

Proposal Stage 0

Agi Rahmawandi as Data Scientist

Shan Ramadhan as Data Analyst

Muhammad Muqorrobin as Business Analyst

I Gusti Ngurah Agung Hari Vijaya Kusuma as PM/DE

August 28, 2025

Submission Links

- Repository: github.com/4Kings-Rakamin
- ManPro & Timeline: [GoogleSheet_Timeline](#)

1 Pendahuluan

1.1 Latar Belakang

Dalam dunia perekrutan karyawan, pengambilan keputusan yang tepat sangat penting untuk memastikan bahwa perusahaan mendapatkan kandidat yang paling sesuai dengan kebutuhan dan budaya organisasi. Dengan kemajuan teknologi dan analisis data, perusahaan kini dapat memanfaatkan data historis untuk meningkatkan proses perekrutan mereka. Dataset yang digunakan dalam proyek ini berisi informasi tentang berbagai kandidat yang telah melamar pekerjaan di sebuah perusahaan, termasuk fitur-fitur seperti umur, jenis kelamin, tingkat pendidikan, pengalaman kerja, dan skor wawancara. Dengan menganalisis data ini, kita dapat mengidentifikasi pola dan faktor-faktor yang mempengaruhi keputusan perekrutan, sehingga membantu perusahaan dalam membuat keputusan yang lebih baik di masa depan.

1.2 Riset Bisnis

Dalam konteks bisnis, proses perekrutan yang efisien dan efektif sangat penting untuk keberhasilan jangka panjang perusahaan. Dengan menggunakan data historis dari proses perekrutan sebelumnya, perusahaan dapat mengidentifikasi karakteristik kandidat yang paling berhasil dan sesuai dengan kebutuhan organisasi. Hal ini tidak hanya membantu dalam mengurangi biaya dan waktu yang dihabiskan untuk proses perekrutan, tetapi juga meningkatkan kualitas karyawan yang direkrut. Dengan demikian, analisis data ini dapat memberikan wawasan berharga bagi tim HR dan manajemen dalam mengoptimalkan strategi perekrutan mereka.

Adapun penggunaan model machine learning dalam proses perekrutan ini dapat membantu perusahaan dalam mengelola biaya secara lebih efisien. Dengan memprediksi kandidat yang memiliki kemungkinan besar untuk diterima berdasarkan data historis, perusahaan dapat mengurangi jumlah pengiriman email penolakan. Hal ini tidak hanya menghemat waktu dan sumber daya, tetapi juga memungkinkan tim HR untuk lebih fokus pada aspek-aspek lain dari proses perekrutan, seperti pengembangan karyawan dan retensi. Dengan demikian, implementasi model machine learning dalam proses perekrutan dapat memberikan manfaat ekonomi yang signifikan bagi perusahaan.

1.3 Problem Statement

Proses rekrutmen sering menghadapi berbagai tantangan, seperti banyaknya jumlah pelamar, keterbatasan waktu dalam melakukan penilaian, serta tingginya tingkat subjektivitas yang dapat memengaruhi konsistensi pengambilan keputusan. Kompleksitas parameter penilaian yang mencakup kompetensi, keterampilan, dan faktor demografi semakin meningkatkan risiko perusahaan gagal merekrut kandidat potensial. Dataset yang digunakan dalam penelitian ini (recruitment_data.csv) berisi 1.500 data historis kandidat yang mencakup usia, gender, tingkat pendidikan, pengalaman kerja, jumlah perusahaan sebelumnya, jarak tempat tinggal dari kantor, skor wawancara, skor keterampilan, skor kepribadian, strategi rekrutmen, serta keputusan akhir perekrutan (HiringDecision).

Dari hasil analisis awal, ditemukan ketidakseimbangan yang signifikan dalam distribusi kelas, di mana sekitar 69% kandidat diterima dan 31% ditolak. Kondisi ini menyebabkan perusahaan harus mengirim banyak email penolakan, yang pada akhirnya meningkatkan biaya operasional dan menurunkan efisiensi rekrutmen. Oleh karena itu, proyek ini bertujuan untuk mengembangkan model machine learning yang mampu memprediksi keputusan perekrutan secara akurat sekaligus menangani ketidakseimbangan kelas, sehingga perusahaan dapat mengoptimalkan proses rekrutmen dan meminimalisasi risiko kegagalan dalam merekrut kandidat potensial.

1.4 Goals, Objectives, and Business Metrics

Tujuan utama dari proyek ini adalah membangun model machine learning yang dapat membantu proses pengambilan keputusan perekrutan karyawan secara lebih cepat, objektif, dan konsisten. Dengan memanfaatkan data historis perekrutan, model ini diharapkan mampu memberikan rekomendasi kandidat potensial serta mengurangi bias subjektif dalam proses seleksi.

Objectives:

1. Mengembangkan model prediktif berdasarkan data historis recruitment untuk mengklasifikasikan kandidat apakah diterima atau tidak.
2. Mengidentifikasi kandidat potensial yang memiliki probabilitas keberhasilan tinggi.
3. Mempercepat proses identifikasi kandidat sebagai penilaian awal HR dalam mengambil keputusan.

4. Mengurangi bias subjektif dengan menyediakan sistem pendukung keputusan berbasis data.
5. Menyediakan interpretasi hasil model melalui analisis feature importance agar perusahaan mengetahui faktor yang paling memengaruhi keputusan perekrutan.

Business Metrics:

1. **Akurasi Model:** Akurasi prediksi $\geq 85\%$ untuk memastikan keputusan perekrutan yang lebih tepat.
2. **Precision:** Precision $\geq 80\%$, agar recruiter dapat mengurangi jumlah kandidat “false positive” yang tidak sesuai.
3. **Recall dan F1-Score:** Recall yang baik untuk meminimalisasi hilangnya kandidat potensial, dengan target F1-Score $\geq 75\%$ agar model seimbang antara precision dan recall.
4. **AUC:** Nilai AUC $\geq 85\%$ untuk menjamin kestabilan performa model pada berbagai threshold dalam proses shortlisting kandidat.
5. **Efisiensi Waktu:** Mengurangi waktu screening kandidat sehingga proses rekrutmen lebih cepat dan efisien.
6. **Reduksi Biaya:** Menurunkan biaya operasional terkait dengan pengiriman email penolakan dan proses administrasi perekrutan.

1.5 Gap Analysis

Saat ini belum tersedia sistem prediksi keputusan perekrutan berbasis data yang cepat dan objektif. Proses pengambilan keputusan HR masih didominasi oleh penilaian subjektif, membutuhkan waktu yang relatif lama, serta tidak memberikan insight yang jelas terkait faktor-faktor yang memengaruhi keputusan, seperti kompetensi, keterampilan, dan aspek demografi kandidat. Kondisi ini berpotensi menurunkan konsistensi dan efektivitas dalam proses seleksi karyawan baru.

Dengan adanya model machine learning yang diusulkan, diharapkan dapat mengisi gap ini dengan menyediakan alat bantu yang mampu menganalisis data historis secara cepat dan akurat. Model ini akan membantu HR dalam mengidentifikasi kandidat potensial berdasarkan fitur-fitur yang relevan, sehingga mengurangi ketergantungan pada penilaian subjektif dan mempercepat proses seleksi. Selain itu, model ini juga akan memberikan wawasan tentang faktor-faktor kunci yang memengaruhi keputusan perekrutan, sehingga perusahaan dapat mengoptimalkan strategi rekrutmen mereka ke depannya.

1.6 Ideal Condition & Expected Impact

Sistem mampu mengidentifikasi kandidat potensial secara cepat, objektif, dan tepat sasaran melalui pemetaan kandidat yang layak direkrut. Dampaknya, proses pengambilan keputusan menjadi lebih efisien sehingga HR dapat memfokuskan waktu pada kandidat yang benar-benar potensial dan sesuai kebutuhan perusahaan.

2 Tinjauan Pustaka

2.1 Proses Rekrutmen

Proses rekrutmen adalah serangkaian langkah yang diambil oleh perusahaan untuk menarik, menilai, dan memilih kandidat yang paling sesuai untuk mengisi posisi yang tersedia. Proses ini biasanya dimulai dengan identifikasi kebutuhan tenaga kerja, diikuti oleh pencarian kandidat melalui berbagai saluran seperti iklan lowongan kerja, agen perekrutan, dan media sosial. Setelah kandidat ditemukan, mereka akan melalui tahap seleksi yang melibatkan penilaian kualifikasi, wawancara, dan tes keterampilan. Akhirnya, keputusan perekrutan dibuat berdasarkan evaluasi menyeluruh dari semua kandidat yang telah melalui proses seleksi. Proses rekrutmen yang efektif tidak hanya memastikan bahwa perusahaan mendapatkan karyawan yang berkualitas, tetapi juga membantu dalam membangun budaya organisasi yang positif dan mendukung tujuan bisnis jangka panjang. (Mathis et al., 2017)

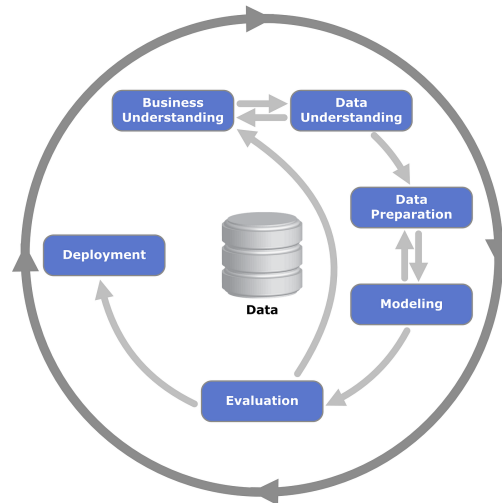
Adapun beberapa faktor yang mempengaruhi keputusan perekrutan meliputi:

1. Demografis: usia, gender, latar belakang pendidikan. (Ng & Burke, 2005)
2. Pengalaman kerja: jumlah tahun pengalaman dan variasi perusahaan sebelumnya. (Ployhart, 2006)
3. Kompetensi teknis dan soft skills: hasil tes keterampilan (skill score), wawancara, dan penilaian kepribadian. (Schmidt & Hunter, 1998)
4. Faktor eksternal: jarak tempat tinggal ke kantor sering menjadi pertimbangan dalam retensi. (Hausknecht et al., 2009)
5. Strategi rekrutmen: pendekatan organisasi (job fairs, rekrutmen online, campus hiring) dapat memengaruhi kualitas kandidat. (Breaugh, 2013).

Berdasarkan studi literatur, faktor-faktor di atas secara signifikan mempengaruhi keputusan perekrutan. Dimana hal tersebut dapat menjadi pertimbangan penting dalam pengembangan model prediktif untuk proses perekrutan.

2.2 CRISP-DM

CRISP-DM (Cross-Industry Standard Process for Data Mining) adalah metodologi standar yang digunakan dalam proyek data mining dan analisis data. Metodologi ini terdiri dari enam fase utama yang membantu dalam mengorganisir dan mengelola proyek data secara sistematis. (Chumbar, 2020)



Gambar 1: Alur Kerja CRISP-DM

Gambar 2 menggambarkan alur kerja CRISP-DM yang terdiri dari enam fase utama, yaitu:

1. Business Understanding: Memahami tujuan bisnis dan kebutuhan proyek.
2. Data Understanding: Mengumpulkan dan memahami data yang tersedia.
3. Data Preparation: Membersihkan dan mempersiapkan data untuk analisis.
4. Modeling: Membangun model prediktif menggunakan teknik machine learning.
5. Evaluation: Mengevaluasi model untuk memastikan bahwa tujuan bisnis tercapai.
6. Deployment: Menerapkan model dalam lingkungan produksi untuk digunakan dalam pengambilan keputusan bisnis.

2.3 EDA (Exploratory Data Analysis)

Exploratory Data Analysis (EDA) adalah proses awal dalam analisis data yang bertujuan untuk memahami karakteristik dan pola dalam dataset. EDA melibatkan berbagai teknik visualisasi dan statistik untuk mengeksplorasi data, mengidentifikasi outlier, dan menemukan hubungan antara variabel. Proses ini sangat penting karena membantu dalam mengarahkan langkah-langkah selanjutnya dalam analisis data, seperti pemilihan fitur dan pemodelan. (Tukey, 1977)

2.4 T-Test

T-Test adalah metode statistik yang digunakan untuk membandingkan rata-rata dari dua kelompok data. Uji ini membantu menentukan apakah perbedaan antara kedua kelompok tersebut signifikan secara statistik atau hanya terjadi secara kebetulan. T-Test dapat digunakan dalam berbagai konteks, seperti membandingkan hasil tes antara dua kelompok siswa

atau mengevaluasi efektivitas dua metode pengajaran yang berbeda. Hasil dari uji T-Test memberikan nilai p-value yang digunakan untuk menilai signifikansi perbedaan antara kedua kelompok. (De Veaux et al., 2011)

Rumus untuk menghitung nilai T-Test (t) adalah sebagai berikut:

$$t = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \quad (1)$$

Dimana:

- t = nilai T-Test
- \bar{X}_1 = rata-rata kelompok pertama
- \bar{X}_2 = rata-rata kelompok kedua
- s_1^2 = varians kelompok pertama
- s_2^2 = varians kelompok kedua
- n_1 = ukuran sampel kelompok pertama
- n_2 = ukuran sampel kelompok kedua

Singkatnya uji t-test ini memudahkan kita dalam membandingkan dua kelompok data untuk menentukan apakah ada perbedaan yang signifikan antara keduanya. Apabila p-value lebih kecil dari tingkat signifikansi (misalnya 0,05), maka kita menolak hipotesis nol dan menyimpulkan bahwa ada perbedaan yang signifikan antara kedua kelompok tersebut. Ini berguna apabila fitur numerik ingin dibandingkan terhadap target kategorikal.

2.5 Chi Square Test

Chi Square Test adalah metode statistik yang digunakan untuk menguji hubungan antara dua variabel kategorikal. Uji ini membandingkan frekuensi yang diamati dalam setiap kategori dengan frekuensi yang diharapkan jika tidak ada hubungan antara variabel. Hasil dari uji Chi Square memberikan nilai p-value yang digunakan untuk menentukan apakah hubungan antara variabel tersebut signifikan secara statistik. (Agresti, 2018)

Rumus untuk menghitung nilai Chi Square (χ^2) adalah sebagai berikut:

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i} \quad (2)$$

Dimana:

- χ^2 = nilai Chi Square
- O_i = frekuensi yang diamati dalam kategori ke-i
- E_i = frekuensi yang diharapkan dalam kategori ke-i

- \sum = penjumlahan untuk semua kategori

Apabila p-value lebih kecil dari tingkat signifikansi (misalnya 0,05), maka kita menolak hipotesis nol dan menyimpulkan bahwa ada hubungan yang signifikan antara kedua variabel tersebut. Ini berguna apabila fitur kategorikal ingin dibandingkan terhadap target kategorikal.

2.6 Standard Scaler

Standard Scaler adalah teknik normalisasi data yang digunakan untuk mengubah fitur numerik sehingga memiliki rata-rata (mean) nol dan standar deviasi (standard deviation) satu. Proses ini membantu dalam mengurangi skala variabilitas antar fitur, sehingga model machine learning dapat belajar lebih efektif. Standard Scaler sangat berguna ketika fitur-fitur dalam dataset memiliki rentang nilai yang berbeda-beda, karena dapat meningkatkan konvergensi dan kinerja model. (Jain & Bhandare, 2016) Rumus untuk menghitung nilai yang telah dinormalisasi (z) menggunakan Standard Scaler adalah sebagai berikut:

$$z = \frac{(X - \mu)}{\sigma} \quad (3)$$

Dimana:

- z = nilai yang telah dinormalisasi
- X = nilai asli dari fitur
- μ = rata-rata (mean) dari fitur
- σ = standar deviasi (standard deviation) dari fitur
- $X - \mu$ = selisih antara nilai asli dan rata-rata
- $\frac{(X-\mu)}{\sigma}$ = hasil pembagian selisih dengan standar deviasi

Dengan menggunakan Standard Scaler, setiap fitur numerik dalam dataset akan diubah sehingga memiliki distribusi yang seragam, yang pada gilirannya dapat meningkatkan performa model machine learning.

2.7 Logistic Regression

Logistic Regression adalah metode statistik yang digunakan untuk memodelkan hubungan antara satu atau lebih variabel independen (fitur) dengan variabel dependen biner (target). Metode ini digunakan untuk memprediksi probabilitas kejadian suatu peristiwa, seperti apakah seorang kandidat akan diterima atau ditolak dalam proses perekrutan. Logistic Regression menggunakan fungsi logit untuk mengubah output linier menjadi probabilitas yang berada dalam rentang 0 hingga 1. (Hosmer et al., 2013) Rumus untuk menghitung probabilitas (p) menggunakan Logistic Regression adalah sebagai berikut:

$$p = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n)}} \quad (4)$$

Dimana:

- p = probabilitas kejadian (misalnya, kandidat diterima)
- e = basis dari logaritma natural (sekitar 2,718)
- β_0 = intercept (konstanta)
- $\beta_1, \beta_2, \dots, \beta_n$ = koefisien regresi untuk masing-masing fitur
- X_1, X_2, \dots, X_n = nilai dari fitur-fitur independen
- $\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n$ = kombinasi linier dari fitur-fitur

Logistic Regression sangat berguna dalam berbagai aplikasi, termasuk analisis risiko kredit, diagnosis medis, dan prediksi perilaku konsumen. Dengan kemampuannya untuk memberikan interpretasi yang jelas melalui koefisien regresi, metode ini memungkinkan pemahaman yang lebih baik tentang faktor-faktor yang mempengaruhi keputusan biner.

2.8 Decision Tree

Decision Tree adalah algoritma machine learning yang digunakan untuk klasifikasi dan regresi. Algoritma ini bekerja dengan membagi data menjadi subset berdasarkan fitur-fitur tertentu, sehingga membentuk struktur pohon yang terdiri dari simpul (nodes) dan cabang (branches). Setiap simpul mewakili keputusan berdasarkan nilai fitur, sementara cabang menghubungkan simpul-simpul tersebut. Proses pembagian data berlanjut hingga mencapai simpul daun (leaf nodes) yang memberikan prediksi akhir. Decision Tree sangat populer karena kemampuannya untuk menangani data kategorikal dan numerik, serta memberikan interpretasi yang mudah dipahami.

Rumus untuk menghitung impurity menggunakan Gini Index (G) adalah sebagai berikut:

$$G = 1 - \sum_{i=1}^C p_i^2 \quad (5)$$

Dimana:

- G = nilai Gini Index
- C = jumlah kelas dalam target
- p_i = proporsi dari kelas ke- i dalam subset data
- $\sum_{i=1}^C p_i^2$ = penjumlahan kuadrat proporsi untuk semua kelas
- $1 - \sum_{i=1}^C p_i^2$ = hasil pengurangan dari 1 dengan penjumlahan kuadrat proporsi

Decision Tree sangat berguna dalam berbagai aplikasi, termasuk diagnosis medis, analisis risiko kredit, dan prediksi perilaku konsumen. Dengan kemampuannya untuk memberikan interpretasi yang jelas melalui struktur pohon, metode ini memungkinkan pemahaman yang lebih baik tentang faktor-faktor yang mempengaruhi keputusan klasifikasi.

2.9 Random Forest

Random Forest adalah algoritma ensemble learning yang digunakan untuk klasifikasi dan regresi. Algoritma ini bekerja dengan membangun sejumlah besar pohon keputusan (decision trees) secara acak dan menggabungkan hasil prediksi dari masing-masing pohon untuk menghasilkan prediksi akhir yang lebih akurat dan stabil. Random Forest mengurangi risiko overfitting yang sering terjadi pada pohon keputusan tunggal dengan cara mengambil rata-rata (untuk regresi) atau mode (untuk klasifikasi) dari hasil prediksi semua pohon dalam hutan. (Breiman, 2001)

Rumus untuk menghitung prediksi akhir (\hat{y}) dalam Random Forest adalah sebagai berikut:

$$\hat{y} = \frac{1}{T} \sum_{t=1}^T h_t(X) \quad (6)$$

Dimana:

- \hat{y} = prediksi akhir dari Random Forest
- T = jumlah pohon keputusan dalam hutan
- $h_t(X)$ = prediksi dari pohon keputusan ke- t untuk input X
- $\sum_{t=1}^T h_t(X)$ = penjumlahan prediksi dari semua pohon keputusan
- $\frac{1}{T} \sum_{t=1}^T h_t(X)$ = hasil pembagian penjumlahan dengan jumlah pohon untuk mendapatkan rata-rata prediksi

Random Forest sangat berguna dalam berbagai aplikasi, termasuk diagnosis medis, analisis risiko kredit, dan prediksi perilaku konsumen. Dengan kemampuannya untuk menangani data yang kompleks dan memberikan interpretasi yang jelas melalui fitur penting (feature importance), metode ini memungkinkan pemahaman yang lebih baik tentang faktor-faktor yang mempengaruhi keputusan klasifikasi atau regresi.

2.10 XGBoost

XGBoost (Extreme Gradient Boosting) adalah algoritma machine learning yang digunakan untuk klasifikasi dan regresi. Algoritma ini merupakan implementasi dari teknik boosting yang menggabungkan beberapa model lemah (weak learners), biasanya pohon keputusan, untuk membentuk model yang lebih kuat dan akurat. XGBoost dikenal karena kemampuannya dalam menangani data besar dan kompleks, serta memberikan performa yang tinggi melalui optimasi paralel dan regularisasi.

Rumus untuk menghitung prediksi akhir (\hat{y}) dalam XGBoost adalah sebagai berikut:

$$\hat{y} = \sum_{k=1}^K f_k(X) \quad (7)$$

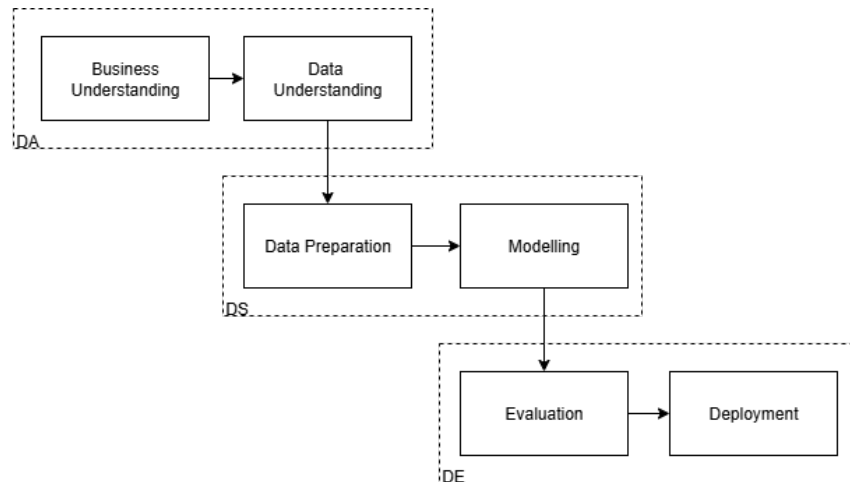
Dimana:

- \hat{y} = prediksi akhir dari XGBoost
- K = jumlah model lemah (pohon keputusan) yang digabungkan
- $f_k(X)$ = prediksi dari model lemah ke- k untuk input X
- $\sum_{k=1}^K f_k(X)$ = penjumlahan prediksi dari semua model lemah

XGBoost sangat berguna dalam berbagai aplikasi, termasuk diagnosis medis, analisis risiko kredit, dan prediksi perilaku konsumen. Dengan kemampuannya untuk menangani data yang kompleks dan memberikan interpretasi yang jelas melalui fitur penting (feature importance), metode ini memungkinkan pemahaman yang lebih baik tentang faktor-faktor yang mempengaruhi keputusan klasifikasi atau regresi.

3 Metodologi Penelitian

Metodologi penelitian yang akan digunakan dalam proyek ini adalah Cross Industry Standard Process for Data Mining (CRISP-DM). CRISP-DM adalah sebuah model proses yang terstruktur dan berulang yang terdiri dari enam fase utama, yaitu Business Understanding, Data Understanding, Data Preparation, Modeling, Evaluation, dan Deployment. Setiap fase memiliki tujuan dan aktivitas spesifik yang membantu dalam mengelola project ini secara efektif. Gambar 2 menunjukkan diagram alur metodologi yang berlandaskan prinsip CRISP-DM yang akan diikuti dalam penelitian ini.



Gambar 2: Diagram Alur Metodologi

Setiap fase dalam metodologi CRISP-DM memiliki peran penting dalam memastikan keberhasilan proyek ini. Fase Business Understanding bertujuan untuk memahami konteks bisnis dan mengidentifikasi tujuan yang ingin dicapai melalui analisis data. Fase Data Understanding melibatkan eksplorasi awal terhadap dataset untuk menilai kualitas data, mengidentifikasi pola, dan mendeteksi potensi masalah seperti data yang hilang atau outlier.

Fase Data Preparation fokus pada pembersihan dan transformasi data agar siap digunakan dalam proses pemodelan. Ini termasuk penanganan data yang hilang, normalisasi, dan encoding variabel kategorikal. Fase Modeling adalah tahap di mana berbagai algoritma machine learning diterapkan untuk membangun model prediktif berdasarkan data yang telah dipersiapkan.

Fase Evaluation melibatkan penilaian kinerja model menggunakan metrik yang relevan untuk memastikan bahwa model memenuhi kebutuhan bisnis yang telah ditetapkan. Terakhir, fase Deployment adalah tahap di mana model yang telah dievaluasi dan disetujui diintegrasikan ke dalam sistem bisnis yang ada, serta dipantau secara berkelanjutan untuk memastikan performa yang optimal di lingkungan nyata.

3.1 Peran Tim dalam Metodologi

Data analyst berperan penting dalam setiap fase pertama, mulai dari memahami kebutuhan bisnis, melakukan eksplorasi data, menyiapkan data untuk analisis, membangun dan mengevaluasi model, hingga memastikan bahwa model yang dihasilkan dapat diimplementasikan secara efektif dalam konteks bisnis.

Data scientist akan lebih fokus pada fase Modeling dan PreProcessing, di mana mereka akan menerapkan teknik machine learning yang lebih kompleks, melakukan tuning hyperparameter, serta mengevaluasi model dengan metrik yang lebih mendalam untuk memastikan bahwa model tidak hanya akurat tetapi juga dapat diinterpretasikan dan diandalkan.

Data engineer akan berperan utama dalam fase Deployment, di mana mereka akan memastikan bahwa model yang telah dikembangkan dapat diintegrasikan dengan lancar ke dalam infrastruktur teknologi yang ada. Mereka juga akan bertanggung jawab untuk membangun pipeline data yang efisien, mengelola penyimpanan data, serta memastikan bahwa sistem dapat menangani beban kerja yang diperlukan untuk menjalankan model secara real-time atau batch processing sesuai kebutuhan bisnis.

3.2 TimeLine Project

Berikut adalah rincian timeline proyek yang direncanakan untuk setiap fase dalam metodologi CRISP-DM, beserta estimasi waktu yang dibutuhkan untuk menyelesaikan masing-masing fase. Tabel 1 merangkum jadwal proyek secara keseluruhan.

Table 1: Timeline Project

Milestone	Aug W4	W1	September W2	W3	W4	Oct W1
Project Initiation & Problem Framing	PM & DA					
Data Acquisition & Preparation		DA & DS				
Model Development & Experimentation			DA & DS			
Model Evaluation & Interpretability				DS & BA		
Deployment & Business Integration					DS & DE	
Final Presentation						All Role

Agar lebih jelas lagi, berikut adalah penjabaran dari setiap fase beserta estimasi waktu yang dibutuhkan untuk per stage sesuai jadwal yang diberikan Rakamin Academy. Gambar 3 menunjukkan Timeline Stage0 Project secara keseluruhan.

No	Nama Aktivitas	Nama Task	Role	start date	due date	PIC	Progress Task
Stage 0	Melakukan riset terkait industri dari dataset yang dipilih	Riset terkait bisnis	Business Analyst	08/23/2025	08/30/2025		Done
		Riset terkait industri	Data Analyst	08/23/2025	08/30/2025		Done
		Identifikasi Kebutuhan data	Project Manager	08/23/2025	08/30/2025		Done
		Riset terkait data	Project Manager	08/23/2025	08/30/2025		Done
	Penyusunan Problem Statement & Business Understanding	Problem Statement	Business Analyst	08/23/2025	08/30/2025		Done
		Tujuan Bisnis	Data Analyst	08/23/2025	08/30/2025		Done
	Penyusunan Project Timeline	Draft TimeLine Project	Business Analyst	08/23/2025	08/30/2025		Done
		Validasi Workflow (Crisp DM)	Data Analyst	08/23/2025	08/30/2025		Done
	Risk & Feasibility Analysis	Identifikasi Resiko	Business Analyst	08/23/2025	08/30/2025		Done
		Review Dampak Bisnis	Business Analyst	08/23/2025	08/30/2025		Done
		Dokumentasi	Project Manager	08/23/2025	08/30/2025		Done
	Sesi Mentoring	Memilih Dataset, Pembahasan EDA, Workflow & Timeline	Data Engineer	08/23/2025	08/30/2025		In Pro...
		Review Mentoring	All	08/23/2025	08/30/2025		In Pro...
		Pembahasan Hasil Stage 0	All	08/23/2025	08/30/2025		In Pro...
	Dokumentasi (Github)	Page Organisasi (4Kings)	Project Manager	08/23/2025	08/30/2025		Done
		Repository untuk Stage 0	Project Manager	08/23/2025	08/30/2025		Done
	Evaluasi & Laporan	Penyusunan Proposal	All	08/23/2025	08/30/2025		In Pro...

Gambar 3: Timeline Stage0 Project

Timeline Stage 0 Project dimulai pada minggu ke-4 bulan Agustus dengan dengan detail seperti pada Gambar 3. Dibuat juga kolom progress task untuk menandai progress mana yang belum dikerjakan, sedang dikerjakan dan belum mulai dikerjakan. Gambar 4 menunjukkan Timeline Stage1 Project dengan progress task.

Stage1	Finalisasi Pemilihan Dataset	Pengecekan Dataset	Data Analyst	08/30/2025	09/06/2025		Not Yet
		Relevansi Dataset dengan Bisnis	Business Analyst	08/30/2025	09/06/2025		Not Yet
		Dokumentasi Dataset	Project Manager	08/30/2025	09/06/2025		Not Yet
	Exploratory Data Analysis (EDA)	Analisis Univariate&Multivariate pada Fitur	Data Analyst	08/30/2025	09/06/2025		Not Yet
		Visualisasi & Insight	Data Scientist	08/30/2025	09/06/2025		Not Yet
	Preprocessing Data	Handling missing values & duplicates	Data Scientist	08/30/2025	09/06/2025		Not Yet
		Outlier detection	Data Scientist	08/30/2025	09/06/2025		Not Yet
		Memilih outlier handling (IQR/ZScore)	Project Manager	08/30/2025	09/06/2025		Not Yet
	Feature Selection & Engineering	Uji Statistik (Chi Square/Anova)	Business Analyst	08/30/2025	09/06/2025		Not Yet
		Pemilihan Fitur yang relevan	Data Analyst	08/30/2025	09/06/2025		Not Yet
		Standarisasi, Encoding & Transform Log (Optional)	Project Manager	08/30/2025	09/06/2025		Not Yet
	Sesi Mentoring	Hasil Riset, Saran Model & Evaluasi Metric	All	08/30/2025	09/06/2025		Not Yet
		Review Mentoring	All	08/30/2025	09/06/2025		Not Yet
		Pembahasan Hasil Stage 1	All	08/30/2025	09/06/2025		Not Yet
	Dokumentasi (Github)	Transisi Proposal ke Report (LaTeX)	Project Manager	08/30/2025	09/06/2025		Not Yet
		Repository untuk Stage 1	Project Manager	08/30/2025	09/06/2025		Not Yet
	Evaluasi & Laporan	Mengisi Final Report	All	08/30/2025	09/06/2025		Not Yet

Gambar 4: Timeline Stage1 Project dengan Progress Task

Gambar 5 menunjukkan Timeline Stage2 Project dengan progress task.

Stage 2	Baseline Model Development	Menentukan Baseline Model (ex : Logistic Regression, DecisionTree dll)	Data Scientist	09/13/2025	09/20/2025	Not Yet
		Model Evaluasi & Interpretasi Bisnis	Data Scientist	09/13/2025	09/20/2025	Not Yet
		Deployment Kasar Coba" Techstack	Data Engineer	09/13/2025	09/20/2025	Not Yet
		Review Baseline model	Project Manager	09/13/2025	09/20/2025	Not Yet
	Experiment & Komparasi Model	Mencoba algoritma alternatif (ex : Random Forest, XGBoost, dll.)	Data Scientist	09/13/2025	09/20/2025	Not Yet
		Analisis & Dokumentasi Hasil Experiment	Project Manager	09/13/2025	09/20/2025	Not Yet
	Hyperparameter Tuning	Menentukan Tuning terbaik (ex : bisa RandomSearch/Grid dll)	Data Scientist	09/13/2025	09/20/2025	Not Yet
		Cross Validation Setup	Data Scientist	09/13/2025	09/20/2025	Not Yet
		Review Hasil Tuning	Project Manager	09/13/2025	09/20/2025	Not Yet
	Feature Selection & Engineering	Uji Statistik (Chi Square/Anova)	Data Analyst	09/13/2025	09/20/2025	Not Yet
		Pemilihan Fitur yang relevan	Project Manager	09/13/2025	09/20/2025	Not Yet
		Standarisasi, Encoding & Transform Log (Opsional)	Data Scientist	09/13/2025	09/20/2025	Not Yet
		Exporting Model	Data Engineer	09/13/2025	09/20/2025	Not Yet
	Pipeline Model	Menggunakan Teknik Processing yang sama untuk data baru	Data Engineer	09/13/2025	09/20/2025	Not Yet
		Review Pipeline	Project Manager	09/13/2025	09/20/2025	Not Yet
	Dokumentasi (Github)	Repository untuk Stage 2	Project Manager	09/13/2025	09/20/2025	Not Yet
	Evaluasi & Laporan	Mengisi Final Report	All	09/13/2025	09/20/2025	Not Yet

Gambar 5: Timeline Stage2 Project dengan Progress Task

Gambar 6 menunjukkan Timeline Stage3 Project dengan progress task.

Stage 3	Model Performance Evaluation	Evaluasi dengan metrik (Accuracy, Precision, Recall, F1, AUC)	Data Scientist	09/20/2025	09/27/2025	Not Yet
		Bandingkan baseline dan model model lain hasil Experiment & Tuning	Project Manager	09/20/2025	09/27/2025	Not Yet
	Feature Importance & Model Tradeoff	Shap Values fitur importance	Data Scientist	09/20/2025	09/27/2025	Not Yet
		Bias & Variance Tradeoff untuk model	Business Analyst	09/20/2025	09/27/2025	Not Yet
		Memilih Sesuai dengan Solusi Bisnis	Business Analyst	09/20/2025	09/27/2025	Not Yet
	Business Impact	Review dampak bisnis dari model	Business Analyst	09/20/2025	09/27/2025	Not Yet
		Penyelarasan hasil dengan tujuan Project	Project Manager	09/20/2025	09/27/2025	Not Yet
		Review Hasil Tuning	Project Manager	09/20/2025	09/27/2025	Not Yet
	Dokumentasi (Github)	Repository untuk Stage 3	Project Manager	09/20/2025	09/27/2025	Not Yet
	Evaluasi & Laporan	Mengisi Final Report	All	09/20/2025	09/27/2025	Not Yet

Gambar 6: Timeline Stage3 Project dengan Progress Task

Gambar 7 menunjukkan Timeline Stage4 Project dengan progress task.

Stage 4	Deployment & API Development	Implementasi model ke API	Data Engineer	09/27/2025	10/04/2025	Not Yet
		Pembuatan dashboard interaktif	Data Scientist	09/27/2025	10/04/2025	Not Yet
		Koordinasi Deployment	Project Manager	09/27/2025	10/04/2025	Not Yet
	Monitoring & Maintenance Strategy	Menentukan strategi monitoring model	Data Scientist	09/27/2025	10/04/2025	Not Yet
		Dokumentasi strategi monitoring	Project Manager	09/27/2025	10/04/2025	Not Yet
		Menyusun Pitch Deck Presentasi	All	09/27/2025	10/04/2025	Not Yet
	Final Presentation Preparation	Visualisasi hasil & insight	All	09/27/2025	10/04/2025	Not Yet
		Dokumentasi source code & pipeline	Project Manager	09/27/2025	10/04/2025	Not Yet
	Dokumentasi (Github)	Repository untuk Stage 4 & Readme Interaktif untuk masing" Repository	Project Manager	09/27/2025	10/04/2025	Not Yet
	Final Presentation Delivery	Presentasi hasil proyek & Penyerahan Laporan	All	09/27/2025	10/04/2025	Not Yet

Gambar 7: Timeline Stage4 Project dengan Progress Task

3.3 Risk & Feasibility Analysis

Dalam menjalankan proyek ini, terdapat beberapa risiko yang perlu diidentifikasi dan dianalisis untuk memastikan kelancaran proses project. Diharapkan dengan memahami potensi risiko yang ada, tim dapat merancang strategi mitigasi yang efektif guna mengurangi dampak negatif yang mungkin timbul. Tabel 2 merangkum berbagai aspek risiko, potensi risiko yang mungkin dihadapi, strategi mitigasi yang dapat diterapkan, serta penilaian kelayakan dari masing-masing risiko tersebut.

Table 2: Risk-Feasibility Analysis

Aspek	Potensi Risiko	Strategi Mitigasi	Kelayakan
Data	Data kandidat bisa tidak lengkap, tidak seimbang (imbalanced), atau mengandung bias (gender, usia).	Lakukan preprocessing, balancing data, feature engineering, serta audit fairness.	Layak jika dilakukan data cleaning & monitoring.
Teknis	Model bisa overfitting atau performa rendah di data baru.	Gunakan cross-validation, regularisasi, dan retraining berkala.	Layak dengan pipeline validasi yang baik.
Operasional	HR sulit mengadopsi sistem baru, lebih percaya screening manual.	Berikan pelatihan, buat antarmuka user-friendly, dan jelaskan transparansi model.	Layak jika ada kolaborasi dengan HR.
Etika & Regulasi	Risiko diskriminasi dalam keputusan perekrutan (misalnya gender bias).	Terapkan fairness metrics, hindari variabel sensitif sebagai faktor utama.	Layak dengan pengawasan etis & regulasi.
Ekonomi	Biaya implementasi dan maintenance model cukup tinggi.	Bandingkan cost vs benefit (efisiensi waktu, cost per hire, kualitas kandidat).	Layak jika ROI positif dalam 1–2 tahun.
Keberlanjutan	Model bisa usang (model drift) seiring perubahan tren pasar tenaga kerja.	Monitoring performa model, retraining dengan data terbaru setiap periode tertentu.	Layak dengan komitmen maintenance rutin.

Dengan melakukan analisis risiko ini, tim proyek dapat lebih siap dalam menghadapi tantangan yang mungkin muncul selama pelaksanaan proyek. Setiap risiko yang diidentifikasi telah diberikan strategi mitigasi yang spesifik, sehingga dapat diatasi dengan cara yang paling efektif. Selain itu, penilaian kelayakan dari setiap risiko membantu dalam menentukan prioritas tindakan yang perlu diambil untuk memastikan keberhasilan proyek secara keseluruhan.

3.4 Penjelasan Dataset

Dataset yang digunakan dalam proyek ini adalah `recruitment_data.csv`, berisi informasi kandidat dan faktor yang dipertimbangkan dalam proses perekrutan. Tujuan pemodelan adalah memprediksi keputusan perekrutan (*HiringDecision*) berdasarkan atribut kandidat.

3.4.1 Ringkasan Dataset

- **Jumlah rekaman (baris):** 1,500
- **Jumlah fitur (prediktor):** 10
- **Target:** *HiringDecision* (biner: 0 = tidak diterima, 1 = diterima)
- **Sifat data:** Sintetis (dibuat untuk tujuan pendidikan/proyek data sains)

3.4.2 Definisi Variabel

Berikut fitur dan target yang tersedia, beserta tipe data, rentang/kategori, dan keterangan singkat. Tabel 3 merangkum definisi variabel dalam dataset.

Table 3: Definisi Variabel Dataset

Nama Fitur	Tipe	Rentang	Keterangan
Age	Numerik	20-50	Umur kandidat
Gender	Kategorikal	0/1	0 = Laki-laki, 1 = Perempuan
EducationLevel	Kategorikal	1/2/3/4	1 = S1 (Tipe 1), 2 = S1 (Tipe 2), 3 = S2, 4 = S3/PhD
ExperienceYears	Numerik	0-15	Lama pengalaman kerja (tahun)
PreviousCompanies	Numerik	1-5	Jumlah perusahaan tempat bekerja sebelumnya
DistanceFromCompany	Numerik	1-50	Jarak dari rumah ke perusahaan
InterviewScore	Numerik	0-100	Skor hasil wawancara
SkillScore	Numerik	0-100	Skor keterampilan teknis
PersonalityScore	Numerik	0-100	Skor aspek kepribadian
RecruitmentStrategy	Kategorikal	1/2/3	1 = Agresif, 2 = Moderat, 3 = Konservatif
HiringDecision	Target	0/1	Target: 0 = tidak diterima, 1 = diterima

3.4.3 Catatan Kodefikasi dan Pra-pemrosesan

- **Gender:** dikodekan sebagai 0 (Laki-laki) dan 1 (Perempuan).
- **EducationLevel:** ordinal 1–4 dengan pemetaan spesifik (S1 Tipe 1, S1 Tipe 2, S2, S3/PhD). Jika korelasi kuat, dapat di *one-hot*.

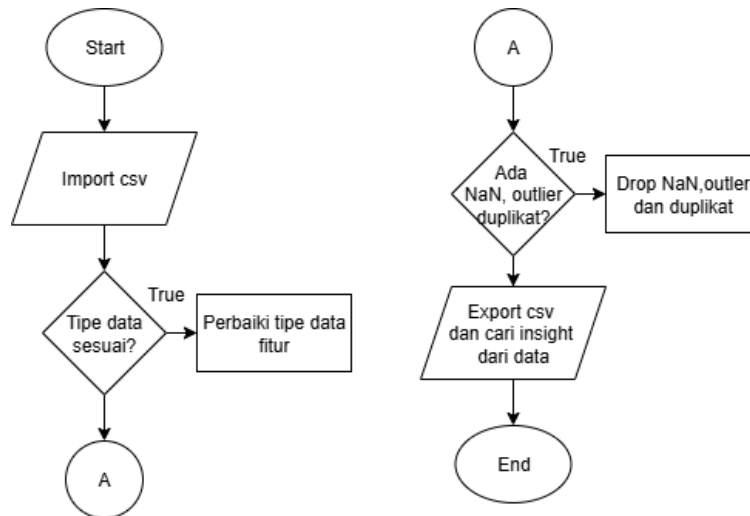
- **RecruitmentStrategy**: kategorikal 1–3. Umumnya di-*one-hot* untuk model linear; model pohon dapat menggunakan kode numeriknya langsung.
- **Skor (Interview/Skill/Personality)**: berada pada skala 0–100; pertimbangkan penskalaan (*standardization/min-max*) untuk model sensitif skala.
- **Fitur numerik lain** (Age, ExperienceYears, PreviousCompanies, DistanceFromCompany): periksa outlier, distribusi, dan lakukan transformasi/penanganan jika diperlukan.

3.4.4 Sumber dan Lisensi

Dataset ini dibagikan oleh **Rabie El Kharoua** dengan lisensi **CC BY 4.0**. Dataset bersifat *exclusive synthetic* dan ditujukan untuk keperluan edukasi/proyek data sains. Penggunaan diperbolehkan dengan mencantumkan atribusi yang tepat kepada pemilik dataset. DOI dan rincian penyedia data tercantum pada kartu data sumbernya. (Kharoua, 2024)

3.5 EDA (Exploratory Data Analysis)

EDA adalah langkah awal yang penting dalam analisis data untuk memahami struktur, pola, dan karakteristik dataset. Gambar 8 menunjukkan flowchart EDA yang akan dilakukan dalam proyek ini.



Gambar 8: Flowchart EDA

Dengan mengikuti langkah EDA yang terstruktur, tim dapat memperoleh wawasan yang mendalam tentang dataset, mengidentifikasi potensi masalah, dan menyiapkan data dengan baik untuk tahap pemodelan selanjutnya. EDA membantu memastikan bahwa model yang dibangun didasarkan pada pemahaman yang kuat tentang data, sehingga meningkatkan peluang keberhasilan proyek secara keseluruhan.

3.5.1 Handle Tipe Data, NaN, & Duplikasi

Data wrangling adalah proses penting dalam persiapan data untuk analisis dan pemodelan. Proses ini melibatkan beberapa langkah kunci yang bertujuan untuk membersihkan, mengubah, dan mengorganisir data agar siap digunakan. Dengan menggunakan `df.info()`, kita dapat memperoleh gambaran umum tentang struktur dataset, termasuk jumlah entri, tipe data setiap kolom, dan informasi tentang nilai yang hilang. Berikut adalah hasil dari `df.info()` pada dataset yang digunakan dalam proyek ini.

Listing 1: Info Dataset

```
1 <class 'pandas.core.frame.DataFrame'>
2 RangeIndex: 1500 entries, 0 to 1499
3 Data columns (total 11 columns):
4 #   Column                               Non-Null Count  Dtype
5 ---  -
6 0   Age                                  1500 non-null   int64
7 1   Gender                              1500 non-null   int64
8 2   EducationLevel                      1500 non-null   int64
9 3   ExperienceYears                     1500 non-null   int64
10 4   PreviousCompanies                  1500 non-null   int64
11 5   DistanceFromCompany               1500 non-null   float64
12 6   InterviewScore                    1500 non-null   int64
13 7   SkillScore                        1500 non-null   int64
14 8   PersonalityScore                  1500 non-null   int64
15 9   RecruitmentStrategy               1500 non-null   int64
16 10  HiringDecision                    1500 non-null   int64
17 dtypes: float64(1), int64(10)
18 memory usage: 129.0 KB
```

Dari hasil `df.info()`, kita dapat melihat bahwa dataset terdiri dari 1500 entri dengan 11 kolom. Semua kolom memiliki tipe data numerik (`int64` dan `float64`), dan tidak ada nilai yang hilang (non-null count sama dengan total entries untuk setiap kolom). Ini menunjukkan bahwa dataset sudah cukup bersih dari segi kelengkapan data, namun masih perlu dilakukan pemeriksaan lebih lanjut terhadap distribusi nilai, outlier, dan potensi inkonsistensi lainnya. Walaupun beberapa fitur numerik memiliki makna kategorikal seperti gender, education level, dan recruitment strategy, hal tersebut tidaklah menjadi masalah karena jika dia bertipe object pada akhirnya akan diubah menjadi numerik juga.

Selanjutnya mengecek apakah ada nilai duplikasi pada dataset. Dengan menggunakan `df.duplicated().sum()`, kita dapat menghitung jumlah baris yang duplikat dalam dataset. Berikut adalah hasil dari pengecekan duplikasi pada dataset yang digunakan dalam proyek ini.

Listing 2: Cek Duplikasi Dataset

```
1 df.duplicated().sum()
2
3 #output
4 np.int64(0)
```

Dari hasil pengecekan duplikasi, kita dapat melihat bahwa tidak ada baris yang duplikat dalam dataset (jumlah duplikasi adalah 0). Ini menunjukkan bahwa setiap entri dalam dataset adalah unik, yang merupakan kondisi ideal untuk analisis data dan pemodelan. Dengan tidak adanya duplikasi, kita dapat melanjutkan ke tahap berikutnya dalam proses data wrangling dengan keyakinan bahwa data yang kita miliki sudah bersih dari masalah duplikasi.

3.5.2 Analisis Fitur Numerik

Fitur numerik dalam dataset ini meliputi `Age`, `ExperienceYears`, `PreviousCompanies`, `DistanceFromCompany`, `InterviewScore`, `SkillScore`, dan `PersonalityScore`. Untuk memahami karakteristik dari fitur-fitur ini, kita dapat melakukan analisis statistik deskriptif dan visualisasi distribusi data. Tabel 4 merangkum statistik deskriptif dari fitur numerik dalam dataset.

Table 4: Statistik Deskriptif Fitur Numerik

Fitur	Count	Mean	Std	Min	25%	50%	75%	Max
Age	1500	35.15	9.25	20.00	27.00	35.00	43.00	50.00
ExperienceYears	1500	7.69	4.64	0.00	4.00	8.00	12.00	15.00
PreviousCompanies	1500	3.00	1.41	1.00	2.00	3.00	4.00	5.00
DistanceFromCompany	1500	25.51	14.57	1.00	12.84	25.50	37.74	50.99
InterviewScore	1500	50.56	28.63	0.00	25.00	52.00	75.00	100.00
SkillScore	1500	51.12	29.35	0.00	25.75	53.00	76.00	100.00
PersonalityScore	1500	49.39	29.35	0.00	23.00	49.00	76.00	100.00

Daftar Pustaka

- Agresti, A. (2018). *Statistical methods for the social sciences* (5th ed.). Pearson.
- Breaugh, J. A. (2013). Employee recruitment. *Annual Review of Psychology*, 64, 389–416. <https://doi.org/10.1146/annurev-psych-113011-143757>
- Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5–32. <https://doi.org/10.1023/A:1010933404324>
- Chumbar, S. (2020). *The crisp-dm process: A comprehensive guide* [Accessed: 2025-08-26]. <https://medium.com/@shawn.chumbar/the-crisp-dm-process-a-comprehensive-guide-4d893aecb151>
- De Veaux, R. D., Velleman, P. F., & Bock, D. E. (2011). *Intro stats* (3rd ed.). Pearson Education.
- Hausknecht, J. P., Rodda, J., & Howard, M. J. (2009). Targeted employee retention: Performancebased and jobrelated differences in reported reasons for staying. *Human Resource Management*, 48(2), 269–288. <https://doi.org/10.1002/hrm.20279>
- Hosmer, D. W., Lemeshow, S., & Sturdivant, R. X. (2013). *Applied logistic regression* (3rd ed.). Wiley. <https://doi.org/10.1002/9781118548387>
- Jain, A. K., & Bhandare, S. (2016). Min max normalization based data perturbation method for privacy protection. *International Journal of Computer Applications*, 136(5), 1–4. <https://doi.org/10.5120/ijca2016908569>
- Kharoua, R. E. (2024). Predicting hiring decisions in recruitment data. <https://doi.org/10.34740/KAGGLE/DSV/8715385>
- Mathis, R., Jackson, J., Valentine, S., & Meglich, P. (2017). *Human resource management*. Cengage Learning.
- Ng, E. S. W., & Burke, R. J. (2005). Person–organization fit and the war for talent: Does diversity management make a difference? *The International Journal of Human Resource Management*, 16(7), 1195–1210. <https://doi.org/10.1080/09585190500144038>
- Ployhart, R. E. (2006). Staffing in the 21st century: New challenges and strategic opportunities. *Journal of Management*, 32(6), 868–897. <https://doi.org/10.1177/0149206306293625>
- Schmidt, F., & Hunter, J. (1998). The validity and utility of selection methods in personnel psychology. *Psychological Bulletin*, 124(2), 262–274.
- Tukey, J. W. (1977). *Exploratory data analysis*. Addison–Wesley.