

# Explaining by Examples

Sargur N. Srihari  
srihari@buffalo.edu

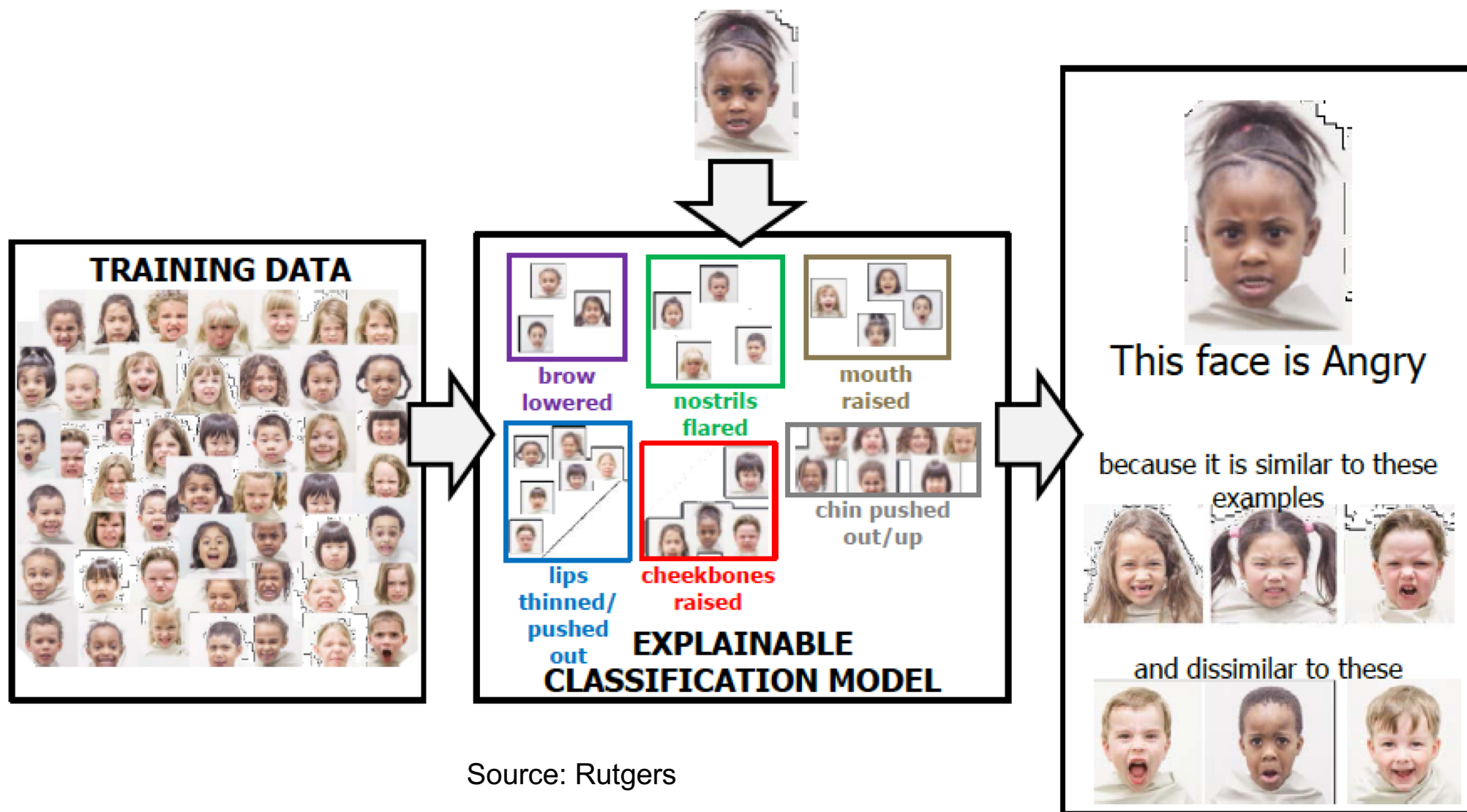
# Topics in Explaining by Examples

1. Two types of Explanation
2. Machine Teaching
3. Bayesian Teaching

# Two types of Explanation

- Explanation by Teaching Examples
  - Help to *understand the data better*
  - Image is of a cat because it looks like another example of a cat (nearest neighbor).
- Explanation by Features
  - It is a cat because it has whiskers and fur

# Explanation by Teaching Examples



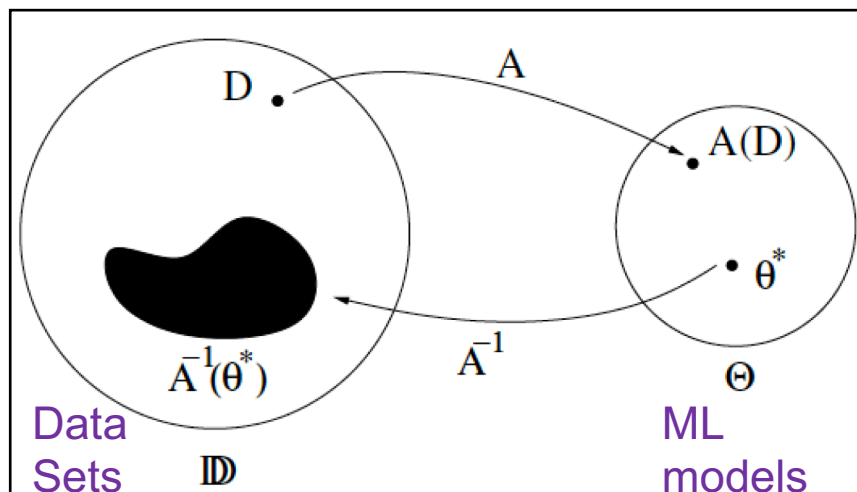
# Explanation by Examples

- Aim to select subset of dataset that leads to similar conclusions as entire dataset
- Intuition
  - *Subsets* of training data that lead a model to the same (or approximately similar) inference as the model trained on *all* the data should be useful for understanding the fitted model

# Motivating explanation by examples

- Explanation is the inverse of modeling
- Based on two fundamental observations
  1. All ML models are trained on data
  2. Data is common language of user and model
- Examples may not capture what philosophers and cognitive scientists call as explanation
  - but harnesses people's proclivity to inductive inference
- Related to *machine teaching*
  - Machine teaching designs the optimal training data to drive learning algorithm to a target model.

# Machine Teaching

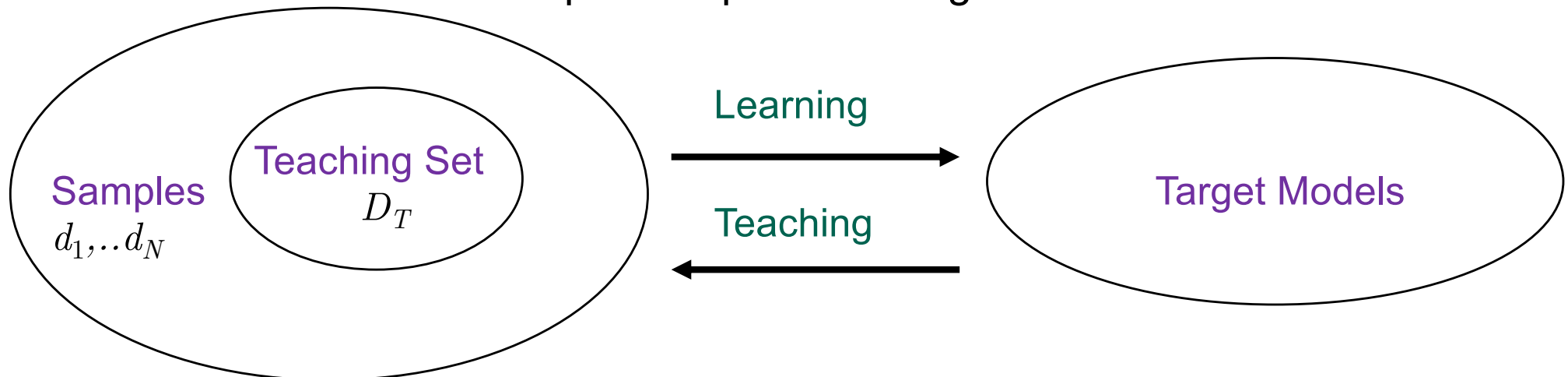


- Given training set  $D \in \mathbb{D}$ , ML returns model  $A(D) \in \Theta$ 
  - $A$  is in general many-to-one
    - Several data sets result in same model
- Conversely, given a target model  $\theta^* \in \Theta$ , the inverse function  $A^{-1}$  returns a set of training examples that will result in  $\theta^*$
- MT aims to identify optimal member among  $A^{-1}(\theta^*)$

# Bayesian Teaching

- Formalization

- Teaching problem is to select a small subset of data that with high probability, lead learner model to the correct inference
- This requires two kinds of inference
  - Teacher's inference
    - Done in the space of possible teaching sets
  - Learner's inference
    - Done in the space of possible target models





# Equation for Bayesian Teaching

$$P_T(x|\Theta) = \frac{P_L(\Theta|x)P(x)}{\int_{x'} P_L(\Theta|x')P(x')dx'},$$

- Where  $x$  can be any subset of the training data
- $\Theta$  denotes the target model, which can be an entire model or particular substructures, such as latent features, relations, grammars, programs, or combinations of these
- $P_T(x|\Theta)$  is the probability of choosing  $x$  as the teaching examples for explaining target model  $\Theta$
- $P_L(\Theta|x)$  is learner's posterior inference after receiving  $x$
- $P(x)$  describes bias for certain kind of examples (e.g., favoring smaller subsets)
- and the integral is over all partitions of the training data (i.e., if the size of  $x$  is  $m$  and the size of the entire training corpus is  $N$ , there are  ${}^N C_m$  partitions)

# A Bayesian Teacher

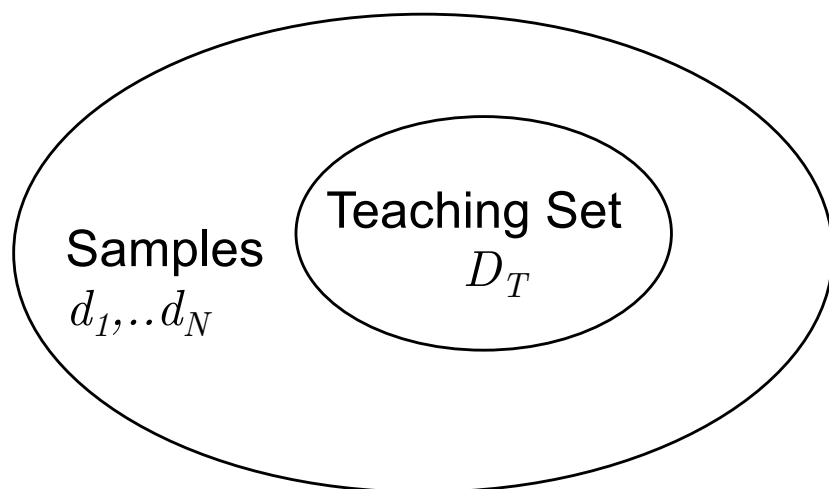
- With training data  $D = \{d_1, d_2, \dots, d_N\}$  and teaching set size  $n < N$  teach a target model  $\Theta^*$  by sampling a teaching set  $D_T \subseteq D$  from  $\mathcal{D} = \{D \mid D \in \mathcal{P}(D) \wedge |D| \leq n\}$  according to

$$\begin{aligned} p(\mathcal{D}_T \mid \Theta^*) &= \frac{p(\mathcal{D}_T) p_L(\Theta^* \mid \mathcal{D}_T)}{p(\Theta^*)} \\ &= \frac{p(\mathcal{D}_T) p_L(\Theta^* \mid \mathcal{D}_T)}{\sum_{\mathcal{D} \in \mathcal{D}} p(\mathcal{D}) p_L(\Theta^* \mid \mathcal{D})}, \end{aligned}$$

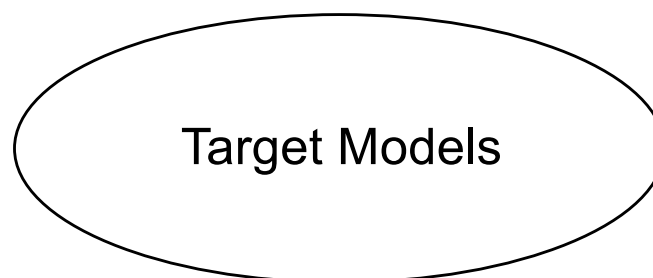
- $\mathcal{P}(D)$  is the power set of  $D$  and  $D$  is the space of teaching sets.  $p_L(\Theta^* \mid D)$  is the probability the learner will infer the target model  $\Theta^*$  given a particular teaching set  $D$  (i.e. the *learner's posterior* probability given that teaching set), and  $p(D)$  is the *teacher's prior* probability on the same teaching set  $D$
- Priors assign higher probabilities to smaller teaching sets

# Teaching Inference for model explanation

Teaching Probabilities



Model Probabilities



1. Pick a teaching set size (e.g. 2) to constrain the search space
2. Perform teaching inference for the *category to be understood*
3. Rank the teaching sets based on the teaching probabilities.

# Explaining a new ML model

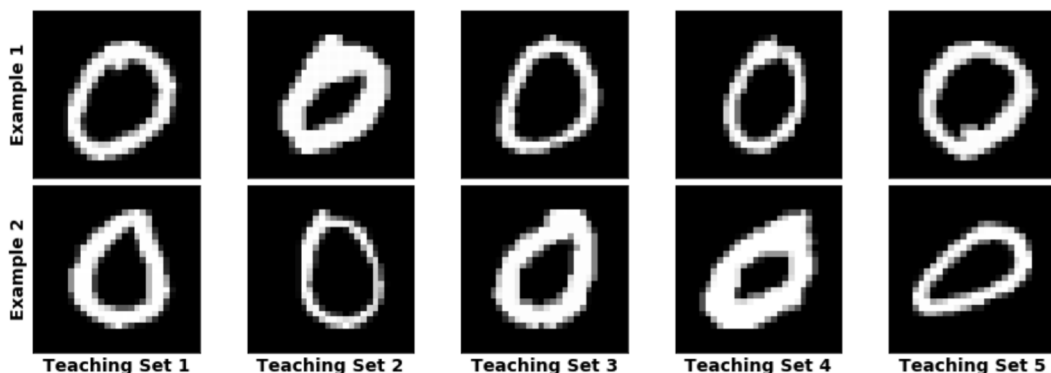
- At a conference, authors show examples classified correctly and incorrectly
  - Mention that humans would find some to be hard
- Rather than cherry-pick, we care about
  1. Those classified *correctly* with *high* confidence
  2. Those classified *correctly* with *low* confidence
  3. Those classified *incorrectly* with *low* confidence
  4. Those classified *incorrectly* with *high* confidence
    - But some methods don't offer certainty estimates
- Bayesian teaching offers finding/ranking such examples

# Using ground truth as labels

Five best teaching sets using ground truth labels.

Each pair of images represents a teaching set

Best Teaching Sets

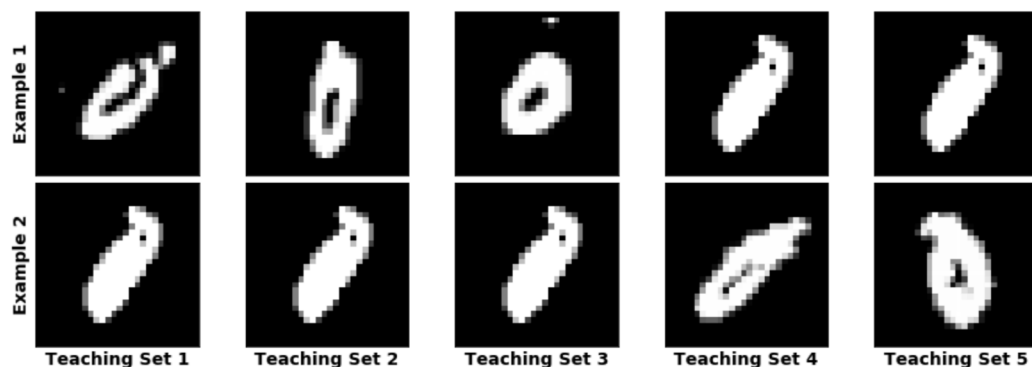


Pairs sorted by teaching probabilities in descending order, from left to right, (leftmost is best)

Five worst teaching sets using ground truth labels.

Each pair of images represents a teaching set

Worst Teaching Sets



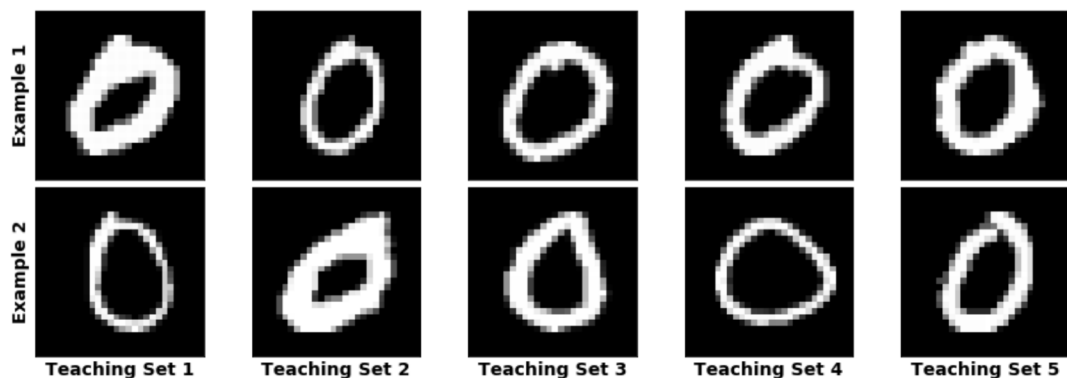
Pairs sorted by teaching probabilities in ascending order, from left to right, (leftmost set is worst)

# Using model predictions as labels

5 best teaching sets using model predictions as labels.

Each pair of images represents a teaching set

Best Teaching Sets

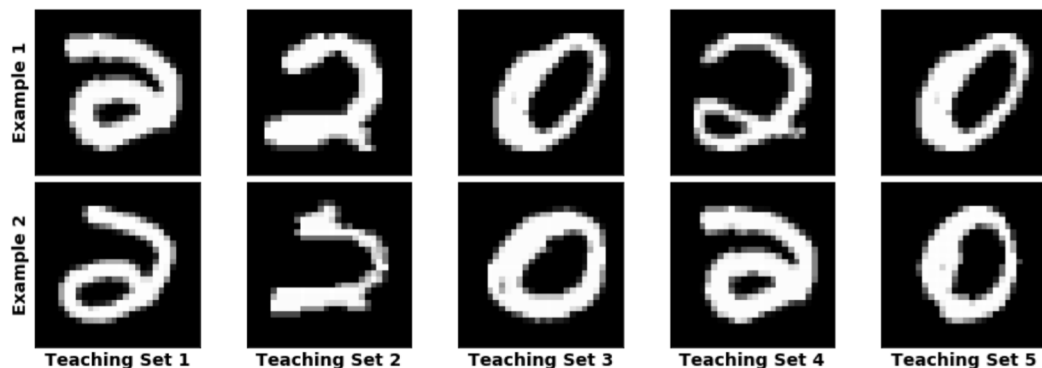


Pairs are sorted by teaching probabilities in descending order, from left to right, (i.e. leftmost set is best).

5 worst teaching sets using model predictions as labels.

Each pair of images represents a teaching set,

Worst Teaching Sets



Pairs are sorted by teaching probabilities in ascending order, from left to right (i.e. leftmost set is worst).

## Explaining the model for category 0

1. Those classified *correctly* with *high* confidence
2. Those classified *correctly* with *low* confidence
3. Those classified *incorrectly* with *low* confidence
4. Those classified *incorrectly* with *high* confidence

# Summary of Bayesian Teaching

- It leverages the natural common language for understanding model behavior—the data—to explain opaque models through the examples from the original data that are most representative of the inference
- In doing so, it integrates learning and explanation by taking the learning model as input into the explanation process, and outputs an explanation in terms of the examples from the original data