

Abstract

Explanation is a simple model (e.g., linear) that approximates the decision boundary of a complex model of interest (e.g., deep neural net).

Questions:

Consistency. Explanations are as good as the features they use to explain predictions. How do feature selection and feature noise affect explanations?

Performance. When explanation is a part of the learning and prediction process, how does that affect performance of the predictive model?

Insights. What insights we can gain by visualizing and inspecting explanations?

1. Background

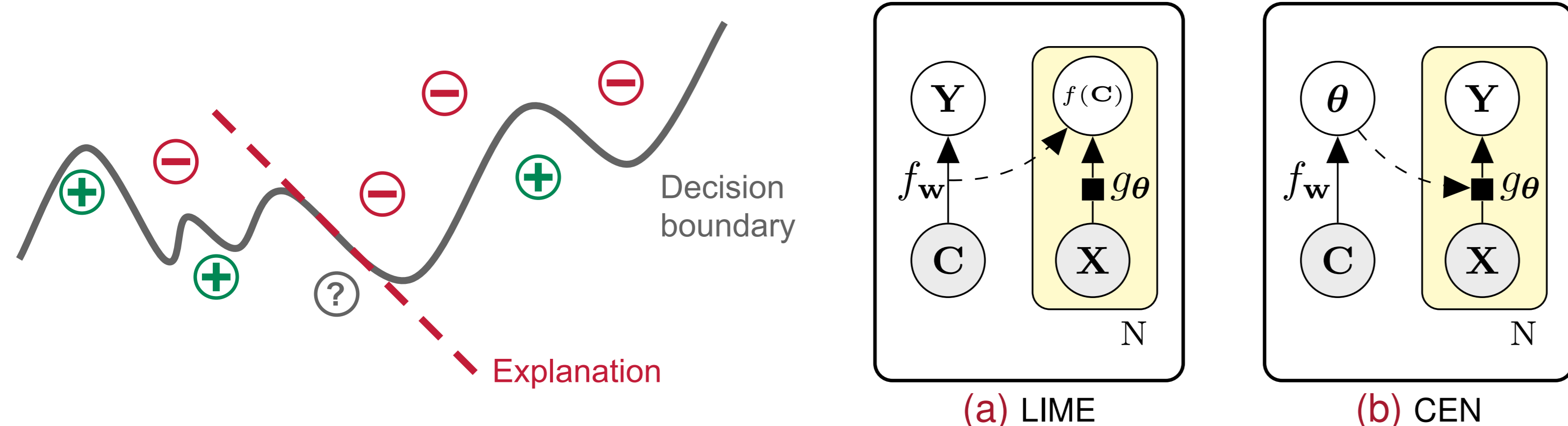


Figure 1: (a) Local interpretable model-agnostic explanations (LIME). (b) Contextual explanation networks.

Data representations:

- C : low-level or unstructured features (e.g., text, image pixels, sensory inputs, etc.)
- X : high-level or human-interpretable features (e.g., categorical variables).

Constructing explanations *post-hoc* [Ribeiro et al., 2016]

$$\text{model: } f : C \mapsto Y, \quad \text{explanation: } g_c := \operatorname{argmin}_{g \in \mathcal{G}} \mathcal{L}(f, g, \pi_c) + \Omega(g)$$

where

- $\mathcal{L}(f, g, \pi_c)$ – loss that measures how well g approximates f in the neighborhood, π_c
- Linear explanation: $g_c(x) := b_c + \theta_c \cdot x$, parametrized by (b_c, θ_c)
- $\Omega(g)$ – explanation complexity penalty (typically, L_1)

2. Contextual Explanation Networks (CENs)

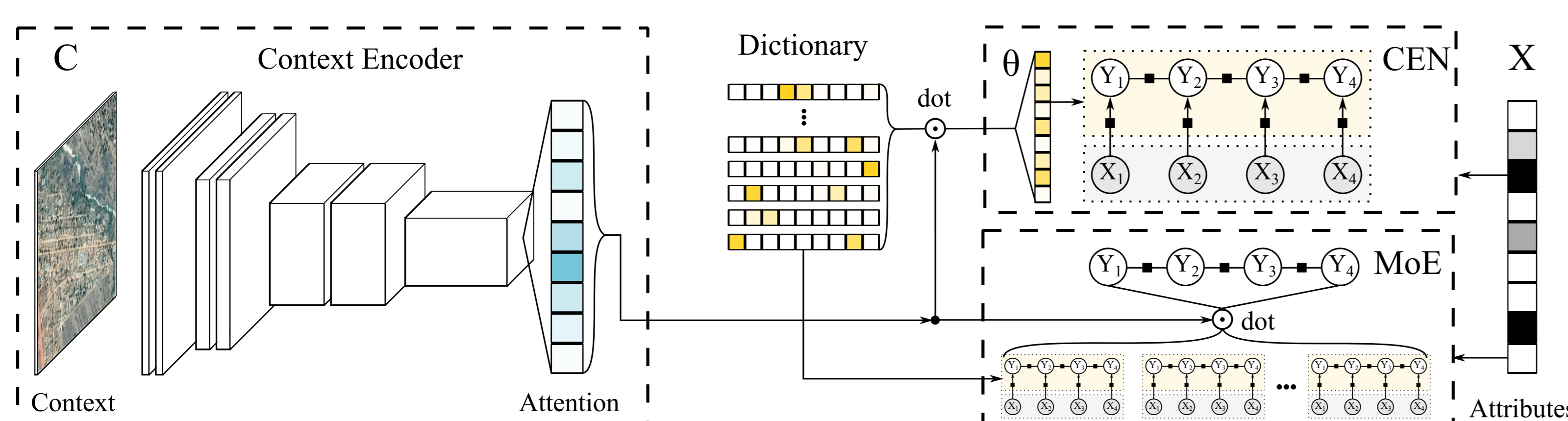


Figure 2: An example of CEN and MoE architectures with CNN context-encoder and structured explanations.

Learning to predict and explain jointly [our work]

$$\text{model: } f(c, x) := g_c(x), \quad \text{explanation: } g_c := \text{enc}(c)$$

where

- Linear explanation: $g_c(x) := b + \theta_c \cdot x$, parametrized by $\theta := (b, \theta_c)$
- $\text{enc}(c)$ – encoder that generates context-specific explanations, g_c
- g_c are applied to x to make a prediction

3. Consistency of Explanations

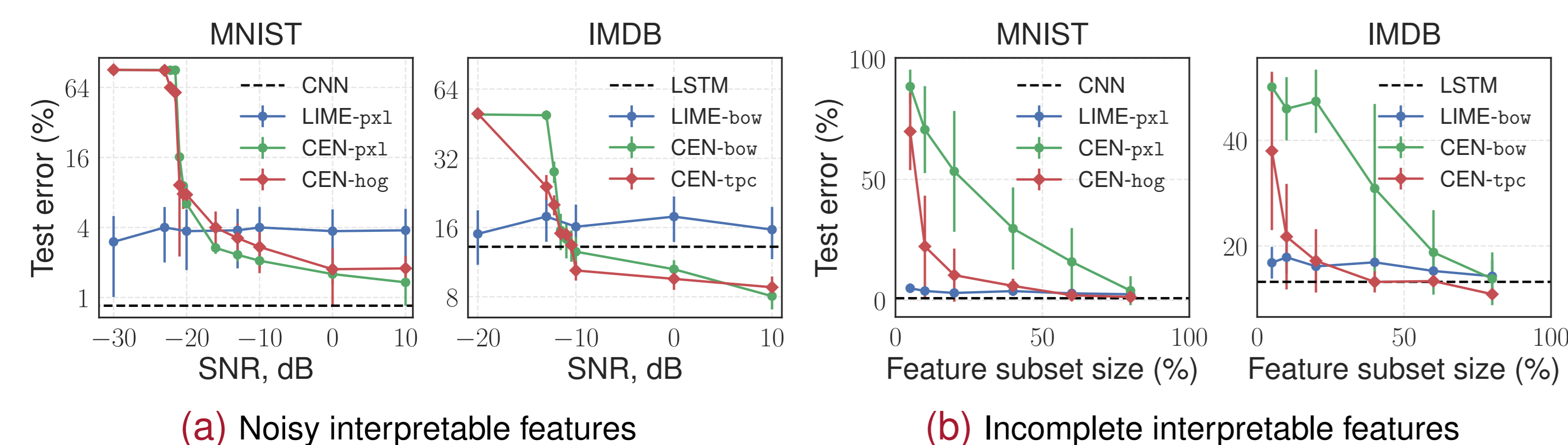


Figure 3: Predictive performance of LIME- and CEN-generated explanations applied to corrupted features.

Setup

- X is subsampled or corrupted with additive noise. C is kept the same.
- LIME and CEN produce explanations using incomplete/corrupted features.
- This simulates a real-life situation when interpretable features are selected imperfectly.

Takeaways

- **LIME overfits explanations.** Despite the loss of information, it almost perfectly approximates the decision boundary of a deep net. *Can we trust such explanations?*
- **CEN-generated explanations are consistent.** Performance of CENs is proportional to the quality of the interpretable features used by explanations.

4. Explanations as a Regularizer

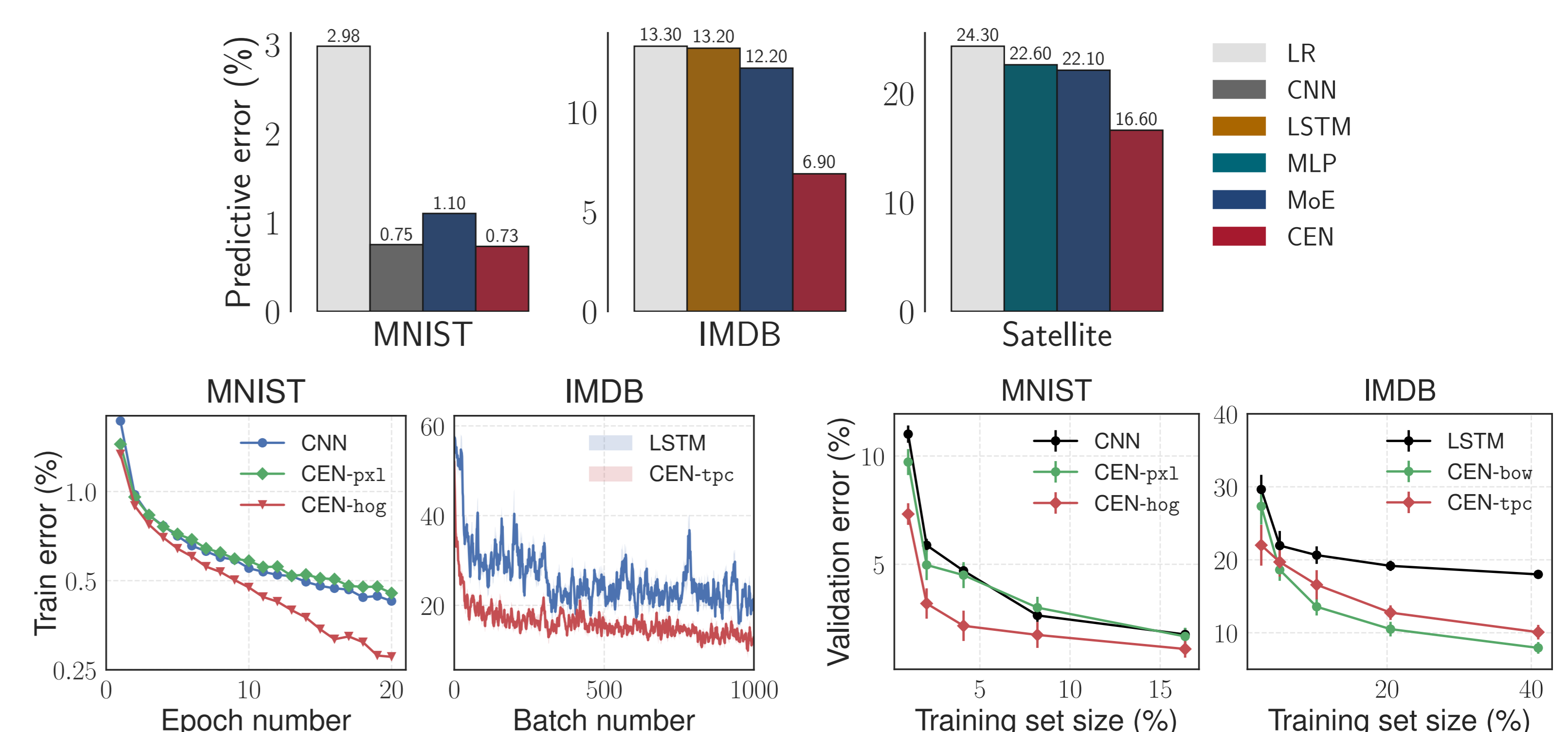


Figure 4: Upper: Predictive errors for the baselines, MoE, and CEN. Lower: Training error vs. iteration (epoch or batch) and validation error for models trained on random subsets of data of different sizes.

Takeaways

- **Large data regime.** When the data is abundant, prediction by explanation does not hurt the performance (CENs match the performance of vanilla deep nets).
- **Low data regime.** When the data is scarce, explanations can play the role of a good regularizer and boost performance.

5. Visualization & Insights

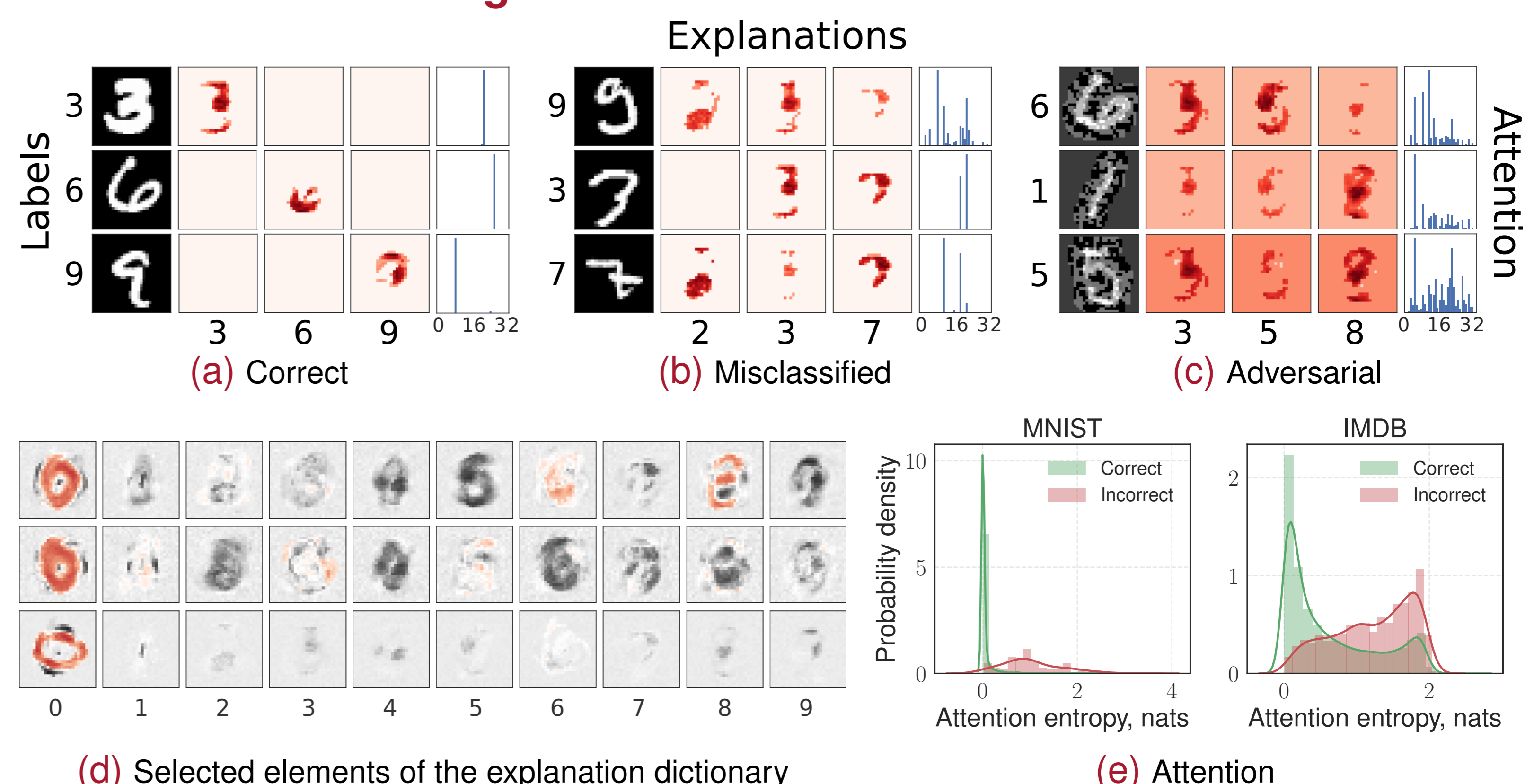


Figure 5: Explanations generated by CEN for the 3 top classes and the corresponding attention vectors for (a) correctly classified, (b) misclassified, and (c) adversarially constructed images. (d) Elements from the learned 32-element dictionary that correspond to different writing styles of 0 digits. (e) Histogram of the attention entropy for correctly and incorrectly classified instances.

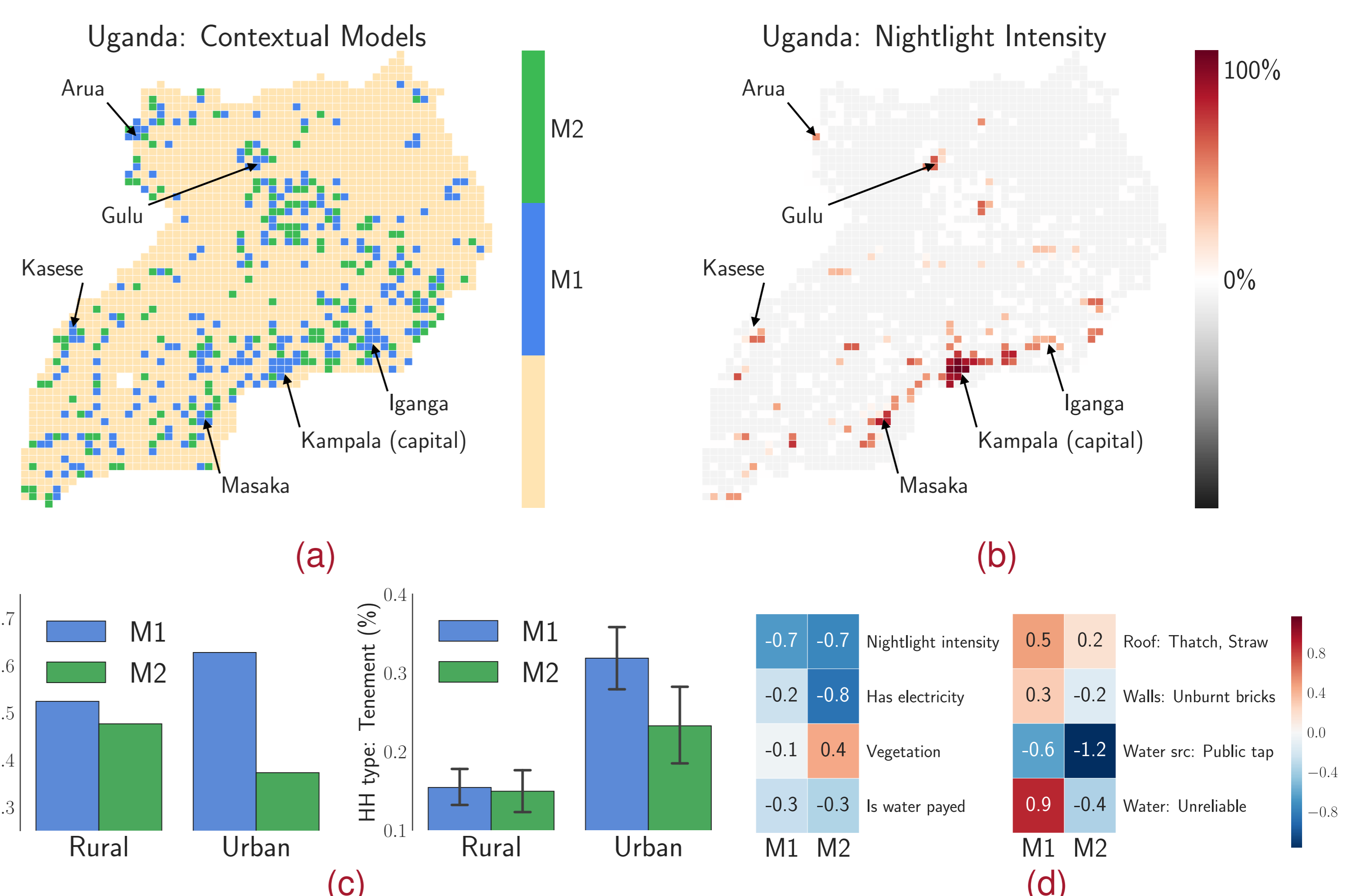


Figure 6: Explanations M1 and M2 (a) and nightlight intensities (b) for different areas in Uganda. (c) How frequently M1 and M2 are selected for rural/urban areas (left) and the average proportion of Tenement-type households in an urban/rural area for which M1 or M2 was selected (right). (d) Weights given to the survey features by the two explanations (M1 and M2) discovered by CEN.

Summary

Consistency. Be careful using post-hoc explanations as they can be silently misleading when the interpretable features are selected imperfectly. Learning to predict and to explain jointly helps to detect such problems.

Performance. Explanations can regularize the model and boost performance.

Insights. Visualizing explanations helps understanding the patterns in the data used by the model to make predictions.

Full paper: <http://arxiv.org/abs/1705.10301>