**HandlingMissingDatain ETL**

**1. WhatarethemostcommonreasonsformissingdatainETLpipelines?**
Ans.
• Sourcesystemissues→Datanotcapturedduetofaultysensors,humanerror,or incomplete forms.
• Schemachangesupstream→Columnnamesorstructureschange,breaking extraction.
• Integrationerrors→APIfailures,connectiontimeouts,ormismatchedformatsduring extraction.
• Datacorruption→Lossduringtransferbecauseofnetworklatencyorhardwarefailures.
• Workflowbottlenecks→Jobsfailmid-processduetoresourcelimitsorpoor orchestration.
• Businessrules→Certainfieldsintentionallyleftblank(e.g.,optionalattributes).

**2. WhyisblindlydeletingrowswithmissingvaluesconsideredabadpracticeinETL?**
Ans.
• Lossofvaluableinformation→Entirerecordsarediscardedevenifmostfieldsare valid.
• Biasinanalysis→Ifmissingnessissystematic(e.g.,certainregionsorcustomer groups), deletion skews results.
• Reducedsamplesize→Smallerdatasetsweakenstatisticalpowerandmodel accuracy.
• Businessimpact→Importantinsights(likewhydataismissing)arelost,whichcould itself be meaningful.

**3. Q3.Explainthedifferencebetween:**
• **Listwisedeletion**
• **Columndeletion**
**Alsomentiononescenariowhereeachisappropriate. Ans.**
ListwiseDeletion
• Removesentirerowsifanyrequiredfieldismissing.
• Preservesschemabutreducesdatasetsize.
• Scenario:Customersurveydatawheremultipleanswersaremissing→bettertodrop the whole record to avoid incomplete profiles.
ColumnDeletion
• Removesentirecolumnsiftheyhavetoomanymissingvalues.
• Preservesrowcountbutlosesthatvariable.

• Scenario: A dataset where "MiddleName" is missing for 95% of customers → dropping the column avoids clutter without harming analysis.

**4. Why is median imputation preferred over mean imputation for skewed data such as income?**

Ans.
• Mean imputation → Sensitive to extreme values (outliers). In skewed distributions like income, a few very high salaries can inflate the mean, making it unrepresentative.
• Median imputation → More robust because the median is the middle value, unaffected by outliers. It better reflects the "typical" case in skewed datasets.
• Example: If most incomes are around 50,000 but a few are above 1,000,000, the mean will be distorted, while the median stays close to the majority.

**5. What is forward fill and in what type of dataset is it most useful?**
Ans.
• Forward Fill (FFILL):
A technique where missing values are replaced with the last valid observation carried forward.
• Usefulness:
o Best suited for time-series or sequential datasets (e.g., monthly sales, stock prices, sensor readings).
o Assumes continuity—the last known value is a reasonable proxy until a new one appears.
• Example:
If sales data for March is missing but February had 12,000, forward fill assigns March = 12,000 until April's actual figure arrives.

**6. Why should flagging missing values be done before imputation in an ETL workflow?**
**Ans.**
• Preserves information about missingness → The fact that a value is missing can itself carry business meaning (e.g., customers not disclosing income).
• Avoids loss of transparency → Once imputation is applied, you can't distinguish between original and filled values unless flagged beforehand.
• Supports better analysis → Analysts can study patterns of missingness separately (e.g., which regions or groups have more missing data).
• Improves model accuracy → Machine learning models can use the flag as an additional feature, helping them account for missingness bias.

**7. Consider a scenario where income is missing for many customers. How can this missingness itself provide business insights?**
Ans.
• Customer reluctance → Missing income may indicate customers are unwilling to disclose sensitive financial details, signaling trust or privacy concerns.
• Segment behavior → Certain demographics or regions may have higher missingness, revealing cultural or socio-economic differences.
• Product targeting → If high-value customers avoid sharing income, it may suggest they don't see relevance in providing it, guiding product design or survey strategy.
• Operational gaps → Consistent missingness could highlight flaws in data collection