# 1. Background Motivation

**How can ML algorithms learn to classify words without exhausting human patience?**

# 1. Background Motivation

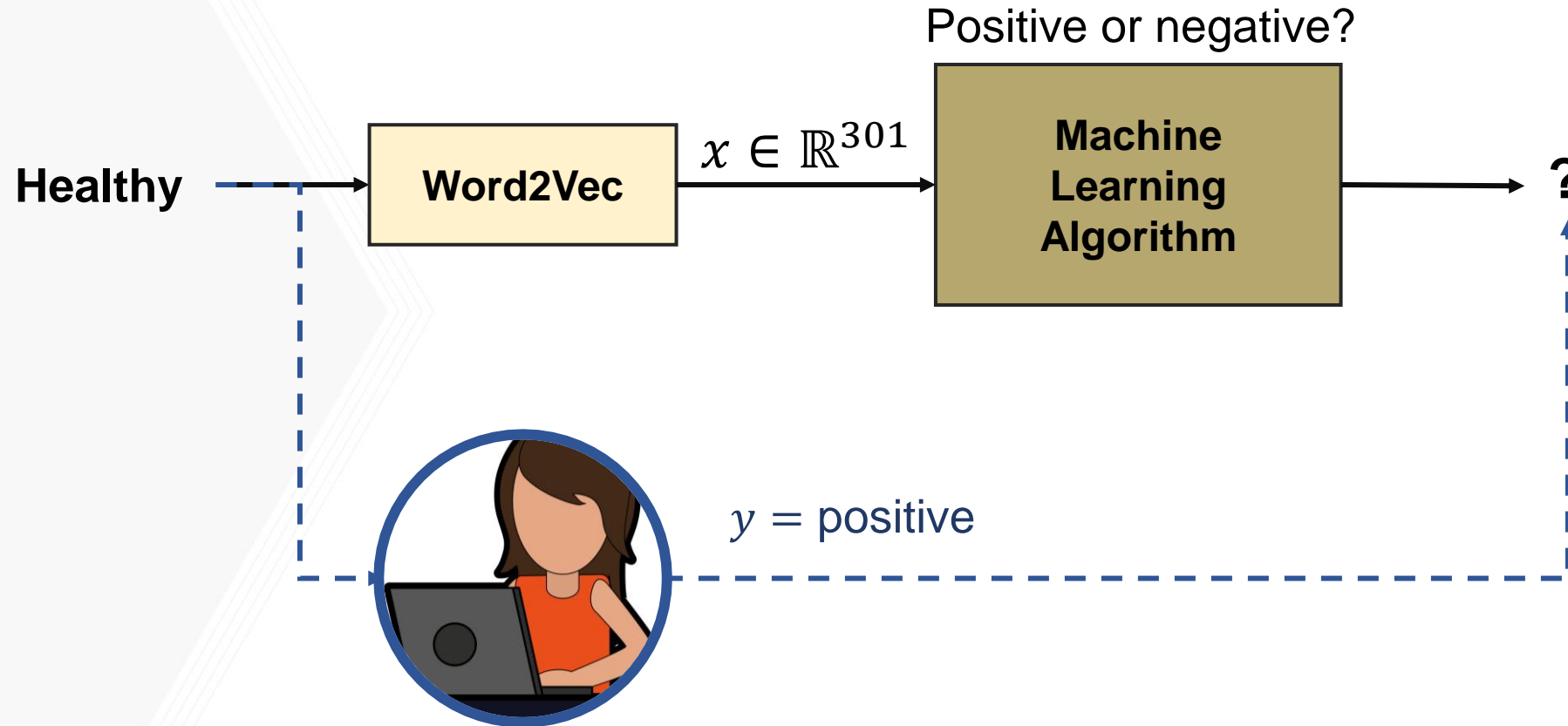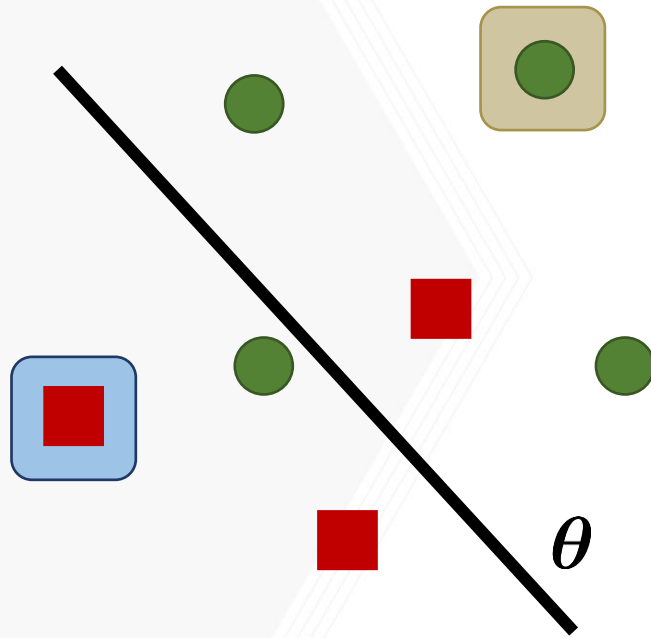**How can ML algorithms learn to classify words without exhausting human patience?**



□ Could we get **more information** from the human?

# 2. More Information

Queries must be <u>human</u> and <u>mathematically</u> interpretable



$\theta$

Legend:
- ● : Positive word embedding
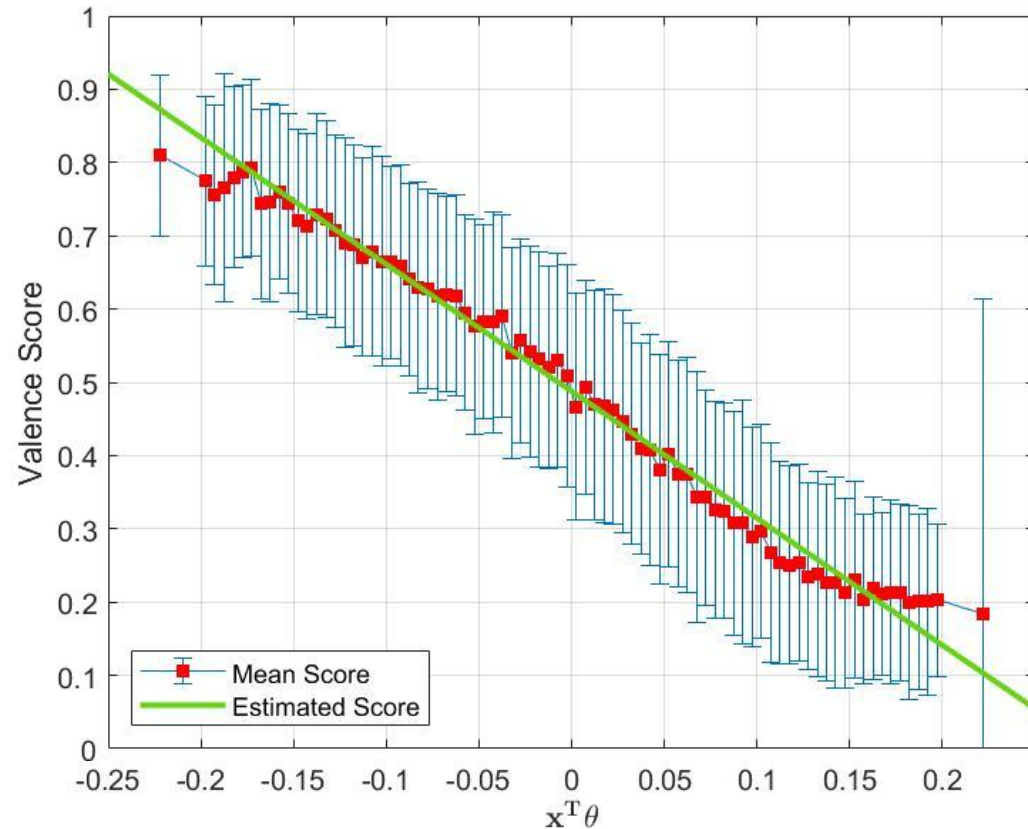- ■ : Negative word embedding
- $\theta$ : Word Classifier

"Select the example that you find…"

- Query: "… most **positive**"

- Query: "… most **negative**"

**Hypothesis:** Human answers depend on the distance to the ground truth

# 2. **More Information** Valence vs Distance to θ

**NRC-VAD Lexicon**



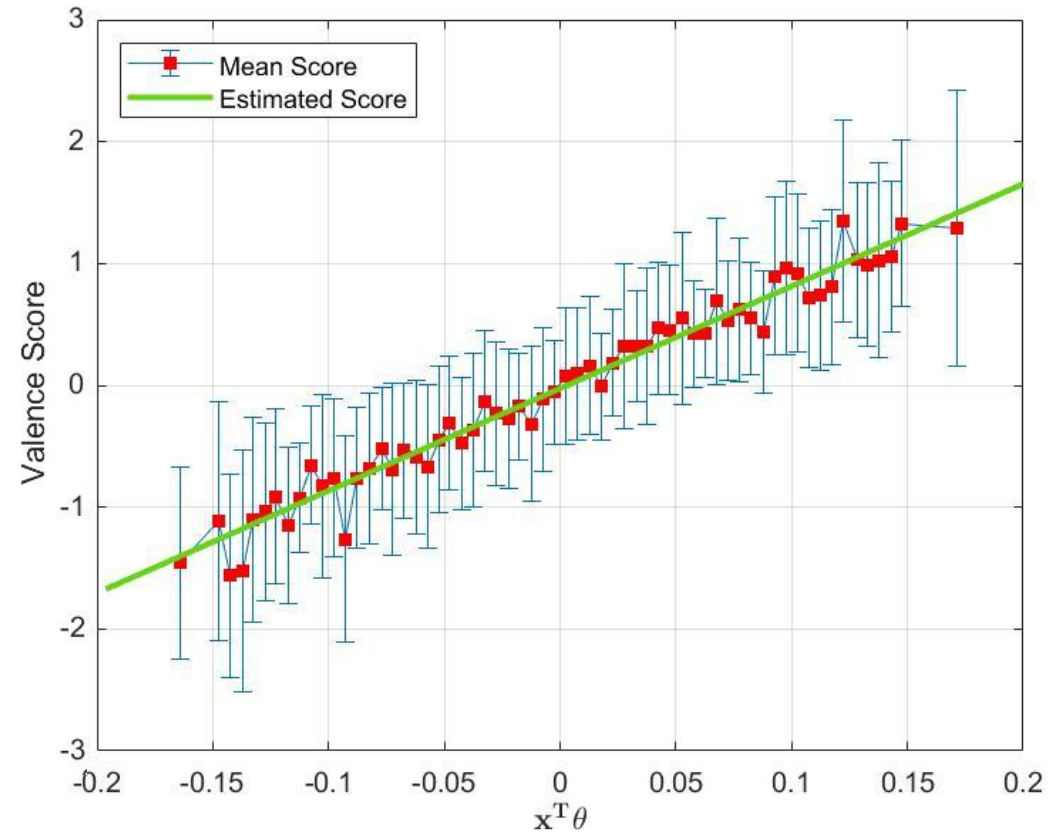**SocialNet**



$$\mathbb{E}\left[\text{score}(\mathbf{x})|\boldsymbol{\theta}\right] = -1.73(\mathbf{x}^T\boldsymbol{\theta}) + 0.49$$

$$\mathbb{E}\left[\text{score}(\mathbf{x})|\boldsymbol{\theta}\right] = 6.96(\mathbf{x}^T\boldsymbol{\theta}) - 0.09$$

$\mathbf{x}$: Word embedding    $\boldsymbol{\theta}$: Ground truth classifier

Georgia
Tech
CREATING THE NEXT

# 2. **More Information** Multinomial Logit Model

Queries must be <u>human</u> and <u>mathematically</u> interpretable
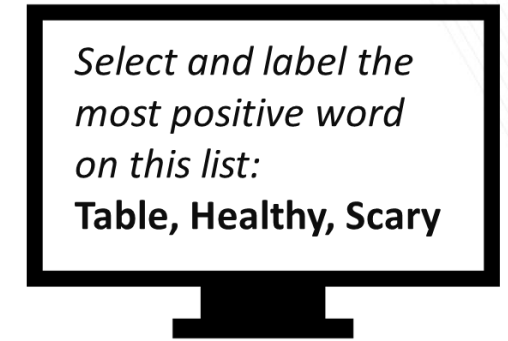
**"Select and label the most positive/negative word"**

$$\mathbb{P}\left[y = 1 | \mathbf{x}\right] = \frac{1}{1 + \exp\left(W(\boldsymbol{\theta}^T \mathbf{x})\right)}, \text{ with } W \in \mathbb{R}.$$

$$\mathbb{P}\left[\mathbf{x}_i | \{\mathbf{x}_j\}_{j=1}^N, \boldsymbol{\theta}\right] = \frac{\exp\left(u(\mathbf{x}_i, \boldsymbol{\theta})\right)}{\sum_{j=1}^N \exp\left(u(\mathbf{x}_j, \boldsymbol{\theta})\right)}$$

**Likelihood**

Select and label the
most positive word
on this list:
**Table, Healthy, Scary**

'Healthy'
positive

- Q1: Most positive

$$u(\mathbf{x}, \boldsymbol{\theta}) = \frac{a}{\sigma} \mathbf{x}^T \boldsymbol{\theta}$$

- Q2: Most negative

$$u(\mathbf{x}, \boldsymbol{\theta}) = -\frac{a}{\sigma} \mathbf{x}^T \boldsymbol{\theta}$$

$\mathbf{x}$: Word embedding
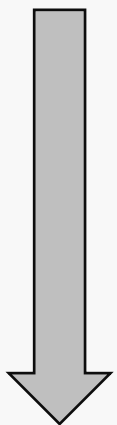
$\boldsymbol{\theta}$: Ground truth classifier

$y$: Word label (positive: 1, negative: 0)

$u$: Utility function (expected valence score)

Georgia
Tech
CREATING THE NEXT

# 2. **More Information** Multinomial Logit Model

$$\mathbb{P}\left[y = 1 | \mathbf{x}\right] = \frac{1}{1 + \exp\left(W(\boldsymbol{\theta}^T \mathbf{x})\right)}, \text{ with } W \in \mathbb{R}.$$

$$\mathbb{P}\left[\mathbf{x}_i | \{\mathbf{x}_j\}_{j=1}^N, \boldsymbol{\theta}\right] = \frac{\exp\left(u(\mathbf{x}_i, \boldsymbol{\theta})\right)}{\sum_{j=1}^N \exp\left(u(\mathbf{x}_j, \boldsymbol{\theta})\right)}$$

**Likelihood**

How can we get the posterior?

$$p_{\boldsymbol{\theta}} = \mathbb{P}\left[\boldsymbol{\theta} | \{\mathbf{x}_t, y_t, q_t, \{\mathbf{x}_j\}_{j=1}^N\}_{t=0}^i\right]$$

→ Approximate $\boldsymbol{\theta}$ as a multivariate Gaussian

**Posterior over** $\theta$

$\mathbf{x}$: Word embedding     $q$: Query     $u$: Utility function (expected valence score)

$\boldsymbol{\theta}$: Ground truth classifier     $y$: Word label

Georgia Tech
CREATING THE NEXT

# 2. **More Information** Update Given Label

$$\mathbb{P}\left[y = 1 | \mathbf{x}\right] = \frac{1}{1 + \exp\left(W(\boldsymbol{\theta}^T \mathbf{x})\right)}, \text{ with } W \in \mathbb{R}.$$

<u>How do we update the posterior given the label?</u>

Jaakkola and Jordan give a closed form approximation*

$$\Sigma_{\text{pos}}^{-1} = \Sigma^{-1} + 2\frac{\tanh(\xi/2)}{4\xi}W^2\mathbf{x}_i\mathbf{x}_i^T$$

$$\boldsymbol{\mu}_{\text{pos}} = \Sigma_{\text{pos}}\left[\Sigma^{-1}\mu + \left(y_i - \tfrac{1}{2}\right)W\mathbf{x}_i\right]$$

$$\xi^2 = W^2\mathbf{x}_i^T\Sigma_{\text{pos}}\mathbf{x}_i + W^2(\mathbf{x}_i^T\boldsymbol{\mu}_{\text{pos}})^2$$

*Bayesian parameter estimation via variational methods – Jaakkola and Jordan, 2000*

$\mathbf{x}$: Word embedding    Prior: $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$    $\boldsymbol{\theta}$: Ground truth classifier
$y$: Word label    Posterior: $\mathcal{N}(\boldsymbol{\mu}_{\text{pos}}, \boldsymbol{\Sigma}_{\text{pos}})$

Georgia Tech
CREATING THE NEXT

# 2. **More Information** Update Given Word

$$\mathbb{P}\left[\mathbf{x}_i | \{\mathbf{x}_j\}_{j=1}^{N}, \boldsymbol{\theta}\right] = \frac{\exp\left(K\mathbf{x}_i^T\boldsymbol{\theta}\right)}{\sum_{j=1}^{N}\exp\left(K\mathbf{x}_j^T\boldsymbol{\theta}\right)}$$

How do we update the posterior given the word selected?

$$ELBO(q) = -\mathrm{KL}(q(\boldsymbol{\theta})\|p(\boldsymbol{\theta})) + \mathbb{E}_{\boldsymbol{\theta}\sim q}\left[K\mathbf{x}_s^T\boldsymbol{\theta}\right] - \mathbb{E}_{\boldsymbol{\theta}\sim q}\left[\log\sum_{j=1}^{|\mathcal{S}|}\exp\left(K\mathbf{x}_j^T\boldsymbol{\theta}\right)\right]$$

$$KL(q\|p) = \frac{1}{2}\left[\log\frac{|\Sigma_p|}{|\Sigma_q|} - d + (\boldsymbol{\mu_q})^T\Sigma_p^{-1}(\boldsymbol{\mu_q}) + (\boldsymbol{\mu_p})^T\Sigma_p^{-1}(\boldsymbol{\mu_p}) - 2(\boldsymbol{\mu_q})^T\Sigma_p^{-1}(\boldsymbol{\mu_p}) + tr\left\{\Sigma_p^{-1}\Sigma_q\right\}\right]$$

$$\mathbb{E}_{\boldsymbol{\theta}\sim q}\left[K\mathbf{x}_s^T\boldsymbol{\theta}\right] = K\mathbf{x}_s^T\boldsymbol{\mu}_q$$

$$\mathbb{E}_{\boldsymbol{\theta}\sim q}\left[\log\sum_{j=1}^{|\mathcal{S}|}\exp\left(K\mathbf{x}_j^T\boldsymbol{\theta}\right)\right] \geq \log\sum_{j=1}^{|\mathcal{S}|}\exp\left(K\mathbf{x}_j^T\boldsymbol{\mu}_q + 0.5\mathbf{x}_j^T\boldsymbol{\Sigma}_q\mathbf{x}_j\right)$$ *[Braun and McAuliffe, 2007]*
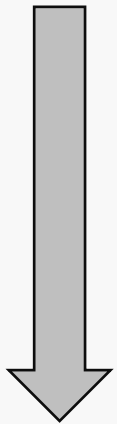
Slide 9 of 19

$p(\boldsymbol{\theta}) \sim \mathcal{N}(\boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p)$: Prior $\qquad q(\boldsymbol{\theta}) \sim \mathcal{N}(\boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q)$: Variational Distribution

Georgia
Tech
CREATING THE NEXT

# 2. **More Information** Multinomial Logit Model

$$\mathbb{P}\left[\mathbf{x}_i \middle| \{\mathbf{x}_j\}_{j=1}^{N}, \boldsymbol{\theta}\right] = \frac{\exp\left(u(\mathbf{x}_i, \boldsymbol{\theta})\right)}{\sum_{j=1}^{N} \exp\left(u(\mathbf{x}_j, \boldsymbol{\theta})\right)}$$

$$\mathbb{P}\left[y = 1 \middle| \mathbf{x}\right] = \frac{1}{1 + \exp\left(W(\boldsymbol{\theta}^T \mathbf{x})\right)}, \text{ with } W \in \mathbb{R}.$$

**Likelihood**

**Variational inference**

**Posterior over** $\theta$
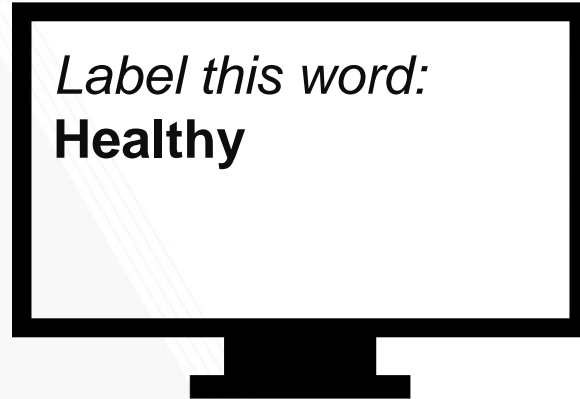
$\mathbf{x}$: Word embedding     $y$: Word label (positive: 1, negative: 0)

$\boldsymbol{\theta}$: Ground truth classifier     $u$: Utility function (expected valence score)

Georgia Tech

CREATING THE NEXT

# 2. **More Information** Experiments

Label:



*Label this word:*
**Healthy**

positive

$s_{\text{healthy}} > 0.5$

Label + Word:

*Select and label the most positive word?*
**Healthy, table, scary, orange**

'Healthy' positive

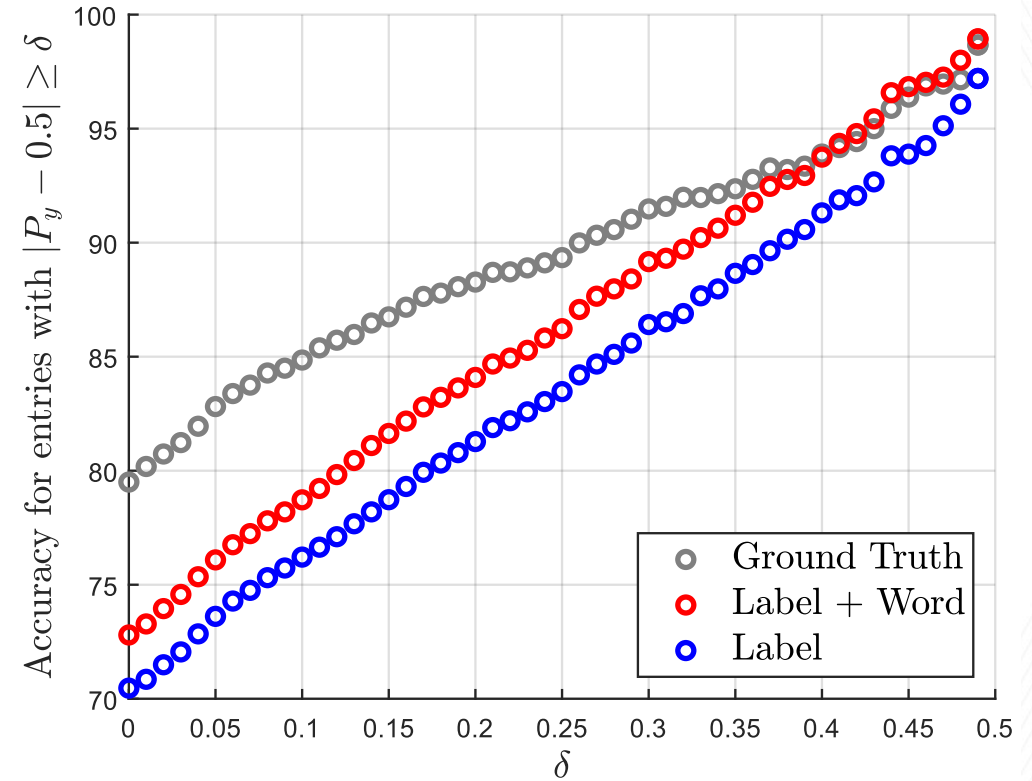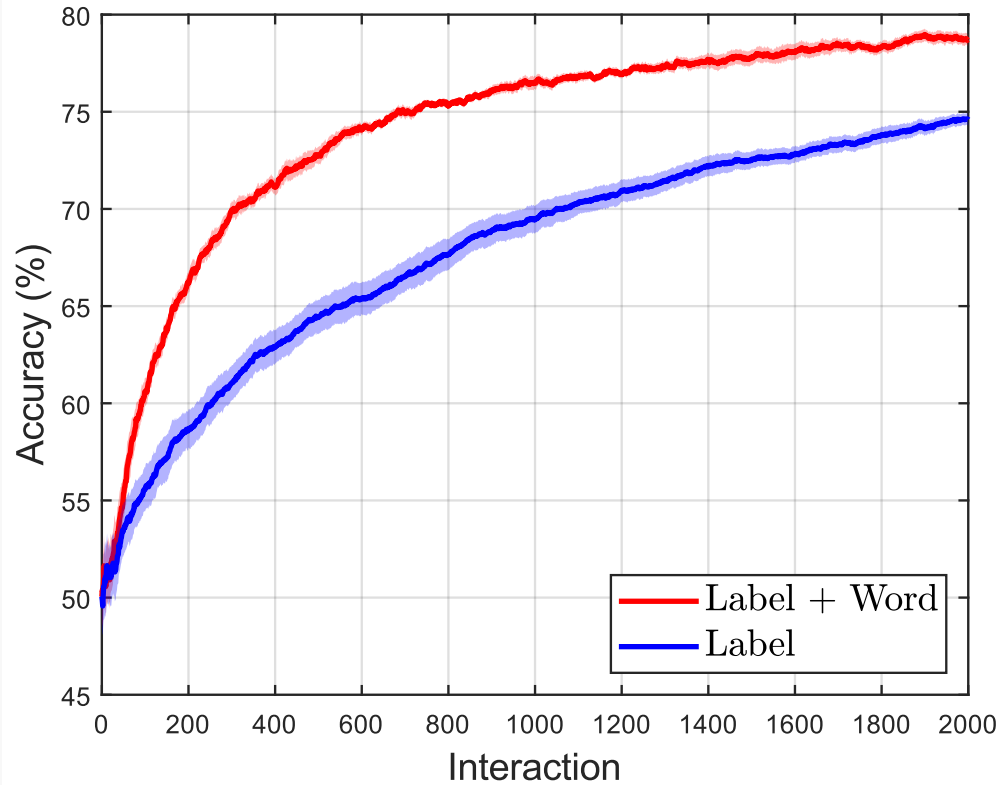$s_{\text{healthy}} > s_{\text{table}}, s_{\text{scary}}, s_{\text{orange}}$
$s_{\text{healthy}} > 0.5$

**NRC-VAD Lexicon dataset:**
Each word $\mathbf{x}$ has a valence score mean $\mu_{\mathbf{x}}$ and variance $\sigma^2_{\mathbf{x}} \rightarrow$ score: $s_{\mathbf{x}} \sim \mathcal{N}(\mu_{\mathbf{x}}, \sigma^2_{\mathbf{x}})$

**Georgia Tech**
CREATING THE NEXT

# 2. **More Information** Empirical Results

# 2. More Information Theoretical Results

Under Assumptions
1. Word and label selected are independent of the history given the classifier

$$p(\mathbf{x}_t, y_t | \boldsymbol{\theta}, q_t, \mathcal{S}_t, \mathcal{F}_{t-1}) = p(\mathbf{x}_t, y_t | \boldsymbol{\theta}, q_t, \mathcal{S}_t)$$

2. The label only depends on the word it is referring to

$$p(y_t | \mathbf{x}_t, q_t, \mathcal{S}_t) = p(y_t | \mathbf{x}_t)$$

3. An answer always provides some information

$$I(\boldsymbol{\theta}; X_t, Y_t | \mathcal{F}_{t-1}) \geq L > 0$$

Simplified Theorem

The expected stopping time $T_\epsilon = \min\{t : \left|\boldsymbol{\Sigma}_{\boldsymbol{\theta}|\mathcal{F}_t}\right|^{1/d} < \epsilon\}$ is bounded as

$$\frac{d}{2} \frac{\log_2 \frac{2}{\pi e \epsilon}}{\log_2 2|\mathcal{S}|} \leq \mathbb{E}[T_\epsilon] \leq \frac{d}{2L} \log_2 \frac{e^4 d^2}{2\sqrt{2}(d+2)\epsilon} - 1.$$

where $\mathbf{x} \in \mathbb{R}^d$ and $\mathcal{S} = \{\mathbf{x}_j\}_{j=1}^{|\mathcal{S}|}$ are the candidate words.

Georgia
Tech
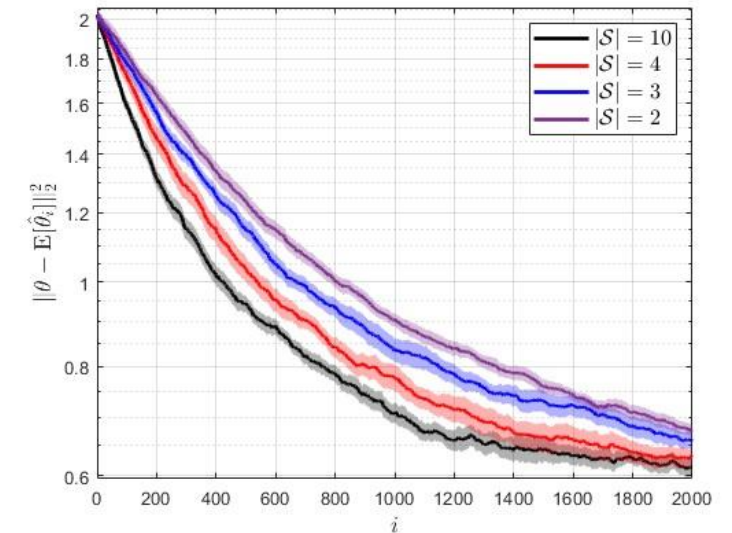
CREATING THE NEXT

# 2. More Information Theoretical Results

Simplified Theorem

The expected stopping time $T_\epsilon = \min\{t : |\Sigma_{\theta|\mathcal{F}_t}|^{1/d} < \epsilon\}$ is bounded as

$$\frac{d}{2}\frac{\log_2 \frac{2}{\pi e \epsilon}}{\log_2 2|\mathcal{S}|} \le \mathbb{E}[T_\epsilon] \le \frac{d}{2L}\log_2 \frac{e^4 d^2}{2\sqrt{2}(d+2)\epsilon} - 1.$$

where $\mathbf{x} \in \mathbb{R}^d$ and $\mathcal{S} = \{\mathbf{x}_j\}_{j=1}^{|\mathcal{S}|}$ are the candidate words.
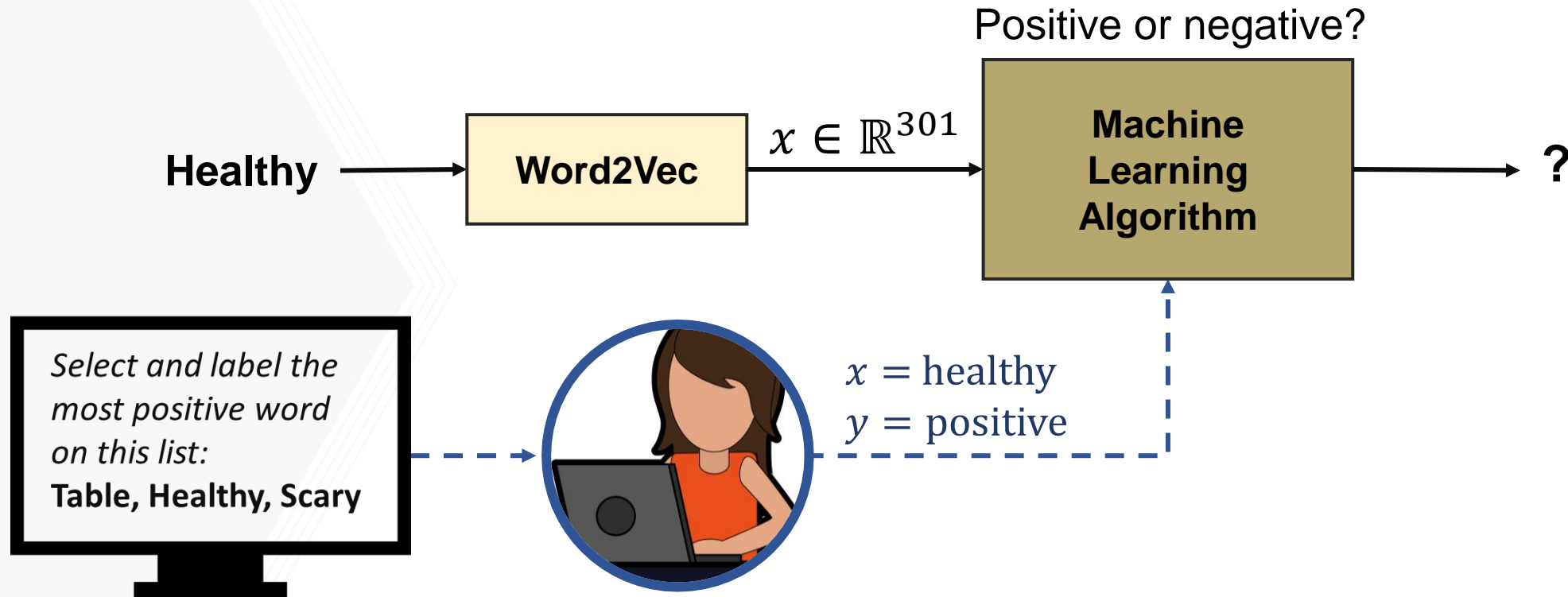
✓ The number of question to ask the human to reach uncertainty $< \epsilon$ is on the order of $\log 1/\epsilon$

✓ Related to the error $\mathrm{MSE}_t = \mathrm{trace}(\Sigma_{\theta|\mathcal{F}_t}) \ge d|\Sigma_{\theta|\mathcal{F}_t}|^{1/d}$

✓ The more words in the list, the faster the error decays

Georgia Tech
CREATING THE NEXT

# 2. More Information Motivation

**How can ML algorithms learn to classify words without exhausting human patience?**
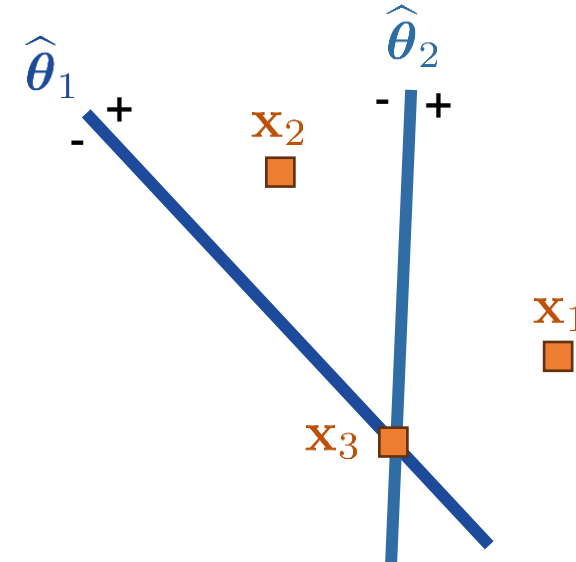


- ✓ Could we get **more information** from the human?
- ☐ Could we get the information **faster**?

# 3. Faster Active Learning Heuristic

💡 Instead of showing the humans a random list of words, could we select them in a smart way?

| Example | $P[+|x, \theta = 1]$ | $P[+|x, \theta = 2]$ |
|:---:|:---:|:---:|
| $x_1$ | 0.99 | 0.99 |
| $x_2$ | 0.99 | 0.01 |
| $x_3$ | 0.5 | 0.5 |



Heuristic in AL:

$$\mathrm{argmax}_{\mathbf{x}} \underbrace{H\left(\mathbb{E}_{\boldsymbol{\theta}}\left[f_{\theta}(\mathbf{x})\right]\right)}_{} - \underbrace{\mathbb{E}_{\theta}\left[H\left(f_{\theta}(\mathbf{x})\right)\right]}_{}$$

→ Maximize uncertainty of the expected output

→ Minimize uncertainty due to noise

$H$: Entropy   $\boldsymbol{\theta}$: Ground truth classifier   $\mathbf{x}$: Word embedding

Georgia Tech
CREATING THE NEXT

# 3. **Faster** Active Word Selection

Heuristic:

$$\mathcal{S} = \text{argmax}_{\mathcal{S} \in \{\mathcal{X}\}^k} \underbrace{H\left(\mathbb{E}_{\boldsymbol{\theta}}\left[\mathbf{x}_i, y_i \mid q, \boldsymbol{\theta}, \mathcal{S}\right]\right)}_{} - \underbrace{\mathbb{E}_{\boldsymbol{\theta}}\left[H\left(\mathbf{x}_i, y_i \mid q, \boldsymbol{\theta}, \mathcal{S}\right)\right]}_{}$$

→ Maximize uncertainty of the expected output

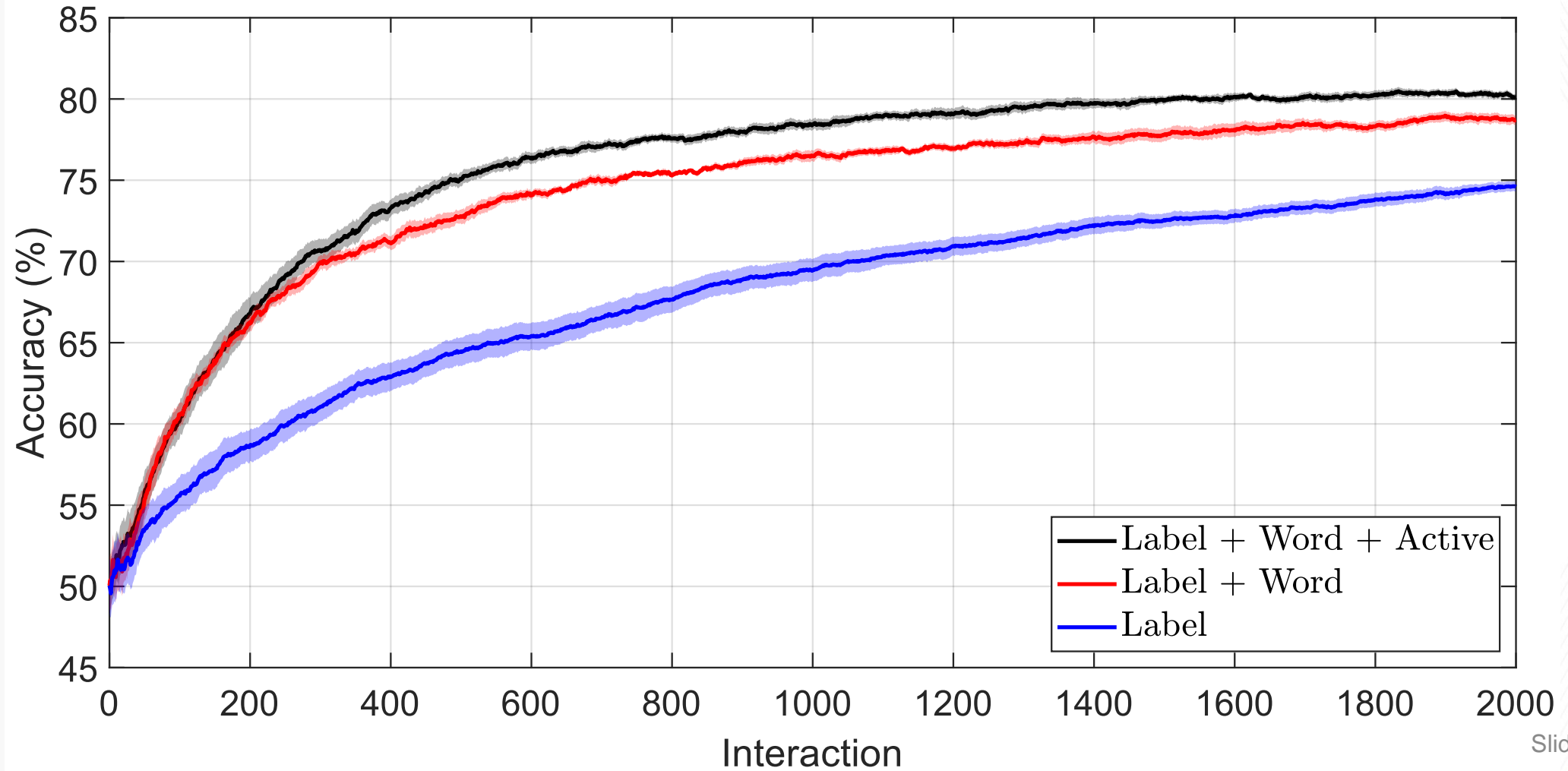→ Minimize uncertainty due to noise

**Problem:** There exist combinatorically many sets $\binom{|\mathcal{X}|}{|\mathcal{S}|}$ to maximize over

**Our Approach:** Greedily select one word at a time

If $|\mathcal{X}| = 3500$ and $|\mathcal{S}| = 4$:
$\binom{|\mathcal{X}|}{|\mathcal{S}|} \sim 10^{15}$
Galaxies $\sim 10^{12}$
If each computation 1ms,
$t \sim 31,7000$ years

$H$: Entropy    $\boldsymbol{\theta}$: Ground truth classifier    $\mathbf{x}$: Word embedding    $y$: Word label

Georgia
Tech
CREATING THE NEXT

# 3. Faster Results

# 4. Summary

1. We introduce a novel **human response model**.

2. We **speed up learning** in sentiment classification
   - By combining label requests with <u>word selection</u>.
   - By <u>active</u> query selection.

3. We **validate** our approach
   - <u>Theoretically</u>: Bounds for expected stopping time.
   - <u>Empirically</u>: Experiments with human data.

Georgia
Tech

CREATING THE NEXT

# 3. Faster Active Learning Heuristic



| Example | $P[+|x, \boldsymbol{\theta} = 1]$ | $P[+|x, \boldsymbol{\theta} = 2]$ |
|---------|-----------------------------------|-----------------------------------|
| $x_1$   | 0.99                              | 0.99                              |
| $x_2$   | 0.99                              | 0.01                              |
| $x_3$   | 0.5                               | 0.5                               |

$$x_1 \rightarrow H(0.99) - \left[0.5\left(H(0.99) + H(0.99)\right)\right] = 0.02 - 0.02 = 0$$

$$x_2 \rightarrow H(0.5) - \left[0.5\left(H(0.01) + H(0.99)\right)\right] = 1 - 0.02 = 0.98$$

$$x_3 \rightarrow H(0.5) - \left[0.5\left(H(0.5) + H(0.5)\right)\right] = 1 - 1 = 0$$

Heuristic in AL:

$$\operatorname{argmax}_{\mathbf{x}} \underbrace{H\left(\mathbb{E}_{\boldsymbol{\theta}}\left[f_\theta(\mathbf{x})\right]\right)}_{} - \underbrace{\mathbb{E}_\theta\left[H\left(f_\theta(\mathbf{x})\right)\right]}_{}$$

→ Maximize uncertainty of the expected output

→ Minimize uncertainty due to noise

$H$: Entropy    $\boldsymbol{\theta}$: Ground truth classifier    $\mathbf{x}$: Word embedding

Georgia
Tech

CREATING THE NEXT