

Alur Pre-processing Data AIS

Langkah	Nama Tahap	Tujuan Singkat	Kriteria Lulus
1	Validasi Skema & Kelengkapan	Pastikan kolom inti tersedia (<i>mmsi</i> , <i>dt_pos_utc</i> , <i>lat</i> , <i>lon</i> , <i>sog</i>) & bertipe benar	Kolom hilang → file ditolak tipe salah → cas jika gagal → drop baris
			Valid MMSI = MMSI < 100000000 or MMSI > 999999999 Valid IMO = 1000000 - 99999999 Valid lat = rentang -90° hingga +90° adalah rentang penuh dari kutub selatan ke kutub utara Valid lon = rentang -180° sampai 180° dan mengukur posisi timur-barat terhadap meridian utama (Greenwich)
2	Deduplikasi Prime Key	Hilangkan rekaman ganda	Definisi PK = (<i>mmsi</i> , <i>dt_pos_utc</i> , <i>message_type</i>); duplikat → simpan satu, log sisanya
3	Filter Domain (Rules Hard)	Bersihkan <i>outlier</i> eksplisit	<ul style="list-style-type: none"> MMSI 9-digit valid <i>draught</i> ≠ 0 Tipe kapal

			relevan (cargo, tanker) • Periode analisis (mis. 2022)
4	Pengujian Range Numerik	Jaga batas fisik & format	$lat \pm 90^\circ, lon \pm 180^\circ; 0 \leq sog \leq 70 \text{ kn}; \text{timestamp} \leq _UTC \text{ now} + 5 \text{ m}$
5	Konsistensi Metadata per MMSI	Satu identitas kapal = satu set metadata	Hitung <code>nunique()</code> <code>imo</code> , <code>vessel_name</code> , <code>length</code> , <code>width</code> ; $>1 \rightarrow$ flag audit
6	Koherensi Spasio-Temporal	Deteksi loncatan & jeda sinyal	• <i>Jump test</i> $> 30 \text{ km/min} \rightarrow$ drop titik • $\Delta t > 60 \text{ m} \rightarrow$ flag gap • <code>nav_status</code> \leftrightarrow <code>sog</code> selaras
7	Enrichment & Spatial Join	Tambahkan konteks	• Join data statis kapal (IMO, tipe) • Spatial join dengan poligon pelabuhan \rightarrow <code>inside_port</code>

QA Quality Dimension	Indicator QA																	
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
Relevance			v			v			v								v	
Accuracy	v	v	v		v		v	v		v	v	v	v	v	v	v		v
Punctuality and Timeliness												v						
Interpretability							v											
Accessibility																		
Coherence		v		v	v	v		v						v		v		v
Comparability																		
Trustworthiness	v									v								

Kerangka Data-Input Quality Metrics (AIS):

	Indikator	Dimensi Kualitas	Goals (penjelasan singkat tapi jelas)	Rumus Indikator*	Urgensi
1	Percent of Missing Values	Accuracy, Trustworthiness	Memastikan kolom-kolom kunci (mmsi, waktu, koordinat, kecepatan) selalu terisi, supaya lintasan kapal tidak “bolong-bolong” dan analisis port-call tidak salah sasaran.	$\text{null_pct}(\text{col}) = (\sum \text{row_null}(\text{col}) / N_{\text{total}}) \times 100 \%$	Critical
2	Invalid MMSI Rate	Accuracy, Coherence	Menyingkirkan pesan AIS yang memakai MMSI di luar 9 digit resmi—karena rekaman palsu dapat “mencemari” statistik per kapal dan pelabuhan.	$\text{invalid_mmsi_pct} = (\sum \text{MMSI} < 100\,000\,000 \vee > 999\,999\,999) / N_{\text{total}} \times 100 \%$	Critical

3	Non-Zero Draught Rate	Relevance, Accuracy	Sarat air 0 artinya nakhoda tidak mengisi / sensor mati. Kita butuh nilai >0 untuk menganalisis bobot muatan dan kedalaman pelabuhan.	$\text{draught_zero_pct} = \frac{(\sum \text{draught} = 0)}{N_total} \times 100 \%$	Important
4	Records & Unique Vessels	Coverage, Coherence	Mengecek apakah jumlah total ping dan kapal harian/-mingguan berada di kisaran normal; penurunan tajam bisa menandakan gangguan satelit atau ingest.	$\text{rows_cnt, mmsi_n} = \text{COUNT(DISTINCT mmsi)}$ dibandingkan baseline rata-rata	Critical
5	Avg. Records per Vessel	Coherence, Accuracy	Mengukur kepadatan ping tiap kapal; kalau terlalu sedikit, kita bisa kehilangan detil manuver di pelabuhan.	$\text{avg_rows_per_mmsi} = \frac{\text{rows_cnt}}{\text{mmsi_n}}$	Important
6	Day-of-Week Pattern	Coherence, Relevance	Melihat distribusi lalu-lintas Senin → Minggu; pola "kosong di Minggu" misalnya bisa ikut memengaruhi jam kerja pelabuhan.	Hitung rows_cnt & mmsi_n per hari (Mon...Sun); cek hari dengan 0 rekor	Nice

7	Schema & Range Pass-Rate	Accuracy, Interpretability	Memastikan semua kolom berada di rentang fisik (lat $\pm 90^\circ$, lon $\pm 180^\circ$, sog 0-70 kn, dst); baris di luar itu dianggap cacat.	$\text{pass_rate} = \frac{\text{rows_good}}{\text{N_total}} \times 100 \%$	Critical
8	Unique H3 Cells (res 7)	Coverage, Accuracy	Menguji apakah cakupan spasial data sudah menutupi area laut yang semestinya; sel res 7 $\approx 1 \text{ km}^2$.	$\text{h3_unique_cnt} = \text{COUNT}(\text{DISTINCT h3_r7})$	Important
9	Ports Covered per Day	Coverage, Relevance	Menjamin data mencakup cukup banyak pelabuhan; penurunan drastis bisa berarti layer poligon salah atau data hilang.	$\text{port_n_daily} = \text{COUNT}(\text{DISTINCT port_id})$	Important
10	Incorrect Coordinates	Accuracy, Trustworthiness	Menolak titik AIS yang "mendarat" di daratan atau di luar peta laut; ini biasanya noise GPS atau spoofing.	$\text{coord_on_land_pct} = \frac{(\sum \text{point_on_land})}{\text{N_total}} \times 100 \%$	Critical

Commented [1]: Penghitungannya mungkin dibagi per orang per port nanti kalau pakai ini takutnya gak dapet ? gimana han ?

11	Spatial Gap (Jump Test)	Accuracy	Mengidentifikasi “teleport” > 30 km per menit yang mustahil secara fisik sehingga wajib dibuang.	<code>jump_rate = distance_km/Δt_min;</code> flag jika > 30	Critical
12	Time Gap between Pings	Timeliness, Accuracy	Mencari jeda sinyal > 60 menit; kalau banyak, artinya ada blank spot satelit atau AIS dimatikan.	Distribusi <code>Δt = t_i - t_{i-1}</code> ; hitung % gap > 60 m	Important
13	Speed Gap Consistency	Accuracy	Memastikan kecepatan terhitung (jarak/Δt) sejalan dengan <code>sog</code> ; gap besar ≈ sensor salah.	<code>speed_calc_kmh = (distance_km / Δt_h)</code> , lalu	calc - sog
14	Metadata Consistency	Coherence, Accuracy	<code>imo</code> , nama, ukuran kapal seharusnya tidak berubah-ubah dalam satu MMSI; perubahan = sinyal kesalahan merge.	<code>meta_conflict = (nunique(col) > 1)</code> per mmsi	Important
15	Ingest Latency p95	Timeliness	Menilai apakah data masuk hampir real-time; latency tinggi mengurangi	<code>latency_p95 = p95(dt_insert - dt_pos)</code>	Nice

manfaat
operasional.

16	Duplicate Primary- Key Rate	Accuracy, Coherence	(mmsi, timestamp) wajib unik; duplikat artinya data direplay atau file dobel.	$\text{dup_pk_pct} = (\sum \text{dup_pk}) / N_total \times 100 \%$	Critical
17	Inside- Port Share	Relevance	Memilih kapal yang betul-betul pernah berada di poligon pelabuhan— hanya mereka yang relevan untuk statistik port-call.	$\text{inside_port_share_m} = \text{rows_inside} / \text{rows_total}$	Important
18	Nav- Status vs Speed	Coherence, Accuracy	Mengecek kecocokan status 'Moored/Anchored' dengan kecepatan < 0,5 kn; inkonsistensi = potensi kesalahan sensor.	$\text{incoherent_nav_pct} = (\sum \text{rows_nav_mismatch}) / N_total \times 100 \%$	Nice