

Roteiro geral de analises, modelagem e dashboards

Documento complementar aos guias existentes em `docs/`. Objetivo: orientar colegas sobre a sequencia minima para atualizar dados, executar as analises/modelos e abastecer dashboards sem se perder na estrutura Bronze/Silver/Gold.

1. Visao rapida do projeto

- Escopo: seis municipios da RMB (Belem, Ananindeua, Marituba, Benevides, Santa Barbara do Para, Santa Izabel do Para) no periodo 2018-2025.
- Hipotese central: melhor serviço (cobertura, eficiencia), boa qualidate da agua e finanças equilibradas reduzem internacoes por agravos hidricos, moduladas por clima.
- Entregaveis: base Gold anual (`data/gold/gold_features_ano.*`), notebooks de EDA/modelagem (pasta `notebooks/`), dashboard no Looker Studio com artefatos exportados em `dashboard/`.

2. Estrutura de dados atualizada

Camada	Uso	Diretorios principais	Formato	Observacoes
Bronze	Dados brutos como baixados	<code>data/bronze/sih</code> , <code>data/bronze/sih/csv</code> , <code>data/bronze/sisagua</code> , <code>data/bronze/siops</code> , <code>data/bronze/inmet</code> , <code>data/bronze/ibge</code> , <code>data/bronze/snus</code>	DBC, CSV, ZIP	Mantener apenas anos/fonte relevantes.
Silver	Dados padronizados (chaves, tipos)	<code>data/silver/sih</code> , <code>data/silver/sisagua</code> , <code>data/silver/siops</code> , <code>data/silver/inmet</code> , <code>data/silver/ibge_populacao</code>	Parquet	Colunas chaves: <code>cod_mun</code> , <code>municipio</code> , <code>ano</code> .
Gold	Conjuntos prontos para analise/dash	<code>data/gold/gold_features_ano.*</code> , <code>data/gold/snus_rmb_indicadores_v2.*</code> , <code>data/gold/gold_qualidade_agua.*</code>	CSV, Parquet	Fonte oficial para notebooks, modelos e dashboard.
Dashboards	Entregaveis finais e artefatos de apresentação	<code>dashboard/</code>	PDFs, imagens, links	Guardar exports do Looker/PowerPoint para referência do grupo.

Referencia cruzada: README.md (secao 4) traz detalhes extras e deve ser atualizado se novos arquivos surgirem.

3. Passo a passo essencial

3.1 Preparacao do ambiente

1. Ativar a virtualenv (`source .venv/bin/activate`).
2. Instalar dependencias se necessário: `pip install -r requirements.txt`.
3. Garantir que PyArrow e Pandas estao presentes (requeridos para Parquet).

3.2 Confirmar camada Bronze

- `data/bronze/sih`: arquivos RDPAyymm.dbc por mes (2018-2025).
- `data/bronze/sih/csv`: CSVs convertidos a partir dos DBCs (um por mes) quando precisar recalcular Silver.
- `data/bronze/sisagua`: pacotes .csv.zip dos controles mensais 2018-2025.
- `data/bronze/inmet`: pastas por ano com CSVs INMET_N_PA_* (estacoes A201, A202, A227).
- `data/bronze/siops`: arquivos consolidados `siops_indicadores_rmb_*.csv`.
- `data/bronze/ibge`: `sidra6579_pop_aaaa.csv`.
- `data/bronze/snus`: planilhas anuais 2018-2023.

3.3 Gerar camada Silver (quando houver atualizacao)

Fonte	Script	Comando base	Saida
SNIS (servico)	<code>scripts/etl_snus_indicadores_rmb.py</code>	<code>python scripts/etl_snus_indicadores_rmb.py</code>	CSV/Parquet direto em Gold (ja aplicado).
SISAGUA	<code>scripts/bronze_to_silver_sisagua_parquet.py</code>	<code>python scripts/bronze_to_silver_sisagua_parquet.py --input-dir data/bronze/sisagua</code>	Parquet particionado em <code>data/silver/sisagua</code> .
SIH	<code>scripts/bronze_to_silver_sih_parquet.py</code>	<code>python scripts/bronze_to_silver_sih_parquet.py --csv-dir data/bronze/sih/csv</code>	Parquet particionado em <code>data/silver/sih</code> (ano/mes).
SIOPS	<code>scripts/bronze_to_silver_siops_parquet.py</code>	<code>python scripts/bronze_to_silver_siops_parquet.py --input data/bronze/siops/siops_indicadores_rmb_2018_2025.csv</code>	<code>data/silver/siops/indicadores.parquet</code> .
INMET	<code>scripts/inmet_to_parquet.py</code>	<code>python scripts/inmet_to_parquet.py --input-dir data/bronze/inmet</code>	Dataset particionado em <code>data/silver/inmet</code> .
IBGE (pop)	<code>scripts/bronze_to_silver_ibge_pop_parquet.py</code>	<code>python scripts/bronze_to_silver_ibge_pop_parquet.py --input-dir data/bronze/ibge</code>	<code>data/silver/ibge_populacao/populacao.parquet</code> .

3.4 Montar camada Gold

1. SNIS v2: `scripts/fix_snus_csv.py` (gera CSV/Parquet corrigidos em `data/gold/`).
2. Qualidate da agua: usar notebook ou script auxiliar (ver `docs/Guia de Preparacao de Dados` – ...) para agregacoes e salvar `data/gold/gold_qualidade_agua.*`.
3. Merge Silver -> Gold final: `python scripts/silver_to_gold_features.py` (usa `config/rmb_municipios.csv` e salva `data/gold/gold_features_ano.*`).
4. Se precisar revisitar versao anterior do ETL SNIS, `scripts/fix_snus_scale_and_ibge.py` continua disponivel.

4. Orientacoes de uso dos dados

- **Chaves padrao:** sempre trabalhar com `cod_mun` (7 dígitos), `municipio` (uppercase) e `ano`. Conferir se há 6 municípios x n anos.
- **Normalização populacional:** utilizar `data/silver/ibge_populacao/populacao.parquet` para gerar indicadores per capita ou por 10k hab.
- **Qualidade da água:** `gold_qualidade_agua` traz tanto contagem de amostras quanto percentuais por parâmetro (`pct_conformes_*`). Para dashboards, priorizar percentuais; para modelos, usar contagens como peso.
- **SNIS v2:** já corrigido (percentuais >100 convertidos) e possui colunas `idx_*` padronizadas. Usar essa versão, não a raw.
- **SIOPS:** indicadores em percentual estão em escala 0-100; despesas per capita em R\$ constantes (ano corrente). Validar se há anos faltantes ao cruzar com Gold.
- **Clima:** dados agregados por estação são replicados aos municípios conforme mapeamento (A201 -> Belém + entorno, A202 -> Castanhal + Sta Izabel + Sta Barbara, A227 -> proxy Barcarena). Ajustar se novos municípios forem incluídos.
- **Saúde (SIH):** `internacoes_total` e `internacoes_hidricas` devem ser agregados antes de calcular taxas (`internacoes_hidricas_10k`). Sempre verificar se CID filtrado continua coerente.

5. Roteiro de análises e modelagem

1. **EDA serviço e qualidade:** notebook `01_eda_servicos.ipynb` utilizando `data/gold/snus_rmb_indicadores_v2.*` e `data/gold/gold_qualidade_agua.*`.
2. **EDA saúde e clima:** notebook `02_eda_saude_clima.ipynb` combinando `data/gold/gold_features_ano.*` (foco em clima + internações).
3. **Modelagem (target `internacoes_10k_hab`):** notebook `03_modelagem.ipynb` carrega `gold_features_ano`.
 - Modelos: Regressão Linear, Lasso, RandomForestRegressor.
 - Métricas: MAE, RMSE, R2 (k-fold). Guardar resultados em `reports/` se atualizar.
 - Explicabilidade: permutation importance + SHAP (opcional).
4. **Classificação opcional:** binarizar não conformidade (p.ex. `pct_conformes_global < 95`) e testar Logistic vs RandomForestClassifier.

5.1 Cobertura das perguntas de negócios

Pergunta (doc MVP)	Entregável/Dado que responde	Ações previstas
1. Quais ações reduzem mais as internações?	Notebook <code>03_modelagem.ipynb</code> (feature importance e análise de elasticidade) + aba "Priorizar & Simular" do dashboard	Rodar modelos explicáveis, calcular efeitos marginais (<code>what-if</code>) e documentar insights na narrativa do dashboard.
2. Qual retorno por município/sistema?	<code>gold_features_ano</code> + cálculo "ganho por ponto" nas notebooks e tabelas do dashboard	Criar tabela ranqueando municípios pelo impacto projetado; publicar no dashboard e no relatório final.
3. Como a chuva modula o risco?	Notebook <code>02_eda_saude_clima.ipynb</code> (series A00-A09 x chuva com defasagens) + cards no dashboard	Incluir análises de correlação/lag e visualizações sazonais (linha + heatmap).
4. Qual ranking de priorização (12 meses)?	<code>03_modelagem.ipynb</code> + dashboard (painel Priorizar)	Gerar score composto e ordenação; exportar versão para dashboard/ .
5. Como medir ODS 6/3/11?	ODS Tracker no dashboard + indicadores em <code>gold_features_ano</code>	Garantir seção dedicada no dashboard com metas e status; anexar ao relatório.

6. Dashboards (Looker Studio)

- **Fonte principal:** `data/gold/gold_features_ano.csv`.
- **Complementos:** `data/gold/gold_qualidade_agua.csv` para visões mensais, `data/gold/snus_rmb_indicadores_v2.csv` para detalhamento de indicadores.
- **Páginas sugeridas:**
 1. Visão RMB (KPI cards + ranking + séries temporais).
 2. Visão municipal (cartões por eixo, barras de não conformidade, explicabilidade do modelo).
 3. Painel de prioridades (tabela com ranking + simulador simples).
- **Bons práticas:** fixar filtros para anos e municípios; destacar meta (ex. `pct_conformes_global >= 95`).
- **Entrega:** salvar versões exportadas (PDF, PNG, links compartilháveis) dentro da pasta `dashboard/`.

6.1 Simulador de cenários e priorização (seção 12 do MVP)

- **Dados de entrada:** `gold_features_ano.csv` (indicadores anuais) + elasticidades estimadas em `03_modelagem.ipynb` para `idx_atend_agua_total`, `idx_tratamento_esgoto`, `idx_perdas_distribuicao`, `pct_conformes_global` e outras variáveis relevantes.
- **Mecanismo:** utilizar os coeficientes/regressões para calcular (Δ) previsto em `internacoes_10k_hab` quando o usuário ajusta sliders (% coleta, % tratado, perdas, conformidade).
- **Saídas:** painel com ranking atualizado (municípios/sistemas) e comparação com baseline, exportado para `dashboard/`.
- **Documentação:** registrar metodologia e supostos nos notebooks e replicar resumo na página "Priorizar & Simular" do dashboard.

7. Checklist rápido antes de compartilhar

- Conferir contagem de linhas das bases Gold (6 municípios x n anos vigentes).
- Validar escalas (percentuais 0-100, valores monetários em R\$).
- Atualizar `README.md` se novas colunas/arquivos surgirem.
- Executar notebooks limpando o kernel para garantir reprodutibilidade.
- Exportar dashboard atualizando a fonte de dados, se necessário, e arquivar em `dashboard/`.
- Validar que as perguntas de negócios (doc MVP, seção 7) estão respondidas em notebooks/dashboards e documentar conclusões.

8. Recursos adicionais

- `docs/Guia de Preparação de Dados -md`: detalhes de cada fonte e ETL.
- `docs/Catálogo De Dados - RMB ...`: glossário completo (versão DOC/PDF).
- `README.md`: visão geral do repositório e dicionário resumido.
- Contato: registrar dúvidas via issue no repositório ou canal Teams do grupo.

Mantemos este roteiro como documento vivo. Ao adicionar novas fontes, modelos ou visualizações, registrar a mudança aqui, no README e nos guias correspondentes para garantir alinhamento da equipe.