

# Conferência Nacional dos Agentes Produtores e Usuários de Dados

## Modernização da produção estatística

**Carlos Torres Freire**

Diretor do Centro de Ciência de Dados para Estatísticas Públicas (CCDEP / SEADE)

03 de Dezembro de 2025

# AGENDA

## 1) Pontos de partida

- Era digital e a produção estatística
- Princípios de qualidade de processos e produtos estatísticos

## 2) Experimentos recentes do CCDEP / SEADE:

- Áudio, texto, imagem, localização e IA Generativa

## 3) Centro de Ciência de Dados para Estatísticas Públicas (CCCDEP)

# Era digital e desafios na produção estatística

## TRADICIONALMENTE

- Estatísticas via coleta primária de dados – censos, *surveys* e registros administrativos

## NA ERA DIGITAL

- Dados em **grande volume, alta frequência e diversidade**
- **Novos métodos e tecnologias** para coletar, armazenar, tratar e analisar dados para produção de estatísticas

## REGISTROS GERADOS DIGITALMENTE POR AÇÕES HUMANAS E INTERAÇÕES SOCIAIS E CAPTURADOS AUTOMATICAMENTE

### REGISTROS ADMINISTRATIVOS DIGITAIS

**Cadastros sociais:**  
educação, saúde e assistência social

**Informações fiscais:**  
pagamentos, notas de transações

**Documentos:**  
Declarações de nascimentos e óbitos

### SUBPRODUTOS DE AÇÕES EM EQUIPAMENTOS

**Uso de celular**  
Transações financeiras

**Bilhetagem em transporte público**  
Visitas a websites

Publicação de notícias e em redes sociais

### COLETA AUTOMATIZADA POR SENSORES

**Presença física:**  
tráfego de veículos e pessoas ou iluminação

**Remotos:** satélites e veículos aéreos

# Princípios de qualidade de processos e produtos estatísticos

## PRODUTOS ESTATÍSTICOS

- **Relevância:** importância para o usuário, demanda da sociedade
- **Precisão:** descrição correta do fenômeno, validade
- **Atualidade:** distância da data de referência
- **Coerência:** comparabilidade, convenções, padrão conceitual, universo
- **Acessibilidade:** disponibilidade e transparência
- **Interpretabilidade:** informações adicionais, como metadados, metodologias, classificações

## PROCESSOS ESTATÍSTICOS

- **Metodologia sólida:** conceitos, classificações, convenções nacionais e internacionais
- **Adequação** de instrumentos, técnicas e sistemas tecnológicos
- **Solicitação não excessiva** em coletas e uso de registros administrativos
- **Relação entre custo e eficiência**

# Monitor Covid-19

## Pesquisa diária sobre temas relacionados à pandemia

- Acompanhamento de sintomas; distanciamento social; atividade econômica e consumo; trabalho e renda; educação; vacina; expectativas da população.
- De 23/Março/2020 a Setembro/2021.
- **Coleta de dados por telefone**, utilizando Unidade de Resposta Audível (**URA**).
- **800 entrevistas diárias**, para constituir um painel móvel de 4.800 entrevistas
- **População do Estado de São Paulo**, com amostra representativa por local de residência, **sexo e idade**
- Desagregação por Estado, RMSP, município SP, demais municípios da RMSP e Interior; e **17 Departamentos Regionais de Saúde**.

## Experiência levou a pesquisas de percepção:

- **Cultura**
- **Atividade física**
- **Turismo**
- **Endividamento**
- **Segurança Alimentar**
- **Cuidados no Domicílio**
- **Hábitos de Consumo por App e Internet**

# Coleta de informações por robô via telefone e aprendizado de máquina para classificar textos

Análise de respostas abertas em formato de áudio a partir de pesquisas de percepção por telefone com a população

## TRANSCRIÇÃO DE ÁUDIO

Processo automatizado que converte mensagens de áudio em texto.

## PREPARAÇÃO DE MENSAGENS

Adequação dos textos substituindo palavras equivalentes e removendo termos comuns.

## CLASSIFICAÇÃO DE TEXTO

Uso de modelo não supervisionado para classificar mensagens em tópicos principais.

## VALIDAÇÃO DE TÓPICOS

Análise de textos outliers e definição de rótulos para cada tópico pela equipe do SEADE.

ATUALIDADE

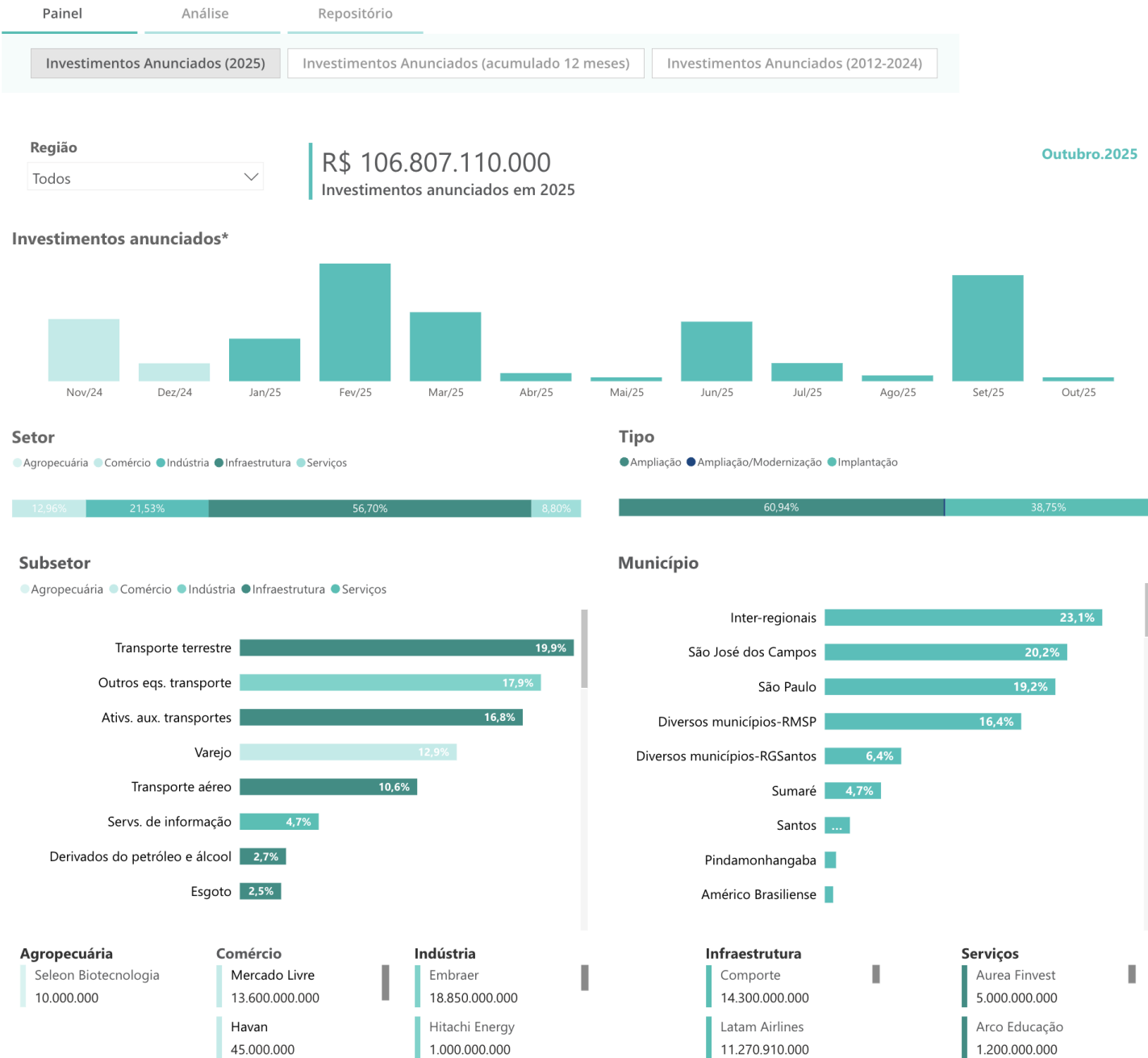
RELEVÂNCIA

CUSTO / BENEFÍCIO

COERÊNCIA e COMPARABILIDADE

**Exemplo internacional:** Statistics Canada utiliza modelos de processamento de linguagem natural para classificar respostas sobre ocupações

- Monitoramento sistemático e automatizado dos anúncios de investimentos de empresas privadas e públicas no Estado
- Captação diária de notícias de anúncios na imprensa
- Equipe de coleta confirma com empresas e completa informações como tipo, período, valor e local do investimento.



# Raspagem de dados para monitorar investimentos

*Web scraping* de anúncios de investimentos em veículos de imprensa, com sistema multiagentes, e classificação de texto para triagem de notícias, com aprendizado de máquina, e produção de indicadores

## CAPTAÇÃO DE NOTÍCIAS

Monitoramento automatizado de anúncios na imprensa sobre investimentos

## PROCESSAMENTO DE INDICADORES

Organização dos dados por tipo, período, valor e localização do investimento



## TRIAGEM POR IA

Modelos de aprendizado de máquina classificam textos sobre investimentos

## VALIDAÇÃO HUMANA

Equipe confirma informações com empresas e completa dados faltantes

**Exemplo internacional:** Projeto ESSnet Big Data do Eurostat classifica setorialmente e por ocupação anúncios de emprego em plataformas

ATUALIDADE

RELEVÂNCIA

PRECISÃO

COMPARABILIDADE



0005072

República Federativa do Brasil  
Ministério da Saúde  
1ª VIA - SECRETARIA DE SAÚDE

Declaração de Nascimento Vivo

30-77617573-6

Número do Cartão Nacional de Saúde do RN

700 5061 4523 0754

I	1 Nome do Recém-nascido (RN)		Número do Cartão Nacional de Saúde do RN 700 5061 4523 0754	
	2 Data e hora do nascimento 2.1 Data _____ Hora _____		3 Sexo M - Masculino <input type="checkbox"/> I - Ignorado <input type="checkbox"/> F - Feminino <input type="checkbox"/>	
II	4 Peso ao nascer _____		5 Índice de Apgar - 1º a 5º minutos 1º _____ 5º _____	
	6 Comprimento _____		7 Perímetro cefálico _____	
III	8 Local de ocorrência 8.1 Hospital <input type="checkbox"/> 8.2 Domicílio <input type="checkbox"/> 8.3 Outro estab. saúde <input type="checkbox"/> 8.4 Outro <input type="checkbox"/>		9 Estabelecimento _____ Código CHES _____	
	10 Endereço da ocorrência, se fora do estado, ou da residência da Mãe (rua, praça, avenida, etc)		11 CEP _____	
IV	12 Bairro/Distrito _____ Código _____		13 Município de ocorrência _____ Código _____	
	14 Nome da Mãe _____		15 Cartão SUS _____	
V	16 Escolaridade (última série concluída) 16.1 Sem escolaridade <input type="checkbox"/> 16.2 Média (certif. 2º grau) <input type="checkbox"/> 16.3 Fundamental I (1ª a 4ª série) <input type="checkbox"/> 16.4 Superior incompleta <input type="checkbox"/> 16.5 Fundamental II (5ª a 8ª série) <input type="checkbox"/> 16.6 Superior completa <input type="checkbox"/>		17 Ocupação habitual (informar atividade ou ocupação/desemprego) _____ Código CBO 2502 _____	
	18 Data nascimento da Mãe _____		19 Situação conjugal 19.1 Solteiro <input type="checkbox"/> 19.2 Casado <input type="checkbox"/> 19.3 Viúvo <input type="checkbox"/> 19.4 Separado judicialmente <input type="checkbox"/> 19.5 União estável <input type="checkbox"/> 19.6 Ignorado <input type="checkbox"/>	
VI	20 Idade (anos) _____		21 Raça / Cor da Mãe 21.1 Branca <input type="checkbox"/> 21.2 Preta <input type="checkbox"/> 21.3 Amarela <input type="checkbox"/> 21.4 Indígena <input type="checkbox"/>	
	22 Maturidade da Mãe _____		23 CEP _____	
VII	24 Nome do Pai _____		25 Idade do Pai _____	
	26 Bairro/Distrito _____ Código _____		27 Município _____ Código _____	
VIII	28 Gestações anteriores 28.1 Histórico gestacional _____		29 Parto _____	
	30 Nº de gestações anteriores _____		31 Nº de partos vaginais _____	
IX	32 Nº de cesáreas _____		33 Nº de nascidos vivos _____	
	34 Nº de perdas fetais / abortos _____		35 Nº de gestações gemelares _____	
X	36 Mês de gestação em que iniciou o pré-natal _____		37 Tipo de gravidez 37.1 Única <input type="checkbox"/> 37.2 Gêmeos <input type="checkbox"/> 37.3 Triplê <input type="checkbox"/> 37.4 Quadrupel <input type="checkbox"/> 37.5 Quintupel <input type="checkbox"/> 37.6 Sextupel <input type="checkbox"/> 37.7 Septupel <input type="checkbox"/> 37.8 Octupel <input type="checkbox"/> 37.9 Outros <input type="checkbox"/>	
	38 Apresentação _____		39 Tipo de parto 39.1 Cefálico <input type="checkbox"/> 39.2 Bacia <input type="checkbox"/> 39.3 Transversal <input type="checkbox"/> 39.4 Outros <input type="checkbox"/>	
XI	40 Trabalho de parto _____		41 Tipo de parto 41.1 Vaginal <input type="checkbox"/> 41.2 Cesárea <input type="checkbox"/> 41.3 Outros <input type="checkbox"/>	
	42 Conduta adotada antes do trabalho de parto (verificar) _____		43 Nascimento assistido por _____ 43.1 Médico <input type="checkbox"/> 43.2 Enfermeiro <input type="checkbox"/> 43.3 Paralelo <input type="checkbox"/> 43.4 Outros <input type="checkbox"/>	
XII	44 Data da última menstruação (DUM) _____		45 Método utilizado para estimar _____	
	46 Nº de semanas de gestação, se DUM ignorado _____		47 Método utilizado para estimar _____	
XIII	48 Descrever todas as anomalias congênicas observadas _____		49 Data do preenchimento _____	
	50 Nome do responsável pelo preenchimento _____		51 Função _____	
XIV	52 Tipo de documento _____		53 Nº do documento _____	
	54 Órgão emissor _____		55 Data _____	
XV	56 Cartório _____ Código _____		57 Registro _____	
	58 Município _____		59 UF _____	

ATENÇÃO: ESTE DOCUMENTO NÃO SUBSTITUI A CERTIDÃO DE NASCIMENTO

O Registro de Nascimento é obrigatório por lei

Para registrar esta criança, o pai ou responsável deverá levar este documento ao cartório de registro civil.

Versão 01/14 - 2º impressão 02/2017

www.ign.com.br

# Dados de telefonia móvel para análise de mobilidade

Metodologia permite análises precisas, desagregadas e atualizadas sobre deslocamentos populacionais para o planejamento de transportes e infraestrutura, especialmente em regiões sem pesquisas de origem e destino

## MATRIZ DE VIAGENS

Por meio da identificação das **origens e destinos** de unidades de celular, é possível mapear **fluxos, duração, distância, número médio, concentrações de viagens** e dinâmicas dentro e entre municípios.

## AVALIAÇÃO DE INTERVENÇÕES

Comparação de matrizes de viagens antes e após intervenções, usando métodos causais para avaliar impacto de nova infraestrutura, equipamentos ou eventos na mobilidade urbana e regional.

## Exemplos internacionais:

**Statistics Netherlands (CBS):** utiliza dados de telefonia móvel para analisar padrões de mobilidade, estimar distâncias médias percorridas, tempos de deslocamento e fluxos populacionais.

**Instituto Nacional de Estadística (INE) da Espanha:** utiliza esses dados para estatísticas de TURISMO; mede presença e deslocamentos de visitantes de forma contínua e com desagregação territorial.

ATUALIDADE

RELEVÂNCIA

ACESSIBILIDADE

PRECISÃO ???

# IA Generativa para disseminação de conteúdo

**Plataforma baseada em IA Generativa para produção automatizada de conteúdo sobre o Estado de São Paulo. Utiliza bases de dados e análises produzidas pelo SEADE como fonte primária.**



Reunião de bancos de dados e textos analíticos e metodológicos do SEADE  
Processamento por IA para consulta, análise e interpretação dos dados por modelos de linguagem



Geração de conteúdo  
Criação automatizada de relatórios, análises e visualizações



Disseminação personalizada às necessidades dos usuários  
Tecnologia para transformar dados e textos em respostas a perguntas de usuários

# Centro de Ciência de Dados para Estatísticas Públicas (CCDEP)



## PROJETO FAPESP

- Início em outubro de 2024 e duração de 5 anos
- SEADE como instituição sede



## PARCERIAS

- Instituições de pesquisa: FFLCH, POLI, IME e ICMC da USP, Unesp, FGV, FBSP
- Órgãos de governo: Sefaz, CPTM, STM, SSP e Conselho Estadual de Educação.



## OBJETIVOS

- Utilizar dados de alta frequência, em grande volume e não estruturados para indicadores de políticas públicas
- Desenvolver novas ferramentas computacionais para coleta, processamento, análise e disseminação.
- Disseminar conhecimento na administração pública.

## LINHAS DE PESQUISA:

- 1 ATIVIDADE ECONÔMICA:** Registros fiscais para indicadores econômicos (NFe), indicadores setoriais e regionais, sistema multiagente para coleta de textos e uso de IA Generativa
- 2 MERCADO DE TRABALHO:** Coleta com URA para caracterizar ocupados, desocupados, inativos, **coleta de anúncios de emprego** e classificação textual de competências profissionais com ML
- 3 MOBILIDADE URBANA:** Indicadores a partir de dados de telefonia móvel, matrizes origem-destino para regiões metropolitanas e integração com dados de bilhetagem eletrônica
- 4 SEGURANÇA PÚBLICA:** Metodologias para **rotinas de PLN de registros criminais** e outros documentos e para **indicadores de população flutuante**, integrando com dados demográficos
- 5 TECNOLOGIAS EM ENGENHARIA DE DADOS E DE SOFTWARE:** P&D de plataforma em Data Mesh, aprimorar modelos de coleta em plataformas web e classificação de conteúdos com novas ferramentas e uso de IA generativa

# Obrigado!

Carlos Torres Freire  
[carlosfreire@seade.gov.br](mailto:carlosfreire@seade.gov.br)

**SEADE**  
Fundação Sistema Estadual  
de Análise de Dados

