



# Learning under algorithmic triage

**Manuel Gomez Rodriguez**

Includes joint work with Nastaran Okati, Abir De, Paramita Koley, Ali Zarezade and Niloy Ganguly

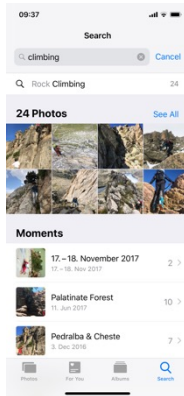


MAX PLANCK INSTITUTE  
FOR SOFTWARE SYSTEMS

# Most stunning developments in computing

Machine learning (ML) has *taught* machines to...

recognize images



translate between languages



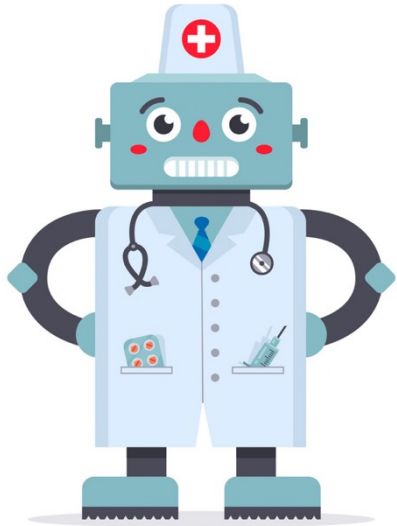
play complex games



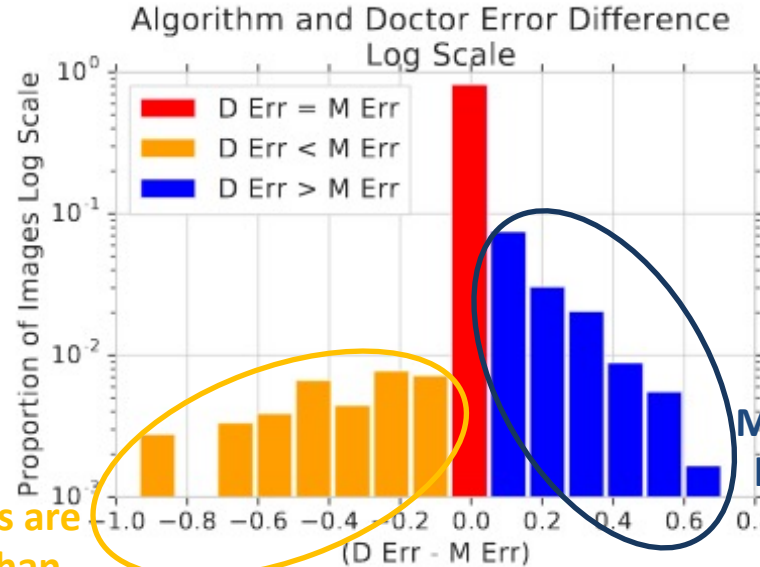
...to name a few

**Machines achieve, or surpass, human performance  
at tasks for which intelligence is required**

# Machines learning is sometimes worse than humans



On some instances, machine predictions are still worse than predictions made by human experts



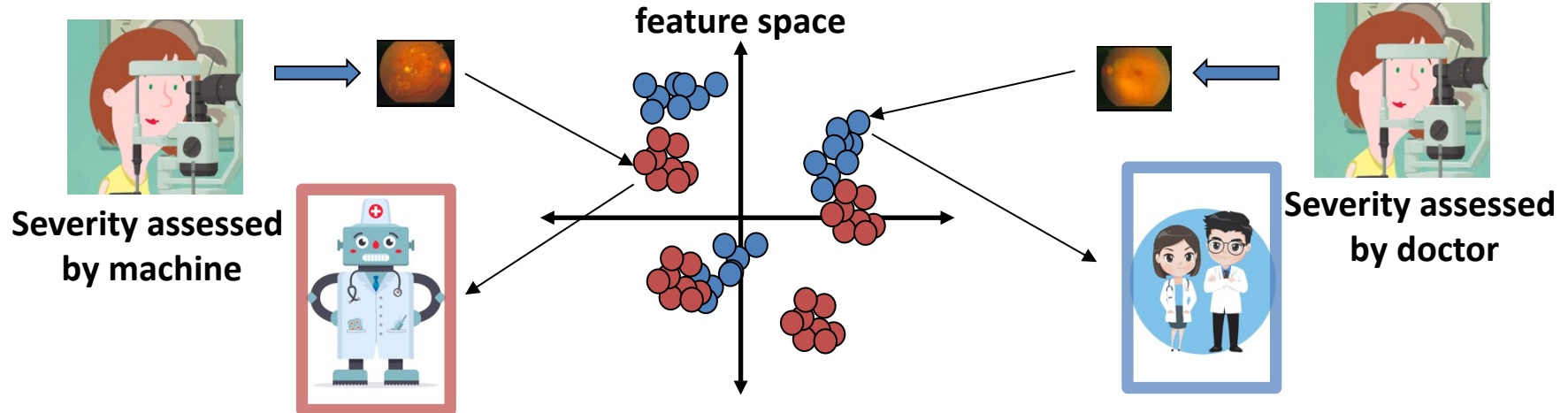
Task: estimating severity of diabetic retinopathy

Machines are worse than humans

Machines are better than humans

# Learning under algorithmic triage

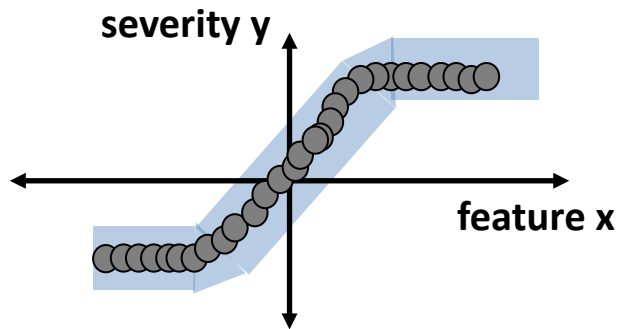
Develop machine learning models that are **optimized** to operate under **different automation levels**



They take **decisions** for a given **fraction of the instances** and leave the **remaining ones** to humans

# Isn't learning under triage just learning to defer?

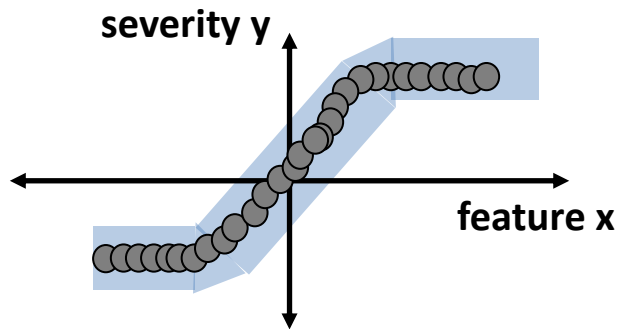
## Regression task



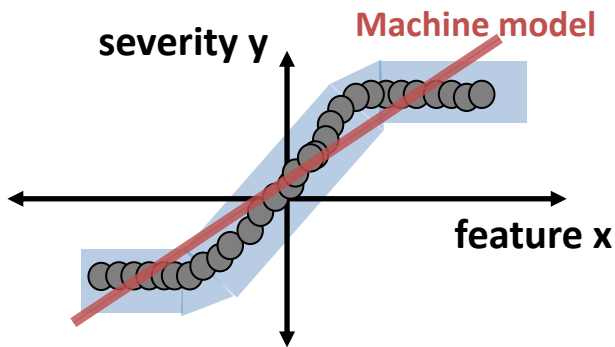
1. **The machine model** is a linear function
2. We can defer some samples to **humans**, as dictated by a *triage policy*

# Isn't learning under triage just learning to defer?

## Regression task



## Learning to defer

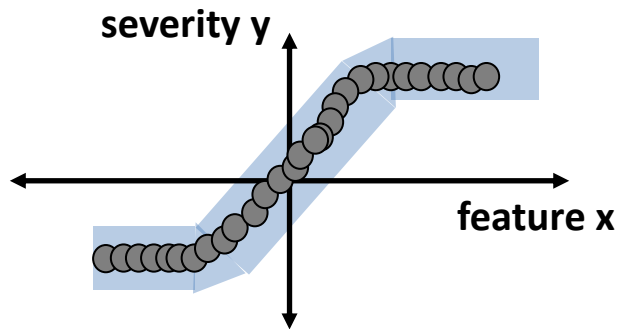


1. The machine model is a linear function
2. We can defer some samples to humans, as dictated by a triage policy

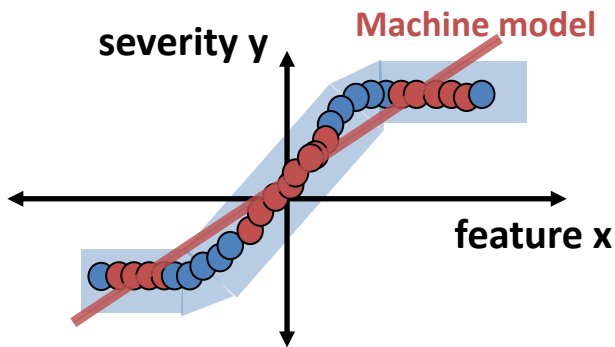
Machine model is not optimized during training

# Isn't learning under triage just learning to defer?

## Regression task



## Learning to defer

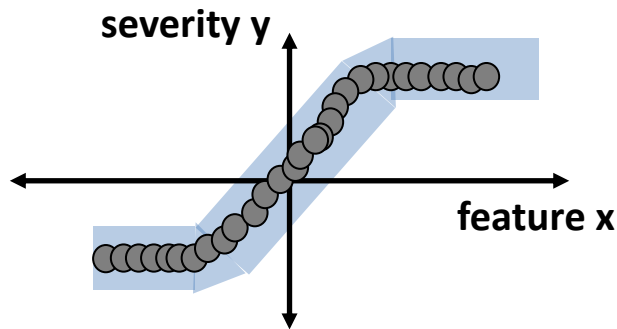


1. **The machine model** is a linear function
2. We can defer some samples to **humans**, as dictated by a *triage policy*

**Machine model**  
is not optimized  
during training

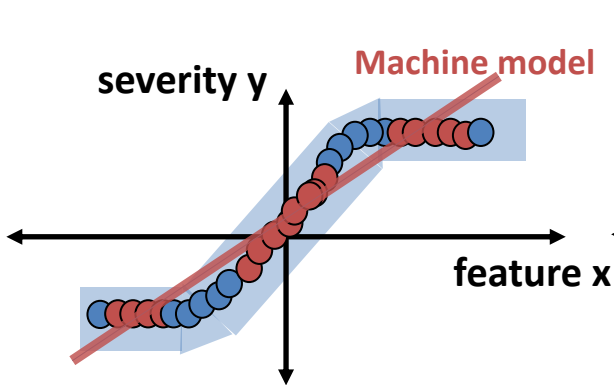
# Isn't learning under triage just learning to defer?

## Regression task



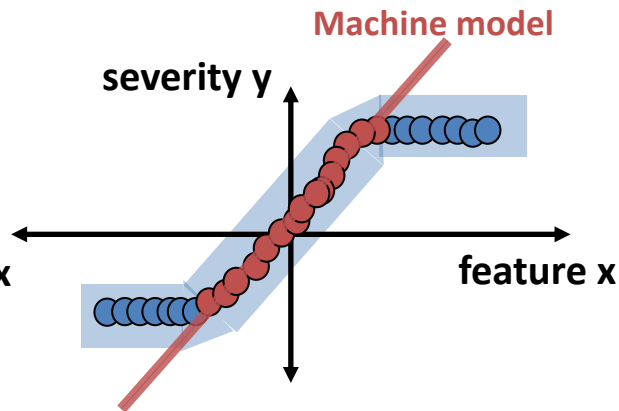
1. The machine model is a linear function
2. We can defer some samples to humans, as dictated by a triage policy

## Learning to defer



Machine model is not optimized during training

## Learning under triage

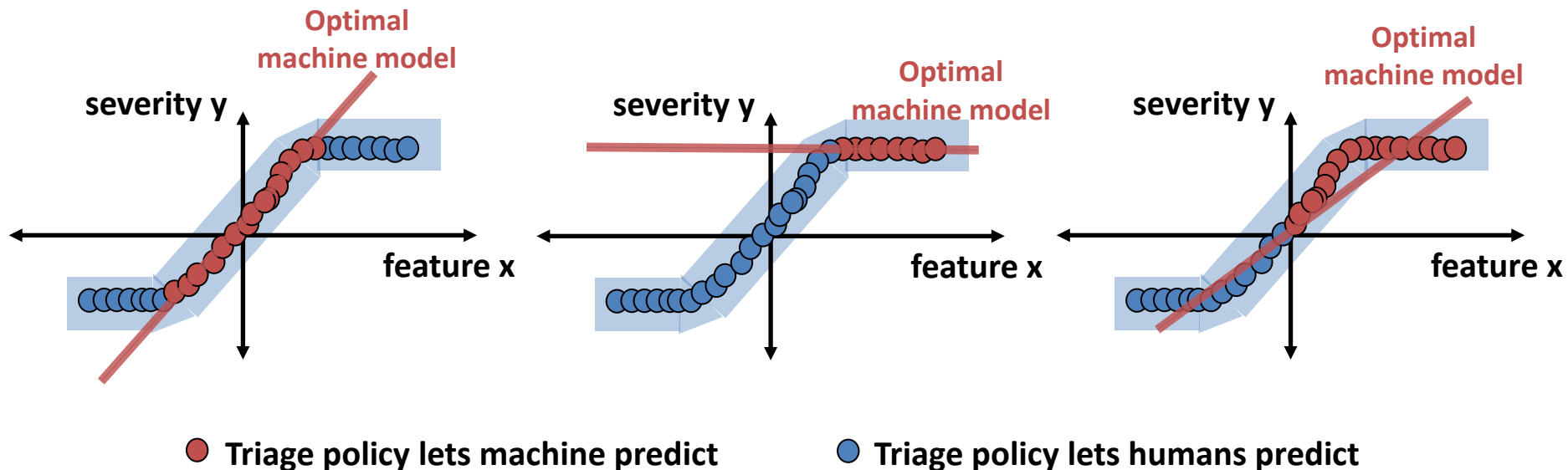


Machine model is optimized during training



# Main challenge of learning under algorithmic triage

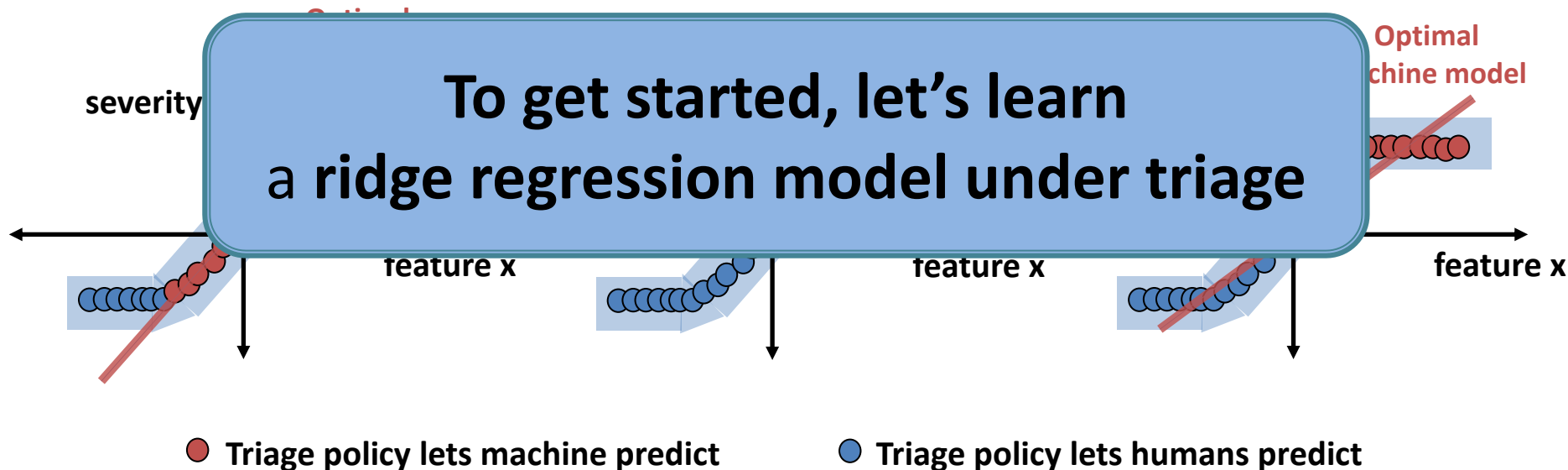
For each **triage policy**, there is an **optimal machine model**.  
However, the **triage policy** is also something **one seeks to optimize**.



# Main challenge of learning under algorithmic triage

For each **triage policy**, there is an **optimal machine model**.  
However, the **triage policy** is also something **one seeks to optimize**.

To get started, let's learn  
a ridge regression model under triage



# Ridge regression, revisited

## Training

$$\underset{w, \mathcal{S}}{\text{minimize}} \quad \sum_{i \in \mathcal{S}} \overbrace{c(\mathbf{x}_i, y_i)}^{\text{Human error per sample}} + \sum_{j \in \mathcal{S}^c} [(y_j - \mathbf{x}_j^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2]$$

subject to  $|\mathcal{S}| \leq n$

Annotations:

- Training samples assigned to humans (blue circles) point to  $i \in \mathcal{S}$ .
- Training samples assigned to machines (red circles) point to  $j \in \mathcal{S}^c$ .
- Machine model points to  $\mathbf{x}_j^\top \mathbf{w}$ .
- Regularization parameter points to  $\lambda \|\mathbf{w}\|_2^2$ .
- Max. number of samples that can be assigned to humans points to  $|\mathcal{S}| \leq n$ .

# Ridge regression, revisited

## Training

minimize  $w, \mathcal{S}$

$$\sum_{i \in \mathcal{S}} \overbrace{c(\mathbf{x}_i, y_i)}^{\text{Human error per sample}} + \sum_{j \in \mathcal{S}^c} [(y_j - \mathbf{x}_j^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2]$$

subject to  $|\mathcal{S}| \leq n$

Training samples assigned to **humans**      Training samples assigned to **machines**

**Machine model**

Regularization parameter

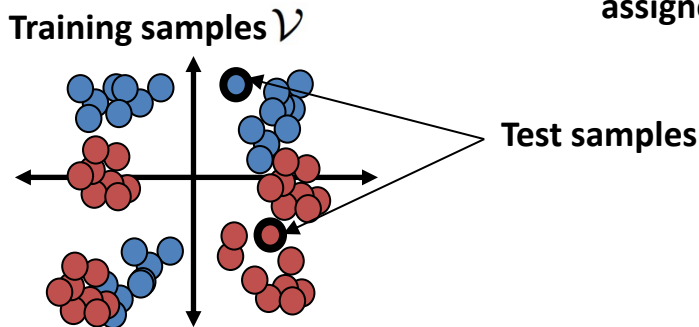
Max. number of samples that can be assigned to humans

## Test

We assign a test sample with features  $x$  to **humans** if

$$\operatorname{argmin}_{i \in \mathcal{V}} \|\mathbf{x}_i - \mathbf{x}\| \in \mathcal{S}^*$$

Training samples assigned to **humans**



# Ridge regression becomes a combinatorial problem

Given a fixed set  $\mathcal{S}$ , the **optimal machine model** is given by

$$\mathbf{w}^*(\mathcal{S}) = (\lambda|\mathcal{S}^c|\mathbb{I} + \mathbf{X}_{\mathcal{S}^c}\mathbf{X}_{\mathcal{S}^c}^\top)^{-1} \mathbf{X}_{\mathcal{S}^c}\mathbf{y}_{\mathcal{S}^c}$$

# Ridge regression becomes a combinatorial problem

Given a fixed set  $\mathcal{S}$ , the **optimal machine model** is given by

$$\mathbf{w}^*(\mathcal{S}) = (\lambda|\mathcal{S}^c|\mathbb{I} + \mathbf{X}_{\mathcal{S}^c}\mathbf{X}_{\mathcal{S}^c}^\top)^{-1} \mathbf{X}_{\mathcal{S}^c}\mathbf{y}_{\mathcal{S}^c}$$

Then, we can **rewrite the ridge regression problem** as a purely **combinatorial maximization problem**

$$\underset{\mathcal{S}}{\text{maximize}} \quad \underbrace{-\log \ell(\mathcal{S})}_{\downarrow} \quad \text{subject to} \quad |\mathcal{S}| \leq n$$

$$\sum_{i \in \mathcal{S}} c(\mathbf{x}_i, y_i) + \mathbf{y}_{\mathcal{S}^c}^\top \mathbf{y}_{\mathcal{S}^c} - \mathbf{y}_{\mathcal{S}^c}^\top \mathbf{X}_{\mathcal{S}^c}^\top (\lambda|\mathcal{S}^c|\mathbb{I} + \mathbf{X}_{\mathcal{S}^c}\mathbf{X}_{\mathcal{S}^c}^\top)^{-1} \mathbf{X}_{\mathcal{S}^c}\mathbf{y}_{\mathcal{S}^c}$$

# Ridge regression under human assistance is hard

Finding the solution to

$$\underset{\mathcal{S}}{\text{maximize}} \quad -\log \ell(\mathcal{S}) \quad \text{subject to} \quad |\mathcal{S}| \leq n$$

is a **NP-hard problem**

# Ridge regression under human assistance is hard

## Finding the solution to

$$\underset{\mathcal{S}}{\text{maximize}} \quad -\log \ell(\mathcal{S}) \quad \text{subject to} \quad |\mathcal{S}| \leq n$$

## is a NP-hard problem

### Proof sketch

Assume  $c(\mathbf{x}_i, y_i) = 0$ ,  $\lambda = 0$  and  $\mathbf{y} = \mathbf{X}^\top \mathbf{w}^* + \mathbf{b}^*$

k-sparse noise vector



Then, **the problem** can be viewed as the **robust least square (RLSR) problem**, which has been shown to be **NP-hard**:

$$\underset{\mathbf{w}, \mathcal{S}}{\text{minimize}} \quad \sum_{i \in \mathcal{S}} (y_i - \mathbf{x}_i^\top \mathbf{w})^2 \quad \text{subject to} \quad |\mathcal{S}| = |\mathcal{V}| - n$$



# A simple greedy algorithm

The **greedy algorithm** proceeds **iteratively**.

At each iteration, it assigns to a human the sample in the training set that provides the **largest marginal gain**, i.e.,

$$k^* \leftarrow \operatorname{argmax}_{k \in \mathcal{V} \setminus \mathcal{S}} \overbrace{-\log \ell(\mathcal{S} \cup k) + \log \ell(\mathcal{S})}^{\text{marginal gain}}$$
$$\mathcal{S} \leftarrow \mathcal{S} \cup \{k^*\} \quad \text{Points not yet assigned to humans}$$

Does this simple greedy algorithm has approximation guarantees? 😊

# The greedy algorithm has approximation guarantees

The function  $-\log \ell(\mathcal{S})$  satisfies an **approximate notion of submodularity**

$$-\log \ell(\mathcal{S} \cup k) + \log \ell(\mathcal{S}) \geq (1 - \alpha) [-\log \ell(\mathcal{T} \cup k) + \log \ell(\mathcal{T})]$$

for all  $\mathcal{S} \subseteq \mathcal{T} \subset \mathcal{V}$

where  $\alpha \leq \alpha^*$  is the **generalized curvature**

↑  
Data dependent constant

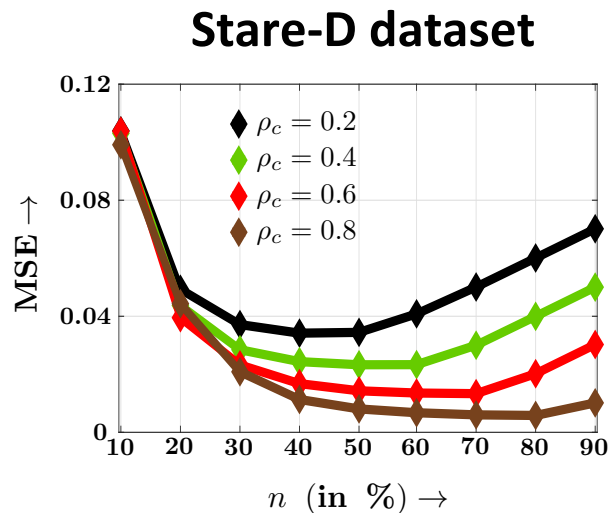
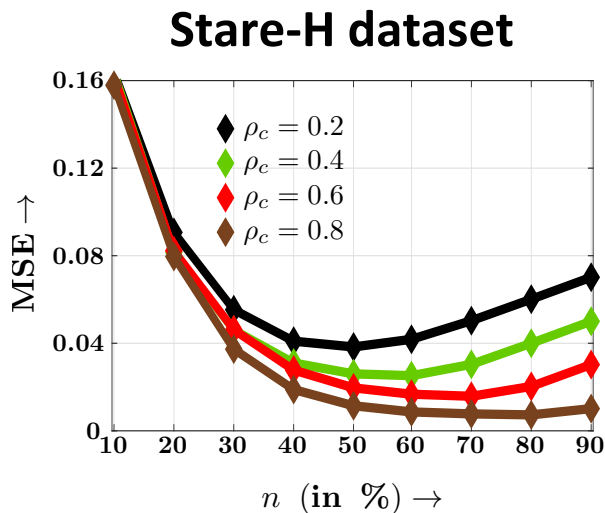
We can conclude that the greedy algorithm will find a set  $\mathcal{S}$  such that

$$-\log \ell(\mathcal{S}) \geq \left(1 + \frac{1}{1 - \alpha}\right)^{-1} OPT$$

Optimal value

[Gatmiry & Gomez Rodriguez, 18  
ACM TOIS 2021]

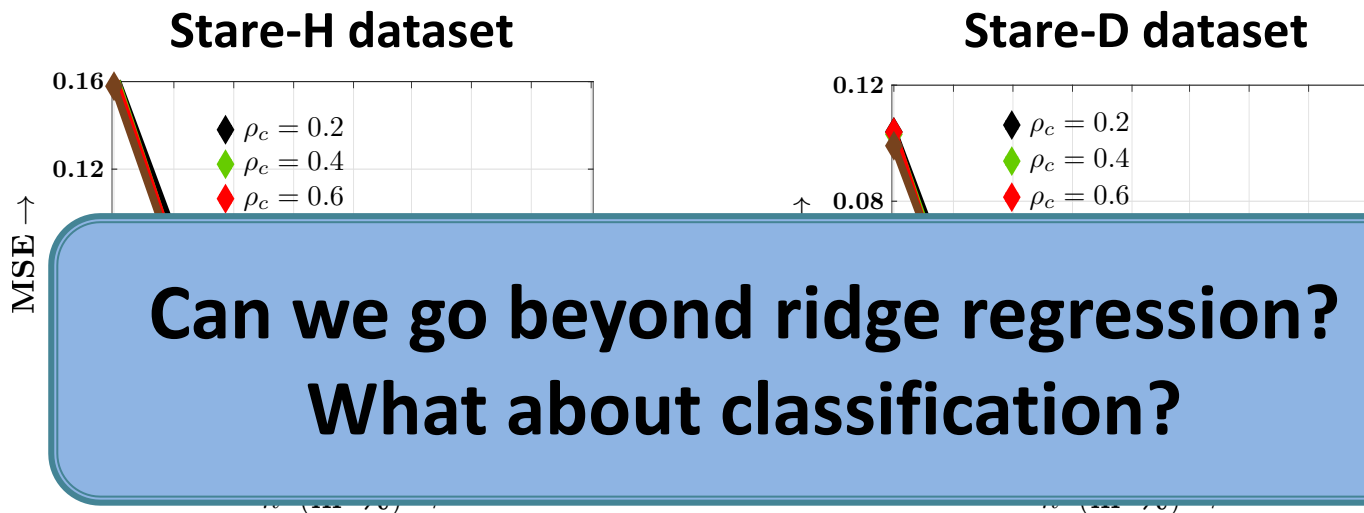
# The greedy algorithm spots samples where humans are accurate



**$\rho_c$ : fraction of samples with low human error**

As long as there are samples that humans can predict with low error, the greedy algorithm outsources them to humans and the performance improves

# The greedy algorithm spots samples where humans are accurate



As long as there are **samples that humans can predict with low error**, the **greedy algorithm outsources them to humans** and the **performance improves**

# Convex margin-based classifiers, revisited

minimize  $\mathcal{S}, \mathbf{w}, b$

Regularization parameter  $\lambda$  points to  $\lambda \|\mathbf{w}\|^2$

Machine model  $\mathbf{w}^\top \Phi(\mathbf{x}_i) + b$  points to  $h_{\mathbf{w}, b}(\mathbf{x}_i)$

Human predictions  $h(\mathbf{x}_i)$  points to  $h(\mathbf{x}_i)$

$$\sum_{i \in \mathcal{V} \setminus \mathcal{S}} [\lambda \|\mathbf{w}\|^2 + [1 - y_i (\mathbf{w}^\top \Phi(\mathbf{x}_i) + b)]_+] + \sum_{i \in \mathcal{S}} [1 - y_i h(\mathbf{x}_i)]_+$$

$\ell(h_{\mathbf{w}, b}(\mathbf{x}_i), y_i)$  (under the first sum)

$c(\mathbf{x}_i, y_i)$  (under the second sum)

subject to  $|\mathcal{S}| \leq n$

Training samples assigned to machines (red dots) points to  $i \in \mathcal{V} \setminus \mathcal{S}$

Training samples assigned to humans (blue dots) points to  $i \in \mathcal{S}$

Max. number of samples that can be assigned to humans points to  $n$

**Key idea:** measure the error in the human predictions (or, scores) using the same (hinge) loss used in the machine model

# Convex optimization + combinatorial optimization

Given a fixed set  $\mathcal{S}$ , we can find the **optimal machine model** using convex programming:

$$\mathbf{w}^*(\mathcal{V} \setminus \mathcal{S}), b^*(\mathcal{V} \setminus \mathcal{S}) = \operatorname{argmin}_{\mathbf{w}, b} \sum_{i \in \mathcal{V} \setminus \mathcal{S}} [\lambda \|\mathbf{w}\|^2 + (1 - y_i(\mathbf{w}^\top \Phi(\mathbf{x}_i) + b))_+]$$

Then, we can **rewrite** the **optimization problem** as the maximization of the difference of two set functions:

$$\underset{\mathcal{S}}{\text{maximize}} \quad g(\mathcal{S}) - c(\mathcal{S}) \quad \text{subject to} \quad |\mathcal{S}| \leq n$$

# Convex optimization + combinatorial optimization

We can show that the **two set functions** satisfy several **desirable properties**:

$$\underset{\mathcal{S}}{\text{maximize}} \quad \underbrace{g(\mathcal{S})}_{\substack{\text{Non-negative} \\ \text{monotone} \\ \gamma\text{-weakly} \\ \text{submodular}}} - \underbrace{c(\mathcal{S})}_{\substack{\text{Non-negative} \\ \text{modular}}} \quad \text{subject to} \quad |\mathcal{S}| \leq n$$

As a result, a **distorted greedy algorithm** is guaranteed to find a set  $\mathcal{S}$  such that:

$$g(\mathcal{S}) - c(\mathcal{S}) \geq (1 - e^{-\gamma})g(\text{OPT}) - c(\text{OPT})$$

We can upper-bound this constant

# Differentiable learning under triage

So far, we “solved” the problem of **learning under triage for ridge regression and SVMs**.

- ↳ The solutions are specialized and it is not easy to extend the resulting methodology to, e.g., deep learning.



# Differentiable learning under triage

So far, we “solved” the problem of **learning under triage for ridge regression and SVMs**.

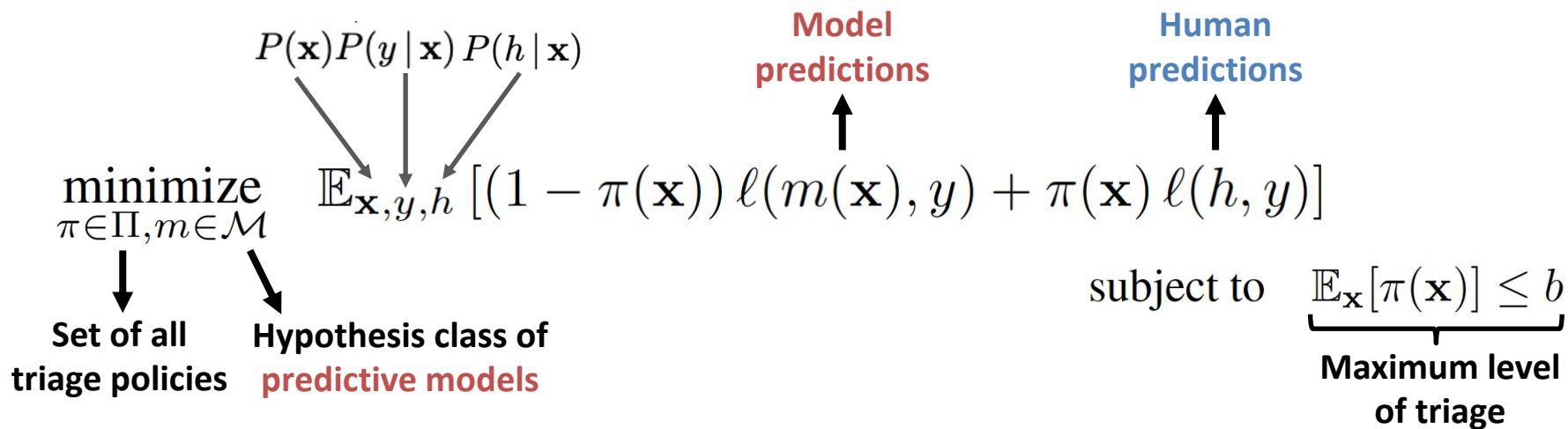
- ↳ The solutions are specialized and it is not easy to extend the resulting methodology to, e.g., deep learning.

Next, we design a **principled & scalable method** that **builds** upon **optimization methods (e.g., SGD)** used in deep learning

- It does not increase the complexity of vanilla SGD
- It is very easy to implement
- It is guaranteed to converge to a local minimum

# Learning under triage, revisited

We look for the **triage policy**  $\pi(\mathbf{x})$  and **predictive model**  $m(\mathbf{x})$  that minimize a loss function  $\ell(\hat{y}, y)$

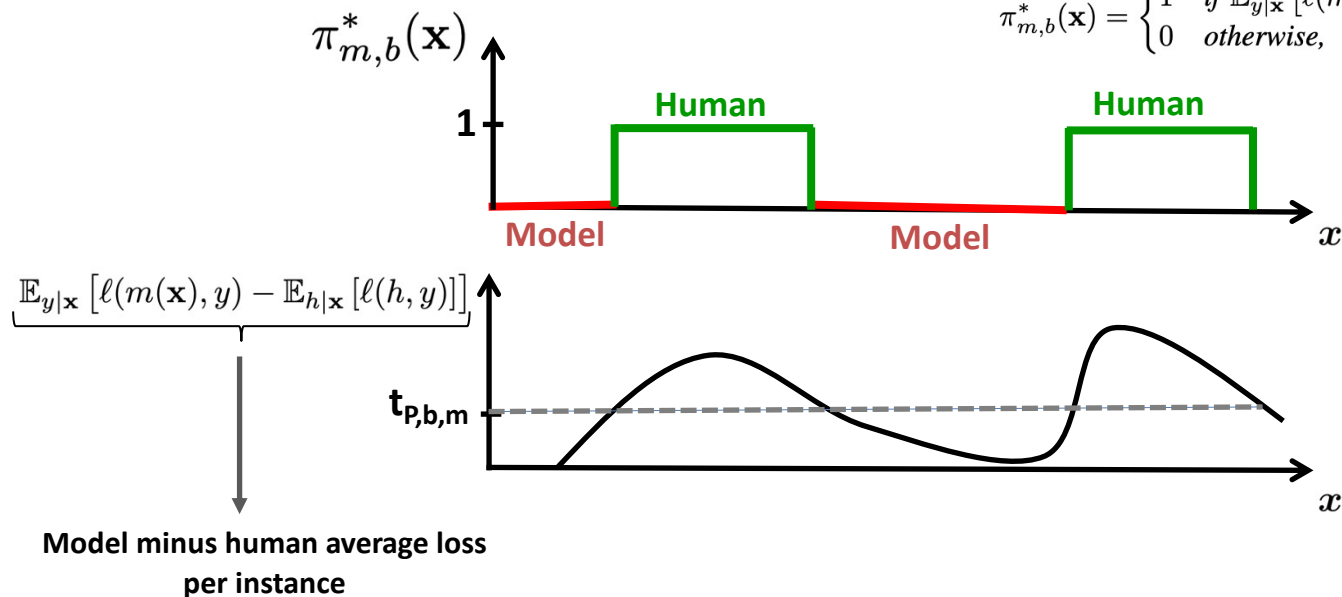


Triage policy  $\begin{cases} \pi(\mathbf{x}) = 0 & \text{Model predicts sample } \mathbf{x} \\ \pi(\mathbf{x}) = 1 & \text{Human predicts sample } \mathbf{x} \end{cases}$

# The optimal triage policy is a threshold rule

Let  $m \in \mathcal{M}$  be any fixed predictive model. Then, the optimal triage policy  $\pi_{m,b}^*(\mathbf{x})$  is a **deterministic threshold rule**:

$$\pi_{m,b}^*(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbb{E}_{y|\mathbf{x}} [\ell(m(\mathbf{x}), y) - \mathbb{E}_{h|\mathbf{x}} [\ell(h, y)]] > t_{P,b,m} \\ 0 & \text{otherwise,} \end{cases}$$



# When is a predictive model suboptimal under triage?

Let  $m_{\theta_0^*}$  the **optimal predictive model under full automation**

Average gradient on all samples

$$\rightarrow \mathbb{E}_{\mathbf{x}, y} \left[ \nabla_{\theta} \ell(m_{\theta}(\mathbf{x}), y) \Big|_{\theta=\theta_0^*} \right] = \mathbf{0}$$

# When is a predictive model suboptimal under triage?

Let  $m_{\theta_0^*}$  the **optimal predictive model under full automation**



$$\overbrace{\mathbb{E}_{\mathbf{x}, y} \left[ \nabla_{\theta} \ell(m_{\theta}(\mathbf{x}), y) \Big|_{\theta=\theta_0^*} \right]}^{\text{Average gradient on all samples}} = \mathbf{0}$$

Let  $\pi_{m_{\theta_0^*}, b}^*$  its **optimal triage policy** under max. triage level  $b$ . Then, if

$$\underbrace{\int_{\mathbf{x} \in \mathcal{V}} \mathbb{E}_{y|\mathbf{x}} \left[ \nabla_{\theta} \ell(m_{\theta}(\mathbf{x}), y) \Big|_{\theta=\theta_0^*} \right] dP}_{\text{Average gradient on the set of samples } \mathcal{V} \text{ that } \pi_{m_{\theta_0^*}, b}^* \text{ lets the model predicts}} \neq \mathbf{0}$$

Average gradient on the set of samples  $\mathcal{V}$  that  $\pi_{m_{\theta_0^*}, b}^*$  lets the model predicts

# When is a predictive model suboptimal under triage?

Let  $m_{\theta_0^*}$  the **optimal predictive model under full automation**

Average gradient on all samples

$$\mathbb{E}_{\mathbf{x}, y} \left[ \nabla_{\theta} \ell(m_{\theta}(\mathbf{x}), y) \Big|_{\theta=\theta_0^*} \right] = \mathbf{0}$$

Let  $\pi_{m_{\theta_0^*}, b}^*$  its **optimal triage policy** under max. triage level  $b$ . Then, if

$$\int_{\mathbf{x} \in \mathcal{V}} \mathbb{E}_{y|\mathbf{x}} \left[ \nabla_{\theta} \ell(m_{\theta}(\mathbf{x}), y) \Big|_{\theta=\theta_0^*} \right] dP \neq \mathbf{0}$$

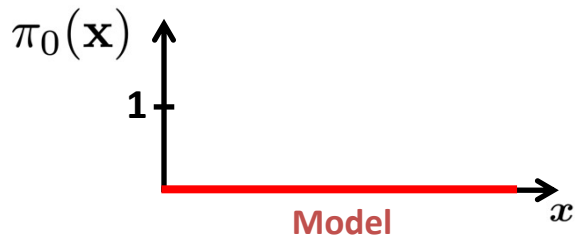
Average gradient on the set of samples  $\mathcal{V}$  that  $\pi_{m_{\theta_0^*}, b}^*$  lets the model predicts

Then, **there exists a better predictive model** under  $\pi_{m_{\theta_0^*}, b}^*$  :

$$L(\pi_{m_{\theta_0^*}, b}^*, m_{\theta_0^*}) > \min_{\theta \in \Theta} L(\pi_{m_{\theta}, b}^*, m_{\theta})$$

# Augmenting SGD to learn under triage

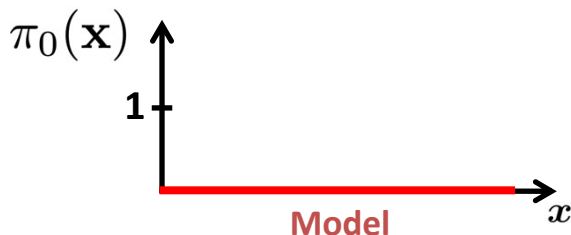
Step  $t = 0$



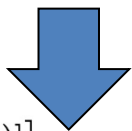
Train predictive model  $m_{\theta_0^*}$   
using SGD with samples for which  
 $\pi_0(\mathbf{x}) = 0$

# Augmenting SGD to learn under triage

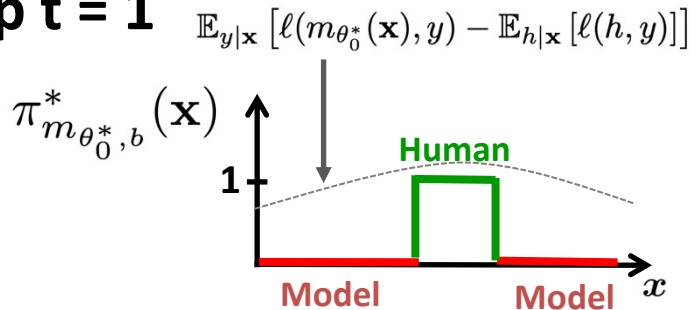
Step  $t = 0$



Train predictive model  $m_{\theta_0^*}$   
using SGD with samples for which  
 $\pi_0(\mathbf{x}) = 0$



Step  $t = 1$

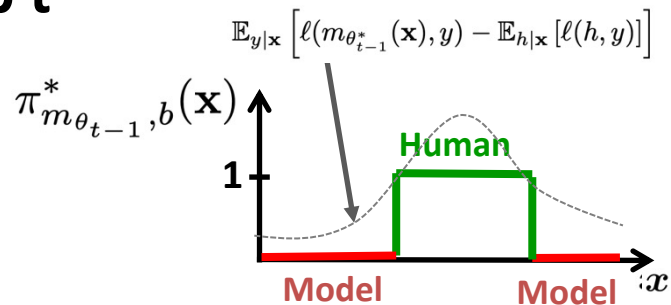


Train predictive model  $m_{\theta_1^*}$   
using SGD with samples for which  
 $\pi_{m_{\theta_0^*, b}^*}(\mathbf{x}) = 0$



# Augmenting SGD to learn under triage

## Step t



Train predictive model  $m_{\theta_t}$   
using SGD with samples for which

$$\pi_{m_{\theta_{t-1}}, b}^*(\mathbf{x}) = 0$$

# Local convergence guarantees of SGD under triage

Under **mild conditions**, it holds that:

$$\underbrace{L(\pi_{m_{\theta_{t-1}}, b}^*, m_{\theta_{t-1}}) < L(\pi_{m_{\theta_t}, b}^*, m_{\theta_t})}$$

**The performance of the triage policies and predictive models improves in each step**

**An implementation of the algorithm with Monte-Carlo estimates converges to a local minimum of the empirical loss**

# Global guarantees for convex losses

Under **convex losses**, it holds that:

$$\lim_{t \rightarrow \infty} L(\pi_{m_{\theta_t}, b}^*, m_{\theta_t}) - L(\pi_{m_{\theta^*}, b}^*, m_{\theta^*}) \leq \frac{4H^2 \Lambda_{\max}}{\Lambda_{\min}^2 (1-b)^2}$$

where

$$\ell(m_{\theta}(\mathbf{x}), y) - \ell(m_{\theta'}(\mathbf{x}), y) \leq H \cdot \|\theta - \theta'\|$$

$$\Lambda_{\min} \mathbb{I} \preceq \nabla_{\theta}^2 \ell(m_{\theta}(\mathbf{x}), y) \preceq \Lambda_{\max} \mathbb{I}$$

↑  
Maximum level  
of triage

# SGD under triage at test time

Let  $m_{\theta_T}$  and  $\pi_{m_{\theta_T}, b}^*(\mathbf{x})$  the last predictive model and its optimal triage policy.

At test time, we have a problem:

we cannot compute the triage value  $\pi_{m_{\theta_T}, b}^*(\mathbf{x})$  for unseen samples. We do not know the value of the model/human loss!

# SGD under triage at test time

Let  $m_{\theta_T}$  and  $\pi_{m_{\theta_T}, b}^*(\mathbf{x})$  the last predictive model and its optimal triage policy.

At test time, we have a problem:

we cannot compute the triage value  $\pi_{m_{\theta_T}, b}^*(\mathbf{x})$  for unseen samples. We do not know the value of the model/human loss!

Solution:

we train a parameterized triage policy  $\hat{\pi}_\gamma(\mathbf{x})$  to approximate the optimal triage policy  $\pi_{m_{\theta_T}, b}^*(\mathbf{x})$

# Learning under triage on real datasets

**Few public datasets with several human predictions per instance, necessary to estimate the human loss per instance.**

# Learning under triage on real datasets

Few public datasets with several human predictions per instance, necessary to estimate the **human loss per instance**.

Experiments with **two datasets** from **applications in content moderation** and **scientific discovery**:

→ **Hatespeech:** 25k tweets containing lexicons used in hate speech



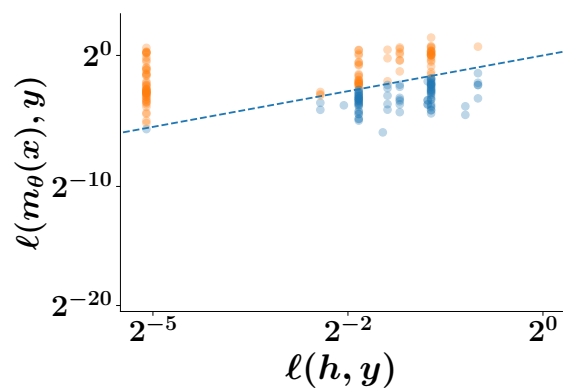
3-5 human predictions per instance  
labels = {hate-speech, offensive, neither}

→ **Galaxy zoo:**

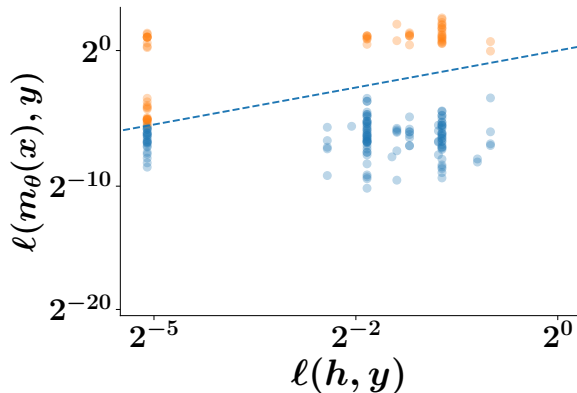


10k galaxy images  
30+ human predictions per instance  
labels = {early-type, spiral}

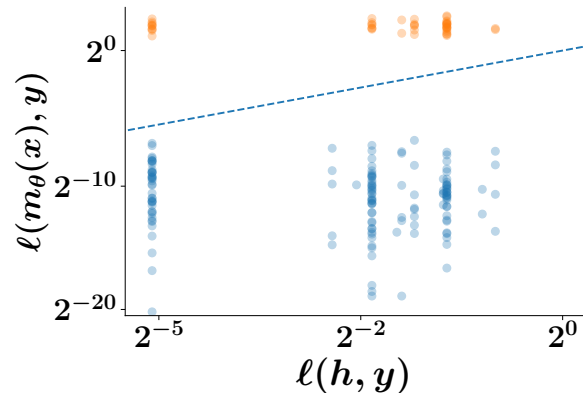
# Model vs human loss on the Hatespeech dataset



$t = 1$



$t = 2$



$t = 3$

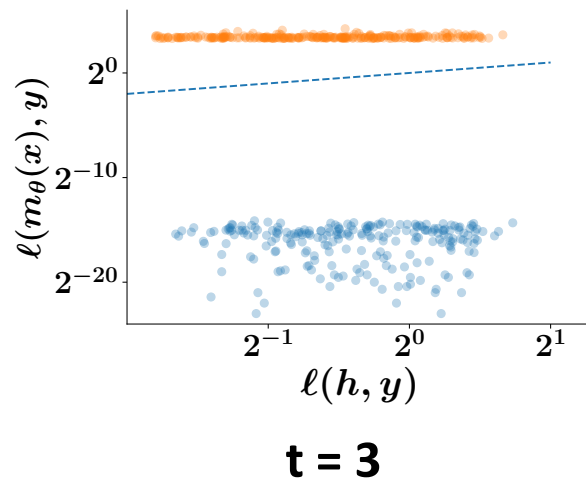
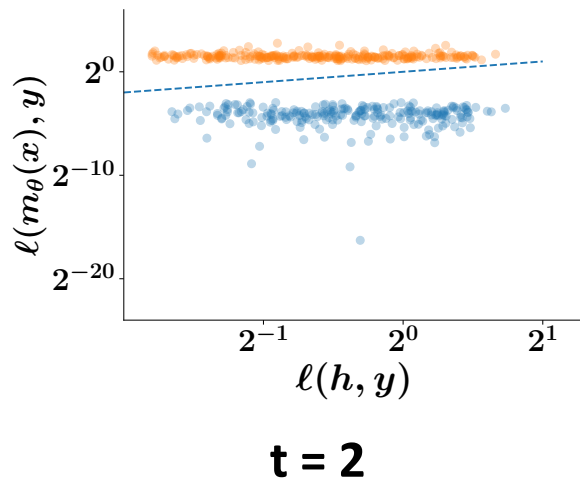
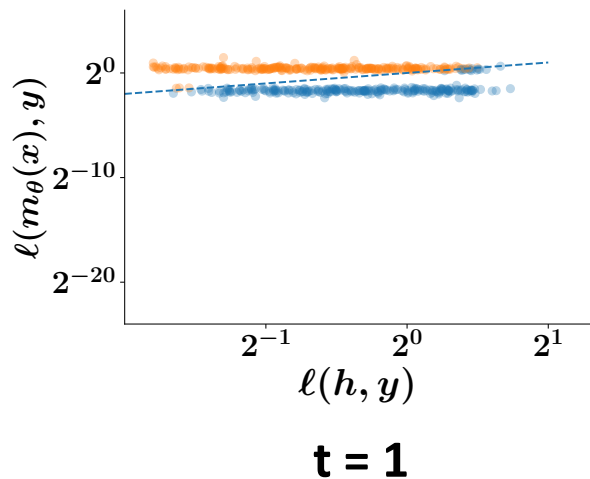
triage level  $\leq 40\%$



The model learns to specialize in a **region of the feature space**, where it achieves **low loss**, and gives up on **the rest**, where the **loss is very high**



# Model vs human loss on the Galaxy zoo dataset



triage level  $\leq 100\%$



The model learns to specialize in a **region of the feature space**, where it achieves **low loss**, and gives up on **the rest**, where the **loss is very high**

# Where do we go from here?

- We have focused on **independent samples**. **Algorithmic triage in sequential decision making (e.g. semi-autonomous driving)**.
- We have assumed **each instance is either predicted by a model or by a human**. **All instances predicted by humans, who are informed by a model**.
- **Validate algorithmic triage using interventional experiments.**

# Thanks!

## **Regression under human assistance, AAAI 2020**

<https://arxiv.org/abs/1909.02963>

<https://github.com/Networks-Learning/regression-under-assistance>

## **Classification under human assistance, AAAI 2021**

<https://arxiv.org/abs/2006.11845>

<https://github.com/Networks-Learning/classification-under-assistance>

## **Differentiable learning under triage, NeurIPS 2021**

<https://arxiv.org/abs/2103.08902>

<https://github.com/Networks-Learning/differentiable-learning-under-triage>