

SKRIPSI

**SISTEM PENILAIAN KEPUASAN PELANGGAN
TERHADAP TRANSAKSI PEMBELIAN
PRODUK E-COMMERCE
MENGUNAKAN ANALISIS SENTIMEN**

EUGENIUS WAHYUDIARTO

D121171526



**DEPARTEMEN TEKNIK INFORMATIKA
FAKULTAS TEKNIK
UNIVERSITAS HASANUDDIN
MAKASSAR
2022**

SKRIPSI

**SISTEM PENILAIAN KEPUASAN PELANGGAN
TERHADAP TRANSAKSI PEMBELIAN
PRODUK E-COMMERCE
MENGUNAKAN ANALISIS SENTIMEN**

TUGAS AKHIR

Sebagai salah satu syarat Mencapai Gelar Sarjana Departemen Teknik
Informatika Fakultas Teknik Universitas Hasanuddin Makassar

Disusun dan diajukan oleh:

EUGENIUS WAHYUDIARTO

D121171526

kepada

**DEPARTEMEN TEKNIK INFORMATIKA
FAKULTAS TEKNIK
UNIVERSITAS HASANUDDIN
MAKASSAR
2022**

LEMBAR PENGESAHAN SKRIPSI
SISTEM PENILAIAN KEPUASAN PELANGGAN TERHADAP
TRANSAKSI PEMBELIAN PRODUK E-COMMERCE MENGGUNAKAN
ANALISIS SENTIMEN


Disusun dan diajukan oleh
EUGENIUS WAHYUDIARTO
D121171526

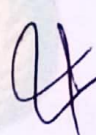
Telah dipertahankan di hadapan Panitia Ujian yang dibentuk dalam rangka Penyelesaian Studi Program Sarjana Program Studi Teknik Informatika Fakultas Teknik Universitas Hasanuddin pada tanggal 20 April 2022 dan dinyatakan telah memenuhi syarat kelulusan.

Menyetujui,

Pembimbing Utama,

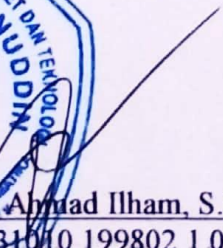
Pembimbing Pendamping,


Dr. Amil Ahmad Ilham, S.T., M.IT.
Nip. 197310101998021001


Anugrayani Bustamin, S.T., M.T.
Nip. 199012012018074001



Kepala Program Studi,


Dr. Amil Ahmad Ilham, S.T., M.IT.
Nip. 19731010 199802 1 002

PERNYATAAN KEASLIAN

Yang bertanda tangan di bawah ini:

Nama: Eugenius Wahyudiarto

NIM: D121171526

Departemen: Teknik Informatika

Jenjang: S1

Menyatakan dengan ini bahwa karya tulisan saya yang berjudul

**SISTEM PENILAIAN KEPUASAN PELANGGAN
TERHADAP TRANSAKSI PEMBELIAN
PRODUK E-COMMERCE
MENGUNAKAN ANALISIS SENTIMEN**

Adalah karya tulisan saya sendiri dan bukan merupakan pengambil alihan tulisan orang lain bahwa skripsi yang saya tulis ini benar-benar merupakan hasil karya saya sendiri.

Apabila di kemudian hari terbukti atau dapat dibuktikan bahwa sebagian atau keseluruhan skripsi ini hasil karya orang lain, maka saya bersedia menerima sanksi atas perbuatan tersebut.

Makassar, 23 Mei 2022

Yang Menyatakan



(Eugenius Wahyudiarto)

KATA PENGANTAR

Puji dan syukur penulis panjatkan kepada Tuhan Yang Maha Esa karena berkat rahmat dan karunia-Nya sehingga skripsi yang berjudul **“Sistem Penilaian Kepuasan Pelanggan Terhadap Transaksi Pembelian Produk E-commerce Menggunakan Analisis Sentimen”** ini dapat terselesaikan dengan baik sebagai salah satu syarat dalam menyelesaikan jenjang Strata-1 pada Departemen Teknik Informatika Fakultas Teknik Universitas Hasanuddin.

Penulis menyadari bahwa banyak kendala yang dihadapi dalam menyelesaikan skripsi ini. Berkat dukungan, bimbingan, serta motivasi yang diberikan, penyusunan skripsi ini dapat terselesaikan dengan baik. Ucapan terima kasih serta penghargaan yang setinggi-tingginya penulis sampaikan kepada:

1. Tuhan Yang Maha Esa yang melalui berkat dan rahmat-Nya sehingga penulis dapat menyelesaikan tugas akhir ini.
2. Kedua orang tua penulis, Bapak Yulius Tandiayu, dan Ibu Anastasia Sunarseh yang selalu memberikan dukungan, doa, semangat dan kasih sayang serta selalu sabar dalam mendidik penulis.
3. Harrison Kinsley & Daniel Kukiela yang telah membuat materi *Neural Networks from Scratch in Python* sehingga membantu penulis untuk lebih memahami secara mendalam secara teori dan praktek mengenai *Neural Network*. Khususnya, Harrison Kinsley atau Sentdex di Youtube

yang juga membimbing penulis untuk menulis program menggunakan Python.

4. Younes Bensouda Mourri sebagai AI instruktur di Stanford University dan Łukasz Kaiser sebagai Staff Research Scientist, Google Brain & Chargé de Recherche, CNRS, yang telah membimbing penulis dalam memahami NLP melalui materinya di Natural Language Processing Specialization.
5. Bapak Dr. Amil Ahmad Ilham, S.T., M.IT., selaku pembimbing I dan ketua Prodi Teknik Informatika Fakultas Teknik Universitas Hasanuddin, serta Ibu Anugrahyani Bustamin, S.T., M.T., selaku pembimbing II yang selalu menyediakan waktu, tenaga, pikiran dan memberikan bimbingan dalam penyusunan skripsi ini.
6. Bapak Ir. Christiforus Yohannes, M.T., selaku penguji I, serta Bapak Dr. Muhammad Niswar, S.T., M.T., selaku penguji II yang memberikan saran dan pelajaran yang sangat berharga dalam hidup penulis.
7. Fitriani Nasir, S.T., yang memberikan arahan mendetail dan praktis bagi penulis dalam menyelesaikan tugas akhir ini.
8. Ayu Lestari, S.Si., yang turut membantu penulis mengoreksi skripsi ini.
9. Teman-teman Prodi Teknik Informatika angkatan 2017, 2021 atas dukungan dan semangat yang telah diberikan.
10. Segenap Dosen dan Staff Departemen Teknik Informatika Fakultas Teknik Universitas Hasanuddin yang telah banyak membantu penulis selama masa perkuliahan.

11. Serta berbagai pihak atas segala dukungan dan bantuannya yang tidak dapat penulis tuliskan satu persatu.

Penulis berharap semoga Tuhan Yang Maha Esa membalas segala kebaikan dari semua pihak yang telah membantu penulis dalam penyusunan skripsi ini dan semoga skripsi ini dapat memberikan hal yang bermanfaat serta menambah wawasan ilmu untuk pembaca dan juga bagi penulis sendiri.

Akhir kata, penulis mengucapkan selamat membaca dan mendalami skripsi ini. Penulis akan membawa pembaca dalam perjalanan untuk mencapai tujuan skripsi ini. Selanjutnya untuk menghubungi penulis, pembaca dapat mengirimkan pesan melalui eugeniusw18@gmail.com. Berbahagialah orang yang haus akan pengetahuan, karena dia akan tahu.

Makassar, Februari 2022

Penulis,

(Eugenius Wahyudiarto)

ABSTRAK

Saat ini semakin banyak produk bermunculan. Berbagai layanan yang menawarkan produk yang mirip dan sama membuat pembeli sulit memutuskan untuk membeli sebelum melihat ulasan dari pengguna lainnya. Hal ini diperparah dengan tumbuhnya platform pemasaran (*e-commerce*) yang berbeda-beda. Pengguna pun harus menghabiskan waktu untuk memilih produk di masing-masing platform dengan banyak pertimbangan, seperti melihat *rating*, harga, dan ulasan dari pembeli lainnya. Optimalisasi pemilihan produk *e-commerce* dilakukan supaya pengguna tidak harus menghabiskan waktu yang lama dengan membaca setiap ulasan ketika ingin membeli produk. Dengan adanya sistem penilaian kepuasan pelanggan terhadap transaksi pembelian produk *e-commerce* menggunakan analisis sentimen diharapkan dapat memberikan penilaian menyeluruh terhadap transaksi pembelian suatu produk dari ulasan yang diberikan. Sistem ini memberikan 91% akurasi dan skor F_1 sebesar 57% dalam memprediksi sentimen menggunakan naive Bayes. Model naive Bayes kemudian digunakan untuk mendapatkan skor kepuasan pelanggan terhadap suatu produk.

Kata kunci: skor kepuasan pelanggan, rekomendasi pembelian produk, analisis sentimen, naive Bayes.

DAFTAR ISI

HALAMAN JUDUL.....	i
HALAMAN PENGANTAR.....	ii
HALAMAN PENGESAHAN.....	iii
PERNYATAAN KEASLIAN.....	iv
KATA PENGANTAR	v
ABSTRAK	viii
DAFTAR ISI.....	ix
DAFTAR TABEL.....	xi
DAFTAR GAMBAR	xii
DAFTAR LAMPIRAN.....	xiv
BAB I PENDAHULUAN.....	1
1.1 Latar belakang	1
1.2 Rumusan masalah.....	2
1.3 Tujuan penelitian	3
1.4 Batasan penelitian.....	3
1.5 Manfaat penelitian	3
1.6 Sistematika penulisan	4
BAB II TINJAUAN PUSTAKA.....	6
2.1 Analisis sentimen	6
2.2 <i>Web scraper</i>	6
2.3 <i>Numeric character references</i>	6
2.4 Google Cloud Natural Language API	7
2.5 Algoritma Naive Bayes	9
2.6 Pelatihan Classifier Naive Bayes	11
2.7 Evaluasi Model Klasifikasi	13
2.8 <i>Neural Network</i>	16
2.9 <i>Continuous Bag-of-Words (CBOW)</i>	21
2.10 Stratified K-Fold Cross-Validation	23

2.11	Konversi nilai ke rentang baru	24
BAB III METODE PENELITIAN.....		25
3.1	Tempat dan waktu	25
3.2	Alat dan bahan.....	25
3.3	Pelaksanaan Penelitian	26
3.4	Metode Penelitian.....	27
3.5	Pengambilan data	29
3.6	<i>Preprocessing</i> I: Persiapan Pelabelan	29
3.7	Pelabelan Data	34
3.8	<i>Preprocessing</i> II: Pembentukan Data Latih, Validasi, dan Uji	35
3.9	Pelatihan, dan Validasi Model.....	36
3.10	Pengujian Model.....	40
3.11	Meningkatkan Kualitas Model Sentimen	41
3.12	Perhitungan Skor Kepuasan Pelanggan.....	47
BAB IV HASIL DAN PEMBAHASAN		49
4.1	Hasil Pengambilan Data	49
4.2	Hasil <i>Preprocessing</i> I: Persiapan Pelabelan.....	51
4.3	Hasil Pelabelan Data	51
4.4	Hasil <i>Preprocessing</i> II: Pembentukan Data Latih, Validasi, dan Uji.....	52
4.5	Hasil Evaluasi Model Sentimen	52
4.6	Contoh Penilaian Kepuasan Pelanggan	57
4.7	<i>Website</i> Penilaian Kepuasan Pelanggan	59
BAB V KESIMPULAN DAN SARAN.....		62
5.1	Kesimpulan.....	62
5.2	Saran	62
DAFTAR PUSTAKA		64

DAFTAR TABEL

Tabel 3.6.1 Pemetaan transformasi simbol	31
Tabel 4.1.1 Jumlah baris setiap tabel	50
Tabel 4.1.2 12 kategori teratas dalam jumlah baris terbanyak	50
Tabel 4.4.1 Komposisi Data	52
Tabel 4.5.1 Perbandingan performa antar model sentimen	57

DAFTAR GAMBAR

Gambar 2.4.1 Contoh interpretasi sentimen	8
Gambar 2.6.1 Algoritma naive Bayes untuk teks	13
Gambar 2.7.1 Confusion matrix pada klasifikasi biner.	14
Gambar 2.7.2 Confusion matrix pada klasifikasi tiga kelas.	15
Gambar 2.8.1 Neural network 3 layers	16
Gambar 2.8.2 Neuron pada neural network.....	16
Gambar 2.8.3 Fungsi aktivasi sigmoid	17
Gambar 2.8.4 Fungsi aktivasi ReLu	18
Gambar 2.8.5 Fungsi aktivasi tanh	18
Gambar 2.8.6 Proses <i>forward propagation</i>	19
Gambar 2.8.7 Proses <i>backpropagation</i>	21
Gambar 2.9.1 Arsitektur CBOW	22
Gambar 2.10.1 Ilustrasi Stratified K-Fold Cross Validation	24
Gambar 2.11.1 Fungsi untuk konversi nilai ke rentang lain.....	24
Gambar 3.3.1 Tahapan Penelitian.....	26
Gambar 3.4.1 Rancangan sistem	27
Gambar 3.4.2 Struktur basis data scraper	28
Gambar 3.6.1 Contoh transformasi simbol ke spasi.....	30
Gambar 3.6.2 Contoh pengurangan emoji.....	31
Gambar 3.6.3 Pengurangan karakter pada kata	31
Gambar 3.6.4 Contoh transformasi simbol.....	32
Gambar 3.6.5 Diagram alur transformasi angka.....	33
Gambar 3.6.6 Contoh transformasi angka	33
Gambar 3.6.7 Contoh pemisahan kalimat/frasa	34
Gambar 3.8.1 Pembagian data latih, validasi, dan uji	35
Gambar 3.9.1 Ekstraksi data dari reviews	36

Gambar 3.11.1 Update algoritma <i>train</i> pada model sentimen.....	41
Gambar 4.1.1 Tabel yang berhasil terbentuk.....	49
Gambar 4.5.1 <i>Confusion matrix</i> model 1 dengan <i>stopwords</i>	53
Gambar 4.5.2 Implementasi perhitungan performa model 1	53
Gambar 4.5.3 <i>Confusion matrix</i> model 1	54
Gambar 4.5.4 <i>Confusion matrix</i> model 2	54
Gambar 4.5.5 Kesalahan prediksi model 2 pada negasi	55
Gambar 4.5.6 <i>Confusion matrix</i> model model 3	55
Gambar 4.5.7 Kalimat negasi berhasil diprediksi oleh model 3.....	56
Gambar 4.5.8 <i>Confusion matrix</i> model GCloud NL	56
Gambar 4.6.1 Skor kepuasan pelanggan dari ulasan produk 1.....	58
Gambar 4.6.2 Skor kepuasan pelanggan dari ulasan produk 2.....	58
Gambar 4.6.3 Rating kedua produk. Produk 1 (atas), produk 2 (bawah).....	59
Gambar 4.7.1 Kode QR dan <i>url website</i> penilaian kepuasan pelanggan.....	60
Gambar 4.7.2 Tampilan <i>website</i> pada komputer	60
Gambar 4.7.3 Tampilan penilaian kepuasan pelanggan di layar komputer.....	61
Gambar 4.7.4 Interpretasi penilaian kepuasan pelanggan	61

DAFTAR LAMPIRAN

Lampiran 1 Diagram alur Tokopedia Scraper	68
Lampiran 2 Diagram alur pelabelan sentimen	69
Lampiran 3 Diagram alur perhitungan skor kepuasan pelanggan (bagian 1).....	70
Lampiran 4 Diagram alur perhitungan skor kepuasan pelanggan (bagian 2).....	71
Lampiran 5 Hasil <i>Preprocessing</i> I	72
Lampiran 6 Hasil prediksi sentimen menggunakan GCloud NL	73
Lampiran 7 Beberapa tes yang salah diprediksi	74
Lampiran 8 <i>Stopwords</i>	75
Lampiran 9 Data Setelah <i>Preprocessing</i> II	80
Lampiran 10 Hasil pelabelan data secara manual	81
Lampiran 11 Performa Model 1 (dengan <i>stopwords</i>) pada SKF terhadap label manual	82
Lampiran 12 Performa Model 1 (tanpa <i>stopwords</i>) pada SKF terhadap label manual	83
Lampiran 13 Performa Model 2 pada SKF terhadap label manual	84
Lampiran 14 Performa Model 3 pada SKF terhadap label manual	85
Lampiran 15 Performa GCloud NL pada SKF terhadap label manual	86

BAB I

PENDAHULUAN

1.1 Latar belakang

Saat ini semakin banyak produk bermunculan. Bahkan di antara produk tersebut terdapat banyak kemiripan dan kesamaan di dalamnya sehingga pembeli sulit memutuskan di mana pembelian produk yang paling cocok untuk mereka. Hal ini diperparah dengan tumbuhnya platform pemasaran (*e-commerce*) yang berbeda-beda. Pengguna pun harus menghabiskan waktu untuk memilih produk di masing-masing platform dengan banyak pertimbangan, seperti melihat *rating*, harga, dan ulasan dari pembeli lainnya.

Ini adalah masalah rekomendasi. Terdapat dua kategori teknik rekomendasi, *content based filtering* (CBF) dan *collaborative filtering* (CF) (Shahbazi & Byun, 2020). Ada juga yang menggabungkan keduanya sehingga menjadi *hybrid*. CBF akan merekomendasikan produk berdasarkan kemiripannya dengan yang telah dibeli sebelumnya dan CF akan merekomendasikan produk berdasarkan yang dibeli entitas lain. Ada beberapa solusi yang dapat dilakukan untuk mengatasi masalah di atas. Pertama, fitur pencarian produk yang lebih baik. Fitur pencarian ini akan masuk dalam model klasifikasi tag dan melakukan *query* ke database. Kedua, merekomendasikan produk yang paling baik. Untuk mendapatkan produk yang paling baik perlu dilakukan perbandingan antara produk dan juga *review* dari setiap produk. *Sentiment analysis* (SA) akan digunakan

untuk mengetahui sikap pengguna terhadap produk yang ditawarkan (Khanvilkar & Deepali Vora, 2018). Terdapat beberapa algoritma yang dapat digunakan dalam SA seperti XGBoost, Random Forest, SVM, KNearest Logistic Regression (Shahbazi & Byun, 2020), dan naive Bayes (Parveen & Pandey, 2016). Bahkan ada yang menggunakan *lexicon-based keyword* yang digabungkan dengan algoritma tersebut (Khanvilkar & Deepali Vora, 2018). Tantangannya adalah *preprocessing* datanya seperti, kesalahan ketik dari pengguna, singkatan Bahasa Indonesia dan *emoticon*. Karena *emoticon* memainkan peranan penting dalam SA (Parveen & Pandey, 2016).

Sistem yang akan dibuat seperti yang dilakukan oleh Kim (Kim dkk., 2019). Perbedaannya adalah sistem yang dibuat oleh Kim hanya menggunakan *keyword* kata positif dan negatif untuk menentukan sentimennya. Sistem yang dibuat oleh Hanni dapat membandingkan spesifikasi produk (Hanni dkk., 2016). Sayangnya tidak semua produk memiliki spesifikasi. Alternatifnya adalah membuat skor penilaian kepuasan pelanggan, sebagai rekomendasi, berdasarkan hasil SA sehingga mendapatkan ringkasan skor mengenai apa yang dirasakan pelanggan ketika melakukan transaksi pembelian produk terkait.

1.2 Rumusan masalah

Adapun rumusan masalah dalam penelitian ini yaitu, bagaimana mengembangkan model analisis sentimen dalam penilaian kepuasan

pelanggan terhadap transaksi produk e-commerce dengan menggunakan algoritma naive Bayes?

1.3 Tujuan penelitian

Tujuan dari penelitian ini adalah untuk mengetahui bagaimana mengembangkan model analisis sentimen menggunakan algoritma naive Bayes dalam penilaian kepuasan pelanggan terhadap transaksi pembelian produk e-commerce.

1.4 Batasan penelitian

1. Data ulasan produk diambil dari Tokopedia sebanyak minimal 3 kategori dan jumlah data minimal 50000 per kategori.
2. Model analisis sentimen menggunakan naive Bayes yang dibuat dengan pemrograman Python versi 3.7.
3. Sistem penilaian ini berdasarkan ulasan dari pelanggan.

1.5 Manfaat penelitian

1. Penelitian dapat memberikan sumbangan pemikiran mengenai penilaian kepuasan pelanggan terhadap transaksi pembelian produk e-commerce.
2. Penelitian dapat mempermudah masyarakat dalam pembelian produk tanpa harus menghabiskan waktu untuk membaca ulasan setiap produk.
3. Penelitian dapat merepresentasikan skor kepuasan pelanggan dengan lebih akurat dibandingkan dengan menggunakan bintang dari pelanggan yang kadangkala tanpa komentar apa-apa.

4. Penelitian ini dapat diterapkan dalam kasus lain yang sejenis, misalnya penilaian video berdasarkan komentar atau penilaian aplikasi berdasarkan ulasan.

1.6 Sistematika penulisan

Sebagai gambaran singkat dari isi tulisan ini, penulis memberikan uraian singkat pada masing-masing tahapan penelitian.

BAB I PENDAHULUAN

Pada tahap ini dijabarkan latar belakang, rumusan masalah, tujuan, batasan, manfaat, dan sistematika penulisan dari penelitian ini.

BAB II TINJAUAN PUSTAKA

Pada tahap ini dituliskan referensi-referensi istilah dan metode yang terkait dengan sistem penilaian kepuasan pelanggan terhadap transaksi produk e-commerce menggunakan analisis sentimen.

BAB III METODOLOGI PENELITIAN

Pada tahap ini dijelaskan secara detail penelitian yang dilakukan. Beberapa poin pembahasannya yaitu, tempat dan waktu, alat dan bahan, metode penelitian, rancangan sistem, serta analisis sistem.

BAB IV HASIL DAN PEMBAHASAN

Pada tahap ini dijelaskan tentang hasil dan pembahasan dari penelitian yang dilakukan.

BAB V PENUTUP

Pada tahap ini diberikan kesimpulan yang didapatkan dari penelitian yang telah dilakukan. Terdapat juga saran-saran untuk pengembangan selanjutnya.

BAB II

TINJAUAN PUSTAKA

2.1 Analisis sentimen

Sentimen adalah pandangan atau perasaan terhadap suatu peristiwa atau sesuatu. Analisis adalah pemeriksaan secara teliti terhadap suatu hal. Analisis sentimen atau *sentiment analysis* bertujuan untuk secara otomatis mengungkap sikap mendasar yang kita pegang terhadap suatu entitas. Analisis sentimen memiliki banyak aplikasi salah satunya adalah mengukur tingkat kepuasan pelanggan terhadap suatu merek terhadap merek lainnya. Saat ini analisis sentimen yang secara luas digunakan berbasis teks (Soleymani dkk., 2017, hlm. 2).

2.2 Web scraper

Web scraper adalah suatu perangkat lunak yang menyimulasikan kunjungan manusia pada web untuk mengumpulkan data informasi terperinci. Keuntungan menggunakan *scraper* yakni dapat diprogram dan diotomatisasikan sehingga pengumpulan informasi menjadi lebih cepat dan informasi dapat disimpan pada format yang terstruktur (Diouf dkk., 2019, hlm. 6040).

2.3 Numeric character references

Numeric character references menentukan posisi kode karakter dalam kumpulan karakter dokumen. Referensi karakter numerik dapat mengambil dua bentuk:

1. Sintaks "**&#D;**". D adalah angka desimal, mengacu pada angka karakter desimal ISO 10646.
2. Sintaks "**&#xH;**" atau "**&#XH;**". H adalah angka heksadesimal, mengacu pada nomor karakter heksadesimal ISO 10646. Angka heksadesimal dalam referensi karakter numerik tidak peka huruf besar-kecil.

Berikut ini beberapa contoh representasi *numeric character references*:

1. **** (dalam desimal) merepresentasikan karakter apostrof.
2. **"**, **"**, **"**, atau **"** (dalam heksadesimal) merepresentasikan karakter apostrof (World Wide Web Consortium / W3C, 2018, Bab 5).

2.4 Google Cloud Natural Language API

Google Cloud Natural Language (GCloud NL) API adalah antarmuka pemrograman aplikasi dari Google yang dikembangkan dengan menggunakan model yang telah dilatih sebelumnya. Model ini memberdayakan pengembang untuk dengan mudah menerapkan Natural Language Understanding (NLU) ke aplikasi mereka dengan fitur-fitur termasuk analisis sentimen, analisis entitas, analisis sentimen entitas, klasifikasi konten, dan analisis sintaksis (Google LLC, 2021a).

Berikut ini contoh fungsi untuk menganalisa sentimen suatu teks menggunakan Google Cloud dalam bahasa pemrograman Python dari <https://cloud.google.com/natural-language/docs/analyzing-sentiment>.

```

from google.cloud import language v1

def sample_analyze_sentiment(text_content):
    """
    Analyzing Sentiment in a String

    Args:
        text_content The text content to analyze
    """

    client = language_v1.LanguageServiceClient()

    # text_content = 'I am so happy and joyful.'

    # Available types: PLAIN TEXT, HTML
    type_ = language_v1.Document.Type.PLAIN_TEXT

    # Optional. If not specified, the language is automatically detected.
    # For list of supported languages:
    # https://cloud.google.com/natural-language/docs/languages
    language = "en"
    document = {"content": text_content, "type_": type_, "language": language}

    # Available values: NONE, UTF8, UTF16, UTF32
    encoding_type = language_v1.EncodingType.UTF8

    response = client.analyze_sentiment(request = {'document': document,
        'encoding type': encoding_type})
    # Get overall sentiment of the input document
    print(u"Document sentiment score: {}".format(response.document_sentiment.score))
    print(
        u"Document sentiment magnitude: {}".format(
            response.document_sentiment.magnitude
        )
    )
    # Get sentiment for all sentences in the document
    for sentence in response.sentences:
        print(u"Sentence text: {}".format(sentence.text.content))
        print(u"Sentence sentiment score: {}".format(sentence.sentiment.score))
        print(u"Sentence sentiment magnitude: {}".format(sentence.sentiment.magnitude))

    # Get the language of the text, which will be the same as
    # the language specified in the request or, if not specified,
    # the automatically-detected language.
    print(u"Language of the text: {}".format(response.language))

```

Sentimen dari suatu teks dapat diinterpretasikan dengan melihat nilai *score* dan *magnitude*-nya. Google Cloud memberikan contoh pada **Gambar 2.4.1** untuk untuk menginterpretasikannya.

Sentiment	Sample Values
Clearly Positive*	"score": 0.8, "magnitude": 3.0
Clearly Negative*	"score": -0.6, "magnitude": 4.0
Neutral	"score": 0.1, "magnitude": 0.0
Mixed	"score": 0.0, "magnitude": 4.0

Gambar 2.4.1 Contoh interpretasi sentimen
 Sumber: (Google LLC, 2021b, bag. Natural Language API Basics)

Nilai *score* adalah skor sentimen mulai dari -1.0 (negatif) sampai 1.0 (positif). Nilai *magnitude* menandakan kekuatan emosional (baik itu positif maupun negatif) dalam teks mulai dari 0.0 sampai +inf. Seandainya nilai *score* adalah 0 bukan berarti itu adalah netral. Bisa saja tingkat emosional negatif dan positif adalah sama sehingga harus dilakukan pengecekan nilai *magnitude*. Penentuan *clearly positive* dan *negative* bervariasi tergantung kasus penggunaan. Misalnya saja dapat didefinisikan *threshold* skor di atas 0.25 sebagai *clearly positive*. Namun setelah dilihat hasilnya skor 0.15 masih positif maka *threshold* dapat diubah (Google LLC, 2021b).

2.5 Algoritma Naive Bayes

Algoritma naive Bayes adalah algoritma probabilitas berdasarkan Teorema Bayes untuk menyelesaikan masalah klasifikasi (Kathuria, 2020; Zhang, 2019). Teorema Bayes adalah cara untuk mencari suatu probabilitas ketika kita mengetahui probabilitas lainnya. Persamaannya adalah pada persamaan (1) sebagai berikut (Pierce, 2020; Zhang, 2019).

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)} \quad (1)$$

P(Y|X) / *posterior* : seberapa sering Y terjadi jika X terjadi;

P(X|Y) / *likelihood* : seberapa sering X terjadi jika Y terjadi;

P(Y) / *prior* : seberapa besar kemungkinan Y terjadi;

P(X) / *evidence* : seberapa besar kemungkinan X terjadi;

Dalam skenario nyata, X atau fitur tidak hanya satu dan selalu ada ketergantungan di antaranya. Misalkan terdapat 5 fitur sehingga secara naif persamaan (1) dimodifikasi menjadi persamaan naive Bayes di persamaan (2) berikut yang mengasumsikan semua fitur independen.

$$P(y|x_1, x_2, x_3, x_4, x_5) = \frac{P(x_1|y)P(x_2|y)P(x_3|y)P(x_4|y)P(x_5|y) * P(y)}{P(x_1)P(x_2)P(x_3)P(x_4)P(x_5)} \quad (2)$$

Persamaan (2) juga dapat dituliskan sebagai persamaan (3) berikut.

$$P(y|x_1, x_2, x_3, x_4, x_5) = \frac{P(y)\prod_{i=1}^5 P(x_i|y)}{P(x_1)P(x_2)P(x_3)P(x_4)P(x_5)} \quad (3)$$

Perhatikan bahwa denominator pada persamaan (3) akan konstan. Sehingga secara proporsi persamaan (3) akan menjadi persamaan (4) berikut.

$$P(y|x_1, x_2, x_3, x_4, x_5) \propto P(y)\prod_{i=1}^5 P(x_i|y) \quad (4)$$

Ini akan memberikan kemungkinan untuk setiap kelas y. Yang lebih tinggi dianggap sebagai prediksinya. Sehingga untuk mendapatkan nilai tersebut digunakan argmax seperti pada persamaan (5) (Kathuria, 2020).

$$\hat{y} = \arg \max_{y \in Y} P(y)\prod_{i=1}^5 P(x_i|y) \quad (5)$$

Untuk mengaplikasikan naive Bayes pada klasifikasi teks, X di sini menjadi kata dalam dokumen. Perhitungan naive Bayes dilakukan di ruang log untuk menghindari *underflow* (kalkulasi yang terlalu kecil untuk direpresentasikan oleh CPU atau memori komputer) dan meningkatkan

kecepatan komputasi. Sehingga persamaan (5) dapat direpresentasikan sebagai persamaan (6) berikut.

$$\hat{y} = \arg \max_{y \in Y} \log P(y) + \sum_{i=1}^5 \log P(x_i|y) \quad (6)$$

Dengan menggunakan ruang log, persamaan (6) akan memprediksi kelas sebagai fungsi linear terhadap fitur yang diberikan (Jurafsky & Martin, 2020).

2.6 Pelatihan Classifier Naive Bayes

Pelatihan model klasifikasi sentimen dengan menggunakan persamaan (6), misalkan $y=c$ dan $x_i=w_i$, dilakukan dengan mencari nilai $P(c)$ dan $P(w_i|c)$. *Prior* atau $P(c)$ adalah berapa persen setiap kelas sentimen c dalam data latih. Untuk menentukannya digunakan persamaan (7) berikut.

$$\hat{P}(c) = \frac{N_c}{N_{doc}} \quad (7)$$

N_c : banyaknya data latih dengan kelas sentimen c ;

N_{doc} : banyaknya data latih

Likelihood atau $P(w_i|c)$ diasumsikan sebagai rasio antara jumlah kemunculan suatu kata dalam *bag of words* dokumen terhadap jumlah kata dengan kelas c , yang dihitung dengan persamaan (8) berikut.

$$\hat{P}(w_i|c) = \frac{\text{count}(w_i, c)}{\sum_{w' \in V} \text{count}(w', c)} \quad (8)$$

V : terdiri dari semua kata dari semua kelas c ;

$count(w_i, c)$: jumlah kata dengan kelas c , contoh: “fantastic”=1;

$count(w', c)$: jumlah semua kata dalam kelas c , contoh: 100 kata positif dari 1000 kata (1/10);

Tetapi apabila suatu kata tidak muncul dalam data latih, misalnya kata “fantastic” maka nominator akan bernilai 0 sehingga seluruh persamaan bernilai 0.

$$\hat{P}(\text{"fantastic"} | \text{positive}) = \frac{count(\text{"fantastic"}, \text{positive})}{\sum_{w' \in V} count(w', \text{positive})} = 0 \quad (9)$$

Hal ini menjadi masalah karena naive Bayes akan mengalikan semua fitur *likelihoods* bersama sehingga satu saja 0 akan membuat nilai *likelihoods* dari fitur lain menjadi 0. Ini dapat diatasi dengan menambahkan 1 (*Laplace smoothing*), seperti pada persamaan (10) berikut.

$$\hat{P}(w_i | c) = \frac{count(w_i, c) + 1}{\sum_{w' \in V} (count(w', c) + 1)} = \frac{count(w_i, c) + 1}{(\sum_{w' \in V} count(w', c)) + |V|} \quad (10)$$

Untuk menyelesaikan masalah negasi pada sentimen, dapat ditambahkan kata “NOT_” di depan setiap kata setelah kata negasi. Penggunaan *stopwords* tidak meningkatkan performa model, sehingga algoritma naive Bayes pada **Gambar 2.6.1** berikut tidak mengikutsertakan *stopwords* (Jurafsky & Martin, 2020).

```

function TRAIN NAIVE BAYES(D, C) returns  $\log P(c)$  and  $\log P(w|c)$ 
for each class  $c \in C$            # Calculate  $P(c)$  terms
   $N_{doc}$  = number of documents in D
   $N_c$  = number of documents from D in class c
   $logprior[c] \leftarrow \log \frac{N_c}{N_{doc}}$ 
   $V \leftarrow$  vocabulary of D
   $bigdoc[c] \leftarrow$  append(d) for d  $\in$  D with class c
  for each word w in V           # Calculate  $P(w|c)$  terms
     $count(w,c) \leftarrow$  # of occurrences of w in  $bigdoc[c]$ 
     $loglikelihood[w,c] \leftarrow \log \frac{count(w,c) + 1}{\sum_{w' \in V} (count(w',c) + 1)}$ 
return  $logprior, loglikelihood, V$ 

function TEST NAIVE BAYES( $testdoc, logprior, loglikelihood, C, V$ ) returns best c
for each class  $c \in C$ 
   $sum[c] \leftarrow logprior[c]$ 
  for each position i in  $testdoc$ 
     $word \leftarrow testdoc[i]$ 
    if  $word \in V$ 
       $sum[c] \leftarrow sum[c] + loglikelihood[word,c]$ 
return  $argmax_c sum[c]$ 

```

Gambar 2.6.1 Algoritma naive Bayes untuk teks
 Sumber: (Jurafsky & Martin, 2020, hlm. 60)

2.7 Evaluasi Model Klasifikasi

Untuk mengevaluasi model klasifikasi diperlukan pembuatan *confusion matrix* dari model tersebut. *Confusion matrix* berfungsi untuk memvisualisasikan bagaimana performa sistem terhadap label yang diberikan (*gold standard label*). Dengan *confusion matrix* dapat dihitung nilai *precision*, *recall*, *accuracy* dan pengukuran *F*. Pada **Gambar 2.7.1** berikut format *confusion matrix* pada klasifikasi biner (Jurafsky & Martin, 2020).

		<i>gold standard labels</i>		
		gold positive	gold negative	
<i>system output labels</i>	system positive	true positive	false positive	precision = $\frac{tp}{tp+fp}$
	system negative	false negative	true negative	
		recall = $\frac{tp}{tp+fn}$		accuracy = $\frac{tp+tn}{tp+fp+tn+fn}$

Gambar 2.7.1 Confusion matrix pada klasifikasi biner.

Sumber: (Jurafsky & Martin, 2020, hlm. 65)

Recall adalah rasio antara semua kelas yang diklasifikasi benar pada kelas positif terhadap total anggota aktual dari kelas positif. Dengan kata lain ini memberi tahu tentang berapa banyak dari jumlah total contoh positif yang diklasifikasikan benar. *Precision* adalah rasio antara semua anggota yang diklasifikasikan benar dalam kelas positif terhadap jumlah total anggota yang diklasifikasikan dalam kelas positif dan benar (Villalobos, 2020). Bagaimana bila ingin mengevaluasi performa model multi kelas? Untuk model multi kelas dapat menggunakan *macroaveraging*, yaitu precision dan recall dihitung seperti pada **Gambar 2.7.2** dan hasil akhirnya dirata-ratakan.

		<i>gold labels</i>			
		urgent	normal	spam	
<i>system output</i>	urgent	8	10	1	$\text{precision}_u = \frac{8}{8+10+1}$
	normal	5	60	50	$\text{precision}_n = \frac{60}{5+60+50}$
	spam	3	30	200	$\text{precision}_s = \frac{200}{3+30+200}$
		$\text{recall}_u = \frac{8}{8+5+3}$	$\text{recall}_n = \frac{60}{10+60+30}$	$\text{recall}_s = \frac{200}{1+50+200}$	

Gambar 2.7.2 Confusion matrix pada klasifikasi tiga kelas.
 Sumber: (Jurafsky & Martin, 2020, hlm. 67)

Nilai *recall* dan *precision* digunakan apabila model menggunakan data yang tidak seimbang sehingga nilai *accuracy* menjadi tidak berguna. Terdapat juga persamaan F apabila ingin menggabungkan nilai *recall* dan *precision* dengan menggunakan parameter β sebagai nilai pementingnya pada persamaan (11). Apabila $\beta > 1$ maka lebih mementingkan *recall* dan bila $\beta < 1$ maka lebih mementingkan *precision*.

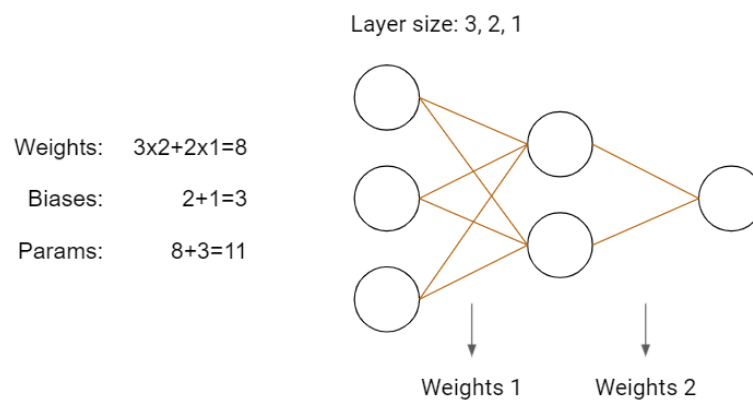
$$F_{\beta} = \frac{(\beta^2 + 1)PR}{\beta^2 P + R} \quad (11)$$

Apabila nilai $\beta = 1$ maka nilai *recall* dan *precision* dianggap seimbang dan disebut sebagai persamaan F_1 yang dirumuskan dalam persamaan (12) berikut (Jurafsky & Martin, 2020).

$$F_1 = \frac{2PR}{P + R} \quad (12)$$

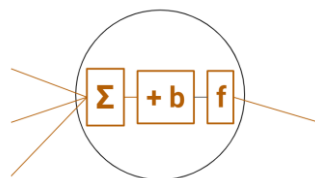
2.8 Neural Network

Artificial neural network (ANN) (yang sekarang orang-orang menyebutnya *neural network* tanpa kata “*artificial*”) terinspirasi dari otak biologis, diterjemahkan ke komputer. ANN terdiri dari syaraf-syaraf (*neurons*) yang saling terhubung satu sama lain. Hal inilah yang membuat ANN seringkali mengalahkan performa metode *machine learning* yang lain. Struktur *neural network* dapat dilihat pada **Gambar 2.8.1** berikut.



Gambar 2.8.1 Neural network 3 layers

Setiap *neuron* pada neural network memiliki struktur seperti pada **Gambar 2.8.2** (Kinsley & Kukiela, 2020).



Gambar 2.8.2 Neuron pada neural network

Neural network, dalam *Weights* 1 (lihat: **Gambar 2.8.1**) atau layer(3) ke layer(2), dapat memproses dalam satu hitungan dengan menggunakan matriks pada *input* X (matriks) dan menghitung setiap *output* Y (matriks) sesuai dengan persamaan **(13)** (Aflak, 2021).

$$Y = XW + B \quad (13)$$

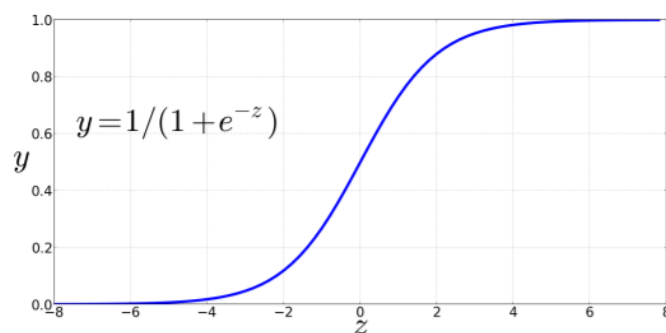
Y : Matriks prediksi;

X : Matriks input;

W : Matriks *weights*;

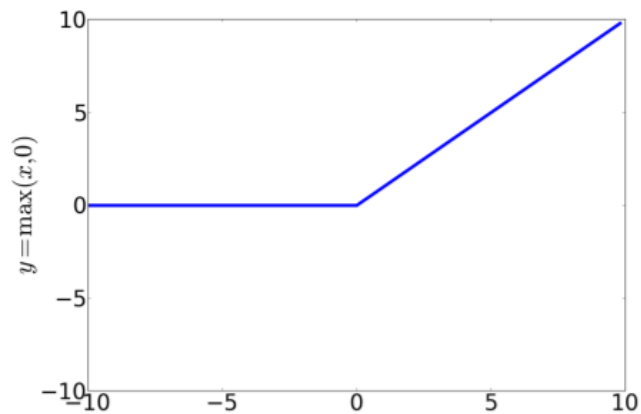
B : Matriks *bias*;

Untuk menambahkan non-linearitas pada model, perlu ditambahkan fungsi non-linear dari output beberapa *layers* (Aflak, 2021). Fungsi non-linear ini seperti sigmoid pada **Gambar 2.8.3**;



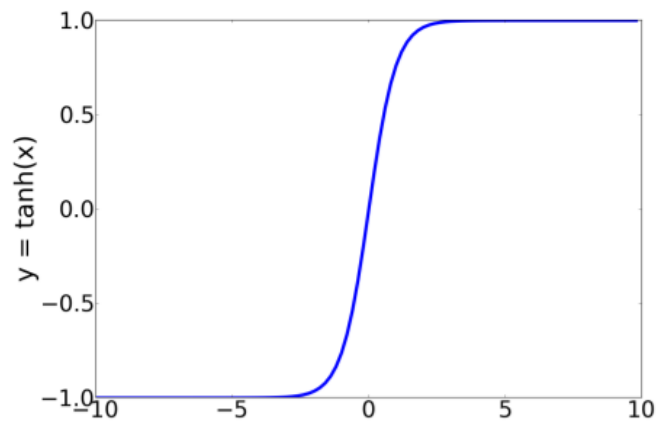
Gambar 2.8.3 Fungsi aktivasi sigmoid
Sumber: (Jurafsky & Martin, 2020, hlm. 128)

ReLU atau *rectified linear unit* pada **Gambar 2.8.4**;



Gambar 2.8.4 Fungsi aktivasi ReLu
 Sumber: (Jurafsky & Martin, 2020, hlm. 130)

dan tanh pada **Gambar 2.8.5**, sebagai variasi dari sigmoid dengan jangkauan output -1 sampai 1 (Jurafsky & Martin, 2020).

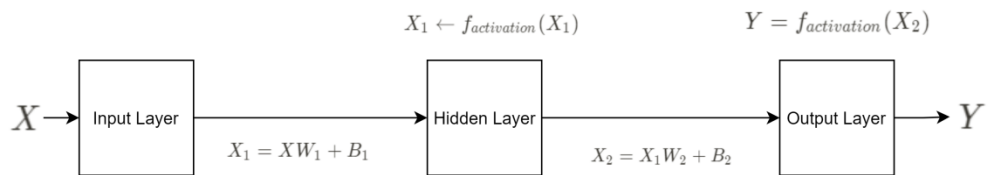


Gambar 2.8.5 Fungsi aktivasi tanh
 Sumber: (Jurafsky & Martin, 2020, hlm. 130)

Fungsi aktivasi di atas dapat dituliskan secara umum sebagai persamaan (14) berikut (Aflak, 2021).

$$Y = f_{activation}(X) \quad (14)$$

Dengan demikian, *forward propagation* dapat dilakukan seperti pada **Gambar 2.8.6** berikut (Aflak, 2021).



Gambar 2.8.6 Proses *forward propagation*

Agar model dapat meminimalkan error antara nilai prediksi dan nilai aktual perlu dilakukan perubahan *weights* dan *biases*. Untuk dapat menentukan *weights* dan *biases* yang optimal digunakan algoritma *gradient descent* (Jurafsky & Martin, 2020). Perubahan nilai *weights* dan *biases* dari *neural network* menggunakan persamaan *gradient descent* **(15)**, **(16)**, **(17)**, dan **(18)** berikut (Aflak, 2021).

$$W \leftarrow W - \alpha \frac{\partial E}{\partial W} \quad (15)$$

$$\frac{\partial E}{\partial W} = X^T \frac{\partial E}{\partial Y} \quad (16)$$

$$B \leftarrow B - \alpha \frac{\partial E}{\partial B} \quad (17)$$

$$\frac{\partial E}{\partial B} = \frac{\partial E}{\partial Y} \quad (18)$$

Perhitungan nilai *error* antara nilai prediksi dan nilai aktual disebut sebagai fungsi *loss*, *cost*, atau, *error* dan ditentukan sendiri oleh pelatih model. Fungsi *loss* merepresentasikan total *error* dalam *neural network* (Aflak, 2021; Taylor, 2017). Terdapat beberapa fungsi *loss* seperti, *Mean Squared Error* (MSE), *Squared Error* (SE), *Root Mean Squared Error* (RMSE) dan *Sum of Squared Errors* (SSE) (Taylor, 2017). Untuk memahami *backpropagation* akan digunakan fungsi *loss* MSE dengan persamaan (19):

$$E = \frac{1}{n} \sum_{i=1}^n (y_{target}^{(i)} - y^{(i)})^2 \quad (19)$$

Turunan dari persamaan (19) adalah persamaan (20):

$$\frac{\partial E}{\partial Y} = \left[\frac{\partial E}{\partial y_1} \dots \frac{\partial E}{\partial y_n} \right] = \frac{2}{n} [y^{(1)} - y_{target}^{(1)} \dots y^{(n)} - y_{target}^{(n)}] = \frac{2}{n} (Y - Y_{target}) \quad (20)$$

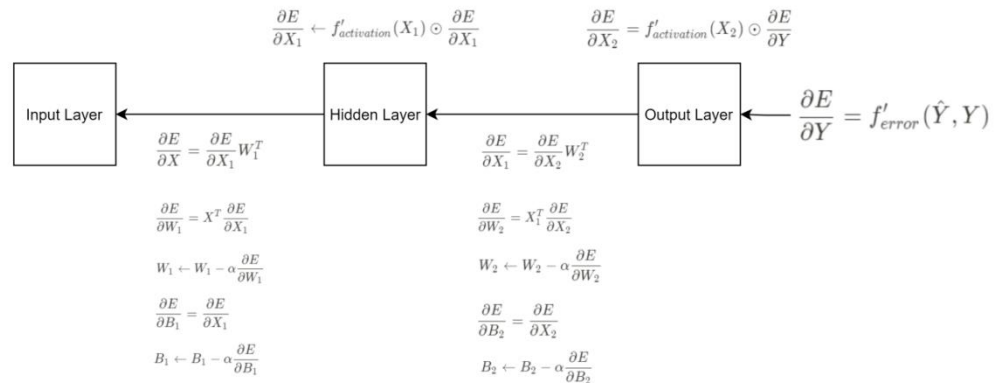
Error dari suatu *layer* dapat diteruskan ke *layer* sebelumnya dengan persamaan (21) berikut.

$$\frac{\partial E}{\partial X} = \frac{\partial E}{\partial Y} W^T \quad (21)$$

Pada fungsi aktivasi, *error* dari *layer* sebelumnya dihitung menggunakan persamaan (22) berikut.

$$\frac{\partial E}{\partial X} \leftarrow f'_{activation}(X) \odot \frac{\partial E}{\partial Y} \quad (22)$$

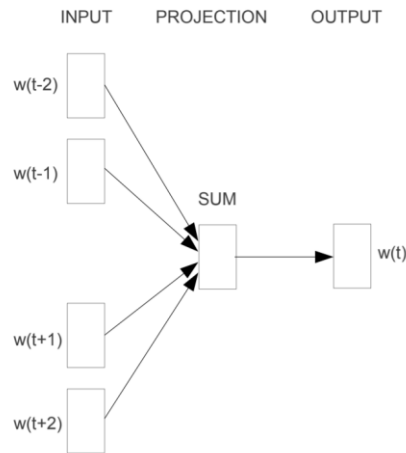
Perhatikan bahwa pada persamaan (22) digunakan *hadamard product*, yaitu perkalian matriks identikal secara elemen. Dengan demikian, proses *backpropagation* terjadi seperti pada **Gambar 2.8.7** (Aflak, 2021).



Gambar 2.8.7 Proses *backpropagation*

2.9 Continuous Bag-of-Words (CBOW)

Continuous Bag-of-Words (CBOW) adalah arsitektur yang mirip *neural network language model* (NNLM) yang mencoba memprediksi kata yang dikelilingi oleh kata-kata lain atau konteks (sebesar C dari kata tengah). CBOW merupakan salah satu arsitektur dari Word2Vec dan yang lainnya adalah skip-gram. Pada CBOW semua input direpresentasikan dengan menggunakan *projection layer* sehingga menghasilkan multi vektor. Multi vektor ini kemudian **dirata-ratakan** menjadi suatu vektor dan masuk ke *output layer*. Fungsi aktivasi *softmax* digunakan pada *output layer* untuk menghasilkan prediksi. Untuk optimalisasi model dengan banyak kosakata digunakan *hierarchical softmax* dan output yang direpresentasikan dengan *huffman tree* yang dibentuk berdasarkan frekuensi kata.



Gambar 2.9.1 Arsitektur CBOW
 Sumber: (Mikolov, Chen, dkk., 2013, hlm. 5)

Setelah proses pelatihan model, vektor yang dihasilkan dari *layer projection* memiliki nilai yang merepresentasikan setiap kata, sehingga vektor tersebut dapat dioperasikan. Misalnya untuk mencari kata yang mirip dengan *small* dalam arti yang sama dengan *biggest* mirip dengan *big*, dapat dihitung dengan cara $X = \text{vektor}(\text{"biggest"}) - \text{vektor}(\text{"big"}) + \text{vektor}(\text{"small"})$. *Embedding* suatu kata yang dihasilkan dicari kemiripannya dengan *embedding* kata lain di dalam vektor menggunakan *cosine distance* (1-*cosine similarity*). Adapun *cosine similarity* dihitung dengan persamaan (23) berikut.

$$\cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}} \quad (23)$$

Bila model telah dilatih maka mungkin menemukan jawaban yang benar, yakni *cosine distance* yang paling kecil, menggunakan metode ini. Untuk mengevaluasi kualitas vektor, digunakan suatu set penalaran analogis yang

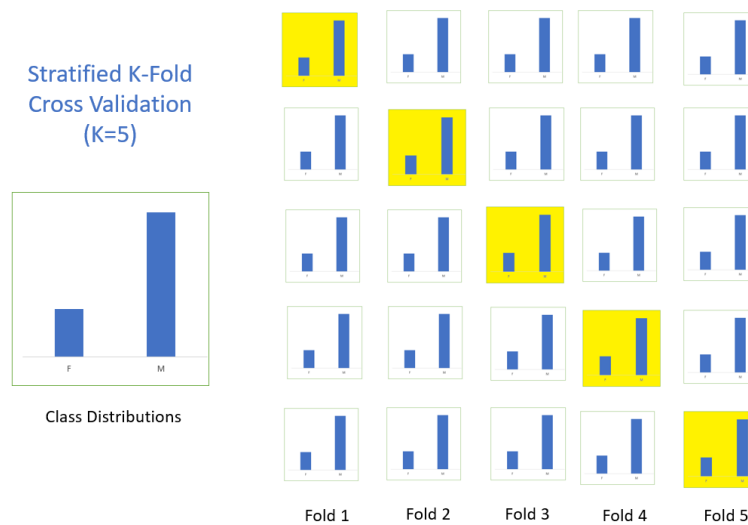
berisi kata dan frasa. Pasangan analogis ini dibentuk dengan cara mengambil frasa yang sering muncul lalu memberikan kata representasinya dari frasa tersebut. Misalnya, "Montreal":"Montreal Canadiens", "Toronto":"Daun Maple Toronto". Jawaban dianggap benar jika representasi terdekat dari vektor("Montreal Canadiens") - vektor("Montreal") + vektor("Toronto") adalah vektor("Toronto Maple Leafs") (Mikolov, Chen, dkk., 2013; Mikolov, Sutskever, dkk., 2013).

Evaluasi dapat juga dilakukan secara intrinsik yakni visualisasi. *Clustering* dapat memberikan visualisasi dengan mengambil n kata yang sering muncul dan memvisualisasikannya. Algoritma k-means digunakan untuk visualisasi pengelompokan karena mempunyai skor *purity* yang tinggi (Zhai dkk., 2016). Dimensi vektor yang tinggi perlu dikurangi menjadi dua dengan menggunakan PCA (*Principal Component Analysis*) atau t-SNE (*t-distributed Stochastic Neighbor Embeddings*) (Liu dkk., 2018). Perbedaannya adalah PCA mengurangi dimensi secara linear dengan memetakan data berdimensi tinggi secara linear ke dimensi yang lebih rendah sambil memaksimalkan variasi dari data. Sedangkan, t-SNE menghitung kemiripan pada dimensi tinggi dan dimensi rendah lalu menggunakan metode optimalisasi, misalnya *gradient descent*, untuk meminimalisir perbedaan antara kedua dimensi tersebut (Winastwan, 2020).

2.10 Stratified K-Fold Cross-Validation

Stratified K-Fold Cross-Validation (SKF) adalah suatu metode yang mengembalikan potongan data (*folds*) sebanyak k (n_splits) dan

memisahkannya menjadi set latih/uji untuk setiap potongan. Potongan tersebut dibuat dengan mempertahankan persentase sampel untuk setiap kelasnya pada masing-masing set (Scikit-Learn, 2021). Pada **Gambar 2.10.1** berikut ini adalah proses SKF untuk membentuk data latih/uji (Sunil, 2018).



Gambar 2.10.1 Ilustrasi Stratified K-Fold Cross Validation
Sumber: (Sunil, 2018, bag. 4)

2.11 Konversi nilai ke rentang baru

Memetakan ulang nomor dari satu rentang ke rentang lainnya dapat dilakukan dengan menggunakan fungsi berikut (Arduino, 2022).

```
long map(long x, long in_min, long in_max, long out_min, long out_max) {
  return (x - in_min) * (out_max - out_min) / (in_max - in_min) + out_min;
}
```

Gambar 2.11.1 Fungsi untuk konversi nilai ke rentang lain