* pengecekan situs per jam sebesar 4700%.

5.2 Saran

Berdasarkan hasil penelitian ini, berikut ini adalah saran untuk pengembangan penelitian berikutnya:

1. Penelitian ini menggunakan *dataset* yang sangat sedikit dan memiliki jumlah kategori yang tidak seimbang (*imbalance data*). Untuk mendapatkan performa yang lebih baik dan stabil, sebaiknya jumlah *dataset* dapat ditingkatkan terutama untuk data berlabel 1.
2. Situs *merchant* merupakan suatu *test data* yang *time-sensitive*. Sebagai contoh, situs X bisa saja pada suatu waktu pengecekan diklasifikasikan sebagai suatu situs yang tidak *fraud*, namun konten situs dapat berubah sewaktu-waktu. Sebuah *scheduler* atau mekanisme pengecekan rutin dapat dikembangkan sebagai solusi untuk selalu mendapatkan prediksi terkini.
3. Selain menggunakan pendekatan ekstraksi fitur-fitur yang statis dan *pre-determined*, metode *content analysis* juga dapat diterapkan. Word vector yang terbentuk dari transformasi *content* yang ada pada suatu halaman situs dapat dimanfaatkan sebagai fitur untuk melakukan prediksi.

DAFTAR REFERENSI

- [1] A. F. Indonesia, “Fintech landscape di indonesia,” Aug 2019.
- [2] M. D. Innovation, “Mobile payments in indonesia,” 2019.
- [3] P. Schueffel, “Taming the beast: A scientific definition of fintech,” *Journal of Information Management*, vol. 4, no. 4, pp. 32–54, 2016.
- [4] “What is fintech?,” Feb 2017.
- [5] “Ini dia empat jenis fintech di indonesia,” Jan 2018.
- [6] “Some insights to kyc, kyb, and what they mean to businesses,” Dec 2019.
- [7] J. Masanès, *Web Archiving: Issues and Methods*, pp. 1–53. 02 2007.
- [8] “What is a crawl frontier?,” 2018.
- [9] “What is selenium webdriver? difference with rc,” Feb 2020.
- [10] R. Fielding, J. Gettys, J. Mogul, H. Nielsen, and T. Berners-Lee, “Hypertext transfer protocol – http1.1,” Jan 1997.
- [11] G. Boeing and P. Waddell, “New insights into rental housing markets across the united states: Web scraping and analyzing craigslist rental listings,” *Journal of Planning Education and Research*, 07 2016.
- [12] F. Blog, “Web crawler 101: What is a web crawler and how do crawlers work?,” 2019.
- [13] “Machine learning faq,” Mar 2020.
- [14] L. Breiman, *Classification and regression trees*. Wadsworth International Group, 1984.
- [15] G. Louppe, *Understanding Random Forests: From Theory to Practice*. PhD thesis, University of Liege, Belgium, 10 2014. arXiv:1407.7502.
- [16] B. Efron, “Bootstrap methods: Another look at the jackknife,” *The Annals of Statistics*, vol. 7, no. 1, p. 1–26, 1979.
- [17] J. Brownlee, “A gentle introduction to xgboost for applied machine learning,” Apr 2020.
- [18] D. Opitz and R. Maclin, “Popular ensemble methods: An empirical study,” *Journal of Artificial Intelligence Research*, vol. 11, p. 169–198, 1999.

- [19] V. Metsis and et al., “Spam filtering with naive bayes – which naive bayes?,” in *THIRD CONFERENCE ON EMAIL AND ANTI-SPAM (CEAS)*, 2006.
- [20] C. D. Manning, P. Raghavan, and H. Schütze, *Introduction to Information Retrieval*. USA: Cambridge University Press, 2008.
- [21] G. C. Cawley and N. L. C. Talbot, “On over-fitting in model selection and subsequent selection bias in performance evaluation,” *Journal of Machine Learning Research*, Oct 2010.
- [22] R. T. T. Hastie, J. Friedman, “The elements of statistical learning: Data mining, inference, and prediction,” p. 193–224, 2001.
- [23] T. Fawcett, “An introduction to roc analysis,” *Pattern Recognition Letters*, vol. 27, no. 8, p. 861–874, 2006.
- [24] “pickle — python object serialization.”
- [25] J. H. Friedman, “Stochastic gradient boosting,” *Computational Statistics amp; Data Analysis*, vol. 38, no. 4, p. 367–378, 2002.

LAMPIRAN

LAMPIRAN 1: JUDUL LAMPIRAN 1

Subbab dari Lampiran 1

Isi subbab dari lampiran 1.