

# Teoria e Métodos no Uso de Registros Administrativos na Pesquisa Aplicada em Políticas Públicas

Ricardo Ceneviva

Métodos Computacionais

Julho 2025

- Definir o que são registros administrativos (RA) e contextualizar seu papel em políticas públicas.
- Apresentar os métodos centrais para tratamento, integração e proteção de RA.
- Ilustrar aplicações com dois estudos de caso paradigmáticos (Chetty *et al.*, 2014; Barreto *et al.*, 2021).
- Discutir governança, ética e sustentabilidade dos dados na administração pública.

- ① Ecossistema de RA
- ② Metodologias de Integração & Proteção
- ③ Estudos de Caso
  - Chetty *et al.*, 2014
  - Barreto *et al.*, 2021
- ④ Governança, Ética & Sustentabilidade
- ⑤ Agenda Estratégica para o Futuro

## Definição

Dados coletados rotineiramente por governos para finalidades operacionais (cadastro, fiscalização, prestação de serviços).

## Principais Tipos

- *Cadastrais*: identificação de pessoas, empresas, imóveis.
- *Transacionais*: eventos recorrentes (benefícios, internações).
- *Financeiros*: arrecadação, folha de pagamento.
- *Serviços*: registros de escola, saúde, segurança.

**Exemplos Brasil:** Cadastro Único, Censo Escolar, RAIS/CAGED, DATASUS.

- ① **Coleta:** digitados na ponta (escolas, postos de saúde, repartições).
- ② **Processamento:** validação e deduplicação em sistemas centrais.
- ③ **Armazenamento:** bases legadas ou *data lakes* na nuvem.
- ④ **Integração:** preparação e *linkage* para pesquisa ou avaliação.
- ⑤ **Arquivamento/Descarte:** períodos de retenção e destruição segura.

- Cobertura quase censitária → poder estatístico elevado.
- Longitudinalidade natural → acompanha indivíduos ao longo do tempo.
- Custo marginal de pesquisa quase nulo → aproveita dados já coletados.
- Precisão em variáveis sensíveis (ex.: renda tributária, óbitos).

- **Vieses de cobertura:** populações rurais e informais sub-representadas.
- **Qualidade heterogênea:** erros de digitação, códigos obsoletos.
- **Rigidez conceitual:** variáveis criadas para gestão, não para pesquisa.
- **Barreiras de acesso:** sigilo, fragmentação institucional.

- **Determinístico:** chave única (CPF, NIS) → simples, mas dependente de completude.
- **Probabilístico:** Fellegi–Sunter (CIDACS-RL) → flexível, requer parâmetros e validação.
- Indicadores de qualidade: taxa de falsos-positivos/negativos, *precision*, *recall*.

- **PPRL:** linkage com hashes/Bloom filters evitando troca de identificadores.
- **Privacidade Diferencial:** garante anonimato matemático (-DP) adicionando ruído.
- **Dados Sintéticos:** conjuntos artificiais que preservam estatísticas essenciais.
- **Aprendizado Federado:** modelos treinados sem mover microdados.

- Padronização de variáveis (códigos, datas, nomes).
- Imputação de valores faltantes (regressão, múltipla imputação).
- Metadados e *data profiling* para auditoria de consistência.
- Documentação reproduzível (R Markdown, Jupyter, data dictionaries).

## Objetivo

Medir mobilidade intergeracional de renda nos EUA e identificar determinantes locais.

**Dados** — Registros fiscais (IRS) + dados de censo; 40 milhões de crianças (1980-1985).

## Metodologia

- Linkage determinístico via Social Security Number.
- Regressão rank-rank (renda pais  $\times$  renda filhos aos 35 anos).

# Insights do Caso 1

- Mobilidade varia até 3× entre condados → fenômeno fortemente local.
- Fatores correlacionados: segregação racial, qualidade escolar, capital social.
- Implicação: políticas territoriais (vouchers de moradia, zonings) podem aumentar oportunidades.

## Coorte 100 Milhões de Brasileiros (CIDACS/Fiocruz)

Integra Cadastro Único (programas sociais) aos sistemas de saúde (SIM, SINASC, SINAN), cobrindo 131 milhões de indivíduos (2001-2018).

### Método

- Record linkage probabilístico de alta escala (CIDACS-RL).
- Avaliação de impacto usando modelagem de painel.

- Programa Bolsa Família associado à redução de mortalidade materna e de tuberculose.
- Demonstra viabilidade de RA em países de renda média para estudos populacionais.
- Geração de infraestrutura de dados permanentes para futuras avaliações.

- Escala nacional possível em contextos de alta (EUA) e média renda (Brasil).
- Identificadores fortes (SSN vs. NIS) e infraestrutura influenciam precisão do linkage.
- Necessário entender vieses de cobertura para interpretar evidências.

- ① Safe Projects — projetos de interesse público legítimo.
- ② Safe People — pesquisadores qualificados e credenciados.
- ③ Safe Data — dados anonimizados/tratados.
- ④ Safe Settings — ambientes controlados (laboratórios seguros, VDI).
- ⑤ Safe Outputs — checagem de resultados antes da divulgação.

**Exemplo BR:** Sala Segura do CIDACS atende às cinco dimensões.

- **Data Stewardship:** papéis claros de custodians e trustees para gerenciar acesso e uso.
- **Espaços de Dados da UE:** EHDS (saúde) → compartilhamento federado entre países.
- Estado como orquestrador: normas, infraestrutura e mecanismos de incentivo ao compartilhamento.

- Custos fixos: servidores, segurança, equipe técnica.
- Modelos de co-financiamento público-privado (data collaboratives).
- Atenção a *data colonialism*: garantir benefícios locais e soberania.

- Padronização nacional de ID único (interoperabilidade CPF-NIS-CNS).
- Portais de consulta segura e aprendizado federado para análises inter-agência.
- Capacitação contínua de servidores e pesquisadores em ciência de dados.
- Incentivos regulatórios (Evidence Act BR) para compartilhamento responsável.

- RA são pilar para políticas baseadas em evidências.
- Desafios: governança, qualidade, ética e financiamento.
- Perguntas, comentários e próximos passos.