

**ANALISIS TINGKAT  
KEJAHATAN MENGGUNAKAN  
MODEL RANDOM FOREST**



**OLEH :**

**PARHAN HAMBALI**

**RONI HABIBI S.KOM.,MT.,SFPC**

Judul :

Analisis Tingkat Kejahatan Menggunakan Model Random Forest  
2021,

Penulis

PARHAN HAMBALI

RONI HABIBI., S.Kom., M.T., SFPC

Penyunting : CAHYO PRIANTO, S.Pd., M.T., CDSP, SFPC

Layout : PARHAN HAMBALI

### **Sanksi Pelanggaran Pasal 113**

### **Undang-Undang Nomor 28 Tahun 2014**

### **tentang Hak Cipta**

- (1) Setiap Orang yang dengan tanpa hak melakukan pelanggaran hak ekonomi sebagaimana dimaksud dalam Pasal 9 ayat (1) huruf i untuk Penggunaan Secara Komersial dipidana dengan pidana penjara paling lama 1 (satu) tahun dan/atau pidana denda paling banyak Rp100.000.000 (seratus juta rupiah).
- (2) Setiap Orang yang dengan tanpa hak dan/atau tanpa izin Pencipta atau pemegang Hak Cipta melakukan pelanggaran hak ekonomi Pencipta sebagaimana dimaksud dalam Pasal 9 ayat (1) huruf c, huruf d, huruf e, dan/atau huruf h untuk Penggunaan Secara Komersial dipidana dengan pidana penjara paling lama 3 (tiga) tahun dan/atau pidana denda paling banyak Rp500.000.000,00 (lima ratus juta rupiah).
- (3) Setiap Orang yang dengan tanpa hak dan/atau tanpa izin Pencipta atau pemegang Hak Cipta melakukan pelanggaran hak ekonomi Pencipta sebagaimana dimaksud dalam Pasal 9 ayat (1) huruf a, huruf b, huruf c, dan/atau huruf g untuk Penggunaan Secara Komersial dipidana dengan pidana penjara paling lama 4 (empat) tahun dan/atau pidana denda paling banyak Rp1.000.000.000,00 (satu miliar rupiah).
- (4) Setiap Orang yang memenuhi unsur sebagaimana dimaksud pada ayat (3) yang dilakukan dalam bentuk pembajakan, dipidana dengan pidana

penjara paling lama 10 (sepuluh) tahun dan/atau pidana denda paling banyak Rp4.000.000.000,00 (empat miliar rupiah).

## Kata Sambutan

Puji syukur kehadiran Allah Subhanahu Wa Ta'ala, atas rahmat dan hidayat-Nya penulis dapat menyelesaikan buku yang berjudul “Analisis Tingkat Kejahatan Menggunakan Model Random Forest”. Buku ini disusun dalam konteks mendorong masyarakat khususnya mahasiswa yang sedang mencari referensi kebutuhan laporan, baik laporan magang ataupun tugas akhir. Semoga buku ini bisa memberikan banyak manfaat kepada kalian.

# Kata Pengantar

Puji syukur kehadiran Allah Subhanahu Wa Ta'ala, atas rahmat dan hidayat-Nya penulis dapat menyelesaikan buku yang berjudul “Analisis Tingkat Kejahatan Menggunakan Model Random Forest” .

Dalam kesempatan ini penulis tidak lupa menyampaikan ucapan terima kasih kepada semua pihak yang telah memberikan bantuan moral dan spiritual, langsung maupun tidak langsung dalam menyelesaikan buku ini. Semoga buku ini bisa bermanfaat bagi penulis dan para pembaca.

# Daftar Isi

<b>KATA SAMBUTAN .....</b>	<b>IV</b>
<b>KATA PENGANTAR .....</b>	<b>V</b>
<b>DAFTAR ISI .....</b>	<b>VI</b>
<b>DAFTAR GAMBAR.....</b>	<b>VIII</b>
<b>DAFTAR SINGKATAN.....</b>	<b>XIII</b>
<b>BAB I PENDAHULUAN .....</b>	<b>1</b>
1.1 LATAR BELAKANG .....	1
1.2 IDENTIFIKASI MASALAH .....	2
1.3 TUJUAN .....	2
1.4 MANFAAT PENELITIAN .....	2
<b>BAB II.....</b>	<b>4</b>
<b>LANDASAN TEORI.....</b>	<b>4</b>
2.1 PENELITIAN YANG BERHUBUNGAN .....	4
2.2 LANDASAN TEORI.....	8
2.2.1 <i>Kriminalitas</i> .....	8
2.2.2 <i>Pencurian</i> .....	8
2.2.3 <i>Python</i> .....	9
2.2.4 <i>Instalasi Python Menggunakan Installer</i> .....	10
2.2.5 <i>Instalasi Python dengan Anaconda</i> .....	13
2.2.6 <i>Machine Learning</i> .....	22
2.2.7 <i>Data Science</i> .....	28
2.2.8 <i>Data Analytic</i> .....	30
2.2.9 <i>Scikit Learn</i> .....	30
2.2.10 <i>Confusion Matrix</i> .....	31
2.2.11 <i>Balance Dataset</i> .....	31
2.2.12 <i>Imbalace Dataset</i> .....	31
2.2.13 <i>Visualisasi Data</i> .....	31
2.2.14 <i>Numpy</i> .....	32
2.2.15 <i>Matplotlib</i> .....	34
2.2.16 <i>Pandas</i> .....	38

2.2.17	<i>Classification Report</i> .....	46
2.2.18	<i>Jupyter Notebook</i> .....	50
BAB III .....		53
METODE PENELITIAN .....		53
3.1	METODOLOGI PENELITIAN .....	53
3.2	TAHAPAN – TAHAPAN DIAGRAM ALUR METODOLOGI PENELITIAN .....	54
3.2.1	<i>Melakukan Analisis</i> .....	54
3.2.2	<i>Obtain Data (Mengumpulkan Data)</i> .....	54
3.2.3	<i>Scrub Data (Pembersihan Data)</i> .....	54
3.2.4	<i>Explore</i> .....	55
3.2.5	<i>Model</i> .....	55
3.2.6	<i>Random Forest</i> .....	55
3.3	METODE PENGUMPULAN DATA .....	56
3.4	METODE ANALISIS DATA .....	56
BAB IV .....		58
ANALISIS HASIL DAN PEMBAHASAN .....		58
4.4.1	<i>Data Yang Digunakan</i> .....	58
4.2.1	<i>Mengimport Library</i> .....	65
4.2.2	<i>Melakukan Proses Obtain Data</i> .....	66
4.4.2	<i>Scrub Data</i> .....	73
4.4.3	<i>Explore Data</i> .....	81
4.4.4	<i>Model</i> .....	89
BAB V .....		102
KESIMPULAN DAN SARAN .....		102
5.1	KESIMPULAN .....	102
5.2	SARAN .....	102
DAFTAR PUSTAKA .....		103
LAMPIRAN – LAMPIRAN .....		107

## Daftar Gambar

<i>Gambar 2. 1 Logo Python .....</i>	<i>10</i>
<i>Gambar 2. 2 Website Resmi Python.....</i>	<i>11</i>
<i>Gambar 2. 3 Instalasi Python Dengan Installer .....</i>	<i>12</i>
<i>Gambar 2. 4 Instalasi Python Dengan Installer .....</i>	<i>12</i>
<i>Gambar 2. 5 Instalasi Python Dengan Installer .....</i>	<i>13</i>
<i>Gambar 2. 6 Website Anaconda .....</i>	<i>14</i>
<i>Gambar 2. 7 Hasil Download Anaconda .....</i>	<i>14</i>
<i>Gambar 2. 8 Instalasi Anaconda.....</i>	<i>15</i>
<i>Gambar 2. 9 Instalasi Anaconda.....</i>	<i>16</i>
<i>Gambar 2. 10 Instalasi Anaconda.....</i>	<i>17</i>
<i>Gambar 2. 11 Instalasi Anaconda.....</i>	<i>18</i>
<i>Gambar 2. 12 Instalasi Anaconda.....</i>	<i>19</i>
<i>Gambar 2. 13 Instalasi Anaconda.....</i>	<i>20</i>
<i>Gambar 2. 14 Instalasi Anaconda.....</i>	<i>21</i>
<i>Gambar 2. 15 Instalasi Anaconda Selesai.....</i>	<i>22</i>
<i>Gambar 2. 16 Machine Learning Pada Web Searching .....</i>	<i>26</i>
<i>Gambar 2. 17 (Sumber : <a href="http://www.bing.com">www.bing.com</a>).....</i>	<i>27</i>
<i>Gambar 2. 18 Logo Scikit Learn .....</i>	<i>31</i>
<i>Gambar 2. 19 Mengimport Library Numpy .....</i>	<i>32</i>
<i>Gambar 2. 20 Contoh Numpy Array.....</i>	<i>32</i>
<i>Gambar 2. 21 Cek Type Numpy Array .....</i>	<i>33</i>
<i>Gambar 2. 22 Cek Type Element Array .....</i>	<i>33</i>
<i>Gambar 2. 23 Array 2 Dimensi .....</i>	<i>34</i>
<i>Gambar 2. 24 Logo Matplotlib.....</i>	<i>34</i>



<i>Gambar 2. 25 Import Library Matplotlib .....</i>	<i>34</i>
<i>Gambar 2. 26 Penggunaan Magic matplotlib .....</i>	<i>35</i>
<i>Gambar 2. 27 Perintah Membuat Bar Plot.....</i>	<i>35</i>
<i>Gambar 2. 28 Gambar Bar Plot.....</i>	<i>36</i>
<i>Gambar 2. 29 Perintah Menampilkan Scatter Plot.....</i>	<i>37</i>
<i>Gambar 2. 30 Tampilan Scatter Plot .....</i>	<i>37</i>
<i>Gambar 2. 31 Perintah Menampilkan Line Plot .....</i>	<i>38</i>
<i>Gambar 2. 32 Tampilan Line Plot.....</i>	<i>38</i>
<i>Gambar 2. 33 Mengimport Library Pandas .....</i>	<i>39</i>
<i>Gambar 2. 34 Penggunaan Fungsi Series .....</i>	<i>39</i>
<i>Gambar 2. 35 Tampilan Fungsi Series .....</i>	<i>39</i>
<i>Gambar 2. 36 Cek Tipe Pandas Series .....</i>	<i>40</i>
<i>Gambar 2. 37 Penggunaan DataFrame .....</i>	<i>40</i>
<i>Gambar 2. 38 Tampilan DataFrame .....</i>	<i>41</i>
<i>Gambar 2. 39 Penggunaan Fungsi Shape dan Hasilnya .....</i>	<i>41</i>
<i>Gambar 2. 40 Penggunaan Fungsi Info .....</i>	<i>42</i>
<i>Gambar 2. 41 Tampilan Info DataFrame .....</i>	<i>42</i>
<i>Gambar 2. 42 Penggunaan Fungsi Describe .....</i>	<i>43</i>
<i>Gambar 2. 43 Tampilan Fungsi Describe .....</i>	<i>43</i>
<i>Gambar 2. 44 Penggunaan Value_Counts .....</i>	<i>44</i>
<i>Gambar 2. 45 Tampilan Fungsi Value Counts .....</i>	<i>44</i>
<i>Gambar 2. 46 Mengakses DataFrame .....</i>	<i>44</i>
<i>Gambar 2. 47 Tampilan Akses DataFrame .....</i>	<i>45</i>
<i>Gambar 2. 48 Penggunaan Fungsi Loc .....</i>	<i>45</i>
<i>Gambar 2. 49 Tampilan Fungsi Loc .....</i>	<i>45</i>
<i>Gambar 2. 50 Penggunaan Loc Lebih Dari 1 Baris .....</i>	<i>46</i>
<i>Gambar 2. 51 Tampilan Loc Lebih dari 1 Baris .....</i>	<i>46</i>
<i>Gambar 2. 52 Perintah Menggunakan Seaborn.....</i>	<i>47</i>
<i>Gambar 2. 53 Tampilan Seaborn.....</i>	<i>47</i>

<i>Gambar 2. 54 Bar Plot Dengan Matplotlib .....</i>	<i>48</i>
<i>Gambar 2. 55 Tampilan Bar Plot Dengan Matplotlib .....</i>	<i>48</i>
<i>Gambar 2. 56 Bar Plot Dengan Seaborn .....</i>	<i>49</i>
<i>Gambar 2. 57 Tampilan Bar Plot Dengan Matplotlib .....</i>	<i>49</i>
<i>Gambar 2. 58 Tampilan Bar Plot Seaborn .....</i>	<i>49</i>
<i>Gambar 2. 59 Logo Jupyter Notebook.....</i>	<i>50</i>
<i>Gambar 3. 1 Metodologi Penelitian.....</i>	<i>53</i>
<i>Gambar 4. 1 Dataset Chicago .....</i>	<i>59</i>
<i>Gambar 4. 2 Membaca Dataset .....</i>	<i>59</i>
<i>Gambar 4. 3 Kompres Dataset .....</i>	<i>62</i>
<i>Gambar 4. 4 Kondisi IF Else .....</i>	<i>63</i>
<i>Gambar 4. 5 Import Library .....</i>	<i>65</i>
<i>Gambar 4. 6 Dataset Chicago .....</i>	<i>66</i>
<i>Gambar 4. 7 Membaca Dataset .....</i>	<i>67</i>
<i>Gambar 4. 8 Hasil Pembacaan Dataset.....</i>	<i>67</i>
<i>Gambar 4. 9 Hasil Pembacaan Dataset.....</i>	<i>68</i>
<i>Gambar 4. 10 Melihat Informasi Dataframe .....</i>	<i>68</i>
<i>Gambar 4. 11 Hasil Informasi Dataframe .....</i>	<i>69</i>
<i>Gambar 4. 12 Menghapus Data Kosong.....</i>	<i>69</i>
<i>Gambar 4. 13 Tampilan Data Kosong .....</i>	<i>70</i>
<i>Gambar 4. 14 Data Setelah Dibersihkan .....</i>	<i>70</i>
<i>Gambar 4. 15 Melihat Informasi Data .....</i>	<i>71</i>
<i>Gambar 4. 16 Tampilan Informasi Data .....</i>	<i>72</i>
<i>Gambar 4. 17 Tampilan Informasi Data .....</i>	<i>72</i>
<i>Gambar 4. 18 Melihat Informasi Kolom .....</i>	<i>73</i>
<i>Gambar 4. 19 Tampilan Informasi Kolom .....</i>	<i>73</i>
<i>Gambar 4. 20 Konversi Bulan, Hari, Jam.....</i>	<i>74</i>
<i>Gambar 4. 21 Fungsi Convert .....</i>	<i>74</i>
<i>Gambar 4. 22 Tampilan Informasi Data .....</i>	<i>74</i>

<i>Gambar 4. 23 Hasil Pembersihan Tabel Date.....</i>	<i>75</i>
<i>Gambar 4. 24 Rekayasa Fitur Date .....</i>	<i>75</i>
<i>Gambar 4. 25 Menampilkan Informasi Dataframe .....</i>	<i>76</i>
<i>Gambar 4. 26 Tampilan Informasi Dataframe .....</i>	<i>76</i>
<i>Gambar 4. 27 Tampilan Informasi Dataframe .....</i>	<i>77</i>
<i>Gambar 4. 28 Klasifikasi Kejahatan Pencurian.....</i>	<i>77</i>
<i>Gambar 4. 29 Hasil Kejahatan Pencurian.....</i>	<i>78</i>
<i>Gambar 4. 30 Melihat Informasi Dataframe .....</i>	<i>78</i>
<i>Gambar 4. 31 Tampilan Informasi Dataframe .....</i>	<i>78</i>
<i>Gambar 4. 32 Pengelompokan Data kolom.....</i>	<i>78</i>
<i>Gambar 4. 33 Tampilan Pengelompokan Data.....</i>	<i>79</i>
<i>Gambar 4. 34 Ubah Nama Kolom .....</i>	<i>79</i>
<i>Gambar 4. 35 Tampilan Perubahan Nama Kolom .....</i>	<i>79</i>
<i>Gambar 4. 36 Explore Data.....</i>	<i>81</i>
<i>Gambar 4. 37 Jumlah Data .....</i>	<i>81</i>
<i>Gambar 4. 38 Melihat Rata Rata Kejahatan Per Waktu .....</i>	<i>81</i>
<i>Gambar 4. 39 Kondisi IF Else .....</i>	<i>82</i>
<i>Gambar 4. 40 Tampilan Hasil Explore Data.....</i>	<i>82</i>
<i>Gambar 4. 41 Value Counts.....</i>	<i>83</i>
<i>Gambar 4. 42 Hasil Value Counts .....</i>	<i>83</i>
<i>Gambar 4. 43 Presentase Tingkat Kejahatan .....</i>	<i>83</i>
<i>Gambar 4. 44 Tampilan Presentase .....</i>	<i>83</i>
<i>Gambar 4. 45 Kondisi IF Else .....</i>	<i>84</i>
<i>Gambar 4. 46 Melihat Data Kejadin .....</i>	<i>84</i>
<i>Gambar 4. 47 Data Kejadian Berdasarkan Tahun .....</i>	<i>85</i>
<i>Gambar 4. 48 Berdasarkan Bulan .....</i>	<i>85</i>
<i>Gambar 4. 49 Hasil Kejadian Berdasarkan Bulan .....</i>	<i>86</i>
<i>Gambar 4. 50 Berdasarkan Hari .....</i>	<i>87</i>
<i>Gambar 4. 51 Hasil Kejadian Berdasarkan Hari .....</i>	<i>87</i>

<i>Gambar 4. 52 Berdasarkan Jam .....</i>	<i>88</i>
<i>Gambar 4. 53 Hasil Kejadian Berdasarkan Jam .....</i>	<i>88</i>
<i>Gambar 4. 54 Berdasarkan Lokasi .....</i>	<i>89</i>
<i>Gambar 4. 55 Meload Dataset .....</i>	<i>90</i>
<i>Gambar 4. 56 Menghilangkan Nilai Kosong .....</i>	<i>90</i>
<i>Gambar 4. 57 Feature Engineering .....</i>	<i>90</i>
<i>Gambar 4. 58 Hasil Testing .....</i>	<i>91</i>
<i>Gambar 4. 59 Balanced Dataset .....</i>	<i>92</i>
<i>Gambar 4. 60 Library .....</i>	<i>92</i>
<i>Gambar 4. 61 Proses Prediksi .....</i>	<i>93</i>
<i>Gambar 4. 62 Untuk Menampilkan Hasil Prediksi .....</i>	<i>93</i>
<i>Gambar 4. 63 Hasil Random Forest Array .....</i>	<i>94</i>
<i>Gambar 4. 64 Hasil Random Forest .....</i>	<i>94</i>
<i>Gambar 4. 65 Hasil Prediksi .....</i>	<i>94</i>
<i>Gambar 4. 66 Menampilkan Hasil Accuracy .....</i>	<i>95</i>
<i>Gambar 4. 67 Hasil Accuracy .....</i>	<i>95</i>
<i>Gambar 4. 68 Classification Report .....</i>	<i>95</i>
<i>Gambar 4. 69 Hasil Classification Report .....</i>	<i>95</i>
<i>Gambar 4. 70 Pohon Keputusan .....</i>	<i>97</i>
<i>Gambar 4. 71 Root Node .....</i>	<i>98</i>
<i>Gambar 4. 72 Branches .....</i>	<i>99</i>
<i>Gambar 4. 73 Leaf Node .....</i>	<i>100</i>

## Daftar Singkatan

OSEMN	= Obtain, Scrub, Explore, Model
OOP	= Object Oriented Programming
TP	= True Positif
TN	= True Negatif
FP	= False Positif
FN	= False Negatif
URL	= Uniform Resource Locator
CSV	= Comma Separated Values



# BAB I

## PENDAHULUAN

### 1.1 Latar Belakang

Kriminalitas merupakan tindakan yang bisa merugikan jiwa orang lain dan pelaku itu sendiri. Kejahatan yang dilakukan sangat berdampak pada kehidupan sosial. Tindakan ini bisa menyebabkan seseorang mengalami trauma atau pun hal-hal yang lain.[1]

Kejahatan bisa seperti pencurian di pinggir jalan atau sama berbahayanya dengan pembunuhan, dan sangat sering sebagian besar kejahatan muncul dari berbagai penyebab seperti masa anak – anak yang bermasalah, pengaruh teman seumuran yang nakal, kemiskinan, pengangguran, penyalahgunaan narkoba, dan lain – lain.[2]

Pada tindak kejahatan suatu daerah, bisa dilakukan dengan memprediksi peramalan suatu kejadian, peramalan merupakan upaya memperkirakan kemungkinan peristiwa yang mungkin terjadi selama beberapa jangka waktu. peramalan adalah tentang upaya memprediksi masa depan seakurat mungkin, mengingat semua informasi yang tersedia, dengan menggunakan data historis dan pengetahuan tentang masa depan ataupun segala peristiwa yang mungkin akan mempengaruhi peramalan. Peramalan juga dapat diartikan sebagai sebuah cara untuk mendapatkan interval berdasarkan data saat ini dan masa lalu dimana pada umumnya bertujuan untuk mengetahui tren masa depan. Hasil peramalan biasanya diterapkan untuk membantu proses pengambilan keputusan, perencanaan, dan alokasi sumber daya.[3]

Beberapa penelitian telah dilakukan untuk memprediksi jenis kejahatan, tingkat kejahatan dan hot spot kejahatan dengan menggunakan dataset kejahatan untuk wilayah yang berbeda, misalnya di Korea Selatan, dan Amerika Serikat, India. [4]

Tujuan penelitian ini dalam ruang lingkup jenis kejahatan pencurian menggunakan dataset Chicago, untuk mencari tau range kejahatan dari rendah sampai tinggi berdasarkan jumlah kejadian dan waktu kejadian menggunakan metode *OSEMN*.

Metode *OSEM* adalah tahapan tahapan yang dilakukan mulai dari pengumpulan data sampai dengan penentuan model yang digunakan untuk melakukan tingkat accuracy data yang akan digunakan.

### **1.2 Identifikasi Masalah**

Berdasarkan latar belakang diatas, maka penulis merumuskan permasalahannya sebagai berikut :

Identifikasi tingkat kejahatan menggunakan model Random Forest berdasarkan data pada website [data.cityofchicago.org](http://data.cityofchicago.org).

### **1.3 Tujuan**

Adapun tujuan melakukan penelitian diantaranya :

Untuk mengetahui tingkat kejahatan di suatu tempat tersebut menggunakan model Random Forest.

### **1.4 Manfaat Penelitian**

Adapun manfaat dari penelitian tersebut bermanfaat untuk membantu masyarakat umum dan pihak keamanan dalam melakukan pengawasan terhadap beberapa tempat, terutama untuk tempat yang memiliki tingkat kejahatan yang rendah sampai dengan tempat yang memiliki tingkat kejahatan tinggi





## BAB II

# LANDASAN TEORI

### 2.1 Penelitian Yang Berhubungan

Penelitian ini tidak terlepas dari hasil penelitian yang pernah dilakukan dan dijadikan sebagai tolak ukur serta acuan. Adapun hasil penelitian yang pernah dilakukan dan berhubungan dengan topik penelitian adalah sebagai berikut :

No	Nama Peneliti	Judul Penelitian	Hasil Penelitian
1.	Raihan Virgatama, Andri Suprayogi dan Hana Sugiastu Firdaus	Identifikasi Pengaruh Sistem Keamanan Lingkungan Terhadap Tingkat Kejahatan Pencurian di Kota Surakarta Dengan Metode Sistem Informasi Geografis	Penelitian ini menggunakan metode sistem informasi geografis dimana datanya di <i>clustering</i> terlebih dahulu, kemudian setelah dilakukan clustering dan klasifikasi data, maka hasil pengolahan datanya dan koordinatnya ditampilkan menggunakan software ArcGIS 10.4 Berdasarkan jumlah kejadian tindak kejahatan pencurian yang terjadi di Kota Surakarta selama tahun 2016-2017, terdapat 228 kasus tindak pencurian yang terdiri dari 110 kasus pencurian pada tahun 2016 dan 118 kasus pada tahun 2017. Berdasarkan jenis kejadian tindak pencurian yang terjadi selama tahun 2016-

			<p>2017 terdapat 32 kasus pencurian biasa, dan 196 kasus pencurian dengan pemberatan. Kecamatan Contoh Judul : Pembelajaran Abad 21: dari <i>E-learning</i> sampai <i>Blended Learning</i> Banjarsari adalah wilayah dengan jumlah tindak pencurian terbanyak yaitu 112 kasus yang terdiri pencurian biasa sebanyak 17 kasus, dan pencurian dengan pemberatan sebanyak 95 kasus. Sedangkan Kecamatan Pasar Kliwon dan Serengan merupakan wilayah dengan jumlah tindak pencurian paling sedikit, yaitu sembilan kasus.</p> <p>2. Berdasarkan pengolahan data tindak kejahatan pencurian dengan pemilihan jarak 750m, dihasilkan data raster yang terbagi menjadi empat kelas tingkat kerawanan dimana wilayah yang tingkat kerawanannya sangat tinggi terletak pada Kecamatan Banjarsari di kawasan GOR Manahan. Sedangkan dari data fasilitas keamanan dengan pemilihan jarak 1200m,</p>
--	--	--	--

			<p>dihasilkan data raster yang terbagi menjadi empat kelas ketersediaan fasilitas keamanan dimana wilayah yang fasilitas keamanannya sangat memadai terletak pada wilayah Kecamatan Pasar Bab 1 Pendidikan di Era Digital Kliwon yang merupakan kawasan pusat kota dan perkantoran di Kota Surakarta. 3. Berdasarkan analisis hubungan sistem keamanan lingkungan terhadap tindak kejahatan pencurian di Kota Surakarta dalam waktu dua tahun, dapat ditarik kesimpulan bahwa ketersediaan sistem keamanan lingkungan mempunyai pengaruh terhadap tingkat kerawanan pencurian di Kota Surakarta tetapi tidak mempunyai korelasi yang kuat, yaitu dengan nilai koefisiens korelasi sebesar 0,36982. [5]</p>
2.	Wajiha Safat, Sohail Asghar, Saira Andleeb Gillani	Empirical Analysis for Crime Prediction and Forecasting Using Machine Learning and Deep Learning Techniques.	<p>Berdasarkan hasil penelitian dan pembahasan, maka didapatkan kesimpulan sebagai berikut : Pada model Logistic Regression mendapatkan accuracy</p>

			sebesar 90% pada dataset Chicago, sedangkan pada Los Angeles mendapatkan accuracy sebesar 48%. Pada model XGBoost, model KNN dan model MLP mendapatkan accuracy yang optimal.[6]
3.	Marchell Rianto, Roni Yunis	Analisis Runtun Waktu Untuk Memprediksi Jumlah Mahasiswa Baru Dengan Model Random Forest	<p>Pada penelitian ini menjelaskan tentang peramalan prediksi jumlah mahasiswa yang mendaftar berdasarkan runtun waktu, dengan menggunakan metode random forest dapat memprediksi seberapa besar jumlah mahasiswa yang masuk berdasarkan tahun.</p> <p>Berdasarkan hasil analisis peramalan yang telah dilakukan maka dapat diperoleh kesimpulan sebagai berikut: Berdasarkan hasil penelitian didapatkan hasil forecasting dengan model random forest 3 memiliki variabel penting dalam prediksi yaitu Jumlah.Grade.B, Lalu parameter terbaik dari model terdapat pada maxnodes 100 dan ntree 900 yang menunjukkan parameter tersebut adalah parameter yang sangat</p>

			akurat. Peramalan jumlah mahasiswa baru dalam beberapa tahun ke depan akan menurun dengan tingkat akurasi sebesar 98% karena nilai MSE dan MAE dari model sangat baik, yaitu di bawah 5%. Dari hasil peramalan disarankan bagi penelitian selanjutnya dapat mengkombinasi model Random Forest dengan model lain yang mendukung agar tingkat keakuratan dan visualisasi data yang lebih baik dan jelas lagi. [7]
--	--	--	--

## 2.2 Landasan Teori

### 2.2.1 Kriminalitas

Kriminalitas merupakan tindakan yang bisa merugikan jiwa orang lain dan pelaku itu sendiri. Kejahatan yang dilakukan sangat berdampak pada kehidupan sosial. Tindakan ini bisa menyebabkan seseorang mengalami trauma atau pun hal-hal yang lain. Kejahatan yang dilakukan bisa berupa perjudian, perampokan, pembunuhan, dan lain-lain. Pelaku tindak pidana harus diproses oleh pihak yang berwajib dan dihukum sesuai dengan kejahatan yang dilakukan. Secara hukum, kejahatan didefinisikan sebagai tindakan atau kelalaian yang dilarang oleh hukum yang dapat dihukum dengan pidana penjara dan atau denda. Pembunuhan, perampokan, pencurian, pemerkosaan, mengemudi, mabuk, pembuangan anak dan sebagainya. [1]

### 2.2.2 Pencurian

Pencurian merupakan kejahatan yang ditujukan terhadap harta benda dan paling sering terjadi di dalam masyarakat. Kejahatan ini merupakan tindakan kejahatan yang dapat mengguncangkan stabilitas keamanan baik terhadap

harta maupun terhadap jiwa masyarakat. Oleh karena itu, baik dalam Kitab Undang - Undang Hukum Pidana (KUHP) maupun dalam Kitab Suci melarang keras tindakan kejahatan tersebut dan menegaskan ancaman hukuman secara rinci dan berat atas diri pelanggarnya. Hal ini dapat dilihat dari bentuk hukuman dan ancaman hukuman yang dijatuhkan.[8]

Adapun mengenai ancaman hukuman tentang kejahatan pencurian dalam hukum pidana positif di Indonesia diatur dalam Kitab Undang-Undang Hukum Pidana (KUHP) Buku Kedua Bab XXII tentang kejahatan terhadap harta benda dari Pasal 362 sampai dengan Pasal 367 KUHP. [8]

### 2.2.3 Python

Python merupakan bahasa pemrograman interpretatif multiguna dengan filosofi perancangan yang berfokus pada tingkat keterbacaan kode. Python diklaim sebagai bahasa yang menggabungkan kapabilitas, kemampuan, dengan sintak kode yang sangat jelas, dan dilengkapi dengan fungsionalitas pustaka standar yang besar serta komprehensif. Python juga didukung oleh komunitas yang besar. Python mendukung konsep paradigma *object oriented programming* (OOP). pemrograman *imperatif*, dan pemrograman fungsional. Saat ini kode python dapat dijalankan di berbagai platform sistem operasi, beberapa diantaranya adalah:

- a. Linux/Unix
- b. Windows
- c. Mac OS
- d. Java Virtual Machine
- e. Amiga
- f. Palm
- g. Symbian (untuk produk-produk Nokia)

Python didistribusikan dengan beberapa lisensi yang berbeda dari beberapa versi, namun pada prinsipnya python dapat diperoleh dan dipergunakan secara bebas, bahkan untuk kepentingan komersial. Lisensi Python tidak bertentangan baik menurut definisi *open source* maupun *General Public License* (GPL). [9]



*Gambar 2. 1 Logo Python*

Selain itu, python mempunyai beberapa kelebihan dibandingkan dengan bahasa pemrograman lain diantaranya :

a. *Readability*

Kelebihan yang pertama yaitu *readability* atau keterbacaan. Artinya bahasa pythob mudah dibaca karena memiliki *cript code* yang sederhana dan mudah ditulis. Hal ini akan memudahkan dalam tahap *development* aplikasi dan juga mudah dalam hal pemeliharaan.[10]

b. Efisien

Selain memiliki kelebihan keterbacaan yang baik, python juga memiliki banyak *library* yang lengkap sehingga membuat bahasa python menjadi efisien ketika membutuhkan sebuah *library* yang dibutuhkan.[10]

c. Multifungsi

Python merupakan bahasa pemrograman multifungsi, yaitu python bisa digunakan untuk berbagai keperluan dalam pembuatan dan pengembangan aplikasi, seperti pembuatan *website* , pembuatan aplikasi *machine learning*, dan pembuatan *internet of think's*.[10]

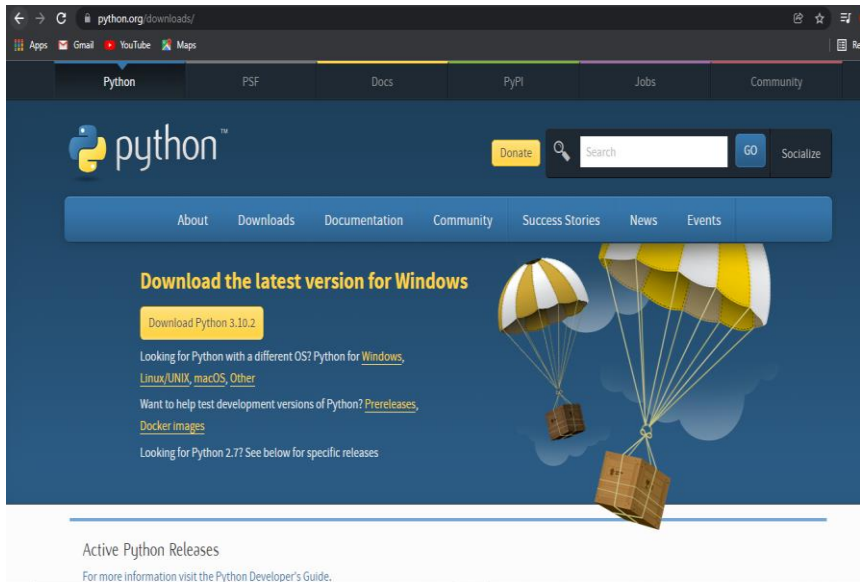
## **2.2.4 Instalasi Python Menggunakan Installer**

Cara melakukan instalasi python dengan installer yaitu :

1. Kunjungi situs <https://www.python.org/downloads/>. Kemudian unduh installer python tersebut, untuk versi terbaru pada saat buku ini dibuat terdapat pada versi 3.10.2.

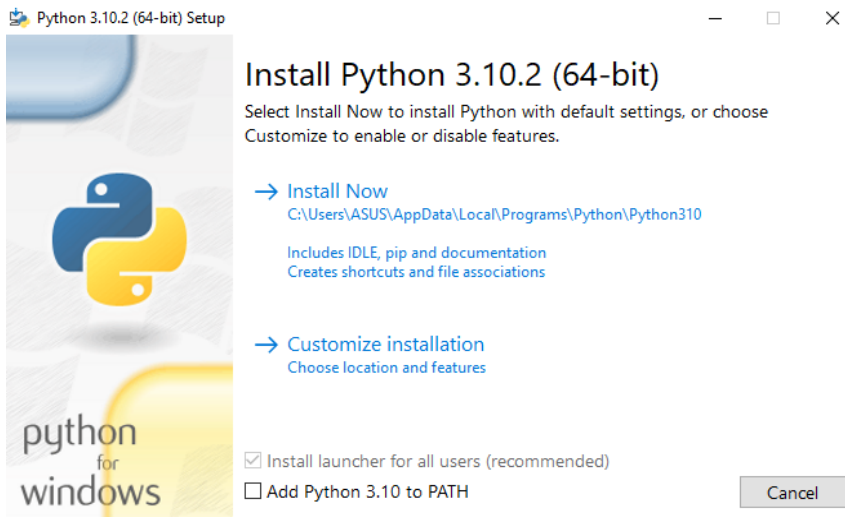


## Analisis Tingkat Kejahatan Menggunakan Model Random Forest



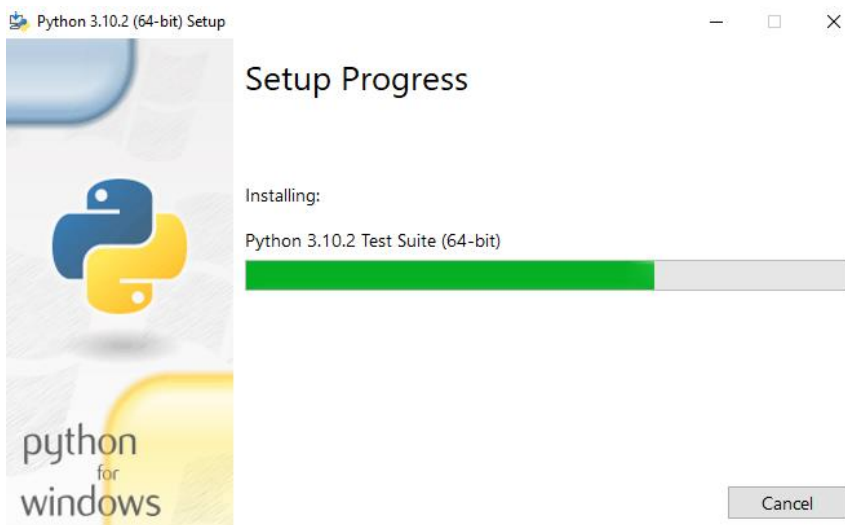
*Gambar 2. 2 Website Resmi Python*

2. Setelah di download, klik hasil download python dan tampilannya akan seperti ini. Kemudian klik install now.



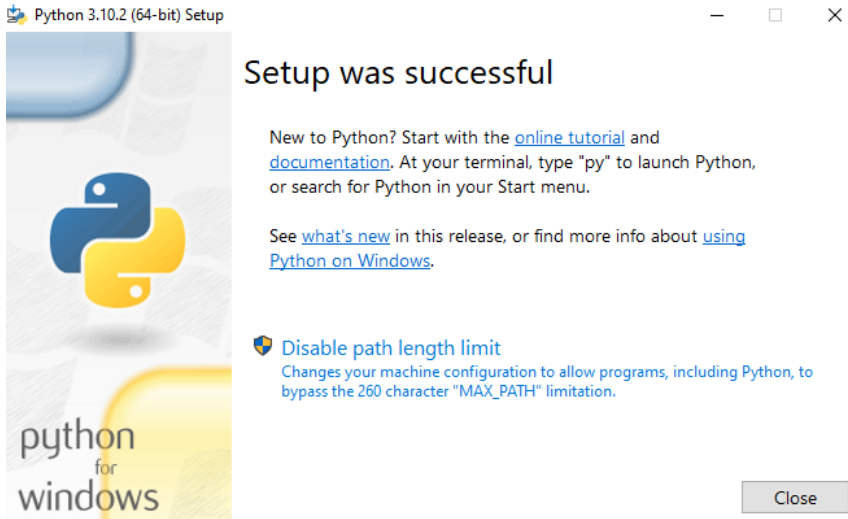
*Gambar 2. 3 Instalasi Python Dengan Installer*

3. Kemudian akan muncul tampilan seperti ini, tunggu hingga instalasi python selesai.



*Gambar 2. 4 Instalasi Python Dengan Installer*

4. Instalasi python telah berhasil, kemudian klik close.



*Gambar 2. 5 Instalasi Python Dengan Installer*

### **2.2.5 Instalasi Python dengan Anaconda**

Cara selanjutnya yaitu dengan menggunakan paket distribusi Anaconda, keuntungan dengan menggunakan Anaconda yaitu telah tersedianya beberapa tools yang digunakan untuk menjalankan python seperti Spyder dan Jupyter Notebook. Sehingga tidak harus menginstalnya secara terpisah.

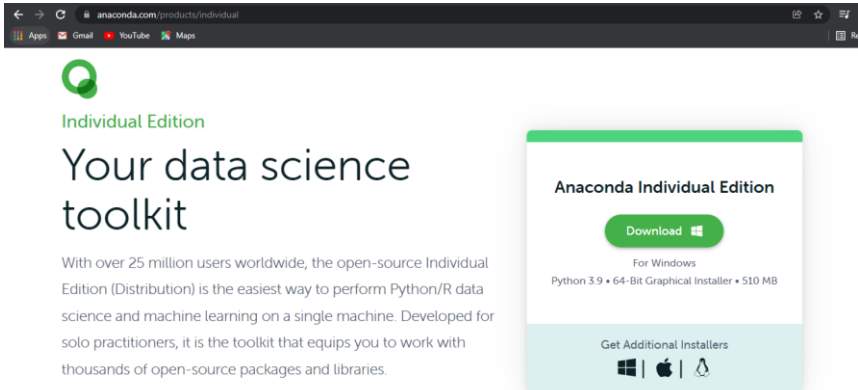
Untuk Anaconda sendiri tersedia dalam bentuk gratis sehingga memudahkan dalam penggunaan python. Selain itu, Anaconda tersedia dalam beberapa sistem operasi seperti Windows, Linux dan Mac OS.

Cara – cara melakukan instalasi Anaconda diantaranya :

1. Masuk ke situs <https://www.anaconda.com/products/individual>, kemudian klik *Download*.

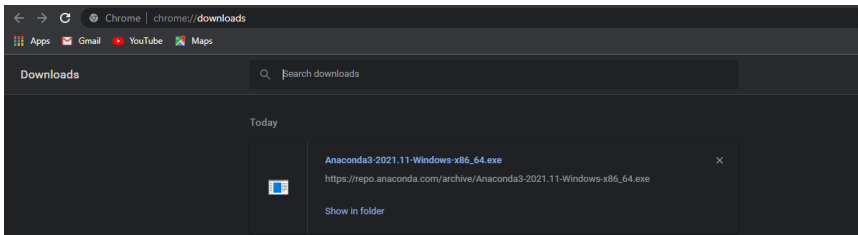
## Analisis Tingkat Kejahatan Menggunakan Model Random Forest

---



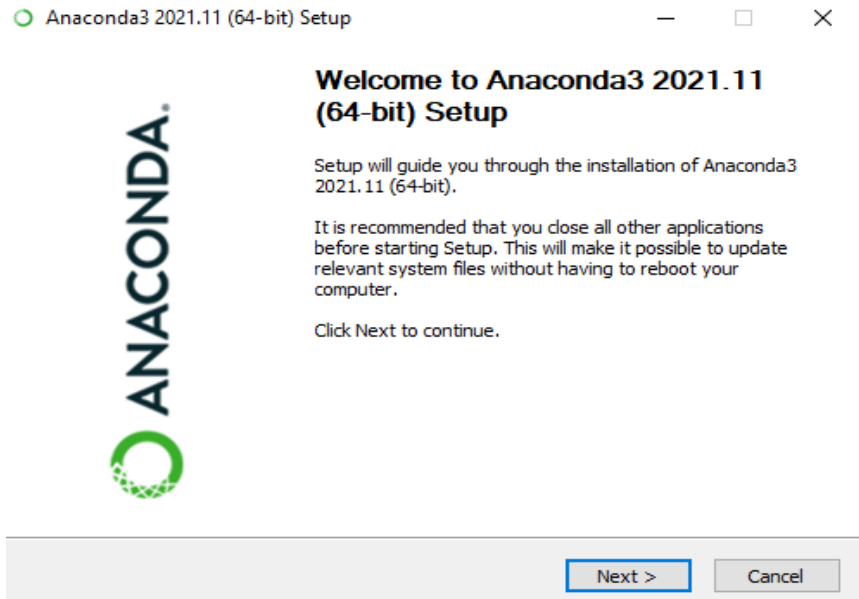
*Gambar 2. 6 Website Anaconda*

2. Selanjutnya cek hasil download pada folder Downloads, atau bisa juga dengan mengklik **ctrl + j** untuk melihat hasil download nya.



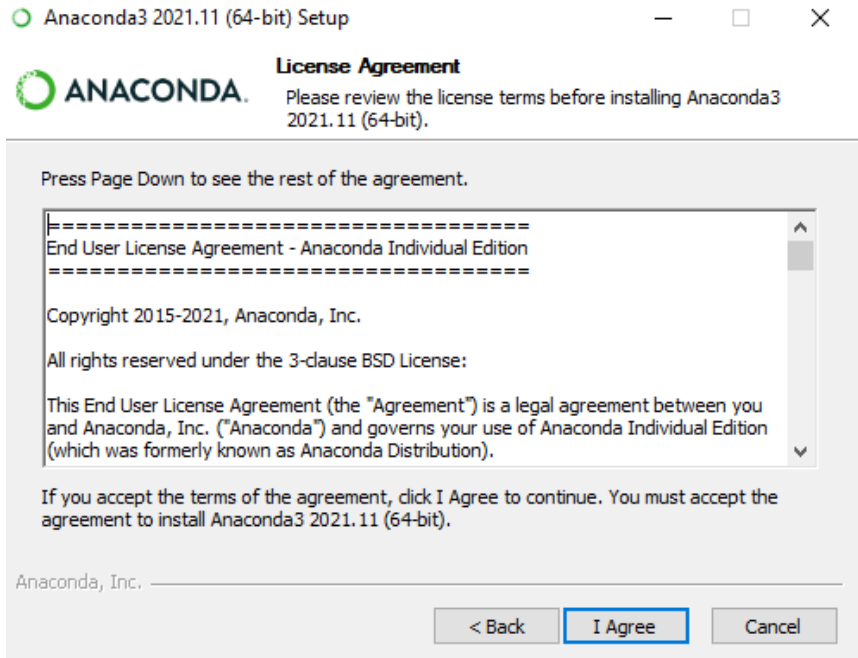
*Gambar 2. 7 Hasil Download Anaconda*

3. Selanjutnya klik hasil download Anaconda nya sampai muncul tampilan seperti ini. Kemudian klik **Next**.



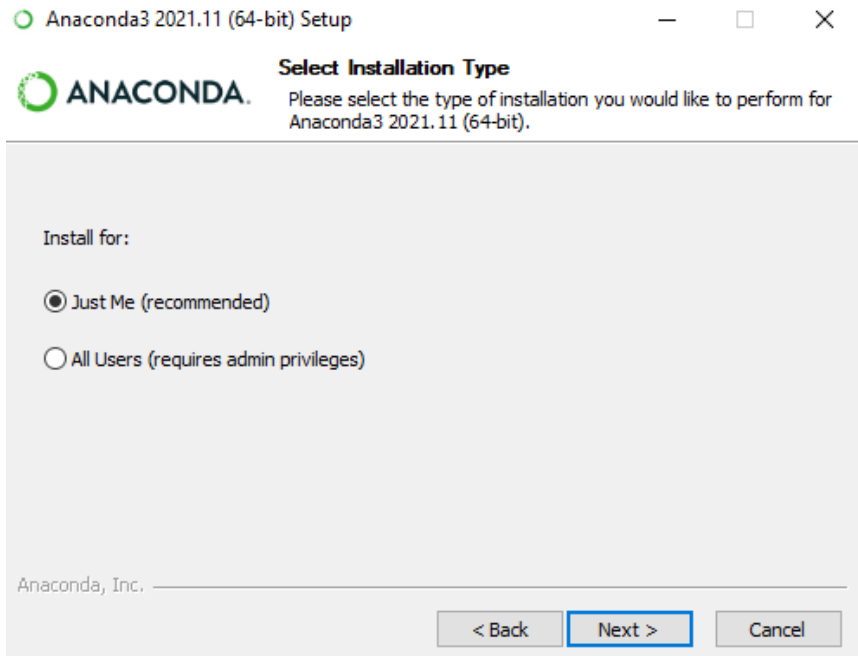
*Gambar 2. 8 Instalasi Anaconda*

4. Kemudian bacalah poin-poin terlebih dahulu, jika setuju, klik ***I Agree***



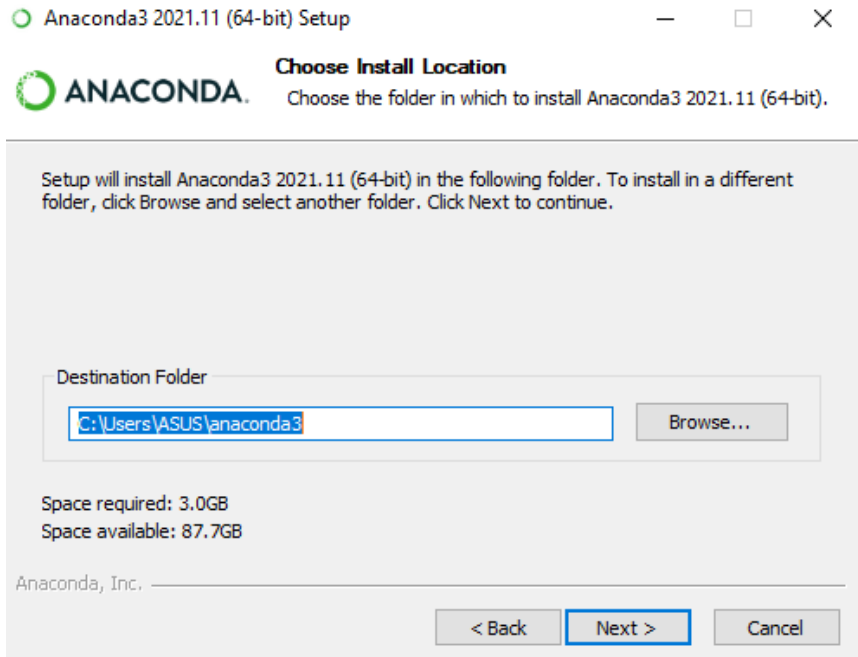
Gambar 2. 9 Instalasi Anaconda

5. Selanjutnya disini terdapat 2 kotak dialog, di proses instalasi disini menggunakan ***Just Me (recommended)*** sesuai dengan rekomendasi pada Anaconda tersebut. klik **Next**.



Gambar 2. 10 Instalasi Anaconda

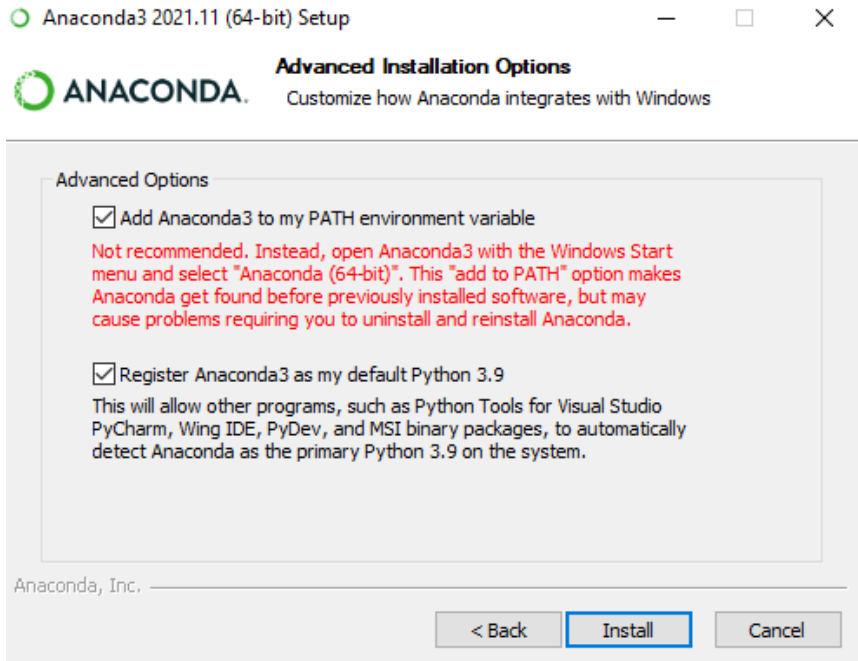
6. Selanjutnya disini muncul **Destination Folder** yang bisa digunakan untuk mengatur tempat folder itu disimpan, tetapi disini saya menggunakan versi defaultnya saja seperti pada gambar. Kemudian klik **Next**.



Gambar 2. 11 Instalasi Anaconda

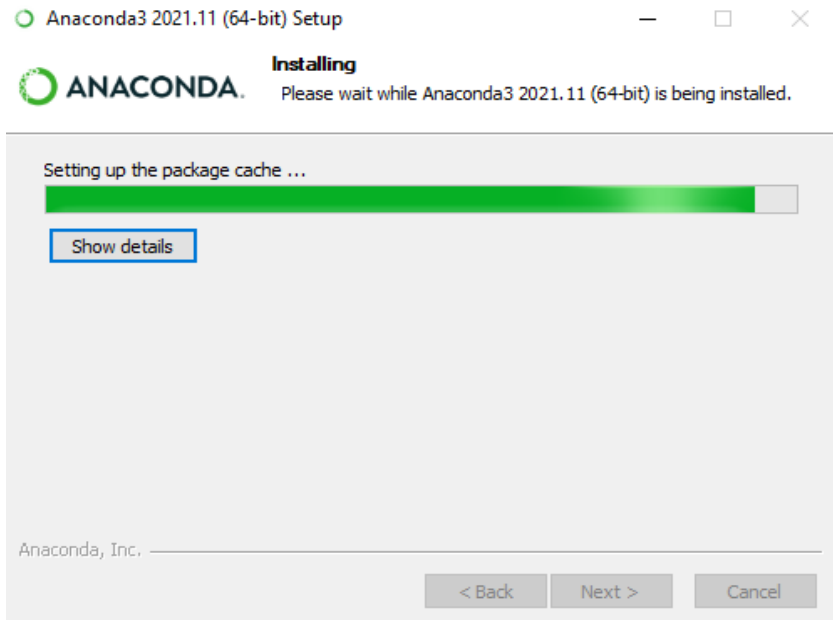
7. Selanjutnya akan muncul **Advanced Installation Options**, disini berilah tanda ceklis untuk kedua form tersebut, tujuannya agar tidak perlu menginstall lagi path yang dibutuhkan, kecuali path tersebut tidak tersedia pada saat proses instalasi. Klik **Install**.





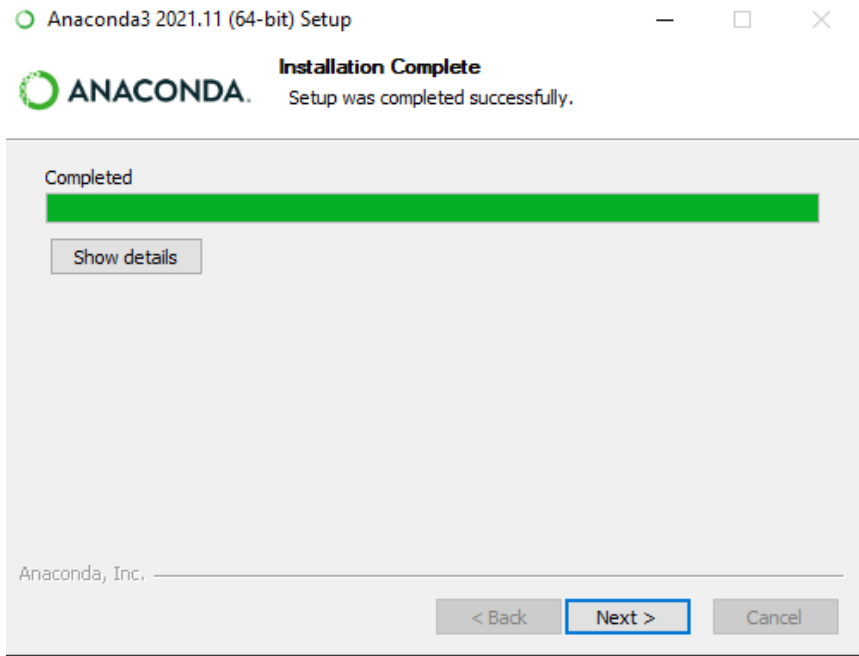
*Gambar 2. 12 Instalasi Anaconda*

8. Proses instalasi berjalan, tunggu hingga instalasi selesai.



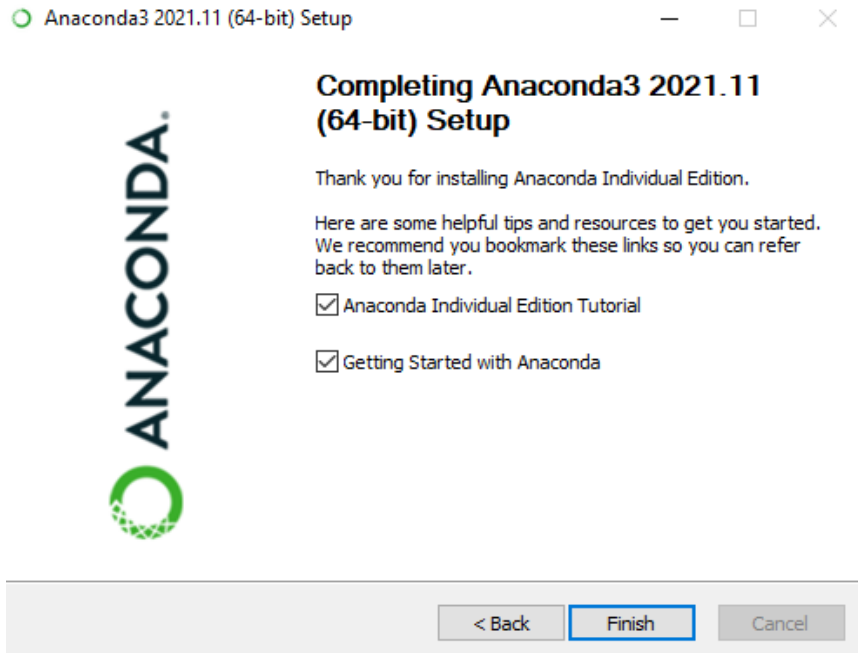
*Gambar 2. 13 Instalasi Anaconda*

9. Setelah prosesnya **Completed**. Klik **Next**.



*Gambar 2. 14 Instalasi Anaconda*

10. Proses Instalasi Anaconda selesai. Klik ***Finish***



Gambar 2. 15 Instalasi Anaconda Selesai

### 2.2.6 Machine Learning

*Machine Learning* adalah cabang aplikasi dari *Artificial Intelligence* (Kecerdasan Buatan) yang fokus pada pengembangan sebuah sistem yang mampu belajar "sendiri" tanpa harus berulang kali di program oleh manusia. Aplikasi *Machine Learning* membutuhkan data sebagai bahan belajar (*training*) sebelum mengeluarkan *output*. Aplikasi sejenis ini juga biasanya berada dalam domain spesifik alias tidak bisa diterapkan secara general untuk semua permasalahan. [11]

*Machine learning* mempunyai fokus pada pengembangan sistem yang mampu belajar sendiri untuk memutuskan sesuatu tanpa harus berulang kali di program oleh manusia. Hal tersebut menjadikan mesin tidak hanya mampu berperilaku mengambil keputusan, namun juga dapat beradaptasi dengan perubahan yang terjadi. *Machine learning* bekerja apabila tersedia data sebagai

input untuk dilakukan analisis terhadap kumpulan data besar (*big data*) sehingga menemukan pola tertentu. [12]

### Sejarah *Machine Learning*

Arthur Samuel, seorang pionir dalam pengembangan permainan komputer dan kecerdasan buatanlah yang pertama kali mengeluarkan istilah "*Machine Learning*" ke publik pada tahun 1959. Perkembangan pembelajaran mesin tumbuh berkat berkembangnya bidang kecerdasan buatan atau *Artificial Intelligence* (AI). Banyak peneliti di bidang AI tertarik untuk memiliki mesin yang dapat belajar dari data. Para peneliti berusaha untuk mendekati masalah dengan berbagai metode simbolik, serta apa yang kemudian disebut *Neural Network*, Penalaran probabilistik dan berbagai model statistik. [13]

Era sebelum 1920an Thomas Bayes, Adrien-Marie Legendre, Andrey Markov, dan matematikawan lainnya memulai penelitian yang mejadi teknik dasar pada konsep-konsep ML. Sejarah *Machine learning* dapat dikatakan dimulai oleh penelitian Thomas Bayes pada tahun 1763 yang dipublikasikan oleh temannya Richard Prince. Penelitian tersebut disempunakan oleh Pierre Simon Laplace sehingga sekarang dikenal dengan Theorema Bayes (1812). Selanjutnya penelitian penting lainnya dilakukan oleh Adrien-Marie Legendre matematikawan asal Prancis, dimana dia mengembangkan metode Least Squares yang bertujuan untuk melakukan data fitting menggunakan pendekatan aljabar (1805), serta Andrey Markov mendeskripsikan sebuah teknik yang ia gunakan untuk menganalisis puisi. Teknik tersebut dikenal dengan Markov Chains (1913). Teori-teori mereka tersebut menjadi landasan utama untuk *machine learning*. [13]

Era 1940-1950an Pada tahun 1943, McMulloh dan Pitts mengembangkan model matematis yang terinspirasi dari bagaimana otak manusia bekerja, dimana dalam model tersebut terdapat neuron-neuron yang dapat menjadi aktif atau nonaktif seperti saklar on-off. Selain itu neuron-neuron tersebut memiliki kemampuan untuk belajar dan menghasilkan respons yang berbeda berdasarkan input yang diberikan. Sumbangan besar lainnya diberikan oleh Alan Turing, pada tahun 1950 yang menciptakan mesin turing untuk menjawab pertanyaan "Dapatkan Komputer Berfikir?". Paper Alan Turing pada tahun 1950 berjudul "*Computing Machinery and Intelligence*" mendiskusikan syarat sebuah mesin dianggap cerdas. Dia beranggapan bahwa jika mesin dapat dengan sukses berperilaku seperti manusia, kita dapat menganggapnya cerdas. Pengembangan selanjutnya, pada tahun 1959, Allen

Newell dan Herbert Simon mengembangkan sebuah proyek disebut General Problem Solver (GPS). GPS telah berhasil menyelesaikan permasalahan manusia menggunakan teknik mean-ends analysis. [13]

Era 1960-1990 an Selanjutnya penemuan fenomenal lainnya pada tahun 1962, dimana Frank Rosenblatt membuktikan teorema Perceptron dan berhasil mengklasifikasi objek dari dua kelompok menggunakan classification rules. Perkembangan selanjutnya dibuat oleh Rumelhart (1986) dengan mencoba mengembangkan sistem layar tunggal (single layer) pada Perceptron menjadi sistem layar jamak (multilayers), yang kemudian disebut dengan sistem back-propagation. Konsep back-propagation telah ada sebelumnya, namun Rumelhart dan Hinton-lah yang pertama kali mempopulerkan dan membuktikan metode ini dapat dilakukan. Setelah itu, muncul beberapa model jaringan saraf tiruan lain yang dikembangkan oleh Kohonen (1972), Hopfield (1982), dan lain-lain. Selain algoritma berbasis Neural Network, di Rusia berkembang algoritma-algoritma Generalized Portrait pada tahun 1960-an (1963, 1964). Algoritma tersebut berakar dari teori pembelajar statistika yang telah dikembangkan selama 3 dekade oleh Vapnik dan Chervonenkis. SVM (Support Vector Machines) diperkenalkan oleh Vapnik, Boser, dan Guyer dalam konferensi COLT tahun 1992 dalam sebuah jurnal dan berkembang dengan sangat cepat. SVM pertama kali diperkenalkan ke dunia sejak akhir tahun 1970-an (Vapnik, 1979), Support Vector Machines tidak mendapat banyak perhatian oleh para peneliti. Kemudian Vladimir Vapnik kembali menerbitkan bukunya pada tahun 1990-an (Vapnik, 1995; Vapnik, 1998). Sejak saat itu, SVM menjadi metode yang populer dan aktif diteliti pada bidang Machine Learning. [13]

Era 2000 an di akhir 90 - an hingga pertengahan tahun 2000 an, neural networks sempat “nyaris dilupakan” dikarenakan muncul berbagai algoritma seperti Support Vector Machines, AdaBoost yang dapat dieksekusi lebih cepat dengan performa yang lebih baik pada waktu itu. Neural networks kembali mendapatkan perhatian ketika Deep Belief Networks (DBN) membuat terobosan dengan menjadi model handwritten digit recognition yang paling akurat, yang pada akhirnya memunculkan istilah Deep Learning. Istilah Deep Learning makin populer dengan diperkenalkannya Convolution Neural Network (CNN). CNN dengan arsitektur tertentu yang dipadukan dengan berbagai trik (misalnya, dropout regularization, pemanfaatan *Rectified Linear Unit* (ReLU) sebagai fungsi aktivasi, data augmentation) sehingga telah mampu melakukan klasifikasi pada data gambar yang berjumlah sangat besar

(ImageNet). ImageNet yang memiliki 1000 kategori objek dan dengan jumlah 1 juta gambar, hal ini melebihi performa manusia. [13]

*Machine learning* memiliki 3 jenis teknik yaitu teknik *supervised learning*, *unsupervised learning* dan *reinforcement learning*. Beberapa praktisi dari *machine learning* menggunakan *supervised learning*. *Supervised learning* adalah salah satu tipe *machine learning* yang menggunakan dataset yang dikenal (*training dataset*) untuk membuat prediksi. *Unsupervised learning* adalah salah satu tipe *machine learning* yang digunakan untuk menarik kesimpulan dari datasets yang terdiri dari *input data labeled response*. [14].

*Reinforcement learning* yaitu teknik ini bekerja dalam lingkungan yang dinamis di mana konsepnya harus menyelesaikan tujuan tanpa adanya pemberitahuan dari komputer secara eksplisit jika tujuan tersebut telah tercapai. [15]

*Machine learning* berdasarkan bagaimana cara kerjanya dapat dikelompokkan menjadi dua yaitu :

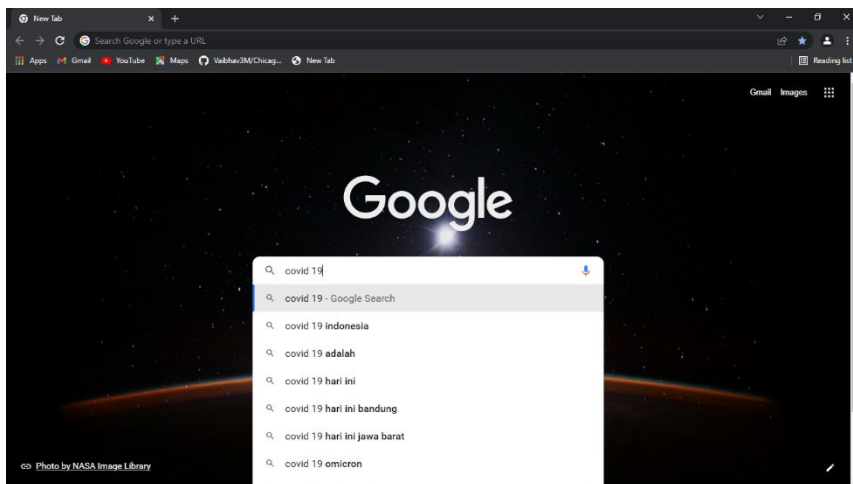
- a. *Instance-based learning* (atau sering disebut *memory-based learning*) adalah sebuah kelompok model *machine learning* yang bekerja dengan membandingkan data *testing* dengan data yang telah dipelajari pada proses *training*.
- b. *Model based learning* adalah kebalikan dari *instance-based* dimana menggunakan memori untuk melakukan pemecahan masalah, algoritma ini membuat sebuah model yang bersifat generik. [13]

*Machine learning* dan ilmu statistik sangat erat kaitannya, suatu ide dalam *machine learning* berawal dari prinsip – prinsip metodologis ke alat teoritis, dan memiliki sejarah keterkaitan yang erat dengan ilmu statistik. [16]

*Machine learning* tersebut dimanfaatkan dalam beberapa aplikasi yang digunakan diantaranya :

### a. *Web Searching*

Pada saat melakukan pencarian pada *website* pencarian, seperti google, Mozilla, opera dan sebagainya. Maka *website* tersebut akan menampilkan halaman web yang paling relevan dan mendekati dengan pencarian yang diinginkan. Dalam hal ini *machine learning* bekerja membantu menangani data yang besar dengan cara yang cerdas, dalam hal ini untuk memberikan hasil pencarian yang tepat kepada pengguna web tersebut..[10]



Gambar 2. 16 Machine Learning Pada Web Searching

### b. Filter Spam

*Machine learning* berperan sangat besar untuk melakukan penyaringan (filter) pesan spam (menganalisa, menilai dan menyaring spam berdasarkan isi ) terutama yang sering di terjadi pada email. model pohon keputusan (*decision tree*) merupakan cikal bakal dari algoritma spam filtering, untuk menentukan suatu pesan termasuk spam atau bukan. [13]

### c. Penerapan pada *Marketplace*

Pada penerapannya dengan memprediksi pergerakan dari pelanggan, dari mulai melakukan pencarian sampai dengan mendeteksi validasi dari sebuah transaksi. Sehingga pada saat kegiatan tersebut bisa memunculkan sebuah



rekomendasi produk sehingga sistem bisa memberikan rekomendasi produk kepada setiap akun *marketplace*. [13]

### d. *Virtual Assistant*

Berbagai teknologi seperti *handphone* dan laptop telah dilengkapi dengan *virtual assistant*, misalnya: Cortana di Microsoft Windows, Siri di Iphone, dan Google Now di Android. *Virtual assistant* dapat membantu pengguna untuk melakukan pencarian di internet, menanyakan jalan, melihat prakiraan cuaca, melakukan panggilan telepon dan lain – lain. Di Windows, Cortana bahkan mempelajari pengguna untuk memberikan rekomendasi perintah yang akan dilakukan. [13]

### e. *Image Recognition*

*Machine learning* juga bisa berfungsi sebagai fitur untuk mengenali gambar yang diberikan. Pada *image recognition* bisa mengenali berbagai macam objek yang tersedia dalam sebuah gambar seperti berikut. [10]



Gambar 2. 17 (Sumber : [www.bing.com](http://www.bing.com))

### f. Mobil Kendali Otomatis

Mobil kendali otomatis merupakan penerapan serta pengembangan dari machine learning yaitu machine vision. Mobil kendali otomatis merupakan penerapan machine learning yang kompleks dan dengan resiko langsung yang tinggi. Banyak hal yang harus dipelajari oleh mobil, mulai dari rambu-rambu

lalulintas, arah dan tujuan, kondisi jalan, traffic light, kondisi manusia sekitarnya, dan sensor lainnya yang terintegrasi. [17]

g. *Videos Surveillance*

*Videos surveillance* atau pengawasan video merupakan teknologi baru yang merupakan penerapan dari *machine learning* yang disematkan pada CCTV untuk mendeteksi suatu tindak kejahatan atau kecelakaan. Di negara-negara maju, CCTV sudah digunakan untuk melakukan pencarian penjahat yang masih buron. [17]

h. *Online Fraud Detection*

*Online fraud detection* adalah metode yang digunakan untuk mendeteksi suatu transaksi digital sah atau tidak. *Online fraud detection* digunakan oleh semua bank baik bank umum maupun bank virtual seperti paypal. *Online fraud detection* menggunakan *machine learning* untuk melakukan perlindungan (*cybersecurity*) terhadap pencucian uang, pendeteksian transaksi palsu, hingga deteksi pembobolan akun bank *digital*. [17]

### 2.2.7 Data Science

*Data science* atau ilmu data adalah suatu disiplin ilmu yang khusus mempelajari data, khususnya data kuantitatif (data numerik), baik yang terstruktur maupun tidak terstruktur. Bidang ilmu data telah muncul dalam menanggapi peningkatan jumlah data. Sejumlah besar data telah tersedia untuk orang-orang di semua lapisan masyarakat, melalui jejaring sosial, perangkat seluler, dan berbagai perangkat sensor Internet of Things. [18]

Berikut beberapa pilar yang ada pada *data science*

a. Bisnis

Seorang data scientist yang mengolah data berdasarkan ilmu data science harus bisa mengolah data menjadi informasi yang bisa dipahami untuk membantu perancangan strategi guna menyelesaikan masalah bisnis. Untuk bisa melakukan ini, keahlian data *science* pun harus disertai pemahaman bisnis sehingga penyelesaian yang diusulkan berdasarkan data mungkin untuk dilakukan sebuah bisnis untuk mencapai tujuannya. [19]

### b. Matematika dan Statistika

Data science sangat membutuhkan ilmu matematika, karena data harus diolah secara kuantitatif. Banyak permasalahan dalam bisnis yang dapat diselesaikan dengan membuat model analitik dengan dasar matematika. Untuk membuatnya, dibutuhkan pemahaman matematika yang mendalam.

Contohnya, model untuk merancang *machine learning* sebagai salah satu aplikasi ilmu data science sangat lekat dengan matematika. Statistik untuk data science adalah hal yang tak kalah penting. Tidak hanya mengerti statistika klasik, seorang data scientist juga perlu memahami statistika Bayes. [19]

### c. Teknologi

Data science tidak bisa lepas dari teknologi dan kreativitas serta kecerdasan dalam menggunakan keahlian teknis untuk menyelesaikan suatu permasalahan. Data science adalah keilmuan yang menggunakan data dalam jumlah besar dan model yang rumit, sehingga butuh keahlian ilmu komputer yang mendalam. Seorang data scientist perlu menguasai bahasa pemrograman seperti SQL, Python, R, SAS, Java, Scala, Julia, dan masih banyak lagi. Seorang data scientist harus mampu berpikir layaknya algoritma dalam memecahkan permasalahan yang paling sulit sekalipun. [19]

Selanjutnya ada beberapa penerapan yang digunakan dalam data science diantaranya :

#### a. Bidang Kesehatan

Pada tahun 2008, data science membuat tanda besar pertamanya di industri perawatan kesehatan. Staff Google menemukan bahwa mereka dapat memetakan wabah flu secara real time dengan melacak data lokasi pada pencarian terkait flu. Peta CDC yang ada tentang kasus flu yang terdokumentasi, FluView, diperbarui hanya sekali seminggu. Google dengan cepat meluncurkan alat yang bersaing dengan pembaruan yang lebih sering: Google Flu Trends. Tapi itu tidak berhasil. Pada tahun 2013, Google memperkirakan sekitar dua kali kasus flu yang benar-benar diamati. Meski begitu, ini menunjukkan potensi serius ilmu data dalam perawatan kesehatan.[20]

### b. Bidang Keuangan

perusahaan berbasis *finance* atau keuangan kini menggunakan data *science* untuk mengklasifikasikan, mengelompokkan, dan menyegmentasikan data yang mungkin menandakan pola penipuan. Hal ini diperlukan guna menghindari terjadinya proses kriminalisasi berkelanjutan di dalam sistem milik perusahaan. Sistem pakar dalam perusahaan finansial juga bisa mengkodekan data yang mampu mendeteksi penipuan dalam bentuk-bentuk yang tak terduga. [19]

### 2.2.8 Data Analytic

*Data Analytics* merupakan proses memeriksa kumpulan data untuk menarik kesimpulan tentang informasi yang dikandungnya. *Data Analytic* digunakan dalam industri komersial untuk memungkinkan organisasi membuat keputusan bisnis yang lebih tepat dan cepat. Oleh para ilmuwan atau peneliti untuk memverifikasi atau menyangkal model, teori, dan atau hipotesis ilmiah.[21]

### 2.2.9 Scikit Learn

Sklearn merupakan library open-source yang dapat digunakan sebagai Machine Learning pada bahasa pemrograman python. scikit-learn atau sklearn adalah modul untuk bahasa pemrograman python yang dibangun diatas numpy, scipy, dan matplotlib, fungsinya dapat membantu melakukan processing data ataupun melakukan *training* data untuk kebutuhan *machine learning* Sklearn merupakan *library* simple dan efisien untuk melakukan data mining serta analisis data. Sklearn memiliki fitur regresi, klasifikasi dan pengelompokan model contohnya seperti support vector machine, nearest neighbors, naive bayes dan lain – lain. library Scikit-learn dibangun menggunakan library lainnya seperti Numpy, Scipy dan matplotlib. Ada banyak fitur yang dapat digunakan dengan sklearn ini, seperti *Classification*, *Regression*, *Clustering*, *Dimensionality reduction*, Model selection, dan *Preprocessing* data [11]



Gambar 2. 18 Logo Scikit Learn

#### 2.2.10 Confusion Matrix

*Confusion Matrix* adalah tabel yang digunakan untuk menggambarkan kinerja model klasifikasi pada set data pengujian yang nilainya telah diketahui.

*Confusion matrix* sendiri relatif sederhana dan mudah untuk mengerti, tetapi terminologi yang terkait dapat membingungkan apabila belum memahami konsep dari *confusion matrix*. [22]

#### 2.2.11 Balance Dataset

*Balance Dataset* adalah *dataset* yang berisi jumlah sampel yang sama atau hampir sama dari kelas positif dan negative. [23]

#### 2.2.12 Imbalance Dataset

*Imbalance Dataset* adalah salah satu dari dua kelas lebih tinggi dari yang lain, dengan cara lain, jumlah pengamatan tidak sama untuk semua kelas dalam dataset klasifikasi. [23]

#### 2.2.13 Visualisasi Data

Visualisasi adalah suatu metode untuk memrepresentasikan suatu data atau permasalahan ke dalam format grafik atau bentuk gambar yang mudah untuk dipahami. Visualisasi data dalam bentuk gambar dan grafik akan memberikan kemudahan dalam membaca dan memahami data tersebut. Salah satu bentuk visualisasi yaitu pemetaan. [24]

Visualisasi data merupakan era baru yang merupakan sumber yang muncul dari kecerdasan, perkembangan teoretis dan kemajuan dalam pencitraan

multidimensi dengan membentuk kembali nilai potensial yang analitik dan wawasan dapat memberikan peran visualisasi. [25]

#### 2.2.14 Numpy

Numpy merupakan *library* python yang berfokus pada *scientific computing*. Numpy mempunyai kemampuan untuk membentuk objek N-dimensional array, yang mirip dengan list pada Python. Keunggulan Numpy array dibandingkan dengan list pada python yaitu konsumsi memory yang lebih kecil serta runtime yang lebih cepat. [26]

Numpy mempunyai beberapa fungsi diantaranya :

##### a. Numpy Array

Untuk menggunakan numpy, pertama – tama *import* terlebih dahulu *library* yang dibutuhkan.

```
import numpy as np
```

Gambar 2. 19 Mengimport Library Numpy

Kemudian untuk membuat array, gunakan fungsi `array()`. Didalam *library* Numpy, terdapat *upcasting*, yaitu ketika tipe data *element* array tidak sama, dilakukan penyesuaian tipe data pada yang lebih tinggi.

Berikut contoh numpy array.

```
import numpy as np
x = np.array([[1,2,3], [7,8,9]])
x

array([[1, 2, 3],
       [7, 8, 9]])
```

Gambar 2. 20 Contoh Numpy Array

Selanjutnya untuk mengecek tipe array gunakan perintah `type()` yang diikuti oleh nama variable yang digunakan.

Berikut untuk memanggil tipe numpy array.

```
import numpy as np
x = np.array([[1,2,3], [7,8,9]])
type(x)

numpy.ndarray
```

Gambar 2. 21 Cek Type Numpy Array

Numpy array adalah objek ndarray, yang merupakan singkatan dari *n-dimensional array*.

Selanjutnya untuk mengecek sebuah tipe data element pada array, gunakan fungsi dtype().

```
import numpy as np
x = np.array([[1,2,3], [7,8,9]])
x.dtype

dtype('int32')
```

Gambar 2. 22 Cek Type Element Array

Disini terlihat bahwa untuk tipe data array tersebut merupakan *integer*.

Numpy mempunyai fungsi shape yang berfungsi untuk menghasilkan sebuah tuple yang berisikan panjang sebuah array pada tiap dimensi. Jadi artinya fungsi shape ini bisa menghitung jumlah baris pada array.

Numpy bisa digunakan untuk membuat array multi dimensi, contohnya membuat array 2 dimensi seperti ini.

```
import numpy as np
x = np.array([[1,2,3], [7,8,9]])
x

array([[1, 2, 3],
       [7, 8, 9]])
```

*Gambar 2. 23 Array 2 Dimensi*

### 2.2.15 Matplotlib

Matplotlib merupakan *library* python yang berfokus pada visualisasi data seperti membuat plot grafik. Matplotlib pertama kali diciptakan oleh John D. Hunter dan sekarang telah dikelola oleh tim developer yang besar. Matplotlib dapat digunakan dalam code python, IPython shell, server aplikasi web, dan beberapa toolkit graphical user interface (GUI) lainnya.[26]

Visualisasi dari matplotlib merupakan sebuah gambar grafik yang terdapat satu sumbu atau lebih. Setiap sumbu memiliki sumbu horizontal (x) dan sumbu vertikal (y), dan data yang direpresentasikan menjadi warna dan glyphs seperti marker (contohnya bentuk lingkaran) atau lines (garis) atau poligon [26].



*Gambar 2. 24 Logo Matplotlib*

Untuk menggunakan *library* pada matplotlib gunakan perintah seperti berikut.

```
import matplotlib.pyplot as plt
```

*Gambar 2. 25 Import Library Matplotlib*



Matplotlib mempunyai *magic command* `%matplotlib inline`, untuk pengaturan pada backend matplotlib agar setiap grafik ditampilkan secara *inline*, yaitu akan ditampilkan langsung pada cell notebook.

Berikut perintahnya seperti ini.

```
import matplotlib.pyplot as plt
%matplotlib inline
```

*Gambar 2. 26 Penggunaan Magic matplotlib*

Berikut beberapa diagram yang bisa dibuat pada matplotlib.

a. Membuat Bar Plot

Bar plot berfungsi untuk membandingkan perubahan tiap waktu pada beberapa kelompok data. Bar plot sangat bagus digunakan dalam visualisasi ketika perubahan data sangat besar dibandingkan dengan line plot. Bar plot biasanya memiliki dua sumbu yaitu sumbu x untuk jenis kelompok dan sumbu y untuk proporsi kelompok. Matplotlib menyediakan fungsi `bar()` untuk mempermudah dalam visualisasi bar plot.

Berikut perintah untuk menampilkan bar plot.

```
import matplotlib.pyplot as plt
%matplotlib inline

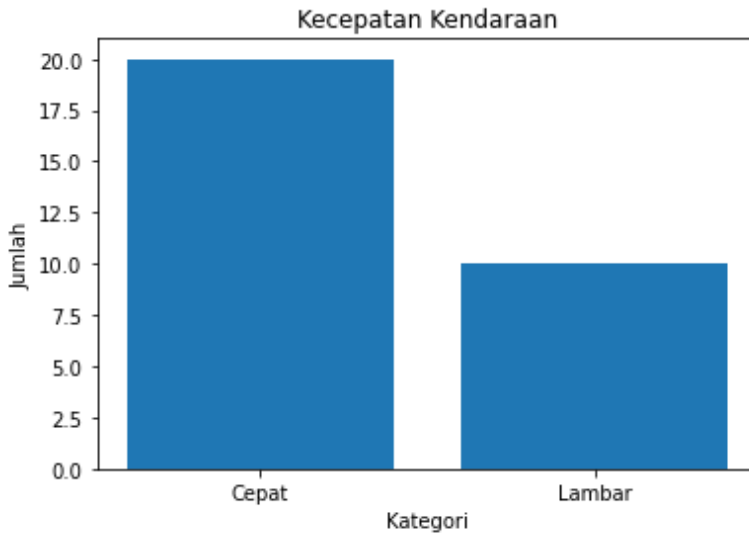
x = ['Cepat', 'Lambat']
y = [20,10]

fig, ax = plt.subplots()
ax.bar(x, y)
ax.set_xlabel('Kategori')
ax.set_ylabel('Jumlah')
ax.set_title('Kecepatan Kendaraan')
```

*Gambar 2. 27 Perintah Membuat Bar Plot*

Maka akan menampilkan bar plot seperti ini.

```
Text(0.5, 1.0, 'Kecepatan Kendaraan')
```



*Gambar 2. 28 Gambar Bar Plot*

b. Scatter Plot

Scatter plot berfungsi untuk melakukan observasi dan menunjukkan hubungan relasi antara dua variabel *numeric*. Titik-titik pada scatter plot juga dapat menggambarkan pola dari data secara keseluruhan. Matplotlib menyediakan fungsi `scatter()` untuk mempermudah dalam visualisasi scatter plot.

Berikut perintah untuk menampilkan scatter plot.

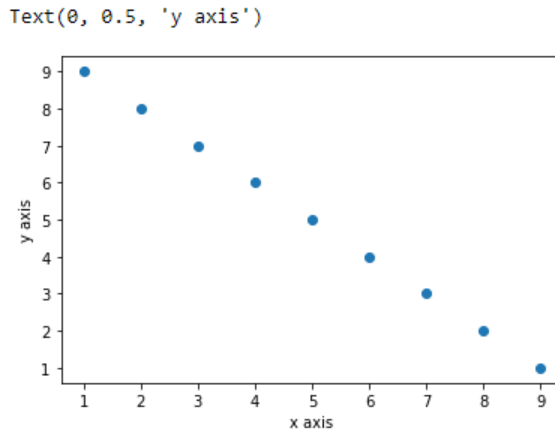
```
import matplotlib.pyplot as plt
%matplotlib inline

x = np.array([1,2,3,4,5,6,7,8,9])
y = np.array([9,8,7,6,5,4,3,2,1])

fig, ax = plt.subplots()
ax.scatter(x,y)
ax.set_xlabel('x axis')
ax.set_ylabel('y axis')
```

Gambar 2. 29 Perintah Menampilkan Scatter Plot

Berikut tampilan dari scatter plot.



Gambar 2. 30 Tampilan Scatter Plot

### c. Line Plot

*Line plot* berguna untuk melacak perubahan pada periode waktu pendek dan panjang. Ketika terdapat perubahan kecil, line plot lebih baik dalam melakukan visualisasi dibandingkan grafik bar. [26]

Berikut perintah untuk menampilkan *line plot*.

```
import matplotlib.pyplot as plt
import numpy as np
%matplotlib inline

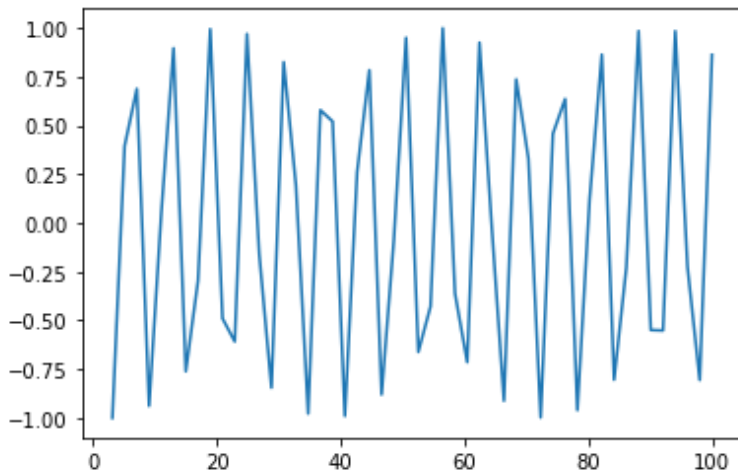
x = np.linspace(1*np.pi, 100)
y = np.cos(x)

fig, ax = plt.subplots()
ax.plot(x, y)
```

Gambar 2. 31 Perintah Menampilkan Line Plot

Berikut tampilan dari *line plot*.

```
[<matplotlib.lines.Line2D at 0x1e16e2388e0>]
```



Gambar 2. 32 Tampilan Line Plot

### 2.2.16 Pandas

Pandas merupakan *library* python yang berfokus untuk proses analisis data seperti manipulasi data, persiapan data, dan pembersihan data. Pandas

menyediakan struktur data dan fungsi high-level untuk membuat pekerjaan dengan data terstruktur/tabular lebih cepat, mudah, dan ekspresif. [26]

Pandas memadukan *library* numpy yang memiliki kemampuan manipulasi data yang fleksibel dengan database relasional (seperti SQL). Sehingga memudahkan kita untuk melakukan *reshape*, *slice* dan *dice*, agregasi data, dan mengakses subset dari data.

Untuk menggunakan *library* pandas gunakan perintah seperti ini.

```
import pandas as pd
```

Gambar 2. 33 Mengimport Library Pandas

Series adalah objek 1-dimensi yang berisi *sequence* nilai dan berasosiasi dengan label data, yang disebut indeks. Berikut contoh penggunaan series.

```
import pandas as pd
x = pd.Series([1,2,3,4,5,6,7,8,9])
x
```

Gambar 2. 34 Penggunaan Fungsi Series

Berikut tampilan fungsi series.

```
0    1
1    2
2    3
3    4
4    5
5    6
6    7
7    8
8    9
dtype: int64
```

Gambar 2. 35 Tampilan Fungsi Series

Selanjutnya untuk mengecek type dari series pandas gunakan perintah `type()`.

```
import pandas as pd
x = pd.Series([1,2,3,4,5,6,7,8,9])
type(x)

pandas.core.series.Series
```

Gambar 2. 36 Cek Tipe Pandas Series

DataFrame merupakan objek yang memiliki struktur data tabular, berorientasi pada kolom dengan label baris dan kolom. Series adalah objek array 1-dimensi yang memiliki label. DataFrame merupakan tabel data yang terdapat kolom dan baris, dimana nilai-nilai yang terdapat di dalamnya dapat berupa tipe berbeda seperti *numeric*, *string*, *boolean*, dan sebagainya.. Dataframe mirip dengan data 2 dimensi dengan adanya baris dan kolom. Selain itu, dataframe bisa dikatakan gabungan dari *dictionary* objek series yang memiliki *indeks* yang sama. Terdapat berbagai macam cara untuk membentuk objek DataFrame. Salah satu cara yang biasa dilakukan untuk membentuk objek DataFrame dengan menggunakan data masukan berupa *dictionary* [26]

Berikut contoh dalam penggunaan DataFrame.

```
import pandas as pd
data = { 'warna' : ['merah', 'kuning', 'hijau', 'biru', 'ungu'],
        'sepatu' : ['nike', 'adidas', 'puma', 'kompas', 'swallow'],
        'tahun' : [2017, 2018, 2019, 2020, 2021]
}

main_data = pd.DataFrame(data)
main_data.info()
```

Gambar 2. 37 Penggunaan DataFrame

Maka hasil tampilan dari dataframe nya seperti berikut.

	warna	sepatu	tahun
0	merah	nike	2017
1	kuning	adidas	2018
2	hijau	puma	2019
3	biru	kompas	2020
4	ungu	swallow	2021

*Gambar 2. 38 Tampilan DataFrame*

Selanjutnya ada fungsi shape yaitu fungsi untuk mengetahui jumlah baris dan kolom dari DataFrame.

Berikut contoh penggunaan shape.

```
import pandas as pd
data = { 'warna' : ['merah', 'kuning', 'hijau', 'biru', 'ungu'],
        'sepatu' : ['nike', 'adidas', 'puma', 'kompas', 'swallow'],
        'tahun' : [2017, 2018, 2019, 2020, 2021]
}

main_data = pd.DataFrame(data)
main_data.shape
```

(5, 3)

*Gambar 2. 39 Penggunaan Fungsi Shape dan Hasilnya*

Selanjutnya ada fungsi info() yang berfungsi untuk mengetahui keterangan dari objek DataFrame yang dibuat seperti index dari DataFrame lengkap dengan range dari index, jumlah kolom beserta informasi tiap kolom untuk null data dan tipe data, dan jumlah total penggunaan memory pada tiap kolom dalam satuan bytes.

Berikut contoh dalam penggunaan fungsi `info()`.

```
import pandas as pd
data = { 'warna' : ['merah', 'kuning', 'hijau', 'biru', 'ungu'],
        'sepatu' : ['nike', 'adidas', 'puma', 'kompas', 'swallow'],
        'tahun' : [2017, 2018, 2019, 2020, 2021]
}

main_data = pd.DataFrame(data)
main_data.info()
```

*Gambar 2. 40 Penggunaan Fungsi Info*

Berikut tampilan hasil fungsi `info()` pada `DataFrame`.

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5 entries, 0 to 4
Data columns (total 3 columns):
#   Column  Non-Null Count  Dtype
---  -
0   warna   5 non-null         object
1   sepatu  5 non-null         object
2   tahun   5 non-null         int64
dtypes: int64(1), object(2)
memory usage: 248.0+ bytes
```

*Gambar 2. 41 Tampilan Info DataFrame*

Maka hasilnya bisa didapat seperti pada diatas, dimana terdapat info tipe data apa saja yang digunakan, diatas didapat bahwa tipe data yang digunakan yaitu object dan integer.

Selanjutnya fungsi `describe()` yang berfungsi mengetahui statistika data untuk data *numeric* seperti *count*, *mean*, *standard deviation*, *maximum*, *minimum*, dan *quartile*. Untuk data *string*, misalkan data tersebut adalah kategori, maka dapat menggunakan fungsi `value_counts()` untuk mengetahui jumlah tiap kategori pada data.



Berikut contoh penggunaan fungsi *describe*.

```
import pandas as pd
data = { 'warna' : ['merah', 'kuning', 'hijau', 'biru', 'ungu'],
        'sepatu' : ['nike', 'adidas', 'puma', 'kompas', 'swallow'],
        'tahun' : [2017, 2018, 2019, 2020, 2021]
}

main_data = pd.DataFrame(data)
main_data.describe()
```

*Gambar 2. 42 Penggunaan Fungsi Describe*

Berikut tampilan dari fungsi *describe()*.

tahun	
count	5.000000
mean	2019.000000
std	1.581139
min	2017.000000
25%	2018.000000
50%	2019.000000
75%	2020.000000
max	2021.000000

*Gambar 2. 43 Tampilan Fungsi Describe*

Selanjutnya fungsi *value\_counts()* yang berfungsi untuk mengetahui jumlah tiap kategori pada data.

Berikut contoh penggunaan `value_counts()`.

```
import pandas as pd
data = { 'warna' : ['merah', 'kuning', 'hijau', 'biru', 'ungu'],
        'sepatu' : ['nike', 'adidas', 'puma', 'kompas', 'swallow'],
        'tahun' : [2017, 2018, 2019, 2020, 2021]
}

main_data = pd.DataFrame(data)
main_data.value_counts()
```

*Gambar 2. 44 Penggunaan Value\_Counts*

Maka tampilan dari fungsi `value_counts` seperti berikut.

warna	sepatu	tahun	
biru	kompas	2020	1
hijau	puma	2019	1
kuning	adidas	2018	1
merah	nike	2017	1
ungu	swallow	2021	1

dtype: int64

*Gambar 2. 45 Tampilan Fungsi Value Counts*

Selanjutnya untuk mengakses DataFrame pada suatu data, alamat data tersebut adalah pada nama kolom sebagai petunjuk lokasi kolom dan indeks sebagai petunjuk lokasi baris.

Berikut contoh untuk mengakses DataFrame.

```
import pandas as pd
data = { 'warna' : ['merah', 'kuning', 'hijau', 'biru', 'ungu'],
        'sepatu' : ['nike', 'adidas', 'puma', 'kompas', 'swallow'],
        'tahun' : [2017, 2018, 2019, 2020, 2021]
}

main_data = pd.DataFrame(data)
main_data['sepatu']
```

*Gambar 2. 46 Mengakses DataFrame*

Maka hasilnya seperti berikut.

```
0      nike
1     adidas
2      puma
3     kompas
4     swallow
Name: sepatu, dtype: object
```

*Gambar 2. 47 Tampilan Akses DataFrame*

Sedangkan mengakses data pada baris tertentu, kita menggunakan fungsi `loc[indeks]`. Indeks disini menunjukan baris pada DataFrame.

Berikut contoh penggunaan fungsi `loc`.

```
import pandas as pd
data = { 'warna' : ['merah', 'kuning', 'hijau', 'biru', 'ungu'],
        'sepatu' : ['nike', 'adidas', 'puma', 'kompas', 'swallow'],
        'tahun' : [2017, 2018, 2019, 2020, 2021]
}

main_data = pd.DataFrame(data)
main_data.loc[4]
```

*Gambar 2. 48 Penggunaan Fungsi Loc*

Berikut tampilan fungsi `loc`

```
warna      ungu
sepatu     swallow
tahun      2021
Name: 4, dtype: object
```

*Gambar 2. 49 Tampilan Fungsi Loc*

DataFrame bisa diakses lebih dari 1 baris dengan menggunakan titik dua ':', seperti ingin mengakses indeks 2–3, maka menggunakan perintah `loc[2:3]`.

Berikut contoh penggunaan fungsi loc dengan akses lebih dari 1 baris.

```
import pandas as pd
data = { 'warna' : ['merah', 'kuning', 'hijau', 'biru', 'ungu'],
        'sepatu' : ['nike', 'adidas', 'puma', 'kompas', 'swallow'],
        'tahun' : [2017, 2018, 2019, 2020, 2021]
}

main_data = pd.DataFrame(data)
main_data.loc[2:3]
```

*Gambar 2. 50 Penggunaan Loc Lebih Dari 1 Baris*

Berikut tampilan penggunaan loc lebih dari 1 baris.

	warna	sepatu	tahun
2	hijau	puma	2019
3	biru	kompas	2020

*Gambar 2. 51 Tampilan Loc Lebih dari 1 Baris*

### 2.2.17 Seaborn

Seaborn berfungsi untuk memproduksi visualisasi dengan Python dan memiliki beberapa kelebihan dibandingkan dengan Matplotlib yaitu, hasil visualisasi Seaborn diklaim lebih bagus dan indah juga menggunakan serangkaian kode yang lebih mudah.[27]

Untuk menggunakan seaborn yaitu harus menggunakan library numpy, pandas dan matplotlib.

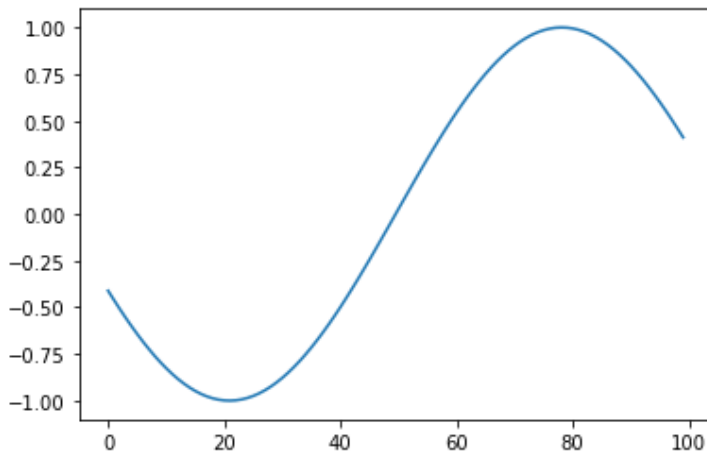
```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

data = np.linspace(-np.e, np.e, 100, endpoint=True)

plt.plot(np.sin(data))
plt.show()
```

*Gambar 2. 52 Perintah Menggunakan Seaborn*

Berikut tampilan dari seaborn.



*Gambar 2. 53 Tampilan Seaborn*

Selanjutnya bandingkan tampilan bar plot dengan matplotlib dan seaborn, maka hasilnya seperti berikut.

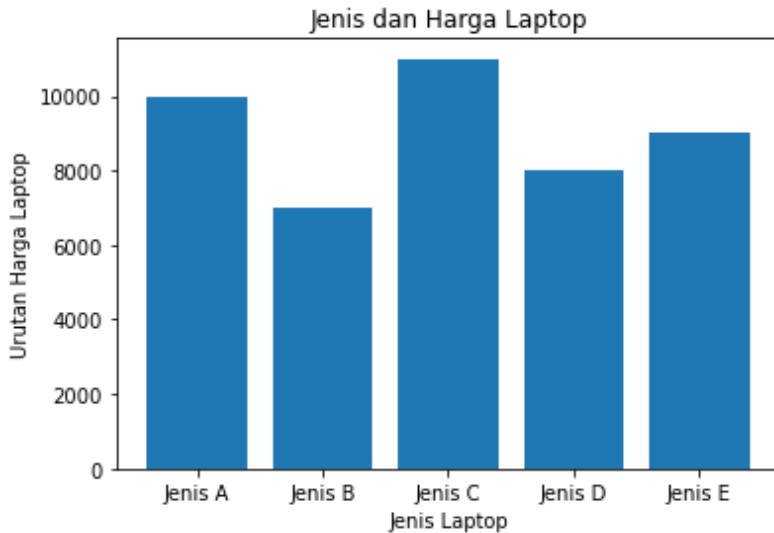
```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

laptop = ["Jenis A", "Jenis B", "Jenis C", "Jenis D", "Jenis E"]
harga_laptop = [10000, 7000, 11000, 8000, 9000]

plt.bar(laptop, harga_laptop)
plt.title("Jenis dan Harga Laptop")
plt.xlabel("Jenis Laptop")
plt.ylabel("Urutan Harga Laptop")
plt.show()
```

*Gambar 2. 54 Bar Plot Dengan Matplotlib*

Maka hasil tampilan bar plot dengan matplotlib seperti berikut.



*Gambar 2. 55 Tampilan Bar Plot Dengan Matplotlib*

Selanjutnya perintah bar plot dengan seaborn maka seperti ini.

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

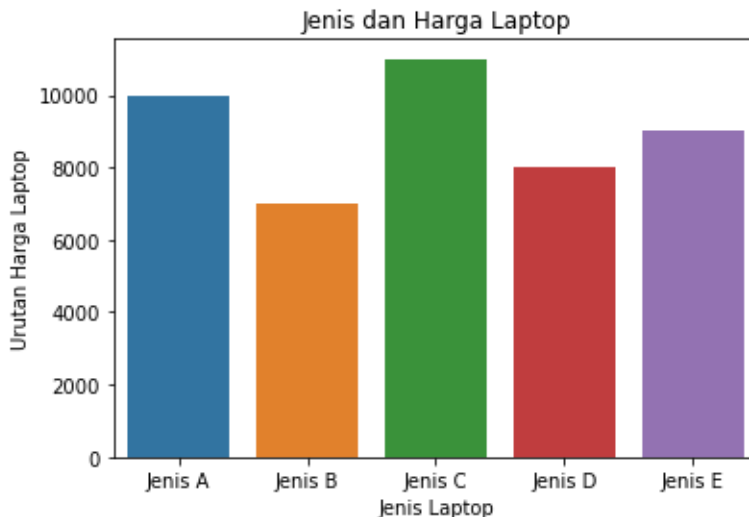
laptop = ["Jenis A", "Jenis B", "Jenis C", "Jenis D", "Jenis E"]
harga_laptop = [10000, 7000, 11000, 8000, 9000]

sns.barplot(x=laptop, y=harga_laptop)
plt.title("Jenis dan Harga Laptop")
plt.xlabel("Jenis Laptop")
plt.ylabel("Urutan Harga Laptop")
plt.show()
```

*Gambar 2. 56 Bar Plot Dengan Seaborn*

Berikut tampilan dari bar plot dengan seaborn.

*Gambar 2. 57 Tampilan Bar Plot Dengan Matplotlib*



*Gambar 2. 58 Tampilan Bar Plot Seaborn*

Terdapat perbedaan antara seaborn dengan matplotlib dimana jika menggunakan matplotlib, tampilan bar plot hanya berwarna biru sedangkan

dengan seaborn terlihat bahwa tampilan bar plotnya lebih berwarna, tidak hanya warna biru tetapi terdapat dengan warna lain.

### 2.2.18 Classification Report

*Classification Report* berfungsi untuk mengukur kualitas prediksi dari model yang digunakan. Ada prediksi yang benar dan ada juga prediksi yang salah. Lebih lanjut, ada prediksi dengan hasil *True Negatif*, *True Positif*, *False Negatif*, dan *False Positif* digunakan untuk memprediksi *metric classification report*. Beberapa penjelasan terkait *True Negatif*, *True Positif*, *False Negatif* dan *False Positif*.

- a. *TN / True Negative* = ketika sebuah kasus negatif dan diprediksi negatif
- b. *TP / True Positive*: ketika kasus positif dan diprediksi positif
- c. *FN / Negatif Palsu*: ketika sebuah kasus positif tetapi diprediksi negatif
- d. *FP / False Positive*: ketika suatu kasus negatif tetapi diprediksi positif. [28]

### 2.2.19 Jupyter Notebook



*Gambar 2. 59 Logo Jupyter Notebook*

Jupyter notebook merupakan aplikasi berbasis web yang bisa digunakan untuk membuat dan membagikan dokumen. Dokumen tersebut berisi kode, persamaan matematika, visualisasi data maupun text. Jupyter notebook ini di kelola oleh orang-orang yang tergabung pada Project Jupyter.



Jupyter notebook merupakan project spin-off dari IPython, yang pada mulanya memiliki proyek tersendiri yaitu Notebook IPython. Memiliki nama Jupyter karena dapat mendukung bahasa pemrograman Julia, Python dan R. Jupyter disajikan dengan kernel IPython, sehingga memungkinkan untuk menulis program dengan menggunakan Python. Namun, pada saat ini ada lebih dari 100 kernel lainnya yang dapat digunakan. [29].

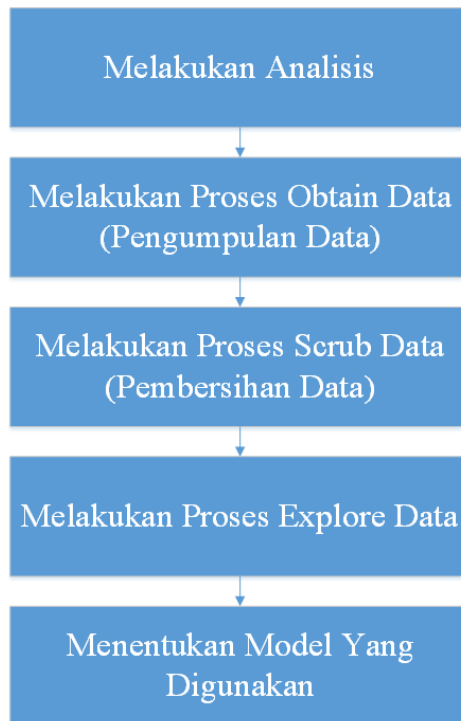


## BAB III

# METODE PENELITIAN

### 3.1 Metodologi Penelitian

Pada penelitian ini, penulis menggunakan metode *OSEMN* yang digunakan untuk mendapatkan analisis terbaik dari data yang disajikan. Untuk menyelesaikan permasalahan yang ada, diperlukan suatu metode penelitian sebagai berikut :



*Gambar 3. 1 Metodologi Penelitian*

### 3.2 Tahapan – Tahapan Diagram Alur Metodologi Penelitian

Tahapan – tahapan dari metodologi penelitian dapat diuraikan dan dijelaskan lebih detail seperti berikut :

#### 3.2.1 Melakukan Analisis

Pada tahap ini, akan dibahas mengenai penelitian yang dilakukan, dari analisis kebutuhan data. Untuk data disini menggunakan data yang disediakan dari *website* Chicago. Berikut *link url* nya (<https://data.cityofchicago.org/Public-Safety/Crimes-2001-to-Present/ijzp-q8t2/data>).

Kemudian melakukan proses pengumpulan data dengan cara mengunduh data pada *website* tersebut. setelah itu masuk ke proses *scrub data* (pembersihan data) , proses pembersihan data dilakukan dengan 3 tahapan yaitu menghilangkan data kosong, melakukan rekayasa fitur dan melakukan kompres data.

Kemudian melakukan explore data pada data yang ada untuk digunakan sesuai dengan kebutuhan.

Kemudian melakukan penentuan model untuk mencari tingkat akurasi dari model yang digunakan. Pada tahap ini menggunakan model *random forest*.

#### 3.2.2 Obtain Data (Mengumpulkan Data)

Dalam bidang *data science*, proses mengumpulkan data wajib dan harus menggunakan data yang usable atau dapat digunakan untuk tujuan analisis dengan melihat jenis dan bentuk datanya, cara memperoleh datanya dan sumber data.[30]

Dalam tahap awal melakukan penelitian. penulis mengidentifikasi dan mengumpulkan data yang tepat dalam penelitian. Jika tidak memiliki data maka tidak ada yang dapat dilakukan. Banyak cara untuk mendapatkan dataset yaitu dari situs kaggle.com, data.cityofchicago.org, dan masih banyak lagi repository dataset lainnya. Dalam penelitian ini, penulis menggunakan dataset yang disediakan oleh data.cityofchicago.org yang membahas tentang criminal.

#### 3.2.3 Scrub Data (Pembersihan Data)

*Scrubbing Data* adalah kegiatan melakukan konversi atau merapikan data dari satu format ke format lain dan menggabungkan semuanya ke dalam satu format standar di agar memberikan gambar yang akurat pada hasil akhir.[7]

Tahapan selanjutnya yaitu melakukan *scrub data* (pembersihan data) dimana pada bagian ini penulis menguraikan dan beberapa subbab diantaranya :

- a. *Getting rid of missing value* (menghilangkan nilai yang kosong)  
Pada tahap ini akan memeriksa data dan membersihkan data dari data yang tidak diperlukan dengan menggunakan Bahasa pemrograman Python.
- b. *Feature Engineering* (Rekayasa Fitur)  
Pada tahap ini peneliti akan merekayasa beberapa fitur yang sesuai dengan bahasa yang digunakan, disini menggunakan Python sebagai bahasa pemrograman yang digunakan.
- c. *Compressing the Dataset* (Melakukan kompres data)  
Tahapan selanjutnya yaitu melakukan kompres dataset yang didapat menjadi dataset yang baru.[2]

### 3.2.4 Explore

*Explore* data adalah langkah di mana kita membiasakan diri dengan data dan memungkinkan kita untuk mencari tahu sub kumpulan data mana yang akan digunakan untuk pemodelan lebih lanjut serta membantu dalam pembuatan hipotesis untuk dieksplorasi.[7]

### 3.2.5 Model

Pada Tahap akhir ini dilakukan proses *modeling* data dengan menggunakan beberapa model yang digunakan untuk mengklasifikasi data. Tahap *Modelling* data adalah proses memvisualisasikan, mengelompokkan, dan melakukan pengurangan dimensi model dari data. Dengan melihat pola dan tren data yang unik dapat dijadikan model terbaik di mana model terbaik adalah model yang prediktif hasilnya.[30]

### 3.2.6 Random Forest

Random Forest merupakan salah satu *ensemble learning* yang dibangun dari pohon keputusan . Beberapa keuntungan menggunakan pendekatan *ensemble learning* adalah dapat digunakan untuk kasus klasifikasi dan regresi, mampu memperoleh akurasi tinggi, cocok untuk analisis ukuran dataset yang besar

dengan banyak dimensi. Selain itu ensemble learning seperti Random Forest juga cocok digunakan untuk menangani *imbalanced* data.[31]

Random forest mempunyai beberapa keunggulan,yaitu dapat meningkatkan akurasi apabila terdapat data yang hilang serta untuk resisting outliers, dan juga efisien untuk penyimpanan data. Tidak hanya itu, pada Random Forest terdapat proses seleksi fitur dimana mampu mengambil fitur terbaik sehingga meningkatkan performa pada model klasifikasi. Dengan adanya fitur seleksi tentunya Random Forest mampu bekerja pada data yang besar dengan parameter yang kompleks secara efektif. Selain itu, Random Forest juga mampu bekerja secara parallel yang dikenal dengan multiple random forest. Namun, Random Forest terkadang memiliki nilai yang tidak diharapkan dan juga tidak memprediksi rangedari responsenilai pada data latih. [32]

Random Forest adalah model yang dapat meningkatkan hasil akurasi karena dalam membangkitkan simpul anak untuk setiap node dilakukan secara acak. Algoritma klasifikasi Random Forest merupakan pengembangan dari Decision Tree, dimana menghasilkan pohon gabungan yang memberikan tingkat akurat yang lebih tinggi dibandingkan dengan pohon tunggal. [33]

### **3.3 Metode Pengumpulan Data**

Metode pengumpulan data yang dilakukan pada penelitian ini, penulis menggunakan data yang terdapat pada web Chicago. Dataset yang berisi tentang data criminal yang dilakukan pada daerah tersebut.

### **3.4 Metode Analisis Data**

Analisis data kuantitatif merupakan sebuah teknik analisis yang digunakan pada data kuantitatif. Data kuantitatif merupakan data yang dapat dibentuk dengan simbol angka atau bilangan. Metode ini merupakan pendekatan pengolahan data melalui metode statistik atau matematik yang terkumpul dari data sekunder. Kelebihan dari metode ini adalah kesimpulan yang lebih terukur dan komprehensif. Hasil dari analisis kuantitatif biasanya dalam bentuk angka yang kemudian akan diinterpretasikan dalam uraian-uraian kalimat yang dapat dipahami oleh pengguna. Analisis data dalam penelitian ini biasanya menggunakan dua macam teknik analisis statistik, yaitu statistik deskriptif dan statistik inferensial.

Metode analisis data dalam penelitian kuantitatif menggunakan statistik. Metode analisis data kuantitatif adalah metode yang bergantung kepada kemampuan untuk menghitung data secara akurat. Selain itu, metode ini juga memerlukan kemampuan untuk menginterpretasikan data yang kompleks. Beberapa contoh metode analisis kuantitatif, seperti analisis deskriptif, regresi, dan faktor. Metode analisis data kuantitatif mempunyai berbagai macam jenis analisis seperti teknik korelasional, regresi, komparasi, deskriptif dan sejenisnya. [34]

## BAB IV

# ANALISIS HASIL DAN PEMBAHASAN

### 4.1 Analisis

Pada analisis ini akan menjelaskan tentang Analisis tentang tingkat kejahatan menggunakan model *random forest*. Pada tahapan penelitian dilakukan analisis kebutuhan data, selanjutnya melakukan analisis data yang digunakan pada penelitian. Pada analisis data yang digunakan, dilakukan proses klasifikasi data dari urutan jenis kejahatan rendah – tinggi, menentukan data apa saja yang digunakan.

*obtain data* (pengumpulan data) yang dimana data tersebut didapatkan dalam *website* Chicago yaitu (<https://data.cityofchicago.org/Public-Safety/Crimes-2001-to-Present/ijzp-q8t2/data>).

Setelah itu masuk ke proses *scrub data* (pembersihan data) , proses pembersihan data dilakukan dengan 3 tahapan yaitu menghilangkan data kosong, melakukan rekayasa fitur dan melakukan kompres data.

Kemudian melakukan explore data pada data yang ada untuk digunakan sesuai dengan kebutuhan.

Kemudian melakukan penentuan model untuk mencari tingkat akurasi dari model yang digunakan. Pada tahap ini menggunakan model *random forest*.

#### 4.4.1 Data Yang Digunakan

Data yang digunakan pada penelitian ini menggunakan dataset kejahatan yang berada di daerah chicago. Untuk mendapatkan dataset tersebut bisa mengunjungi web <https://data.cityofchicago.org/Public-Safety/Crimes-2001-to-Present/ijzp-q8t2/data> . Data yang digunakan penulis data rentan tahun 2019 – 2021 yang sudah dalam bentuk kompres data. Berikut gambaran dari tampilan data yang digunakan.





shape = digunakan untuk untuk mengetahui jumlah baris dan kolom dari DataFrame

print = fungsi untuk menampilkan hasil dari suatu parameter.

Berikut penjelasan tentang dataset Chicago sebagai berikut :

1. *ID* : Merupakan nomor id dari dataset
2. *Case Number* : Merupakan case number atau number perkara pada dataset
3. *Date* : Merupakan tanggal kejadian dari sebuah peristiwa pada dataset.
4. *Primary Type* : Merupakan daftar kejahatan yang tersedia pada dataset.
5. *Description* : Merupakan deskripsi kejahatan tersebut.
6. *Location Description* : Merupakan deskripsi lokasi kejadian.
7. *Arrest* : Merupakan kode penangkapan, apakah dilakukan penangkapan atau tidak, dengan kode *TRUE* atau *FALSE*
8. *Domestic* : Merupakan kode domestic
9. *Beat* : Merupakan kode jumlah kejahatan
10. *District* : Merupakan kode distrik atau tempat
11. *Ward* : merupakan kode bangsal atau tempat tinggal
12. *Community Area* : Merupakan kode area dari suatu wilayah
13. *FBI Code* : Merupakan kode kejahatan yang kode tersebut dibawah fbi
14. *X Coordinate* : Merupakan titik koordinat X
15. *Y Coordinate* : Merupakan titik koordinat Y
16. *Year* : Merupakan tahun kejadian
17. *Update On* : Merupakan Update tanggal kejadian.
18. *Latitude* : Merupakan kode Lintang.
19. *Longitude* : Merupakan kode bujur.
20. *Location* : Merupakan kode lokasi.

## Analisis Tingkat Kejahatan Menggunakan Model Random Forest

---

Setelah didapatkan data, maka selanjutnya data tersebut di bersihkan dan direkayasa beberapa fitur, sehingga data yang tersedia terdapat *Month*, *Day*, *Hour*, *District*, *Jumlah Kejahatan*, *Alarm* (Jenis Kejahatan).

	Month	Day	Hour	District	Jumlah_Kejahatan	Alarm
25709	7	6	10	31	1	0
36724	10	6	3	31	1	0
35145	10	3	23	31	1	0
28336	8	4	0	31	1	0
36199	10	5	2	31	1	0
36200	10	5	7	31	1	0
22568	7	0	12	31	1	0
13676	4	4	15	31	1	0
528	1	0	0	31	1	0
7360	2	6	12	31	1	0
3156	1	5	11	31	1	0
42461	12	3	11	31	1	0
40377	11	6	12	31	1	0
23090	7	1	9	31	1	0
21517	6	5	7	31	1	0

11568	4	0	16	25	10	0
27269	8	2	3	25	2	0
11569	4	0	17	25	18	1
11567	4	0	15	25	11	0

Berikut *script code* untuk melakukan proses kompres dataset.

```
def crime_rate_assign(x):  
    if(x<=14):  
        return 0  
    elif(x>14 and x<=33):  
        return 1  
    else:  
        return 2  
cri6['Alarm'] = cri6['Jumlah_Kejahatan'].apply(crime_rate_assign)  
cri6 = cri6[['Month','Day','Hour','District','Jumlah_Kejahatan','Alarm']]  
cri6.head()
```

*Gambar 4. 3 Kompres Dataset*

def = merupakan sebuah fungsi dalam python diikuti nama fungsinya, disini nama fungsi nya yaitu crime\_rate\_assign.

If else = merupakan kondisi perulangan ya dan tidak.

Cri6 = variable yang digunakan untuk menampung nilai array.

Cri6.head = merupakan variable yang mempunyai fungsi head yang digunakan untuk menampilkan *value* dari variable cri6.

Selanjutnya untuk melakukan proses klasifikasi, penulis membuat pelabelan tingkat kejahatan menjadi 3 level.

0 = tingkat kejahatan rendah.

1 = tingkat kejahatan sedang.

2 = tingkat kejahatan tinggi.

Untuk penentuan tingkat kejahatan, diukur dari seberapa banyak jumlah kejahatan yang terdapat pada waktu tersebut. Jumlah kejahatan dibagi menjadi 3 bagian.

Kejahatan rendah memiliki jumlah kejahatan sebanyak 0 – 14 kejadian.

Kejahatan sedang memiliki jumlah kejahatan sebanyak 15 – 33 kejadian.

Kejahatan tinggi memiliki jumlah kejahatan sebanyak 34 kejadian keatas.

Untuk kategori kejahatan, penulis menggunakan jenis kejahatan pencurian untuk penelitian ini, karena kejahatan tersebut adalah yang tertinggi dari kejahatan yang lain.

Untuk memprediksi tingkat kejahatan antara jumlah kejahatan dan jenis kejahatan, maka disini menggunakan percabangan if else. [2]

```
if(x<=14):  
    return 0  
elif(x>14 and x<=33):  
    return 1  
else:  
    return 2  
end if
```

Berikut script code untuk melakukan kondisi perulangan.

```
def crime_rate_assign(x):  
    if(x<=14):  
        return 0  
    elif(x>14 and x<=33):  
        return 1  
    else:  
        return 2
```

*Gambar 4. 4 Kondsi IF Else*

Untuk menentukan accuracy dari model yang digunakan, disini menggunakan beberapa rumus.

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{TN} + \text{FN}}$$

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{F1} = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$\text{UAR} = \frac{\text{R1} + \text{R2} + \text{R3}}{3}$$

Beberapa penjelasan diatas diantaranya :

TP = True Positif ( Mengacu pada jumlah sampel positif yang telah diidentifikasi dengan benar sebagai positif). [2]

FP = False Positif ( Mengacu pada jumlah sampel negatif yang salah diidentifikasi sebagai positif) [2]

TN = True Negatif ( Mengacu pada jumlah sampel negatif yang telah diidentifikasi dengan benar sebagai negatif). [2]

FN = False Negatif ( Mengacu pada jumlah sampel positif yang salah diidentifikasi sebagai negatif). [2]

Accuracy = Akurasi merupakan tingkat ketepatan atau seberapa dekat suatu pengukuran dengan nilai sebenarnya. [2]

Presisi = Mendeskripsikan kemungkinan suatu sampel berada di kelas X jika sudah diprediksi masuk ke dalam kelas X. [2]

Recall = Menjelaskan kemungkinan model untuk mengklasifikasikan sampel dengan benar sebagai milik kelas X jika itu benar benar milik kelas X. [2]

F1 = Ukuran yang menggambarkan hubungan antara presisi dan mengingat. Semakin tinggi skor F1 untuk suatu kelas, semakin baik kinerja model dalam mengklasifikasi kelas tersebut. [2]

## 4.2 Hasil dan Pembahasan

### 4.2.1 Mengimport Library

Pada tahap ini, persiapkan library yang akan digunakan.

```
import pandas as pd
import sys
sys.path.append("utils")
import matplotlib.pyplot as plt
import numpy as np
from datetime import datetime
import seaborn as sns

import util_script as us
```

*Gambar 4. 5 Import Library*

Pada code `util_script` tersebut, buat sebuah fungsi dengan menggunakan bahasa python untuk melakukan proses load data yang akan digunakan.

```
import pandas as pd

def create_df(filenamees):

    main_df = pd.read_csv(filenamees[0])
    print("Proses Loading Dataset "+filenamees[0][-8:4]+".")
    main_df = main_df[list(main_df.columns[:22])]
    for file in filenamees[1:]:
        print("Proses Loading Dataset "+file[-8:-4]+".")
        df_temp = pd.read_csv(file)
        df_temp = df_temp[list(df_temp.columns[:22])]
        main_df
    main_df.append(df_temp,ignore_index=True)
    print("Prses Loading Dataframe Selesai.\n\n")
```

```
return main_df
```

Year	Case Number	Time	Class	Block	LCR	Primary	Description	Location	Descend	Alt	Domestic	Seat	Dates	Flt	Community	Alt	FB	CR	X Coordinate	Y Coordinate	Altitude	Year	Updated On	Latitude	Longitude	Location	
2020	120147	01/12/20	12.00	000	5	BATTERY	DOMESTIC BATT APPOINTMENT	FALSE	TRUE	752	7	7	2020/01/01	4	08	08	08	174406	158621			2020/01/01	41.764345	-87.053487	41.764345	41.830607	
2020	120148	01/10/17	12.00	000	5	BATTERY	\$000 AND UNDERBATTERY	FALSE	TRUE	623	6	17	2020/01/01	4	17	44	17	17011	105178			2020/01/01	41.744718	-87.060388	41.744718	41.877407	
2020	120150	01/10/17	12.00	0.0	000	W/BATT	BATT BATTERY	DOMESTIC BATT APPOINTMENT	FALSE	TRUE	153	15	20	2020/01/01	4	28	28	10872	108325			2020/01/01	41.782222	-87.768156	41.782222	41.760716	
2020	120169	01/10/18	12.00	000	000	12.00	POSSIBLY - CASH	TRUE	TRUE	710	7	16	2020/01/01	4	16	16	16	16926	168621			2020/01/01	41.770938	-87.064768	41.769994	41.769994	
2020	120172	01/10/18	12.00	000	000	000	12.00	HOMECIDE	FIRST DEGREE STREET	FALSE	TRUE	302	3	3	3	3	3	4318	168621			2020/01/01	41.758647	-87.058617	41.758647	41.862781	
2020	120173	01/10/18	12.00	000	000	000	12.00	HOMECIDE	CASH APPOINTMENT	FALSE	TRUE	302	3	3	3	3	3	19184	168621			2020/01/01	41.758647	-87.058617	41.758647	41.862781	
2020	120184	01/10/18	12.00	000	000	000	12.00	BATT	\$000 AND UNDERBATTERY	TRUE	TRUE	139	13	20	2020/01/01	4	6	16	16167	105757			2020/01/01	41.760123	-87.074884	41.760123	41.898248
2020	120193	01/10/18	12.00	000	000	000	12.00	W/BATT	\$000 AND UNDERBATTERY	FALSE	TRUE	321	3	5	43	15	15	10872	105960			2020/01/01	41.738078	-87.061307	41.738078	41.877407	
2020	120193	01/10/18	12.00	000	000	000	12.00	W/BATT	\$000 AND UNDERBATTERY	FALSE	TRUE	321	3	5	43	15	15	10872	105960			2020/01/01	41.738078	-87.061307	41.738078	41.877407	
2020	120200	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	243	24	48	77	14	16500	104543			2020/01/01	41.754452	-87.066120	41.754452	41.898248		
2020	120220	01/10/18	12.00	000	000	000	12.00	CRIMINAL	TRIAL TO VEHICLE STREET	FALSE	TRUE	243	24	5	41	26						2020/01/01					
2020	120132	01/10/02	12.00	000	000	000	12.00	HOMECIDE	FIRST DEGREE STREET	FALSE	TRUE	434	4	10	40	18	19360	1047157			2020/01/01	41.746833	-87.754445	41.746833	41.844421		
2020	120169	01/10/18	12.00	000	000	000	12.00	HOMECIDE	CASH TO PROPERTY SMALL RETAIL	FALSE	TRUE	102	10	22	30	14	15276	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120118	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120118	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.00	000	000	000	12.00	CRIMINAL	CASH TO PROPERTY RESCUE	FALSE	TRUE	142	14	36	21	14	13437	105960			2020/01/01	41.758647	-87.058617	41.758647	41.862781		
2020	120140	01/10/18	12.0																								



## Analisis Tingkat Kejahatan Menggunakan Model Random Forest

---

```
DATA_PATH = "E:/Poltekpos/TI/Tingkat 4/Intership/# Source Code Khusus Jurnal/Data Crime Internship/data/"

file_names = ['crimes_2019.csv', 'crimes_2020.csv', 'crimes_2021.csv']
file_names = [DATA_PATH+x for x in file_names]

main_df = us.create_df(file_names)
orig_shape = main_df.shape
print("Jumlah Kejahatan: "+ str(main_df.shape[0]))
print("\n Jumlah Kolom: "+ str(main_df.shape[1]))
```

*Gambar 4. 7 Membaca Dataset*

DATA\_PATH = digunakan untuk mengetahui posisi tempat penyimpanan data csv yang digunakan.

file\_names = variable yang digunakan untuk menampung data array.

main\_df = variable yang digunakan untuk mengeksekusi parameter file\_names.

orig\_shape = variable yang digunakan untuk menampung fungsi shape, disini main\_df.shape.

print = fungsi untuk menampilkan hasil dari suatu parameter.

Berikut tampilan dari perintah tersebut seperti ini.

```
Proses Loading Dataset 2019.
Proses Loading Dataset 2020.
Proses Loading Dataset 2021.
Prses Loading Dataframe Selesai.
```

```
Jumlah Kejahatan: 492007
```

```
Jumlah Kolom: 22
```

*Gambar 4. 8 Hasil Pembacaan Dataset*

Pada tampilan gambar tersebut menjelaskan bahwa :

1. Tiga **Proses Loading** diatas menjelaskan bahwa proses dataset sedang di load untuk ditampilkan.
2. **Proses Loading Dataframe Selesai** menjelaskan bahwa proses untuk meload dataset telah selesai.

```
Proses Loading Dataset 2019.  
Proses Loading Dataset 2020.  
Proses Loading Dataset 2021.  
Prses Loading Dataframe Selesai.
```

```
Jumlah Kejahatan: 492007
```

```
Jumlah Kolom: 22
```

*Gambar 4. 9 Hasil Pembacaan Dataset*

3. **Jumlah Kejahatan** menjelaskan tentang total kasus kejahatan yang terjadi dari tahun 2019 – 2021.
4. **Jumlah Kolom** menjelaskan total kolom yang tersedia pada dataset tersebut.
5. Untuk melihat info tipe data yang digunakan pada dataset tersebut menggunakan fungsi `main_df.info()`.

```
main_df.info()
```

*Gambar 4. 10 Melihat Informasi Dataframe*

Berikut tampilan dari fungsi tersebut.

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 492007 entries, 0 to 492006
Data columns (total 22 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   ID                                     492007 non-null  int64
1   Case Number                           492007 non-null  object
2   Date                                  492007 non-null  object
3   Block                                 492007 non-null  object
4   IUCR                                  492007 non-null  object
5   Primary Type                           492007 non-null  object
6   Description                             492007 non-null  object
7   Location Description                    489790 non-null  object
8   Arrest                                492007 non-null  bool
9   Domestic                              492007 non-null  bool
10  Beat                                  492007 non-null  int64
11  District                              492007 non-null  int64
12  Ward                                  491985 non-null  float64
13  Community Area                         492006 non-null  float64
14  FBI Code                              492007 non-null  object
15  X Coordinate                           486901 non-null  float64
16  Y Coordinate                           486901 non-null  float64
17  Year                                   492007 non-null  int64
18  Updated On                             492007 non-null  object
19  Latitude                               486901 non-null  float64
20  Longitude                              486901 non-null  float64
21  Location                               486901 non-null  object
dtypes: bool(2), float64(6), int64(4), object(10)
memory usage: 76.0+ MB
```

*Gambar 4. 11 Hasil Informasi Dataframe*

Terdapat beberapa tipe data yang digunakan pada dataset tersebut diantaranya integer, object, boolean, dan float.

6. Untuk menghapus data NaN atau kosong, null data tidak ada, menggunakan fungsi `dropna()`, maka akan menghapus data yang kosong.

```
main_df = main_df.dropna()
main_df.isna().sum()
```

*Gambar 4. 12 Menghapus Data Kosong*

Maka tampilannya seperti ini :

```
ID                                0
Case Number                      0
Date                            0
Block                           0
IUCR                            0
Primary Type                    0
Description                     0
Location Description             0
Arrest                          0
Domestic                        0
Beat                            0
District                        0
Ward                            0
Community Area                  0
FBI Code                        0
X Coordinate                    0
Y Coordinate                    0
Year                            0
Updated On                      0
Latitude                        0
Longitude                      0
Location                        0
dtype: int64
```

*Gambar 4. 13 Tampilan Data Kosong*

Hasil yang diperoleh 0 karena kekosongan data tidak terlalu banyak.

7. Untuk melihat data yang telah dibersihkan

```
print("Data Setelah Di Cleaning / Dibersihkan:",round((((main_df.shape[0]/orig_shape[0]) * 100),2),"%")
```

*Gambar 4. 14 Data Setelah Dibersihkan*

Untuk melihatnya hasilnya menggunakan perintah print, kemudian gunakan fungsi round, fungsi round ini digunakan untuk mengembalikan nilai dalam bentuk decimal. Kemudian lakukan pembagian pada variable

`((main_df.shape[0]/orig_shape[0]) * 100),2)“%”,` berfungsi untuk membagi nilai array, seperti nilai dari data dari dataset, dimana data yang tersedia, dibagi dengan data yang dibutuhkan pada saat proses pembersihan data tersebut. Angka 2 disini berfungsi untuk menampilkan berapa banyak angka dibelakang koma, disini ditampilkan yaitu 2 angka, symbol persen disini sebagai string untuk menambahkan symbol di akhir pada saat program telah selesai di jalankan. maka akan mendapatkan hasil 98, 57%.

**Data Setelah Di Cleaning / Dibersihkan: 98.57 %**

8. Untuk melihat data pada dataset menggunakan fungsi `head()`, secara default fungsi ini akan menampilkan 5 data pada dataset, bilangannya dimulai dari 0 – 4.

```
main_df.head()
```

*Gambar 4. 15 Melihat Informasi Data*

Berikut tampilannya seperti ini :

## Analisis Tingkat Kejahatan Menggunakan Model Random Forest

	ID	Case Number	Date	Block	IUCR	Primary Type	Description	Location Description	Arrest	Domestic	...	Ward	Community Area	FBI Code
0	11675533	JC248984	05/03/2019 11:40:00 PM	027XX W NORTH AVE	1812	NARCOTICS	POSS CANNABIS MORE THAN 30GMS	STREET	True	False	...	1.0	24.0	18
1	11675453	JC248966	05/03/2019 11:33:00 PM	051XX W MADISON ST	1330	CRIMINAL TRESPASS	TO LAND	PARKING LOT/GARAGE(NON.RESID.)	True	False	...	28.0	25.0	26
2	11675494	JC248999	05/03/2019 11:30:00 PM	011XX W LAWRENCE AVE	0810	THEFT	OVER \$500	MOVIE HOUSE/THEATER	False	False	...	46.0	3.0	06
3	11675481	JC248956	05/03/2019 11:30:00 PM	121XX S HARVARD AVE	0810	THEFT	OVER \$500	RESIDENTIAL YARD (FRONT/BACK)	False	False	...	34.0	53.0	06
4	11675910	JC248965	05/03/2019 11:25:00 PM	025XX N MILWAUKEE AVE	5011	OTHER OFFENSE	LICENSE VIOLATION	BAR OR TAVERN	False	False	...	32.0	22.0	26

5 rows × 22 columns

Gambar 4. 16 Tampilan Informasi Data

X Coordinate	Y Coordinate	Year	Updated On	Latitude	Longitude	Location
1157682.0	1910520.0	2019	05/10/2019 04:20:42 PM	41.910225	-87.696173	(41.910224908, -87.696172663)
1142285.0	1899542.0	2019	05/10/2019 04:20:42 PM	41.880400	-87.753009	(41.880399914, -87.753008553)
1167886.0	1932034.0	2019	05/10/2019 04:20:42 PM	41.969046	-87.658065	(41.969046206, -87.658064717)
1176145.0	1824335.0	2019	05/10/2019 04:20:42 PM	41.673328	-87.630934	(41.67332819, -87.630934077)
1155343.0	1916744.0	2019	05/10/2019 04:20:42 PM	41.927351	-87.704598	(41.92735143, -87.704597631)

Gambar 4. 17 Tampilan Informasi Data

#### 4.4.2 Scrub Data

Tahapan selanjutnya yaitu melakukan scrub data (pembersihan data) pada dataset, tahapan tersebut diantaranya :

1. Untuk melihat nama kolom pada tabel tersebut menggunakan fungsi **columns**.

```
print(main_df.columns)
```

*Gambar 4. 18 Melihat Informasi Kolom*

Berikut tampilan informasi kolom tersebut.

```
Index(['ID', 'Case Number', 'Date', 'Block', 'IUCR', 'Primary Type',  
      'Description', 'Location Description', 'Arrest', 'Domestic', 'Beat',  
      'District', 'Ward', 'Community Area', 'FBI Code', 'X Coordinate',  
      'Y Coordinate', 'Year', 'Updated On', 'Latitude', 'Longitude',  
      'Location'],  
      dtype='object')
```

*Gambar 4. 19 Tampilan Informasi Kolom*

2. *Getting rid of missing value* (menghilangkan nilai yang kosong)  
Pada tahap ini akan memeriksa data dan membersihkan data dari data yang tidak diperlukan dengan menggunakan Bahasa pemrograman Python. Data yang dibersihkan terletak pada kolom Date.  
Berikut fungsi yang digunakan untuk membersihkan data pada kolom *Date* tersebut.

```
def time_convert(date_time):
    s1 = date_time[:11]
    s2 = date_time[11:]

    month = s1[:2]
    date = s1[3:5]
    year = s1[6:10]

    hr = s2[:2]
    mins = s2[3:5]
    sec = s2[6:8]
    time_frame = s2[9:]
    if(time_frame == 'PM'):
        if (int(hr) != 12):
            hr = str(int(hr) + 12)
    else:
        if(int(hr) == 12):
            hr = '00'

    final_date = datetime(int(year), int(month), int(date), int(hr), int(mins), int(sec))
    return final_date
```

*Gambar 4. 20 Konversi Bulan, Hari, Jam*

Kemudian pada kolom **Date** tersebut dilakukan convert dengan menggunakan fungsi **apply(time\_convert)**

```
main_df['Date'] = main_df['Date'].apply(time_convert)
```

*Gambar 4. 21 Fungsi Convert*

Untuk menampilkan tanggal hasil *convert* tersebut bisa menggunakan perintah **head**, maka akan ditampilkan 5 data teratas.

```
main_df['Date'].head()
```

*Gambar 4. 22 Tampilan Informasi Data*



Berikut hasil dari *convert* tersebut.

```
0    2019-05-03 23:40:00
1    2019-05-03 23:33:00
2    2019-05-03 23:30:00
3    2019-05-03 23:30:00
4    2019-05-03 23:25:00
Name: Date, dtype: datetime64[ns]
```

Gambar 4. 23 Hasil Pembersihan Tabel Date

3. *Feature Engineering* (Rekayasa Fitur)

Pada tahap ini akan merekayasa beberapa fitur yang sesuai dengan bahasa yang digunakan, disini akan melakukan *convert* field **Date (Tanggal)** pada dataset yang dimana yang direkayasa agar dipisahka tiap datanya yaitu **Month (Bulan)**, **Day (Hari)** dan **Hour (Jam)** menggunakan library *pandas* dan *datetime* dari *python*.

Proses *feature engineering* pada *pyhon* dilakukan seperti ini:

```
# Feature Engineering 1 : Month
def month_col(x):
    return int(x.strftime("%m"))
main_df['Month'] = main_df['Date'].apply(month_col)

# Feature Engineering 2 : Day
def day_col(x):
    return int(x.strftime("%w"))
main_df['Day'] = main_df['Date'].apply(day_col)

# Feature Engineering 3 : Hour
def hour_col(x):
    return int(x.strftime("%H"))
main_df['Hour'] = main_df['Date'].apply(hour_col)
```

Gambar 4. 24 Rekayasa Fitur Date

## Analisis Tingkat Kejahatan Menggunakan Model Random Forest

Pada tahap ini terdapat penambahan kolom sebanyak 3 kolom untuk menampung nilai dari Bulan (*Month*), Hari (*Day*) dan Jam (*Hour*).

Untuk melihat hasilnya gunakan fungsi **head**.

```
main_df.head()
```

Gambar 4. 25 Menampilkan Informasi Dataframe

Pada hasilnya terdapat penambahan kolom, dimana yang pertama terdapat 22 kolom, dan untuk sekarang bertambah menjadi 25 kolom seperti ini.

	ID	Case Number	Date	Block	IUCR	Primary Type	Description	Location Description	Arrest	Domestic	...	X Coordinate	Y Coordinate
0	11675533	JC248984	2019-05-03 23:40:00	027XX W NORTH AVE	1812	NARCOTICS	POSS. CANNABIS MORE THAN 30GMS	STREET	True	False	...	1157682.0	1910520.0
1	11675453	JC248966	2019-05-03 23:33:00	051XX W MADISON ST	1330	CRIMINAL TRESPASS	TO LAND	PARKING LOT/GARAGE(NON RESID.)	True	False	...	1142285.0	1899542.0
2	11675494	JC248999	2019-05-03 23:30:00	011XX W LAWRENCE AVE	0810	THEFT	OVER \$500	MOVIE HOUSE/THEATER	False	False	...	1167886.0	1932034.0
3	11675481	JC248956	2019-05-03 23:30:00	121XX S HARVARD AVE	0810	THEFT	OVER \$500	RESIDENTIAL YARD (FRONT/BACK)	False	False	...	1176145.0	1824335.0
4	11675910	JC248965	2019-05-03 23:25:00	025XX N MILWAUKEE AVE	5011	OTHER OFFENSE	LICENSE VIOLATION	BAR OR TAVERN	False	False	...	1155343.0	1916744.0

5 rows × 25 columns

Gambar 4. 26 Tampilan Informasi Dataframe

Disini terlihat ada 3 kolom yang bertambah yaitu kolom *Month*, *Day*, *Hour*.

Year	Updated On	Latitude	Longitude	Location	Month	Day	Hour
2019	05/10/2019 04:20:42 PM	41.910225	-87.696173	(41.910224908, -87.696172663)	5	5	23
2019	05/10/2019 04:20:42 PM	41.880400	-87.753009	(41.880399914, -87.753008553)	5	5	23
2019	05/10/2019 04:20:42 PM	41.969046	-87.658065	(41.969046206, -87.658064717)	5	5	23
2019	05/10/2019 04:20:42 PM	41.673328	-87.630934	(41.67332819, -87.630934077)	5	5	23
2019	05/10/2019 04:20:42 PM	41.927351	-87.704598	(41.92735143, -87.704597631)	5	5	23

Gambar 4. 27 Tampilan Informasi Dataframe

Selanjutnya memilih salah satu kejahatan yang digunakan untuk penelitian ini, peneliti memilih kejahatan pencurian sebagai tingkat kejahatan yang akan dianalisis, kemudian memfilter jenis kejahatan hanya tindak pidana pencurian.

```
top_10 = list(main_df['Primary Type'].value_counts().head(10).index)

def filter_top_10(df):
    df2=df[df['Primary Type']=='THEFT']
    for crime in top_10[1:]:
        temp=df[df['Primary Type']==crime]
        df2 = df2.append(temp, ignore_index=True)
    return df2

df2=filter_top_10(main_df) # the dataframe with all the data of only the top 10 crimes
df2.shape
```

Gambar 4. 28 Klasifikasi Kejahatan Pencurian

Tampilan hasil kejahatan dengan kejahatan pencurian (**THEFT**) sebagai berikut.

(330221, 22)

Gambar 4. 29 Hasil Kejahatan Pencurian

Untuk menampilkan dataset baru untuk pencurian, gunakan fungsi head.

df2.head()

Gambar 4. 30 Melihat Informasi Dataframe

Hasil dari dataset untuk pencurian (**THEFT**) seperti ini.

	ID	Case Number	Date	Block	IUCR	Primary Type	Description	Location Description	Arrest	Domestic	...	Ward	Community Area	FBI Code	X Coordinate
0	11675494	JC248999	2019-05-03 23:30:00	011XX W LAWRENCE AVE	0810	THEFT	OVER \$500	MOVIE HOUSE/THEATER	False	False	...	46.0	3.0	06	1167886.0
1	11675481	JC248956	2019-05-03 23:30:00	121XX S HARVARD AVE	0810	THEFT	OVER \$500	RESIDENTIAL YARD (FRONT/BACK)	False	False	...	34.0	53.0	06	1176145.0
2	11675513	JC248946	2019-05-03 23:11:00	066XX W GRAND AVE	0810	THEFT	OVER \$500	RESTAURANT	True	False	...	29.0	18.0	06	1131423.0
3	11675615	JC249145	2019-05-03 23:00:00	079XX S ST LAWRENCE AVE	0820	THEFT	\$500 AND UNDER	RESIDENCE-GARAGE	False	False	...	6.0	44.0	06	1181614.0
4	11676084	JC249692	2019-05-03 23:00:00	032XX N LAMON AVE	0820	THEFT	\$500 AND UNDER	STREET	False	False	...	31.0	15.0	06	1143136.0

Gambar 4. 31 Tampilan Informasi Dataframe

Kemudian kelompokkan beberapa kolom menjadi sebuah group. Ini berfungsi sebagai patokan dalam penelitian, dalam menentukan bulan, Hari dan Jam.

```
cri5 = df2.groupby(['Month', 'Day', 'District', 'Hour'], as_index=False).agg({"Primary Type": "count"})
cri5 = cri5.sort_values(by=['District'], ascending=False)
cri5.head()
```

Gambar 4. 32 Pengelompokan Data kolom

Berikut hasil tampilan dari pengelompokan yang telah dibuat.

	Month	Day	District	Hour	Primary Type
25709	7	6	31	10	1
36724	10	6	31	3	1
35145	10	3	31	23	1
28336	8	4	31	0	1
36199	10	5	31	2	1

*Gambar 4. 33 Tampilan Pengelompokan Data*

```
cri6=cri5.rename(index=str, columns={"Primary Type":"Jumlah_Kejahatan"})
cri6.head()
```

*Gambar 4. 34 Ubah Nama Kolom*

Berikut tampilan hasil dari perubahan nama kolom tersebut.

	Month	Day	District	Hour	Jumlah_Kejahatan
25709	7	6	31	10	1
36724	10	6	31	3	1
35145	10	3	31	23	1
28336	8	4	31	0	1
36199	10	5	31	2	1

*Gambar 4. 35 Tampilan Perubahan Nama Kolom*

4. *Compressing the Dataset* (Melakukan kompres pada dataset)  
Tahapan selanjutnya yaitu melakukan kompres dataset yang didapat menjadi dataset yang baru, dimana data yang sebelumnya tidak dibutuhkan, dihilangkan dan data ini yang akan digunakan

## Analisis Tingkat Kejahatan Menggunakan Model Random Forest

---

untuk melakukan proses prediksi tersebut. data yang di kompres merupakan data dalam bentuk csv.

Hasil dari compress data tersebut.

	Month	Day	Hour	District	Jumlah_Kejahatan	Alarm
25709	7	6	10	31	1	0
36724	10	6	3	31	1	0
35145	10	3	23	31	1	0
28336	8	4	0	31	1	0
36199	10	5	2	31	1	0
36200	10	5	7	31	1	0
22568	7	0	12	31	1	0
13676	4	4	15	31	1	0
528	1	0	0	31	1	0
7360	2	6	12	31	1	0
3156	1	5	11	31	1	0
42461	12	3	11	31	1	0
40377	11	6	12	31	1	0
23090	7	1	9	31	1	0
21517	6	5	7	31	1	0
11568	4	0	16	25	10	0
27269	8	2	3	25	2	0
11569	4	0	17	25	18	1
11567	4	0	15	25	11	0

Data hasil kompres yang ditampilkan hanya sebagian data. Data yang berhasil dikompres sebanyak **44015 data** dan mempunyai **7 kolom**.

### 4.4.3 *Explore Data*

Pada tahap ini dilakukan *explore* data dimana pada bagian ini jumlah kejahatan dibagi berdasarkan bulan, hari, tempat dan waktu.

Untuk pertama mengetahui jumlah dataset yang telah hasil di kompres

```
cri6 = cri6[['Month', 'Day', 'District', 'Hour', 'Jumlah_Kejahatan']]
cri6.head()
print("Jumlah Dataset:", cri6.shape)
```

*Gambar 4. 36 Explore Data*

Untuk tampilan jumlah hasil nya seperti berikut

```
Jumlah Dataset: (44015, 5)
```

*Gambar 4. 37 Jumlah Data*

Untuk melihat rata rata jumlah kejahatan yang dilakukan.

```
print("Rata-rata kejahatan per distrik per titik waktu  
: ", round(cri6['Jumlah_Kejahatan'].sum()/cri6.shape[0],  
2), ".")
```

Untuk melihat jumlah hasil rata rata seperti berikut.

```
Rata-rata kejahatan per distrik per titik waktu : 7.59 .
```

*Gambar 4. 38 Melihat Rata Rata Kejahatan Per Waktu*

Selanjutnya pada tahap ini dilakukan percabangan if else, pada tahap analisis dijelaskan bahwa pada proses terdapat 3 label, yaitu 0,1 dan 2 yang artinya tingkat kejahatan rendah, sedang dan tinggi.

Untuk percabangannya seperti ini.

```
def crime_rate_assign(x):  
    if(x<=14):  
        return 0  
    elif(x>14 and x<=33):  
        return 1  
    else:  
        return 2  
cri6['Alarm'] = cri6['Jumlah_Kejahatan'].apply(crime_rate_assign)  
cri6 = cri6[['Month', 'Day', 'Hour', 'District', 'Jumlah_Kejahatan', 'Alarm']]  
cri6.head()
```

*Gambar 4. 39 Kondisi IF Else*

Buat fungsi terlebih dahulu, disini terdapat fungsi `crime_rate_assign` dengan nilai `x`, selanjutnya didalam fungsi berikan beberapa kondisi, disini menggunakan kondisi `if else`.

Untuk hasil nya seperti ini.

	Month	Day	Hour	District	Jumlah_Kejahatan	Alarm
21329	6	5	7	31	1	0
28101	8	4	0	31	1	0
22895	7	1	9	31	1	0
7306	2	6	12	31	1	0
35861	10	5	2	31	1	0

*Gambar 4. 40 Tampilan Hasil Explore Data*

Terdapat penambahan kolom yaitu `district` untuk menjelaskan kode daerah tersebut.

Untuk melihat hasil kompres dataset tersebut dapat menggunakan perintah.

```
cri6.to_csv("E:/Poltekpos/TI/Tingkat 4/Intership/# Source Code Khusus  
Jurnal/Data Crime Internship/data/Hasil_Kompres.csv")
```



Terdapat variable yang digunakan untuk menampung data penyimpanan ke lokasi yang dituju. Untuk penyimpanan file tersebut bisa berbeda beda tergantung selera anda untuk menyimpan file tersebut.

Kemudian untuk mengecek seberapa banyak data yang digunakan untuk klasifikasi gunakan fungsi **value\_counts()**, berfungsi untuk melihat nilai hitungannya.

```
cri6['Alarm'].value_counts()
```

*Gambar 4. 41 Value Counts*

Berikut hasil dari fungsi tersebut.

```
0    39933
1     3589
2         14
Name: Alarm, dtype: int64
```

*Gambar 4. 42 Hasil Value Counts*

Kemudian untuk melihat presentase tingkat kejahatan, baik rendah sedang dan tinggi.

```
print("Persentase Kejahatan Rendah:", round(cri6['Alarm'].value_counts()[0]/cri6['Alarm'].value_counts().sum()*100,2))
print("Persentase Kejahatan Sedang:", round(cri6['Alarm'].value_counts()[1]/cri6['Alarm'].value_counts().sum()*100,2))
print("Persentase Kejahatan Tinggi:", round(cri6['Alarm'].value_counts()[2]/cri6['Alarm'].value_counts().sum()*100,2))
```

*Gambar 4. 43 Presentase Tingkat Kejahatan*

Untuk hasilnya seperti berikut.

```
Persentase Kejahatan Rendah: 91.72
Persentase Kejahatan Sedang: 8.24
Persentase Kejahatan Tinggi: 0.03
```

*Gambar 4. 44 Tampilan Presentase*

Untuk melihat kejahatan berdasarkan tempat kejadian, dilakukan perintah seperti berikut :

Selanjutnya melakukan explore data berdasarkan nama lokasi seperti berikut:

```
def crime_rate_assign(x):  
    if(x<=14):  
        return 0  
    elif(x>14 and x<=33):  
        return 1  
    else:  
        return 2  
cri6['Alarm'] = cri6['Jumlah_Kejahatan'].apply(crime_rate_assign)  
cri6 = cri6[['Month','Day','Hour','Location Description','Jumlah_Kejahatan','Alarm']]  
cri6.head(60)
```

*Gambar 4. 45 Kondisi IF Else*

Untuk melihat kejadian berdasarkan tahun

```
%matplotlib inline  
  
main_df.resample('M').size().plot(legend=True)  
plt.title('Jumlah Kejadian')  
plt.xlabel('Data Kejadian Bulan ')  
plt.ylabel('Bulan dan Tahun Kejadian')  
plt.show()
```

*Gambar 4. 46 Melihat Data Kejadin*

Berikut penjelasannya ,

%matplotlib inline = berfungsi agar hasil visualisasi bisa langsung tercetak di Jupyter Notebook.

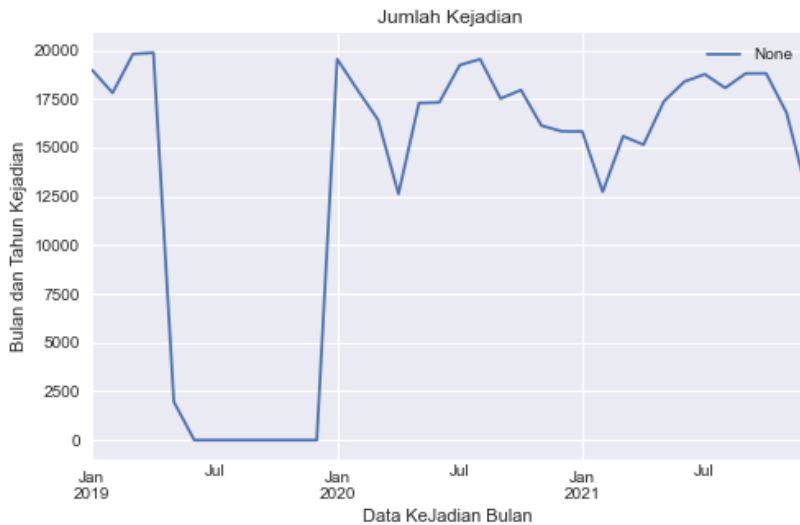
plt.title = digunakan untuk menentukan judul visualisasi yang digambarkan dan menampilkan judul menggunakan berbagai atribut.title()

plt.xlabel = digunakan untuk memberikan nama label untuk garis x.

plt.ylabe = digunakan untuk memberikan nama label untuk garis y.

plt.show = digunakan untuk menampilkan plot.

Beriku tampilan kejadian berdasarkan tahun.



*Gambar 4. 47 Data Kejadian Berdasarkan Tahun*

Pada gambar diatas, yang paling banyak kejadian kejahatan adalah tahun 2019 dari januari sampai juli dengan kejadian 20000 kejadian.

Selanjutnya untuk melihat kejadian berdasarkan bulan

```
%matplotlib inline
main_df['Date'].apply(lambda x: x.month).value_counts(sort=False).plot(kind='barh')
plt.xlabel("Jumlah Kejadian Per Bulan")
plt.ylabel("Bulan Kejadian")
plt.title("Data Kejadian Per Bulan")

plt.savefig("E:/Poltekpos/TI/Tingkat 4/jumlahdataperbulan.png")
plt.show()
```

*Gambar 4. 48 Berdasarkan Bulan*

Berikut penjelasannya ,

`%matplotlib inline` = berfungsi agar hasil visualisasi bisa langsung tercetak di Jupyter Notebook.

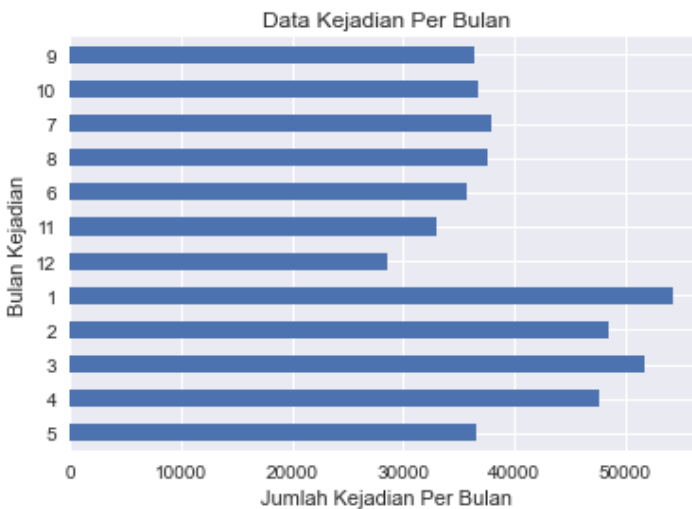
`plt.title` = digunakan untuk menentukan judul visualisasi yang digambarkan dan menampilkan judul menggunakan berbagai atribut `title()`

`plt.xlabel` = digunakan untuk memberikan nama label untuk garis x.

`plt.ylabel` = digunakan untuk memberikan nama label untuk garis y.

`plt.show` = digunakan untuk menampilkan plot.

Berikut tampilan kejadian berdasarkan bulan



*Gambar 4. 49 Hasil Kejadian Berdasarkan Bulan*

Selanjutnya untuk melihat kejadian berdasarkan hari

Pada gambar diatas, yang paling banyak kejadian kejahatan yaitu pada bulan 1 dengan kejadian sebanyak 50000 lebih kejadian.

## Analisis Tingkat Kejahatan Menggunakan Model Random Forest

---

```
%matplotlib inline  
  
main_df['Date'].apply(lambda x: x.weekday()).value_counts(sort=False).plot(kind='barh')  
plt.xlabel("Jumlah Kejadian Per Per Hari")  
plt.ylabel("Hari Kejadian")  
plt.title("Data Kejadian Per Hari")  
plt.show()
```

*Gambar 4. 50 Berdasarkan Hari*

%matplotlib inline = berfungsi agar hasil visualisasi bisa langsung tercetak di Jupyter Notebook.

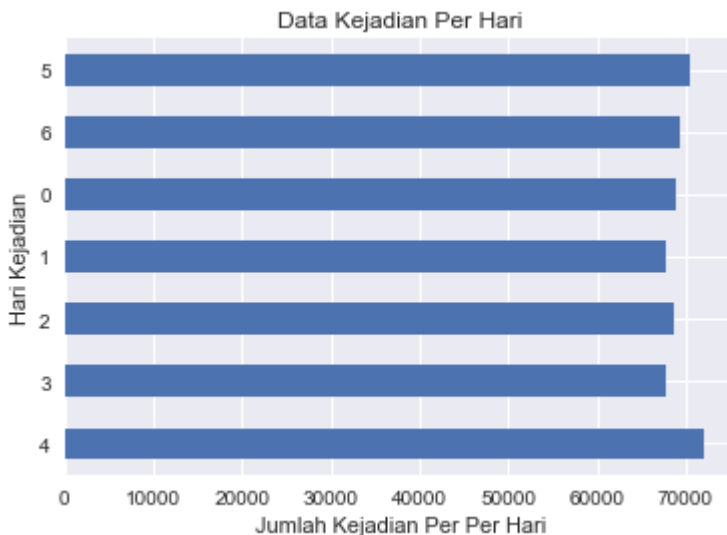
plt.title = digunakan untuk menentukan judul visualisasi yang digambarkan dan menampilkan judul menggunakan berbagai atribut.title()

plt.xlabel = digunakan untuk memberikan nama label untuk garis x.

plt.ylabe = digunakan untuk memberikan nama label untuk garis y.

plt.show = digunakan untuk menampilkan plot.

Berikut tampilan kejadian berdasarkan hari



*Gambar 4. 51 Hasil Kejadian Berdasarkan Hari*

Pada gambar diatas, yang paling banyak kejadian terdapat pada hari ke 4 dengan kejadian sebanyak 70000 lebih kejadian.

Selanjutnya untuk melihat kejadian berdasarkan Jam.

```
%matplotlib inline
main_df['Date'].apply(lambda x: x.hour).value_counts(sort=False).plot(kind='bar')
plt.xlabel("Jumlah Kejadian Per Jam")
plt.ylabel("Jam Kejadian")
plt.title("Data Kejadian Per Jam")
plt.show()
```

*Gambar 4. 52 Berdasarkan Jam*

%matplotlib inline = berfungsi agar hasil visualisasi bisa langsung tercetak di Jupyter Notebook.

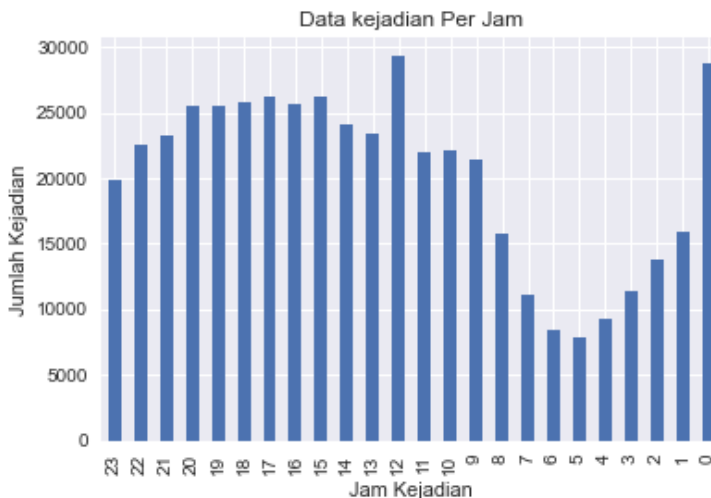
plt.title = digunakan untuk menentukan judul visualisasi yang digambarkan dan menampilkan judul menggunakan berbagai atribut.title()

plt.xlabel = digunakan untuk memberikan nama label untuk garis x.

plt.ylabe = digunakan untuk memberikan nama label untuk garis y.

plt.show = digunakan untuk menampilkan plot.

Berikut tampilan kejadian berdasarkan Jam



*Gambar 4. 53 Hasil Kejadian Berdasarkan Jam*

Pada gambar diatas, yang paling banyak kejadian terdapat pada jam ke 12 dengan kejadian sebanyak 30000 kejadian.

Selanjutnya untuk melihat kejadian berdasarkan lokasi.

```
%matplotlib inline
main_df['Location Description'].value_counts()
```

*Gambar 4. 54 Berdasarkan Lokasi*

%matplotlib inline = berfungsi agar hasil visualisasi bisa langsung tercetak di Jupyter Notebook.

Main\_df = digunakan untuk menampung nilai array.

Value\_counts() = merupakan fungsi untuk menghitung nilai dari array.

Berikut tampilan kejadian berdasarkan lokasi.

```
STREET          118128
APARTMENT       89202
RESIDENCE       81858
SIDEWALK        30433
SMALL RETAIL STORE 12604
...
CTA PROPERTY    1
CHA ELEVATOR    1
DRIVEWAY        1
ELEVATOR        1
STAIRWELL       1
Name: Location Description, Length: 172, dtype: int64
```

Pada gambar diatas, yang paling banyak kejadian yaitu di tempat *street* (Jalan) dengan kejadian sebanyak 118128.

#### 4.4.4 Model

Pada tahap model ini, akan dilakukan proses pemodelan dengan menggunakan metode random forest untuk menentukan tingkat akurasi dari penelitian ini, untuk menentukan accuracy berikut.

Untuk pengujian dataset dengan model, sebelumnya akan meload terlebih dahulu data data yang akan dibutuhkan. Disini akan melakukan pengujian dari tahapan yang sebelumnya telah dilakukan.

```
# Load the Dataset
test_files = ['crimes_2019.csv', 'crimes_2020.csv', 'crimes_2021.csv']
test_files = [DATA_PATH+x for x in test_files]
test_df = us.create_df(test_files)
```

*Gambar 4. 55 Meload Dataset*

Pada tahap ini berfungsi untuk meload dataset yang digunakan. Pada tahap ini bernama test\_file dimana nama variable ini berfungsi untuk mengetes dataset yang digunakan.

Tahap selanjutnya untuk menghilangkan nilai kosong.

```
# Drop missing values
test_df = test_df.dropna()
```

*Gambar 4. 56 Menghilangkan Nilai Kosong*

Pada tahap ini, nilai yang kosong akan dihapus menggunakan fungsi dropna().

Tahap selanjutnya yaitu *feature engineering*.

```
# Feature Engineering our columns
test_df['Month'] = test_df['Date'].apply(month_col)
test_df['Day'] = test_df['Date'].apply(day_col)
test_df['Hour'] = test_df['Date'].apply(hour_col)
```

*Gambar 4. 57 Feature Engineering*

Pada tahap ini dilakukan rekayasa fitur, dimana pada kolom *date* dipecah menjadi 3 bagian, yaitu bulan, hari dan jam.

Tahap selanjutnya yaitu compress dataset



## Analisis Tingkat Kejahatan Menggunakan Model Random Forest

```
# Compressing
df7 = filter_top_10(test_df)
cri7 = df7.groupby(["Month", "Day", "District", "Hour"], as_index=False).agg({"Primary Type" : "count"})
cri7 = cri7.sort_values(by=["District"], ascending=False)
cri8 = cri7.rename(index=str, columns={"Primary Type" : "Jumlah_Kejahatan"})
cri8 = cri8[["Month", "Day", "District", "Hour", "Jumlah_Kejahatan"]]
cri8['Alarm'] = cri8['Jumlah_Kejahatan'].apply(crime_rate_assign)
cri8 = cri8[['Month', 'Day', 'Hour', 'District', 'Jumlah_Kejahatan', 'Alarm']]
print(cri8.head(90))
```

Pada tahap ini dilakukan kompres data, artinya hanya menggunakan kolom-kolom yang akan digunakan.

Berikut tampilan hasil dari test data diatas tersebut.

```
Proses Loading Dataset 2019.
Proses Loading Dataset 2020.
Proses Loading Dataset 2021.
Prses Loading Dataframe Selesai.

   Month  Day  Hour  District  Crime_Count  Alarm
21329    6    5    7        31           1     0
28101    8    4    0        31           1     0
22895    7    1    9        31           1     0
7306     2    6   12        31           1     0
35861   10    5    2        31           1     0
Class Imbalance

0    39933
1    3589
2      14
Name: Alarm, dtype: int64
```

Gambar 4. 58 Hasil Testing

Membuat *balance* dataset, membuat data seimbang.

## Analisis Tingkat Kejahatan Menggunakan Model Random Forest

---

```
from sklearn.utils import resample # for upsampling

# Set individual classes
cri6_low = cri6[cri6['Alarm']==0]
cri6_medium = cri6[cri6['Alarm']==1]
cri6_high = cri6[cri6['Alarm']==2]

# Upsample the minority classes to size of class 1 (medium)
cri6_low_upsampled = resample(cri6_low,
                              replace=True,      # sample with replacement
                              n_samples=22640,    # to match majority class
                              random_state=101)

cri6_high_upsampled = resample(cri6_high,
                               replace=True,      # sample with replacement
                               n_samples=22640,    # to match majority class
                               random_state=101)

# Combine majority class with upsampled minority class
cri6_upsampled = pd.concat([cri6_medium, cri6_low_upsampled, cri6_high_upsampled])
```

*Gambar 4. 59 Balanced Dataset*

Untuk melakukan prediksi, diperlukan beberapa library seperti berikut

```
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.ensemble import RandomForestClassifier
from sklearn import metrics
from sklearn.metrics import confusion_matrix, classification_report
import joblib
import matplotlib.pyplot as plt
```

*Gambar 4. 60 Library*

`Train_test_split` = merupakan *library* yang digunakan untuk membagi dataset menjadi train dan test.

`Standarscaler` = untuk mengubah data sedemikian rupa sehingga distribusinya akan memiliki nilai rata-rata 0 dan deviasi standar 1.

`RandomForestClassfier` = merupakan model yang digunakan dalam prediksi ini.

`Metrics` = merupakan ukuran penilaian kuantitatif untuk menilai, membandingkan, dan melacak kinerja.

Confusion matrix = tabel yang digunakan untuk menggambarkan kinerja model klasifikasi pada set data pengujian yang nilainya telah diketahui.

Classification\_report = untuk memberikan laporan hasil prediksi.

Matplotlib = digunakan untuk membuat diagram.

Selanjutnya lakukan proses prediksi dengan model random forest.

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=101)

classifier = RandomForestClassifier(n_estimators = 1000, criterion = 'entropy', random_state = 0)
classifier.fit(X_train, y_train)

y_pred = classifier.predict(X_test)
```

*Gambar 4. 61 Proses Prediksi*

X\_train = adalah variable dataset yang digunakan untuk proses training. Pada proses training dilakukan sebanyak 75 % dataset yang digunakan.

X\_test = adalah variable dataset yang digunakan untuk proses test dataset. Pada test dataset tersebut dilakukan sebanyak 25% dari data yang digunakan.

Train\_test\_split = adalah memecah dataset training dan test.

Test\_size = adalah ukuran dataset yang diuji, dimana disini ada 0.25 berarti 25% dataset yang diuji.

n\_estimators = merupakan jumlah pohon yang diuji, disini ada 1000 pohon.

random\_state = adalah status acak dari random forest tersebut.

kemudian untuk menampilkan hasil dari prediksi.

```
y_pred = classifier.predict(X_test)
print(y_pred)
```

*Gambar 4. 62 Untuk Menampilkan Hasil Prediksi*

Nilai y\_pred diambil dari nilai classifier, selanjutnya classifier tersebut diberi fungsi predict untuk memberikan hasil prediksi. Hasil prediksi sendiri diambil dari classifier.predict(X\_test), jadi data yang ditampilkan adalah data dari nilai X\_test. Untuk menampilkan hasil prediksi tersebut, maka disini menggunakan

perintah *print* diikuti dengan variable *y\_pred* untuk mengeluarkan hasil dari nilai prediksi *y\_pred* tersebut.

Pada saat tampil hasil random forest tersebut berbentuk array seperti ini. Hasil ini didapat setelah mengeluarkan hasil prediksi dari *y\_pred*.

```
[2 0 0 ... 2 0 0]
```

Gambar 4. 63 Hasil Random Forest Array

Selanjutnya untuk menghasilkan nilai hasil prediksi tersebut.

```
cm = pd.crosstab(y_test, y_pred, rownames=['Kejahatan Sebenarnya'], colnames=['Prediksi Kejahatan'])  
print(cm)
```

Gambar 4. 64 Hasil Random Forest

*cm* = merupakan variable untuk menampung nilai.

*pd* = merupakan library dari *pandas*.

*Crosstab* = merupakan fungsi yang tersedia di dalam *library* *pandas* untuk menghitung tabulasi silang dari dua (atau lebih) factor.

*y\_test* = merupakan variabel untuk menampung nilai test.

*y\_pred* = merupakan variable untuk menampung nilai prediksi.

*print* = digunakan untuk menampilkan nilai. Disini nilai yang ditampilkan yaitu *cm*.

berikut tampilan hasil dari prediksi nilai *cm*.

Prediksi Kejahatan	0	1	2
Kejahatan Sebenarnya			
0	5499	177	0
1	478	422	2
2	0	0	5640

Gambar 4. 65 Hasil Prediksi

Untuk menampilkan nilai *accuracy* dari hasil model yang diprediksi.

```
print("Accuracy:",(metrics.accuracy_score(y_test, y_pred)*100),"\n")
```

*Gambar 4. 66 Menampilkan Hasil Accuracy*

Berikut tampilan hasil *accuracy*

```
Accuracy: 94.62268783761662
```

*Gambar 4. 67 Hasil Accuracy*

Untuk menampilkan hasil *classification report* tersebut.

```
print("\n-----Classification Report-----")
print(classification_report(y_test,y_pred))
```

*Gambar 4. 68 Classification Report*

Berikut tampilan hasil *classification report*.

```
-----Classification Report-----
              precision    recall  f1-score   support

     0           0.92       0.97       0.94       5676
     1           0.70       0.47       0.56        902
     2           1.00       1.00       1.00       5640

 accuracy          0.87       0.81       0.84      12218
  macro avg         0.87       0.81       0.84      12218
 weighted avg         0.94       0.95       0.94      12218
```

*Gambar 4. 69 Hasil Classification Report*

Disini terdapat beberapa keterangan.

*Accuracy* = Akurasi merupakan tingkat ketepatan atau seberapa dekat suatu pengukuran dengan nilai sebenarnya.

*Precision* = Mendeskripsikan kemungkinan suatu sampel berada di kelas X jika sudah diprediksi masuk ke dalam kelas X.

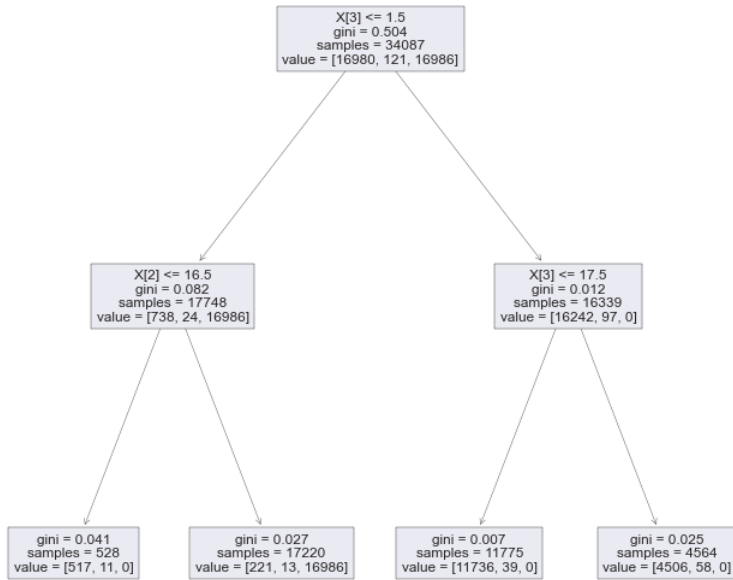
*Recall* = Menjelaskan kemungkinan model untuk mengklasifikasikan sampel dengan benar sebagai milik kelas X jika itu benar benar milik kelas X.

F1 = Ukuran yang menggambarkan hubungan antara presisi dan mengingat. Semakin tinggi skor F1 untuk suatu kelas, semakin baik kinerja model dalam mengklasifikasi kelas tersebut.

*macro avg* = merupakan nilai rata rata macro

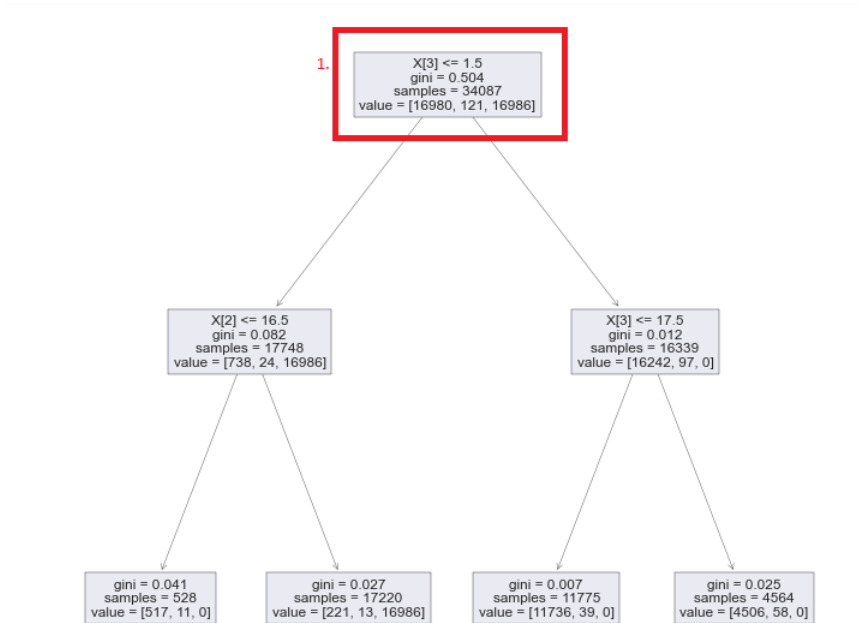
*weight avg* = rata-rata yang memperhitungkan pentingnya setiap angka dalam pembuatan rata-rata.

Selanjutnya Berikut adalah tampilan pohon keputusan



*Gambar 4. 70 Pohon Keputusan*

Berikut penjelasan dari pohon keputusan

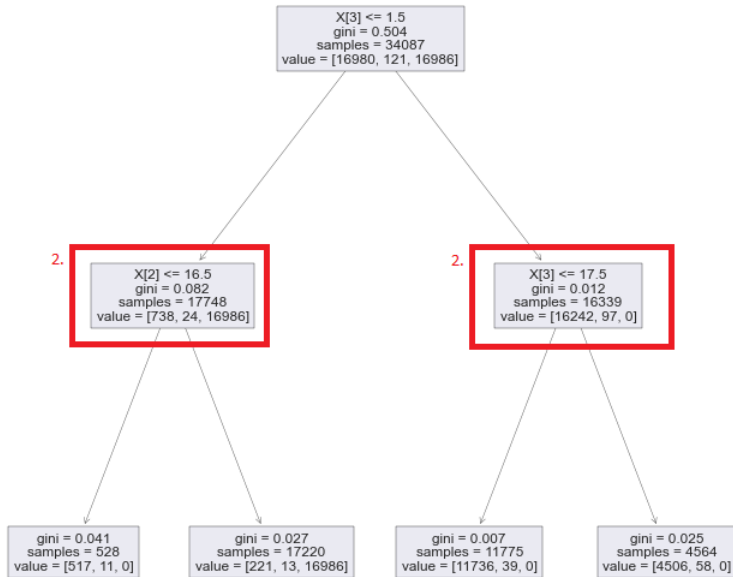


*Gambar 4. 71 Root Node*

Pada bagian no 1 yang ditandai dinamakan dengan *root node* (akar), *root node* (akar) merupakan tujuan akhir atau keputusan besar yang ingin diambil. [35]



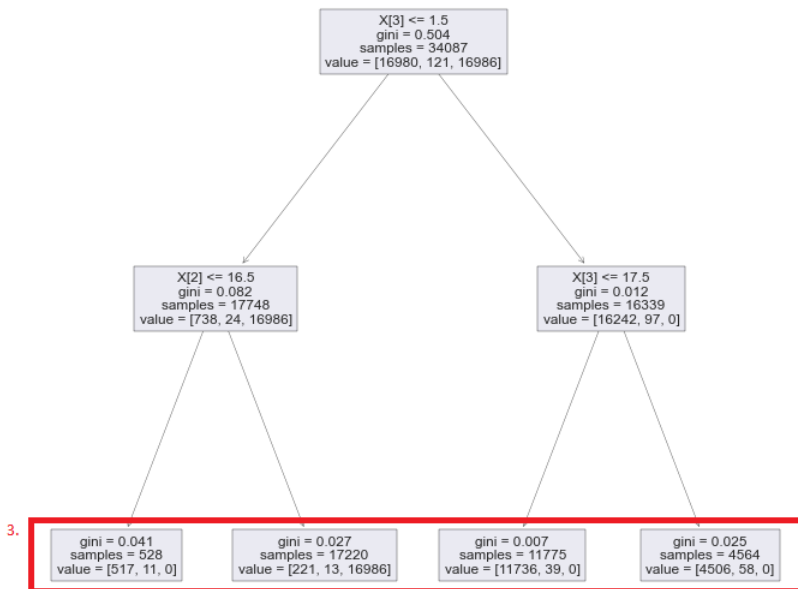
Selanjutnya terdapat bagian *branches* (ranting)



Gambar 4. 72 Branches

*Branches* (ranting) merupakan berbagai pilihan tindakan. Tindakan yang diambil dalam pohon keputusan setelah dari akar. [35]

Selanjutnya ada bagian *leaf node* (daun)



Gambar 4. 73 Leaf Node

*Leaf node* (daun) merupakan berbagai pilihan tindakan. Setelah dari *branches* (ranting) selanjutnya dalam pohon keputusan memilih tindakan apa yang dilakukan untuk proses selanjutnya.[35]



## BAB V

# KESIMPULAN DAN SARAN

### 5.1 Kesimpulan

Pada hasil evaluasi menggunakan model random forest, didapatkan *accuracy* pengujian model tersebut sebesar 94% dari hasil pengujian data.

Pada kejadian diatas, didapatkan kesimpulan bahwa tempat yang sering kejadian berada di *street* (jalan), pada bulan ke 1 (Januari), hari ke 4 (Kamis) dan jam 12.00.

### 5.2 Saran

Berikut beberapa saran yang bisa dilakukan pada penelitian ini diantaranya :

Pada sistem ini dapat dikembangkan untuk melakukan prediksi pada tingkat lebih dari satu jenis kejahatan yang lain dan ruang lingkup yang lebih luas dalam konteks jenis kejahatan.

## DAFTAR PUSTAKA

- [1] D. Winarti, M. Kom, E. Revita, dan M. Kom, “Penerapan Data Mining untuk Analisa Tingkat Kriminalitas Dengan Algoritma Association Rule Metode FP-Growth,” *J. SIMTIKA*, vol. 4, no. 3, hal. 8–22, 2021.
- [2] R. Yadhunath, S. Srikanth, A. Sudheer, dan S. Palaniswamy, “Identification of Criminal Activity Hotspots using Machine Learning to Aid in Effective Utilization of Police Patrolling in Cities with High Crime Rates,” *CSITSS 2019 - 2019 4th Int. Conf. Comput. Syst. Inf. Technol. Sustain. Solut. Proc.*, 2019, doi: 10.1109/CSITSS47250.2019.9031057.
- [3] N. O. Syamsiah dan I. Purwandani, “Penerapan Ensemble Stacking untuk Peramalan Laba Bersih Bank Syariah Indonesia (BSI),” *Build. Informatics, Technol. Sci.*, vol. 3, no. 3, hal. 295–301, 2021, doi: 10.47065/bits.v3i3.1017.
- [4] K. Kumari, M. Bhardwaj, dan S. Sharma, “OSEMN Approach for Real Time Data Analysis,” *Int. J. Eng. Manag. Res.*, vol. 10, no. 02, hal. 107–110, 2020, doi: 10.31033/ijemr.10.2.11.
- [5] I. Pengaruh *et al.*, “Identifikasi Pengaruh Sistem Keamanan Lingkungan Terhadap Tingkat Kejahatan Pencurian Di Kota Surakarta Dengan Metode Sistem Informasi Geografis,” *J. Geod. Undip*, vol. 8, no. 1, hal. 398–407, 2019.
- [6] W. Safat, S. Asghar, dan S. A. Gillani, “Empirical Analysis for Crime Prediction and Forecasting Using Machine Learning and Deep Learning Techniques,” *IEEE Access*, vol. 9, hal. 70080–70094, 2021, doi: 10.1109/ACCESS.2021.3078117.
- [7] M. Rianto dan R. Yunis, “Analisis Runtun Waktu Untuk Memprediksi Jumlah Mahasiswa Baru Dengan Model Random Forest,” *Paradig. - J. Komput. dan Inform.*, vol. 23, no. 1, 2021, doi: 10.31294/p.v23i1.9781.
- [8] R. Rusmiati, S. Syahrizal, dan M. Din, “Konsep Pencurian Dalam

- Kitab Undang-Undang Hukum Pidana dan Hukum Pidana Islam,” *Syiah Kuala Law J.*, vol. 1, no. 1, hal. 339–352, 2018, doi: 10.24815/sklj.v1i1.12318.
- [9] A. N. Syahrudin dan T. Kurniawan, “Input Dan Output Pada Bahasa,” *J. Dasar Pemrograman Python STMIK*, no. January, hal. 1–7, 2018.
- [10] T. Wahyono, *Fundamental Of Python For Machine Learning*, Cetakan I. Yogyakarta: Gava Media, 2018.
- [11] R. I. F. Ibadurrohman, D. R. Wijaya, dan E. Hernawati, “Pengembangan Aplikasi Machine Learning Menggunakan Algoritma Support Vector Regression Dan Statistical-based Feature Selection Untuk Memprediksi Kemiskinan,” *eProceedings Appl. Sci.*, vol. 6, no. 2, hal. 1910–1917, 2020, [Daring]. Tersedia pada: <https://openlibrarypublications.telkomuniversity.ac.id/index.php/applidscience/article/view/12271>.
- [12] E. Retnoningsih dan R. Pramudita, “Mengenal Machine Learning Dengan Teknik Supervised Dan Unsupervised Learning Menggunakan Python,” *Bina Insa. Ict J.*, vol. 7, no. 2, hal. 156, 2020, doi: 10.51211/biict.v7i2.1422.
- [13] I. D. Id, *MACHINE LEARNING : Teori, Studi Kasus dan Implementasi Menggunakan Python*, Edisi I. Riau: UR PRESS, 2021.
- [14] N. Buslim dan R. P. Iswara, “Pengembangan Algoritma Unsupervised Learning Technique Pada Big Data Analysis di Media Sosial sebagai media promosi Online Bagi Masyarakat,” *J. Tek. Inform.*, vol. 12, no. 1, hal. 79–96, 2019, doi: 10.15408/jti.v12i1.11342.
- [15] A. Roihan, P. A. Sunarya, dan A. S. Rafika, “Pemanfaatan Machine Learning dalam Berbagai Bidang: Review paper,” *IJCIT (Indonesian J. Comput. Inf. Technol.)*, vol. 5, no. 1, hal. 75–82, 2020, doi: 10.31294/ijcit.v5i1.7951.
- [16] A. R. Muslikh, H. A. Santoso, A. Marjuni, P. Teknik, I. Universitas, dan D. Nuswantoro, “Klasifikasi Data Time Series Arus Lalu Lintas,” vol. 14, hal. 24–38, 2018.
- [17] Adversenia, “10 Manfaat dan Contoh Penerapan Machine Learning,” [www.advernesia.com](https://www.advernesia.com/blog/data-science/10-manfaat-dan-contoh-penerapan-machine-learning/). <https://www.advernesia.com/blog/data-science/10-manfaat-dan-contoh-penerapan-machine-learning/>.

- [18] M. Harahap, F. Rozi, Y. Yennimar, dan S. D. Siregar, “Analisis Wawasan Penjualan Supermarket dengan Data Science,” *Data Sci. Indones.*, vol. 1, no. 1, hal. 1–7, 2021, doi: 10.47709/dsi.v1i1.1173.
- [19] N. Rahmalia, “Data Science: Arti, Manfaat, Proses, dan Contoh Penerapannya,” *glints.com*, 2020. <https://glints.com/id/lowongan/data-science-adalah/>.
- [20] “17 Data Science Applications and Examples,” *builtin.com*, 2019. <https://builtin.com/data-science/data-science-applications-examples>.
- [21] J. Sanjaya, E. Renata, V. E. Budiman, F. Anderson, dan M. Ayub, “Prediksi Kelalaian Pinjaman Bank Menggunakan Random Forest dan Adaptive Boosting,” *J. Tek. Inform. dan Sist. Inf.*, vol. 6, no. 1, hal. 50–60, 2020, doi: 10.28932/jutisi.v6i1.2313.
- [22] A. Rosadi *et al.*, “Analisis Sentimen Berdasarkan Opini Pengguna pada Media Twitter Terhadap BPJS Menggunakan Metode Lexicon Based dan Naïve Bayes Classifier,” *J. Ilm. Komputasi*, vol. 20, no. 1, hal. 39–52, 2021, doi: 10.32409/jikstik.20.1.401.
- [23] S. Ghosh, “Imbalanced vs Balanced Dataset in Machine Learning,” *medium.com*, 2019. <https://medium.com/open-datascience/imbalanced-vs-balanced-dataset-in-machine-learning-4faec5629b7e>.
- [24] K. Kurniawan dan D. Antoni, “Visualisasi Data Penduduk Dalam Membangun E-government Berbasis Sistem Informasi Geografis (GIS),” *J. Sisfokom (Sistem Inf. dan Komputer)*, vol. 9, no. 3, hal. 310–316, 2020, doi: 10.32736/sisfokom.v9i3.828.
- [25] R. E. Nalawati dan D. Y. Liliana, “Visualisasi Data Program Vaksinasi Covid-19 di Kota Depok dengan Big Data Analytics,” *J. Media Inform. ...*, vol. 5, hal. 1570–1579, 2021, doi: 10.30865/mib.v5i4.3330.
- [26] Yasir Abdur Rohman, “Pengenalan NumPy, Pandas, Matplotlib,” *medium.com*, 2019. <https://medium.com/@yasirabd/pengenalan-numpy-pandas-matplotlib-b90bafd36c0>.
- [27] H. Mccaslin, “Visualisai Data Menggunakan Seaborn dan Matplotlib | Python,” *medium.com*, 2020. <https://medium.com/@HakimMccaslin/apa-itu-seaborn-234a224d946f>.

- [28] Muthukrishnan, “Understanding the Classification report through sklearn,” *medium.com*, 2018. <https://muthu.co/understanding-the-classification-report-in-sklearn/>.
- [29] M. Z. Asy’ari, “Cara Menggunakan Jupyter Notebook dengan Mudah,” *auftechnique.com*, 2020. <https://auftechnique.com/cara-menggunakan-jupyter-notebook/>.
- [30] A. U. Jamila, B. M. Siregar, dan R. Yunis, “Analisis Runtun Waktu Untuk Memprediksi Jumlah Mahasiswa Baru Dengan Model Arima,” *Paradig. - J. Komput. dan Inform.*, vol. 23, no. 1, hal. 99–105, 2021, doi: 10.31294/p.v23i1.9758.
- [31] F. Hamami dan I. A. Dahlan, “Klasifikasi Cuaca Provinsi Dki Jakarta Menggunakan Algoritma Random Forest Dengan Teknik Oversampling,” *J. Teknoinfo*, vol. 16, no. 1, hal. 87, 2022, doi: 10.33365/jti.v16i1.1533.
- [32] S. Devella, Y. Yohannes, dan F. N. Rahmawati, “Implementasi Random Forest Untuk Klasifikasi Motif Songket Palembang Berdasarkan SIFT,” *JATISI (Jurnal Tek. Inform. dan Sist. Informasi)*, vol. 7, no. 2, hal. 310–320, 2020, doi: 10.35957/jatisi.v7i2.289.
- [33] F. Dan, R. Forest, dan R. F. Method, “( Cow Weight Classification Based on Digital Image Using Fractal and,” vol. 8, no. 2, hal. 1472–1480, 2021.
- [34] S. M. Rezkia, “Macam-Macam Metode Analisis Data: 2 Macam Metode Penting dalam Mengolah Data,” 2020. <https://dqlab.id/data-analisis-pahami-2-metode-analisis-data>.
- [35] G. N. Arviana, “Coba Metode Decision Tree bagi Kamu yang Sulit Ambil Keputusan,” *glints.com*, 2020. <https://glints.com/id/lowongan/decision-tree-adalah/#.YiBk-KtBzIU>.



# LAMPIRAN – LAMPIRAN

FORMAT PENILAIAN INTERNSHIP PROGRAM STUDI D4 TEKNIK INFORMATIKA POLITEKNIK POS INDONESIA				
N A M A		N P M	Tempat Tgl Lahir	
PARHAN HAMBALI		1184042	Bandung, 15 Oktober 2000	
JUDUL INTERNSHIP :		Rancang Bangun Sistem Sinkronisasi Surat Keputusan Direksi (Studi Kasus PT Pos Indonesia)		
PEMBIMBING EKSTERNAL : Suprpto				
NO	KOMPONEN YANG DI NILAI	NILAI MAKS	PENILAIAN (ANGKA)	RATA - RATA
1.	PENAMPILAN INDIVIDUAL			
	A. Penampilan Berpakaian	7	6	
	B. Sikap Terhadap Orang Lain	8	6	
	C. Semangat Bekerja	7	6	
	D. Kematangan Dalam Bertindak	6	6	
	E. Adaptasi Tempat Kerja	6	6	
	F. Pengetahuan Yang Mendukung Pekerjaan	6	6	
	G. Kehadiran Ditempat Kerja	8	6	
2.	KINERJA PKL			
	A. Ketelitian & Ketepatan Dalam Bekerja	8	6	
	B. Kualitas Produk / Kerja	8	6	
	C. Kemandirian Dalam Melaksanakan Pekerjaan	7	6	
	D. Kemampuan Bekerjasama	7	6	
	E. Kemampuan Berkomunikasi	8	6	
	F. Manajemen Waktu	7	6	
	G. Kemampuan Menganalisa Masalah	7	6	
TOTAL			84	

BANDUNG, 31 Desember 2021

BEKTERIAT PERUSAHAAN  
POS INDONESIA  
7.1

(SUPRPTO)  
NIPPOS : 969349427

### Sinopsis

Pencurian merupakan kejahatan yang ditujukan terhadap harta benda dan paling sering terjadi di dalam masyarakat. Kejahatan ini merupakan tindakan kejahatan yang dapat mengguncangkan stabilitas keamanan baik terhadap harta maupun terhadap jiwa masyarakat. Oleh karena itu, baik dalam Kitab Undang - Undang Hukum Pidana (KUHP) maupun dalam Kitab Suci melarang keras tindakan kejahatan tersebut dan menegaskan ancaman hukuman secara rinci dan berat atas diri pelanggarnya. Hal ini dapat dilihat dari bentuk hukuman dan ancaman hukuman yang dijatuhkan. Untuk meminimalisir pencurian, dari daerah yang mempunyai tingkat kejahatan tinggi. Maka pada tahap ini membuat klasifikasi pelabelan untuk mengetahui tingkat keamanan suatu daerah, selanjutnya dilakukan prediksi seberapa akurat prediksi tersebut menggunakan model OSEMN. Penentuan model menggunakan Random Forest dipilih karena Random Forest merupakan salah satu metode yang digunakan untuk klasifikasi dan regresi. Random Forest dapat digunakan untuk time series dengan cara meningkatkan keakuratan metode klasifikasi dengan menggabungkan metode klasifikasi dalam kata lain Random Forest dilakukan secara ensemble.