measuring attribution in natural language generation models

Measuring Attribution in Natural Language Generation Models: Understanding the Why Behind the Words

measuring attribution in natural language generation models is becoming increasingly important as AI-generated text finds its way into more aspects of our digital lives. From chatbots assisting customers to algorithms drafting news articles, these models are shaping how we communicate and consume information. But as these systems grow in complexity, questions arise about how they produce specific outputs. How can we determine which parts of the input data influence a generated response? What role do certain tokens or training examples play in shaping the final text? This is where attribution measurement steps in, offering transparency and interpretability in natural language generation (NLG).

In this article, we'll dive into the nuances of measuring attribution in natural language generation models, exploring the techniques, challenges, and implications of understanding the inner workings of these powerful AI tools.

Why Measuring Attribution Matters in Natural Language Generation

When a language model generates text, it's essentially predicting the next word based on patterns it learned from vast datasets. However, the "black box" nature of these models means it's not always clear why a particular phrase or sentence was produced. Measuring attribution helps shed light on this process, allowing researchers, developers, and users to:

- **Understand model decision-making:** Attribution reveals which input features, tokens, or data points most influenced the output.
- **Improve model transparency:** Knowing why a model generated certain content fosters trust and accountability.
- **Debug and refine models:** Attribution highlights biases or errors by showing unexpected influences on the output.
- **Meet regulatory requirements:** For applications in sensitive fields like healthcare or finance, interpretability through attribution is crucial for compliance.

By focusing on how different parts of the input contribute to the generation, we gain valuable insights into model behavior, paving the way for more responsible AI deployment.

Core Techniques for Measuring Attribution in NLG Models

Attribution methods in natural language generation borrow ideas from explainable AI (XAI) and

interpretability research. Some techniques are adapted specifically to handle the challenges posed by sequential text generation.

Gradient-Based Attribution

One common approach uses gradients, calculating how changes in input tokens affect the model's output probabilities. By backpropagating through the model, gradients highlight which words or phrases have the greatest influence on a prediction.

- **Integrated Gradients:** This method accumulates gradients along a path from a baseline input (like an empty sentence) to the actual input, producing a more stable attribution score.
- **SmoothGrad:** By averaging gradients over multiple noisy versions of the input, SmoothGrad reduces sensitivity to small perturbations, yielding clearer attributions.

These gradient-based methods are computationally efficient and compatible with transformer architectures like GPT and BERT, making them practical for real-world NLG attribution.

Attention Mechanisms as Attribution Proxies

Transformers rely heavily on attention mechanisms to weigh input tokens. Some researchers argue that attention weights themselves can serve as a form of attribution, indicating which parts of the input the model "focused" on when generating output.

However, attention isn't always a perfect explanation. While it provides useful clues, attention weights don't always correlate with feature importance and should be interpreted with caution. Combining attention analysis with other attribution techniques often yields more comprehensive insights.

Perturbation and Occlusion Methods

Another intuitive way to measure attribution is to systematically alter or remove parts of the input and observe changes in the output:

- **Token occlusion:** Deleting or masking individual words to see how the output probability shifts.
- **Input perturbation: ** Introducing noise or replacing phrases to test sensitivity.

These approaches are model-agnostic and easy to understand but can be computationally expensive, especially with long sequences, since multiple modified inputs must be processed.

Shapley Values and Game-Theoretic Attribution

Inspired by cooperative game theory, Shapley values assign an attribution score to each input feature by considering all possible combinations. This method provides a fair and theoretically

grounded measure of contribution.

Although calculating exact Shapley values is often infeasible for large language models due to combinatorial explosion, approximations like SHAP (SHapley Additive exPlanations) have been adapted for NLP tasks, offering meaningful insights into input influence.

Challenges in Measuring Attribution for NLG

Despite progress, measuring attribution in natural language generation models is still fraught with difficulties that researchers continue to tackle.

Complexity of Language and Context

Language is inherently complex and context-dependent. The influence of a single word can vary dramatically depending on surrounding words or the overall sentence structure. Capturing these intricate dependencies in attribution scores is a non-trivial challenge.

Non-Linear and Distributed Representations

Deep language models encode information in high-dimensional, distributed embeddings. Attribution methods must navigate these non-linear transformations to accurately trace how input features propagate through layers. This complexity can obscure straightforward cause-effect relationships.

Sequence Length and Computational Cost

Natural language inputs can be lengthy, and generating explanations often requires multiple forward and backward passes through the model. This leads to significant computational overhead, especially when applying perturbation or Shapley-based methods.

Attribution Granularity

Determining the right level of granularity for attribution is tricky. Should explanations focus on individual tokens, phrases, sentences, or even higher-level concepts? Different applications demand different granularities, complicating the development of universal attribution frameworks.

Practical Tips for Measuring Attribution in NLG Projects

If you're working with natural language generation models and want to incorporate attribution measurements, here are some useful strategies to consider:

- **Combine multiple attribution methods:** Relying on a single technique can be misleading. Blend gradient-based, attention, and perturbation approaches to get a holistic view.
- **Use baseline inputs wisely:** For methods like Integrated Gradients, choose meaningful baselines (e.g., empty text or neutral sentences) to improve attribution stability.
- **Visualize attributions:** Heatmaps or highlight overlays on generated text help communicate attribution results effectively to non-technical stakeholders.
- **Consider downstream impact:** Attribution is not just about technical explanations but also about how it affects user trust, model fairness, and ethical AI deployment.
- **Iterate and validate:** Attribution techniques should be evaluated for consistency and alignment with human intuition through user studies or benchmark datasets.

Emerging Trends and Future Directions

The field of measuring attribution in natural language generation models is rapidly evolving. Some exciting developments on the horizon include:

- **Causal attribution approaches:** Moving beyond correlation-based methods, causal inference techniques aim to establish cause-effect links between inputs and outputs in NLG.
- **Multimodal attribution:** As language models integrate with images, audio, and other modalities, attribution methods are expanding to handle richer input types.
- **Real-time explainability:** Developing efficient algorithms that provide attribution feedback during text generation, enabling interactive AI systems.
- **User-centric explanations:** Tailoring attribution outputs to different audiences, from developers to end-users, to enhance comprehension and trust.

As AI-generated text becomes more pervasive, these advances will be crucial in ensuring that natural language generation models are not just powerful but also transparent and accountable.

Understanding the "why" behind the words generated by AI is no longer a luxury—it's a necessity. Measuring attribution in natural language generation models opens a window into the model's reasoning process, helping us harness these technologies responsibly and effectively.

Frequently Asked Questions

What is attribution in the context of natural language generation (NLG) models?

Attribution in NLG models refers to the process of identifying and tracing the source or origin of the generated content, determining which parts of the input data or training information contributed to specific outputs.

Why is measuring attribution important for natural language generation models?

Measuring attribution is important to ensure transparency, improve model interpretability, verify factual accuracy, and mitigate issues like misinformation or bias by understanding how and why a model produces certain outputs.

What are common methods used to measure attribution in NLG models?

Common methods include attention analysis, gradient-based attribution techniques (e.g., Integrated Gradients), perturbation methods, and using external tools like SHAP or LIME adapted for language generation.

How does attention mechanism help in attribution measurement for NLG?

Attention mechanisms highlight which parts of the input the model focuses on when generating each token, providing a form of attribution by showing the alignment between input tokens and generated output.

What challenges exist in measuring attribution for large-scale NLG models?

Challenges include the complexity and opacity of models, the distributed nature of knowledge in parameters, difficulty in interpreting attention weights as true attribution, and the lack of ground truth for validating attribution methods.

Can attribution measurement improve the factual accuracy of NLG outputs?

Yes, by identifying the sources influencing specific outputs, attribution measurement can help detect when a model generates hallucinated or unsupported information and guide improvements to enhance factual accuracy.

Are there benchmark datasets or frameworks for evaluating attribution in NLG models?

While benchmarks specifically for attribution in NLG are emerging, some datasets like FEVER and

fact-checking corpora, combined with explainability frameworks, are used to evaluate attribution and factual grounding.

How do perturbation-based methods work for measuring attribution in NLG?

Perturbation-based methods involve systematically modifying or removing parts of the input and observing the impact on generated outputs to infer which input components are most influential for specific outputs.

What role does human evaluation play in assessing attribution in NLG?

Human evaluation helps validate the quality and reliability of attribution methods by assessing whether the identified sources or explanations align with human understanding and domain knowledge.

Additional Resources

Measuring Attribution in Natural Language Generation Models: Challenges and Approaches

measuring attribution in natural language generation models has become a critical area of research and application, as these models increasingly influence decision-making, content creation, and customer engagement across industries. Attribution, in this context, refers to the process of identifying which parts of the input data, model components, or training signals contribute to a specific output generated by a natural language generation (NLG) system. Understanding and quantifying attribution is essential for improving transparency, trustworthiness, and interpretability of language models, especially as their complexity continues to grow.

As NLG models like GPT, BERT-derived generators, and other transformer-based architectures achieve remarkable fluency and adaptability, the challenge lies not only in their ability to produce coherent text but also in explaining how and why certain outputs are produced. This article delves into the methodologies, challenges, and practical implications of measuring attribution in natural language generation models, while exploring the latest research trends and industry practices.

The Importance of Attribution in NLG Systems

Natural language generation models operate by predicting sequences of words based on input data and learned parameters. However, the black-box nature of deep learning architectures often obscures the rationale behind specific outputs, raising concerns in domains where accountability and accuracy are paramount — such as healthcare, finance, and legal services. Measuring attribution helps stakeholders:

• Identify influential input features or tokens that steer generation outcomes.

- Diagnose errors and biases embedded in model predictions.
- Enhance model explainability for regulatory compliance and user trust.
- Guide model refinement and feature engineering for improved performance.

Attribution also plays a pivotal role in combating misinformation and ensuring that generated content aligns with factual sources, especially in applications involving summarization, question-answering, and content recommendation.

Key Techniques for Measuring Attribution in Natural Language Generation Models

Several analytical frameworks and computational methods have been proposed to measure attribution within NLG systems. These approaches vary in their focus—some emphasize input attribution (which parts of the input text influence output), while others explore internal model workings or training data influence.

Gradient-Based Attribution Methods

One of the foundational approaches to measuring attribution relies on gradient-based techniques. These methods examine how changes in input tokens or embeddings affect the output by calculating gradients of the output with respect to inputs.

- Saliency Maps: By computing the gradient of the output logit for a particular word with respect to each input token, saliency maps highlight which input tokens have the most significant influence on that output word.
- **Integrated Gradients:** This method improves upon basic saliency by integrating gradients along a path from a baseline input to the actual input, reducing noise and providing more stable attribution scores.
- Layer-wise Relevance Propagation (LRP): LRP attributes output decisions back through the layers of the neural network, distributing relevance scores to input features in a way that respects the network's architecture.

Gradient-based methods are computationally efficient and model-agnostic but can be sensitive to model non-linearities and may not always provide human-interpretable explanations.

Perturbation-Based Attribution Approaches

Perturbation methods measure attribution by systematically modifying or masking parts of the input and observing the impact on the output. This approach is intuitive and often more aligned with human reasoning.

- **Occlusion:** Removing or masking one token or phrase at a time to evaluate its effect on the generation outcome.
- **Shapley Values:** Borrowed from cooperative game theory, Shapley values estimate the contribution of each input token by considering all possible subsets of inputs, offering a theoretically sound attribution measure.
- **Counterfactual Generation:** Generating alternative outputs by changing input features to understand causal relationships between inputs and outputs.

These methods tend to be computationally expensive, especially for large models and long input sequences, but often yield more interpretable and robust attribution insights.

Attention Weights as Attribution Indicators

Given that most modern NLG models use attention mechanisms, some researchers propose leveraging attention weights as proxies for attribution. Essentially, the attention scores reflect the importance assigned to different input tokens during generation.

While attention visualization can offer intuitive explanations, studies show that attention weights do not always correlate directly with true attribution, leading to debates regarding their reliability as interpretability tools.

Challenges in Measuring Attribution for NLG Models

Despite progress, measuring attribution in natural language generation models faces several intrinsic difficulties:

Complexity and Scale of Models

State-of-the-art NLG models often consist of billions of parameters and multiple attention layers. This complexity makes it challenging to trace output decisions back to specific inputs or internal mechanisms without incurring prohibitive computational costs.

Contextual Dependencies and Non-Linearity

Language generation relies heavily on context, with outputs influenced by subtle interactions between distant tokens. Attribution methods must account for these dependencies and the nonlinear transformations performed by deep networks, complicating straightforward attribution.

Evaluation and Benchmarking of Attribution Methods

A major hurdle is the lack of universally accepted benchmarks or ground truth for attribution in NLG. Unlike classification tasks where feature importance can sometimes be directly validated, text generation outputs are inherently subjective, making attribution evaluation subjective and nuanced.

Bias and Ethical Considerations

Misattribution can obscure biases in training data or model architecture. Without accurate measurement, models may perpetuate harmful stereotypes or misinformation, underscoring the ethical imperative of reliable attribution techniques.

Emerging Trends and Tools in Attribution Measurement

The research community is actively exploring novel frameworks and tools to improve attribution measurement in NLG models.

Explainability Frameworks and Libraries

Open-source initiatives such as Captum (by Facebook AI) and AllenNLP Interpret offer integrated toolkits for applying gradient- and perturbation-based attribution methods to transformer models, facilitating experimentation and deployment.

Attribution for Training Data Influence

Recent work extends attribution beyond inputs to analyze the impact of specific training examples on generation outputs, using techniques like influence functions. This helps identify problematic data points and improves dataset curation.

Hybrid Attribution Approaches

Combining multiple attribution methods to cross-validate results is gaining traction. For instance, integrating gradient-based and perturbation methods can provide both computational efficiency and interpretability robustness.

Human-in-the-Loop Attribution

Interactive tools that allow domain experts to guide and refine attribution analysis are becoming valuable, particularly in sensitive applications like medical report generation or legal document drafting.

Practical Implications for Industry and Research

Accurately measuring attribution in natural language generation models is no longer just an academic exercise. Companies deploying chatbots, automated content generators, or recommendation engines increasingly demand transparency to maintain user trust and comply with emerging AI regulations.

Moreover, improved attribution techniques can accelerate debugging and model improvement cycles, leading to more reliable and context-aware language systems. As regulatory bodies focus on AI explainability, organizations that invest in robust attribution measurement will be better positioned to navigate compliance requirements.

In research, attribution analysis deepens our understanding of model behavior, enabling the design of architectures that are not only more performant but also inherently interpretable. This shift toward explainable AI aligns with broader trends emphasizing responsible and ethical AI deployment.

The journey toward comprehensive and reliable attribution measurement in natural language generation models continues to unfold, driven by advances in computational methods, growing awareness of AI transparency, and the ever-expanding role of NLG in modern technology ecosystems.

Measuring Attribution In Natural Language Generation Models

Find other PDF articles:

https://espanol.centerforautism.com/archive-th-117/Book?trackid=PNb27-3215&title=by-nicholas-gi ordano-college-physics-reasoning-and-relationships-1st-first-edition.pdf

Machine Learning and Knowledge Discovery in Databases. Research Track Albert Bifet, Jesse Davis, Tomas Krilavičius, Meelis Kull, Eirini Ntoutsi, Indrė Žliobaitė, 2024-09-01 This multi-volume set, LNAI 14941 to LNAI 14950, constitutes the refereed proceedings of the European Conference on Machine Learning and Knowledge Discovery in Databases, ECML PKDD 2024, held in Vilnius, Lithuania, in September 2024. The papers presented in these proceedings are from the following three conference tracks: -Research Track: The 202 full papers presented here, from this track, were carefully reviewed and selected from 826 submissions. These papers are present in the following volumes: Part I, II, III, IV, V, VI, VII, VIII. Demo Track: The 14 papers presented here, from this track, were selected from 30 submissions. These papers are present in the following volume: Part VIII. Applied Data Science Track: The 56 full papers presented here, from this track, were carefully reviewed and selected from 224 submissions. These papers are present in the following volumes: Part IX and Part X.

measuring attribution in natural language generation models: Information Access in the Era of Generative AI Ryen W. White, Chirag Shah, 2024-12-24 Generative Artificial Intelligence (GenAI) has emerged as a groundbreaking technology that promises to revolutionize many industries as well as people's personal and professional lives. This book discusses GenAI and its role in information access - often referred to as Generative Information Retrieval (GenIR) - or more broadly, information interaction. The role of GenAI in information access is complex and dynamic, with many dimensions. To address this, following a brief introduction to GenAI and GenIR, the remainder of the book provides eight chapters, each targeting a different dimension or sub-topic. These cover foundations of GenIR, interactions with GenIR systems, adapting them to users, tasks, and scenarios, improving them based on user feedback, GenIR evaluation, the sociotechnical implications of GenAI for information access, recommendations within GenIR, and the future of information access with GenIR. The book is targeted at graduate students and researchers interested in issues of information retrieval, access, and interactions, as well as applications of GenAI in various informational contexts. While some of the parts assume prior background in IR or AI, most others do not, making this book suitable for adoption in various classes as a primary source or as a supplementary material.

measuring attribution in natural language generation models: Advances in Information Retrieval Nazli Goharian, Nicola Tonellotto, Yulan He, Aldo Lipani, Graham McDonald, Craig Macdonald, Iadh Ounis, 2024-03-19 The six-volume set LNCS 14608, 14609, 14609, 14610, 14611, 14612 and 14613 constitutes the refereed proceedings of the 46th European Conference on IR Research, ECIR 2024, held in Glasgow, UK, during March 24-28, 2024. The 57 full papers, 18 finding papers, 36 short papers, 26 IR4Good papers, 18 demonstration papers, 9 reproducibility papers, 8 doctoral consortium papers, and 15 invited CLEF papers were carefully reviewed and selected from 578 submissions. The accepted papers cover the state of the art in information retrieval focusing on user aspects, system and foundational aspects, machine learning, applications, evaluation, new social and technical challenges, and other topics of direct or indirect relevance to search.

measuring attribution in natural language generation models: Understanding Machine Understanding Ken Clements, 2024-10-15 This is a comprehensive and thought-provoking exploration of the nature of machine understanding, its evaluation, and its implications. The book proposes a new framework, the Multifaceted Understanding Test Tool (MUTT), for assessing machine understanding across multiple dimensions, from language comprehension and logical reasoning to social intelligence and metacognition. Through a combination of philosophical analysis, technical exposition, and narrative thought experiments, the book delves into the frontiers of machine understanding, raising fundamental questions about the cognitive mechanisms and representations that enable genuine understanding in both human and machine minds. By probing the boundaries of artificial comprehension, the book aims to advance our theoretical grasp on the elusive notion of understanding and inform responsible development and deployment of AI technologies. In an era where Artificial Intelligence systems are becoming integral to our daily lives,

a pressing question arises: Do these machines truly understand what they are doing, or are they merely sophisticated pattern matchers? Understanding Machine Understanding delves into this profound inquiry, exploring the depths of machine cognition and the essence of comprehension. Join Ken Clements and Claude 3 Opus on an intellectual journey that challenges conventional benchmarks like the Turing Test and introduces the innovative Multifaceted Understanding Test Tool (MUTT). This groundbreaking framework assesses AI's capabilities across language, reasoning, perception, and social intelligence, aiming to distinguish genuine understanding from mere imitation. Through philosophical analysis, technical exposition, and engaging narratives, this book invites readers to explore the frontiers of AI comprehension. Whether you're an AI researcher, philosopher, or curious observer, Understanding Machine Understanding offers a thought-provoking guide to the future of human-machine collaboration. Discover what it truly means for a machine to understand--and the implications for our shared future.

measuring attribution in natural language generation models: Artificial Neural Networks and Machine Learning - ICANN 2025 Walter Senn, Marcello Sanguineti, Ausra Saudargiene, Igor V. Tetko, Alessandro E. P. Villa, Viktor Jirsa, Yoshua Bengio, 2025-10-11 The four-volume set LNCS 16068-16071 constitutes the proceedings of the 34th International Conference on Artificial Neural Networks and Machine Learning, ICANN 2025, held in Kaunas, Lithuania, September 9-12, 2025. The 170 full papers and 8 abstracts included in these conference proceedings were carefully reviewed and selected from 375 submissions. The conference strongly values the synergy between theoretical progress and impactful real-world applications, and actively encourages contributions that demonstrate how artificial neural networks are being used to address pressing societal and technological challenges.

measuring attribution in natural language generation models: Artificial Intelligence in Education Alexandra I. Cristea, Erin Walker, Yu Lu, Olga C. Santos, Seiji Isotani, 2025-07-17 This six-volume set LNAI 15877-15882 constitutes the refereed proceedings of the 26th International Conference on Artificial Intelligence in Education, AIED 2025, held in Palermo, Italy, during July 22–26, 2025. The 130 full papers and 129 short papers presented in this book were carefully reviewed and selected from 711 submissions. The conference program comprises seven thematic tracks: Track 1: AIED Architectures and Tools Track 2: Machine Learning and Generative AI: Emphasising datadriven Track 3: Learning, Teaching, and Pedagogy Track 4: Human-Centred Design and Design-Based Research Track 5: Teaching AI Track 6: Ethics, Equity, and AIED in Society Track 7: Theoretical Aspects of AIED and AI-Based Modelling for Education

measuring attribution in natural language generation models: Information Retrieval Xiangnan He, Zhaochun Ren, Ruiming Tang, 2025-01-31 This book constitutes the refereed proceedings of the 30th China Conference on Information Retrieval, CCIR 2024, held in Wuhan, China, during October 18–20, 2024. The 11 full papers presented in this volume were carefully reviewed and selected from 26 submissions. As the flagship conference of CIPS, CCIR focuses on the development of China's internet industry and provides a broad platform for the exchange of the latest academic and technological achievements in the field of information retrieval.

Measuring attribution in natural language generation models: Explainable Natural Language Processing Anders Søgaard, 2022-06-01 This book presents a taxonomy framework and survey of methods relevant to explaining the decisions and analyzing the inner workings of Natural Language Processing (NLP) models. The book is intended to provide a snapshot of Explainable NLP, though the field continues to rapidly grow. The book is intended to be both readable by first-year M.Sc. students and interesting to an expert audience. The book opens by motivating a focus on providing a consistent taxonomy, pointing out inconsistencies and redundancies in previous taxonomies. It goes on to present (i) a taxonomy or framework for thinking about how approaches to explainable NLP relate to one another; (ii) brief surveys of each of the classes in the taxonomy, with a focus on methods that are relevant for NLP; and (iii) a discussion of the inherent limitations of some classes of methods, as well as how to best evaluate them. Finally, the book closes by providing a list of resources for further research on explainability.

measuring attribution in natural language generation models: Customer-Centric AI: Conversational Technologies, Personalization, and Ethical Innovation Zahara, Mahwish, 2025-08-13 Artificial intelligence (AI) is transforming how businesses engage with customers, with conversational technologies like chatbots and voice assistants enabling more responsive and personalized experiences. By leveraging data-driven insights, organizations can tailor interactions to individual preferences, enhancing satisfaction and loyalty. However, this shift toward hyper-personalization also raises ethical concerns related to privacy, transparency, and algorithmic bias. Addressing these challenges is essential to fostering responsible innovation that respects user autonomy while maximizing the benefits of AI. As AI continues to evolve, placing the customer at the center of technological development is key to building trust and long-term value in digital interactions. Customer-Centric AI: Conversational Technologies, Personalization, and Ethical Innovation explores the transformative impact of AI on customer engagement, focusing on how technologies are reshaping marketing, service, and personalization strategies. It addresses the ethical implications of AI-driven interactions, highlighting issues of privacy, transparency, and trust in digital environments. Covering topics such as conversational AI, customer service, and social media, this book is an excellent resource for marketing professionals, customer experience and service managers, business leaders and strategists, AI developers, data scientists, graduate and postgraduate students, policymakers, researchers, and more.

measuring attribution in natural language generation models: The Semantic Web – ISWC 2023 Terry R. Payne, Valentina Presutti, Guilin Qi, María Poveda-Villalón, Giorgos Stoilos, Laura Hollink, Zoi Kaoudi, Gong Cheng, Juanzi Li, 2023-11-01 This book constitutes the proceedings of the 22nd International Semantic Web Conference, ISWC 2023, which took place in October 2023 in Athens, Greece. The 58 full papers presented in this double volume were thoroughly reviewed and selected from 248 submissions. Many submissions focused on the use of reasoning and query answering, witha number addressing engineering, maintenance, and alignment tasks for ontologies. Likewise, there has been a healthy batch of submissions on search, query, integration, and the analysis of knowledge. Finally, following the growing interest in neuro-symbolic approaches, there has been a rise in the number of studies that focus on the use of Large Language Models and Deep Learning techniques such as Graph Neural Networks.

measuring attribution in natural language generation models: Legal Knowledge and **Information Systems** E. Francesconi, G. Borges, C. Sorge, 2022-12-21 In recent years, interest within the research community and the legal industry regarding technological advances in legal knowledge representation and processing has been growing. This relates to areas such as computational models of legal reasoning, cybersecurity, privacy, trust and blockchain methods, among other things. This book presents the proceedings of JURIX 2022, the 35th International Conference on Legal Knowledge and Information Systems, held from 14-16 December in Saarbrücken, Germany, under the auspices of the Dutch Foundation for Legal Knowledge Based Systems and hosted by Saarland University. The annual JURIX conference has become an international forum for academics and professionals to exchange knowledge and experiences at the intersection of law and artificial intelligence (AI). For this edition, 62 submissions were received from 163 authors in 24 countries. Following a rigorous review process, carried out by a programme committee of 72 experts recognised in the field, 14 submissions were selected for publication as long papers, 22 as short papers and 5 as demo papers, making a total of 41 papers altogether and representing a 22.5% acceptance rate for long papers (66.1% overall). The broad array of topics covered includes argumentation and legal reasoning, legal ontologies and the semantic web, machine and deep learning and natural language processing for legal knowledge extraction, as well as argument mining, translation of legal texts, defeasible logic, legal compliance, explainable AI, alternative dispute resolution, legal drafting and smart contracts. Providing an overview of recent advances, the book will be of interest to all those working at the interface between the law and AI.

measuring attribution in natural language generation models: *Models of Reference* Kees van Deemter, Emiel Krahmer, Albert Gatt, Roger P.G. van Gompel, 2017-04-28 To communicate,

speakers need to make it clear what they are talking about. Referring expressions play a crucial part in achieving this, by anchoring utterances to things. Examples of referring expressions include noun phrases such as "this phenomenon", "it" and "the phenomenon to which this Topic is devoted". Reference is studied throughout the Cognitive Sciences (from philosophy and logic to neuro-psychology, computer science and linguistics), because it is thought to lie at the core of all of communication. Recent years have seen a new wave of work on models of referring, as witnessed by a number of recent research projects, books, and journal Special Issues. The Research Topic "Models of Reference" in Frontiers in Psychology is a new milestone, focusing on contributions from Psycholinguistics and Computational Linguistics. The articles in it are concerned with such issues as audience design, overspecification, visual perception, and variation between speakers.

measuring attribution in natural language generation models: Flair for Natural Language Processing William Smith, 2025-08-19 Flair for Natural Language Processing Flair for Natural Language Processing presents a comprehensive and authoritative guide to the design, implementation, and deployment of cutting-edge NLP solutions using the Flair framework. Addressing readers ranging from advanced practitioners to curious researchers, this book delves into the architectural foundations of Flair, exploring its modular design, extensibility, and robust workflow orchestration. Detailed expositions of data abstractions, performance optimization, and plugin integration provide all the necessary tools for building bespoke solutions tailored to research or production environments. Spanning sequence labeling, document classification, and information extraction, the text offers an in-depth analysis of modern embedding architectures—including contextual, static, character-level, and transformer-based methods. The reader is guided through advanced sequence modeling, multi-task learning, and cross-lingual adaptation, as well as practical strategies for classifying complex, noisy, and imbalanced datasets. A particular emphasis is placed on constructing, evaluating, and optimizing custom NLP pipelines, with concrete best practices for benchmarking, diagnostics, and explainability. Recognizing the critical importance of scalability, ethical governance, and operational excellence, the book covers every facet of deploying NLP systems at scale—from distributed training and cloud-native deployment, to security, privacy, and responsible AI considerations. It concludes with a forward-looking exploration of emerging trends such as large language model integration, green NLP, and interactive human-in-the-loop systems. For those seeking a rigorous yet accessible resource on the application of contemporary NLP, Flair for Natural Language Processing stands as a definitive reference.

measuring attribution in natural language generation models: Explainable Artificial Intelligence Luca Longo, Sebastian Lapuschkin, Christin Seifert, 2024-07-09 This four-volume set constitutes the refereed proceedings of the Second World Conference on Explainable Artificial Intelligence, xAI 2024, held in Valletta, Malta, during July 17-19, 2024. The 95 full papers presented were carefully reviewed and selected from 204 submissions. The conference papers are organized in topical sections on: Part I - intrinsically interpretable XAI and concept-based global explainability; generative explainable AI and verifiability; notion, metrics, evaluation and benchmarking for XAI. Part II - XAI for graphs and computer vision; logic, reasoning, and rule-based explainable AI; model-agnostic and statistical methods for eXplainable AI. Part III - counterfactual explanations and causality for eXplainable AI; fairness, trust, privacy, security, accountability and actionability in eXplainable AI. Part IV - explainable AI in healthcare and computational neuroscience; explainable AI for improved human-computer interaction and software engineering for explainability; applications of explainable artificial intelligence.

measuring attribution in natural language generation models: Neural Information Processing Mufti Mahmud, Maryam Doborjeh, Kevin Wong, Andrew Chi Sing Leung, Zohreh Doborjeh, M. Tanveer, 2025-07-25 The eleven-volume set LNCS 15286-15296 constitutes the refereed proceedings of the 31st International Conference on Neural Information Processing, ICONIP 2024, held in Auckland, New Zealand, in December 2024. The 318 regular papers presented in the proceedings set were carefully reviewed and selected from 1301 submissions. They focus on four main areas, namely: theory and algorithms; cognitive neurosciences; human-centered

computing; and applications.

measuring attribution in natural language generation models: Interpretability and Explainability in AI Using Python Aruna Chakkirala, 2025-04-15 TAGLINE Demystify AI Decisions and Master Interpretability and Explainability Today KEY FEATURES

Master Interpretability and Explainability in ML, Deep Learning, Transformers, and LLMs

Implement XAI techniques using Python for model transparency • Learn global and local interpretability with real-world examples DESCRIPTION Interpretability in AI/ML refers to the ability to understand and explain how a model arrives at its predictions. It ensures that humans can follow the model's reasoning, making it easier to debug, validate, and trust. Interpretability and Explainability in AI Using Python takes you on a structured journey through interpretability and explainability techniques for both white-box and black-box models. You'll start with foundational concepts in interpretable machine learning, exploring different model types and their transparency levels. As you progress, you'll dive into post-hoc methods, feature effect analysis, anchors, and counterfactuals—powerful tools to decode complex models. The book also covers explainability in deep learning, including Neural Networks, Transformers, and Large Language Models (LLMs), equipping you with strategies to uncover decision-making patterns in AI systems. Through hands-on Python examples, you'll learn how to apply these techniques in real-world scenarios. By the end, you'll be well-versed in choosing the right interpretability methods, implementing them efficiently, and ensuring AI models align with ethical and regulatory standards—giving you a competitive edge in the evolving AI landscape. WHAT WILL YOU LEARN ● Dissect key factors influencing model interpretability and its different types. ● Apply post-hoc and inherent techniques to enhance AI transparency. ● Build explainable AI (XAI) solutions using Python frameworks for different models. • Implement explainability methods for deep learning at global and local levels. • Explore cutting-edge research on transparency in transformers and LLMs. • Learn the role of XAI in Responsible AI, including key tools and methods. WHO IS THIS BOOK FOR? This book is tailored for Machine Learning Engineers, AI Engineers, and Data Scientists working on AI applications. It also serves as a valuable resource for professionals and students in AI-related fields looking to enhance their expertise in model interpretability and explainability techniques. TABLE OF CONTENTS 1. Interpreting Interpretable Machine Learning 2. Model Types and Interpretability Techniques 3. Interpretability Taxonomy and Techniques 4. Feature Effects Analysis with Plots 5. Post-Hoc Methods 6. Anchors and Counterfactuals 7. Interpretability in Neural Networks 8. Explainable Neural Networks 9. Explainability in Transformers and Large Language Models 10. Explainability and Responsible AI Index

measuring attribution in natural language generation models: PROCEEDINGS OF THE 24TH CONFERENCE ON FORMAL METHODS IN COMPUTER-AIDED DESIGN - FMCAD 2024 Nina Narodytska, Philipp Rümmer, 2024-10-01 Die Proceedings zur Konferenz "Formal Methods in Computer-Aided Design 2024" geben aktuelle Einblicke in ein spannendes Forschungsfeld. Zum fünften Mal erscheinen die Beiträge der Konferenzreihe "Formal Methods in Computer-Aided Design" (FMCAD) als Konferenzband bei TU Wien Academic Press. Der aktuelle Band der seit 2006 jährlich veranstalteten Konferenzreihe präsentiert in 35 Beiträgen neueste wissenschaftliche Erkenntnisse aus dem Bereich des computergestützten Entwerfens. Die Beiträge behandeln formale Aspekte des computergestützten Systemdesigns einschließlich Verifikation, Spezifikation, Synthese und Test. Die FMCAD-Konferenz findet im Oktober 2024 in Prag, Tschechische Republik, statt. Sie gilt als führendes Forum im Bereich des computer-aided design und bietet seit ihrer Gründung Forschenden sowohl aus dem akademischen als auch dem industriellen Umfeld die Möglichkeit, sich auszutauschen und zu vernetzen.

measuring attribution in natural language generation models: Deep Learning for Natural Language Processing Stephan Raaijmakers, 2022-12-06 Explore the most challenging issues of natural language processing, and learn how to solve them with cutting-edge deep learning! Deep learning has advanced natural language processing to exciting new levels and powerful new applications! For the first time, computer systems can achieve human levels of summarizing, making connections, and other tasks that require comprehension and context. Deep Learning for Natural

Language Processing reveals the groundbreaking techniques that make these innovations possible. Stephan Raaijmakers distills his extensive knowledge into useful best practices, real-world applications, and the inner workings of top NLP algorithms. Deep learning has transformed the field of natural language processing. Neural networks recognize not just words and phrases, but also patterns. Models infer meaning from context, and determine emotional tone. Powerful deep learning-based NLP models open up a goldmine of potential uses. Deep Learning for Natural Language Processing teaches you how to create advanced NLP applications using Python and the Keras deep learning library. You'll learn to use state-of the-art tools and techniques including BERT and XLNET, multitask learning, and deep memory-based NLP. Fascinating examples give you hands-on experience with a variety of real world NLP applications. Plus, the detailed code discussions show you exactly how to adapt each example to your own uses!

measuring attribution in natural language generation models: Internet of Things, Smart Spaces, and Next Generation Networks and Systems Olga Galinina, Sergey Andreev, Sergey Balandin, Yevgeni Koucheryavy, 2020-12-22 This book constitutes the joint refereed proceedings of the 20th International Conference on Next Generation Teletraffic and Wired/Wireless Advanced Networks and Systems, NEW2AN 2020, and the 13th Conference on Internet of Things and Smart Spaces, ruSMART 2020. The conference was held virtually due to the COVID-19 pandemic. The 79 revised full papers presented were carefully reviewed and selected from 225 submissions. The papers of NEW2AN address various aspects of next-generation data networks, with special attention to advanced wireless networking and applications. In particular, they deal with novel and innovative approaches to performance and efficiency analysis of 5G and beyond systems, employed game-theoretical formulations, advanced queuing theory, and stochastic geometry, while also covering the Internet of Things, cyber security, optics, signal processing, as well as business aspects. ruSMART 2020, provides a forum for academic and industrial researchers to discuss new ideasand trends in the emerging areas.

measuring attribution in natural language generation models: MEDINFO 2019: Health and Wellbeing e-Networks for All L. Ohno-Machado, B. Séroussi, 2019-11-12 Combining and integrating cross-institutional data remains a challenge for both researchers and those involved in patient care. Patient-generated data can contribute precious information to healthcare professionals by enabling monitoring under normal life conditions and also helping patients play a more active role in their own care. This book presents the proceedings of MEDINFO 2019, the 17th World Congress on Medical and Health Informatics, held in Lyon, France, from 25 to 30 August 2019. The theme of this year's conference was 'Health and Wellbeing: E-Networks for All', stressing the increasing importance of networks in healthcare on the one hand, and the patient-centered perspective on the other. Over 1100 manuscripts were submitted to the conference and, after a thorough review process by at least three reviewers and assessment by a scientific program committee member, 285 papers and 296 posters were accepted, together with 47 podium abstracts, 7 demonstrations, 45 panels, 21 workshops and 9 tutorials. All accepted paper and poster contributions are included in these proceedings. The papers are grouped under four thematic tracks: interpreting health and biomedical data, supporting care delivery, enabling precision medicine and public health, and the human element in medical informatics. The posters are divided into the same four groups. The book presents an overview of state-of-the-art informatics projects from multiple regions of the world; it will be of interest to anyone working in the field of medical informatics.

Related to measuring attribution in natural language generation models

MEASURING | **English meaning - Cambridge Dictionary** MEASURING definition: 1. present participle of measure 2. to discover the exact size or amount of something: 3. to be a. Learn more **MEASURING Definition & Meaning - Merriam-Webster** The meaning of MEASURE is an adequate or due portion. How to use measure in a sentence

Measurement | Definition, Types, Instruments, & Facts | Britannica Measurement is fundamental to the sciences; to engineering, construction, and other technical fields; and to almost all everyday activities. For that reason the elements, conditions,

Measurement - Wikipedia The use of the word measure, in the sense of a measuring instrument, only survives in the phrase tape measure, an instrument that can be used to measure but cannot be used to draw straight

Units of Measurement - List, Chart, Length, Mass, Examples In this article, we shall explore the concept of metric and imperial units of measurement. We will also discuss the various measurement units used for measuring length, mass, time,

MEASURING definition and meaning | Collins English Dictionary If possible, invest in some proper measuring spoons - a teaspoon and tablespoon are most commonly needed

Measuring - definition of measuring by The Free Dictionary e. A device used for measuring. f. The act of measuring: By measure the picture was four feet tall. 2. An evaluation or a basis of comparison: "the final measure of the worth of a society" (Joseph

measuring - Dictionary of English Also, measure off, to mark off or deal out by measuring: [\sim + out/off + object] to measure out a cup of flour. [\sim + object + out/off] He measured it out and handed it to her

Measuring - Definition, Meaning & Synonyms | Whether you're a teacher or a learner, Vocabulary.com can put you or your class on the path to systematic vocabulary improvement measuring: Explore its Definition & Usage | RedKiwi Words 'Measuring' means ascertaining the size, amount, or degree of something by using an instrument or device marked in standard units or by comparing it with an object of known size, or taking

MEASURING | English meaning - Cambridge Dictionary MEASURING definition: 1. present participle of measure 2. to discover the exact size or amount of something: 3. to be a. Learn more **MEASURING Definition & Meaning - Merriam-Webster** The meaning of MEASURE is an adequate or due portion. How to use measure in a sentence

Measurement | Definition, Types, Instruments, & Facts | Britannica Measurement is fundamental to the sciences; to engineering, construction, and other technical fields; and to almost all everyday activities. For that reason the elements, conditions,

Measurement - Wikipedia The use of the word measure, in the sense of a measuring instrument, only survives in the phrase tape measure, an instrument that can be used to measure but cannot be used to draw straight

Units of Measurement - List, Chart, Length, Mass, Examples In this article, we shall explore the concept of metric and imperial units of measurement. We will also discuss the various measurement units used for measuring length, mass, time,

MEASURING definition and meaning | Collins English Dictionary If possible, invest in some proper measuring spoons - a teaspoon and tablespoon are most commonly needed

Measuring - definition of measuring by The Free Dictionary e. A device used for measuring. f. The act of measuring: By measure the picture was four feet tall. 2. An evaluation or a basis of comparison: "the final measure of the worth of a society" (Joseph

measuring - Dictionary of English Also, measure off, to mark off or deal out by measuring: [\sim + out/off + object] to measure out a cup of flour. [\sim + object + out/off] He measured it out and handed it to her

Measuring - Definition, Meaning & Synonyms | Whether you're a teacher or a learner, Vocabulary.com can put you or your class on the path to systematic vocabulary improvement measuring: Explore its Definition & Usage | RedKiwi Words 'Measuring' means ascertaining the size, amount, or degree of something by using an instrument or device marked in standard units or by comparing it with an object of known size, or taking

MEASURING | **English meaning - Cambridge Dictionary** MEASURING definition: 1. present participle of measure 2. to discover the exact size or amount of something: 3. to be a. Learn more **MEASURING Definition & Meaning - Merriam-Webster** The meaning of MEASURE is an

adequate or due portion. How to use measure in a sentence

Measurement | Definition, Types, Instruments, & Facts | Britannica Measurement is fundamental to the sciences; to engineering, construction, and other technical fields; and to almost all everyday activities. For that reason the elements, conditions,

Measurement - Wikipedia The use of the word measure, in the sense of a measuring instrument, only survives in the phrase tape measure, an instrument that can be used to measure but cannot be used to draw straight

Units of Measurement - List, Chart, Length, Mass, Examples In this article, we shall explore the concept of metric and imperial units of measurement. We will also discuss the various measurement units used for measuring length, mass, time,

MEASURING definition and meaning | Collins English Dictionary If possible, invest in some proper measuring spoons - a teaspoon and tablespoon are most commonly needed

Measuring - definition of measuring by The Free Dictionary e. A device used for measuring. f. The act of measuring: By measure the picture was four feet tall. 2. An evaluation or a basis of comparison: "the final measure of the worth of a society"

measuring - Dictionary of English Also, measure off, to mark off or deal out by measuring: [\sim + out/off + object] to measure out a cup of flour. [\sim + object + out/off] He measured it out and handed it to her

Measuring - Definition, Meaning & Synonyms | Whether you're a teacher or a learner, Vocabulary.com can put you or your class on the path to systematic vocabulary improvement measuring: Explore its Definition & Usage | RedKiwi Words 'Measuring' means ascertaining the size, amount, or degree of something by using an instrument or device marked in standard units or by comparing it with an object of known size, or taking

MEASURING | English meaning - Cambridge Dictionary MEASURING definition: 1. present participle of measure 2. to discover the exact size or amount of something: 3. to be a. Learn more **MEASURING Definition & Meaning - Merriam-Webster** The meaning of MEASURE is an adequate or due portion. How to use measure in a sentence

Measurement | Definition, Types, Instruments, & Facts | Britannica Measurement is fundamental to the sciences; to engineering, construction, and other technical fields; and to almost all everyday activities. For that reason the elements, conditions,

Measurement - Wikipedia The use of the word measure, in the sense of a measuring instrument, only survives in the phrase tape measure, an instrument that can be used to measure but cannot be used to draw straight

Units of Measurement - List, Chart, Length, Mass, Examples In this article, we shall explore the concept of metric and imperial units of measurement. We will also discuss the various measurement units used for measuring length, mass, time,

MEASURING definition and meaning | Collins English Dictionary If possible, invest in some proper measuring spoons - a teaspoon and tablespoon are most commonly needed

Measuring - definition of measuring by The Free Dictionary e. A device used for measuring. f. The act of measuring: By measure the picture was four feet tall. 2. An evaluation or a basis of comparison: "the final measure of the worth of a society" (Joseph

measuring - Dictionary of English Also, measure off, to mark off or deal out by measuring: [\sim + out/off + object] to measure out a cup of flour. [\sim + object + out/off] He measured it out and handed it to her

Measuring - Definition, Meaning & Synonyms | Whether you're a teacher or a learner, Vocabulary.com can put you or your class on the path to systematic vocabulary improvement measuring: Explore its Definition & Usage | RedKiwi Words 'Measuring' means ascertaining the size, amount, or degree of something by using an instrument or device marked in standard units or by comparing it with an object of known size, or taking

MEASURING | English meaning - Cambridge Dictionary MEASURING definition: 1. present participle of measure 2. to discover the exact size or amount of something: 3. to be a. Learn more

MEASURING Definition & Meaning - Merriam-Webster The meaning of MEASURE is an adequate or due portion. How to use measure in a sentence

Measurement | Definition, Types, Instruments, & Facts | Britannica Measurement is fundamental to the sciences; to engineering, construction, and other technical fields; and to almost all everyday activities. For that reason the elements, conditions,

Measurement - Wikipedia The use of the word measure, in the sense of a measuring instrument, only survives in the phrase tape measure, an instrument that can be used to measure but cannot be used to draw straight

Units of Measurement - List, Chart, Length, Mass, Examples In this article, we shall explore the concept of metric and imperial units of measurement. We will also discuss the various measurement units used for measuring length, mass, time,

MEASURING definition and meaning | Collins English Dictionary If possible, invest in some proper measuring spoons - a teaspoon and tablespoon are most commonly needed

Measuring - definition of measuring by The Free Dictionary e. A device used for measuring. f. The act of measuring: By measure the picture was four feet tall. 2. An evaluation or a basis of comparison: "the final measure of the worth of a society" (Joseph

measuring - Dictionary of English Also, measure off, to mark off or deal out by measuring: [\sim + out/off + object] to measure out a cup of flour. [\sim + object + out/off] He measured it out and handed it to her

Measuring - Definition, Meaning & Synonyms | Whether you're a teacher or a learner, Vocabulary.com can put you or your class on the path to systematic vocabulary improvement measuring: Explore its Definition & Usage | RedKiwi Words 'Measuring' means ascertaining the size, amount, or degree of something by using an instrument or device marked in standard units or by comparing it with an object of known size, or taking

Back to Home: https://espanol.centerforautism.com