

Seminar II

**KOMBINASI TEKNIK RESAMPLING DAN ALGORITMA
MACHINE LEARNING PADA KELAS TAK SEIMBANG**



Oleh :

ZINEDINE KAHLIL GIBRAN ZIDANE

H13116304

Pembimbing Utama : Dr. Amran, S.Si., M.Si.
Pembimbing Pertama : Supri Bin Hj Amir, S.Si., M.Eng.
Penguji : 1. Dr. Anna Islamiyati, S.Si, M.Si.
2. Nur Hilal A Syahrir, S.Si, M.Si.

PROGRAM STUDI ILMU KOMPUTER
DEPARTEMEN MATEMATIKA
FAKULTAS MATEMATIKA DAN ILMU PENGETAHUAN ALAM
UNIVERSITAS HASANUDDIN
2019

Daftar Isi

BAB I	PENDAHULUAN	1
1.1	Latar Belakang	1
1.2	Rumusan Masalah	6
1.3	Batasan Masalah	6
1.4	Tujuan Penelitian.....	6
1.5	Manfaat Penelitian.....	6
BAB II	TINJAUAN PUSTAKA	8
2.1	Machine Learning.....	8
2.2	Imbalanced Class.....	8
2.3	Confusion Matrix	9
2.4	Jarak Euklid	10
2.5	Principal Common Analysis (PCA)	10
2.6	LASSO dan Ridge Regression	12
2.7	Teknik Resampling.....	13
2.7.1	Random Oversampling.....	13
2.7.2	Synthetic Minority Oversampling Technique (SMOTE).....	14
2.7.3	Borderline – Synthetic Minority Oversampling Technique (Borderline-SMOTE).....	14
2.7.4	Adaptive Synthetic (ADASYN).....	16
2.7.5	Random Undersampling.....	17
2.7.6	Tomek Links	17

2.8	Algoritma Klasifikasi Machine Learning	18
2.8.1	Logistic Regression	18
2.8.2	Decision Tree	21
2.8.3	Support Vector Machine	24
2.8.4	Multilayer Perceptron (MLP).....	32
2.8.5	K-Nearest Neighbor	39
BAB III	METODE PENELITIAN	40
3.1	Waktu dan Tempat	40
3.2	Tahapan Penelitian	40
3.3	Deskripsi Data	41
3.4	Alur Penelitian.....	42
BAB IV	HASIL DAN PEMBAHASAN	43
4.1	Eksplorasi dan Preprocessing Data	43
4.1.1	Image Segmentation Dataset.....	43
4.1.2	Spambase Dataset.....	46
4.1.3	Credit Card Fraud Dataset.....	53
4.2	Model Tuning & Fitting	57
4.3	Analisis Hasil	60
BAB V	KESIMPULAN DAN SARAN	62
5.1	Kesimpulan.....	62
5.2	Saran	63
	Daftar Pustaka	64

BAB I

PENDAHULUAN

1.1 Latar Belakang

Dalam beberapa tahun terakhir, perkembangan yang pesat dalam sains dan teknologi telah berdampak pada pertumbuhan data mentah secara eksponensial (He & Garcia, Learning from Imbalanced Data, 2008). Berdasarkan World Economic Forum (Desjardins, 2019), data digital di dunia akan mencapai 44 zettabyte, atau 44 triliun gigabyte pada tahun 2020. Jumlah ini akan terus berkembang hingga lebih dari dua kali lipat setiap dua tahunnya (Chen, Mao, & & Liu, 2014; Lohr, 2012).

Dari pertumbuhan data tersebut, kebutuhan untuk menganalisis data terus meningkat (Elgendy & Elragal, 2014). Menurut laporan dari LinkedIn pada tahun 2018, permintaan pekerjaan yang membutuhkan analisis data berkembang hingga 12 kali lipat lebih banyak dari tahun 2014 (LinkedIn Economic Graph Team, 2018). Ini disebabkan karena data merupakan kunci dari setiap organisasi, institusi atau perusahaan untuk mengambil keputusan yang lebih cermat dan efektif (McAfee, Brynjolfsson, Davenport, Patil, & Barton, 2012).

Namun data mentah atau *raw data* tidak dapat memberikan informasi yang berguna. Data tersebut perlu diolah dan dianalisis menjadi informasi yang berguna. *Raw data* yang dimaksud adalah data yang belum diproses, yang terstruktur dan yang tidak terstruktur. Data yang terstruktur adalah data yang disimpan dengan format yang telah ditentukan seperti database (Beal, 2019), di mana atribut data dapat dibedakan dengan jelas sehingga dapat langsung diproses oleh peralatan komputasi (Baars & Kemper, 2008). Sedangkan data tidak terstruktur adalah data yang disimpan dalam format yang tidak terstruktur, sehingga membutuhkan campur tangan manusia agar dapat diinterpretasi oleh mesin; seperti dokumen, gambar, video dan audio (Weglarz, 2004). Menurut laporan Beal (Beal, 2019), 80% hingga 90% data di dunia tidak terstruktur.

Dengan jumlah data yang sangat banyak, tidak mungkin oleh manusia untuk menganalisis dan membuat perhitungan mengenai data secara manual. Maka dari itu diperlukan bidang khusus untuk mengolah dan menganalisis data. *Data Science* adalah bidang yang mempelajari bagaimana meng-ekstraksi *raw data* menjadi *meaningful information* atau informasi yang berguna (Berman, et al., 2018; Dhar, 2012). Data Science melibatkan prinsip, proses, dan teknik untuk memahami fenomena melalui data (Provost & Fawcett, 2013). Data Science merupakan bidang yang sangat luas dan sedang dikembangkan (Provost & Fawcett, 2013), namun salah satu bidang khusus dari data science adalah *machine learning* yang merupakan paduan antar *computer science* (ilmu komputer) dan *statistics* (statistika) (Jordan & Mitchell, 2015). Machine learning membahas mengenai bagaimana membangun sistem komputer yang dapat belajar melalui pengalaman tanpa harus diprogram secara spesifik dan manual (Jordan & Mitchell, 2015) (Domingos, 2011).

Untuk menyelesaikan suatu masalah pada komputer, dibutuhkan algoritma. Namun untuk beberapa masalah, tidak ada algoritma yang memadai. Contohnya adalah bagaimana komputer mengenali angka dalam bentuk tulisan tangan (Bishop, Neural Networks for Pattern Recognition, 1995) dan mengklasifikasi suatu email menjadi spam atau bukan spam (Ethem, 2009). Dalam masalah seperti ini, seluruh objek dengan label atau kelas berbeda dikumpulkan dan diidentifikasi ciri-ciri yang membedakan mereka. Di dalam hal ini lah machine learning bekerja.

Data yang dikumpulkan biasanya dalam bentuk dataset atau tabel, di mana setiap kolomnya adalah atribut atau ciri-ciri dan setiap barisnya adalah instansi atau observasi. Dataset tersebut ada yang memiliki kolom label, kelas, atau kolom yang berisi informasi mengenai kategori dari setiap observasi (contohnya, spam atau bukan spam), dan ada juga dataset yang tidak memiliki kolom label, di mana isinya hanyalah atribut atau ciri-ciri dari setiap observasi, tanpa mengindikasikan kategori dari tiap observasi (Bishop, Neural Networks for Pattern Recognition, 1995). Pembelejaran pada dataset berlabel disebut *supervised learning*. Kasus di mana tujuannya adalah

mengklasifikasikan input data ke suatu kategori diskrit tertentu disebut *klasifikasi*, dan kasus di mana outputnya adalah suatu variabel kontinu disebut *regresi*. Selain itu, pembelajaran pada dataset tanpa label atau acuan kategori yang benar disebut *unsupervised learning*. Kasus *unsupervised learning* di mana tujuannya adalah mengelompokkan observasi-observasi yang mirip disebut *clustering*, jika menentukan distribusi data pada input disebut *estimasi kepadatan*. Dan yang terakhir, pembelajaran di mana mesin dilatih untuk membuat keputusan tertentu dengan cara *trial and error* disebut *reinforcement learning* (Bishop, Pattern Recognition and Machine Learning, 2006). Masing-masing jenis pembelajaran memiliki banyak algoritma yang telah dikembangkan dengan berbagai pendekatan yang berbeda-beda (Jordan & Mitchell, 2015). Berdasarkan dokumentasi dari *scikit-learn* (Pedregosa, et al., 2011), terdapat lebih dari 100 algoritma machine learning yang ada.

Untuk distribusi data pada suatu dataset, terdapat istilah kelas yang terdistribusi secara seimbang (*balanced*) dan secara tak seimbang (*imbalanced*). Dataset dengan kelas yang seimbang berarti jumlah observasi untuk setiap kelas tidak jauh dari kelas-kelas yang lain (Galar, Fernandez, Barrenechea, Bustince, & Herrera, 2011). Sedangkan untuk dataset dengan distribusi kelas yang tak seimbang, jumlah suatu observasi pada kelas tertentu sangat jauh berbeda dengan kelas yang lain. Hal ini berlaku pada dataset dengan kelas biner (dua kelas saja) dan juga *multiclass* (lebih dari dua kelas) (He & Garcia, Learning from Imbalanced Data, 2008). Kelas dengan jumlah observasi sedikit disebut kelas minoritas (*minority class*) dan kelas dengan jumlah observasi yang sangat banyak disebut kelas mayoritas (*majority class*). Tidak jarang suatu dataset terdistribusi secara tak seimbang dengan proporsi antara kelas minoritas dan mayoritasnya adalah 1:100, 1:1000, atau 1:1000 (He & Garcia, Learning from Imbalanced Data, 2008). Sebagian besar data asli di dunia terdistribusi secara tak seimbang (He & Garcia, Learning from Imbalanced Data, 2008; Kotsiantis, Kanellopoulos, & Pintelas, 2006; Kumar & Sheshadri, 2012; Visa & Ralescu, 2005).

Pada umumnya, algoritma-algoritma machine learning, dalam hal ini pada masalah klasifikasi, bekerja dengan tujuan utama memaksimalkan akurasi (Provost F. , 2000). Hal ini sangat masuk akal, karena akurasi yang tinggi menjelaskan bahwa model algoritma tersebut melaksanakan tugasnya dengan baik, mengklasifikasikan kelas data dengan benar dengan sedikit kesalahan. Namun, akurasi hanya memberikan informasi secara umum, bagaimana jika model algoritma tersebut bekerja pada dataset tak seimbang, dan hanya mampu mengklasifikasikan kelas mayoritas dengan benar tetapi tak mampu mengklasifikasikan kelas minoritas? Jika perbandingan antara kelas minoritas dan mayoritas saja satu berbanding seratus, maka akurasi yang akan diperoleh lebih besar dari 99%, dengan kesalahan lebih kecil dari 1% yang hampir seluruhnya adalah kelas minoritas. Masalah ini memberi bias terhadap performa algoritma-algoritma klasifikasi, terutama jika kelas yang lebih utama untuk diklasifikasikan dengan benar adalah kelas minoritas, seperti email spam, diagnosis penyakit di bidang kedokteran, deteksi kartu kredit palsu dan lain-lain (Visa & Ralescu, 2005; Rahman & Davis, 2013). Hal ini menunjukkan bahwa pada kasus dataset tak seimbang, dibutuhkan perhatian lebih terhadap *preprocessing* data sebelum dimasukkan ke model.

Banyak cara yang telah ditemukan untuk mengatasi dataset tak seimbang ini, seperti melakukan *resampling* terhadap data yang ada. Resampling adalah teknik mengambil sampel secara berulang dari sampel data asli (Statistics Solution, 2016). Teknik resampling terdiri dari *oversampling*, yaitu mengambil sampel berulang kali dari kelas minoritas; dan *undersampling*, yaitu mengambil sampel secara acak dari kelas mayoritas (Burnaev, Erofeev, & Papanov, 2015). Kedua teknik ini dapat digunakan secara terpisah ataupun digabung (Burnaev, Erofeev, & Papanov, 2015; Anand, Pugalenthi, Fogel, & Suganthan, 2010; More, 2016; Yen & Lee, Under-Sampling Approaches for Improving Prediction of the Minority Class in an Imbalanced Dataset, 2006). SMOTE adalah teknik oversampling yang terpopuler, dengan Borderline-SMOTE merupakan ekstensi dari SMOTE. Salah satu Teknik resampling

yang cukup populer adalah ADASYN yang mampu menyesuaikan jumlah data sintetiknya.

Dalam beberapa penelitian terkait (More, 2016; Batista, Prati, & Monard, 2004; Amin, et al., 2016; Burnaev, Erofeev, & Papanov, 2015), telah dilakukan berbagai percobaan untuk mengatasi masalah dataset tak seimbang, namun metode-metode resampling maupun algoritma machine learning yang digunakan tidak beragam untuk mengetahui metode terbaik untuk mengatasi masalah ini. Seperti penelitian yang dilakukan Amin (Amin, et al., 2016) hanya meneliti teknik oversampling, Burnaev, More, dan Batista dkk (Burnaev, Erofeev, & Papanov, 2015; More, 2016; Batista, Prati, & Monard, 2004) meneliti teknik oversampling dan undersampling namun hanya menggunakan satu algoritma machine learning, sedangkan Diri (Diri & Albayrak, 2008) hanya meneliti beberapa algoritma machine learning tanpa pertimbangan dataset tak seimbang.

Sedangkan untuk mengetahui metode resampling dan algoritma machine learning terbaik untuk masalah ini, dibutuhkan kombinasi-kombinasi antar teknik resampling, dan juga antar algoritma machine learning. Setiap kombinasi (pasangan) ini, seperti SMOTE dengan Support Vector Machine, atau Tomek Links dengan Regresi Logistik akan diuji performanya terhadap dataset yang diberikan, kemudian dari kombinasi-kombinasi tersebut dapat ditarik kesimpulan mengenai kombinasi algoritma dan teknik resampling yang terbaik, dan algoritma machine learning dengan performa terbaik, dan teknik resampling dengan performa terbaik. Setiap kombinasi atau pasangan dievaluasi hasilnya dengan tidak hanya pada satu dataset tak seimbang saja, melainkan dengan beberapa dataset tambahan untuk mendapatkan hasil yang lebih umum.

Berdasarkan uraian di atas, penulis ingin melakukan penelitian mengenai dataset tak seimbang dengan menggunakan “Kombinasi Teknik Resampling dan Algoritma Machine Learning pada Kelas tak Seimbang”.

1.2 Rumusan Masalah

Berdasarkan uraian pada latar belakang masalah di atas, dapat dirumuskan masalah: Kombinasi algoritma Machine Learning dan teknik Resampling yang mana saja yang berperforma dengan baik untuk mengatasi dataset tak seimbang?

1.3 Batasan Masalah

Batasan masalah pada penelitian ini adalah:

1. Merupakan masalah klasifikasi (*Supervised Learning*).
2. Dataset yang digunakan adalah dataset tak seimbang (*imbalanced ratio* < 3:4).
3. Dataset hanya memiliki dua kelas (biner).
4. Atribut-atribut dataset merupakan tipe numerik.
5. Teknik resampling yang digunakan adalah Random Undersampling, Tomek Links, Random Oversampling, SMOTE, Borderline-SMOTE, dan ADASYN.
6. Algoritma machine learning yang digunakan adalah Regresi Logistik, Decision Tree, Neural Network, Nearest Neighbor, dan Support Vector Machine.

1.4 Tujuan Penelitian

Berdasarkan rumusan masalah, maka tujuan dari penelitian ini adalah:

1. Mengetahui algoritma machine learning yang berperforma dengan baik dalam mengklasifikasi data.
2. Mengetahui teknik resampling yang memiliki performa baik pada dataset tak seimbang.
3. Mengetahui kombinasi algoritma machine learning dan teknik resampling yang berperforma dengan baik dalam mengatasi dataset tak seimbang.

1.5 Manfaat Penelitian

Hasil penelitian ini diharapkan dapat bermanfaat:

1. Sebagai rujukan untuk mengatasi dataset tak seimbang yang sering dijumpai.
2. Menjadi sumber informasi mengenai performa beberapa teknik resampling.

3. Menjadi sumber informasi mengenai performa dari beberapa algoritma machine learning.

BAB II

TINJAUAN PUSTAKA

2.1 Machine Learning

Machine learning adalah memprogram komputer untuk mengoptimalkan suatu ukuran kinerja menggunakan sampel data atau berdasarkan pengalaman (Ethem, 2009). Machine learning menggunakan suatu algoritma untuk menganalisis data.

Dalam pembelajaran yang terawasi atau *supervised learning*, pengklasifikasi akan diberikan suatu input tertentu dan menghubungkannya dengan suatu output. Kasus di mana tujuannya adalah mengklasifikasikan input data ke suatu kategori diskrit tertentu disebut *klasifikasi*, dan kasus di mana outputnya adalah suatu variabel kontinu disebut *regresi*. Dalam pembelajaran tanpa pengawasan atau *unsupervised learning*, pengklasifikasi diberi input dan dibiarkan sendiri untuk menemukan pola pada data tersebut. Kasus *unsupervised learning* di mana tujuannya adalah mengelompokkan observasi-observasi yang mirip disebut *clustering*, jika menentukan distribusi data pada input disebut *estimasi kepadatan*. Dalam *reinforcement learning*, sistem komputer menerima input secara terus menerus dan mencoba memilih keputusan-keputusan yang paling optimal berdasarkan kondisi lingkungannya.

Masing-masing jenis pembelajaran memiliki banyak algoritma yang telah dikembangkan dengan pendekatan yang berbeda-beda (Jordan & Mitchell, 2015).

2.2 Imbalanced Class

Dataset adalah kumpulan data yang berbentuk tabel, di mana setiap kolomnya merepresentasikan suatu ciri-ciri, atribut atau fitur. Setiap barisnya menyatakan observasi suatu individu, record atau sampel (Snijders, Matzat, & Reips, 2012). Suatu dataset biasanya memiliki satu kolom tambahan yang merepresentasikan kelas dari observasi tersebut, kolom ini disebut kolom kelas. Kolom kelas ini juga disebut sebagai variabel dependen terhadap variabel-variabel independen yang merupakan ciri-ciri (atribut) dari suatu observasi tertentu.

Dalam machine learning dikenal istilah dataset dengan class yang tak seimbang. Istilah ini berlaku ketika kelas dari dataset tersebut bersifat kategorik diskrit. Dataset dengan class yang tak seimbang (*imbalanced class*) adalah dataset yang frekuensi kejadian dari kelas tertentu sangat jauh berbeda dengan kelas yang lain. Contohnya seperti suatu dataset dengan jumlah pasien yang berkelas “diabetes” jumlahnya jauh lebih sedikit dibanding pasien yang “tidak diabetes”.

Masalah ketidakseimbangan ini akan memberi bias terhadap performa pengklasifikasi sebab jumlah sampel pada kelas tertentu tidak dapat memberi informasi yang cukup kepada pengklasifikasi berdasarkan ciri-ciri yang diberikan (Bishop, Pattern Recognition and Machine Learning, 2006; Ethem, 2009; Domingos, 2011).

2.3 Confusion Matrix

Di dalam machine learning, mengukur kinerja atau performa dari suatu model adalah hal yang esensial. Model yang diperoleh dari pelatihan melalui data training perlu diuji melalui data testing. Kinerja diukur berdasarkan seberapa baik model tersebut memprediksi dengan benar data yang ada.

Pada klasifikasi biner, kelas positif yang berhasil diprediksi dengan benar disebut *true positive*, jika kelas positif tersebut diprediksi negatif (salah) disebut *false negative*. Kelas negatif yang berhasil diprediksi negatif (benar) disebut *true negative*, dan kelas negatif yang diprediksi positif disebut *false positive*. Jumlah dari kasus-kasus tersebut direpresentasikan dalam suatu tabel kontingensi yang disebut *confusion matrix* (Swets, 1988).

Hasil prediksi	Kelas asli		
		Positif	Negatif
	Positif	TP	FP
	Negatif	FN	TN

Tabel 1 Confusion Matrix

Akurasi adalah ukuran kinerja yang menunjukkan seberapa baik suatu pengklasifikasi dalam mengklasifikasikan seluruh data. Akurasi adalah rasio antara observasi yang diklasifikasikan secara benar dengan total observasi:

$$Accuracy = \frac{(TP + TN)}{(TP + FP + FN + TN)} \quad (II.1)$$

Presisi adalah ukuran kinerja yang menunjukkan seberapa besar kebenaran suatu pengklasifikasi dari seluruh kelas positif yang diprediksi. Presisi adalah rasio antara jumlah kelas positif yang diklasifikasikan secara benar dengan jumlah observasi yang diklasifikasikan positif:

$$Precision = \frac{TP}{(TP + FP)} \quad (II.2)$$

Recall atau sensitivitas adalah ukuran kinerja yang menunjukkan seberapa baik suatu pengklasifikasi dalam mengklasifikasikan kelas positif. Recall adalah rasio antara jumlah observasi positif yang diklasifikasikan secara benar dengan jumlah observasi positif asli:

$$Recall = \frac{TP}{(TP + FN)} \quad (II.3)$$

2.4 Jarak Euklid

Jarak Euklid adalah jarak antara suatu vektor $X = (x_1, \dots, x_n)$ ke suatu vektor $Y = (y_1, \dots, y_n)$ pada ruang euklid berdimensi n :

$$d(X, Y) = \sqrt{(x_1 - y_1)^2 + \dots + (x_n - y_n)^2} \quad (II.4)$$

Pada machine learning, jarak Euklid digunakan untuk menghitung jarak antar dua observasi berdasarkan vektor fitur yang bersifat kontinu (Howard, 2013).

2.5 Principal Common Analysis (PCA)

Principal Common Analysis atau PCA adalah metode feature extraction dari suatu dataset dengan p atribut numerik, untuk setiap atribut memiliki n observasi. Dengan

kata lain didefinisikan sejumlah p vektor X berdimensi n (X_1, \dots, X_p) atau $n * p$ data matrix \mathbf{X} , di mana kolom ke- j adalah vektor X_j . Cari kombinasi linear dari kolom dari matrix \mathbf{X} dengan variansi maksimum. Kombinasi linear tersebut diekspresikan dengan:

$$\sum_{j=1}^p A_j X_j = \mathbf{X}A$$

di mana A adalah vektor konstan a_1, a_2, \dots, a_p .

Variansi dari kombinasi linear tersebut didefinisikan sebagai:

$$\text{var}(\mathbf{X}A) = A^T \mathbf{S}A \quad (\text{II.5})$$

di mana \mathbf{S} adalah matrix kovarians dari \mathbf{X} dan A^T adalah vektor transpose dari A . Maka mengidentifikasi kombinasi linear dengan variansi maksimum ekuivalen dengan memperoleh vektor A berdimensi p yang memaksimalkan ekspresi $A^T \mathbf{S}A$. Untuk memperoleh solusi yang terdefinisi, vektor A harus merupakan vektor satuan, atau $A^T A = 1$:

$$A^T \mathbf{S}A - \lambda(A^T A - 1) \quad (\text{II.6})$$

di mana λ adalah pengganda Lagrange. Menurunkan ekspresi (II.6) terhadap vektor A , dan menyamakan dengan vektor nol, menghasilkan persamaan:

$$\mathbf{S}A - \lambda A = \mathbf{0} \Leftrightarrow \mathbf{S}A = \lambda A \quad (\text{II.7})$$

Dapat disimpulkan bahwa A adalah vektor eigen, dan λ adalah nilai eigen dari matrik kovarians \mathbf{S} . Karena nilai eigen adalah variansi dari transformasi linear pada persamaan (II.5):

$$\begin{aligned} \text{var}(\mathbf{X}A) &= A^T \mathbf{S}A \\ &= \lambda A^T A \\ &= \lambda \end{aligned}$$

maka persamaan (II.7) bernilai benar jika setiap vektor eigen dikali -1.

Seluruh matriks real simetrik berukuran $p * p$, seperti matriks kovarians \mathbf{S} , memiliki tepat p nilai eigen yang real, $\lambda_k (k = 1, \dots, p)$, dan vektor eigennya dapat didefinisikan untuk membentuk satu set vektor ortonormal. Pendekatan pengganda Lagrange, dengan batasan tambahan ortogonal dari vektor koefisien yang berbeda, juga dapat digunakan untuk menunjukkan bahwa seluruh vektor eigen dari \mathbf{S} adalah solusi dari masalah dengan memperoleh p kombinasi linear baru:

$$\mathbf{X}A_k = \sum_{j=1}^p A_{jk}X_j \quad (\text{II.8})$$

yang memaksimalkan varians, tergantung pada korelasi dengan kombinasi linear sebelumnya. Pada kombinasi linear ini, $\mathbf{X}A_k$ adalah principal component atau komponen utama dari dataset.

(Jolliffe & Cadima, 2016)

2.6 LASSO dan Ridge Regression

Regularisasi adalah metode untuk menghindari overfitting dengan memberi nilai penalti terhadap koefisien regresi yang bernilai tinggi. Regularisasi mengurangi parameter dan menyederhanakan model yang berbentuk kompleks. Regularisasi menambah nilai penalti pada model yang lebih kompleks dan mengoptimisasi parameter dengan meminimalkan nilai suatu cost function.

LASSO Regression atau Least Absolute Shrinkage and Selection Operator Regression adalah metode regularisasi yang menambahkan penalti sebesar nilai absolut dari besarnya koefisien dikalikan dengan suatu nilai λ . Regularisasi ini dapat menghasilkan model dengan koefisien yang lebih sedikit dikarenakan beberapa koefisien dapat menjadi nol dan dihilangkan dari model. Nilai penalti λ yang lebih besar menghasilkan nilai-nilai koefisien yang lebih mendekati nol, yang ideal untuk menghasilkan model yang lebih sederhana.

$$\sum_{i=1}^n \left(y_i - \sum_j x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p |\beta_j| \quad (\text{II.9})$$

Ridge Regression adalah metode regularisasi yang menambahkan penalti sebesar nilai kuadrat dari besarnya koefisien dikalikan dengan suatu nilai λ . Regularisasi ini menambahkan bias terhadap regresi dengan hasil nilai variansi yang lebih kecil.

$$\sum_{i=1}^n \left(y_i - \sum_j x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \quad (\text{II.10})$$

Dilihat dari persamaan (II.9) dan persamaan (II.10), jika λ bernilai nol, maka cost function tersebut kembali ke Ordinary Least Squares, jika nilai λ terlalu besar, maka bobot yang diberikan akan terlalu besar dan model dapat mengalami underfitting.

Perbedaan dari LASSO dan Ridge Regression adalah pada LASSO Regression, dengan semakin besarnya nilai λ , koefisien-koefisien parameter dapat dirubah menjadi nol, mengeliminasi fitur yang tidak berperan terhadap model; sedangkan pada Ridge Regression, koefisien parameter tidak akan mencapai nol, namun secara asimtotik mendekati nol. (Tibshirani, 1996; Pereira, Basto, & Silva, 2016)

2.7 Teknik Resampling

2.7.1 Random Oversampling

Random oversampling, atau oversampling secara acak adalah teknik oversampling di mana anggota dari kelas minoritas dipilih secara acak dan diduplikasi ke dataset yang baru hingga tercapai keseimbangan (Liu A. Y.-c., 2004). Data minoritas tersebut dapat diduplikasi beberapa kali. Teknik ini biasanya menyebabkan overfitting pada model (Amin, et al., 2016; Liu A. Y.-c., 2004).

2.7.2 Synthetic Minority Oversampling Technique (SMOTE)

SMOTE atau Synthetic Minority Oversampling Technique adalah teknik oversampling terpopuler yang diproposalkan oleh Chawla (Chawla, Bowyer, Hall, & Kegelmeyer, 2002) pada tahun 2002. Teknik ini membuat data tiruan atau sintetik berdasarkan tetangga-tetangga terdekat dari sampel kelas minoritas.

Teknik ini dimulai dengan menentukan n data yang akan dibuat untuk setiap data pada kelas minoritas dalam dataset \mathbf{D} . Kemudian untuk setiap data kelas minoritas $M_{i,c}$, pilih n tetangga secara acak dari k tetangga terdekat data tersebut, di mana $i = \{1, 2, \dots, m\}$ dengan m adalah jumlah data pada kelas minoritas dan $c \in C$, di mana $C = \{c_1, c_2, \dots, c_z\}$ adalah fitur-fitur pada \mathbf{D} . Lalu untuk setiap fitur c pada $M_{i,c}$, hitung jarak euklid d antara $M_{i,c}$ dengan salah satu tetangga $T_{s,c}$, di mana s adalah bilangan acak $s = \{1, 2, \dots, k\}$ dari k tetangga terdekat $M_{i,c}$. Kemudian suatu bilangan acak $g = [0, 1]$ ditentukan. Data tiruan dibuat berdasarkan:

$$S_{j,c} = M_{i,c} + (g * d) \quad (\text{II.11})$$

di mana $j = \{1, 2, \dots, n * m\}$ bersifat inkremental dan S adalah data kelas minoritas tiruan.

Teknik ini membuat $n * m$ data sintetik pada suatu titik dari jarak antara setiap fitur dari M dengan T (Chawla, Bowyer, Hall, & Kegelmeyer, 2002).

2.7.3 Borderline – Synthetic Minority Oversampling Technique (Borderline-SMOTE)

Terinspirasi oleh SMOTE, Han memproposalkan teknik oversampling baru yang menyerupai SMOTE. Teknik Han tersebut membuat data tiruan di sekitar data kelas minoritas yang berada di *borderline* atau perbatasan antara kelas mayoritas dan kelas minoritas saja, di mana SMOTE membuat N data tiruan pada setiap sampel kelas minoritas (Han, Wang, & Mao, 2005).

Misalkan untuk data training \mathbf{D} terdapat kelas minoritas P dan kelas mayoritas N , diekspresikan:

$$P = \{p_1, p_2, \dots, p_{pnum}\}, N = \{n_1, n_2, \dots, n_{nnum}\}$$

di mana $pnum$ dan $nnum$ adalah jumlah sampel kelas minoritas dan jumlah sampel kelas mayoritas. Borderline-SMOTE bekerja sebagai berikut.

Untuk setiap $p_i (i = 1, 2, \dots, pnum)$ pada kelas minoritas P , hitung m tetangga terdekat dari seluruh data training D . Jumlah kelas mayoritas dari m tetangga terdekat p_i disimbolkan $m' (0 \leq m' \leq m)$.

Jika $m' = m$, atau seluruh tetangga terdekat dari p_i adalah sampel kelas mayoritas, maka p_i dinyatakan sebagai noise dan bukan merupakan sampel perbatasan. Jika $m/2 \leq m' < m$, atau tetangga terdekat p_i lebih banyak merupakan sampel kelas mayoritas dibanding minoritas, maka p_i dianggap sampel perbatasan dan dimasukkan ke dalam himpunan sampel perbatasan B . Jika $0 \leq m' < m/2$, maka p_i dianggap aman dan bukan merupakan sampel perbatasan.

Untuk himpunan sampel perbatasan B , di mana $B \subseteq P$. Maka:

$$B = \{p'_1, p'_2, \dots, p'_{dnum}\}, 0 \leq dnum \leq pnum$$

di mana $dnum$ adalah jumlah sampel perbatasan. Untuk setiap sampel perbatasan di B , tentukan k tetangga terdekat dari P .

Data tiruan berkelas minoritas dibuat sebanyak $s * dnum$ ($s = 1, 2, \dots, k$). Untuk setiap p'_i , Pilih sebanyak s tetangga terdekat secara acak dari P , lalu hitung jarak euklid (Persamaan (II.4)) $d_j (j = 1, 2, \dots, s)$ antara p'_i dengan p_j di mana d_j adalah jarak antara p'_i dengan salah satu tetangga terdekatnya. Kemudian suatu bilangan acak g ($g = [0, 1]$) ditentukan dan data tiruan sebanyak s dari p'_i dibuat.

$$S_j = p'_i + (g * d_j) \quad (II.12)$$

Proses di atas dilakukan untuk setiap p'_i , maka data tiruan akan dibuat sebanyak $s * dnum$ kali (Han, Wang, & Mao, 2005).

2.7.4 Adaptive Synthetic (ADASYN)

Adaptive Synthetic adalah teknik oversampling yang diproposalkan oleh He di mana data tiruan dibuat berdasarkan tingkat kesulitan suatu sampel kelas minoritas untuk dipelajari. Setiap sampel kelas minoritas memiliki alokasi data tiruan sesuai dengan banyaknya tetangga sampel kelas mayoritas dari k tetangga terdekat sampel kelas minoritas tersebut (He, Bai, Garcia, & Li, 2008).

Misalkan untuk data training D dengan m sampel $\{x_i, y_i\}$, $i = 1, \dots, m$, di mana x_i adalah instansi dari n dimensi dari ruang fitur X dan $y_i \in \{1, -1\}$ adalah label kelas dari x_i . Didefinisikan m_s dan m_l sebagai jumlah sampel kelas minoritas, dan jumlah sampel kelas mayoritas secara berurutan. Maka $m_s \leq m_l$ dan $m_s + m_l = m$.

Tentukan derajat ketidakseimbangan d :

$$d = \frac{m_s}{m_l} \quad (\text{II.13})$$

di mana $d \in (0, 1]$. Jika $d < d_{th}$ di mana d_{th} adalah nilai maksimum toleransi derajat ketidakseimbangan, maka hitung jumlah data kelas minoritas tiruan:

$$G = (m_l - m_s) * \beta \quad (\text{II.14})$$

di mana $\beta \in [0, 1]$ adalah parameter yang menunjukkan rasio ketidakseimbangan yang diinginkan setelah data tiruan dibuat. $\beta = 1$ menunjukkan dataset akan seimbang sepenuhnya.

Untuk setiap sampel $x_i \in P$ di mana P adalah kelas minoritas, tentukan K tetangga terdekat dengan jarak Euklid pada n dimensi ruang fitur. Kemudian hitung rasio r_i , yang didefinisikan sebagai:

$$r_i = \frac{\Delta_i}{K}, i = 1, \dots, m_s \quad (\text{II.15})$$

di mana Δ_i adalah jumlah sampel pada K tetangga terdekat x_i yang merupakan sampel kelas mayoritas, maka $r_i \in [0, 1]$. Kemudian normalisasikan r_i menjadi:

$$\hat{r}_i = \frac{r_i}{\sum_{i=1}^{m_s} r_i} \quad (\text{II.16})$$

agar \hat{r}_i menjadi distribusi kepadatan ($\sum_i \hat{r}_i = 1$).

Hitung data tiruan yang akan dibuat untuk setiap sampel x_i :

$$g_i = \hat{r}_i * G \quad (\text{II.17})$$

Untuk setiap sampel $x_i \in P$, buat g_i data tiruan dengan persamaan:

$$s_i = x_i + (x_{zi} - x_i) * \lambda \quad (\text{II.18})$$

di mana $(x_{zi} - x_i)$ adalah selisih vektor pada n dimensi dan λ adalah bilangan acak $\lambda \in [0, 1]$.

2.7.5 Random Undersampling

Random undersampling, atau undersampling secara acak adalah teknik undersampling di mana anggota dari kelas mayoritas dipilih secara acak dan dihapus dari dataset training hingga tercapai keseimbangan. Kekurangan dari teknik ini adalah tidak ada cara untuk mengatur informasi apa saja yang dihilangkan dari dataset tersebut, informasi yang berguna bisa saja hilang (Amin, et al., 2016; Liu A. Y.-c., 2004; Yen & Lee, Under-Sampling Approaches for Improving Prediction of the Minority Class in an Imbalanced Dataset, 2006; More, 2016).

2.7.6 Tomek Links

Tomek Link atau tautan Tomek adalah metode undersampling yang diproposalkan oleh Tomek (Tomek, 1976) untuk memodifikasi CNN (Condensed Nearest Neighbor). Teknik ini menghapus sampel kelas mayoritas jika dan hanya jika tetangga kelas mayoritas tersebut berasal dari kelas berbeda dan merupakan tetangga terdekat satu sama lain. Jika x dan y adalah sampel dengan kelas berbeda, maka suatu Tomek Link didefinisikan untuk setiap sampel z :

$$(d(x, y) < d(x, z)) \wedge (d(x, y) < d(y, z))$$

di mana fungsi $d(\cdot)$ adalah jarak Euklid pada persamaan (II.4).

2.8 Algoritma Klasifikasi Machine Learning

2.8.1 Logistic Regression

Regresi logistik adalah algoritma klasifikasi yang menggunakan kelas untuk membangun dan menggunakan model regresi logistik multinomial tunggal dengan *estimator* tunggal. Regresi logistik menyatakan probabilitas kelas (antara 0 dan 1) tergantung pada jarak dari batas, dengan satu pendekatan. Hasil regresi logistik adalah suatu bilangan antara 0 dan 1 yang menyatakan probabilitas kelas.

Nilai atribut pada data observasi dari regresi logistik dapat berupa nominal, ordinal, interval atau skala rasio, sedangkan untuk atribut kelas harus merupakan kelas biner. Hubungan antara atribut dengan kelas bersifat non-linear. Distribusi data tidak tersebar dalam bentuk distribusi Gauss, melainkan distribusi Bernoulli, yang dikarenakan kelasnya berbentuk biner.

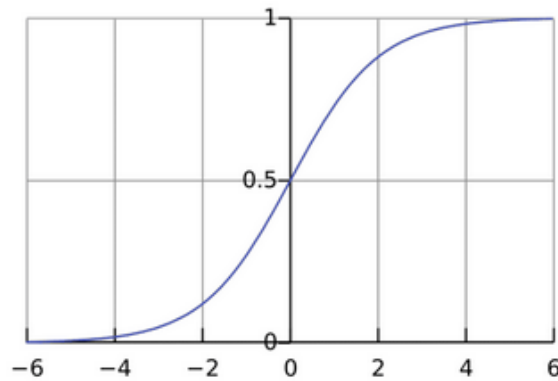
Hipotesis H_θ dari logistik regresi didefinisikan sebagai:

$$H_\theta(x) = g(\theta^T x) \quad (\text{II.19})$$

di mana:

$$g(z) = \frac{1}{1 + e^{-z}}$$

adalah **fungsi logistik** atau **fungsi sigmoid**.



Gambar 1: Fungsi Sigmoid

Terlihat bahwa $g(z)$ cenderung ke 1 jika $z \rightarrow \infty$ dan $g(z)$ cenderung ke 0 jika $z \rightarrow -\infty$. Lebih jauh, $g(z)$ dan $h(z)$ selalu bernilai dalam interval 0 dan 1. Jika dimisalkan $x_0 = 1$, maka:

$$\theta^T x = \theta_0 + \sum_{j=1}^n \theta_j x_j$$

dengan θ_0 adalah nilai **intercept**.

Turunan dari fungsi sigmoid, atau g' adalah:

$$\begin{aligned} g'(z) &= \frac{d}{dz} \frac{1}{1 + e^{-z}} \\ &= \frac{1}{(1 + e^{-z})^2} (e^{-z}) \\ &= \frac{1}{(1 + e^{-z})} \left(1 - \frac{1}{(1 + e^{-z})} \right) \\ &= g(z)(1 - g(z)) \end{aligned} \tag{II.20}$$

Untuk melakukan fitting θ ke dalam regresi logistik, asumsikan beberapa persamaan probabilistik dan tentukan nilai θ melalui maximum likelihood.

Asumsikan:

$$P(y = 1 | x; \theta) = h_{\theta}(x)$$

$$P(y = 0 | x; \theta) = 1 - h_{\theta}(x)$$

dapat ditulis:

$$p(y | x; \theta) = (h_{\theta}(x))^y (1 - h_{\theta}(x))^{1-y}$$

Dengan mengasumsikan m sampel data training yang saling independen, likelihood dari parameter-parameter tersebut dapat ditulis:

$$\begin{aligned} L(\theta) &= p(\vec{y} | X; \theta) \\ &= \prod_{i=1}^m p(y_i | x_i; \theta) \\ &= \prod_{i=1}^m (h_{\theta}(x_i))^{y_i} (1 - h_{\theta}(x_i))^{1-y_i} \end{aligned}$$

Log likelihood-nya adalah:

$$\begin{aligned} \ell(\theta) &= \log L(\theta) \\ &= \sum_{i=1}^m y_i \log h(x_i) + (1 - y_i) \log(1 - h(x_i)) \end{aligned}$$

Log likelihood di atas dapat dimaksimalkan menggunakan **gradient ascent**. Sebelum itu, **gradient descent**, adalah algoritma yang meminimalkan nilai suatu fungsi, berkebalikan dengan gradient ascent yang memaksimalkan nilai suatu fungsi. Pada gradient descent, θ dipilih sedemikian hingga nilai suatu fungsi $J(\theta)$ minimal. Algoritma ini memilih θ secara acak kemudian secara berulang (iterasi) mengubah θ untuk meminimalkan nilai $J(\theta)$ hingga mencapai nilai konvergen untuk θ yang meminimalkan $J(\theta)$. Setiap iterasi dari algoritma ini memilih nilai θ_j yang baru, disebut *update*, yang ditulis:

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad (\text{II.21})$$

Update pada (II.21) dilakukan secara bersamaan untuk setiap nilai $j = 0, \dots, n$. Parameter α disebut **learning rate**, yang mengendalikan seberapa besar update dari θ_j .

Sebagai kontras, pada *gradient ascent*, *update* dari persamaan (II.21) ditulis sebagai:

$$\theta_j := \theta_j + \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad (\text{II.22})$$

Menggunakan persamaan (II.20) Turunan, atau **gradient** dari $\ell(\theta)$ adalah:

$$\begin{aligned} \frac{\partial}{\partial \theta_j} \ell(\theta) &= \left(y \frac{1}{g(\theta^T x)} - (1 - y) \frac{1}{1 - g(\theta^T x)} \right) \frac{\partial}{\partial \theta_j} g(\theta^T x) \\ &= \left(y \frac{1}{g(\theta^T x)} - (1 - y) \frac{1}{1 - g(\theta^T x)} \right) g(\theta^T x) (1 - g(\theta^T x)) \frac{\partial}{\partial \theta_j} \theta^T x \\ &= \left(y(1 - g(\theta^T x)) - (1 - y)g(\theta^T x) \right) x_j \\ &= (y - h_\theta(x))x_j \end{aligned}$$

Substitusi ke persamaan (II.22):

$$\theta_j := \theta_j + \alpha(y_i - h_\theta(x_i))x_{j_i}$$

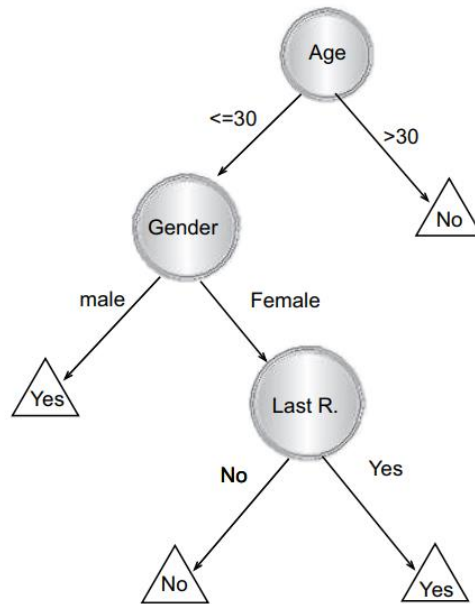
Diberikan input x testing, untuk membuat prediksi kelas, suatu nilai **threshold** r ditentukan antara 0 dan 1 sedemikian hingga $y = 1$ jika $H_\theta(x) \geq r$ dan $y = 0$ jika $H_\theta(x) < r$.

(Tsangaratos & Ilia, 2016; Hosmer & Lemeshow, 2000; Osisanwo, et al., 2017).

2.8.2 Decision Tree

Decision Tree atau pohon keputusan adalah pengklasifikasi yang menentukan output kelas dari suatu sampel berdasarkan keputusan yang diambil dari setiap nilai atribut dari sampel tersebut. Pohon keputusan adalah suatu fungsi Boolean di mana inputnya adalah suatu sampel E dengan vektor fitur X dan outputnya adalah 0 atau 1. Pada pohon keputusan, setiap node pohon yang bukan node daun adalah suatu uji

atribut atau suatu ekspresi boolean, setiap node daun adalah nilai Boolean, dan setiap cabang mewakili salah satu nilai yang mungkin dari atribut yang diuji.



Gambar 2: Decision Tree

Gambar 2 menunjukkan apakah seorang customer merespon surat langsung dari perusahaan atau tidak. Pada gambar di atas, node internal atau atribut disimbolkan sebagai lingkaran dan node daun digambarkan sebagai segitiga.

Jika suatu atribut bersifat kontinu, decision tree mengurutkan nilai kontinu tersebut dan memilih splitting point di antara nilai tersebut berdasarkan kelasnya.

Untuk membuat decision tree, banyak algoritma yang dapat digunakan, seperti ID3, C4.5, dan CART.

Algoritma ID3 menggunakan teori entropi informasi yang menghitung entropi dari masing-masing atribut. Kemudian menghitung Information Gain dari suatu atribut berdasarkan entropinya terhadap target kelas. Di mana **entropi** adalah ukuran ketidakpastian suatu atribut yang terkait dengan suatu kejadian. Semakin kecil entropi dari suatu atribut, semakin tinggi kemurnian informasi yang ada. Semakin besar entropi

dari suatu atribut, semakin besar ketidakpastian informasi tersebut. **Information Gain** adalah jumlah ketidakpastian informasi yang berkurang berdasarkan informasi yang diterima terkait dengan suatu atribut.

Misalkan S adalah data Training dengan jumlah sampel s , dengan m jumlah nilai atribut kelas yang berbeda dari $C_i (i = 1, 2, \dots, m)$. Misalkan S_i adalah jumlah sampel pada C_i , Information Gain minimal dari suatu sampel klasifikasi adalah:

$$I(S_1, \dots, S_m) = - \sum_{i=1}^m p_i \log_2(p_i) \quad (\text{II.23})$$

di mana p_i adalah peluang suatu sampel berasal dari kelas C_i yang diperoleh dari s_i/s . Suatu himpunan atribut A memiliki n nilai berbeda $\{a_1, \dots, a_n\}$. S dibagi menjadi n subset $\{S_1, \dots, S_n\}$, di mana S_j memiliki beberapa sampel di S yang memiliki nilai A_j di A .

Misalkan s_{ij} adalah jumlah sampel di kelas C_i dari subset S_j , entropi subset yang dibagi oleh A adalah:

$$E(A) = - \sum_{j=1}^n \frac{s_{1j} + s_{2j} + \dots + s_{mj}}{s} * I(s_{1j}, s_{2j}, \dots, s_{mj}) \quad (\text{II.24})$$

Berdasarkan persamaan (II.24) di atas, informasi yang dibutuhkan dari suatu subset S_j dihitung berdasarkan:

$$I(s_{1j}, s_{2j}, \dots, s_{mj}) = - \sum_{i=1}^m p_{ij} \log_2 p_{ij} \quad (\text{II.25})$$

Maka Information Gain dari A dapat dihitung dengan:

$$\text{Gain}(A) = I(s_1, s_2, \dots, s_m) - E(A) \quad (\text{II.26})$$

Gini index adalah criterion berbasis impurity yang mengukur divergensi antara distribusi peluang dari nilai atribut target, didefinisikan:

$$Gini(y, S) = \sum_{i=1}^m p_i^2$$

Algoritma ID3 membangun decision tree berdasarkan urutan atribut yang memiliki information gain tertinggi sebagai root node. Node tersebut kemudian bercabang sesuai dengan jumlah nilai berbeda pada atribut tersebut. Cabang dengan entropi nol adalah node daun, dan cabang dengan entropi yang lebih dari nol membutuhkan pecabangan lebih lanjut. Proses tersebut dilakukan secara rekursif pada semua node yang bukan merupakan node daun hingga seluruh data terklasifikasi

(Osisanwo, et al., 2017; Quinlan, 1986; Dai, Zhang, & Wu, 2016; Vafeiadis, Diamantaras, Sarigiannidis, & Chatzisavvas, 55).

2.8.3 Support Vector Machine

Support Vector Machine adalah algoritma pengklasifikasi yang membangun suatu hyperplane yang memisahkan data-data dengan kelas berbeda dengan margin sebesar mungkin.

Misalkan suatu data training pada sampel $(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$ dipisahkan menjadi dua kelas, di mana $x_i \in \mathbb{R}^n$ adalah vektor fitur dan $y_i \in \{-1, 1\}$ adalah label kelasnya. Jika diasumsikan dua kelas tersebut dapat dipisahkan dengan hyperplane $w \cdot x + b = 0$ pada suatu bidang H , dan informasi mengenai distribusi datanya tidak diketahui, maka hyperplane yang optimal adalah hyperplane yang memaksimalkan margin. Model SVM didefinisikan sebagai:

$$H_{w,b}(x) = g(w^T x + b)$$

di mana $g(z) = 1$ jika $z \geq 0$, dan $g(z) = -1$ jika $z < 0$.

Functional margin dari hyperplane (w, b) untuk x_i didefinisikan sebagai:

$$\hat{y}_i = y_i(w^T x_i + b) \quad (\text{II.27})$$

dan *functional margin* dari seluruh dataset didefinisikan sebagai:

$$\hat{\gamma} = \min_{i=1,\dots,m} \hat{\gamma}_i$$

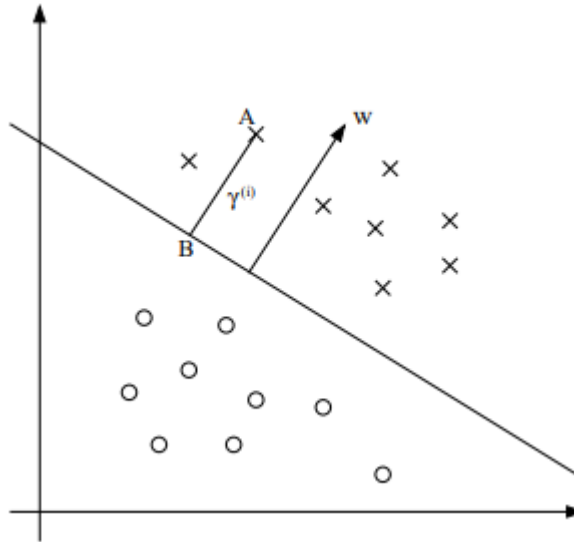
Geometric margin dari hyperplane (w, b) untuk x_i adalah *functional margin* yang dinormalisasikan dengan $\|w\|$:

$$\gamma_i = \frac{y_i(w^T x_i + b)}{\|w\|} \quad (\text{II.28})$$

dan geometric margin dari seluruh dataset didefinisikan sebagai:

$$\gamma = \min_{i=1,\dots,m} \gamma_i$$

Dari persamaan (II.27), jika $y_i = 1$, maka untuk membuat *functional margin* yang besar, $w^T x_i + b$ harus menghasilkan nilai yang besar. Sebaliknya, jika $y_i = -1$, maka untuk membuat *functional margin* yang besar, $w^T x_i + b$ harus menghasilkan nilai negatif yang besar.



Gambar 3: Geometric margin pada SVM

Pada Gambar 3, *decision boundary* atau *hyperplane* digambarkan sebagai garis diagonal pemisah antara kelas X ($y_i = 1$) dan O ($y_i = -1$). Vektor w adalah vektor yang ortogonal terhadap hyperplane. Titik pada A, merupakan suatu vektor x_i dari data

training dengan label $y_i = 1$. Jarak dari titik A dari hyperplane, atau γ_i , digambarkan pada garis AB. Untuk mencari nilai γ_i , $\frac{w}{\|w\|}$ adalah vektor satuan yang mempunyai arah yang sama dengan w . Karena A merepresentasikan x_i , maka titik B adalah $x_i - \gamma_i * \frac{w}{\|w\|}$ yang terletak tepat pada hyperplane. Seluruh titik yang terletak pada hyperplane memenuhi persamaan $w^T x + b = 0$. Maka:

$$w^T \left(x_i - \gamma_i \frac{w}{\|w\|} \right) + b = 0$$

atau dapat ditulis:

$$\begin{aligned} \gamma_i &= \frac{w^T x_i + b}{\|w\|} \\ \gamma_i &= \left(\frac{w}{\|w\|} \right)^T x_i + \frac{b}{\|w\|} \end{aligned} \quad (\text{II.29})$$

Namun persamaan (II.29) diperoleh dari titik A yang memiliki $y_i = 1$ kelas positif. Lebih umumnya, *geometric margin* dari (w, b) terhadap sampel training (x_i, y_i) didefinisikan:

$$\gamma_i = y_i \left(\left(\frac{w}{\|w\|} \right)^T x_i + \frac{b}{\|w\|} \right) \quad (\text{II.30})$$

yang ekuivalen dengan persamaan (II.28).

Jika diasumsikan training set yang diberikan dapat dipisah secara linear (*linearly separable*), *functional margin* tertinggi dapat direpresentasikan sebagai masalah optimisasi berikut:

$$\max_{\gamma, w, b} \gamma$$

dengan syarat:

$$\begin{aligned} y_i(w^T x_i + b) &\geq \gamma, \quad i = 1, \dots, m \\ \|w\| &= 1 \end{aligned}$$

yang berarti γ dimaksimalkan dengan syarat functional margin bernilai lebih besar atau sama dengan γ dan $\|w\|$ bernilai 1. $\|w\| = 1$ memastikan bahwa functional margin akan bernilai sama dengan *geometric margin* yang menjamin seluruh *geometric margin* bernilai lebih besar atau sama dengan γ . Namun karena $\|w\| = 1$ juga masalah optimisasi di atas bersifat *non-convex*, maka masalah optimisasi tersebut dapat transformasikan menjadi:

$$\max_{\hat{\gamma}, w, b} \frac{\hat{\gamma}}{\|w\|}$$

dengan syarat:

$$y_i(w^T x_i + b) \geq \hat{\gamma}, \quad i = 1, \dots, m$$

yang berarti $\frac{\hat{\gamma}}{\|w\|}$ dimaksimalkan dengan syarat seluruh *functional margin* lebih besar atau sama dengan $\hat{\gamma}$. Bentuk ini menghapus syarat $\|w\| = 1$, namun masalah optimisasi ini masih bersifat *non-convex* yang disebabkan oleh fungsi objektif $\frac{\hat{\gamma}}{\|w\|}$. Namun dengan memberi syarat bahwa *functional margin* dari (w, b) harus bernilai 1:

$$\hat{\gamma} = 1$$

Maka fungsi objektif $\frac{\hat{\gamma}}{\|w\|}$ akan berubah menjadi $\frac{1}{\|w\|}$ yang ekuivalen dengan meminimalkan $\|w\|^2$, masalah optimisasi yang baru dapat ditulis:

$$\min_{w, b} \frac{1}{2} \|w\|^2 \tag{II.31}$$

dengan syarat:

$$y_i(w^T x_i + b) \geq 1, \quad i = 1, \dots, m$$

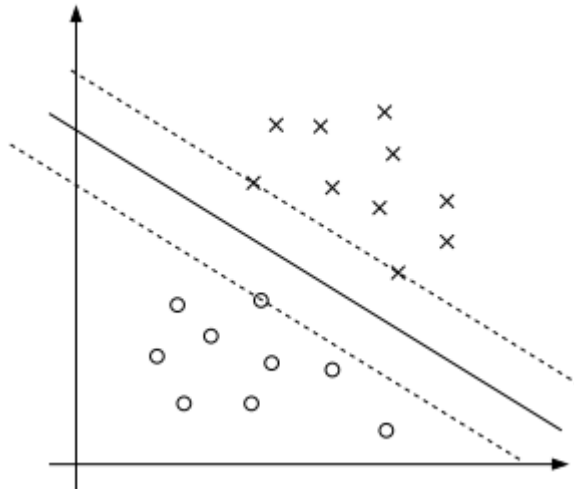
yang merupakan objective function kuadratik yang bersifat *convex* dengan syarat linear. Masalah optimisasi (II.31) dapat diselesaikan menggunakan *Quadratic Programming*. Dari masalah optimisasi (II.31), syarat optimisasinya dapat ditulis sebagai:

$$g_i(w) = -y_i(w^T x_i + b) + 1 \leq 0 \quad (\text{II.32})$$

Didefinisikan kondisi dual komplementer Karush-Kuhn-Tucker (KKT):

$$\alpha_i^* g_i(w^*) = 0, \quad i = 1, \dots, k \quad (\text{II.33})$$

Dari persamaan (II.33), $\alpha_i > 0$ hanya untuk sampel training yang memiliki nilai functional margin 1, atau $g_i(w) = 0$



Gambar 4: Margin dari Support Vector Machines

Pada Gambar 4, titik-titik dengan *margin* terkecil adalah titik-titik yang terdekat dengan *decision boundary*. Terdapat tiga titik yang berada tepat pada garis putus-putus yang paralel terhadap *decision boundary*. Jadi, hanya terdapat tiga α_i yang tidak bernilai nol untuk solusi optimal pada masalah optimisasi di atas.

Masalah optimisasi di atas akan menggunakan inner product $\langle x_i, x_j \rangle$ atau $x_i^T x_j$ antara titik-titik pada ruang input. Inner product ini berguna pada saat menggunakan kernel trick. Mentransformasikan ke bentuk Lagrange, masalah optimisasi (II.31) akan berbentuk:

$$\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i [y_i(w^T x_i + b) - 1] \quad (\text{II.34})$$

Untuk menemukan dual dari masalah di atas, $\mathcal{L}(w, b, \alpha)$ perlu diminimalkan dengan w dan b (dengan α tetap), dan mengubah turunan dari \mathcal{L} ke nol:

$$\nabla_w \mathcal{L}(w, b, \alpha) = w - \sum_{i=1}^m \alpha_i y_i x_i = 0$$

Yang berarti:

$$w = \sum_{i=1}^m \alpha_i y_i x_i \quad (\text{II.35})$$

Dengan menurunkan terhadap b , diperoleh:

$$\frac{\partial}{\partial b} \mathcal{L}(w, b, \alpha) = \sum_{i=1}^m \alpha_i y_i = 0 \quad (\text{II.36})$$

Kemudian w dari persamaan (II.35) disubstitusi ke persamaan (II.34), diperoleh:

$$\mathcal{L}(w, b, \alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y_i y_j \alpha_i \alpha_j x_i^T x_j - b \sum_{i=1}^m \alpha_i y_i$$

Namun dari (II.36), suku terakhir bernilai nol:

$$\mathcal{L}(w, b, \alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y_i y_j \alpha_i \alpha_j x_i^T x_j$$

Persamaan di atas diperoleh dari meminimalkan $\mathcal{L}(w, b, \alpha)$ dengan w dan b . Menggabungkan dengan syarat koefisien Lagrange $\alpha_i \geq 0$ dan (II.36), diperoleh masalah optimisasi dual:

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y_i y_j \alpha_i \alpha_j \langle x_i, x_j \rangle$$

dengan syarat:

$$\alpha_i \geq 0, \quad i = 1, \dots, m$$

$$\sum_{i=1}^m \alpha_i y_i = 0$$

Masalah optimisasi di atas dapat diselesaikan dengan menggunakan persamaan (II.36), memperoleh w yang optimal sebagai fungsi dari α . Setelah memperoleh w^* , mencari nilai optimal untuk b diperoleh dari:

$$b^* = -\frac{\max_{i:y_i=-1} w^{*T} x_i + \min_{i:y_i=1} w^{*T} x_i}{2}$$

Pada persamaan (II.36), yang menghasilkan nilai w optimal dari α , misalkan model yang diperoleh digunakan untuk memprediksi titik x baru. $w^T x + b$ dikalkulasi dan diprediksi $y = 1$ jika dan hanya jika $w^T x + b > 0$. Namun dengan menggunakan persamaan (II.36), $w^T x + b$ dapat ditulis sebagai:

$$w^T x + b = \left(\sum_{i=1}^m \alpha_i y_i x_i \right)^T x + b \quad (\text{II.37})$$

$$w^T x + b = \sum_{i=1}^m \alpha_i y_i \langle x_i, x \rangle + b \quad (\text{II.38})$$

Maka dari itu, ketika seluruh α_i telah diperoleh, untuk membuat prediksi hanya perlu menghitung inner product antara x dengan seluruh sampel training. Namun, dari apa yang dilihat sebelumnya, terdapat hanya beberapa α_i yang tidak bernilai nol. Jadi, banyak iterasi dari (II.38) yang bernilai nol. Dan pada akhirnya model hanya perlu menghitung inner product antara x dengan support vector saja (bukan dari seluruh sampel training).

Untuk membuat hyperplane, terkadang data training yang diberikan sulit untuk dipisahkan. Data training yang diberikan perlu untuk ditransformasikan ke suatu

bentuk lain dalam ruang yang berbeda. Fungsi yang memetakan vektor input ke suatu vektor lain disebut **Kernel**. Sebagai contoh, diberikan fungsi

$$\phi(x) = \begin{bmatrix} x \\ x^2 \\ x^3 \end{bmatrix}$$

Pada beberapa kasus, di mana input x tidak dapat dipisahkan secara linear, input tersebut perlu ditransformasikan menggunakan fungsi $\phi(x)$. Pada persamaan (II.38) seluruh x diganti dengan $\phi(x)$. Karena algoritma SVM ini hanya perlu menghitung *inner product* $\langle x, z \rangle$, maka inner product tersebut berubah menjadi $\langle \phi(x), \phi(z) \rangle$. Lebih spesifik, diberikan fungsi pemetaan ϕ , suatu **Kernel** didefinisikan sebagai:

$$K(x, z) = \phi(x)^T \phi(z) \quad (\text{II.39})$$

Maka dari itu, seluruh *inner product* $\langle x, z \rangle$ pada persamaan (II.38) dapat disubstitusi menjadi persamaan (II.39), dan dengan mengganti $w^T x + b$ menjadi fungsi f yang menerima input x , persamaan (II.38) dapat ditulis menjadi:

$$f(x) = \sum_{i=1}^m \alpha_i y_i K(x_i, x) + b \quad (\text{II.40})$$

di mana α_i dan b diperoleh dari hasil optimisasi di atas, dan α_i tidak nol adalah *support vector*.

Untuk $K(x_i, x)$ adalah *inner product* atau hasil kali dalam antara dua observasi $i \in \{1, \dots, n\}$ dan $i' \in \{1, \dots, n\}$ yang diekspresikan:

$$K(x_i, x_{i'}) = \sum_{j=1}^p x_{ij} x_{i'j} \quad (\text{II.41})$$

yang disebut kernel Linear.

Metode SVM yang menggunakan kernel linear dapat memisahkan kelas-kelas secara linear. Namun secara praktis dataset yang dijumpai sering tidak dapat dipisahkan secara linear. Kernel non-linear lebih fleksibel dikarenakan kernel ini memetakan

variabel p dan x ke dimensi yang lebih tinggi. Salah satu kernel non-linear yang populer digunakan adalah kernel radial atau kernel Gauss:

$$K(x_i, x_{i'}) = \exp \left\{ -\gamma \sum_{j=1}^p (x_{ij} - x_{i'j})^2 \right\} \quad (\text{II.42})$$

di mana $\gamma > 0$ adalah parameter tuning tambahan. Ketika suatu sampel test terletak jauh dari sampel training, eksponensial tersebut akan menjadi negatif, dan $K(x_i, x_{i'})$ mendekati nol. Pada kasus ini, sampel training yang terletak jauh dari sampel test hampir tidak memiliki efek atau kontribusi kepada keputusan yang digunakan untuk mengklasifikasikan sampel test. Parameter γ mempengaruhi seberapa jauh letak suatu observasi agar dapat berkontribusi terhadap keputusan klasifikasi. Kernel ini lebih fleksibel dari kernel linear (Guenther & Schonlau, 2016; Schudt, Laptev, & Caputo, 2004; Vapnik & Cortes, 1995; Andrew, 2000).

2.8.4 Multilayer Perceptron (MLP)

Artificial Neural Network (ANN) atau Jaringan Saraf Tiruan adalah salah satu algoritma machine learning yang mencari suatu fungsi f yang tidak diketahui berdasarkan input vektor X ke suatu output y :

$$y = f(X)$$

Pada tahap pelatihan, fungsi f dioptimalkan sedemikian sehingga hasil yang didapat untuk setiap vektor X yang diberikan sedekat mungkin dengan nilai y .

Istilah neuron atau node digunakan pada ANN untuk menandakan suatu atribut (pada layer input), suatu jenis output (pada layer output) atau suatu ekspresi matematika (pada layer tersembunyi). Neuron atau node ini masing-masing memiliki *weight* (bobot) yang berbeda-beda yang mempengaruhi nilai output dari node tersebut berdasarkan suatu fungsi aktivasi. Untuk node pada layer non-output, output node tersebut akan dikirim ke seluruh node pada layer selanjutnya. Koneksi atau sambungan antara node memiliki nilai bias.

Misalkan diberikan suatu input fitur x_1, x_2, x_3 yang secara kolektif disebut *input layer*, empat unit tersembunyi yang secara kolektif disebut *hidden layer* dan satu output neuron yang disebut *output layer*. Unit tersembunyi yang pertama membutuhkan input x_1, x_2, x_3 dan memberikan output a_1 (a untuk nilai “aktivasi”). Pada contoh ini, hanya terdapat satu hidden layer, namun hidden layer dapat berjumlah lebih dari satu. Misalkan $a_1^{[1]}$ menyatakan nilai output dari hidden unit pertama pada hidden layer pertama. Index berbasis 0 digunakan untuk penomoran index dari layer. Input layer adalah layer 0, hidden layer pertama adalah layer 1, dan output layer adalah layer 2. Dengan notasi matematika, output dari layer 2 adalah $a_1^{[2]}$. Notasi ini dapat disatukan:

$$\begin{aligned}x_1 &= a_1^{[0]} \\x_2 &= a_2^{[0]} \\x_3 &= a_3^{[0]}\end{aligned}$$

di mana $a_j^{[l]}$ menyatakan nilai output dari unit ke j pada layer ke l .

Unit pertama pada hidden layer melakukan komputasi:

$$z_1^{[1]} = W_1^{[1]T} x + b_1^{[1]} \text{ dan } a_1^{[1]} = g(z_1^{[1]})$$

di mana W adalah matrix parameter dan W_1 menyatakan baris pertama dari matrix W . Parameter yang terkait dengan hidden unit pertama adalah vektor $W_1^{[1]} \in \mathbb{R}^3$ dan skalar $b_1^{[1]} \in \mathbb{R}$.

Multilayer Perceptron (MLP) adalah salah satu jenis Artificial Neural Network yang memiliki layer tersembunyi atau hidden layer. MLP terdiri dari setidaknya 3 layer, yaitu input layer, hidden layer dan output layer. Secara praktis, MLP biasanya memiliki lebih dari satu hidden layer. Unit tersembunyi ke dua dan ketiga dari hidden layer pertama, komputasinya didefinisikan sebagai:

$$\begin{aligned}z_2^{[1]} &= W_2^{[1]T} x + b_2^{[1]} \text{ dan } a_2^{[1]} = g(z_2^{[1]}) \\z_3^{[1]} &= W_3^{[1]T} x + b_3^{[1]} \text{ dan } a_3^{[1]} = g(z_3^{[1]})\end{aligned}$$

di mana setiap unit tersembunyi masing-masing memiliki parameter W dan b tersendiri. Output layer melakukan komputasi:

$$z_1^{[2]} = W_1^{[2]T} a^{[1]} + b_1^{[2]} \text{ dan } a_1^{[2]} = g(z_1^{[2]})$$

di mana $a^{[1]}$ didefinisikan sebagai vektor dari seluruh nilai aktivasi dari layer 1:

$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix}$$

Nilai aktivasi $a_1^{[2]}$ dari layer 2, yang berupa nilai skalar yang didefinisikan oleh $a_1^{[2]} = g(z_1^{[2]})$, melambangkan prediksi nilai output terakhir dari neural network tersebut.

Menggunakan matrix aljabar, nilai aktivasi dapat dikomputasi:

$$\underbrace{\begin{bmatrix} z_1^{[1]} \\ \vdots \\ z_4^{[1]} \end{bmatrix}}_{z^{[1]} \in \mathbb{R}^{4 \times 1}} = \underbrace{\begin{bmatrix} \dots & W_1^{[1]T} & \dots \\ \dots & W_2^{[1]T} & \dots \\ \dots & \vdots & \dots \\ \dots & W_4^{[1]T} & \dots \end{bmatrix}}_{W^{[1]} \in \mathbb{R}^{4 \times 3}} \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}}_{x \in \mathbb{R}^{3 \times 1}} + \underbrace{\begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ \vdots \\ b_4^{[1]} \end{bmatrix}}_{b^{[1]} \in \mathbb{R}^{4 \times 1}}$$

Di mana $\mathbb{R}^{d \times n}$ mengindikasikan dimensi matrix.

Berdasarkan operasi matrix di atas, diberikan suatu input $x \in \mathbb{R}^3$, nilai aktivasi dari hidden layer dikomputasikan dengan $z^{[1]} = W^{[1]}x + b^{[1]}$ dan $a^{[1]} = g(z^{[1]})$. Sedangkan untuk nilai aktivasi dari layer output dikomputasi:

$$\underbrace{z^{[2]}}_{1 \times 1} = \underbrace{W^{[2]}}_{1 \times 4} \underbrace{a^{[1]}}_{4 \times 1} + \underbrace{b^{[2]}}_{1 \times 1} \text{ dan } \underbrace{a^{[2]}}_{1 \times 1} = g(\underbrace{z^{[2]}}_{1 \times 1})$$

Jika diberikan suatu training set dengan tiga sampel. Nilai aktivasi dari setiap sampel adalah:

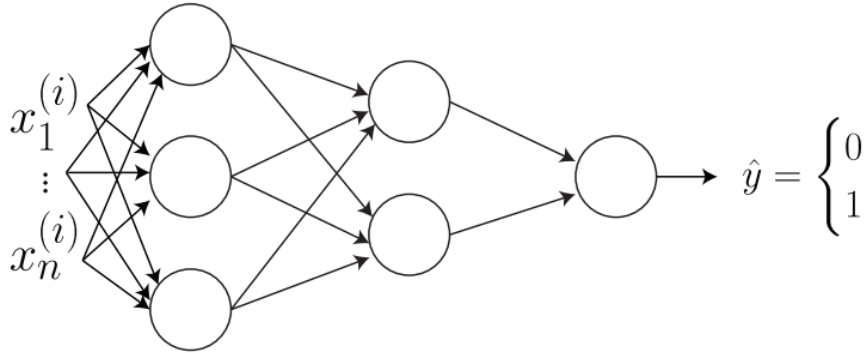
$$z^{1} = W^{[1]}x^{(1)} + b^{[1]}$$

$$z^{[1](2)} = W^{[1]}x^{(2)} + b^{[1]}$$

$$z^{[1](3)} = W^{[1]}x^{(3)} + b^{[1]}$$

Dengan $[i]$ melambangkan layer ke- i dan (j) melambangkan sampel training ke- j .

Dilihat dari Gambar 5, seluruh input adalah vektor $x_1^{(i)}, \dots, x_n^{(i)}$. Pada hidden layer pertama, seluruh input tersambung ke seluruh neuron pada layer berikutnya. Layer ini disebut *fully connected layer*.



Gambar 5: Neural Network

Langkah berikutnya adalah menghitung berapa banyak parameter pada jaringan ini:

$$z^{[1]} = W^{[1]}x^{(i)} + b^{[1]}$$

$$a^{[1]} = g(z^{[1]})$$

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

$$a^{[2]} = g(z^{[2]})$$

$$z^{[3]} = W^{[3]}a^{[2]} + b^{[3]}$$

$$\hat{y}^{(i)} = a^{[3]} = g(z^{[3]})$$

Diketahui bahwa $z^{[1]}, a^{[1]} \in \mathbb{R}^{3 \times 1}$ dan $z^{[2]}, a^{[2]} \in \mathbb{R}^{2 \times 1}$ dan $z^{[3]}, a^{[3]} \in \mathbb{R}^{1 \times 1}$. Untuk $W^{[1]}$ dapat dideduksi dari:

$$z^{[1]} = W^{[1]}x^{(i)}$$

yang size-nya dapat ditulis:

$$\mathbb{R}^{3 \times 1} = \mathbb{R}^{2 \times 3} \times \mathbb{R}^{n \times 1}$$

Menggunakan perkalian matriks, dapat disimpulkan bahwa size $W^{[1]}$ adalah $3 \times n$. Dengan cara yang serupa untuk setiap hidden layer. Disimpulkan:

$$W^{[2]} \in \mathbb{R}^{2 \times 3}, b^{[2]} \in \mathbb{R}^{3 \times 1} \text{ dan } W^{[3]} \in \mathbb{R}^{3 \times 1}, b^{[3]} \in \mathbb{R}^{3 \times 1}$$

Secara keseluruhan, ada $3n + 3$ pada layer pertama, $2 * 3 + 2$ pada layer kedua dan $2 + 1$ pada layer terakhir, dengan total $3n + 14$ parameter. Parameter-parameter tersebut diinisialisasikan secara acak dengan distribusi normal sekitar nol; $N(0, 0.1)$. Kemudian parameter-parameter ini diestimasi menggunakan *gradient descent*.

Untuk memperbarui (*update*) parameter-parameter ini, input vektor harus melewati seluruh neuron pada neural network dan outputnya dinotasikan dengan suatu nilai prediksi \hat{y} . Error atau loss dari \hat{y} , dinotasikan \mathcal{L} , dapat dihitung:

$$\mathcal{L}(\hat{y}, y) = -[(1 - y) \log(1 - \hat{y}) + y \log \hat{y}]$$

Fungsi *loss* $\mathcal{L}(\hat{y}, y)$ menghasilkan suatu nilai skalar. Seluruh parameter pada setiap layer di neural network diperbaharui berdasarkan nilai ini. Untuk setiap layer pada index ℓ , perbaharui parameter:

$$W^{[\ell]} = W^{[\ell]} - \alpha \frac{\partial \mathcal{L}}{\partial W^{[\ell]}}$$

$$b^{[\ell]} = b^{[\ell]} - \alpha \frac{\partial \mathcal{L}}{\partial b^{[\ell]}}$$

di mana α adalah learning rate.

Pada tahap optimisasi, turunan (*gradient*) untuk $W^{[3]}$ diutamakan untuk dikomputasi terlebih dahulu. Hal ini karena pengaruh $W^{[1]}$ terhadap *loss* lebih kompleks dari $W^{[3]}$ yang terletak lebih dekat ke output \hat{y} , dalam hal jumlah komputasi. Dengan menggunakan sifat turunan sigmoid pada persamaan (II.21), turunan $W^{[3]}$ adalah:

$$\begin{aligned}
\frac{\partial \mathcal{L}}{\partial W^{[3]}} &= -\frac{\partial}{\partial W^{[3]}} ((1-y) \log(1-\hat{y}) + y \log \hat{y}) \\
&= -(1-y) \frac{\partial}{\partial W^{[3]}} \log(1 - g(W^{[3]}a^{[2]} + b^{[3]})) \\
&\quad - y \frac{\partial}{\partial W^{[3]}} \log(g(W^{[3]}a^{[2]} + b^{[3]})) \\
&= -(1-y) \frac{1}{1 - g(W^{[3]}a^{[2]} + b^{[3]})} (-1) g'(W^{[3]}a^{[2]} + b^{[3]}) a^{[2]T} \\
&\quad - y \frac{1}{g(W^{[3]}a^{[2]} + b^{[3]})} g'(W^{[3]}a^{[2]} + b^{[3]}) a^{[2]T} \\
&= (1-y) \sigma(W^{[3]}a^{[2]} + b^{[3]}) a^{[2]T} \\
&\quad - y(1 - \sigma(W^{[3]}a^{[2]} + b^{[3]})) a^{[2]T} \\
&= (1-y) a^{[3]} a^{[2]T} - y(1 - a^{[3]}) a^{[2]T} \\
&= (a^{[3]} - y) a^{[2]T}
\end{aligned} \tag{II.43}$$

dengan $a^{[3]} = \sigma(W^{[3]}a^{[2]} + b^{[3]})$.

Turunan $W^{[2]}$ dapat dihitung melalui aturan rantai (*chain rule*). Diketahui bahwa \mathcal{L} bergantung pada $\hat{y} = a^{[3]}$:

$$\frac{\partial \mathcal{L}}{\partial W^{[2]}} = \frac{\partial \mathcal{L}}{\partial a^{[3]}} \frac{\partial a^{[3]}}{\partial W^{[2]}}$$

$\partial a^{[3]}/\partial a^{[3]}$ dapat disisipkan:

$$\frac{\partial \mathcal{L}}{\partial W^{[2]}} = \frac{\partial \mathcal{L}}{\partial a^{[3]}} \frac{\partial a^{[3]}}{\partial z^{[3]}} \frac{\partial z^{[3]}}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial W^{[2]}}$$

lebih lanjut:

$$\frac{\partial \mathcal{L}}{\partial W^{[2]}} = \frac{\partial \mathcal{L}}{\partial a^{[3]}} \frac{\partial a^{[3]}}{\partial z^{[3]}} \frac{\partial z^{[3]}}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial W^{[2]}} \tag{II.44}$$

Dengan menggunakan persamaan (II.43), persamaan (II.44) menjadi:

$$\begin{aligned}
\frac{\partial \mathcal{L}}{\partial W^{[2]}} &= \frac{\partial \mathcal{L}}{\partial a^{[3]} - y} \underbrace{\frac{\partial a^{[3]}}{\partial z^{[3]}}}_{W^{[3]}} \underbrace{\frac{\partial z^{[3]}}{\partial a^{[2]}}}_{g'(z^{[2]})} \underbrace{\frac{\partial a^{[2]}}{\partial z^2}}_{a^{[1]}} \frac{\partial z^2}{\partial W^{[2]}} \\
&= (a^{[3]} - y) W^{[3]} g'(z^{[2]}) a^{[1]} \\
&= W^{[3]T} \cdot g'(z^{[2]}) (a^{[3]} - y) a^{[1]T}
\end{aligned} \tag{II.45}$$

Parameter $W^{[1]}$ dapat dihitung melalui aturan rantai menggunakan persamaan (II.43) dan persamaan (II.45).

Dengan seluruh parameter telah diestimasi oleh gradient descent, model akhir MLP dapat diekspresikan sebagai berikut:

$$\hat{y} = v_0 + \sum_{j=1}^{NH} b_j g(W_j^T x') \tag{II.46}$$

di mana x' adalah vektor input x dengan $x' = (1, x^T)^T$. Fungsi g adalah fungsi aktivasi pada node tersembunyi. Fungsi aktivasi adalah fungsi yang memberi output dari suatu node berdasarkan input yang diberikan untuk diteruskan ke node selanjutnya.

Fungsi aktivasi Sigmoid logistik adalah fungsi aktivasi yang memberi nilai antara 0 dan 1:

$$f(x) = \frac{1}{1 + e^{-x}} \tag{II.47}$$

Fungsi aktivasi Sigmoid tangensial (tanh) adalah fungsi aktivasi yang memberi nilai antara -1 dan 1

$$f(x) = \frac{2}{1 + e^{-2x}} - 1 \tag{II.48}$$

Fungsi aktivasi relu adalah fungsi aktivasi yang mengembalikan nilai nol untuk seluruh input negatif namun untuk seluruh input positif x , ReLU mengembalikan nilai x itu sendiri.

$$f(x) = f(x) = \max(0, x) \quad (\text{II.49})$$

(Ahmed, Atiya, Gayar, & El-Shishiny, 2010; Zurada, 1992; Haykin, 1994; Zhou & Liu, 2006; Bishop, Pattern Recognition and Machine Learning, 2006; Bishop, Neural Networks for Pattern Recognition, 1995; MULTILAYER, 1998; Ng & Katanforoosh, 2018)

2.8.5 K-Nearest Neighbor

K-Nearest Neighbor adalah algoritma prediksi non-parametrik di mana hasil prediksi kelas dari suatu titik didasarkan oleh mayoritas kelas dari k tetangga terdekatnya. Diberikan suatu titik, hitung jarak Euklid sesuai dengan persamaan (II.4) antara titik tersebut dengan semua titik pada data training. Kemudian pilih k tetangga terdekat berdasarkan jarak Euklid tersebut, prediksi kelas dari titik tersebut adalah modus kelas dari k tetangga terdekatnya.

Untuk suatu sampel x dengan k tetangga terdekat dari B , prediksi kelas ditentukan oleh:

$$\hat{y} = \underset{r}{\operatorname{argmax}} \left(\sum_{i=1}^k y_i (y_i = r) \right) \quad (\text{II.50})$$

di mana y_i adalah label kelas dari x_i (Ahmed, Atiya, Gayar, & El-Shishiny, 2010; Osisanwo, et al., 2017).

Weighted-KNN, atau KNN berbobot adalah jenis KNN yang memberi bobot w kepada k tetangga terdekat berdasarkan invers jarak d . Prediksi kelas dipilih berdasarkan jumlah bobot tertinggi dari tiap kelas yang ada, diekspresikan:

$$\hat{y} = \underset{r}{\operatorname{argmax}} \left(\sum_{i=1}^k w_i * I(y_i = r) \right) \quad (\text{II.51})$$

(Hechenbichler, 2004)

BAB III

METODE PENELITIAN

3.1 Waktu dan Tempat

Penelitian ini dilaksanakan dari bulan Agustus 2019 sampai dengan bulan Oktober 2019. Lokasi penelitian dilakukan di Laboratorium Rekayasa Perangkat Lunak Fakultas Matematika dan Ilmu Pengetahuan Alam, Universitas Hasanuddin.

3.2 Tahapan Penelitian

Untuk menyelesaikan penelitian ini, peneliti akan melewati beberapa tahap penelitian, yaitu: Pra-penelitian, eksplorasi dan preprocessing data, model *tuning and fitting*, dan analisis hasil.

Pada tahap pra-penelitian, peneliti menentukan tema penelitian, masalah yang akan diteliti, mengumpulkan sumber referensi atau literatur seperti jurnal dan buku yang mendukung dalam penelitian, dan menentukan metode yang digunakan beserta batasan masalahnya. Kemudian peneliti mencari data yang sesuai dengan tema penelitian sebagai objek penelitian.

Pada tahap eksplorasi dan preprocessing data, peneliti mencoba menguraikan karakteristik-karakteristik setiap dataset sebagai informasi untuk mengambil keputusan pada tahap preprocessing. Peneliti akan mengidentifikasi masalah-masalah yang terdapat pada dataset tersebut kemudian mengambil pendekatan untuk menyelesaikan masalah yang terkait. Normalisasi dan attribute reduction termasuk pada tahap ini.

Pada tahap model *tuning and fitting*, peneliti akan mencari parameter-parameter terbaik untuk model yang akan digunakan berdasarkan hasil eksplorasi data dan *trial and error* untuk mendapatkan hasil terbaik. Tuning juga dilakukan terhadap beberapa teknik resampling yang membutuhkan parameter. Kemudian model akan memberi hasil prediksi yang akan dianalisis pada tahap selanjutnya.

Pada tahap analisis hasil, peneliti akan merangkum hasil yang diperoleh dari metode-metode yang digunakan ke dalam bentuk tabel dan diagram, kemudian menyimpulkan hasilnya sebagai output dari penelitian ini.

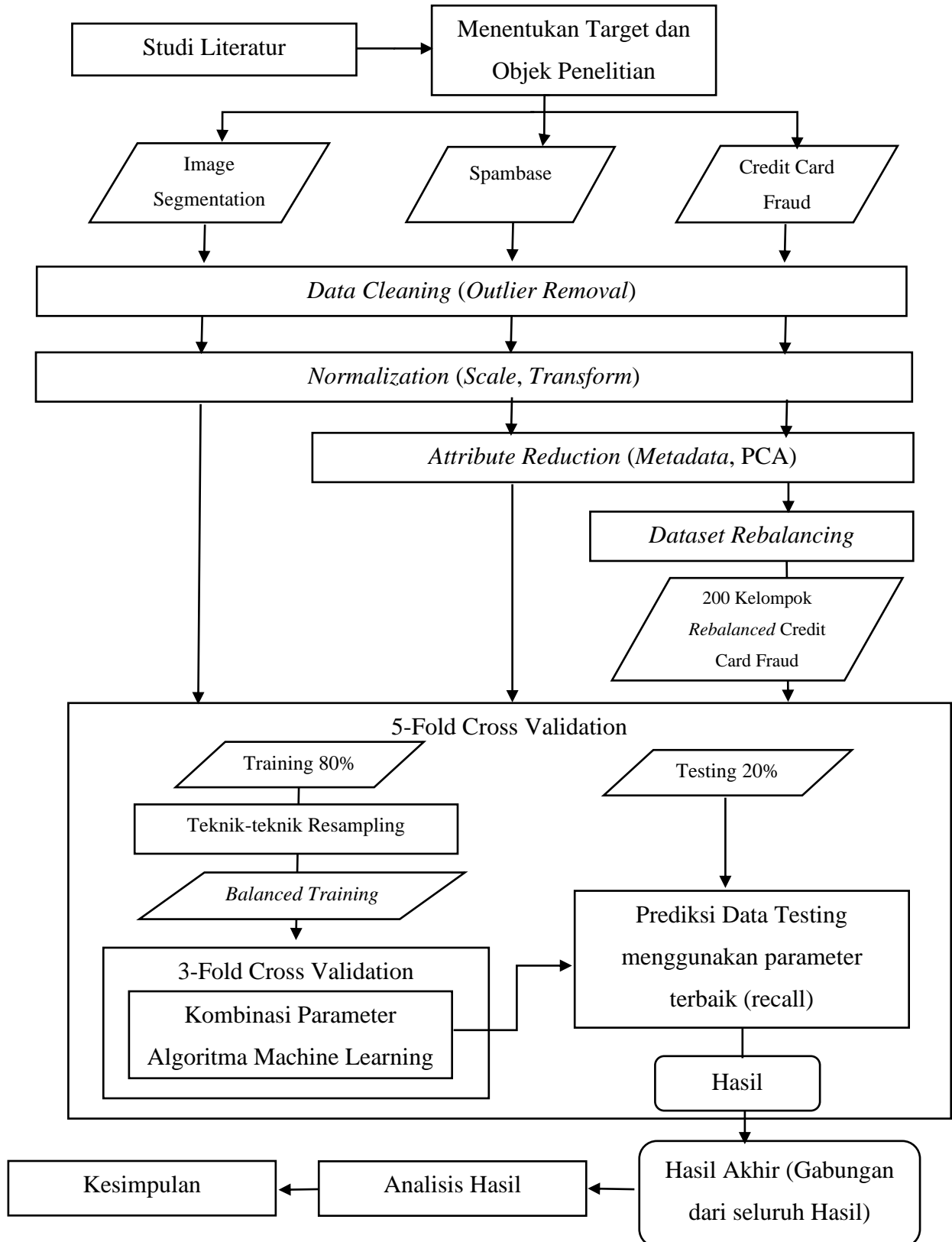
3.3 Deskripsi Data

Data diambil dari Website resmi Kaggle (kaggle.com), UCI Machine Learning Repository (archive.ics.uci.edu/ml/) dan KEEL (sci2s.ugr.es/keel/imbalanced.php). Data tersebut berupa tiga dataset, yaitu:

1. Credit Card Fraud Dataset (Kaggle), yang terdiri dari 30 kolom atribut dengan 1 kolom kelas, 284.807 baris. 284.315 jumlah sampel kelas mayoritas dan 492 jumlah sampel kelas minoritas dengan *imbalanced ratio* sebesar 577:1. Dataset ini merupakan dataset terpopuler di Kaggle sebab jumlah data yang besar dengan imbalanced ratio yang sangat tinggi.
2. Spambase Dataset (UCI), yang terdiri dari 57 kolom atribut dengan 1 kolom kelas, 4.601 baris. 2788 jumlah sampel kelas mayoritas dan 1813 jumlah sampel kelas minoritas dengan *imbalanced ratio* sebesar 1,5:1
3. Image Segmentation Dataset (KEEL), yang terdiri dari 19 kolom atribut dengan 1 kolom kelas, 2308 baris. 1962 jumlah sampel kelas mayoritas dan 346 jumlah sampel kelas minoritas dengan *imbalanced ratio* sebesar 6:1.

Seluruh dataset hanya memiliki atribut kontinu dengan label kelas biner yang sesuai dengan tema penelitian dan metode-metode yang digunakan.

3.4 Alur Penelitian



BAB IV

HASIL DAN PEMBAHASAN

4.1 Eksplorasi dan Preprocessing Data

4.1.1 Image Segmentation Dataset

4.1.1.1 Data Cleaning & Normalization

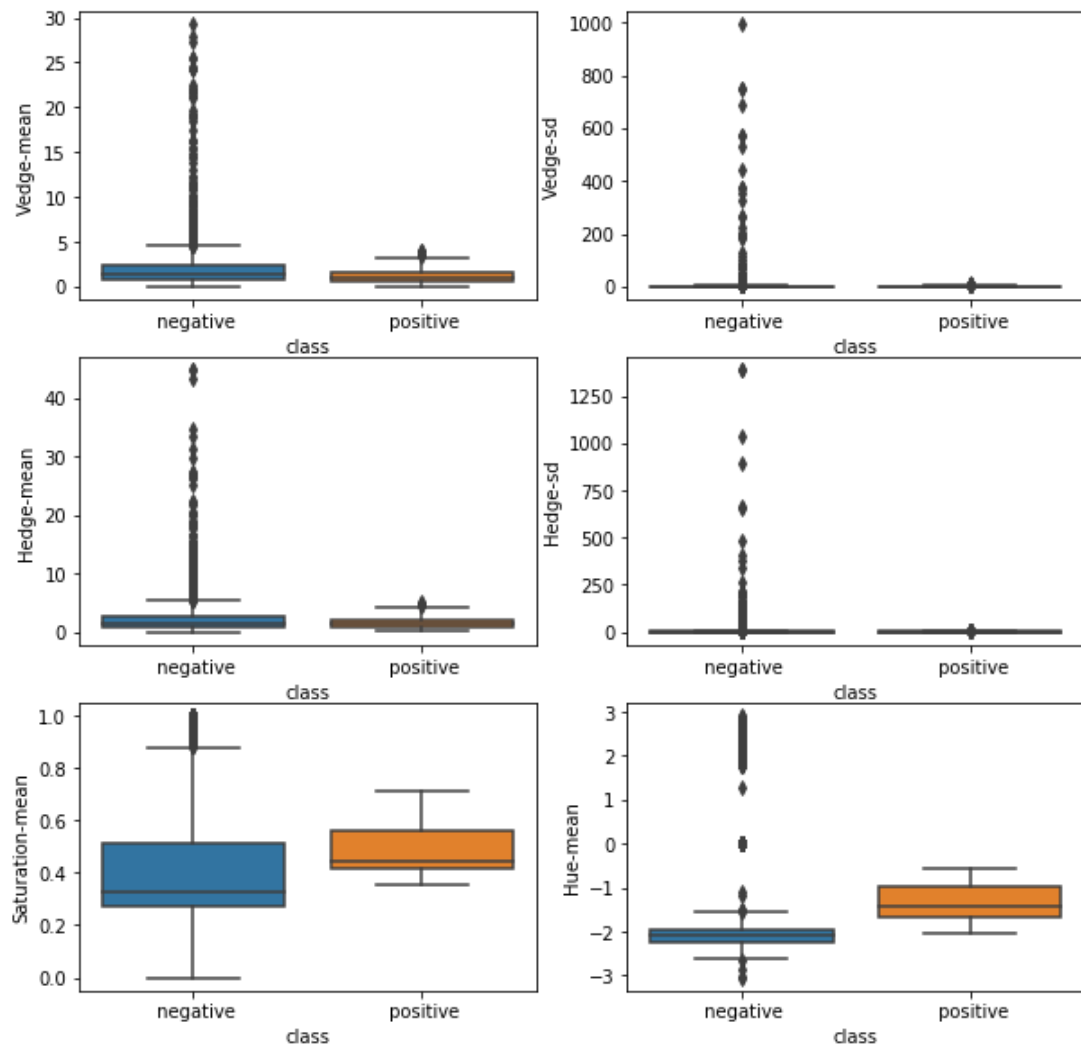
Image Segmentation dataset merupakan dataset mengenai citra di mana setiap pikselnya memiliki kelas berdasarkan hasil segmentasi manual dari citra outdoor. Setiap observasi dari dataset ini adalah gabungan piksel yang berukuran 3x3 (disebut region). Pada dataset asli terhadap 6 kelas berbeda. Namun dataset ini diubah oleh Keel di mana hanya terdapat 2 kelas, yaitu positif dan negatif. Kelas positif adalah kelas 0 pada dataset asli, dan kelas negatif adalah kelas 1, 2, 3, 4 dan 5 pada dataset asli. Perubahan ini dilakukan untuk menciptakan dataset yang tidak seimbang. Terdapat 20 atribut pada dataset ini:

1. 'region_centroid_row' dan 'region_centroid_col' menunjukkan lokasi centroid dari region observasi. region_centroid_row menunjukkan lokasi baris dari centroid, dan region_centroid_col menunjukkan lokasi kolom dari centroid. Atribut-atribut ini bersifat bilangan bulat [0, 255].
2. 'region_pixel_count' berisi informasi tentang jumlah piksel dari region observasi. Atribut ini bersifat bilangan bulat dan hanya terdapat satu jenis nilai pada atribut ini [9].
3. 'short_line_density-5' dan 'short_line_density-2' berisi informasi mengenai hasil dari algoritma ekstraksi garis yang melewati region tersebut. short_line_density-5 menghitung garis dengan kontras rendah (lebih kecil dari atau sama dengan 5) dan short_line_density-2 menghitung garis dengan kontras tinggi (lebih besar dari 5). Atribut-atribut ini bersifat riil kontinu [0, 1]
4. 'vedge_mean', 'vedge_sd', 'hedge_mean', 'hedge_sd' berisi informasi mengenai piksel-piksel tetangga yang berseblahan dengan region tersebut.

Terdapat 6 piksel pada masing-masing vedge dan hedge, yaitu 6 piksel tetangga secara horizontal dan 6 piksel tetangga secara vertikal. Rata-rata dari 6 piksel tersebut dimasukkan dalam vedge_mean dan hedge_mean, kemudian standar deviasi dari 6 piksel tersebut dimasukkan dalam vedge_sd dan hedge_sd. Atribut-atribut ini bersifat riil kontinu [0,...]

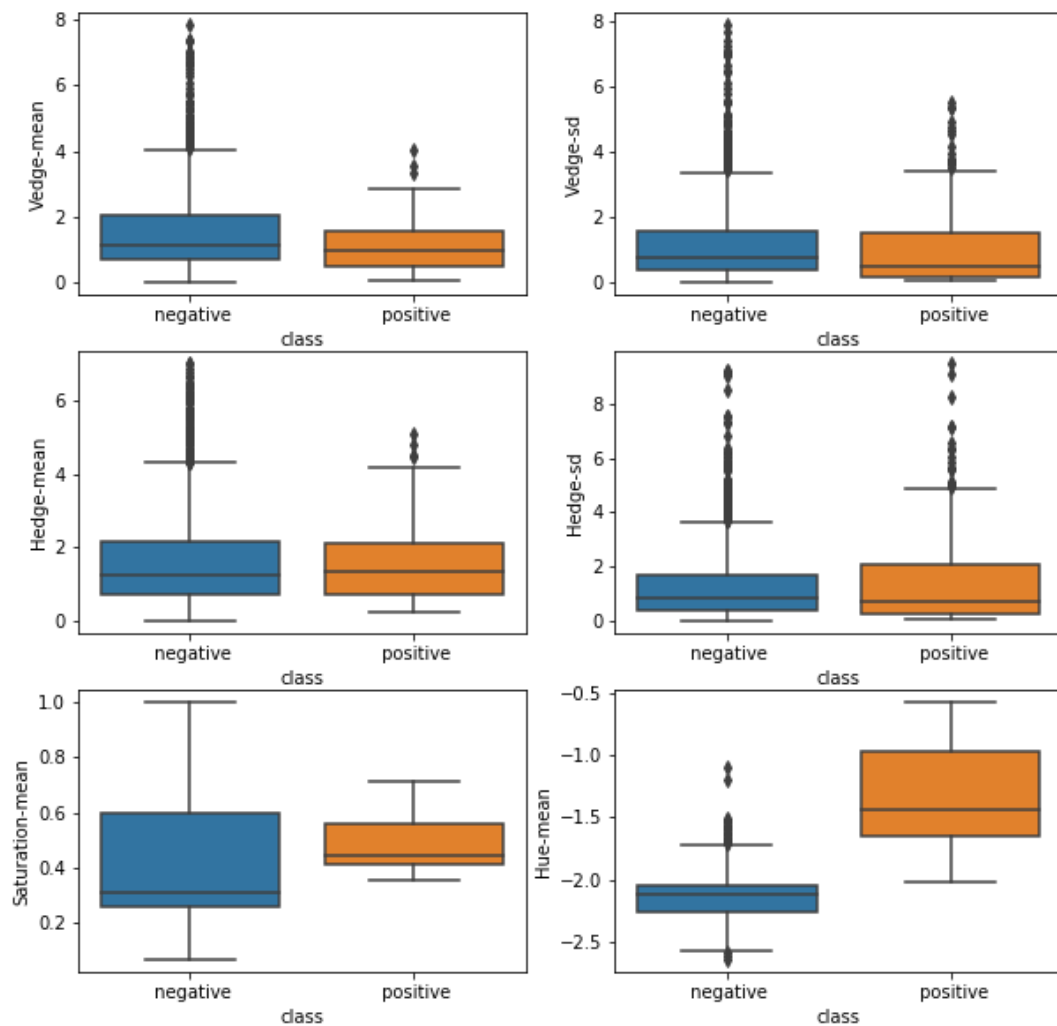
5. 'intensity_mean' adalah rata-rata dari seluruh warna pada region tersebut, dikalkulasi dari $(R+G+B)/3$. Atribut ini bersifat riil kontinu [0, 255]
6. 'rawred_mean', 'rawgreen_mean', dan 'rawblue_mean', berisi informasi mengenai rata-rata dari masing-masing warna merah, hijau, dan biru dari region tersebut. Atribut-atribut ini bersifat riil kontinu [0, 255]
7. 'exred_mean', 'exgreen_mean', dan 'exblue_mean', berisi informasi mengenai kelebihan (excess) dari masing-masing warna merah, hijau, dan biru. exred_mean dikalkulasi dari $2R-(G+B)$, exgreen_mean dikalkulasi dari $2G-(R+B)$, dan exblue_mean dikalkulasi dari $2B-(R+G)$. Atribut-atribut ini bersifat riil kontinu [-256, 255]
8. 'value_mean', 'hue_mean', dan 'saturation_mean', yaitu hasil transformasi linear 3D dari RGB (menggunakan algoritma dari buku Foley and VanDam, Fundamentals of Interactive Computer Graphics). Atribut-atribut ini bersifat riil kontinu, value_mean memiliki range [0, 255], hue_mean memiliki range [-100,100], dan saturation_mean memiliki range [0, 1].
9. Atribut class menunjukkan region adalah kelas positif atau negatif. Atribut ini bersifat kategorik {'positive', 'negative'}.

Image Segmentation dataset ini termasuk kategori dataset tak seimbang yang di mana kelas mayoritasnya adalah region berkategori negatif dengan jumlah 1962 dan kelas minoritas adalah region berkategori positif dengan jumlah 346.



Gambar 6: Boxplot sebelum Outlier Removal - Image Segmentation Dataset

Dilihat dari seluruh boxplot pada Gambar 6, terdapat banyak outlier ekstrim yang tersebar jauh dari 90% sebaran data. Namun outlier ekstrim tersebut hanya tersebar pada kelas negatif (mayoritas) saja, yang berarti bahwa menghapus sebagian outlier ini tidak berdampak terhadap kelas positif (minoritas).



Gambar 7: Boxplot setelah Outlier Removal – Image Segmentation Dataset

Gambar 7 menunjukkan dengan jelas distribusi masing-masing kelas pada 6 atribut di atas setelah *outlier removal*.

4.1.2 Spambase Dataset

4.1.2.1 Data Cleaning & Normalization

Spambase dataset merupakan dataset yang menunjukkan email berlabel spam atau ham. Email berlabel spam pada dataset ini diperoleh dari kantor pos dan orang-orang yang mengajukan spam pada emailnya. Email berlabel ham (non-spam) diperoleh dari

email kerja dan email pribadi seseorang yang bernama George. Indikator non-spam adalah kata ‘George’ dan kode area ‘650’ yang hanya diketahui oleh orang-orang yang mengenal George. Dengan menggunakan *keyword* tertentu seperti ‘George’ dan ‘650’, *personalized spam filter*. Terdapat 58 atribut dalam dataset ini:

1. 48 atribut dalam format `word_freq_WORD` menunjukkan persentase dari frekuensi kemunculan kata `WORD` dari total kata pada email. Contoh: `word_freq_george` dan `word_freq_order`. Atribut-atribut ini bersifat riil kontinu $[0, 100]$.
2. 6 atribut dalam format `char_freq_CHAR` menunjukkan persentase dari frekuensi kemunculan karakter `CHAR` dari total karakter pada email. Contoh: `char_freq_`(dan `char_freq_!`. Atribut-atribut ini bersifat riil kontinu $[0, 100]$.
3. 1 atribut `capital_run_length_average` menunjukkan rata-rata panjang huruf kapital berurutan. Atribut ini bersifat riil kontinu $[1, \dots]$.
4. 1 atribut `capital_run_length_longest` menunjukkan jumlah terbanyak panjang huruf kapital berurutan. Atribut ini bersifat bilangan bulat $[1, \dots]$.
5. 1 atribut `capital_run_length_total` menunjukkan jumlah dari huruf kapital pada email. Atribut ini bersifat bilangan bulat $[1, \dots]$.
6. 1 atribut `class` menunjukkan email termasuk kategori spam (1) atau ham (0). Atribut ini bersifat nominal $\{0, 1\}$.

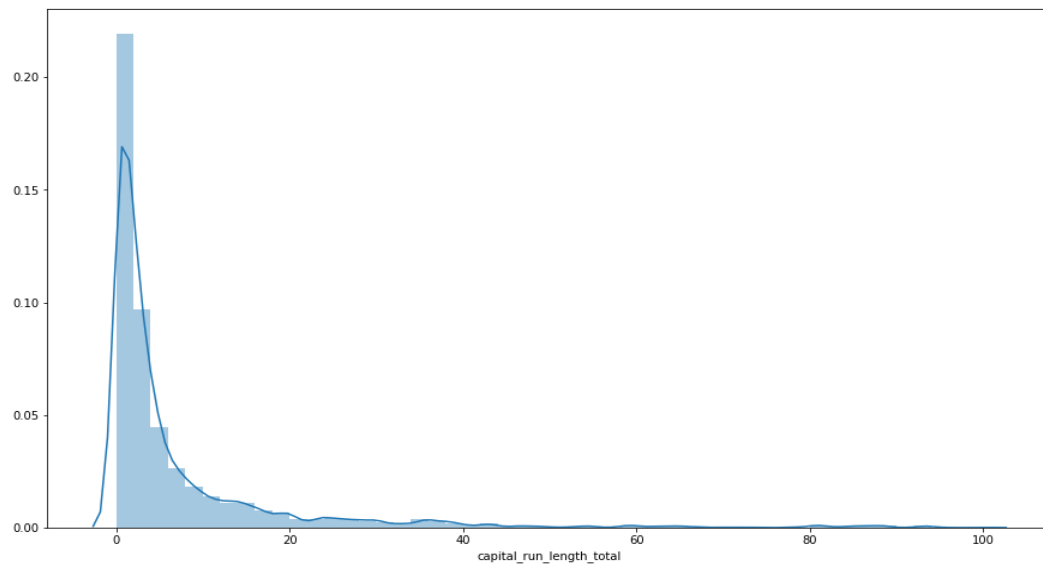
Spambase dataset ini termasuk kategori dataset tak seimbang yang di mana kelas mayoritasnya adalah email berkategori ham dengan jumlah 2788 dan kelas minoritas adalah email berkategori spam dengan jumlah 1813.

Sebagian besar atribut pada dataset ini memiliki distribusi yang sama, yaitu 48 atribut `word_freq_WORD` dan 6 atribut `char_freq_CHAR` yang memiliki rentang nilai $[1, 100]$. Tiga atribut `capital_run_length_average/longest/total` memiliki distribusi yang jauh berbeda dengan 54 atribut sebelumnya.

	capital_run_length_average	capital_run_length_longest	capital_run_length_total
count	4601.000000	4601.000000	4601.000000
mean	5.191515	52.172789	283.289285
std	31.729449	194.891310	606.347851
min	1.000000	1.000000	1.000000
25%	1.588000	6.000000	35.000000
50%	2.276000	15.000000	95.000000
75%	3.706000	43.000000	266.000000
max	1102.500000	9989.000000	15841.000000

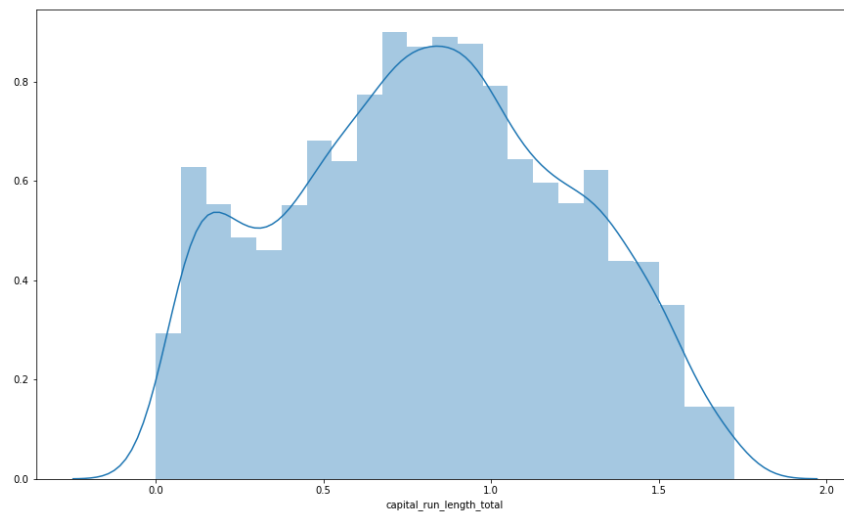
Tabel 2 Distribusi capital_run_length_average, capital_run_length_longest dan capital_run_length_total

Dari Tabel 2 dapat dilihat bahwa nilai maksimal dari masing-masing atribut berbeda jauh dari rentang nilai [1, 100]. Maka atribut-atribut ini perlu dinormalkan ke skala [1, 100] menggunakan MinMax. Namun terdapat lonjakan yang sangat besar dari kuartil ketiga (75%) ke nilai maksimal dari masing-masing atribut di atas. Hal ini mengindikasikan terdapat outlier ekstrim yang memberi bias terhadap standar deviasi atribut, dan memberi dampak negatif jika data diskalakan dengan outlier. Maka normalisasi akan dilakukan setelah outlier ekstrim telah dihapus dari dataset.



Gambar 8: Distribusi `capital_run_length_total` sebelum normalisasi

Sebagian besar atribut dari Spambase dataset condong ke kiri (*left-skewed*) seperti Gambar 8. Normalisasi dilakukan menggunakan transformasi log pada setiap atribut yang condong ke kiri.



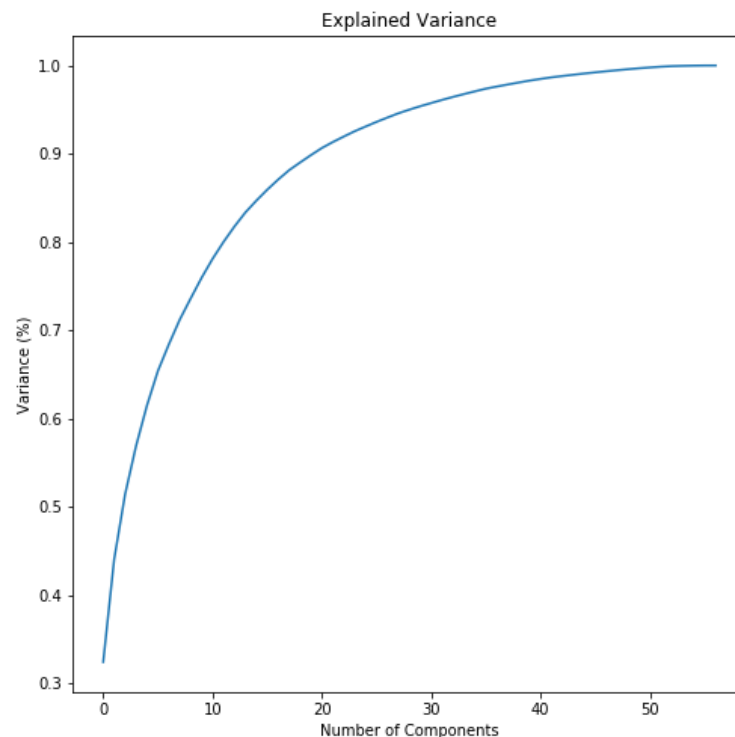
Gambar 9: Distribusi `capital_run_length_total` setelah normalisasi

Gambar 9 menunjukkan distribusi dari `capital_run_length_total` setelah normalisasi. Seluruh atribut telah diuji menggunakan Saphiro-Wilk test dengan p-value lebih kecil dari 0.05.

4.1.2.2 Attribute Reduction

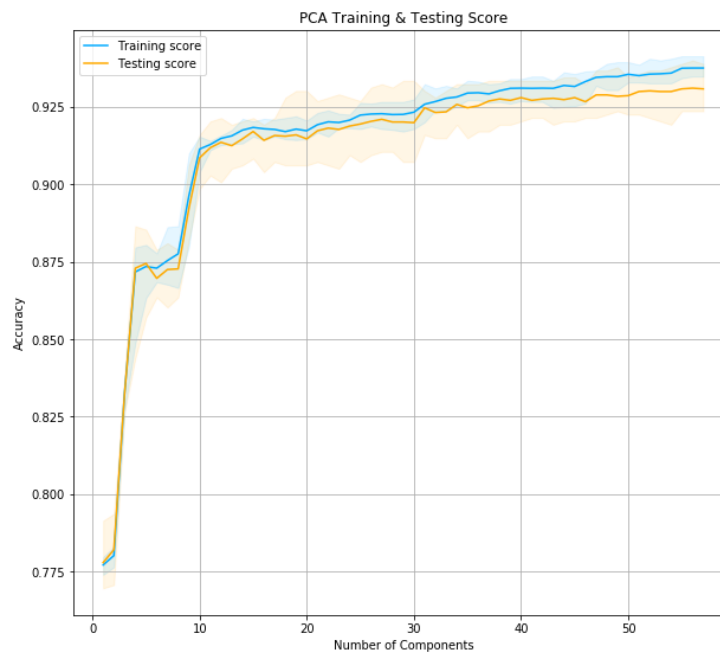
Atribut atau dimensi dataset perlu dikurangi untuk memberi hasil yang lebih umum, mengurangi overfitting, dan mengurangi running time. Terlebih jika dataset tersebut memiliki atribut yang sangat banyak atau data yang sangat besar.

Pada Spambase dataset terdapat 58 atribut independen yang akan ditransformasi dan direduksi menggunakan PCA. Atribut yang ditransformasi disebut komponen. Setiap jumlah komponen dilakukan uji sebanyak 5 kali cross validation yang dirata-ratakan.

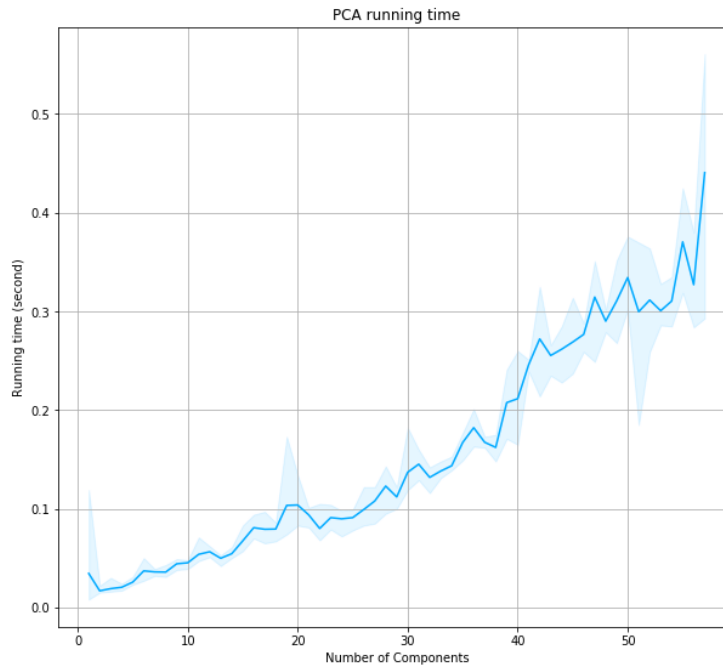


Gambar 10: Scree plot pada Spambase Dataset

Gambar 10 menunjukkan variansi yang dijelaskan oleh komponen-komponen utama PCA. Sumbu X adalah jumlah komponen dengan sumbu Y adalah total variansi yang dikandung untuk setiap jumlah komponen secara kumulatif. Memilih jumlah komponen dapat dilakukan dengan mengandalkan scree plot seperti Gambar 10. Namun memilih jumlah komponen untuk digunakan tidak cukup hanya dengan mereferensikan scree plot tersebut, konsekuensi lain seperti performa pada tahap *training*, *testing* dan *running time* juga perlu dipertimbangkan.



Gambar 11: Training & Testing Score spambase dataset berdasarkan jumlah komponen utama PCA



Gambar 12: Running time dari satu fit berdasarkan jumlah komponen utama PCA

Dilihat dari gambar Gambar 11, menggunakan 58 komponen menjamin akurasi training dan testing yang tertinggi, namun dengan running time yang tertinggi pula (Gambar 12). Sedangkan 10 komponen pertama yang menjelaskan 77% dari variansi data berkontribusi besar terhadap akurasi model, dengan rata-rata akurasi sebesar 91% hanya dengan perbedaan kurang dari 2% dibanding dengan menggunakan seluruh komponen. Dari sisi *running time*, menggunakan 10 komponen pertama memiliki perbedaan running time yang sangat besar hingga 8 kali lebih cepat. Dilihat dari sisi overfitting, terlihat bahwa model dapat memperoleh testing score yang tinggi dengan hanya sedikit bias dari training score-nya, namun perbedaan dari training score dan testing score dapat dilihat jelas setelah 30 komponen, dengan perbedaan atau bias terbesar terjadi ketika seluruh komponen digunakan.

Dari hasil analisis di atas, dapat disimpulkan bahwa jumlah komponen terbaik yang digunakan adalah 10 komponen.

4.1.3 Credit Card Fraud Dataset

4.1.3.1 Data Cleaning & Normalization

Credit Card Fraud dataset merupakan dataset berisi transaksi yang dilakukan dengan kartu kredit pada September 2013 di Eropa yang terjadi dalam dua hari. Dataset ini berisi informasi mengenai transaksi yang bersifat fraud dan non-fraud. Terdapat 31 atribut dalam dataset ini:

1. 'Time', yang berisi informasi mengenai waktu yang telah berlalu (detik) mengenai masing-masing transaksi sejak transaksi pertama pada dataset ini.
2. 28 Fitur V1 hingga V28, merupakan fitur hasil transformasi PCA yang dilakukan oleh instansi yang mengembangkan dataset ini untuk melindungi informasi sensitif dari kartu kredit.
3. 'Amount', yang berisi informasi mengenai jumlah hasil transaksi dalam USD.
4. Atribut class yang menunjukkan apakah suatu transaksi bersifat fraud atau non-fraud. Transaksi fraud berlabel 1 dan non-fraud berlabel 0. Atribut ini bersifat nominal {0, 1}.

Dataset ini berisi 284.807 total transaksi dengan 284.315 transaksi non-fraud dan 492 transaksi fraud. Dataset ini sangat tidak seimbang dengan kelas minoritas hanya berkisar sebesar 0,17% dari total transaksi.

Sebagian besar atribut pada dataset ini memiliki distribusi yang sama, yaitu 28 atribut V1 hingga V28 yang memiliki mean 0 dan standar deviasi mendekati 1. Dua atribut Amount dan Time memiliki distribusi yang jauh berbeda dengan 28 atribut tersebut.

	<i>Time</i>	<i>V1</i>	<i>V2</i>	...	<i>V27</i>	<i>V28</i>	<i>Amount</i>
<i>mean</i>	94813	0	0	...	0	0	88.349
<i>min</i>	0.000	-5.641e+01	-7.271e+01	...	-2.256e+01	-1.543e+01	0.000
<i>25%</i>	54201.500	-9.203e-01	-5.985e-01	...	-7.083e-02	-5.295e-02	5.600
<i>50%</i>	84692.000	1.810e-02	6.548e-02	...	1.342e-03	1.124e-02	22.000
<i>75%</i>	139320.500	1.315e+00	8.037e-01	...	9.104e-02	7.827e-02	77.165
<i>max</i>	172792.000	2.454e+00	2.205e+01	...	3.161e+01	3.384e+01	25691.160

Tabel 3: Karakteristik atribut dari Credit Card Fraud Dataset

Tabel 3 menunjukkan bahwa Time dan Amount memiliki mean dan nilai max yang sangat besar dibanding atribut V1 hingga V28. Atribut yang memiliki karakteristik berbeda ini dinormalisasi dengan cara normalisasi standar. Namun sebelum itu outlier ekstrim perlu dihilangkan terlebih dahulu.

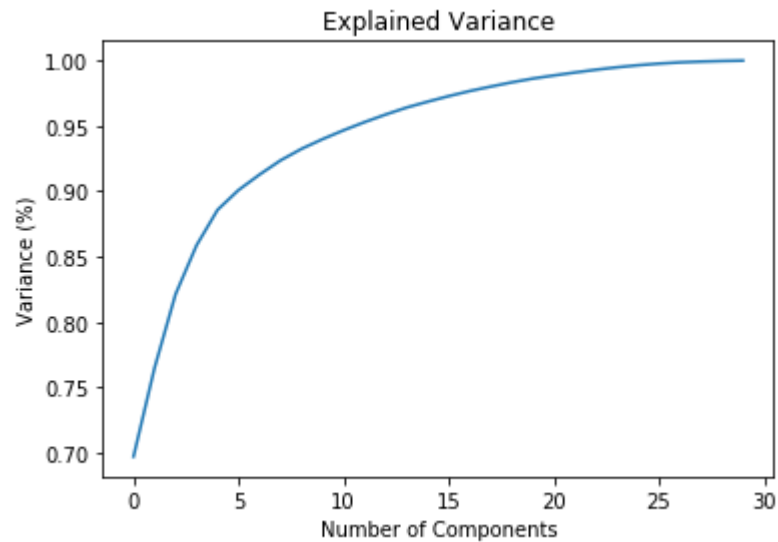
	<i>Time</i>	<i>V1</i>	<i>V2</i>	...	<i>V27</i>	<i>V28</i>	<i>Amount</i>
<i>mean</i>	0.118965	0	0	...	0	0	0.907330
<i>min</i>	-0.995030	-5.641e+01	-7.271e+01	...	-2.256e+01	-1.543e+01	-0.307865
<i>25%</i>	-0.358151	-9.203e-01	-5.985e-01	...	-7.083e-02	-5.295e-02	-0.229639
<i>50%</i>	0.000000	1.810e-02	6.548e-02	...	1.342e-03	1.124e-02	0.000000
<i>75%</i>	0.641849	1.315e+00	8.037e-01	...	9.104e-02	7.827e-02	0.770361
<i>max</i>	1.035070	2.454e+00	2.205e+01	...	3.161e+01	3.384e+01	69.099496

Tabel 4: Karakteristik atribut dari Credit Card Fraud Dataset setelah normalisasi

Tabel 4 menunjukkan bahwa atribut Time dan Amount telah berhasil dinormalisasi.

4.1.3.2 Attribute Reduction

Dilihat dari jumlah data, sangat perlu dilaksanakan atribut reduction pada dataset ini. Menggunakan seluruh atribut akan memakan waktu yang sangat lama untuk memproses data besar, terlebih jika memproses data ini dilakukan berkali-kali pada model *tuning* dan model *fitting*.



Gambar 13: Scree Plot Credit Card Fraud Dataset

Berdasarkan *scree plot* pada Gambar 13, sebesar 90% dari variansi data dicapai hanya dengan menggunakan 6 komponen utama, maka dataset ini berhasil direduksi menjadi 20% dengan tetap menyimpan 90% dari total informasi pada dataset asli (sebelum direduksi). Testing-training plot untuk tiap komponen tidak dapat dilakukan seperti spambase dataset dikarenakan besarnya data yang memakan waktu yang sangat lama untuk diproses.

4.1.3.3 Dataset Rebalancing

Walaupun *attribute reduction* telah dilaksanakan, data masih terlalu besar untuk melakukan *tuning*, *fitting*, dan cross validation berkali-kali. Maka sub-dataset perlu dibangun dari dataset asli.

Dataset dipisah berdasarkan kelasnya, yaitu dataset kelas minoritas dan dataset kelas mayoritas. Dataset yang lebih besar, atau dataset kelas mayoritas akan dibagi menjadi 200 kelompok, dengan masing-masing kelompok memiliki size yang sama, yaitu sekitar 1420 sampel per kelompoknya. Dataset baru dibangun berdasarkan dataset kelas minoritas dengan satu kelompok dataset kelas mayoritas, dengan total 200 dataset baru dengan kelas minoritas yang sama di setiap datasetnya namun dengan sampel mayoritas yang berbeda. Setiap dataset akan memiliki Imbalance Ratio sebesar 1:3 (dibanding 1:577 pada dataset asli). Model *tuning* dan *fitting* dilakukan terhadap masing-masing dataset baru ini. Pendekatan ini memiliki keuntungan:

1. Tidak ada information loss.
2. Berupa cross validation untuk kelas mayoritas.
3. Mengurangi overfitting.

Walaupun jumlah data yang diproses sama besarnya, memproses dataset dengan size kecil ini memiliki running time yang jauh lebih kecil dibanding memproses dataset dengan size besar walaupun jumlah iterasinya 200 kali lebih banyak.

	<i>original dataset</i>	<i>200 datasets approach</i>
<i>time(s)</i>	>15000	~6000

Tabel 5: Perbedaan waktu dari total fitting pada original dataset dan 200 datasets approach

Tabel 5 menunjukkan bahwa pendekatan 200 dataset (*200 datasets approach*) dapat menyelesaikan 30 kombinasi pada dataset Credit Card Fraud dalam 6.000 detik, sedangkan tanpa *dataset rebalancing* membutuhkan waktu lebih dari 15.000 detik.

4.2 Model Tuning & Fitting

Setelah melakukan data preprocessing, setiap dataset melewati 5-fold cross validation, dengan pengecualian Credit Card Fraud Dataset yang harus melewati dataset rebalancing terlebih dahulu. 5-fold cross validation memberi distribusi 80% data training dan 20% data testing dengan data testing yang berbeda pada tiap iterasinya. Untuk setiap iterasi cross validation, data training diseimbangkan dengan melakukan resampling. Hasil resampling adalah data training yang telah diseimbangkan. Data training yang telah diseimbangkan ini akan diklasifikasikan melalui algoritma machine learning dengan berbagai jenis parameter:

1. Regresi Logistik: C dan regularization. C parameter berguna sebagai pengontrol dari regularisasi Ridge dan Lasso Regression. C parameter adalah invers lambda atau $C = \frac{1}{\lambda}$.
 - C: {0.01, 0.1, 0.5, 0.75, 1, 10, 100}
 - Regularization: Lasso Regression dan Ridge Regression
2. Support Vector Machines: C, kernel trick, dan penalty function. Kernel trick adalah fungsi pemetaan yang digunakan untuk mengubah ruang input ke ruang lain dengan tujuan mempermudah pemisahan class dengan hyperplane.
 - C: {0.01, 0.1, 0.5, 1, 10, 100}
 - Kernel trick: Radial Basis Function, Polynomial, Sigmoid, Linear
 - Regularization: Lasso Regression dan Ridge Regression

3. MultiLayer Perceptron: fungsi aktivasi, jumlah hidden layer dan jumlah neuron pada hidden layer.
 - Fungsi aktivasi: Logistik, tanh, Rectified Linear Unit (ReLU)
 - Jumlah hidden layer: { 1, 2 }
 - Jumlah neuron pada hidden layer: $\{|X|, \frac{|X|}{2}\}$ dengan $|X|$ adalah panjang fitur vektor X .
4. Decision Tree: Sampel minimal untuk leaf dan kedalaman maksimal.
 - Sampel minimal untuk leaf: { 5, 10, 15 }
 - Kedalaman maksimal : { 3, 4, 5, 6, 7 }
5. KNN: K atau jumlah tetangga terdekat, dan bobot untuk masing-masing tetangga.
 - K: { 3, 5, 7 }
 - Bobot: Uniform, weighted (Weighted-KNN).

Setiap pengklasifikasi melakukan 3-fold cross validation pada data training yang telah diseimbangkan (*balanced dataset*) dengan setiap kombinasi parameter yang ada (Contoh kombinasi parameter untuk MLP: Fungsi aktivasi Logistik, jumlah hidden layer 1, dan jumlah neuron pada hidden layer $\frac{|X|}{2}$). Kombinasi parameter dengan recall tertinggi pada 3-fold cross validation tersebut dipilih untuk melakukan uji performa pada 20% data testing asli.

Nilai ukur yang digunakan untuk setiap hasil testing dari tiap kombinasi adalah akurasi (secara umum), precision-0 (akurasi prediksi kelas mayoritas), precision-1 (akurasi prediksi kelas minoritas), recall-0 (tingkat pengenalan kelas mayoritas), recall-1 (tingkat pengenalan kelas minoritas), f1-0 (nilai f1 dari kelas mayoritas), f1-1 (nilai f1 dari kelas minoritas). Namun pada kasus imbalanced class, nilai ukur yang terpenting adalah recall-1.

Merata-ratakan hasil dari ketiga dataset, algoritma machine learning memiliki performa sebagai berikut:

<i>clf</i>	<i>accuracy</i>	<i>precision-0</i>	<i>precision-1</i>	<i>recall-0</i>	<i>recall-1</i>	<i>f1-0</i>	<i>f1-1</i>
<i>Lr</i>	0.948	0.972	0.902	0.948	0.942	0.960	0.919
<i>svm</i>	0.955	0.975	0.916	0.957	0.946	0.966	0.930
<i>mlp</i>	0.956	0.976	0.921	0.958	0.946	0.966	0.931
<i>dt</i>	0.921	0.947	0.893	0.935	0.870	0.938	0.870
<i>knn</i>	0.949	0.977	0.895	0.944	0.951	0.960	0.921

Tabel 6: Performa algoritma machine learning

Tabel 6 menunjukkan bahwa berdasarkan recall-1, seluruh algoritma machine learning memiliki performa yang serupa (0.942 – 0.951) dengan pengecualian pada algoritma Decision Tree yang memiliki recall-1 0.870.

Dengan merata-ratakan hasil dari ketiga dataset, teknik resampling memiliki performa sebagai berikut:

<i>Res</i>	<i>accuracy</i>	<i>precision-0</i>	<i>precision-1</i>	<i>recall-0</i>	<i>recall-1</i>	<i>f1-0</i>	<i>f1-1</i>
<i>Smote</i>	0.952	0.969	0.920	0.958	0.931	0.963	0.923
<i>bsmote</i>	0.931	0.976	0.854	0.919	0.949	0.945	0.895
<i>adasyn</i>	0.928	0.978	0.844	0.914	0.953	0.944	0.891
<i>Ros</i>	0.951	0.969	0.919	0.957	0.929	0.963	0.922
<i>Rus</i>	0.948	0.969	0.904	0.951	0.930	0.960	0.914
<i>TL</i>	0.955	0.963	0.946	0.969	0.914	0.966	0.927

Tabel 7: Performa teknik resampling

Tabel 7 menunjukkan bahwa berdasarkan recall-1, ADASYN adalah teknik resampling yang memiliki *detection rate* tertinggi untuk kelas minoritas dengan nilai 0.953.

Dengan merata-ratakan hasil dari ketiga dataset, berikut adalah 10-tertinggi kombinasi dari algoritma machine learning diurut berdasarkan recall-1 (*detection rate* kelas minoritas):

<i>res</i>	<i>clf</i>	<i>accuracy</i>	<i>precision-0</i>	<i>precision-1</i>	<i>recall-0</i>	<i>recall-1</i>	<i>f1-0</i>	<i>f1-1</i>
<i>adasyn</i>	mlp	0.936	0.984	0.852	0.921	0.969	0.951	0.904
<i>adasyn</i>	Lr	0.927	0.981	0.819	0.909	0.967	0.943	0.885
<i>adasyn</i>	svm	0.939	0.983	0.855	0.927	0.966	0.953	0.905
<i>adasyn</i>	knn	0.931	0.984	0.840	0.913	0.966	0.946	0.897
<i>bsmote</i>	Lr	0.929	0.981	0.824	0.912	0.965	0.945	0.888
<i>bsmote</i>	mlp	0.943	0.982	0.870	0.932	0.964	0.956	0.912
<i>bsmote</i>	svm	0.944	0.982	0.869	0.934	0.963	0.957	0.913
<i>bsmote</i>	knn	0.941	0.982	0.869	0.928	0.963	0.953	0.912
<i>rus</i>	knn	0.945	0.979	0.875	0.939	0.954	0.958	0.912
<i>smote</i>	knn	0.951	0.977	0.900	0.948	0.951	0.962	0.924

Tabel 8: 10-tertinggi kombinasi algoritma machine learning dan teknik resampling diurut berdasarkan recall-1

4.3 Analisis Hasil

Borderline-SMOTE dan ADASYN memiliki recall tertinggi dalam memprediksi kelas minoritas, namun dengan mengorbankan akurasi. Hal ini disebabkan karena Borderline-SMOTE dan ADASYN membuat banyak data sintetik di sekitar kelas mayoritas, memberi batas yang jelas antara kelas minoritas dan mayoritas. Pada ketiga dataset, Borderline-SMOTE dan ADASYN secara konsisten memiliki recall tertinggi.

Tomek Links memiliki recall terendah dalam memprediksi kelas minoritas jika dibandingkan dengan teknik resampling yang lain. Hal ini disebabkan karena Tomek Links tidak dapat menyeimbangkan jumlah kelas mayoritas dan minoritas. Secara praktis, teknik ini hanya merupakan teknik penghapusan outlier ketika teknik resampling lain telah dilakukan.

KNN memiliki recall yang baik dikarenakan jumlah tetangga minoritas akan sangat dominan ketika class telah diseimbangkan dengan teknik resampling. Namun berdampak sangat buruk terhadap presisi dari kelas minoritas, menghasilkan banyak *false positive*.

Performa MLP dan SVM adalah yang terbaik pada algoritma klasifikasi machine learning, namun dengan running time yang tertinggi pula. Kedua metode ini secara konsisten memiliki recall terbaik pada ketiga dataset yang diberikan, menjadikannya pengklasifikasi yang ideal pada masalah ketidakseimbangan data.

Decision Tree adalah algoritma machine learning yang memiliki performa terburuk pada ketiga dataset ini, hal ini dikarenakan decision tree adalah algoritma klasifikasi yang lebih baik digunakan pada dataset dengan atribut kategorik, sedangkan ketiga dataset yang digunakan seluruhnya memiliki atribut kontinu.

BAB V

KESIMPULAN DAN SARAN

5.1 Kesimpulan

Berdasarkan penelitian yang telah dilakukan, peneliti menarik beberapa kesimpulan sebagai berikut:

1. ADASYN dan Borderline-SMOTE adalah teknik resampling terbaik untuk mengenali kelas minoritas, namun dengan menurunkan sedikit akurasi. SMOTE berada di bawah ADASYN dan Borderline-SMOTE pada recall-1, namun SMOTE jauh lebih unggul pada precision-1 dan akurasi pada umumnya, yang menghasilkan nilai f1-1 yang lebih tinggi dari ADASYN dan Borderline-SMOTE. SMOTE juga lebih baik pada recall-0, yang ideal jika *false positive* juga sangat tidak diinginkan pada suatu dataset.
2. KNN memiliki recall terbaik dibanding seluruh pengklasifikasi lain ketika class telah diseimbangkan dengan teknik resampling, namun MLP dan SVM memiliki akurasi yang lebih tinggi dengan perbedaan recall yang sangat tipis dari KNN (2%). MLP dan SVM juga memiliki precision-1 yang jauh lebih tinggi dari KNN, menghasilkan nilai f1-1 yang lebih tinggi. Secara umum, MLP dan SVM memiliki performa yang lebih baik pada ketiga dataset yang diberikan.
3. Kombinasi algoritma machine learning dan teknik resampling yang terbaik adalah *adasyn_mlp* untuk recall tertinggi, *smote_knn* untuk nilai f1 tertinggi (keseimbangan antara recall dan precision), dan *tl_mlp* untuk akurasi pada umumnya.
4. Dataset dengan class yang tingkat ketidakseimbangannya sangat tinggi, seperti Credit Card Fraud Dataset, lebih mudah untuk diproses jika dataset dikelompokkan seperti pada bagian **4.1.3.3 Dataset Rebalancing**. Pendekatan ini menunjukkan jauhnya perbedaan pada recall-1 jika dibandingkan data diolah langsung sebelum dikelompokkan.

5.2 Saran

Penelitian ini hanya membahas tentang ketidakseimbangan class biner, untuk ketidakseimbangan multiclass, di mana class lebih dari dua, tidak dibahas pada penelitian ini. Teknik undersampling seperti SBC dan NearMiss juga tidak dibahas pada penelitian ini. Maka saran untuk penelitian selanjutnya adalah membahas mengenai ketidakseimbangan class di mana class-nya lebih dari dua (multiclass), dan menggunakan teknik-teknik resampling atau algoritma machine learning yang belum dibahas pada penelitian ini.

Daftar Pustaka

- Ahmed, N. K., Atiya, A. F., Gayar, N. E., & El-Shishiny, H. (2010). An Empirical Comparison of Machine Learning Models for Time Series Forecasting. *Econometric Reviews*, 29(5-6), 594-621.
- Ali, A., Shamsuddin, S. M., & Ralescu, A. L. (2015). Classification with class imbalance problem: A Review. *Int. J. Advance Soft Compu. Appl*, 7(3), 176-204.
- Amin, A., Anwar, S., Adnan, A., Nawaz, M., Howard, N., Qadir, J., & Hawalah, A. (2016). Comparing Oversampling Techniques to Handle the Class Imbalance Problem: A Customer Churn Prediction Case Study. *IEEE Access*, 4, 7940-7957.
- Anand, A., Pugalenth, G., Fogel, G. B., & Suganthan, P. N. (2010). An approach for classification of highly imbalanced data using weighting and undersampling. *Amino acids*, 39(5), 1385-1391.
- Andrew, N. (2000). CS229 Lecture Notes.
- Baars, H., & Kemper, H. G. (2008). Management support with structured and unstructured data—an integrated business intelligence framework. *Information Systems Management*, 25(2), 132-148.
- Batista, G. E., Prati, R. C., & Monard, M. C. (2004). A Study of the Behavior of Several Methods for Balancing Machine Learning Training Data. *ACM SIGKDD explorations newsletter*, 6(1), 20-29.
- Beal, V. (2019). *Unstructured Data*. Retrieved 6 20, 2019, from https://www.webopedia.com/TERM/U/unstructured_data.html
- Berman, F., Rutenbar, R., Hailpern, B., Christensen, H., Davidson, S., Estrin, D., . . . Szalay, A. S. (2018). Realizing the Potential of Data Science. *Communications Of The Acm*, 61(4), 67-72.

- Bishop, C. M. (1995). *Neural Networks for Pattern Recognition*. Oxford university press.
- Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Springer.
- Blumberg, R., & Atre, S. (2003). The Problem with Unstructured Data. *DM Review*, 13(42-49), 62.
- Bradley, A. P. (1997). The Use of the Area Under the ROC Curve in the Evaluation of Machine Learning Algorithms. *Pattern recognition*, 30(7), 1145-1159.
- Breiman, L., Friedman, J. H., Olshen, R. A., & Stone, C. J. (2017). *Classification and Regression Trees*. Routledge: CRC Press LLC.
- Burnaev, E., Erofeev, P., & Papanov, A. (2015). Influence of Resampling on Accuracy of Imbalanced Classification. *In Eighth International Conference on Machine Vision (ICMV 2015)*, 9875, 987521.
- Chang, C.-C., & Lin, C.-J. (2011). LIBSVM: A Library for Support Vector Machines. *ACM Transactions on Intelligent Systems and Technology*.
- Chawla, N. V. (2009). *DATA MINING FOR IMBALANCED DATASETS: AN OVERVIEW*. Boston: Springer.
- Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: Synthetic Minority Over-sampling Technique. *Journal of artificial intelligence research*, 16, 321-357.
- Chen, H., Chiang, R. H., & Storey, V. C. (2012). Business intelligence and analytics: From big data to big impact. *MIS quarterly*, 36(4).
- Chen, M., Mao, S., & Liu, Y. (2014). Big Data: A survey. *Mobile networks and applications*, 19(2), 171-209.

- Dai, Q.-y., Zhang, C.-p., & Wu, H. (2016). Research of Decision Tree Classification Algorithm in Data Mining. *International Journal of Database Theory and Application*, 9(5), 1-8.
- Desjardins, J. (2019, April 17). *How much data is generated each day?* World Economic Forum. Retrieved June 14, 2019, from <https://www.weforum.org/agenda/2019/04/how-much-data-is-generated-each-day-cf4bddf29f/>
- Dhar, V. (2012). Data Science and Prediction. *Communications of the ACM*, 56(12), 64-73.
- Díaz-Uriarte, R. &. (2006). Gene selection and classification of microarray data using random forest. *BMC Bioinformatics* 7 no. 1, 3.
- Diri, B., & Albayrak, S. (2008). Visualization and analysis of classifiers performance in multi-class medical data. *Expert Systems with Applications*, 34(1), 628-634.
- Domingos, P. (2011). A Few Useful Things to Know about Machine Learning. 78-87.
- Donoho, D. L., & Tanner, J. (2010). Precise Undersampling Theorems. *Proceedings of the IEEE*, 98(6), 913-924.
- Drummond, C., & Holte, R. C. (2003, August). C4.5, Class Imbalance, and Cost Sensitivity: Why Under-Sampling beats Over-Sampling. *Workshop on learning from imbalanced datasets II*, 11, 1-8.
- Elgendy, N., & Elragal, A. (2014). Big data analytics: a literature review paper. *Industrial Conference on Data Mining*, 214-227.
- Elrahman, S. M., & Abraham, A. (2013). A Review of Class Imbalance Problem. *Journal of Network and Innovative Computing*, 332-340.
- Ethem, A. (2009). *Introduction to Machine Learning*. MIT press.

- Fernandez, A., Garcia, S., Herrera, F., & Chawla, N. V. (2018). SMOTE for Learning from Imbalanced Data: Progress and Challenges, Marking the 15-year Anniversary. *Journal of Artificial Intelligence Research*, 61, 863-905.
- Galar, M., Fernandez, A., Barrenechea, E., Bustince, H., & Herrera, F. (2011). A review on ensembles for the class imbalance problem: bagging-, boosting-, and hybrid-based approaches. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 42(4), 463-484.
- García, S., & Herrera, F. (2009). Evolutionary Undersampling for Classification with Imbalanced Datasets: Proposals and Taxonomy. *Evolutionary computation*, 17(3), 275-306.
- Guenther, N., & Schonlau, M. (2016). Support vector machines. *The Stata Journal*, 16(4), 917-937.
- Guo, H., & Viktor, H. L. (2004). Learning from Imbalanced Data Sets with Boosting and Data Generation: The DataBoost-IM Approach. *ACM Sigkdd Explorations Newsletter*, 6(1), 30-39.
- Han, H., Wang, W. Y., & Mao, B. H. (2005, August). Borderline-SMOTE: A New Over-Sampling Method in Imbalanced Data Sets Learning. *International conference on intelligent computing*, 878-887.
- Haykin, S. (1994). *Neural Networks: A Comprehensive Foundation*. Prentice Hall PTR.
- He, H., & Garcia, E. A. (2008). Learning from Imbalanced Data. *IEEE Transactions on Knowledge & Data Engineering*(9), 1263-1284.
- He, H., Bai, Y., Garcia, E. A., & Li, S. (2008, June). ADASYN: Adaptive synthetic sampling approach for imbalanced learning. *2008 IEEE International Joint Conference on Neural Networks (IEEE World Congress on Computational Intelligence)*, 1322-1328.

- Hechenbichler, S. (2004). Weighted k-Nearest-Neighbor Techniques and Ordinal. *Sonderforschungsbereich 386*.
- Hosmer, D. W., & Lemeshow, S. (2000). *Applied Logistic Regression (Second Edition)*. Canada: Wiley-Interscience Publication.
- Howard, A. (2013). *Elementary Linear Algebra, Binder Ready Version: Applications Version*. John Wiley & Sons.
- Hu, S. L. (2009, October). MSMOTE: Improving Classification Performance When Training Data is Imbalanced. *2009 second international workshop on computer science and engineering*, 2, 13-17.
- IBM Corporation. (n.d.). *What is a data set?* (IBM Corporation) Retrieved June 14, 2019, from https://www.ibm.com/support/knowledgecenter/zosbasics/com.ibm.zos.zconcepts/zconc_datasetintro.htm
- Jolliffe, I. T., & Cadima, J. (2016). Principal component analysis: a review and recent developments. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*.
- Jordan, I. M., & Mitchell, M. T. (2015). Machine learning: Trends, perspectives, and prospects. *Science*, 349(6245), 255-260.
- Kittler, J., Hater, M., & Duin, R. P. (1996). Combining classifiers. *Proceedings of 13th international conference on pattern recognition*, 897-901.
- Kotsiantis, S., Kanellopoulos, D., & Pintelas, P. (2006). Handling imbalanced datasets: A review. *GESTS International Transactions on Computer Science and Engineering*, 30(1), 25-36.
- Kubat, M., & Matwin, S. (1997, July). Addressing the Curse of Imbalanced Training Sets: One-Sided Selection. *Icml*, 97, 179-186.

- Kumar, A., & Sheshadri, H. (2012). On the Classification of Imbalanced Datasets. *International Journal of Computer Applications* (0975 – 8887), 44(8), 1-7.
- LeDell, E., Petersen, M., & Laan, M. v. (2015). Computationally efficient confidence intervals for cross-validated area under the ROC curve estimates. *Electronic Journal of Statistics*, 9, 1583-1607.
- Lewis, D. D. (1998). Naive (Bayes) at Forty: The Independence Assumption in Information Retrieval. *European conference on machine learning*, 4-15.
- Liaw, A., & Wiener, M. (2002). Classification and regression by randomForest. *R news* 2 no. 3, 18-22.
- Lin, W. C., Tsai, C. F., Hu, Y. H., & Jhang, J. S. (2017). Clustering-based undersampling in class-imbalanced data. *Information Sciences*, 409, 17-26.
- LinkedIn Economic Graph Team. (2018). *Linkedin 2018 Emerging Jobs Report*. LinkedIn.
- Liu, A. Y.-c. (2004). The effect of oversampling and undersampling on classifying imbalanced text datasets. *The University of Texas at Austin*.
- Liu, X. Y., Wu, J., & Zhou, Z. H. (2008). Exploratory Undersampling for Class-Imbalance Learning. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 39(2), 539-550.
- Lohr, S. (2012). *The age of big data*. New York: New York Times.
- McAfee, A., Brynjolfsson, E., Davenport, T. H., Patil, D. J., & Barton, D. (2012). Big data: the management revolution. *Harvard business review*, 90(10), 60-68.
- More, A. (2016). Survey of resampling techniques for improving classification performance in unbalanced datasets. *arXiv preprint arXiv:1608.06048*.

- MULTILAYER, A. N. (1998). Artificial Neural Networks (The Multilayer Perceptron) - A Review of Applications in the Atmospheric Sciences. *Atmospheric Environment*, 32(14-15), 2627-2636.
- Ng, A., & Katanforoosh, K. (2018). CS 229 Lecture Notes. In A. NG. California: Stanford Edu.
- Osisanwo, F., Akinsola, J., Awodele, O., Hinmikaiye, Olakanmi, & Akinjobi. (2017). Supervised Machine Learning Algorithms: Classification and Comparison. *International Journal of Computer Trends and Technology (IJCTT)*, 48(3), 128-138.
- Pal, M. (2005). Random forest classifier for remote sensing classification. *International Journal of Remote Sensing* 26 no.1, 217-222.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., . . . Duchesnay, M. (2011). Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* 12.
- PENG, C.-Y. J., LEE, K. L., & INGERSOLL, G. M. (2002). An introduction to logistic regression analysis and reporting. *The journal of educational research*, 96(1), 3-14.
- Pereira, J. M., Basto, M., & Silva, A. F. (2016). The Logistic Lasso and Ridge Regression in Predicting Corporate Failure. *Procedia Economics and Finance*, 634-641.
- Provost, F. (2000, July). Machine Learning from Imbalanced Data Sets 101. *Proceedings of the AAAI'2000 workshop on imbalanced data sets*, 68, 1-3.
- Provost, F., & Fawcett, T. (2013). Data Science and its Relationship to Big Data and Data-Driven Decision Making. *Big data*, 1(1), 51-59.
- Quinlan, J. (1986). Induction of Decision Trees. *Machine Learning*, 1, 81-106.

- Rahman, M. M., & Davis, D. N. (2013). Addressing the Class Imbalance Problem in Medical Datasets. *International Journal of Machine Learning and Computing*, 3(2), 224.
- Rish, I. (2001). An Empirical Study of the Naive Bayes Classifier. *IJCAI 2001 workshop on empirical methods in artificial intelligence*, 41-46.
- Rish, I. (2011). An empirical study of the naive Bayes classifier. *Watson Research Center*.
- Russell, S. J., & Norvig, P. (2016). *Artificial Intelligence: A Modern Approach*. Malaysia: Pearson Education Limited.
- Safavian, S. R., & Landgrebe, D. (1991). A survey of decision tree classifier methodology. *IEEE transactions on systems, man, and cybernetics* 21 no.3, 660-674.
- Scholkopf, B., & Smola, A. J. (2001). *Learning with kernels: support vector machines, regularization, optimization, and beyond*. MIT press.
- Schuldt, C., Laptev, I., & Caputo, B. (2004). Recognizing human actions: a local SVM approach. *Proceedings of the 17th International Conference on Pattern Recognition*, 3, 32-36.
- Snijders, C., Matzat, U., & Reips, U.-D. (2012). 'Big Data'. Big gaps of knowledge in the field of Internet. *International Journal of Internet Science*, 7, 1-5.
- Statistics Solution. (2016). *Statistics Solution*. Retrieved August 14, 2019, from <https://www.statisticssolutions.com/sample-size-calculation-and-sample-size-justification-resampling/>
- Suykens, J., & Vandewalle, J. (1999). Least Squares Support Vector Machine Classifiers. *Neural Processing Letters* 9.
- Swets, J. A. (1988). Measuring the Accuracy of Diagnostic Systems. *Science*, 240(4857), 1285-1293.

- Tibshirani, R. (1996). Regression Shrinkage and Selection via the Lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 267-288.
- Tomek, I. (1976). Two Modifications of CNN. *IEEE Trans. Systems, Man and Cybernetics*, 6, 769-772.
- Tsangaratos, P., & Ilia, I. (2016). Comparison of a logistic regression and Naïve Bayes classifier in landslide susceptibility assessments: The influence of models complexity and training dataset size. *Catena*, 145, 164-179.
- Vafeiadis, T., Diamantaras, K., Sarigiannidis, G., & Chatzisavvas, K. (55). A comparison of machine learning techniques for customer. *Simulation Modelling Practice and Theory*, 55(1), 1-9.
- Vapnik, V., & Cortes, C. (1995). Support-vector networks. *Machine Learning*, 20(3), 273-297.
- Visa, S., & Ralescu, A. (2005, April). Issues in Mining Imbalanced Data Sets - A Review Paper. *Proceedings of the sixteen midwest artificial intelligence and cognitive science conference, 2005*, 67-73.
- Wang, Y., Kung, L., & Byrd, T. A. (2018). Big data analytics: Understanding its capabilities and potential benefits for healthcare organizations. *Technological Forecasting and Social Change*, 126, 3-13.
- Weglarz, G. (2004). Two Worlds of Data - Unstructured and Structured. *DM Review*, 14, 19-23.
- Yen, S. J., & Lee, Y. S. (2006). Under-Sampling Approaches for Improving Prediction of the Minority Class in an Imbalanced Dataset. *Intelligent Control and Automation*, 731-740.
- Yen, S. J., & Lee, Y. S. (2009). Cluster-based under-sampling approaches for imbalanced data distributions. *Expert Systems with Applications*, 36(3), 5718-5727.

- Zhang, G. P. (2000). Neural Networks for Classification: A Survey. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 30(4), 451-462.
- Zhang, H. (2004). The Optimality of Naive Bayes. *American Association for Artificial Intelligence*.
- Zhou, Z. H., & Liu, X. Y. (2006). Training Cost-Sensitive Neural Networks with Methods Addressing the Class Imbalance Problem. *IEEE Transactions on Knowledge & Data Engineering*(1), 63-77.
- Zurada, J. M. (1992). *Introduction to artificial neural systems*. St. Paul: West publishing company.