

Efficient PAC Learning from the Crowd with Pairwise Comparisons



Shiwei Zeng

Stevens Institute of Technology



Jie Shen

Stevens Institute of Technology

Crowdsourced PAC Learning

Crowdsourced PAC Learning

Standard PAC learning [Valiant 84]

Crowdsourced PAC Learning

Standard PAC learning [Valiant 84]

- Given samples (x, y)

Crowdsourced PAC Learning

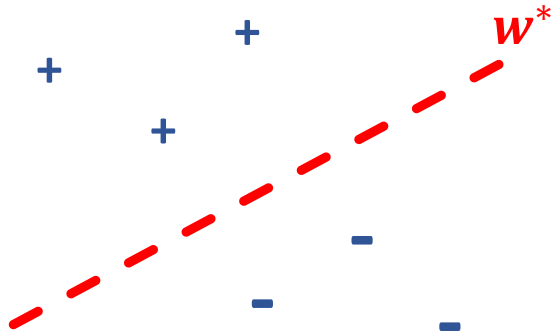
Standard PAC learning [Valiant 84]

- Given samples (x, y)
- Assume $y = h(x)$ for some unknown hypothesis $h \in H$

Crowdsourced PAC Learning

Standard PAC learning [Valiant 84]

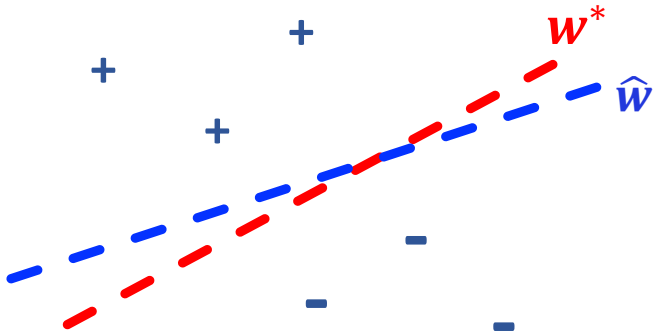
- Given samples (x, y)
- Assume $y = h(x)$ for some unknown hypothesis $h \in H$



Crowdsourced PAC Learning

Standard PAC learning [Valiant 84]

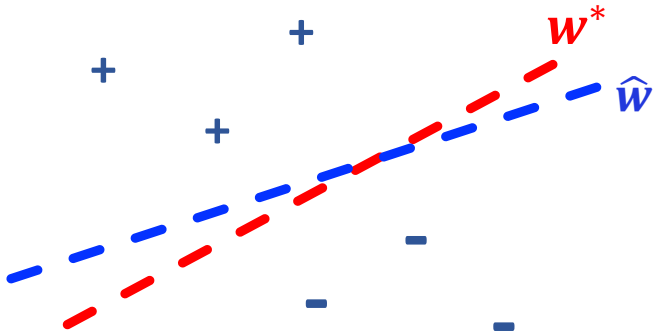
- Given samples (x, y)
- Assume $y = h(x)$ for some unknown hypothesis $h \in H$



Crowdsourced PAC Learning

Standard PAC learning [Valiant 84]

- Given samples (x, y)
- Assume $y = h(x)$ for some unknown hypothesis $h \in H$

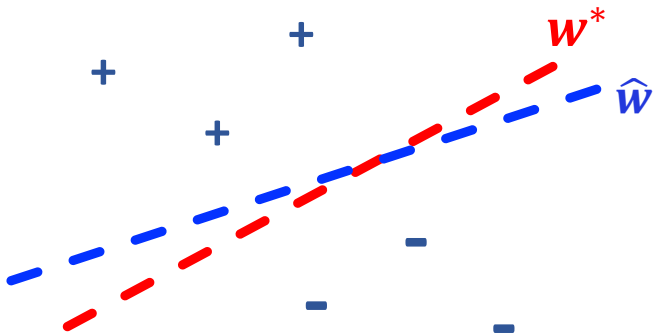


- Goal: \hat{w} is **P**robably **A**pproximately **C**orrect.

Crowdsourced PAC Learning

Standard PAC learning [Valiant 84]

- Given samples (x, y)
- Assume $y = h(x)$ for some unknown hypothesis $h \in H$



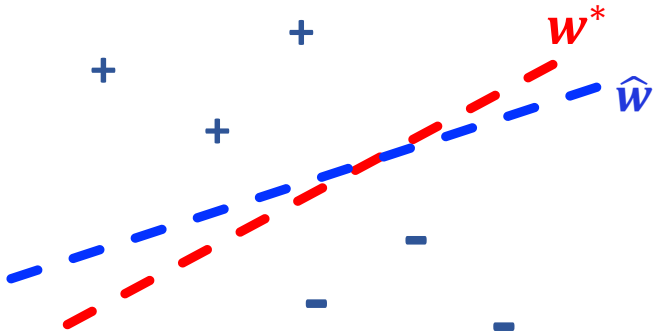
- Goal: \hat{w} is **P**robably **A**pproximately **C**orrect.

Crowdsourced PAC learning [ABHM17]

Crowdsourced PAC Learning

Standard PAC learning [Valiant 84]

- Given samples (x, y)
- Assume $y = h(x)$ for some unknown hypothesis $h \in H$



- Goal: \hat{w} is **P**robably **A**pproximately **C**orrect.

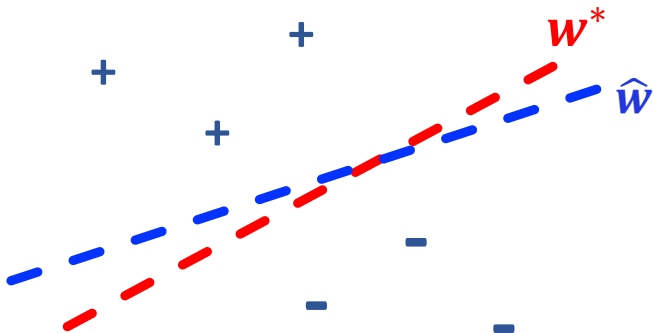
Crowdsourced PAC learning [ABHM17]

- y is not given

Crowdsourced PAC Learning

Standard PAC learning [Valiant 84]

- Given samples (x, y)
- Assume $y = h(x)$ for some unknown hypothesis $h \in H$



- Goal: \hat{w} is **P**robably **A**pproximately **C**orrect.

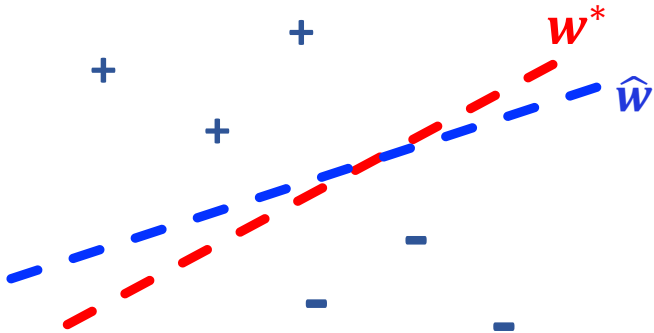
Crowdsourced PAC learning [ABHM17]

- y is not given
- But can collect $\{y_1, \dots, y_k\}$ from crowd

Crowdsourced PAC Learning

Standard PAC learning [Valiant 84]

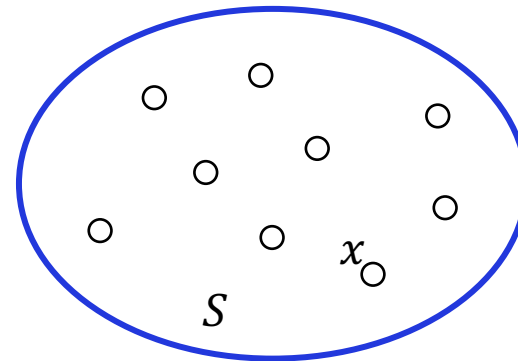
- Given samples (x, y)
- Assume $y = h(x)$ for some unknown hypothesis $h \in H$



- Goal: \hat{w} is **P**robably **A**pproximately **C**orrect.

Crowdsourced PAC learning [ABHM17]

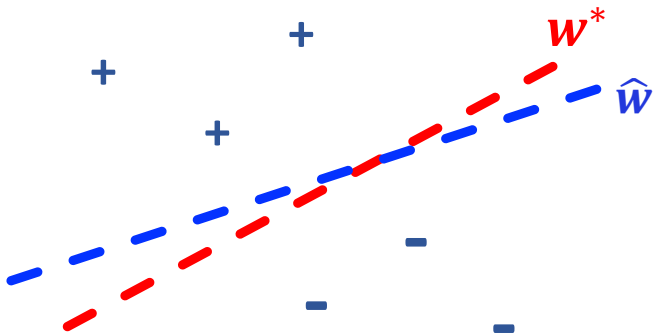
- y is not given
- But can collect $\{y_1, \dots, y_k\}$ from crowd



Crowdsourced PAC Learning

Standard PAC learning [Valiant 84]

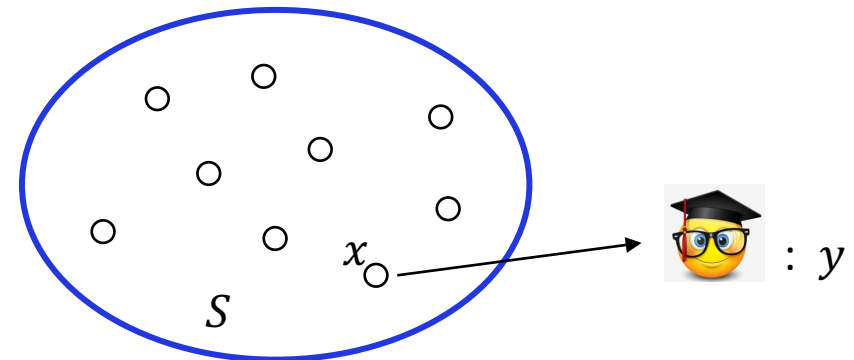
- Given samples (x, y)
- Assume $y = h(x)$ for some unknown hypothesis $h \in H$



- Goal: \hat{w} is **P**robably **A**pproximately **C**orrect.

Crowdsourced PAC learning [ABHM17]

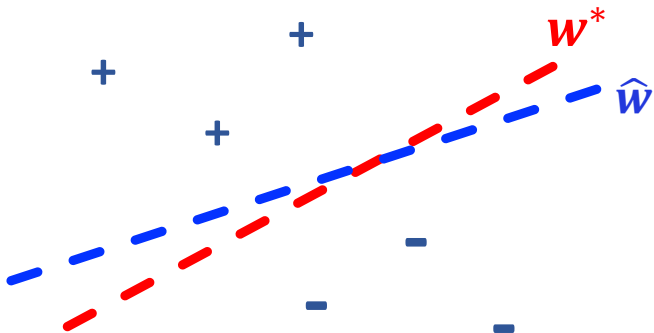
- y is not given
- But can collect $\{y_1, \dots, y_k\}$ from crowd



Crowdsourced PAC Learning

Standard PAC learning [Valiant 84]

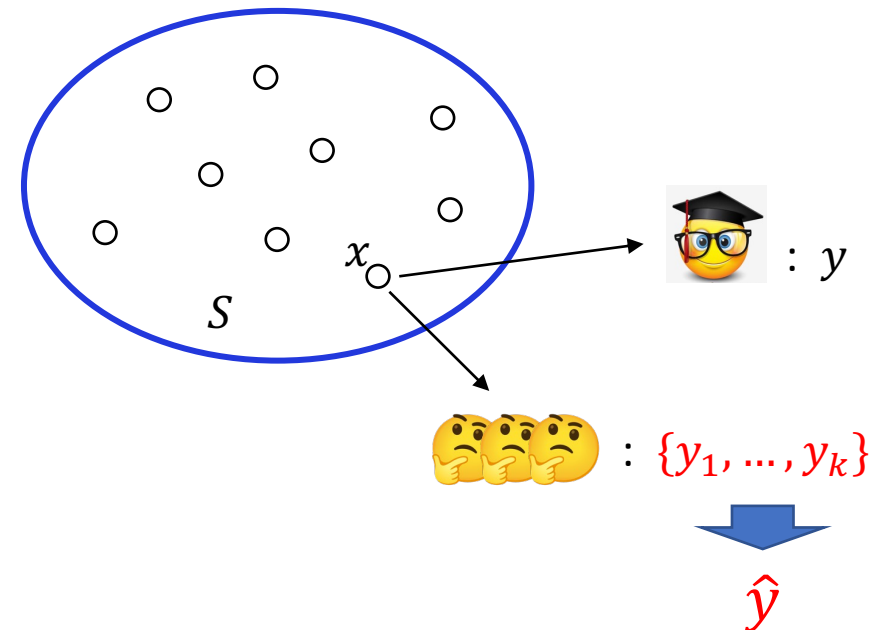
- Given samples (x, y)
- Assume $y = h(x)$ for some unknown hypothesis $h \in H$



- Goal: \hat{w} is **P**robably **A**pproximately **C**orrect.

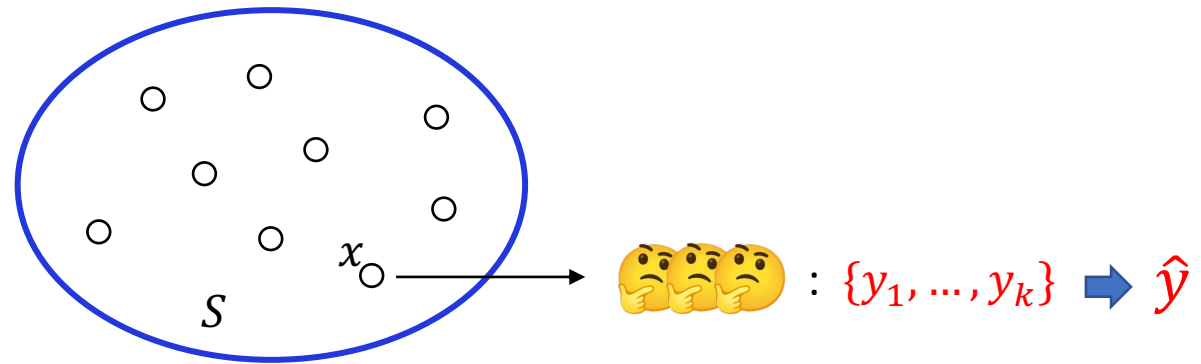
Crowdsourced PAC learning [ABHM17]

- y is not given
- But can collect $\{y_1, \dots, y_k\}$ from crowd

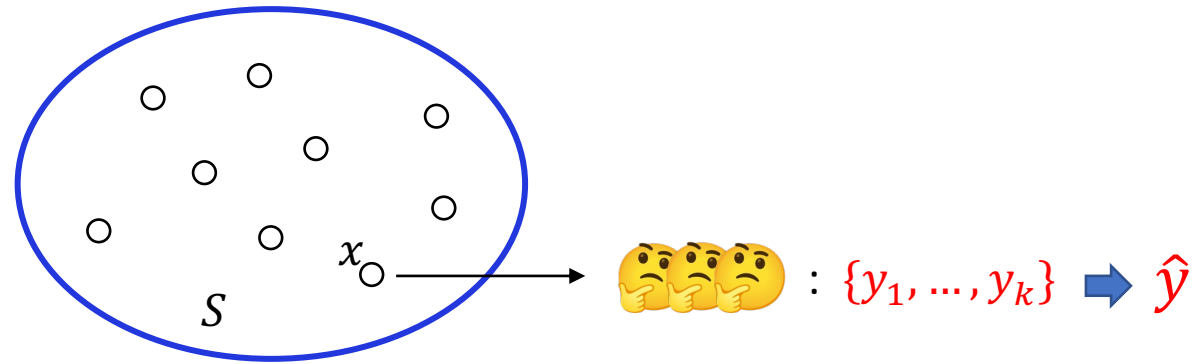


Motivation

Motivation

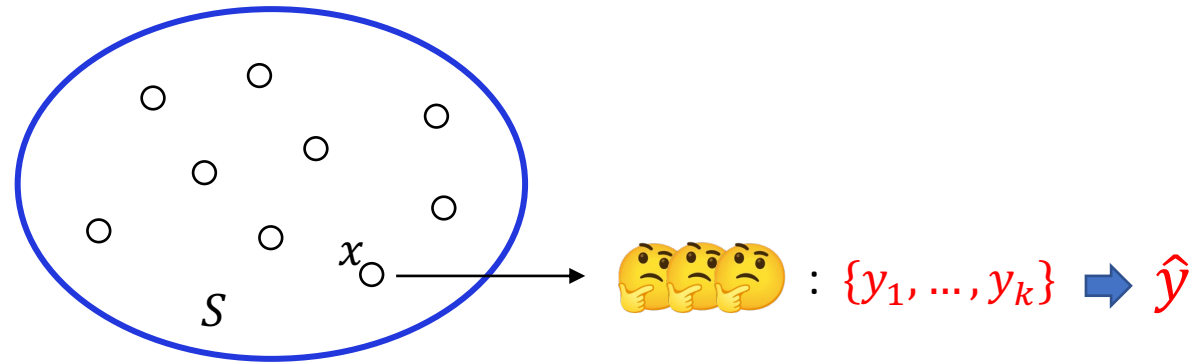


Motivation



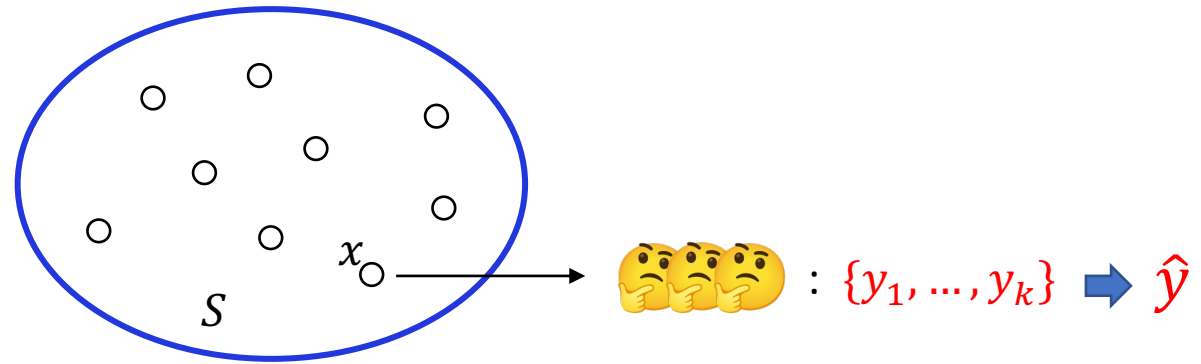
- Label-and-train: $k = \log m$ for each x .

Motivation



- Label-and-train: $k = \log m$ for each x .
- Can we achieve $k = \mathbf{O(1)}$?

Motivation



- Label-and-train: $k = \log m$ for each x .
- Can we achieve $k = \mathbf{O(1)}$? \rightarrow **Query-efficient.**

Motivation

Motivation

- *Recommendation system:*

Motivation

- *Recommendation system:*



Motivation

- *Recommendation system:*



Prefer which one?

Motivation

- *Recommendation system:*



Prefer which one?

- *Covid-19 diagnosis:*

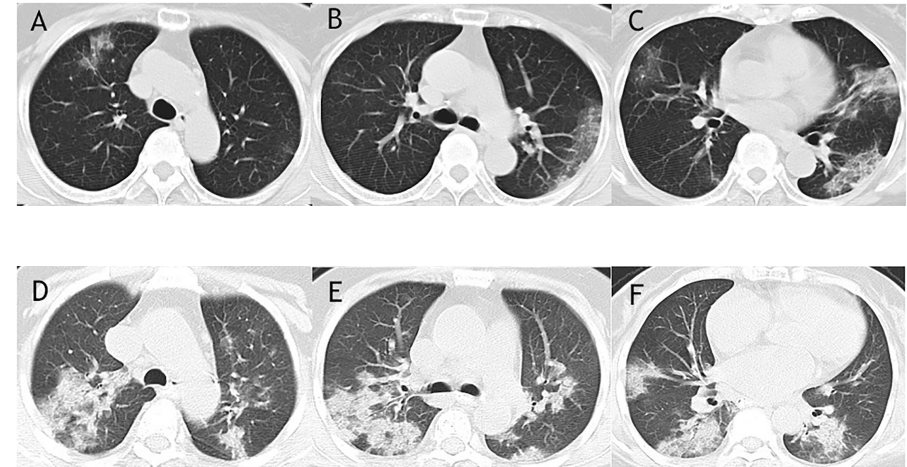
Motivation

- *Recommendation system:*



Prefer which one?

- *Covid-19 diagnosis:*



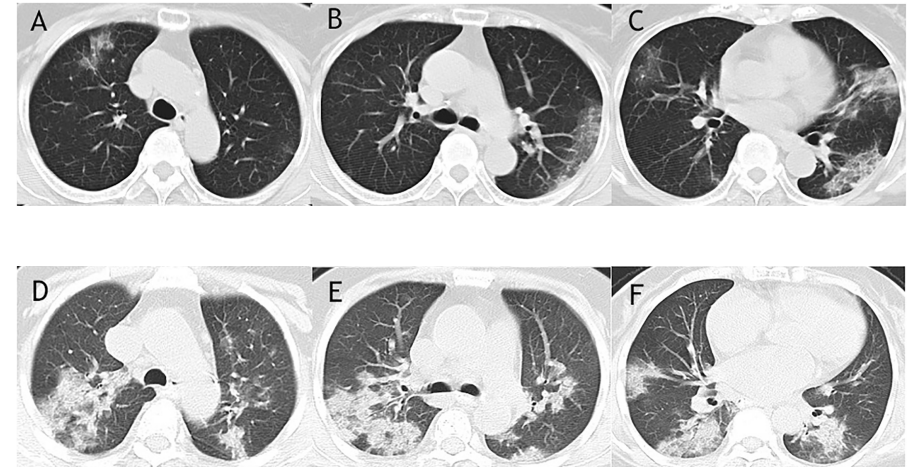
Motivation

- *Recommendation system:*



Prefer which one?

- *Covid-19 diagnosis:*



Which one is healthier?

Motivation

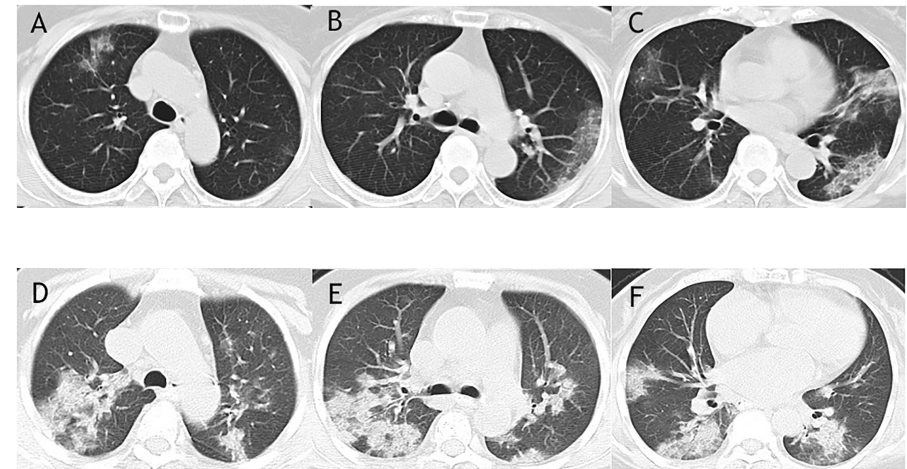
- *Recommendation system:*



Prefer which one?

- **Pairwise comparisons:** $f(x) >^? f(x')$

- *Covid-19 diagnosis:*



Which one is healthier?

Our results

Our results

- PAC guarantees: w.p. $1-\delta$, error $\leq \epsilon$.

Our results

- PAC guarantees: w.p. $1-\delta$, error $\leq \epsilon$.
- No distributional assumptions.

Our results

- PAC guarantees: w.p. $1-\delta$, error $\leq \epsilon$.
- No distributional assumptions.
- **Query** and **label-efficient**.

Our results

- PAC guarantees: w.p. $1-\delta$, error $\leq \epsilon$.
- No distributional assumptions.
- **Query** and **label-efficient**.
- #Labels = $\tilde{O}(\log \frac{d+1/\delta}{\epsilon})$, #pairwise comparisons = $\tilde{O}(\frac{d+(1/\delta)^{1/1000}}{\epsilon})$.

Our results

- PAC guarantees: w.p. $1-\delta$, error $\leq \epsilon$.
- No distributional assumptions.
- **Query** and **label-efficient**.
- #Labels = $\tilde{O}(\log \frac{d+1/\delta}{\epsilon})$, #pairwise comparisons = $\tilde{O}(\frac{d+(1/\delta)^{1/1000}}{\epsilon})$.

#Label per instance:

$$\Lambda_L = \frac{\epsilon}{d} \tilde{O}(\log \frac{d}{\epsilon})$$

Our results

- PAC guarantees: w.p. $1-\delta$, error $\leq \epsilon$.
- No distributional assumptions.
- **Query** and **label-efficient**.
- #Labels = $\tilde{O}(\log \frac{d+1/\delta}{\epsilon})$, #pairwise comparisons = $\tilde{O}(\frac{d+(1/\delta)^{1/1000}}{\epsilon})$.

#Label per instance:

$$\Lambda_L = \frac{\epsilon}{d} \tilde{O}(\log \frac{d}{\epsilon})$$

#Comparison per instance:

$$\Lambda_C = \tilde{O}(\sqrt{\epsilon} \log^2 \frac{d}{\epsilon} + 1)$$

Our results

- PAC guarantees: w.p. $1-\delta$, error $\leq \epsilon$.
- No distributional assumptions.
- **Query** and **label-efficient**.
- #Labels = $\tilde{O}(\log \frac{d+1/\delta}{\epsilon})$, #pairwise comparisons = $\tilde{O}(\frac{d+(1/\delta)^{1/1000}}{\epsilon})$.

#Label per instance:

$$\Lambda_L = \frac{\epsilon}{d} \tilde{O}(\log \frac{d}{\epsilon})$$

#Comparison per instance:

$$\Lambda_C = \tilde{O}(\sqrt{\epsilon} \log^2 \frac{d}{\epsilon} + 1)$$

- When $\epsilon \rightarrow 0$, $\Lambda_L = o(1)$, $\Lambda_C = O(1)$.

Main techniques

Main techniques

- Comparison-based *sorting* algorithm.

Main techniques

- Comparison-based *sorting* algorithm.
- Query and label-efficient *filtering* algorithm.

Main techniques

- Comparison-based *sorting* algorithm.
- Query and label-efficient *filtering* algorithm.

See you at **poster session!**