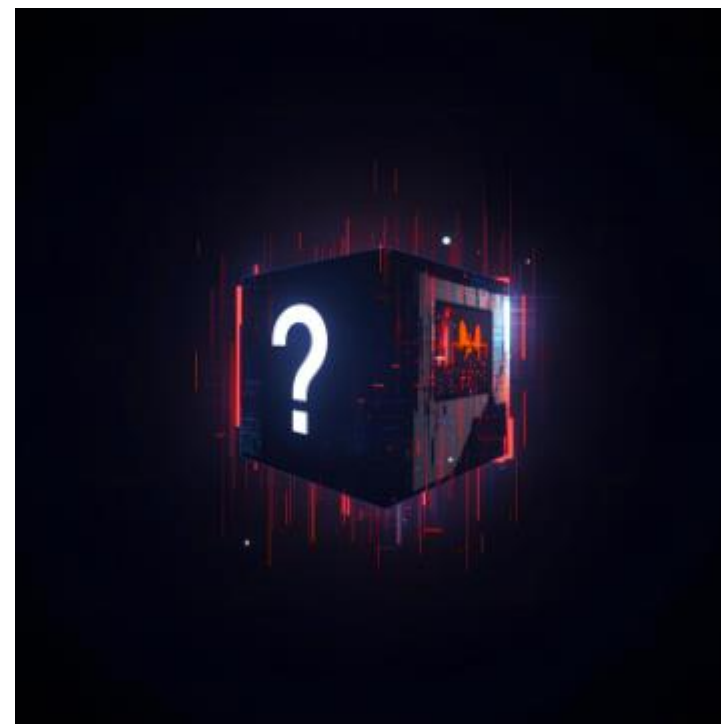




Not Just Another Black Box: AI in the Age of Trusted Statistics

Sónia Quaresma
Instituto Nacional de Estatística





What Has Generative AI Brought to the Table for Official Statistics? Large Language Models

These technologies offer new capabilities that can:

- transform how statistical agencies process data;
- build infrastructure;
- interact with users;





Not Just Another Black Box: AI in the Age of Trusted Statistics

Powerful capabilities Support innovation and efficiency

Extracting Structure from Unstructured Data:

- GenAI can analyze free-text inputs (like survey comments or administrative records) and extract relevant structured variables, enabling statistical use of previously untapped data sources





Not Just Another Black Box: AI in the Age of Trusted Statistics

Powerful capabilities Support innovation and efficiency

Processing Multiple Data Modalities:

- Modern GenAI tools can interpret not just text, but also audio, images, and video, potentially expanding the scope of data collection and integration (e.g., image classification, speech-to-text processing in interviews, video metadata analysis)





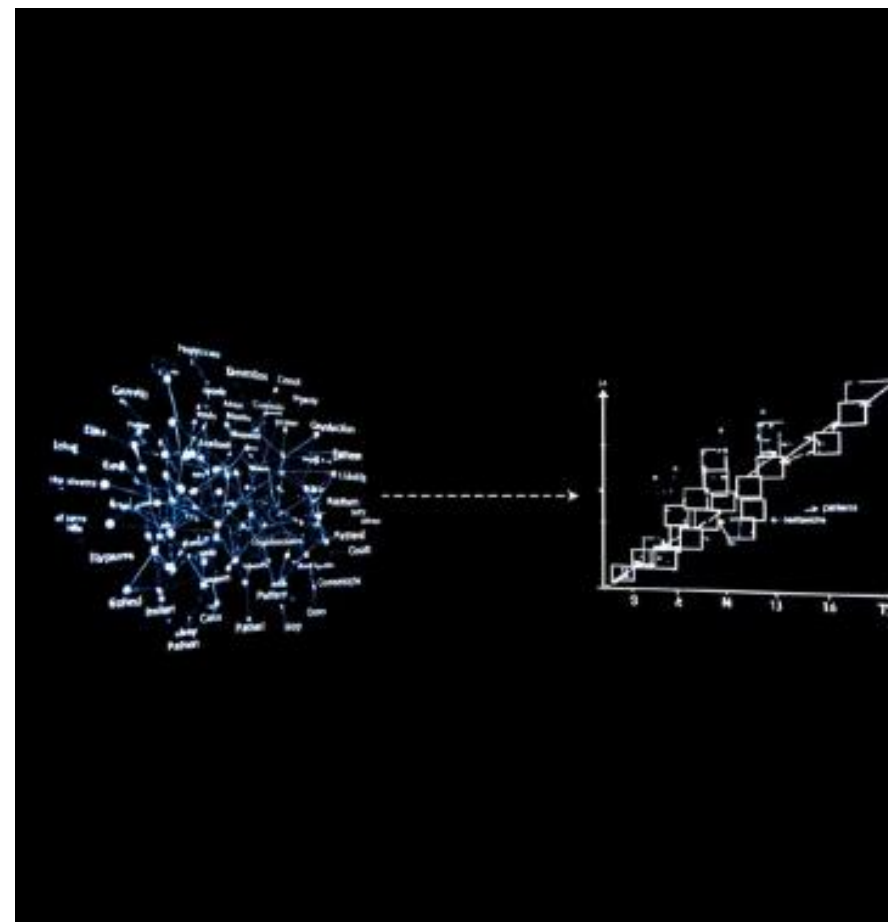
Not Just Another Black Box: AI in the Age of Trusted Statistics



Powerful capabilities Support innovation and efficiency

Support for Causality and Pattern Detection:

- While not a substitute for formal causal modelling, LLMs can surface relationships or hypotheses from large datasets that may guide further statistical analysis or exploratory work





Not Just Another Black Box: AI in the Age of Trusted Statistics



Powerful capabilities Support innovation and efficiency

Semantic Interoperability:

- GenAI helps harmonize data from different sources by aligning schema definitions, resolving naming conflicts, and disambiguating contextual meanings, key for cross-country or multi-source statistical projects.

id	Data	Data namespaces	Atber:icng	tydichght	usveogereets
1	1.aJehgpatinemo:(11:	Cuncalbalinfeadicesta.	1.leffecoplistentib1:	ty.tents	(o2e(10.9.381)
2	8.aJebulapfates:(11:	suneleadApatisfarkiug-	3:6.1{topplisessii:	ty.tents	(2.85)
5	T.aJenderfesterf(11:	suseleadAbleleFarkiug-	2:5.1{steplisteall:	ty.tents	(6.42)
9	T.aJebdrtostert(12:	subeleadApcisierkiug-	2:5.A{topplistentib1	ty.tents	(1.88)
6	7.aJebubatestses(1):	suoeleadAbesitierkiug-	2:5.A{topplizessil:	ty.tents	(3.100)
4	7.aJebulatesterf(21:	byeeleadAbesiofarkiug-	2:5.2{topplistesil:	ty.tents	(2.508)
9	8.aJebulatesterf(21:	byueleadAhsersierkiug-	3:5.4{leiparAdbl:	ty.tents	(5.93)
9	A.aJebulataster(111:	exueleadAAbdelG: Acee	igliates(2legi{tcntoliite	ty.tents	(16.01)
7	5.aJebdentesseet(11:	..Eaunedfestoidlieerice	1.3.A.Renplalieililugs-	ty.tents	(120(11.55)
1	8:027 GS telreb: -#P/Tvate	*alcorep.nal:	Getl naldesdayyereTrSesariatem	>	
2	0.aJeeoarcalsises(0)	evveleuccaser(116:	7:5.8.CZrac{csbes(1):	ty.tents	(2.981)
7	A.adubiegroeatcco(1:	adepliyeras(1011)	1:6.A.CCoulisises(11:	ty.tents	(14700)
8	A.aJebulienner#2-A-tessu	liptfecer(006)	1:4:2.CcacbVatesl:	ty.tent2	(1.925)
y	A.aJepully fecrdsee{0et	peoliprager:000}	2.7.lraIiseVioa 1.291	tiolIsram	(2.2083)
q	A.aJepulidemets <esote	oveuIntfecer(008)	1:5.3.temolii gyoet2	liireabs	1.3589.60)
y	A.aJepuliteasnent(1:	evveIiptfecer(001)	1:5.A.Idellis(isneas:	cy.tent2	(1.3647)
v	A.aJeeoliasges-(2)	evveIiptfecst(001)	1:5.4.leoplie(teale:	cy.tent2	(21.3.328)
y	A.aJepully *icrdseEnatt	vumdfosdo11065	1:6.5.iterseto:	ty.tent2	(1.3231)
d	Q.recoollisgeen(0)	evveIiseIacri:1015	1:8.4.Jcoplis(ishage:	cy.tent2	(1.806)
5	T.aJepully #ecrdse>Eart	verIigeraeer:017)	9:4.A.Idellieffsmeest:	ty.tent2	(10.507)
9	A.aJebulligotoade-fter	ceulsbenhier(005)	2:6.4.Idellit(eerent:	ty.tent1	(1.02.54)
8	5.2:ades L.P-4lsuebs	pecki(7)urelo 066vive:	3:2.A.Jonliat(smeent:	ty.tent2	(2.9911)
9	5.asSCatice (STA)	9yaalouhofantigraven	5:6.4.ltgplio{irfage:	esecten2:8o0byRiesE:	
9	8.atScootle (STAT)	Garsiontierpistfocadse	7:6.2.luplal{f{spepe	>> eisito Trep{luralateri<	
7	T.asCChrosen*LuA:	aqueleatetevisferaser	0:7.1.3jeat{/sl.ver:1:	tybes (8>.268)	
8	1:6repurcfection(1:	Titelalletinfeenter	2:3.3.teInficities/12:	carces(4)	
2	5.aZuphlice (R)	1lteleaZbul3ca5	8:7.2.Irelp)/ss ysef.2.	S.Irient.Passur >>	
3	9.Ovepblioel (R)	BatotiooACheraiertfer	3:5.3.lcetnirsccutRL:	c:litent\$ (1.980-	





Not Just Another Black Box: AI in the Age of Trusted Statistics



Powerful capabilities Support innovation and efficiency

Contextual Aggregation and Disambiguation:

- By understanding the broader linguistic context, LLMs can aggregate values intelligently or clarify ambiguous terms.





Not Just Another Black Box: AI in the Age of Trusted Statistics



Powerful capabilities Support innovation and efficiency

Synthetic Data Generation:

- GenAI can produce realistic synthetic datasets useful for testing methods, system development, or training purposes, while helping protect the privacy of real individuals.





Not Just Another Black Box: AI in the Age of Trusted Statistics

Significant Risks For Official Statistics

particularly sensitive in this field are:

- Trust
- Transparency
- Accuracy





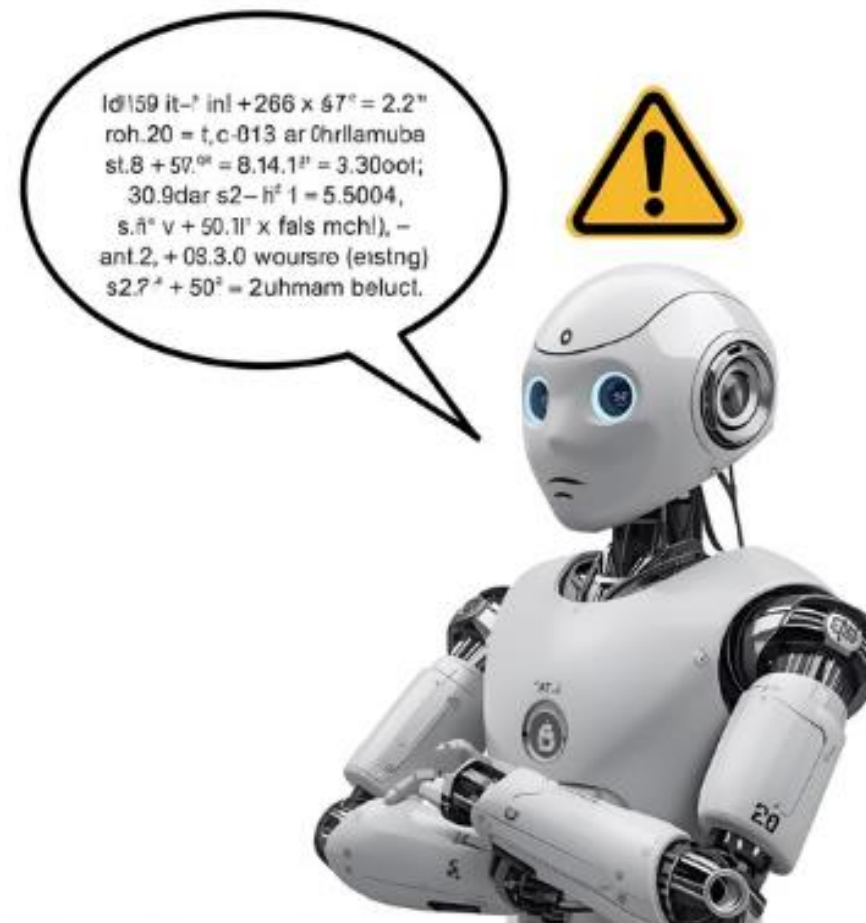
Not Just Another Black Box: AI in the Age of Trusted Statistics



Significant Risks For Official Statistics

Uncontrolled Uncertainty and Hallucinations:

- LLMs may produce plausible sounding but entirely incorrect or invented content, known as hallucinations, posing a risk of misinformation if outputs are not validated.





Not Just Another Black Box: AI in the Age of Trusted Statistics

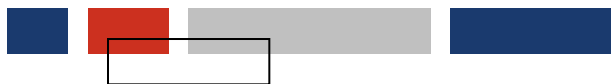


Significant Risks For Official Statistics

Bias:

- Inherited from their training data, LLMs may produce biased or unfair outputs, which can conflict with the principles of impartiality and fairness that underpin official statistics





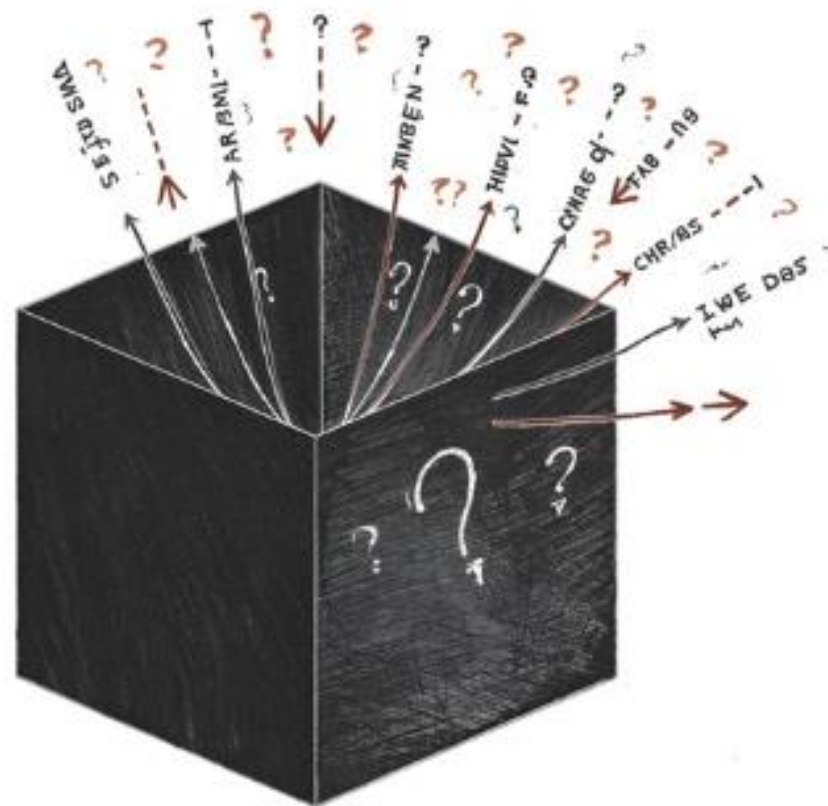
Not Just Another Black Box: AI in the Age of Trusted Statistics



Significant Risks For Official Statistics

Lack of Transparency:

- LLMs are often considered “black boxes,” making it difficult to understand or trace the rationale behind specific outputs





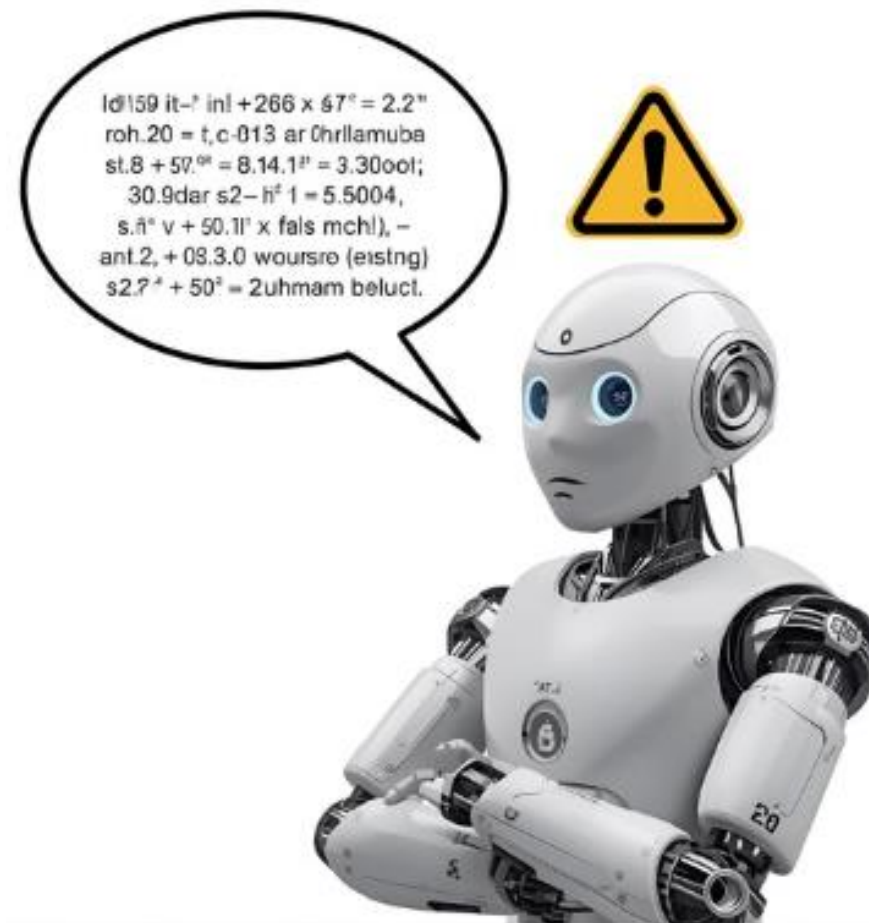
Not Just Another Black Box: AI in the Age of Trusted Statistics



Significant Risks For Official Statistics

Numerical Inaccuracy:

- GenAI models are poor at performing calculations or interpreting numeric data accurately, making them unsuitable for tasks requiring exact values or arithmetic precision.





Not Just Another Black Box: AI in the Age of Trusted Statistics



Significant Risks For Official Statistics

Limited Reasoning Capabilities:

- While LLMs can reproduce patterns seen in training data, they often fail at logical reasoning, deduction, and generalising to unseen situations, limiting their utility in analytical tasks requiring inference or abstraction





Not Just Another Black Box: AI in the Age of Trusted Statistics



Significant Risks For Official Statistics

Suggestive but Not Self-Solving:

- LLMs can suggest code structures or outline problem-solving steps but cannot execute code or evaluate its correctness. They are assistive, not autonomous decision-makers.





Not Just Another Black Box: AI in the Age of Trusted Statistics



Significant Risks For Official Statistics

High Computational Costs:

- Running GenAI models, especially at national scale, requires significant infrastructure, energy, and maintenance, raising questions of cost-efficiency and sustainability.





Not Just Another Black Box: AI in the Age of Trusted Statistics



Significant Risks For Official Statistics

Lack of Theoretical Guarantees:

- Unlike traditional statistical models, GenAI tools typically lack the formal mathematical foundations needed for reproducibility, error bounds, or inference guarantees.





Not Just Another Black Box: AI in the Age of Trusted Statistics



Significant Risks For Official Statistics

Scalability Challenges with Big Data:

- GenAI models are not inherently optimized for large-scale data processing. Integrating them into big data pipelines can be costly and technically demanding.





Not Just Another Black Box: AI in the Age of Trusted Statistics



Significant Risks For Official Statistics

Vulnerability to Prompt Hacking:

- LLMs can be manipulated through carefully crafted inputs to produce misleading or inappropriate outputs, raising concerns about reliability, especially in open-ended tools.





Not Just Another Black Box: AI in the Age of Trusted Statistics



Significant Risks For Official Statistics

Privacy and Confidentiality Risks:

- If not properly controlled, LLMs may unintentionally reproduce sensitive information seen during training, which could violate confidentiality policies fundamental to official statistics.





Not Just Another Black Box: AI in the Age of Trusted Statistics



Should We Still Make the Case for Using LLMs?

Aligning Use with Statistical Principles

Accuracy and Reliability:

- Outputs must be verified, contextualized, and never used as the sole source of truth for statistical indicators.





Not Just Another Black Box: AI in the Age of Trusted Statistics



Should We Still Make the Case for Using LLMs?

Aligning Use with Statistical Principles

Relevance and Timeliness:

- LLMs can help speed up interaction and document processing but must rely on curated, up-to-date sources to avoid spreading outdated information.





Not Just Another Black Box: AI in the Age of Trusted Statistics



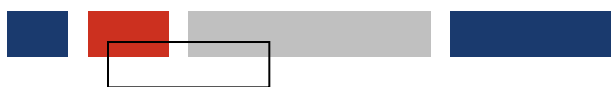
Should We Still Make the Case for Using LLMs?

Aligning Use with Statistical Principles

Transparency and Impartiality:

- Clear documentation of the models, training data, and limitations is essential to maintain trust and ensure fair use.





Not Just Another Black Box: AI in the Age of Trusted Statistics



Should We Still Make the Case for Using LLMs?

Aligning Use with Statistical Principles

Explainability:

- Mechanisms should be implemented to trace how conclusions were reached, particularly when used to assist decision-making or public communication.





Not Just Another Black Box: AI in the Age of Trusted Statistics

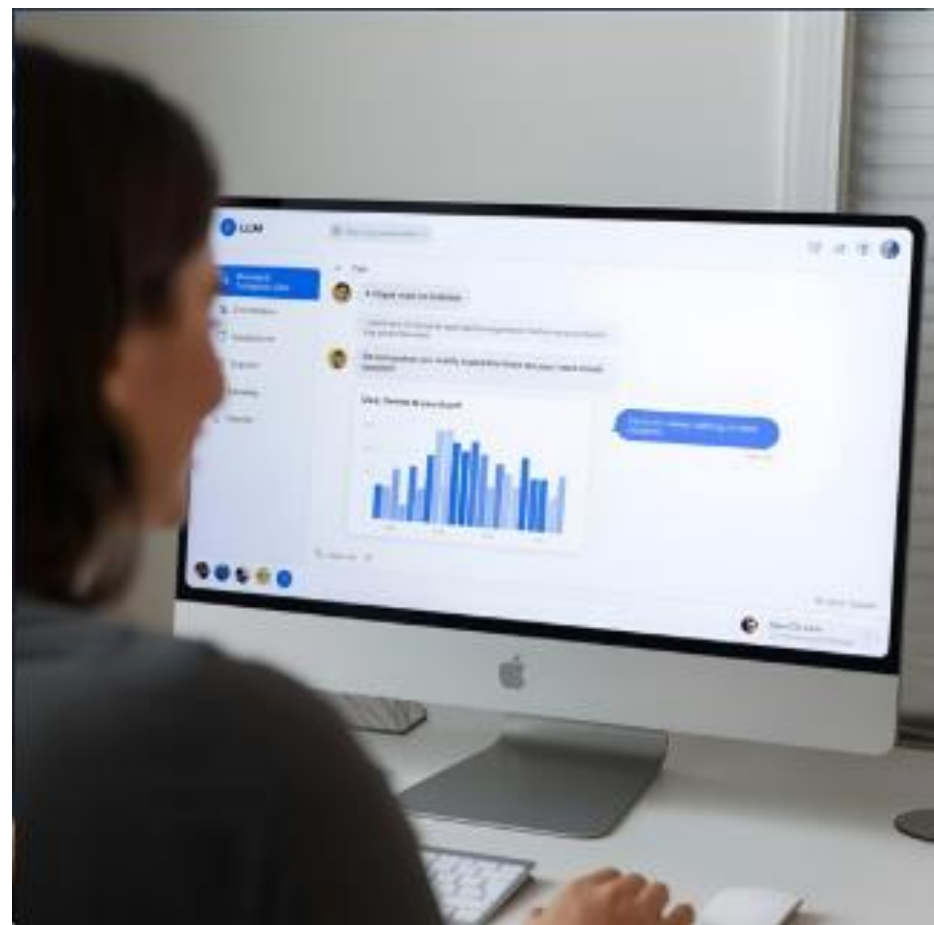


Should We Still Make the Case for Using LLMs?

Aligning Use with Statistical Principles

Accessibility:

- Well-designed LLM interfaces can make statistical information more approachable, especially for non-expert users.





Not Just Another Black Box:
AI in the Age of Trusted Statistics



Should We Still Make the Case for Using LLMs?

Aligning Use with Statistical Principles

Confidentiality:

- All implementations must rigorously safeguard private data, preventing the leakage of sensitive information from training or interactions.





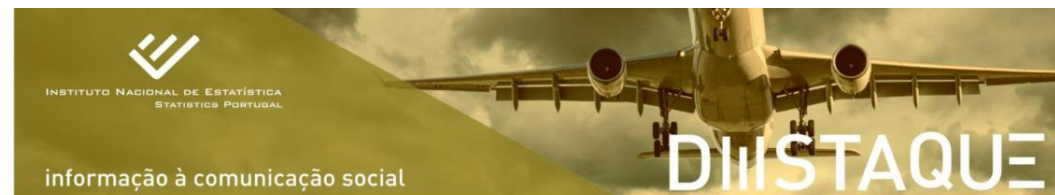
Not Just Another Black Box: AI in the Age of Trusted Statistics



From Potential to Practice Prototyping with LLMs at INE

Promising Use Cases:

- **Summarizing historical documents** - LLMs can assist in synthesizing large volumes of previously published material, such as statistical reports, metadata, and methodological notes, making it easier to extract relevant insights and provide contextual overviews.



12 de junho de 2025
ATIVIDADE DOS TRANSPORTES
Estatísticas rápidas do transporte aéreo
abril 2025

MOVIMENTO DE PASSAGEIROS NOS AEROPORTOS NACIONAIS MANTEVE TENDÊNCIA DE CRESCIMENTO

Em **abril de 2025**, nos aeroportos nacionais movimentaram-se 6,5 milhões de passageiros e 21,4 mil toneladas de carga e correio, correspondendo a variações de +8,1% e +1,7%, respetivamente, face a abril de 2024 (+2,1% e +3,9% no mês anterior, pela mesma ordem).

Em abril de 2025 registou-se o desembarque médio diário de 109,9 mil passageiros, valor superior ao registado em abril de 2024 (101,5 mil; +8,3%).

O Reino Unido passou a ser o principal país de origem e de destino dos voos, considerando os primeiros quatro



Not Just Another Black Box: AI in the Age of Trusted Statistics

From Potential to Practice Prototyping with LLMs at INE

Promising Use Cases:

- **Helping users navigate published content** - With the increasing amount of information available on INE's portal, LLMs could help users find relevant tables, publications, or concepts by interpreting natural language queries and mapping them to structured resources.



The screenshot shows the INE Statistics Portugal website interface. The header includes the logo and name 'INSTITUTO NACIONAL DE ESTATÍSTICA | STATISTICS PORTUGAL' and a navigation menu with 'Estatísticas', 'Território', 'Produtos', and 'WebInq'. The breadcrumb trail is 'Início / Produtos / Base de Dados'. The main content area shows a table titled 'Taxa de desemprego (Série 2021 - %) por Local de residência (NUTS - 2024)'. The table is filtered for '2.º Trimestre de 2025' and 'Sexo'. The data is presented in a grid format with columns for 'Local de residência (NUTS - 2024)', 'HM', and 'Sexo'. The rows include 'Portugal', 'Continente', 'Região Autónoma dos Açores', and 'Região Autónoma da Madeira'.

Local de residência (NUTS - 2024)	Taxa de desemprego (Série 2021 - %) por Local de residência (NUTS - 2024)	
	Período de referência dos dados (1)	
	2.º Trimestre de 2025	
	Sexo	
	HM	H
	%	%
Portugal	5,9	5,3
Continente	6,0	5,4
Região Autónoma dos Açores	3,9 §	4,0 §
Região Autónoma da Madeira	4,0 §	4,2 §



Not Just Another Black Box: AI in the Age of Trusted Statistics



From Potential to Practice Prototyping with LLMs at INE

Promising Use Cases:

- Supporting metadata management - LLMs may support harmonization of terminology, the completion of metadata fields, or the alignment of definitions across sources, improving coherence and reducing manual work.

Interface: Conteúdo: Tem conta

INSTITUTO NACIONAL DE ESTATÍSTICA
STATISTICS PORTUGAL

Sistema de Metainformação

Palavra a pesquisar

Geral **Conceitos** **Classificações** **Documentação metodológica** **Variáveis** **Suportes de recolha**

SMI

Conceitos
Módulo que integra e disponibiliza termos e definições dos conceitos utilizados no âmbito do Sistema Estatístico Nacional que fundamentam a consistência e comparabilidade dos dados recolhidos e difundidos. [Mais...](#)

Cód.	Designação
2051	ATIVIDADE ECONÓMICA

Variáveis
Módulo que integra e disponibiliza as variáveis observadas e os indicadores divulgados no âmbito do Sistema Estatístico Nacional. [Mais...](#)

Cód.	Designação
10009	Remunerações (E) das Empresas por Afili...
11900	Nados-vivos (N.º) por Local de residênci...

Sistema de Metainformação
O Sistema de Metainformação integra e disponibiliza conceitos, classificações, variáveis, suportes de recolha de informação e documentação metodológica com aplicação no âmbito do Sistema Estatístico Nacional (SEN). Os vários componentes do sistema estão interrelacionados, visam apoiar a produção estatística e documentar a difusão de Estatísticas Oficiais. O Sistema de Metainformação é um instrumento de coordenação e harmonização no seio do SEN. [Mais...](#)

Documentação Metodológica
Módulo que integra e disponibiliza (designados Documentos Metodológicos) descrições, com atualidade, da estatística realizada no âmbito do Sistema Estatístico Nacional. Obedece estrutura comum, com tópicos: [Mais...](#)

Cód.	Designação
------	------------

Suportes de Recolha
Módulo que integra e disponibiliza sobre os instrumentos de notação no INE e respetiva imagem digital de questionários (suportes assíncronos) e ficheiros (suportes eletrónicos). [Mais...](#)

Cód.	Designação
------	------------

Classificações
Módulo que integra e disponibiliza as classificações nacionais, comunitárias e internacionais utilizadas para fins estatísticos no âmbito do Sistema Estatístico Nacional. [Mais...](#)

Código	Designação	Sigla
V05497	Classificação portuguesa das atividades económicas, revisão 4	CAE Rev 4
V00017	Código da divisão administrativa (distritos/municípios/freguesias)	

Cód.	Designação
------	------------

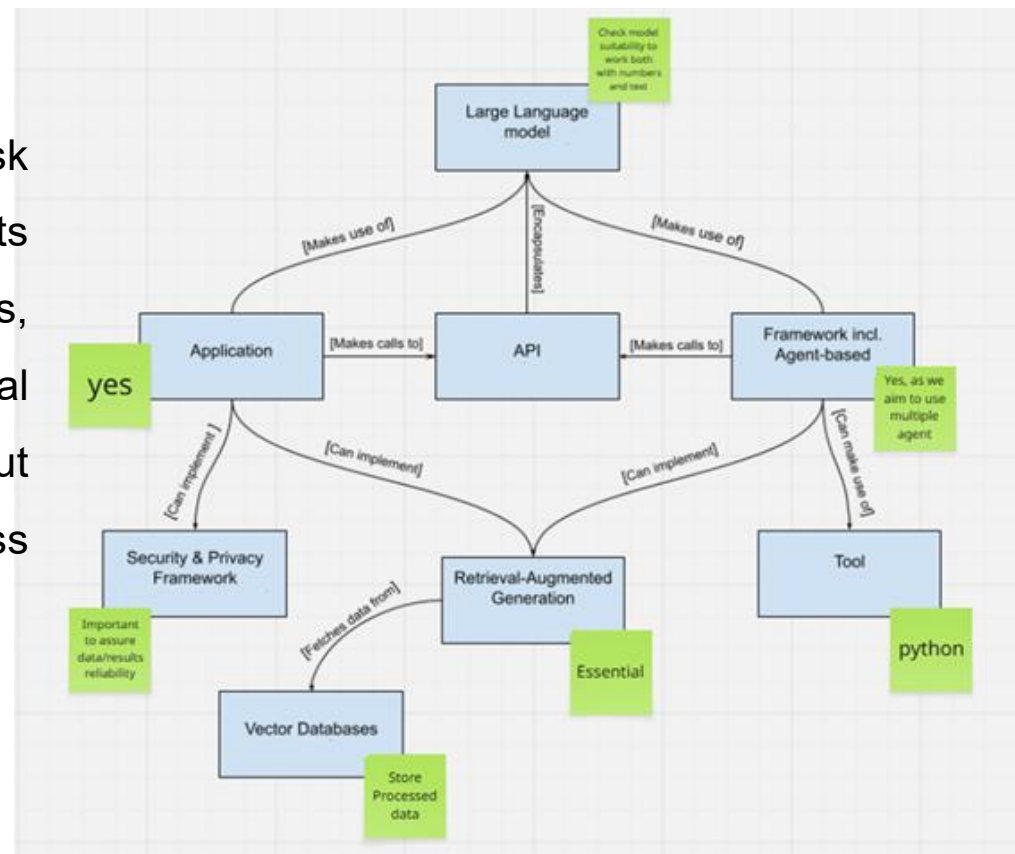




Not Just Another Black Box: AI in the Age of Trusted Statistics

Bringing It All Together A Prototype for Smarter Dissemination

At the heart of this prototype is a common yet time-consuming task faced by many statistical offices, summarizing analytical reports into clear, concise overviews for publication. These summaries, sometimes of historical data, are typically required in both the local language and English. The challenge is not just linguistic; it's about ensuring consistency, clarity, and contextual fidelity across domains.





Not Just Another Black Box: AI in the Age of Trusted Statistics



Input

The system begins by ingesting a report, usually a PDF authored by a subject-matter expert, and extracting its content for semantic processing.



12 de junho de 2025
ATIVIDADE DOS TRANSPORTES
Estatísticas rápidas do transporte aéreo
abril 2025

MOVIMENTO DE PASSAGEIROS NOS AEROPORTOS NACIONAIS MANTEVE TENDÊNCIA DE CRESCIMENTO

Em **abril de 2025**, nos aeroportos nacionais movimentaram-se 6,5 milhões de passageiros e 21,4 mil toneladas de carga e correio, correspondendo a variações de +8,1% e +1,7%, respetivamente, face a abril de 2024 (+2,1% e +3,9% no mês anterior, pela mesma ordem).

Em abril de 2025 registou-se o desembarque médio diário de 109,9 mil passageiros, valor superior ao registado em abril de 2024 (101,5 mil; +8,3%).

O Reino Unido passou a ser o principal país de origem e de destino dos voos, considerando os primeiros quatro



Not Just Another Black Box: AI in the Age of Trusted Statistics

Interface





Not Just Another Black Box:

AI in the Age of Trusted Statistics



Code

A Retrieval-Augmented Generation (RAG) pipeline is then deployed, enabling the LLM to work with the actual content of the report rather than relying on pretraining alone.

```
wp12_hackathon > dissemination_summary_prototype > summarizer_unified.py
129 class PDFSummarizer:
207     def process_pdf(self, pdf_path: str, use_vector_store:bool = False, document_loader:str
242     else:
252
253         # 3 - create the summary chain
254         chain = self._create_summary_chain()
255
256         result = chain.invoke({
257             "content": context,
258             "max_keywords": max_keywords,
259             "max_tags": max_tags,
260             "max_words": max_words,
261             "out_lang": out_lang
262         })
263
264         return result
265
266
267     def _create_summary_chain(self):
268         """
269         Create a runnable sequence for direct text summarization (non-vector store mode).
270         """
271         chat_prompt, _ = self._get_shared_prompts()
272         return chat_prompt | self.llm | self._safe_json_parse
273
```





Not Just Another Black Box: AI in the Age of Trusted Statistics



Output

```
summarizer/.venv_pdf/bin/python /home/ivo/Desktop/Trabalho/WP12/proof-of-concepts/pdf_summarizer/testing_pdf_llm.py
```

Com base nos relatórios apresentados, podemos observar que há uma tendência crescente no movimento de passageiros nos aeroportos nacionais. De acordo com o Documento 1 (página 4), entre janeiro e abril de 2025, os aeroportos nacionais movimentaram-se 6,5 milhões de passageiros, o que representa um crescimento de +8,1% em relação ao mesmo período do ano anterior. Além disso, o desembarque médio diário de 109,9 mil passageiros em abril de 2025 é superior ao registado em abril de 2024 (101,5 mil; +8,3%).

O Documento 2 (página 1) também confirma essa tendência, indicando que o movimento de passageiros nos aeroportos nacionais manteve uma tendência de crescimento. Em abril de 2025, os aeroportos nacionais registaram um desembarque médio diário de 109,9 mil passageiros, o que é superior ao registado em abril de 2024.

Além disso, o Documento 3 (página 2) apresenta uma análise detalhada do movimento de passageiros nos aeroportos nacionais, indicando que os passageiros desembarcados e embarcados em voos internacionais aumentaram significativamente. Em abril de 2025, 82,5% dos passageiros desembarcados correspondiam a tráfego internacional, atingindo 2,7 milhões de passageiros (+8,5%), e os passageiros embarcados também apresentaram um crescimento significativo.

Em resumo, a tendência recente no movimento de passageiros nos aeroportos nacionais é um crescimento significativo, com aumento tanto no número de passageiros desembarcados quanto no número de passageiros embarcados em voos internacionais.



h Flag

Ln 77, Col 23 Spac





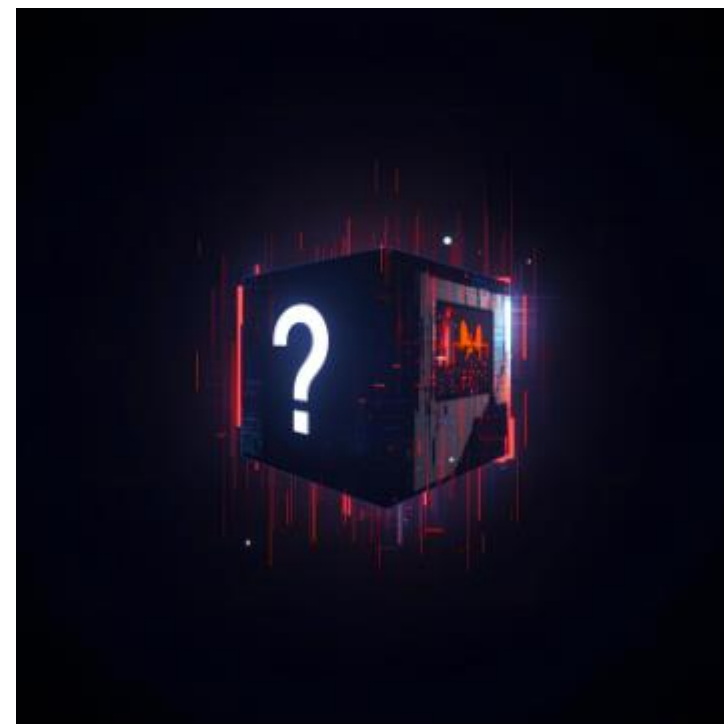
Not Just Another Black Box: AI in the Age of Trusted Statistics

Bringing It All Together A Prototype for Smarter Dissemination

More than just a tech demo, this prototype offers a glimpse into how AI can be embedded thoughtfully into statistical processes, augmenting expert work, not replacing it.

It reflects a pragmatic approach to innovation at INE:

- solving concrete problems,
- testing real tools, and
- doing so in a way that's scalable across national contexts.





Not Just Another Black Box: AI in the Age of Trusted Statistics



Bringing It All Together A Prototype for Smarter Dissemination

As we continue to explore AI's place in official statistics, prototypes like this remind us that progress doesn't always require revolution. Sometimes, it's about making the everyday easier, smarter, and more consistent.

Thank you for your attention!
Questions?

