




Integrative urban AI to expand coverage, access, and equity of urban data

Bill Howe^{1,a}, Jackson Maxfield Brown^{1,b}, Bin Han^{1,c}, Bernease Herman^{1,d}, Nic Weber^{1,e}, An Yan^{2,f}, Sean Yang^{1,g}, and Yiwei Yang^{1,h}

¹ University of Washington, Seattle, USA

² Meta, Seattle, USA

Received 22 November 2021 / Accepted 31 January 2022 / Published online 9 April 2022
© The Author(s), under exclusive licence to EDP Sciences, Springer-Verlag GmbH Germany, part of Springer Nature 2022

Abstract We consider the use of AI techniques to expand the coverage, access, and equity of urban data. We aim to enable holistic research on city dynamics, steering AI research attention away from profit-oriented, societally harmful applications (e.g., facial recognition) and toward foundational questions in mobility, participatory governance, and justice. By making available high-quality, multi-variate, cross-scale data for research, we aim to link the macrostudy of cities as complex systems with the reductionist view of cities as an assembly of independent prediction tasks. We identify four research areas in AI for cities as key enablers: interpolation and extrapolation of spatiotemporal data, using NLP techniques to model speech- and text-intensive governance activities, exploiting ontology modeling in learning tasks, and understanding the interaction of fairness and interpretability in sensitive contexts.

1 Introduction

Cities are complex systems: collections of interacting agents that exhibit non-trivial collective behavior [2, 19]. This observation has guided research in general principles of city planning that can govern the behavior of the complex adaptive system the city manifests. Early work by Jacobs proposed ideal sizes and specific guidelines for city neighborhoods [25], and more recently, researchers have begun to empirically validate these ideas using mobile phone data [46]. West and Kempes model scaling behavior for cities as balancing sublinear growth in resource consumption (as a function of population) against the superlinear growth of socioeconomic effects, both positive (per capita wages) and negative (inequity, disease) as a power law with an exponent of about 0.15 (as opposed to, for example, the exponent of 0.25 identified in many biological processes) [28, 67]. More recently, the rate of COVID-19 spread has been shown to approximate the same power law [57]. As West and Kempes argue “Cities are machines we evolved to facilitate, accelerate, amplify, and densify social interactions.”

The holistic study of cities as complex systems complements the rapid (yet ultimately opportunistic) proliferation of artificial intelligence technology in the public sector. Although conventional machine learning techniques are common in urban applications [44, 65, 71], neural architectures are opening new opportunities by adapting convolutional, recurrent, and transformer architectures to spatiotemporal data [35, 38, 56, 70, 73, 75, 76]; see Grekousis 2020 for a recent survey [17].

These two lines of inquiry—top-down modeling of cities as complex systems and bottom-up modeling of specific urban systems using deep learning—are difficult to reconcile. Complex systems are not amenable to reductionist statistical experiments: comparing the results of an agent-based model with observed data (e.g., for autonomous vehicle research [11]) is often the best we can do, despite the challenges of addressing the inverse problems implied [8, 59]. The central issue is that observational micro-data for cities are inconsistent in availability and quality, limiting the opportunity for validation of sophisticated models.

This inconsistency persists despite significant investments in open data. Over the last 2 decades, cities have increasingly released datasets publicly on the web, proactively, in response to transparency regulation. For example, in the US, all 50 states and the District of Columbia have passed some version of the federal Freedom of Information (FOI) Act. While this first wave of open data was driven by FOI laws and made national government data available primarily to journalists, lawyers, and activists, a second wave of open

^a e-mail: billhowe@uw.edu (corresponding author)

^b e-mail: jmxbrown@uw.edu

^c e-mail: bh193@uw.edu

^d e-mail: bernease@uw.edu

^e e-mail: nmweber@uw.edu

^f e-mail: annieyan9033@fb.com

^g e-mail: seanyang38@uw.edu

^h e-mail: yanyiwei@uw.edu

data, enabled by the advent of open source and web 2.0 technologies, was characterized by an attempt to make data “open by default” to civic technologists, government agencies, and corporations [64]. While open data have indeed made significant data assets available online, their uptake and use have been weaker than anticipated [64], an effect may attribute to inconsistent availability of high-value data across cities [32]. Ultimately, open data exhibit convenience sampling effects.

In this paper, we consider four research thrusts all aimed at using AI techniques to improve the coverage, access, and equity of urban data, and thereby reduce barriers and attract attention to the study of critical questions in city dynamics and socioeconomic interactions. Machine learning research is broadly recognized to be too narrow in applications and datasets, focusing on opportunistic, discriminatory, and profit-oriented applications [50, 53, 68]. By making high-quality urban data available across cities, across variables, across time-scales, and at multiple resolutions, we aim to make AI research on societally important problems the path of least resistance. However, to accomplish this long-term goal, we need to address specific challenges in working with urban data.

Expanding existing sources By simultaneously modeling multiple heterogeneous datasets [69], we aim to identify the underlying relationships and interactions between urban systems as a middle path between reductionist, application-specific prediction tasks and holistic, simulation-oriented inference. However, where our earlier work assumed uniform data coverage, we now need to apply advanced learning techniques to interpolate and extrapolate dense spatiotemporal datasets to account for inconsistent coverage (Sect. 2). These techniques help expand the utility and reach of data-hungry predictive models to counteract the sparse and inconsistent availability of public and private data. As an example, we show how the interpolation of urban transportation data is remarkably amenable to deep learning architectures developed for image inpainting on the web.

Developing new sources Open data in urban contexts are typically either spatiotemporal (vector or raster) or administrative (structured). However, by investing in infrastructure, we can develop and make available data sources around governance, economics, decision-making, and public participation. As an example, we show how transcripts from public meetings are amenable to computational processing to increase oversight and participation, if we can first establish an infrastructure and appropriate standards to collect and manage this data (Sect. 5).

Exploiting rich ontologies The use of large, noisy, and heterogeneous data motivates investment in data curation: associating contextual information with the data to mediate its collection and use. However, manual curation activities (e.g., human labeling of data) scale poorly. In complex domains, human expertise is better invested developing richer labeling schemes than actu-

ally labeling data. For example, ontologies have been developed for electric mobility [55], humanitarian [3], and smart city applications [1, 9, 13, 18, 61]. However, categorizing public data (e.g., social media posts) using these ontologies requires new techniques in hierarchical multi-label classification. As an example, we show how graph encoding techniques can be used to significantly improve performance in these contexts (Sect. 4).

Incorporating fairness and interpretability In every application of urban machine learning, prediction and modeling carries enormous risk of exacerbating inequity and opacity [21, 42, 43, 49]. Building on recent advances in fair and explainable AI, we consider the interactions between accuracy, fairness, and explainability in urban applications. We then propose new methods for controlling these tradeoffs in response to emerging regulation (Sect. 3).

2 Interpolation of spatiotemporal data using deep learning

Image inpainting is a task of synthesizing missing pixels in images. In computer vision, there are two board branches to inpaint images. The first branch contains diffusion-based or patch-based methods that utilize low-level image features to reconstruct the missing regions. The second branch contains learning-based methods that involve the training of deep learning models. Traditional diffusion-based methods transfer information from the valid regions to the missing regions, which are convenient to apply but limited to small missing regions only. Learning-based approaches aim to recover the images based on the patterns learned from large amount of training data. Such methods include context encoder by Pathak et al. [47], global and local image inpainting by Lizuka et al. [22], partial convolution method by Liu et al. [36], etc.

Image inpainting techniques have wide application potentials, including the geospatial domain that works frequently with satellite images. Zhang et al. [77] proposed a unified spatial–temporal–spectral deep convolutional neural network (CNN) image inpainting architecture to recover information obscured by poor atmospheric conditions in satellite images. Kang et al. [27] modified the architecture from [72] to restore the missing patterns of sea surface temperature from satellite images. Tasnim and Mondal [60] also applied the inpainting architecture from [72] to remove redundancies in satellite images and restore the imagery.

We build on prior work from our group in learning fair integrations of heterogeneous urban data [69]. We originally assumed uniform spatial and temporal coverage data, but in practice, urban datasets are spatially imbalanced: one neighborhood may be missing a variable of interest defined everywhere else, undermining trust in the results. Conventional statistical approaches to impute missing data, such as global/local mean imputing, interpolations, and spatial regression models, are limited in their ability to capture non-linear

interactions, where deep learning methods, including image inpainting techniques in geospatial imputation, excel.

Given the similar nature of images and gridded urban data, we conjecture that image inpainting techniques can be adapted to impute missing urban data, improving coverage and quality, and therefore usability. As far as we know, no prior work that has exploited image inpainting techniques to reconstruct missing values in raster urban data. In this section, we present our preliminary experiments and results of utilizing an image inpainting technique to compute missing values in gridded urban data.

2.1 Example: interpolating urban mobility data

We use taxi trip data as a representative example of urban data, though the coverage of urban data is much broader. We used NYC taxi trip data from 2011 to 2016 from NYC Open Data Portal [10]. The years 2011–2014 cover the trips throughout the entire year, while 2015 and 2016 only cover half of the year. The raw data are collected tabular format, where each record/row contains the information of each taxi trip, including the longitude and latitude of the starting location. We considered the demand prediction problem, interpreting each record as an indicator of demand following Mooney et al. [43]. We processed the tabular data into raster format given the following steps:

- We defined a rectangular subset of the greater metropolitan area of New York City representing lower Manhattan. We only consider the taxi trips that began within this rectangular region.
- We imposed a 32×32 grid over our selected region. This choice of dimensions is somewhat arbitrary, balancing fidelity (reducing the need to upsample or downsample datasets too much), computational efficiency, and interpretability (1 grid cell is approximately 1 km^2 .) For each year, each unique date, and each unique hour, we count how many taxi trips are within each grid and interpret these values as an estimate of taxi demand in that cell, at that time.
- In total, we have 32,616 samples to model with, each having 32×32 dimension. 70% of samples (23,482) are used as training data, 10% (2610) as validation set and the rest 20% (6524) as test data.

2.2 Modeling and results

We implemented the architecture from Liu et al. [36]. Many prior works in image inpainting only considered rectangular-shaped missing regions, but rarely are the patterns of missing data so regular. In urban data, the missing regions could be scattered or in irregular shapes corresponding to irregular political boundaries or inconsistent data collection. Therefore, Liu et al.'s work fits well into the urban scheme. Liu et al. used the summation of four different losses as the objective function, to account for different factors related to the percep-

tion of the resulting image, which was appropriate for web images but less appropriate for quantitative urban data. We only used the ℓ_1 regularization loss between the original data and inpainted data. The model hyperparameters are set to be consistent with Liu's work. The learning rate is set to $1e^{-4}$ flat. The batch size is set to 32. The maximum iteration is 10,000 and we evaluated the model on validation set every 100 iterations.

Five inpainted examples are visually presented in Fig. 1. We can see that the inpainting technique can be naturally applied to gridded urban data and yield promising results. Imputing the missing values in urban setting could also be viewed as a type of synthesis. The synthesis of partial urban data could improve the applicability and usability of urban data, but will require future work in multiple areas:

- Though deep learning methods are powerful, we need rigorous evaluation against traditional imputation techniques to see if these complex methods are warranted. Additionally, visual similarities are subjective, which is appropriate for web images but not if we intend to use these datasets for quantitative analysis. Additional quantitative measurements should be incorporated.
- The region of the experiment (NYC) and the dimensions of the urban grid (32×32) are both limited. Expanding the region to cover more area would capture more urban dynamics, while evaluating the effect of different grid sizes, will be necessary to test generalizability.
- We treat each date and each hour as a unique sample. However, in reality, the current hour timestamp is closely related to the demand from the previous hour. Modeling each sample separately ignores such dependency. Therefore, we hypothesize that modifying the architecture to work with temporal blocks would help improve performances.

3 Trade-off among distributive and procedural fairness

Real-world datasets often contain societal biases, which are perpetuated in the machine learning models, leading to discriminatory decisions in high-stake domains. In response, many methods were developed to mitigate fairness by achieving some statistical measure of equity between majority and minority groups (e.g., equalized odds and equality of opportunity) [4, 7, 30, 39, 74]. This line of work is guided primarily by the notion of distributive fairness, which emphasizes on a fair allocation of resources.

However, prior work has shown that *procedural* fairness, the perceived fairness of the process that leads to the outcome, is equally as important as distributive fairness [5, 63]. For example, in court systems, studies have shown that “most people care more about procedural fairness ... than they do about winning or

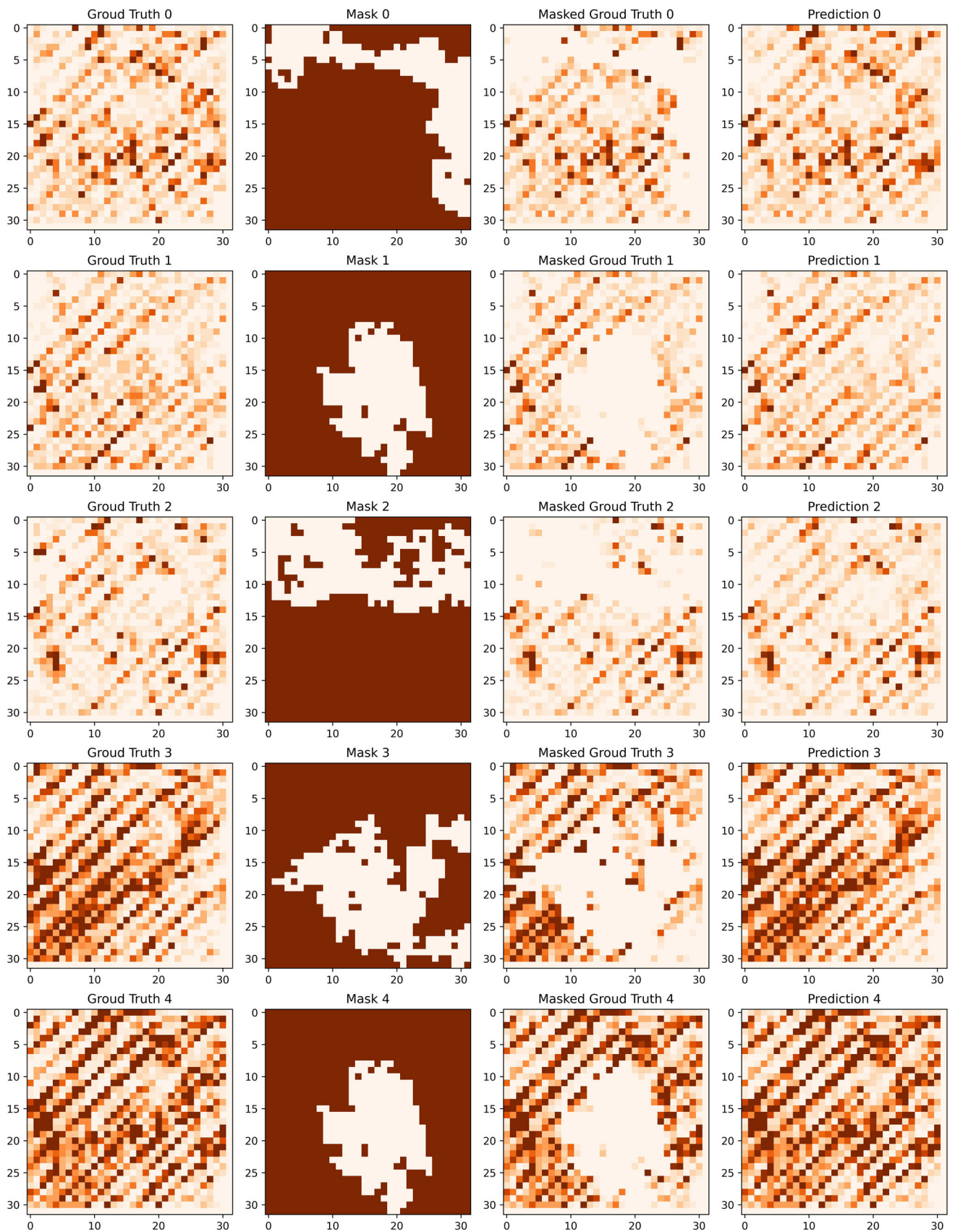


Fig. 1 Inpainting results of taxi trip data. From left to right, the columns are: ground truth images; the irregular masks; the masked ground truth; the final inpainting results

losing the particular case.” [63] Recent studies have also shown that procedural fairness is critical to automated decision systems [33, 40]. For instance, through a cross-sectional survey study at a large German university, Marsinkowski et al. found that both distributive and procedural fairness have significant implications on higher education admission that uses an automated decision system [40].

The interaction between distributive fairness, interpretability, and procedural fairness are rapidly becoming a compliance issue. In April 2021, the EU released a proposal for sweeping regulation of algorithmic bias [14]. In the same week, the Federal Trade Commission released a blog post [26] that described a legal framework for evaluating AI bias, foreshadowing enforcement. In the California, a bill regulating automated decision systems is in committee [23].

Drawing from procedural fairness theory, we propose Explanation Loss (see Eq. 1), a novel fairness metric that measures procedural fairness and a method to optimize for it [34]. In particular, this metric measures the neutrality of the decision process to different demographic groups. Since complex black-box models (e.g., deep neural networks and tree-based ensemble models) are often used due to their high predictive power, we use interpretability methods to generate explanations that reveal the model decision process for each datum. The metric then computes the average absolute differences of the explanations between all possible pairs of input samples, one from the minority group and one from the majority group. The intuition is that the difference in the model’s explanation for two groups can be approximated by the average differences of all pairs of individual explanations. Therefore, Explanation Loss measures how far away the decision process is from being perfectly neutral. In the following sections, we describe the data we used, the method that optimizes for the metric, and the preliminary results we obtained.

3.1 Data

We used the COMPAS dataset [31] for the preliminary study. COMPAS dataset, which contains attributes of criminal defendants, is often being used to study (deeply flawed) recidivism models: whether a person will reoffend within 2 years). It is known to exhibit severe biases against minority groups. Specifically, studies have shown that models trained on COMPAS tend to overpredict recidivism for black defendants and underpredict recidivism for white defendants [31].

We preprocessed the dataset following Rieger et al. [52]. The dataset contains a total of 7214 samples. We filtered 1042 due to missing information about the recidivism. We categorized age into under 25, between 25 and 45, inclusively, and above 45. We categorized sex into Male and Female. We also categorized the crime description based on matching words, resulting in categories Possession (of drugs), Driving, Violence, Theft, and No Charge. For example, descriptions that are matched with “theft” or “burglary” are categorized as Theft. We then one-hot encoded all categorical vari-

ables, and used the numeric variables as is. We focused on equalizing explanations of the Black and Caucasian records, since these two are the predominant groups of the data.

We split the data into train, validation, and test sets, with a ratio of 80/10/10.

3.2 Method

Interpretability techniques aim to generate explanations for a model’s individual predictions. A popular class of such techniques is known as feature attribution, which, given an input, the model, and prediction, assigns a number to each feature of the input to represent how much it contributes to the prediction [37, 45, 51, 58]. There are two reasons that feature attribution methods are appropriate for our study: (1) they allow us to compare model’s explanation for each prediction at the feature level, which is especially important for fairness, since certain features are more sensitive than others; (2) feature attribution vectors can be interpreted as attribution priors to incorporate the notion of procedural fairness in the model.

We propose the following regularization to achieve procedural fairness:

$$\hat{\theta} = \arg \min_{\theta} \sum_{x_i, y_i \in D} \mathcal{L}(f_{\theta}(x_i), y_i) + \lambda \frac{1}{|D_{s1}| |D_{s2}|} \sum_{x_j \in D_{s1}, x_k \in D_{s2}} |expl(x_j) - expl(x_k)|. \quad (1)$$

The regularization computes the average L1 norm of difference between explanations of every pair of instances (one from each group), $x \in D_{s1}$, $x' \in D_{s2}$. Each explanation, $expl(x)$, is a vector of feature importance scores with dimension equal to the number of features of the input. This vector is generated using any feature attribution method. In this study, we used Contextual Decomposition (CD) as the feature attribution method [45]. This regularization term takes the exact form of our proposed metric for procedural fairness.

We trained a simple multi-layer neural network model with one hidden fully connected layer of 100 neurons and ReLU activation, and varied a weight for the regularization term of 0 (no explanation loss), 0.2, 0.4, 0.6, 0.8, and 1.0. The model was trained with a batch size of 256 and a learning rate of 0.001 for 5 seeds, and the average results were reported. While the regularization equation refers to explanations of the entire dataset, in practice, this term is computed per batch for faster convergence. Specifically, for each batch, we partition the instances into two groups, then for every possible pair (one from each group), we compute the L1 norm of the difference of the feature attributions, and finally, we average the differences. An ablation study of the effect of batch size on results remains future work.

3.3 Metrics

In addition to our proposed metric of procedural fairness (Explanation Loss), the metrics we used to evaluate the model include 1) accuracy and 2) fairness. We considered two popular fairness metrics: equality of opportunity [20] and equalized odds [54]. A model is said to satisfy equality of opportunity if the false-positive rates are equivalent across two demographic groups. Similarly, equalized odds require false-negative rates to be equivalent across two groups in addition to false-positive rates.

The loss term represents the notion of fairness *distance*: how far away the model is from perfectly fair. The fairness distance we consider is based on equality of opportunity [], and measures the absolute difference between the false-positive rates of one demographic group (FPR1) and another (FPR2): $|FPR1 - FPR2|$. On the other hand, fairness distance is based on equalized odds measures [] $|FPR1 - FPR2| + |FNR1 - FNR2|$, which adds an additional absolute difference between the false-negative rates.

3.4 Results

The results are summarized in Table 1. From the first two columns of the table, we can see that the regularization effectively encourages the model to predict with similar explanations across two demographic groups, which we interpret as improved procedural fairness by penalizing the tendency for a model to essentially learn two separate submodels, one for each group. Second, adding the regularization term does not reduce the accuracy of the model. Third, equalizing the explanations has a minor effect on fairness of the outcomes, causing a slight increase on fairness distances.

4 Hierarchical multi-label classification

We demonstrate hierarchical multi-label classification (HMC) in the urban domain. HMC tasks involve a large set of labels organized into parent-child relationships, typically representing increasing specificity or isa relationships. Each input record is associated with multiple labels in the hierarchy, representing the uncertainty

Table 1 The effect of equal explanations on accuracy and fairness distances of the model

| Reg rate | Explanation loss | Accuracy | Equality of opportunity | Equalized odds |
|----------|------------------|----------|-------------------------|----------------|
| 0 | 1.68 | 70.40 | 0.19 | 0.50 |
| 0.2 | 0.05 | 70.11 | 0.22 | 0.56 |
| 0.4 | 0.03 | 70.16 | 0.22 | 0.58 |
| 0.6 | 0.02 | 70.90 | 0.21 | 0.56 |
| 0.8 | 0.01 | 70.65 | 0.23 | 0.56 |
| 1 | 0.02 | 69.70 | 0.24 | 0.60 |

associated with a large label space in a complex domain. HMC has received increasing attention with the adoption of neural networks [15, 16, 66, 78], often in contexts requiring significant human expertise, making large-scale labeling exercise infeasibly expensive. In other words, human expertise is invested in modeling the world through a complex ontology rather than labeling data using that ontology. As a result, the machine learning tasks represent a different set of requirements: the number of labels can be large relative to the number of labeled items, but there is structure among the labels that algorithms can exploit for distance supervision.

Ontology development is common in urban planning, where the complexity of the domain and multiplicity of perspectives require building consensus around a universe of discourse. For example, ontologies have been developed by teams of experts to describe electric mobility [55], humanitarian efforts [3], and smart city applications [1, 9, 13, 18, 61]¹. The HMC literature rarely considers these urban applications, instead favoring biological and scientific domains where public data are more readily available.

We are exploring new approaches for HMC that involve learning reusable representations of the ontology itself (using graph encoding techniques) to tame the complexity, then using these learned representations as the labels when training a classifier. We show that using these ontologies as a source of supervision can significantly improve the classification performance over other HMC techniques, motivating greater investment in developing comprehensive ontologies to represent the complex urban domain as a whole rather than expending resources on creating expensive labeled datasets for myriad specific applications.

4.1 Case study: community listener

We worked with a local non-profit organization to identify the community needs from several sources of the data, such as social media (Twitter, Reddit, and Facebook conversations) and long-form survey responses. We classify these discourses into the Sustainable Development Goals Ontology (SDG) [3]² and the Social Progress Index (SPI)³. The data and the predicted labels are then aggregated and visualized on an online dashboard serving policymakers and entrepreneurs. The Sustainable Development Goals Interface Ontology (SDG) was developed by United National Environment Programme to support the achievement of the 17 United National Sustainable Development Goals to promote human rights and equity. The ontology includes

¹ <https://techcommunity.microsoft.com/t5/internet-of-things/smart-cities-ontology-for-digital-twins/ba-p/2166585>.

² <https://www.unep.org/explore-topics/sustainable-development-goals/what-we-do/monitoring-progress/sdg-interface-ontology>.

³ <https://www.socialprogress.org/2020-Social-Progress-Index-Methodology.pdf>.

169 nodes with 3 levels. **Social Progress Index (SPI)** was introduced by Social Progress Imperative to promote improvement and actions for social progress. They define Social Progress as “the capacity of a society to meet the basic human needs of its citizens, establish the building blocks that allow citizens and communities to enhance and sustain the quality of their lives, and create the conditions for all individuals to reach their full potential.” SPI includes three levels with 124 nodes.

4.2 Experimental settings

There are two datasets used in this experiment, Programs and Organizations. Programs is a list of descriptions of humanitarian programs; the task is to determine which areas of humanitarian need are intended by the Program. The description typically mentions the mission and areas of focus for the program, which we anticipated would make Programs relatively easy to classify. The Organizations dataset is a list of companies and non-profits that may work in areas of interest for humanitarian causes. In this case, the descriptions are less likely to explicitly mention areas of humanitarian need. For both datasets, we associate each record with zero or more labels from the SDG and SPI ontologies. Statistics of the two datasets are shown in Table 2. We split each dataset into training, validation, and test set with 8:1:1 ratio. The models are optimized with validation set and the experimental results are reported from the test set.

We experimented with different text embedders and classification models to find the best combinations. Because the organization did not have abundant computation resources, we limited our choices within computation efficient models. We chose TF-IDF and Glove [48] as our text embedders. For the classification model, we adopted two frameworks for classification: one considered the hierarchical structure with graph encoding (named **Ontology**) within the labels and the other did not (named **naive**). The **naive** model considered the labels as a flat list. The model consisted of two fully connected layers and was optimized with Binary cross entropy, which is often used for multi-label classification. The diagram of the **ontology** framework is shown in Fig. 2. The framework learned a representation for the label ontology using a graph autoencoder [29]. Then, the model considered the node embeddings and mapped the input instances onto the

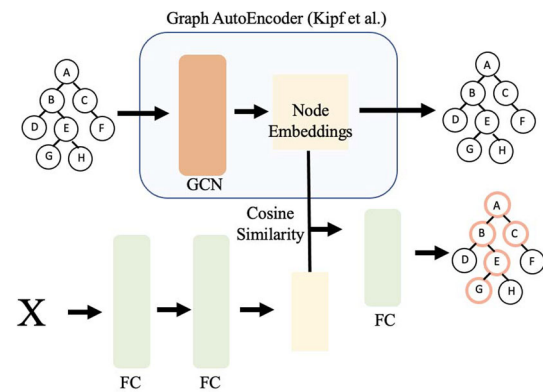


Fig. 2 Illustration of our framework

node embedding space with cosine similarity. Finally, the model was optimized with binary cross entropy and produce probability confidence as output. The threshold for classification is set to be 0.5. Finally, we evaluated all models with Precision (P), Recall (R), and F1 score which are commonly used in multi-label classification community. Following the literature, a data record is considered correctly classified when the predicted leaves match the ground truth exactly: there is no partial credit for siblings, for example.

Experimental results

We demonstrate our experimental results in Table 3. Because these two datasets are custom and not publicly available, we provide results from two baseline methods—majority vote and random selection. We can observe that considering the label ontology significantly improve the results. The trained model then allows us to tag the discourses on social media based on the humanitarian ontologies from SPI and SDG and to visualize within an online dashboard serving policymakers and entrepreneurs. As a result, we can organize public discourse and participation to capture levels of interest in various topics.

This approach is potentially critical for addressing data scarcity in practice. As we have argued, in complex domains, obtaining labeled data is expensive and requires significant human expertise. For example, determining whether a potential project is related to a goal to *enhance inclusive and sustainable urbanization* (SDG 11.3), *achieve sustainable management of resources* (SDG 12.2), *encourage adoption of sustainable practices* (SDG 12.6), or all three, requires significant expertise with the SDG ontology, municipal government practices, and the data being labeled. Moreover, labeled datasets can be rendered obsolete with only minor changes to the ontology, requiring an expensive re-labeling exercise. To enable comprehensive cross-sector models that can be deployed in a variety of contexts, we need to make efficient use of the human attention invested in creating the ontology.

Table 2 Dataset statistics

| Programs | | | Organization | | |
|----------|-----|------|--------------|-----|------|
| Train | Val | Test | Train | Val | Test |
| 6412 | 801 | 802 | 4558 | 570 | 570 |

Table 3 Experimental results on the Program and Organization datasets

| | Programs | | | | Organization | | | |
|-------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | Acc | P | R | F1 | Acc | P | R | F1 |
| Majority | 0.075 | 0.006 | 0.075 | 0.01 | 0.056 | 0.003 | 0.056 | 0.006 |
| Random | 0.003 | 0.016 | 0.003 | 0.003 | 0.006 | 0.018 | 0.006 | 0.008 |
| TF-IDF + naive | 0.102 | 0.085 | 0.132 | 0.087 | 0.090 | 0.032 | 0.090 | 0.030 |
| Glove + naive | 0.202 | 0.143 | 0.201 | 0.158 | 0.249 | 0.141 | 0.249 | 0.171 |
| TF-IDF + Ontology | 0.145 | 0.134 | 0.167 | 0.159 | 0.143 | 0.095 | 0.123 | 0.117 |
| Glove + Ontology | 0.245 | 0.175 | 0.254 | 0.219 | 0.286 | 0.176 | 0.287 | 0.256 |

Bold indicates best scores (our method)

Acc accuracy, P precision, R recall, F1 F1 score. Because these are custom datasets that are not publicly available, we also provide results from two baseline methods, majority vote and random selection. We can observe that considering ontology improves the results significantly

5 Modeling governance behaviors

The source of municipal democracy in the United States is found in city halls across the country. Even as our collective work in the analysis of urban space is used to create, debate, and ultimately enact urban policy, there is a lack of large-scale quantitative studies on municipal government. Comparative research into municipal governance in the USA is often prohibitively difficult due to a broad federal system where states, counties, and cities divide legislative powers differently. This power distribution has contributed to the lack of necessary research into the procedural elements of administrative and legislative processes, because it affords each municipality to each have their own standards for archival and publishing of municipal data [62].

To better study the complexities of municipal councils across the county, multiple tools are needed to standardize and aggregate data into large research databases and access portals. The data from municipal government meetings (videos, transcripts, voting records, etc.) must be made more accessible to both the general public and to researchers, and, such tools must be deployed in multiple municipalities across the nation, so that data can be used in aggregate to study the spread of policy, topic coverage, public sentiment, and more.

Once this infrastructure is available, it becomes possible to conduct large-scale quantitative studies on the dynamics of discourse in policy deliberation and enactment, quantifying how much of policy is decided upon using community sentiment as the policy basis, how such policy is supported or not from the public, and how similar policy proposals in different municipalities (or levels of government) are discussed and either enacted or rejected.

5.1 Council data project

To enable such large-scale studies, we have begun work on “Council Data Project,” [6] a suite of tools for deploying and managing infrastructure for rapidly generating, archiving, and analyzing transcript datasets

of municipal council meeting content. Council Data Project (CDP) is easily deployable and generalizes to many different meeting venues, but is specifically built with municipal council meetings in mind.

For each meeting a CDP deployment processes, our tools generate a transcript of timestamped sentences, and archives the produced transcript and all attached metadata (minutes items, presentations and attachments, voting records, etc.). CDP deployments additionally create a keyword-based index multiple times a week to enable plain-text search of events.

To further the utility of the CDP produced corpus, we are creating audio classification models for labeling each sentence with the classified speaker, aligning sentences in the transcript to the provided list of minutes item, re-using the generated keyword-based index for a municipality level n-gram viewer [41], and much more. Such work will enable the creation of datasets such as a dataset of discussions where only a set of specific councilmembers are present, or a dataset of discussions regarding specific pieces of legislation (minutes items).

Council Data Project enables large-scale quantitative studies by generating standardized municipal governance corpora—including legislative voting records, timestamped transcripts, and full legislative matter attachments (related reports, presentations, amendments, etc.). CDP enables the reproduction of political science research such as studying the effects of gender, ideology, and seniority in council deliberation [24], and studying the effects that adopting information communication technologies have on the civic participation process [12].

In constructing CDP to be as easily deployable as possible, we enable studies to understand how these behaviors generalize (or act as outliers) in different municipalities and settings.

Effective use of this new source of data motivates research in adapting deep learning techniques to multi-speaker, structured settings. Tasks include identifying speakers, topic and sentiment labeling by speaker to understand political positions, labeling speech by agenda topic, summarizing public sentiment to guide outreach, and communication investments. These prob-

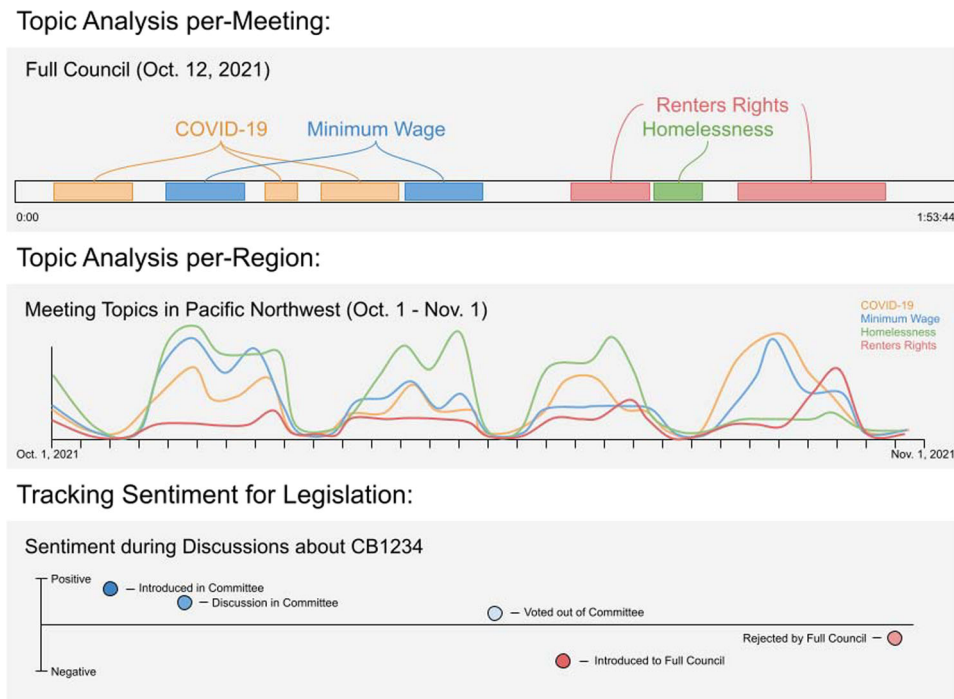


Fig. 3 Examples of analysis made possible through CDP infrastructure. Using the produced transcripts, we can build topic models to tag topics both in a single meeting’s transcript and track topic trends over time. With multiple CDP instances, we can show how these trends hold (and spread—

lems appear to be within the capabilities of emerging deep learning techniques, but require research attention in formulating the problem and evaluating competing techniques, which in turn require access to high-quality labeled datasets. Moreover, linking public discourse over social media with formal discourse in public hearings, administrative data collected through municipal service delivery, and geospatial data collected through sensing technologies will be required to meet our goal of a holistic study of the science of cities.

6 Conclusions

We aim to improve the coverage, access, and equity of urban data to advance understanding of city dynamics, unifying a top-down, holistic view of cities as a complex system and bottom-up, application-oriented view of cities as an assembly of independent subsystems. We aim to combat the disproportionate attention received by online advertising, face recognition, image labeling, and NLP tasks that dominate the machine learning literature by making high-quality, comprehensive urban datasets available for research. We identify four areas of research, with promising preliminary results, that involve the application of AI in urban contexts: spatiotemporal interpolation of data, unifying fairness, and interpretability in the context of emerging regulation of algorithms, accommodating the complex domain mod-

investigating the topical latency between municipalities over entire regions. Additionally, building models for tracking the sentiment of discussions regarding specific pieces of legislation as they move through council

els that are necessary to describe cities holistically, and engaging with new sources of data at the intersection of public discourse and policymaking.

References

1. T. Abid, H. Zarzour, M.R. Laouar, M.T. Khadir, Towards a smart city ontology. In *2016 IEEE/ACS 13th International Conference of Computer Systems and Applications (AICCSA)*, IEEE (2016), pp. 1–6
2. P.M. Allen, in *Cities and Regions as Self-Organizing Systems: Models of Complexity* (Routledge, 2012)
3. K. Ban, Sustainable development goals (2016). <https://sdgs.un.org/goals>. Accessed Oct 2021
4. R. Berk, H. Heidari, S. Jabbari, M. Joseph, *Michael Kearns (Seth Neel, and Aaron Roth. A convex framework for fair regression* (Jamie Morgenstern, 2017)
5. S.L. Blader, T.R. Tyler, A four-component model of procedural justice: defining the meaning of a fair process. *Person. Soc. Psychol. Bull.* **29**(6), 747–758 (2003). (PMID: 15189630)
6. J.M. Brown, T. Huynh, I. Na, B. Ledbetter, H. Ticehurst, S. Liu, E. Gilles, K.M.F. Greene, S. Cho, S. Ragoler, N. Weber, Council data project: software for municipal data collection, analysis, and publication. *J. Open Source Softw.* **6**(68), 3904 (2021)
7. P. Flavio, in *Calmon, Dennis Wei, Karthikeyan Natesan Ramamurthy, and Kush R (Varshney, Optimized data pre-processing for discrimination prevention)* (2017)

8. K. Cranmer, J. Brehmer, G. Louppe, The frontier of simulation-based inference. *Proc. Natl. Acad. Sci.* **117**(48), 30055–30062 (2020)
9. A. De Nicola, M.L. Villani, Smart city ontologies and their applications: a systematic literature review. *Sustainability* **13**(10), 5578 (2021)
10. B. Donovan, DB. Work, New York city taxi trip data (2010–2013). University of Illinois Urbana-Champaign, Champaign, IL, USA, Technical Report (2014)
11. A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, V. Koltun, An open urban driving simulator, Carla (2017)
12. K.L. Einstein, D. Glick, L.G. Ptug, M. Palmer, Evidence from public meeting minutes, Zoom does not reduce unequal participation (2021)
13. P. Espinoza-Arias, M. Poveda-Villalón, R. García-Castro, O. Corcho, Ontological representation of smart city data: from devices to cities. *Appl. Sci.* **9**(1), 32 (2019)
14. EUR-Lex, Proposal for a regulation of the European parliament and of the council. laying down harmonised rules on artificial intelligence (artificial intelligence act) and amending certain union legislative acts. <https://eur-lex.europa.eu/legal-content/EN/TXT/?qid=1623335154975&uri=CELEX%3A52021PC0206> (2021). Accessed Oct 2021
15. S. Feng, F. Ping, W. Zheng, A hierarchical multi-label classification method based on neural networks for gene function prediction. *Biotechnol. Equip.* **32**(6), 1613–1621 (2018)
16. E. Giunchiglia, T. Lukasiewicz, Coherent hierarchical multi-label classification networks. *NeurIPS* **2020**, 33 (2020)
17. G. Grekousis, Artificial neural networks and deep learning in urban geography: a systematic review and meta-analysis. *Comput. Environ. Urb. Syst.* **74**, 244–256 (2019)
18. A. Gyrard, A. Zimmermann, A. Sheth, Building IOT-based applications for smart cities: how can ontology catalogs help? *IEEE Internet Things J.* **5**(5), 3978–3990 (2018)
19. S. Ha, H. Jeong, Unraveling hidden interactions in complex systems with deep learning. *Sci. Rep.* **11**(1), 1–13 (2021)
20. M. Hardt, E. Price, N. Srebro, Equality of opportunity in supervised learning (2016)
21. H. Heidari, C. Ferrari, K. Gummadi, A. Krause, Fairness behind a veil of ignorance: a welfare analysis for automated decision making. In *Advances in Neural Information Processing Systems* (2018), pp. 1265–1276
22. S. Iizuka, E. Simo-Serra, H. Ishikawa, Globally and locally consistent image completion. *ACM Trans. Graph.* **36**(4), 1–14 (2017). <https://doi.org/10.1145/3072959.3073659>
23. California Legislative Information, Ab-13 public contracts: automated decision systems. (2021–2022). https://leginfo.ca.gov/faces/billTextClient.xhtml?bill_id=20210220AB13 (2021). Accessed Oct 2021
24. T. Jacobi, D. Schweers, Justice, interrupted: the effect of gender, ideology, and seniority at supreme court oral arguments. *Va. L. Rev.* **103**, 1379 (2017)
25. J. Jacobs, *The Death and Life of Great American Cities*. (Vintage, 2016)
26. E. Jillson, Aiming for truth, fairness, and equity in your company’s use of AI (2021)
27. S.-H. Kang, Y. Choi, J.Y. Choi, Restoration of missing patterns on satellite infrared sea surface temperature images due to cloud coverage using deep generative inpainting network. *J. Mar. Sci. Eng.* **9**(3), 310 (2021)
28. C.P. Kempes, G.B. West, The bridge on complex unifiable systems. *Bridge* **50**(4) (2020)
29. T.N. Kipf, M. Welling, Variational graph auto-encoders. In *Bayesian Deep Learning Workshop (NeurIPS 2016)* (2016)
30. P. Lahoti, A. Beutel, J. Chen, K. Lee, F. Prost, N. Thain, X. Wang, E.H. Chi, Fairness without demographics through adversarially reweighted learning (2020)
31. J. Angwin, J. Larson, S. Mattu, L. Kirchner, How We Analyzed the COMPAS Recidivism Algorithm (Pro Publica, 2016). <https://www.propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm>
32. M. Lee, E. Almirall, J. Wareham, Open data and civic apps: first-generation failures, second-generation improvements. *Commun. ACM* **59**(1), 82–89 (2015)
33. M.K. Lee, A. Jain, H.J. Cha, S. Ojha, D. Kusbit, Procedural justice in algorithmic fairness: leveraging transparency and outcome control for fair algorithmic mediation. In *Proceedings of ACM Human–Computer Interactions, 3(CSCW)* (2019)
34. G.S. Leventhal, *What Should Be Done with Equity Theory* (Springer US, Boston, 1980), pp. 27–55
35. B. Liao, J. Zhang, C. Wu, D. McIlwraith, T. Chen, S. Yang, Y. Guo, F. Wu, Deep sequence learning with auxiliary information for traffic prediction. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. ACM* (2018), pp. 537–546
36. G. Liu, F.A. Reda, K.J. Shih, T.-C. Wang, A. Tao, B. Catanzaro, Image inpainting for irregular holes using partial convolutions (2018)
37. S.M. Lundberg, S.-I. Lee, A unified approach to interpreting model predictions. In ed. by I. Guyon, U.V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, R. Garnett, *Advances in Neural Information Processing Systems*, vol. 30 (Curran Associates, Inc., 2017)
38. X. Ma, Z. Dai, Z. He, J. Ma, Y. Wang, Y. Wang, Learning traffic as images: a deep convolutional neural network for large-scale transportation network speed prediction. *Sensors* **17**(4), 818 (2017)
39. D. Madras, E. Creager, T. Pitassi, R.S. Zemel, in *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, July 10–15, 2018*, ed. by J.G. Dy, A. Krause. Learning adversarially fair and transferable representations. *Proceedings of Machine Learning Research*, vol. 80 (PMLR, 2018), pp. 3381–3390. <http://proceedings.mlr.press/v80/madras18a.html>
40. F. Marcinkowski, K. Kieslich, C. Starke, M. Lünich, Implications of ai (un-)fairness in higher education admissions: the effects of perceived ai (un-)fairness on exit, voice and organizational reputation. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency, FAT* ’20, New York, NY, USA* (Association for Computing Machinery, 2020), pp. 122–130

41. J.-B. Michel, Y.K. Shen, A.P. Aiden, A. Veres, M.K. Gray, Google Books Team, J.P. Pickett, D. Hoiberg, D. Clancy, P. Norvig, et al. Quantitative analysis of culture using millions of digitized books. *Science* **331**(6014):176–182 (2011)
42. C.C. Miller, When algorithms discriminate. *N. Y. Times* **9** (2015)
43. S.J. Mooney, K. Hosford, B. Howe, A. Yan, M. Winters, A. Bassok, J.A. Hirsch, Spatial equity in access to dockless bike share, Freedom from the station. *J. Transp. Geogr.* **74**, 91–96 (2019)
44. L. Moreira-Matias, J. Gama, M. Ferreira, J. Mendes-Moreira, L. Damas, Predicting taxi-passenger demand using streaming data. *IEEE Trans. Intell. Transport. Syst.* **14**(3), 1393–1402 (2013)
45. W. James Murdoch, P.J. Liu, B. Yu, Beyond word importance: contextual decomposition to extract interactions from lstms (2018)
46. M.D. Nadai, J. Staiano, R. Larcher, N. Sebe, D. Quercia, B. Lepri, The death and life of great Italian cities: a mobile phone data perspective. *CoRR abs/1603.04012* (2016)
47. D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, A.A. Efros, Context encoders: feature learning by inpainting (2016)
48. J. Pennington, R. Socher, C.D. Manning, Glove: global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)* (2014), pp. 1532–1543
49. K. Petrasic, B. Saul, J. Greig, M. Bornfreund, K. Lamberth, Algorithms and bias: What lenders need to know. White & Case (2017). <https://www.whitecase.com/publications/insight/algorithms-and-bias-what-lenders-need-know>. Accessed Oct 2021
50. I.D. Raji, E. Denton, E.M. Bender, A. Hanna, A. Paullada, AI and the everything in the whole wide world benchmark. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)* (2021)
51. M.Y. Ribeiro, S. Singh, C. Guestrin, "why should i trust you?": explaining the predictions of any classifier (2016)
52. L. Rieger, C. Singh, W. James Murdoch, B. Yu, *Proceedings of the 37th International Conference on Machine Learning*, vol. 119 (PMLR, 2020), pp. 8116–8126
53. C. Rudin, K.L. Wagstaff, Machine learning for science and society. *Mach. Learn.* **95**(1), 1–9 (2014)
54. C. Schumann, X. Wang, A. Beutel, J. Chen, H. Qian, E.H. Chi, Transfer of machine learning fairness across domains. *CoRR* (2019). [arXiv: 1906.09688](https://arxiv.org/abs/1906.09688)
55. M. Scrocca, I. Baroni, I. Celino, Urban IOT ontologies for sharing and electric mobility. *Semant. Web Interoper. Usab. Appl.* (2021). <http://www.semantic-web-journal.net/content/urban-iot-ontologies-sharing-and-electric-mobility-0>
56. B. Shen, X. Liang, Y. Ouyang, M. Liu, W. Zheng, K.M. Carley, Stepdeep: a novel spatial-temporal mobility event prediction framework based on deep neural network. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. ACM* (2018), pp. 724–733
57. A.J. Stier, M.G. Berman, L.M.A. Bettencourt, COVID-19 attack rate increases with city size. *MedRxiv* (2020). <https://doi.org/10.1101/2020.03.22.20041004>
58. M. Sundararajan, A. Taly, Q. Yan, in *Proceedings of the 34th International Conference on Machine Learning - Volume 70. Axiomatic attribution for deep networks, ICML'17.* (JMLR.org, Sydney, 2017), pp. 3319–3328
59. M. Taddy, H.K.H. Lee, B. Sansó, Fast Bayesian inference for computer simulation inverse problems. *AMS2008-3*, UC Santa Cruz (2008). <https://tr.soe.ucsc.edu/research/technical-reports/AMS2008-3>
60. J. Tasnim, D. Mondal, Data reduction and deep-learning based recovery for geospatial visualization and satellite imagery. In *2020 IEEE International Conference on Big Data (Big Data)* (2020), pp. 5276–5285
61. J. Teller, A.K. Keita, C. Roussey, R. Laurini, Urban ontologies for an improved communication in urban civil engineering projects. Presentation of the cost urban civil engineering action c21 "towntology". *Cybergeog. Eur. J. Geogr.* (2007). <https://journals.openedition.org/cybergeog/8322>
62. J. Trounstine, All politics is local: the reemergence of the study of city politics. *Perspect. Politics* **7**(3), 611–618 (2009)
63. T.R. Tyler, Procedural justice and the courts. *Court Rev.* **26** (2007)
64. S. Verhulst, A. Young, A. Zahuranec, A. Calderon, M. Gee, S.A. Aaronson, The emergence of a third wave of open data: How to accelerate the re-use of data for public interest purposes while ensuring data rights and community flourishing (2020). <https://doi.org/10.2139/ssrn.3937638>
65. P. Vogel, T. Greiser, D.C. Mattfeld, Exploring activity patterns, Understanding bike-sharing systems using data mining. *Procedia Soc. Behav. Sci.* **20**, 514–523 (2011)
66. J. Wehrmann, R. Cerri, R. Barros, Hierarchical multi-label classification networks. In *ICML*, vol. 2018 (2018), pp. 5075–5084
67. G.B. West, Scale: the universal laws of growth, innovation, sustainability, and the pace of life in organisms, cities, economies, and companies. Penguin (2017)
68. M. Whittaker, The steep cost of capture. *Interactions* **28**(6), 50–55 (2021)
69. A. Yan, B. Howe, Equitensors: learning fair integrations of heterogeneous urban data. In *ACM SIGMOD* (2021), pp. 2338–2347
70. H. Yao, X. Tang, H. Wei, G. Zheng, Y. Yu, Z. Li, Modeling spatial-temporal dynamics for traffic prediction. [arXiv:1803.01254](https://arxiv.org/abs/1803.01254) (2018)
71. J.W. Yoon, F. Pinelli, F. Calabrese, Cityride: a predictive bike sharing journey advisor. In *2012 IEEE 13th International Conference on Mobile Data Management. IEEE* (2012), pp. 306–311
72. J. Yu, Z. Lin, J. Yang, X. Shen, X. Lu, T.S. Huang, in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*. Generative image inpainting with contextual attention (2018), pp. 5505–5514. <https://doi.org/10.1109/CVPR.2018.00577>
73. Z. Yuan, X. Zhou, T. Yang, Hetero-convlstm: a deep learning approach to traffic accident prediction on heterogeneous spatio-temporal data. In *Proceedings of the 24th ACM SIGKDD International Conference on*

- Knowledge Discovery & Data Mining ACM* (2018), pp. 984–992
74. R. Zemel, Y. Wu, K. Swersky, T. Pitassi, Cynthia Dwork, Learning fair representations. In ed. by S. Dasgupta, D. McAllester, *Proceedings of the 30th International Conference on Machine Learning, volume 28 of Proceedings of Machine Learning Research, Atlanta, Georgia, USA, 17–19 (Jun 2013)*. PMLR, pp. 325–333
 75. J. Zhang, Y. Zheng, D. Qi, R. Li, X. Yi, Dnn-based prediction model for spatio-temporal data. In *Proceedings of the 24th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems, ACM* (2016), p. 92
 76. Y. Junbo Zhang, D.Q. Zheng, R. Li, X. Yi, T. Li, Predicting citywide crowd flows using deep spatio-temporal residual networks. *Artif. Intell.* **259**, 147–166 (2018)
 77. Q. Zhang, Q. Yuan, C. Zeng, X. Li, Y. Wei, Missing data reconstruction in remote sensing image with a unified spatial-temporal-spectral deep convolutional neural network. *IEEE Trans. Geosci. Remote Sens.* **56**(8), 4274–4288 (2018)
 78. Z. Zou, S. Tian, X. Gao, Y. Li, mldeepr: multi-functional enzyme function prediction with hierarchical multi-label deep learning. *Front. Genet.* **9**, 714 (2019)