

Promises and Pitfalls of Threshold-based Auto-labeling

Harit Vishwakarma

hvishwakarma@cs.wisc.edu

Ph.D. Student

Dept. of Computer Sciences
University of Wisconsin-Madison



Appearing in

NEURAL INFORMATION PROCESSING SYSTEMS 2023



Huguang Lin

hglin@seas.upenn.edu



Frederic Sala

fredsala@cs.wisc.edu



Ramya Korlakai Vinayak

ramya@ece.wisc.edu

Roadmap

What & Why auto-labeling?

Data labeling problem

Wide adoption of
auto-labeling

How does it work?

Workflow of TBAL

Finding the
auto-labeling region

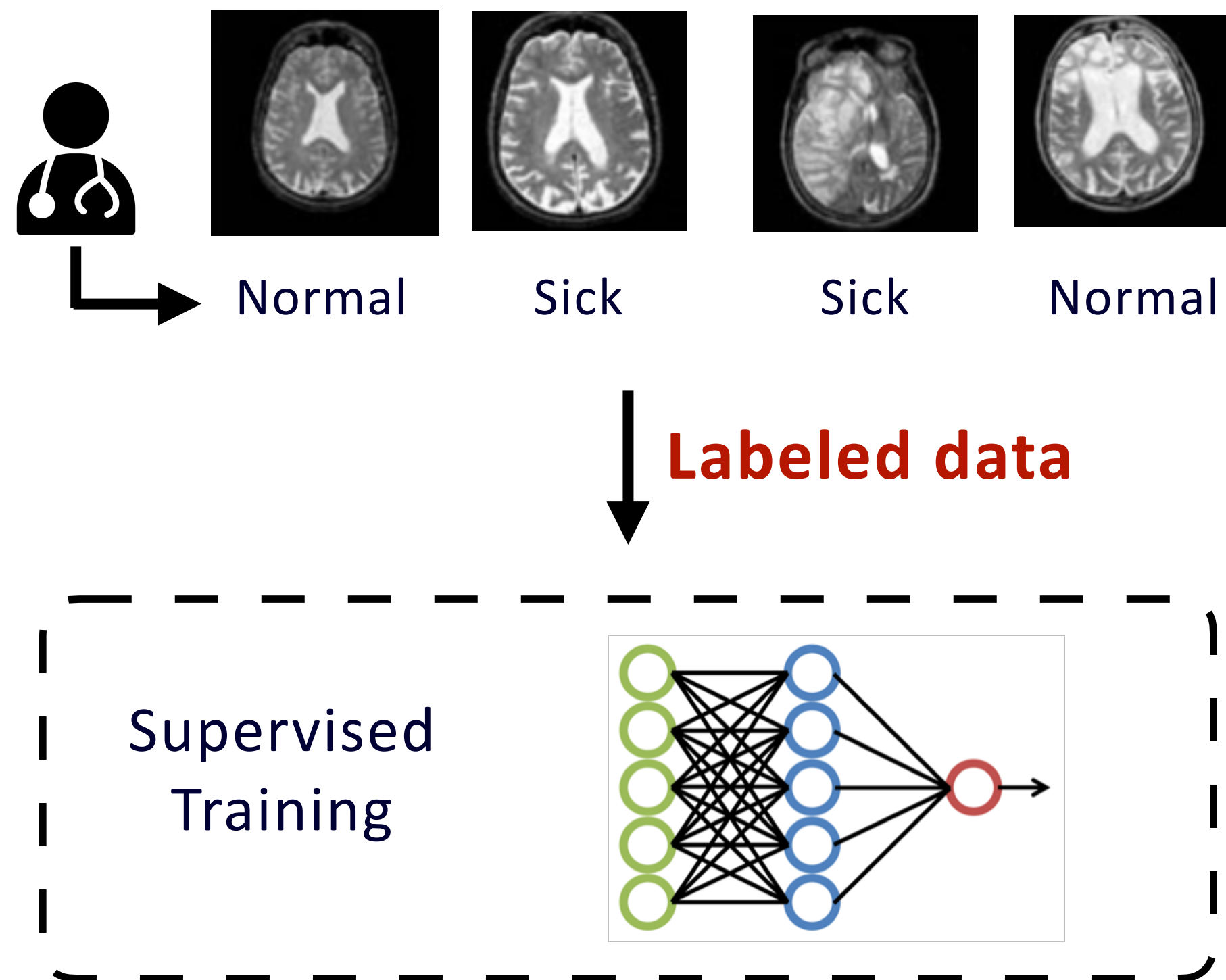
Analysis & Results

Conditions when TBAL works.

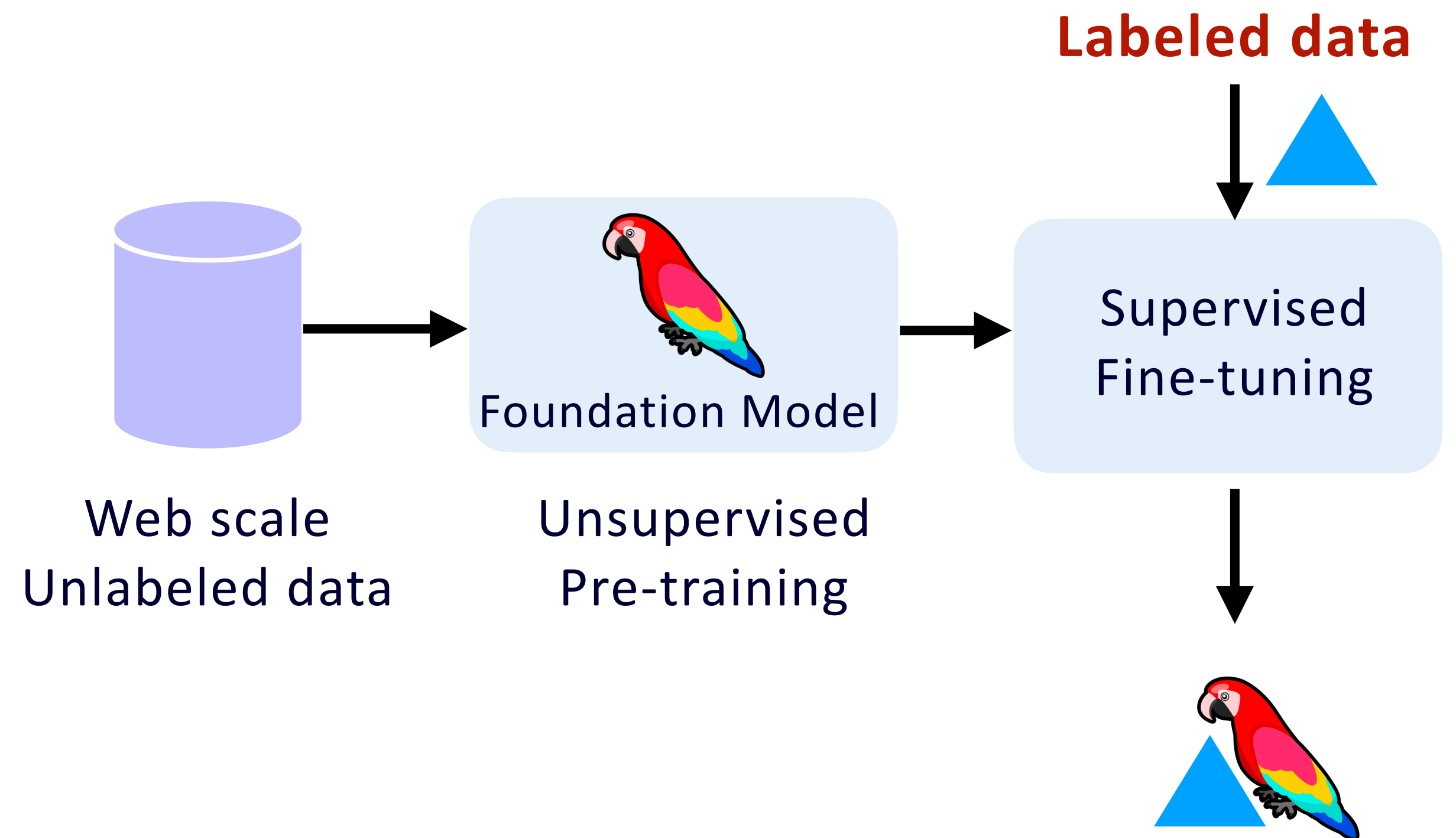
Comparison with
Active Learning, Selective
Classification

We need labeled data and often a lot of it!

Diagnosing a novel disease using brain scans



Fine-tuning Foundation models or Aligning LLMs



Data Labeling costs a lot of **time and money**

Crowdsourcing is widely used to get labels

Wisdom of Crowd



amazon

mechanical turk

and many others...

Takes a lot of time and money to get labels.

IMAGENET

Deng et. Al. 2009

Took multiple years and a lot of human effort

Geological formation, formation (geology) the geological features of the earth

14M Images, 20K Classes.

1808 pictures 86.24% Popularity Percentile Wordnet IDs

Numbers in brackets: (the number of synsets in the subtree).

ImageNet 2011 Fall Release (32326)

- plant, flora, plant life (4486)
- geological formation, formation (1:
 - aquifer (0)
 - beach (1)
 - cave (3)
 - cliff, drop, drop-off (2)
 - delta (0)
 - diapir (0)
 - follum (0)
 - foreshore (0)
 - ice mass (10)
 - lakefront (0)
 - massif (0)
 - monocline (0)
 - mouth (0)
 - natural depression, depression (
 - natural elevation, elevation (41
 - oceanfront (0)
 - range, mountain range, range of
 - relict (0)
 - ridge, ridgeline (2)
 - ridge (0)
 - shore (7)
 - slope, incline, side (17)
 - spring, fountain, outflow, outpo
 - talus, scree (0)
 - vein, mineral vein (1)
 - volcanic crater, crater (2)
 - wall (0)

Treemap Visualization Images of the Synset Downloads

ImageNet 2011 Fall Release Geological formation, formation

Natural	Slope	Shore
Ice	Water	Vein
Delta	Foreshore	
Massif	Talus	Volcanic
Beach		
Mouth	Lakefront	Range
Diapir	Cliff	
Wall	Oceanfront	Aquifer
Cave	Spring	
Monocline		Ridge

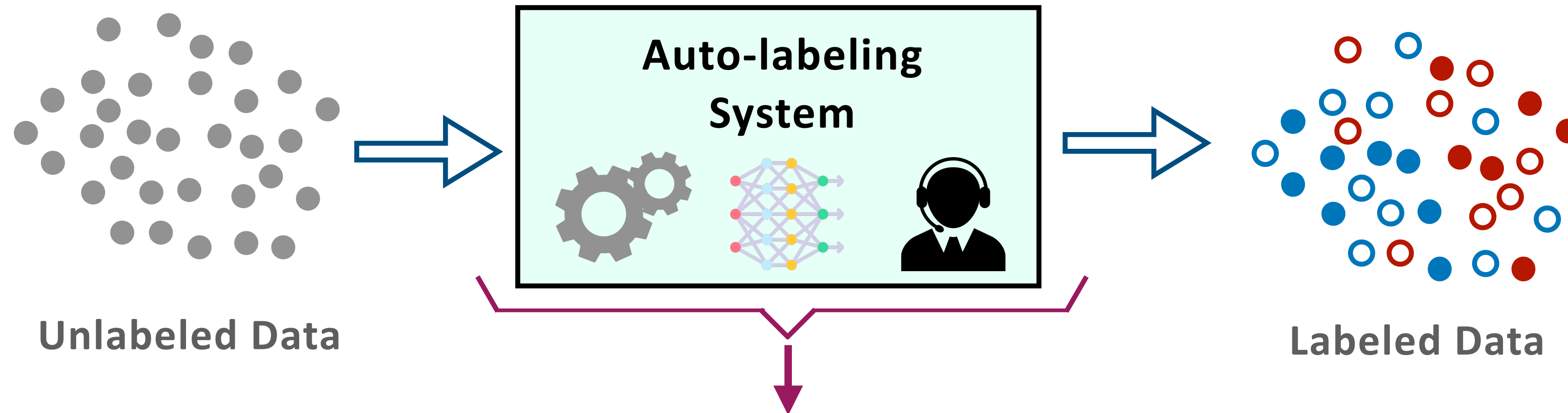
A screenshot of the ImageNet database online

IMAGENET

How do we get **accurately labeled** data, while spending **less time and money**?

Automatically label datasets with minimal human feedback

● ● Human-labeled
○ ○ Auto-labeled



Get labels for “minimal” points from human



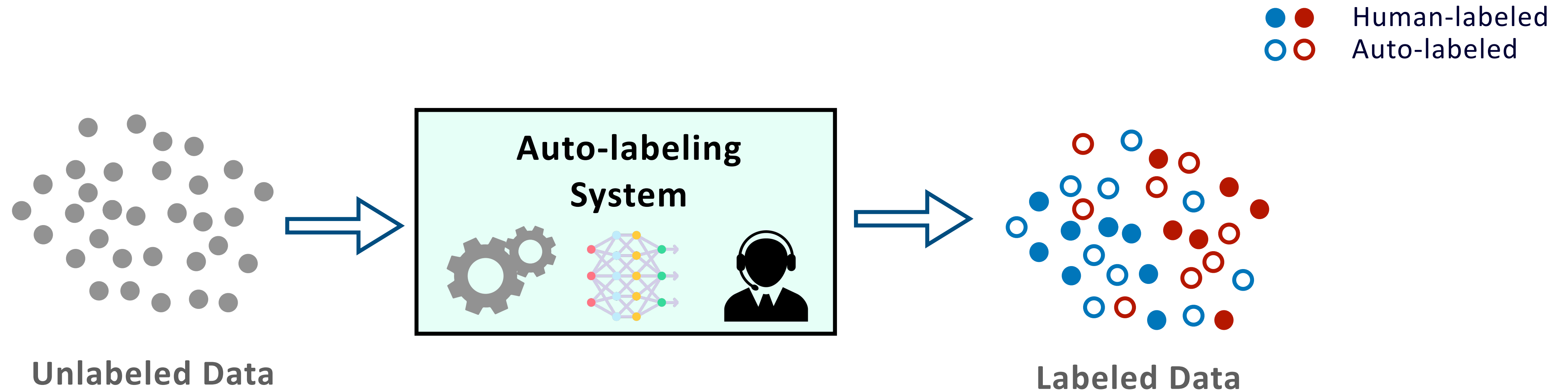
Train a model on these labeled points and



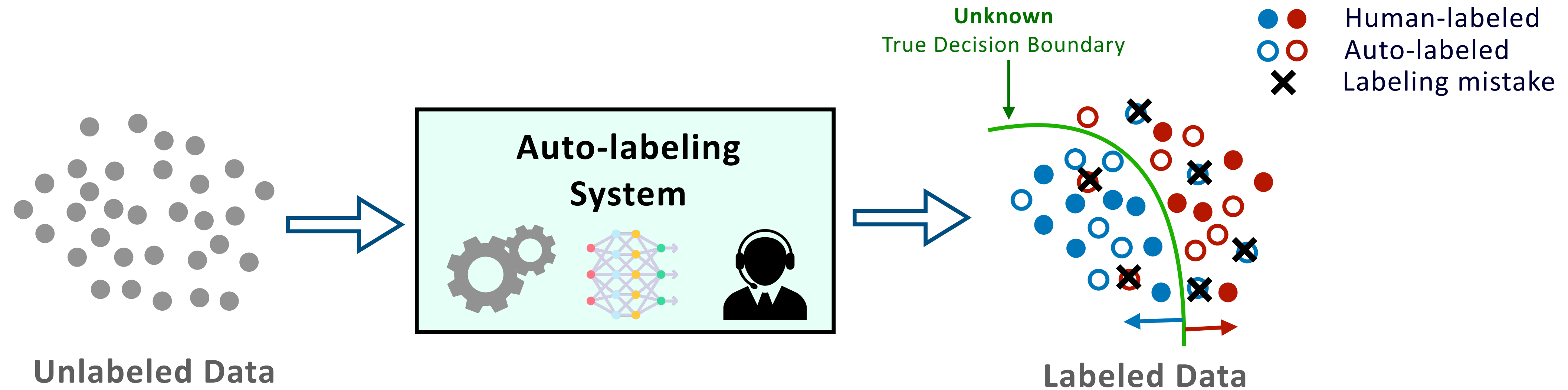
Auto-label using the model



Auto-Labeling Errors and Their Impact



Auto-Labeling Errors and Their Impact



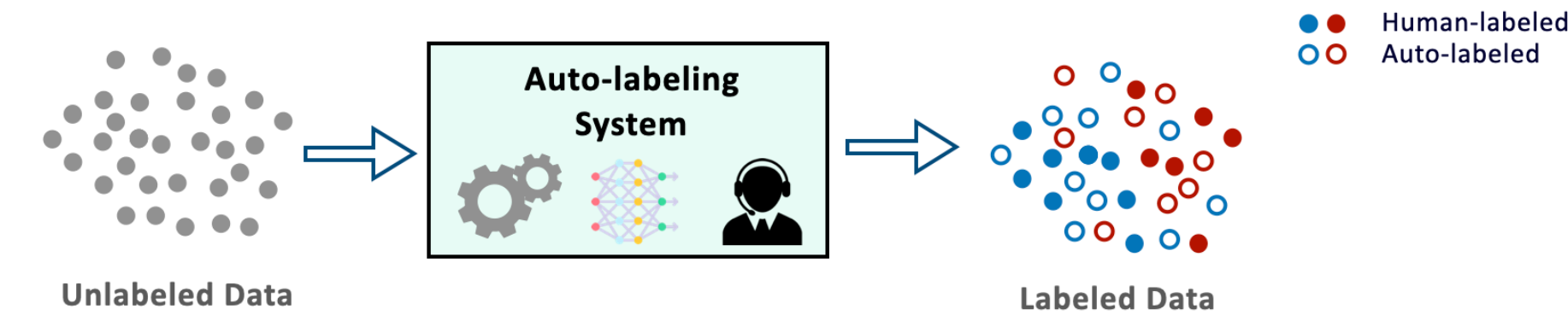
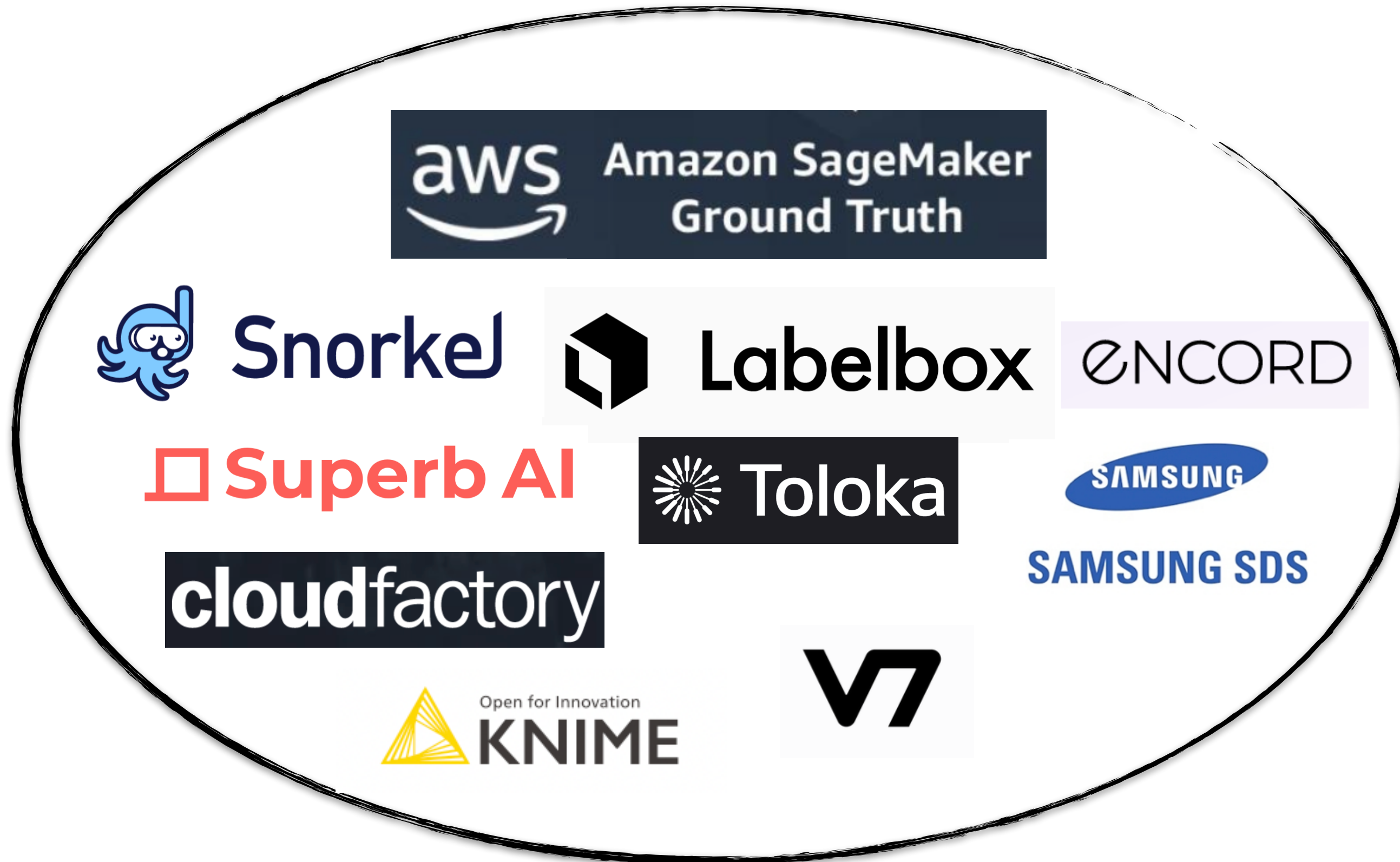
1. The output dataset may have labeling errors

2. The impact of errors in datasets is more severe

- a) Multiple downstream applications
- b) Longer shelf-life than models.

Auto-labeling systems are widely used

Auto-labeling Platforms



Auto-labeling is heavily used commercially.

Even in **high risk** applications

health care, telecom, recruiting...

So we need to understand them.

Despite wide adoption, our **understanding of auto-labeling systems is limited!**

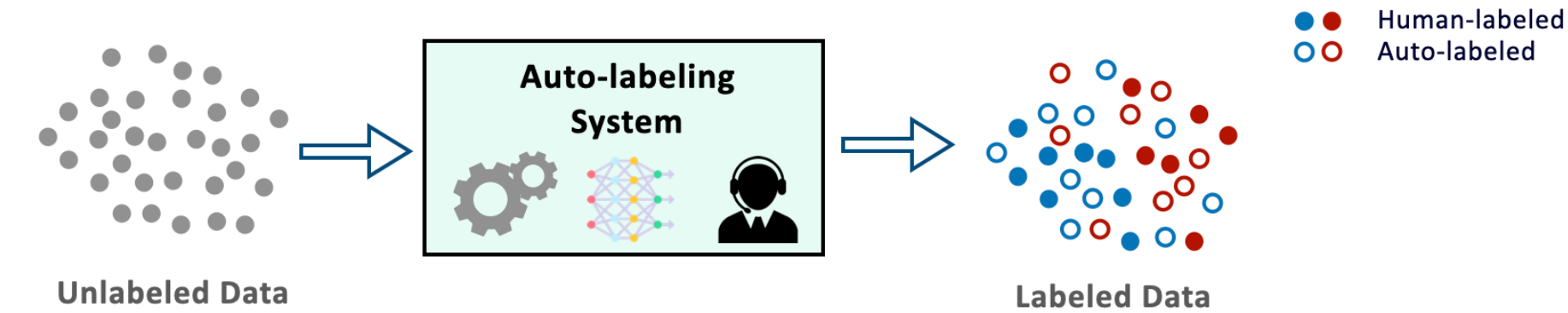
Despite wide adoption, our **understanding of auto-labeling systems is limited!**

To address this gap we **develop a theoretical understanding** of auto-labeling systems.

Auto-labeling systems are widely used

Auto-labeling Platforms

Study a w/f inspired from it.



Auto-labeling is heavily used commercially.

Even in **high risk** applications

health care, telecom, recruiting...

So we need to understand them.

Roadmap

What & Why auto-labeling?

Data labeling problem

Adoption of auto-labeling

How does it work?

Workflow of TBAL

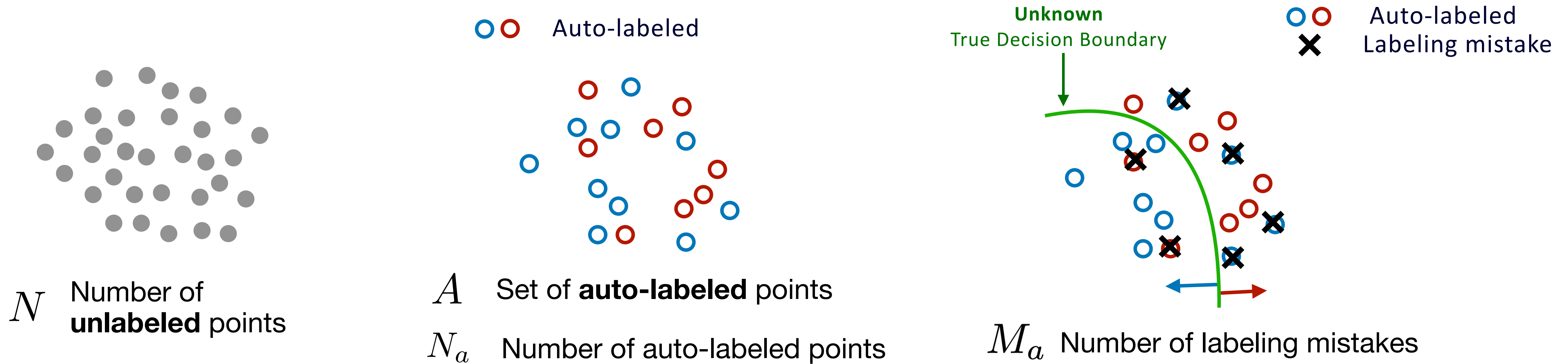
Finding the
auto-labeling region

Analysis & Results

Conditions when TBAL works.

Comparison with
Active Learning, Selective
Classification

Quality and Quantity of Auto-labeled Data



Quantity

Auto-labeling Coverage

$$\hat{\mathcal{P}} = \frac{N_a}{N}$$

Good Stuff
maximize this



Quality

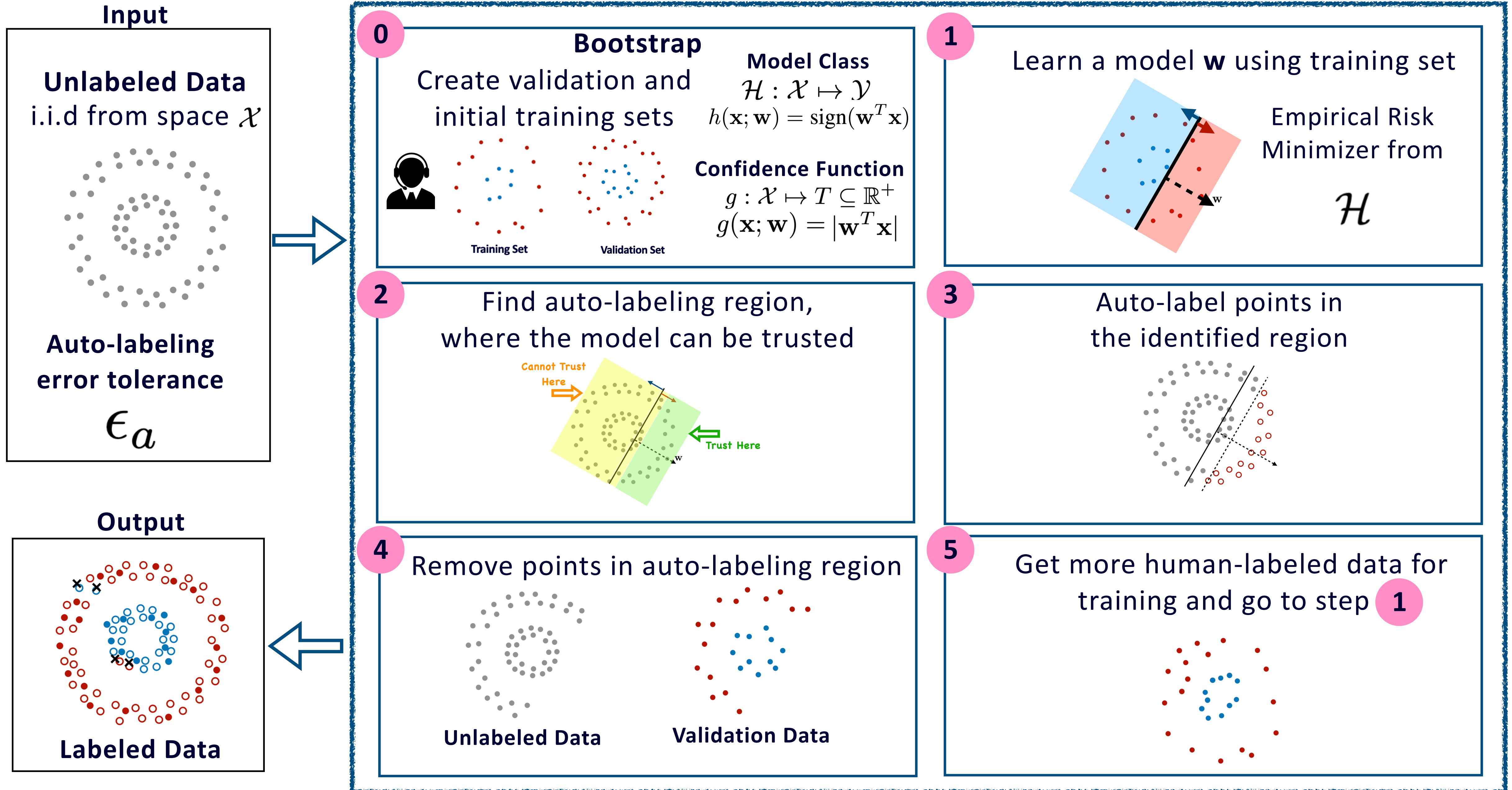
Auto-labeling Error

$$\hat{\mathcal{E}} = \frac{M_a}{N_a}$$

Bad Stuff
minimize this



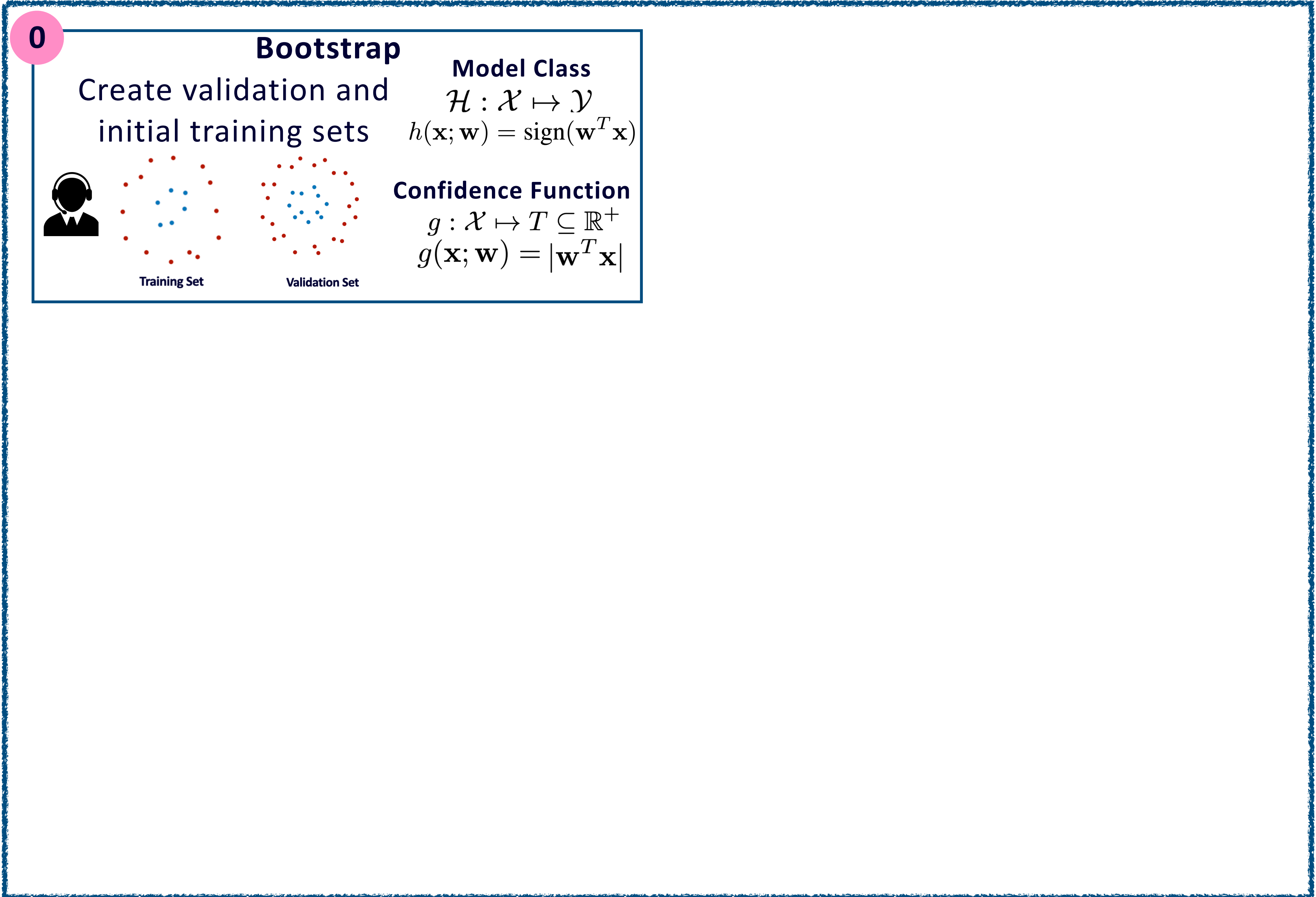
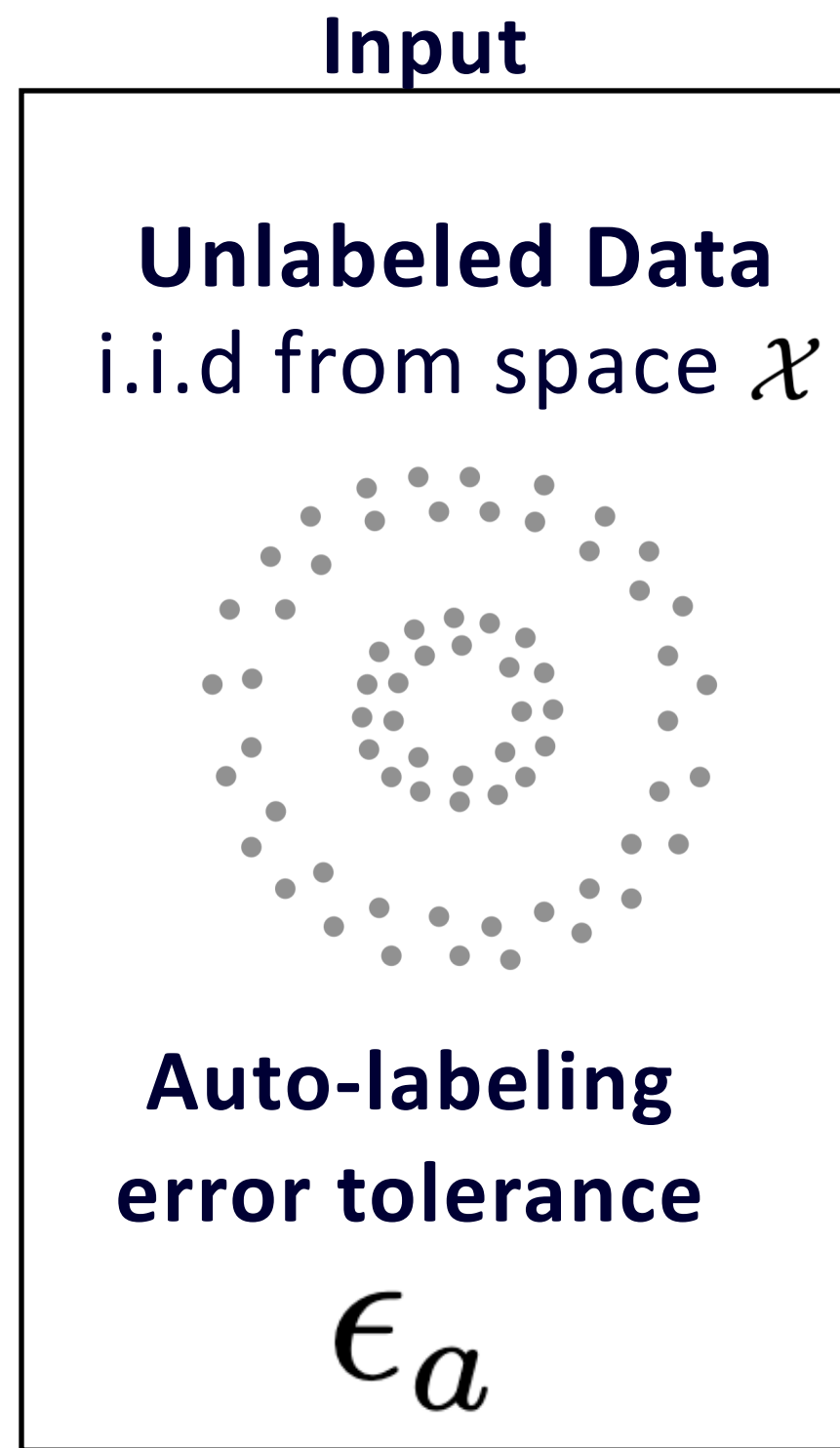
Threshold-based Auto-labeling Workflow(TBAL)



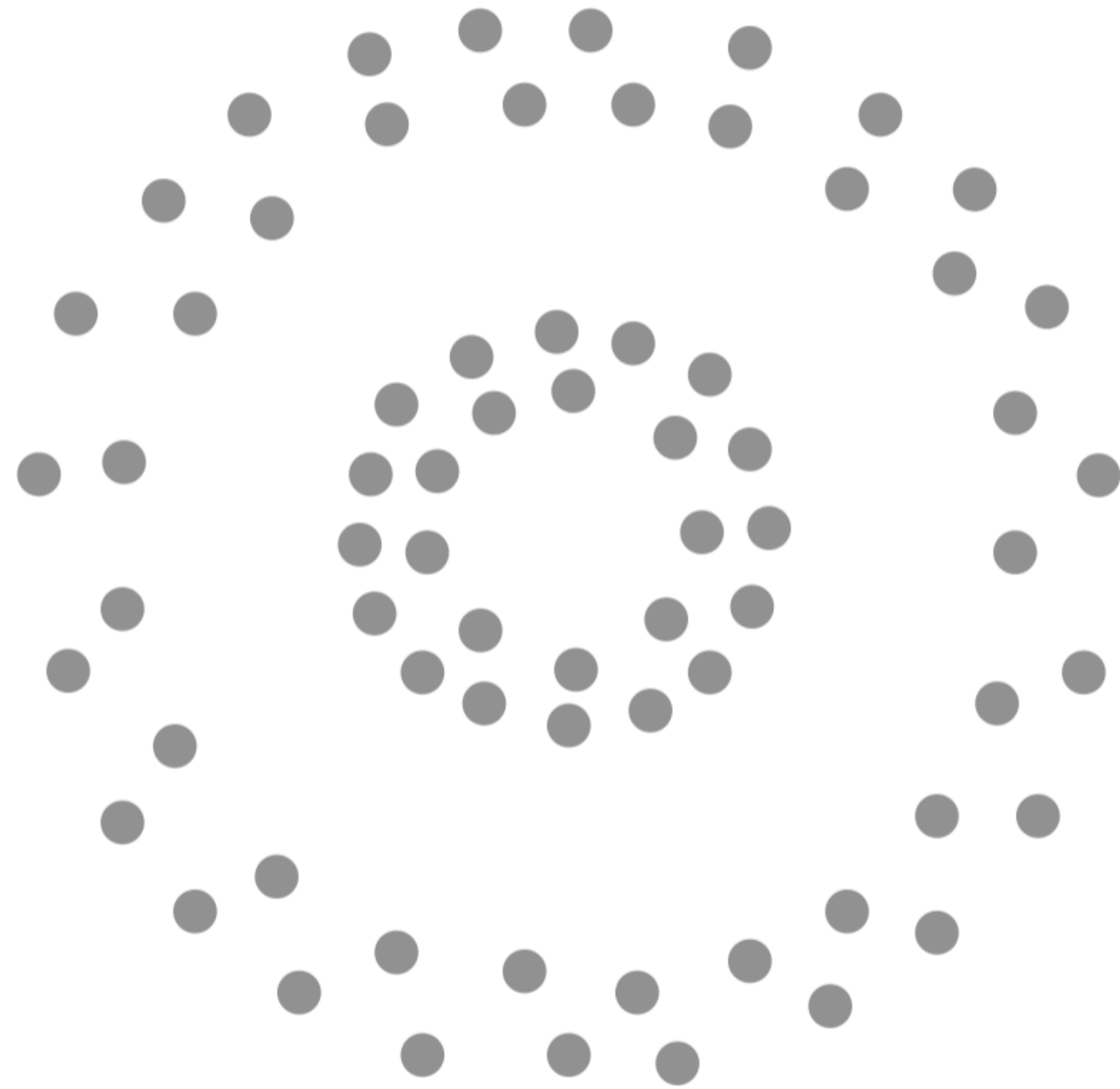
~~Pretend~~ we are LLMs and

Let's think step by step with an example

Threshold-based Auto-labeling Workflow(TBAL)

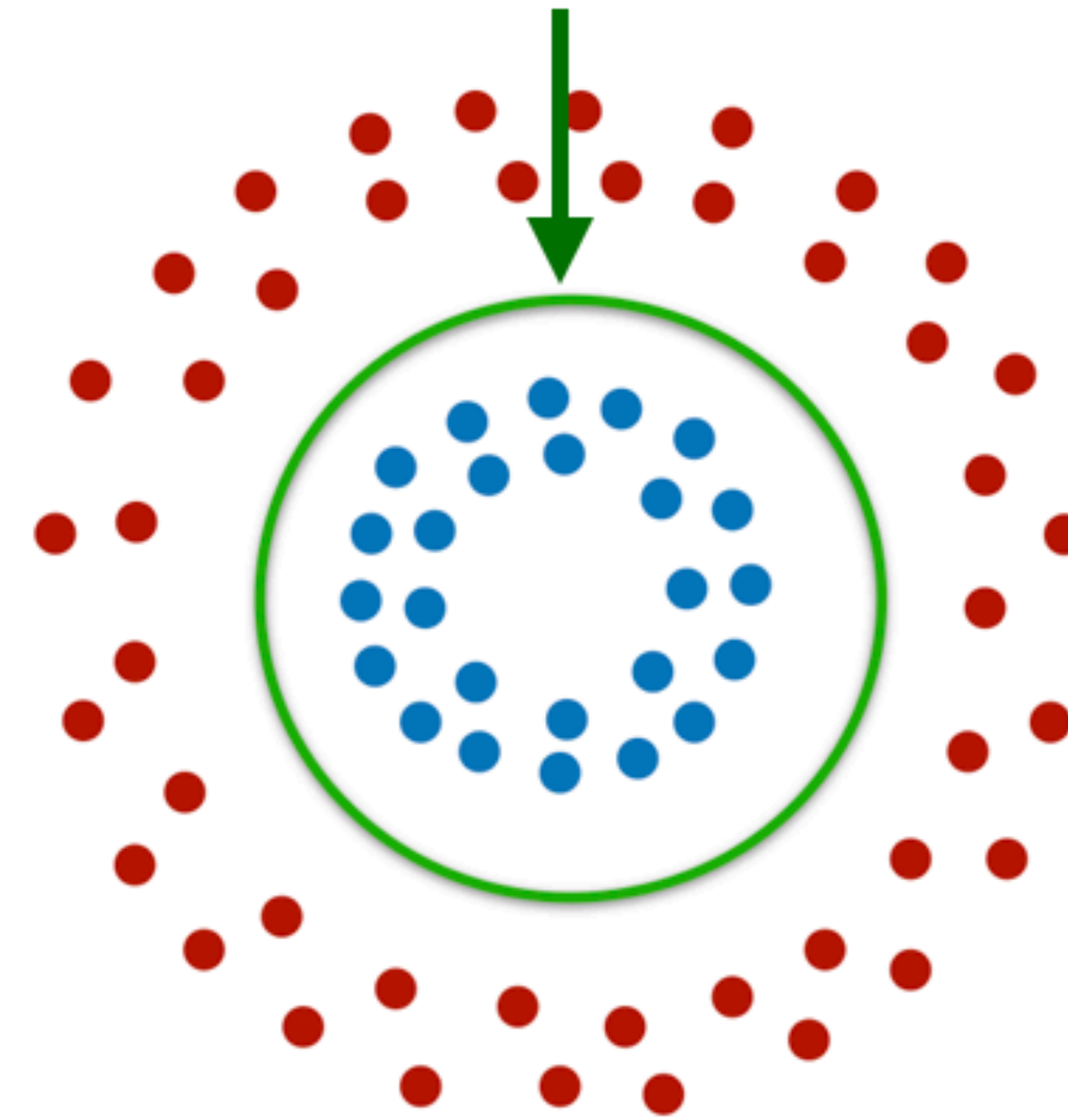


Input



Unlabeled Data
i.i.d from space \mathcal{X}

Unknown
True Decision Boundary f^*

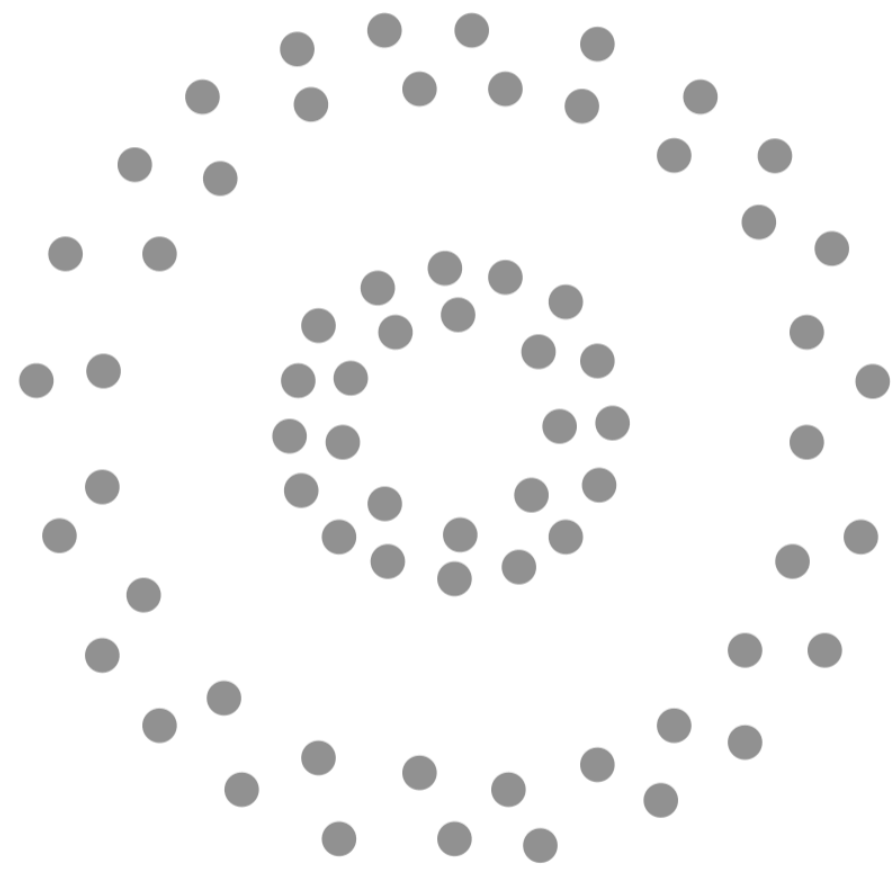


**Known to
the oracle**



Learning f^* is NOT the goal.

Input

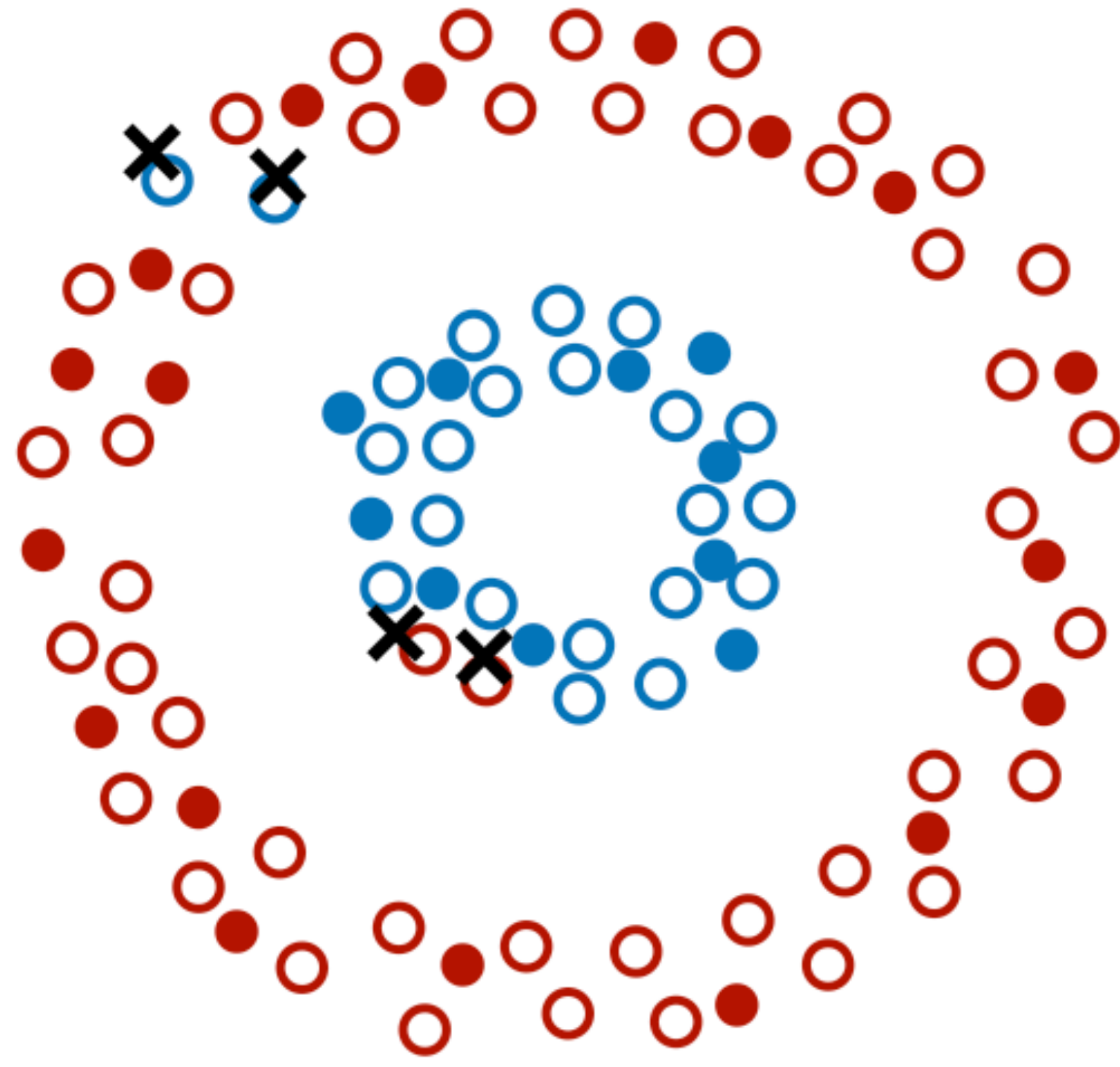


Unlabeled Data

Auto-labeling error tolerance

$$\epsilon_a$$

Expected Output



Labeled Data

- ● Human-labeled
- ○ Auto-labeled
- ✕ Labeling mistake

Auto-labeling Error

$$\hat{\epsilon} = \frac{M_a}{N_a} = \frac{\# \times}{\# \circ + \# \bullet} \leq \epsilon_a$$

Coverage

$$\hat{p} = \frac{N_a}{N} = \frac{\# \circ + \# \bullet}{\# \bullet}$$

TBAL Workflow : Bootstrap (Step 0)

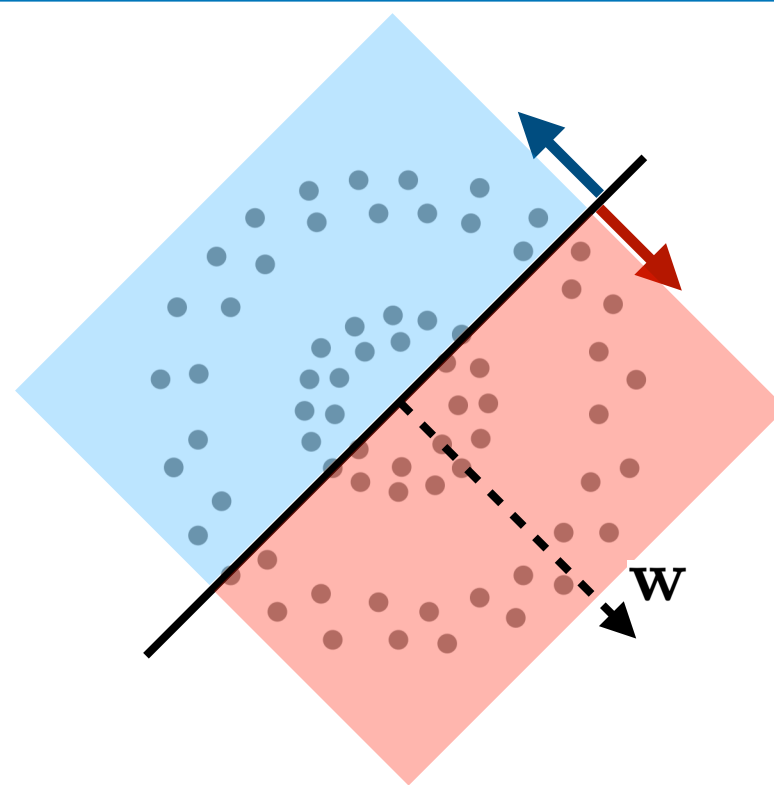
Pick a Model class and Confidence function

Model/Hypothesis Class

$$\mathcal{H} : \mathcal{X} \mapsto \mathcal{Y}$$

$$\mathcal{X} = \{\mathbf{x} \in \mathbb{R}^2 : \|\mathbf{x}\|_2 \leq 1\}$$

$$\mathcal{Y} = \{-1, +1\}$$



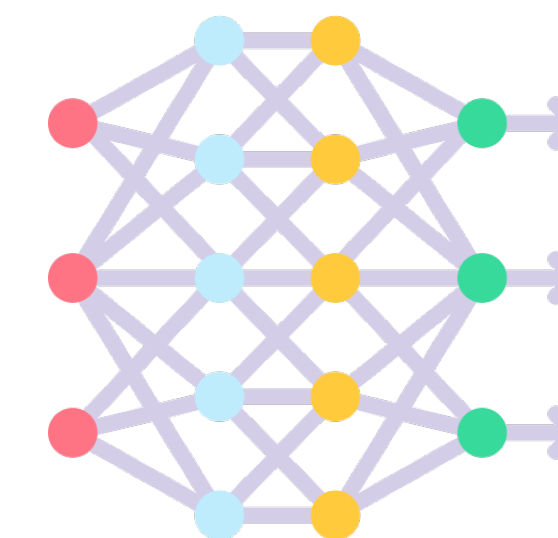
Linear Classifiers

$$\mathcal{W} = \{\mathbf{w} \in \mathbb{R}^2 : \|\mathbf{w}\|_2 \leq 1\}$$

$$h(\mathbf{x}; \mathbf{w}) = \text{sign}(\mathbf{w}^T \mathbf{x})$$

$$f^* \notin \mathcal{H}$$

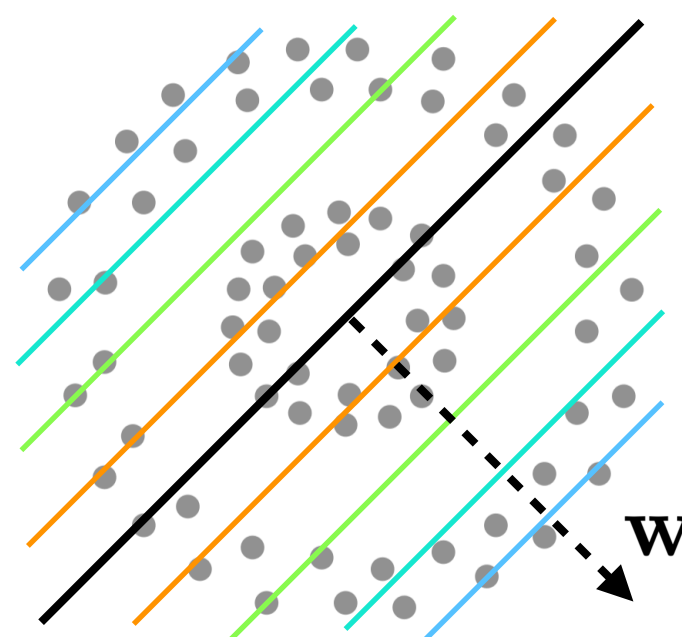
Neural Nets



Confidence/Scoring Function

$$g : \mathcal{X} \mapsto T \subseteq \mathbb{R}^+$$

$$T = [0, 1]$$

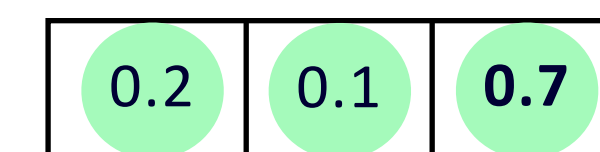


Linear Confidence Function

$$g(\mathbf{x}; \mathbf{w}) = \frac{1}{1 + e^{-|\mathbf{w}^T \mathbf{x}|}}$$

$$\equiv |\mathbf{w}^T \mathbf{x}|$$

Softmax Score



TBAL Workflow : Bootstrap (Step 0)

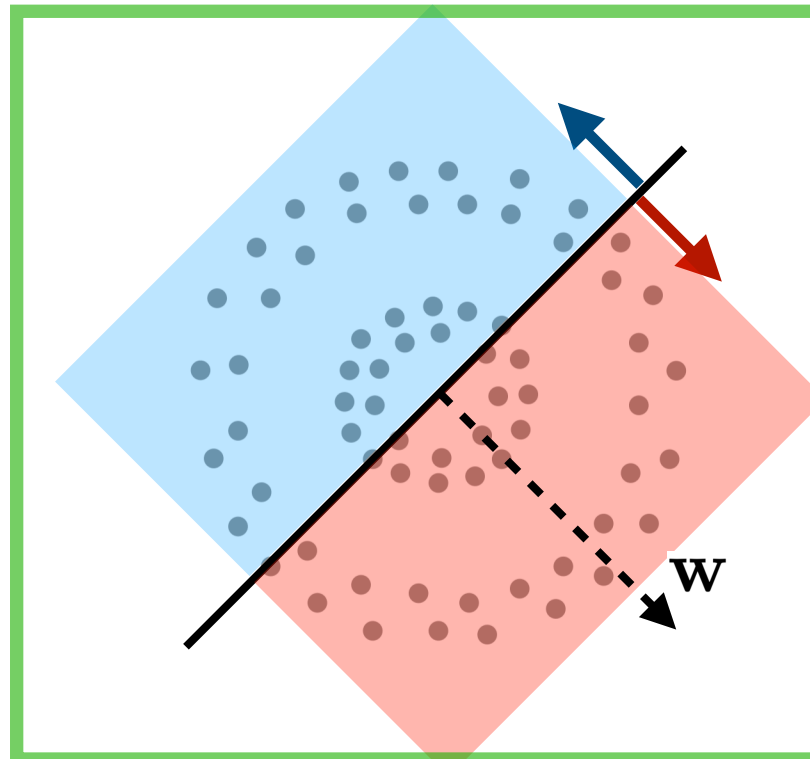
Pick a Model class and Confidence function

Model/Hypothesis Class

$$\mathcal{H} : \mathcal{X} \mapsto \mathcal{Y}$$

$$\mathcal{X} = \{\mathbf{x} \in \mathbb{R}^2 : \|\mathbf{x}\|_2 \leq 1\}$$

$$\mathcal{Y} = \{-1, +1\}$$



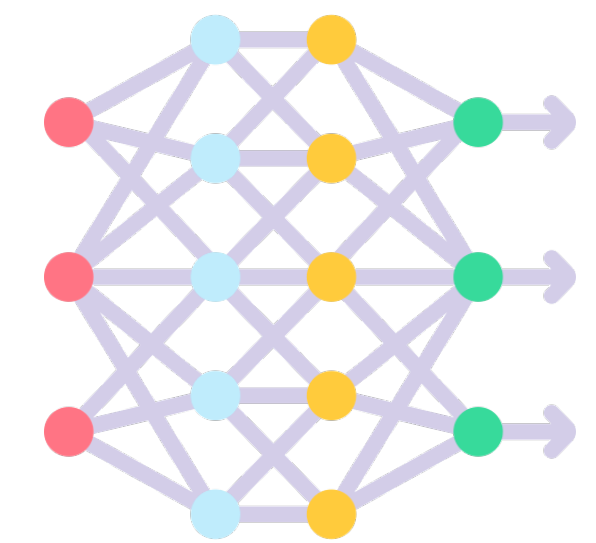
Linear Classifiers

$$\mathcal{W} = \{\mathbf{w} \in \mathbb{R}^2 : \|\mathbf{w}\|_2 \leq 1\}$$

$$h(\mathbf{x}; \mathbf{w}) = \text{sign}(\mathbf{w}^T \mathbf{x})$$

$$f^* \notin \mathcal{H}$$

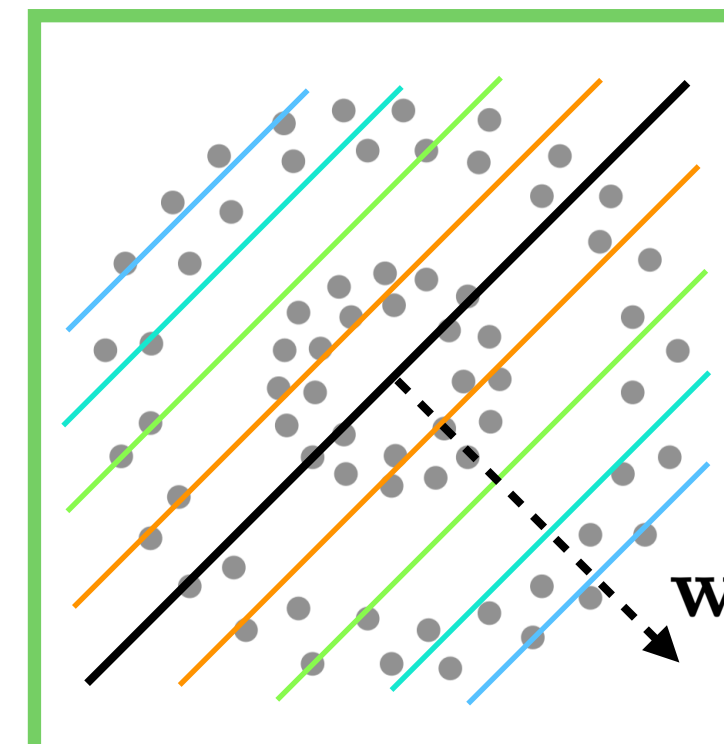
Neural Nets



Confidence/Scoring Function

$$g : \mathcal{X} \mapsto T \subseteq \mathbb{R}^+$$

$$T = [0, 1]$$



Linear Confidence Function

$$g(\mathbf{x}; \mathbf{w}) = \frac{1}{1 + e^{-|\mathbf{w}^T \mathbf{x}|}}$$

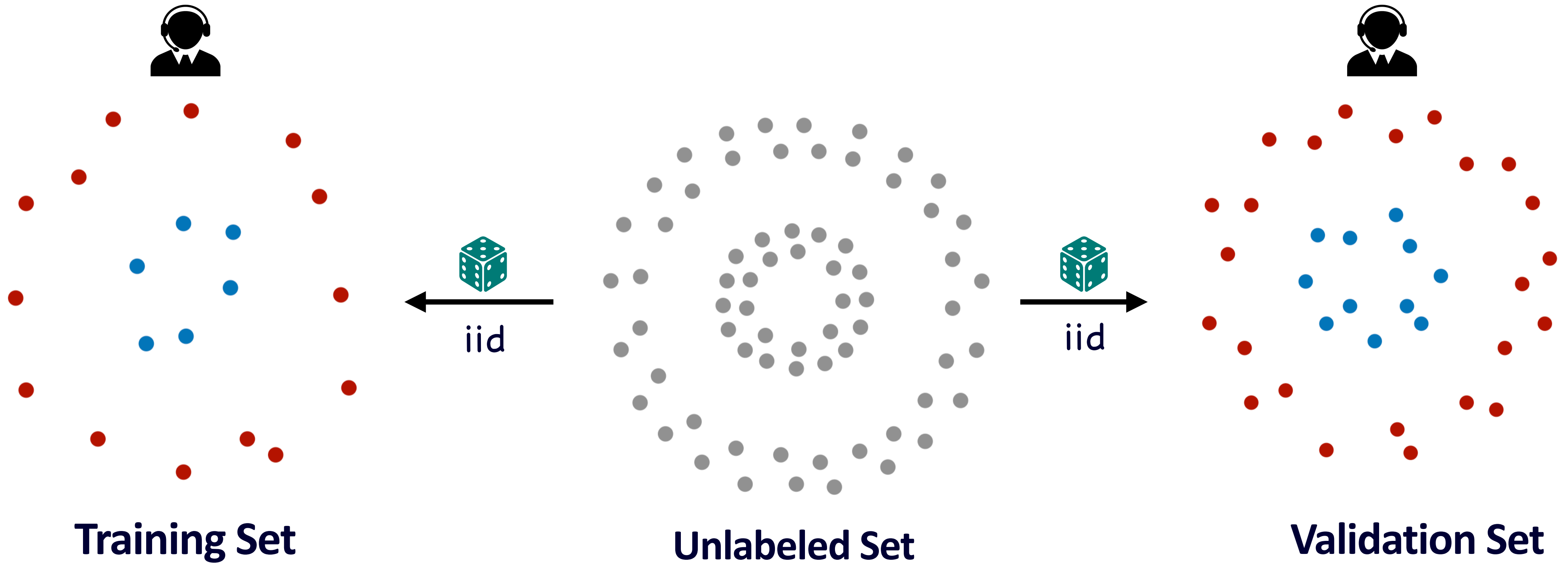
$$\equiv |\mathbf{w}^T \mathbf{x}|$$

Softmax Score



TBAL Workflow : Bootstrap (Step 0)

Get some labeled data for training and validation



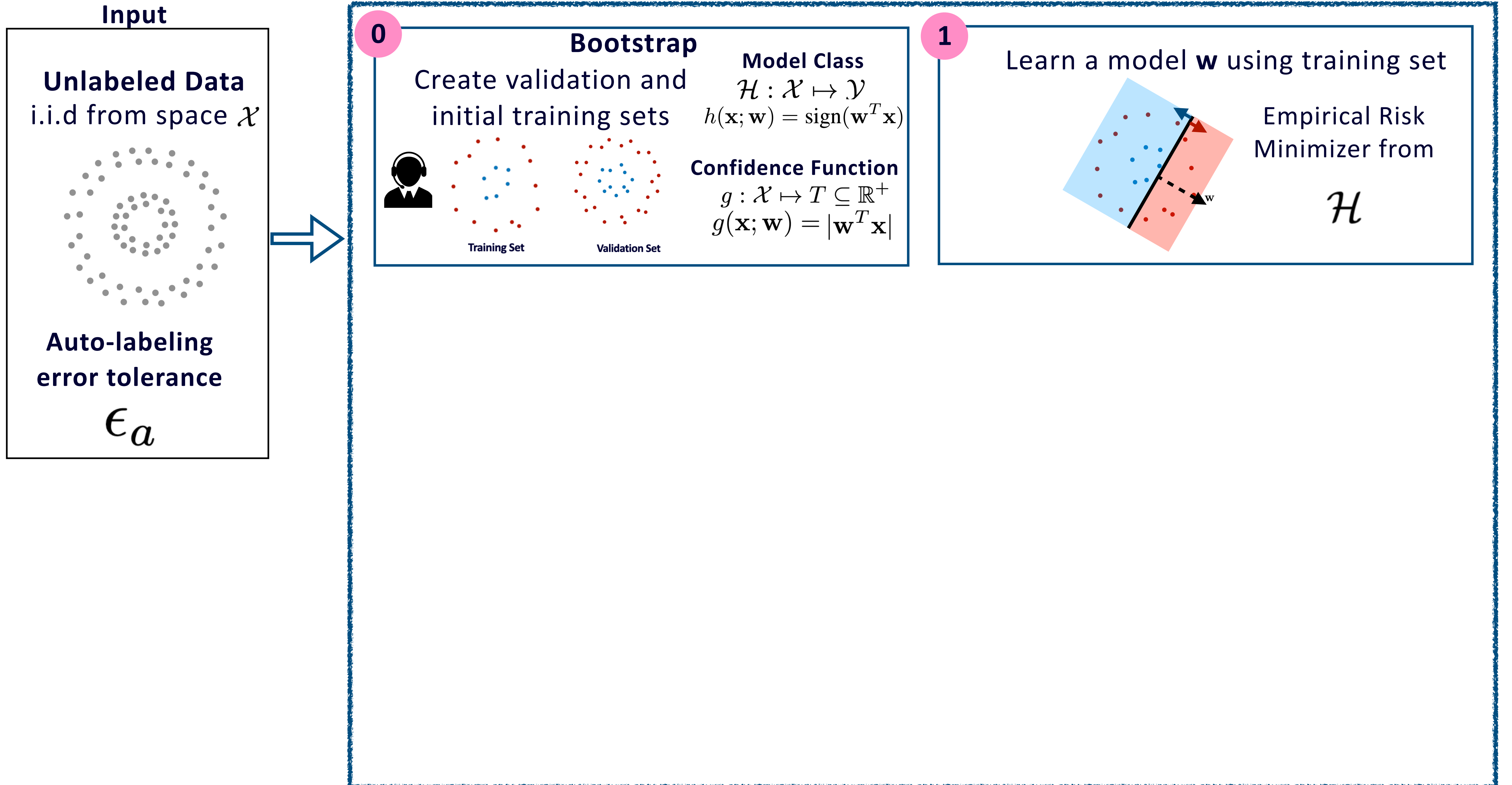
$$D_{train} = \{(\mathbf{x}_i, y_i) : i \in I_{train}\}$$

Start small and gradually add more

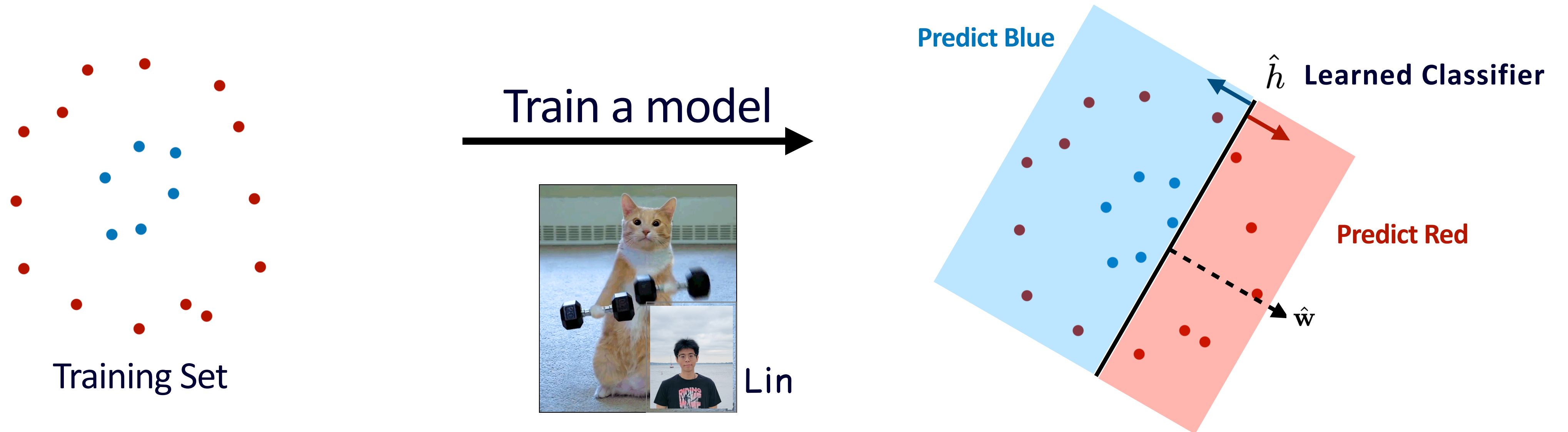
$$D_{val} = \{(\mathbf{x}_i, y_i) : i \in I_{val}\}$$

Get “sufficiently” large amount of it.

Threshold-based Auto-labeling Workflow(TBAL)



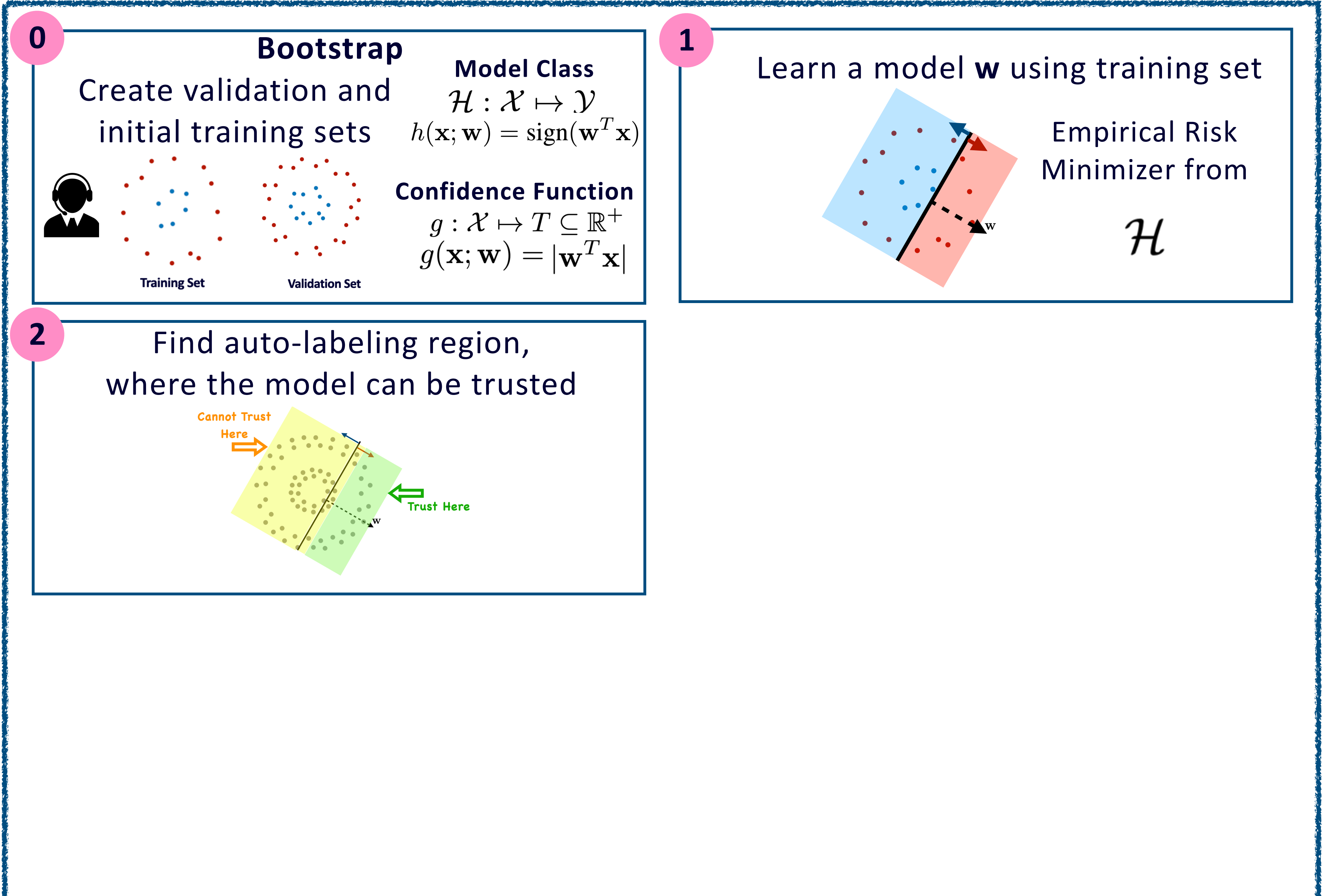
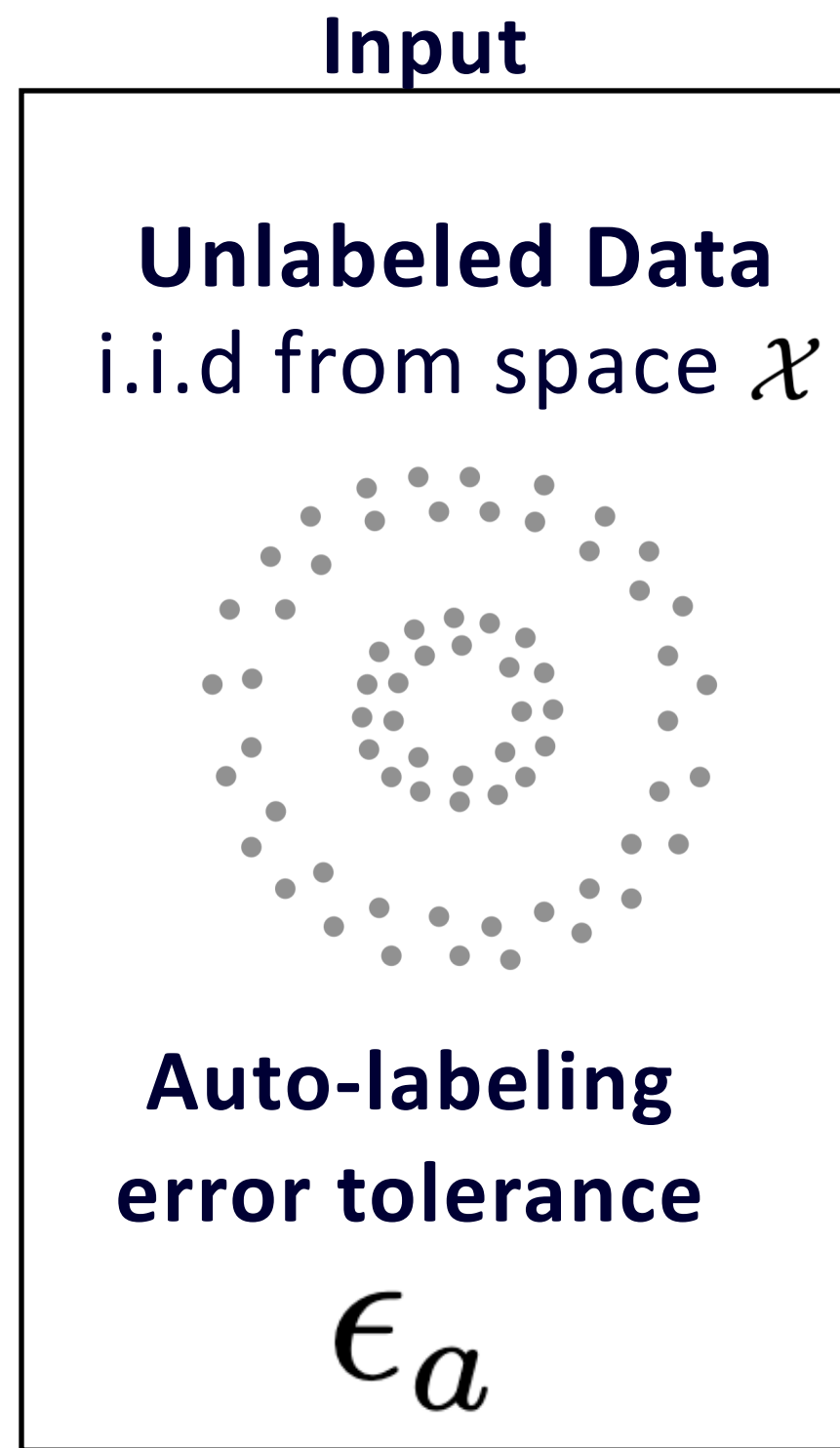
TBAL Workflow : Step 1 Model training



$$\hat{h} = \text{EmpiricalRiskMinimizer}(\mathcal{H}, D_{train})$$
$$\hat{h} = \arg \min_{h \in \mathcal{H}} \frac{1}{|D_{train}|} \sum_{(\mathbf{x}_i, y_i) \in D_{train}} \mathbb{1}\{h(\mathbf{x}_i) \neq y_i\}$$

In practice, usually some surrogate loss is minimized

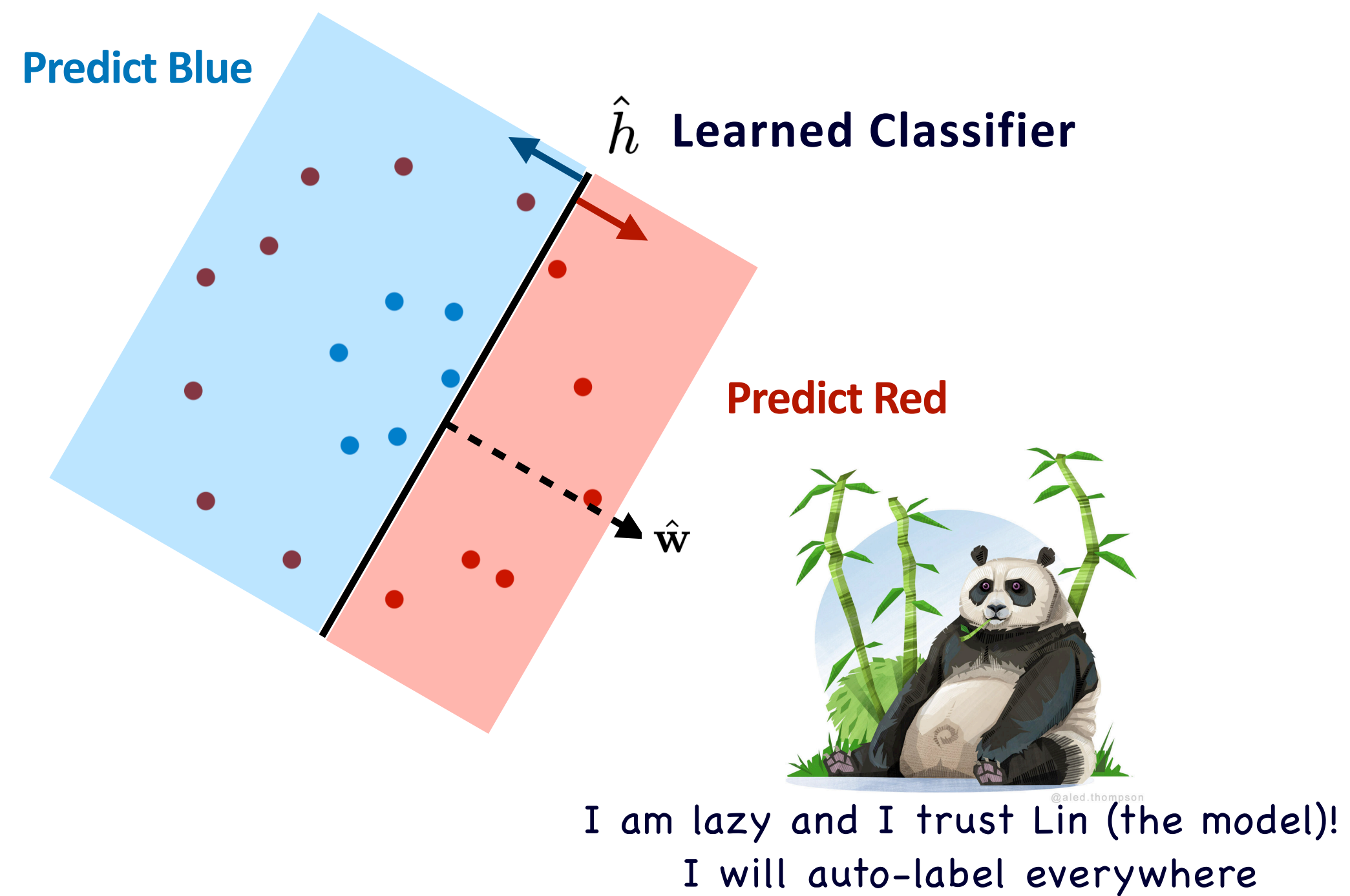
Threshold-based Auto-labeling Workflow(TBAL)



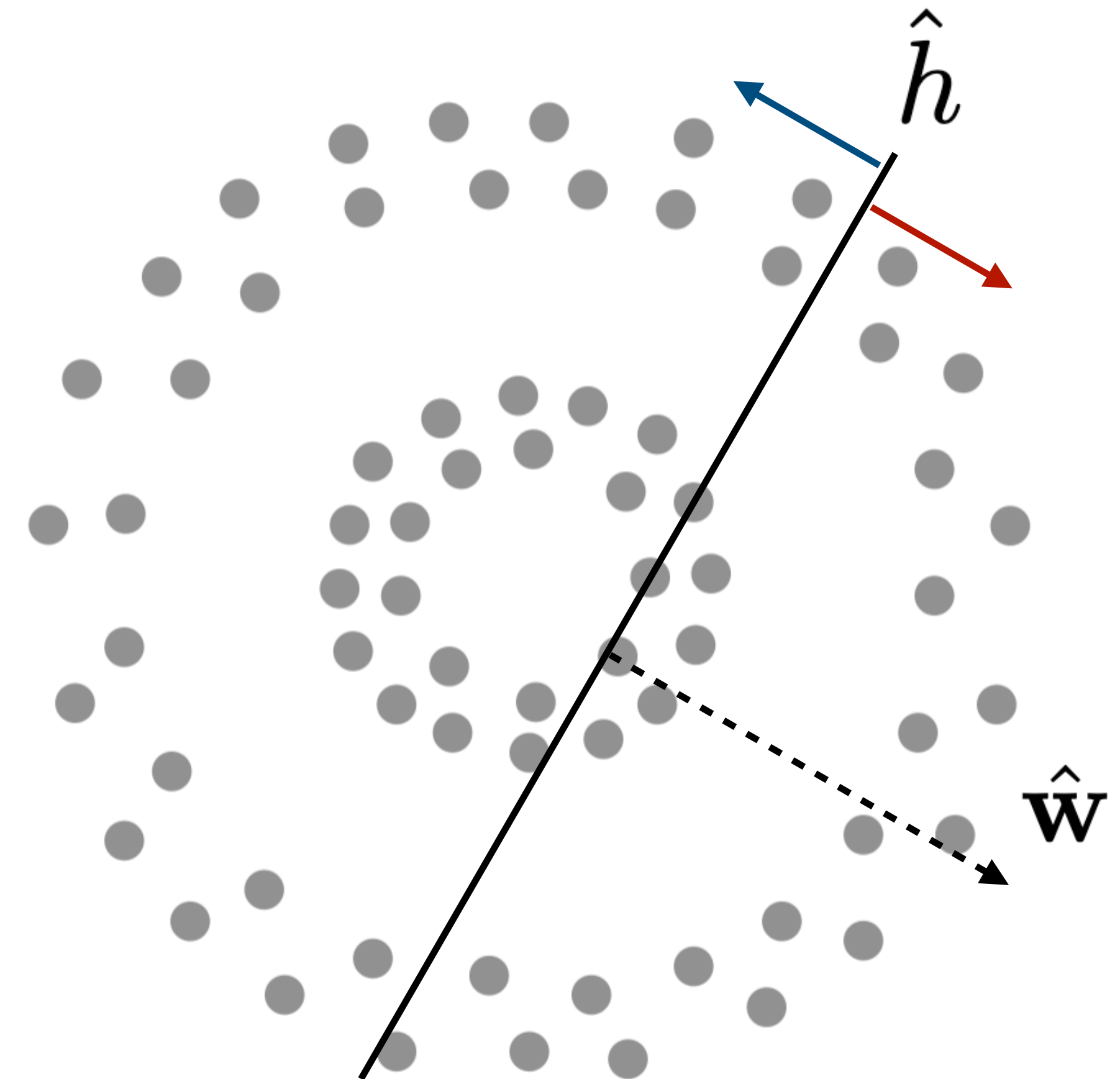
TBAL Workflow: Step 2

Find the Auto-labeling region

Idea 1: Auto-label everywhere.



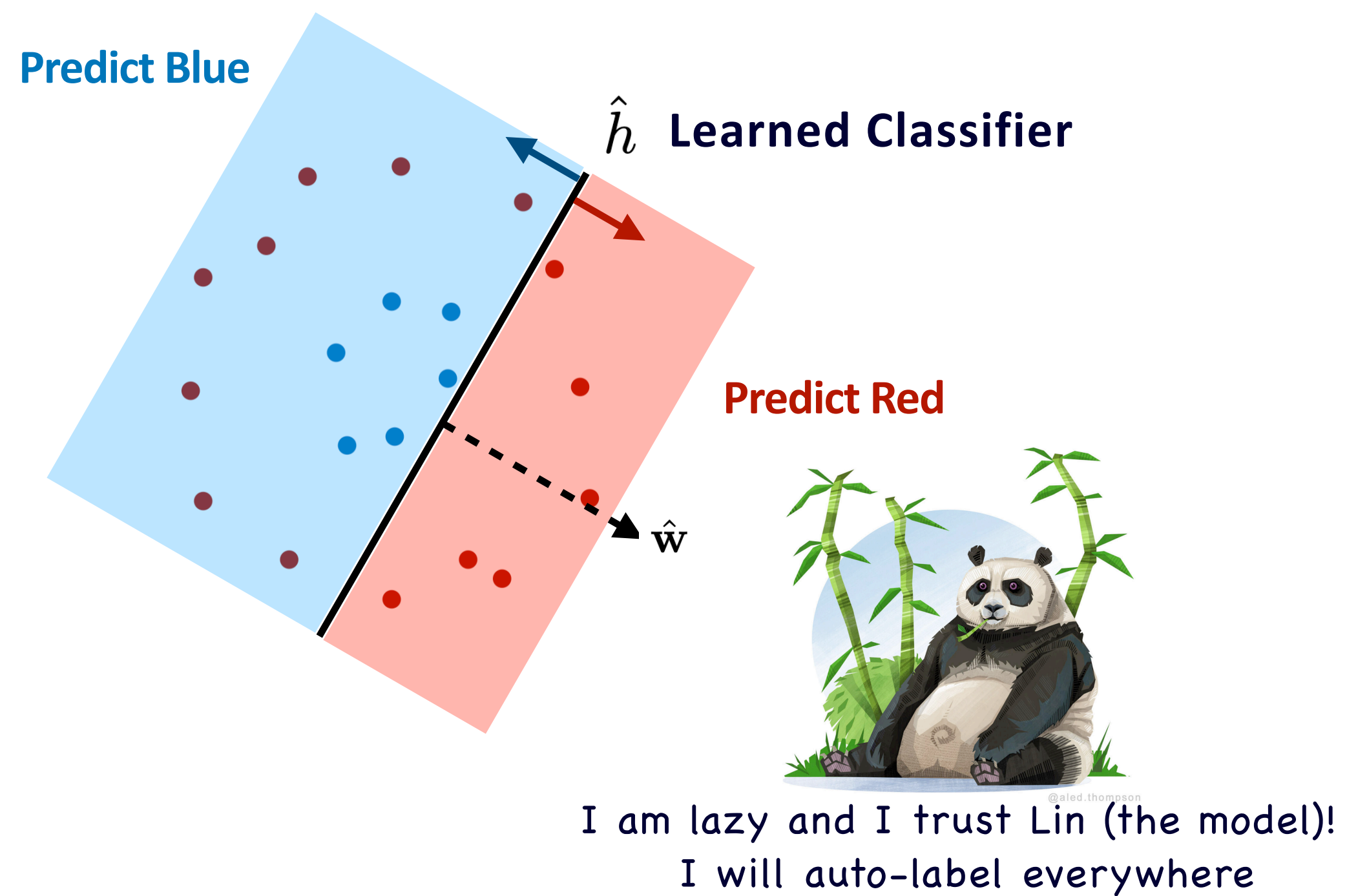
- ● Human-labeled
- ○ Auto-labeled
- ✕ Labeling mistake



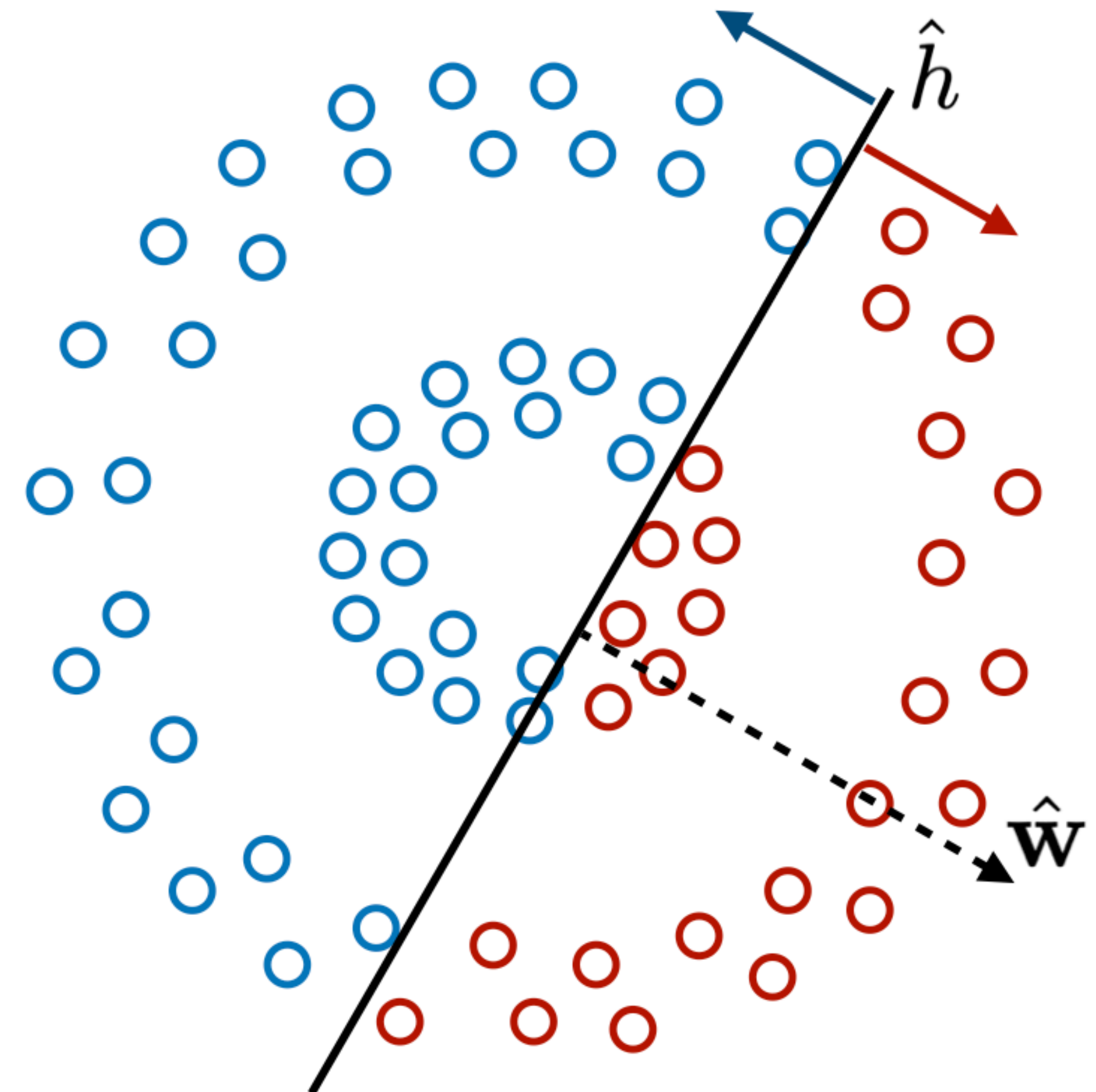
TBAL Workflow: Step 2

Find the Auto-labeling region

Idea 1: Auto-label everywhere.



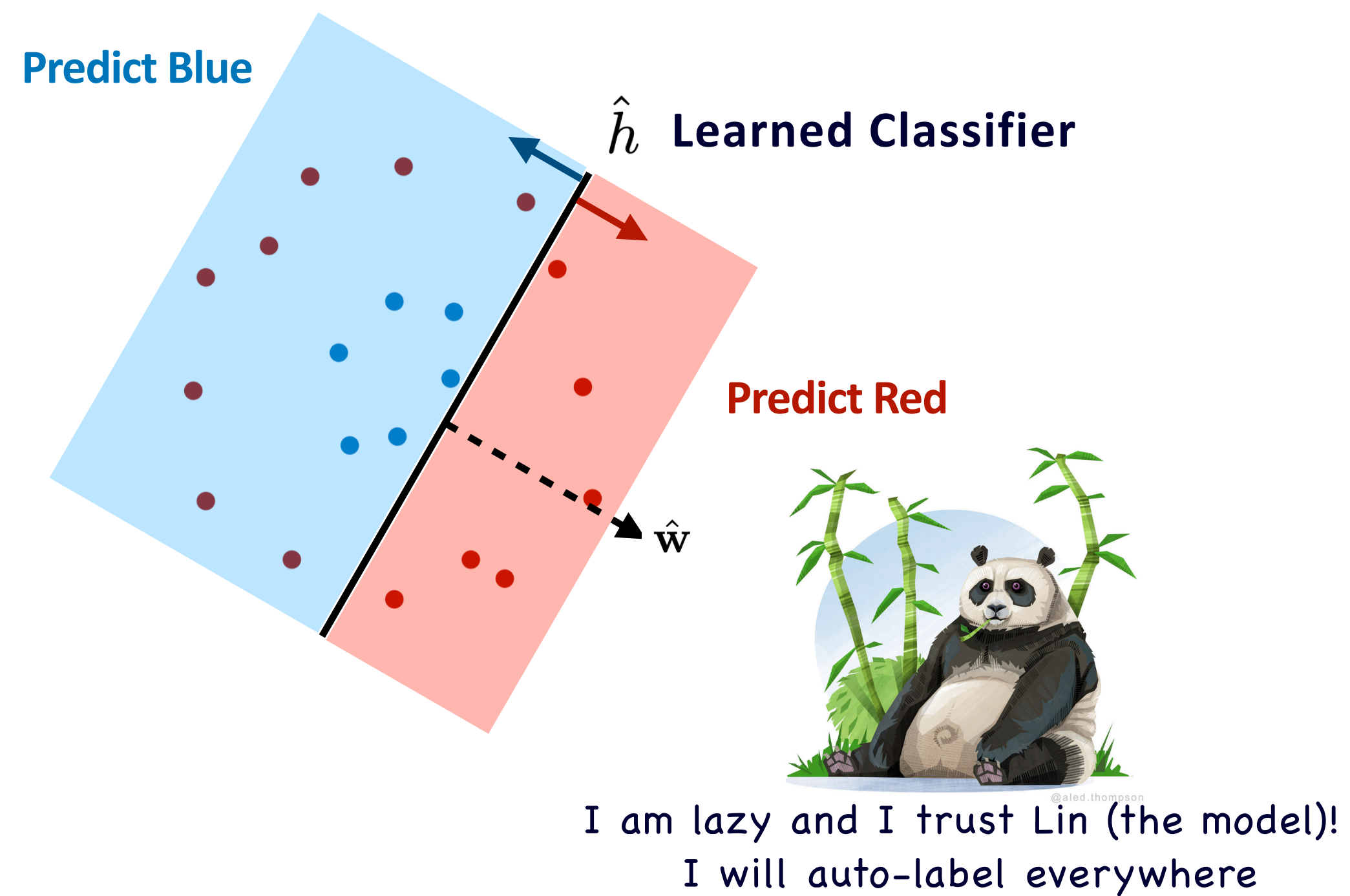
- ● Human-labeled
- ○ Auto-labeled
- ✗ Labeling mistake



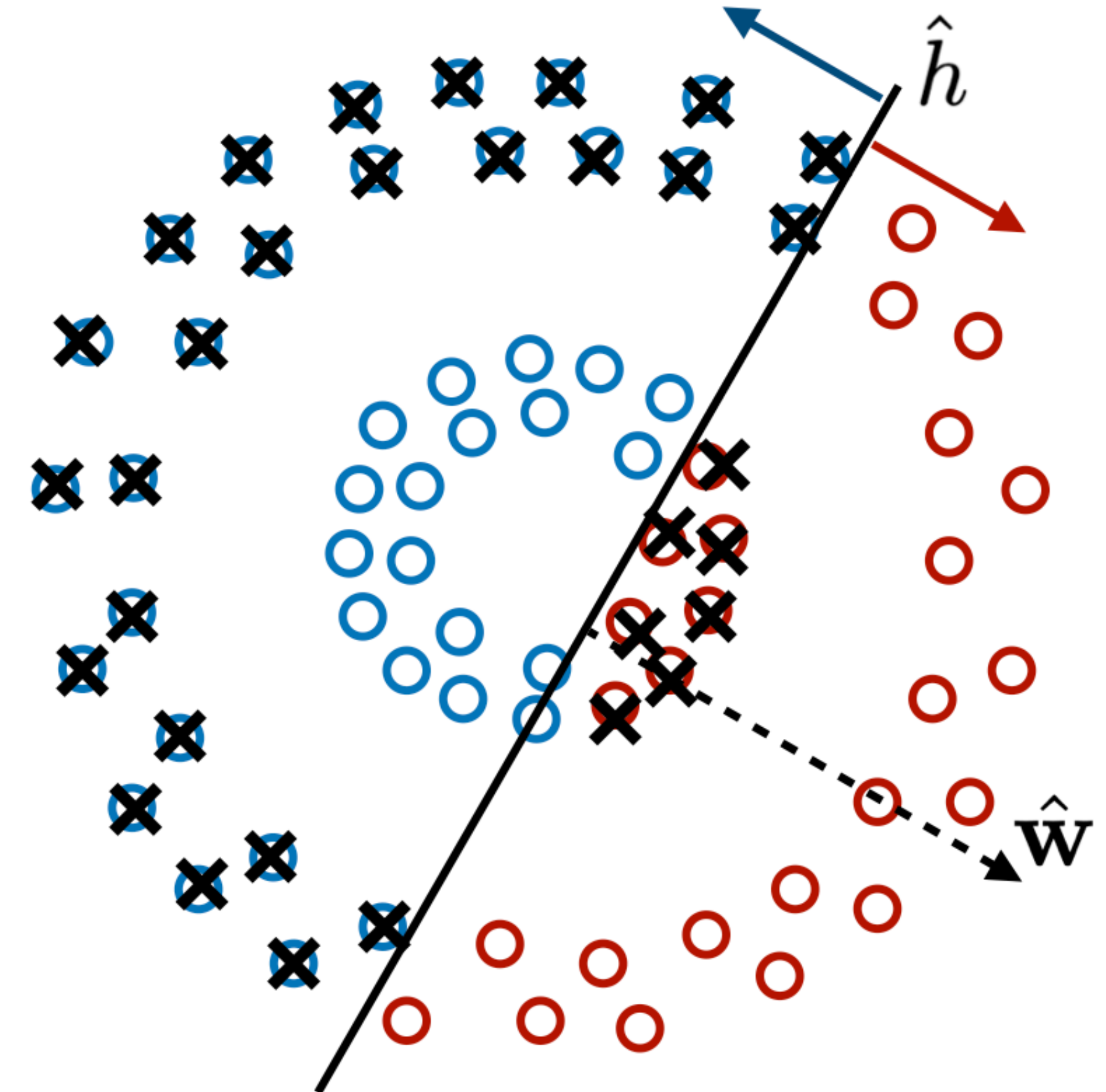
TBAL Workflow: Step 2

Find the Auto-labeling region

Idea 1: Auto-label everywhere.



- Human-labeled
- Auto-labeled
- ✕ Labeling mistake



Could lead to high auto-labeling errors!

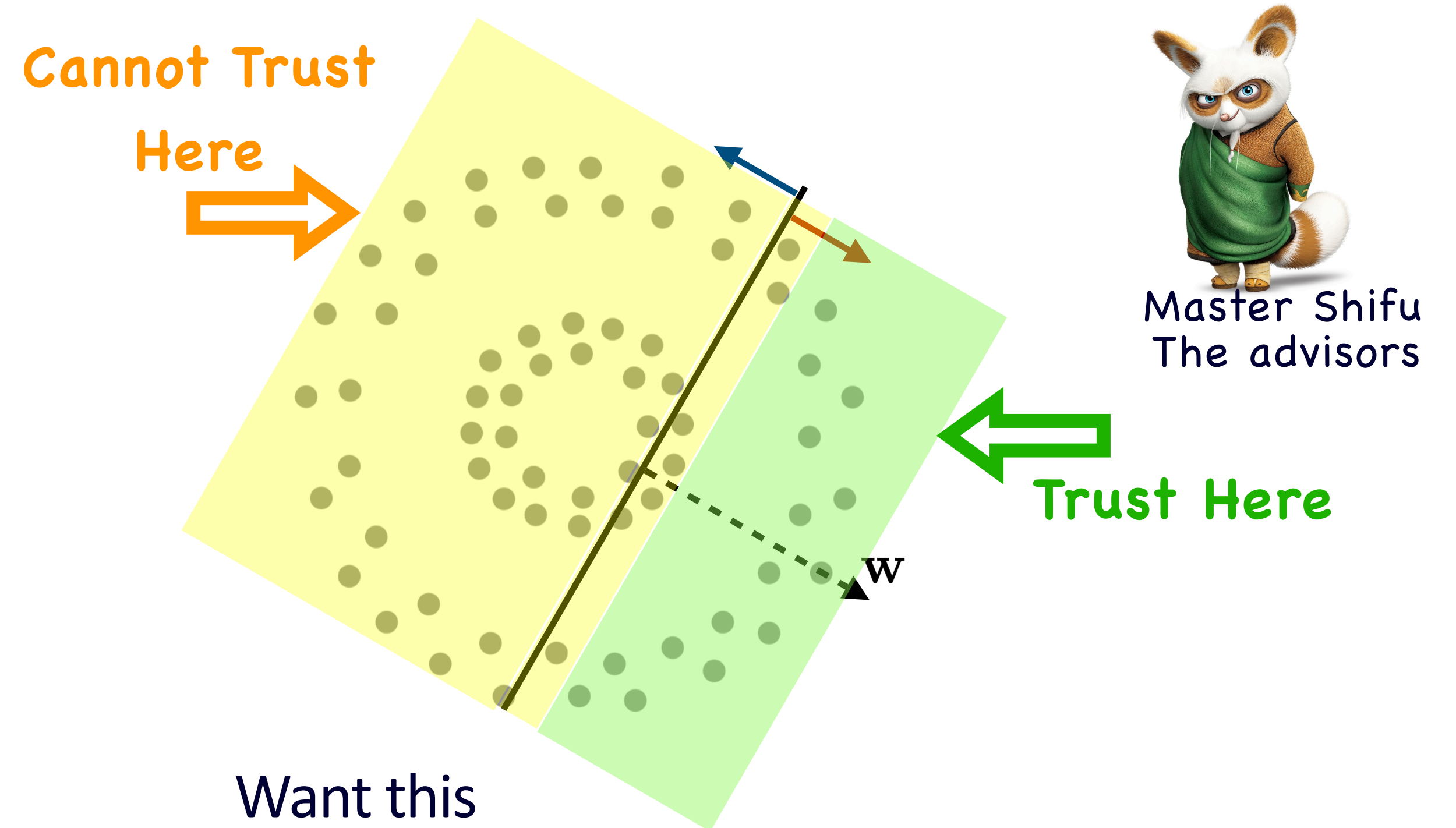
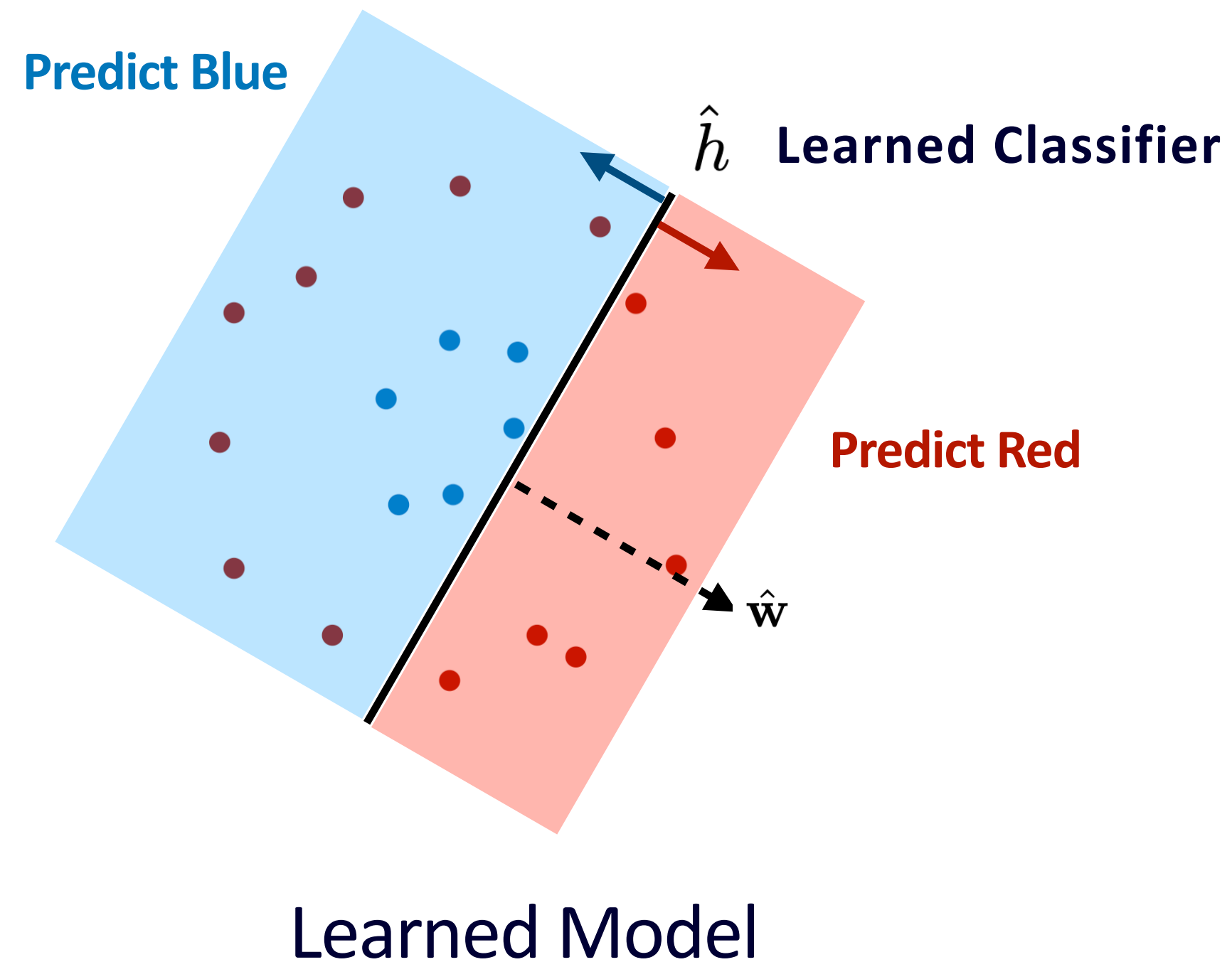


Panda's strategy does not work,
he goes to Master Shifu for advice.

TBAL Workflow: Step 2

Find the Auto-labeling region

Idea 2: Auto-label where the model is accurate (or trustworthy?)

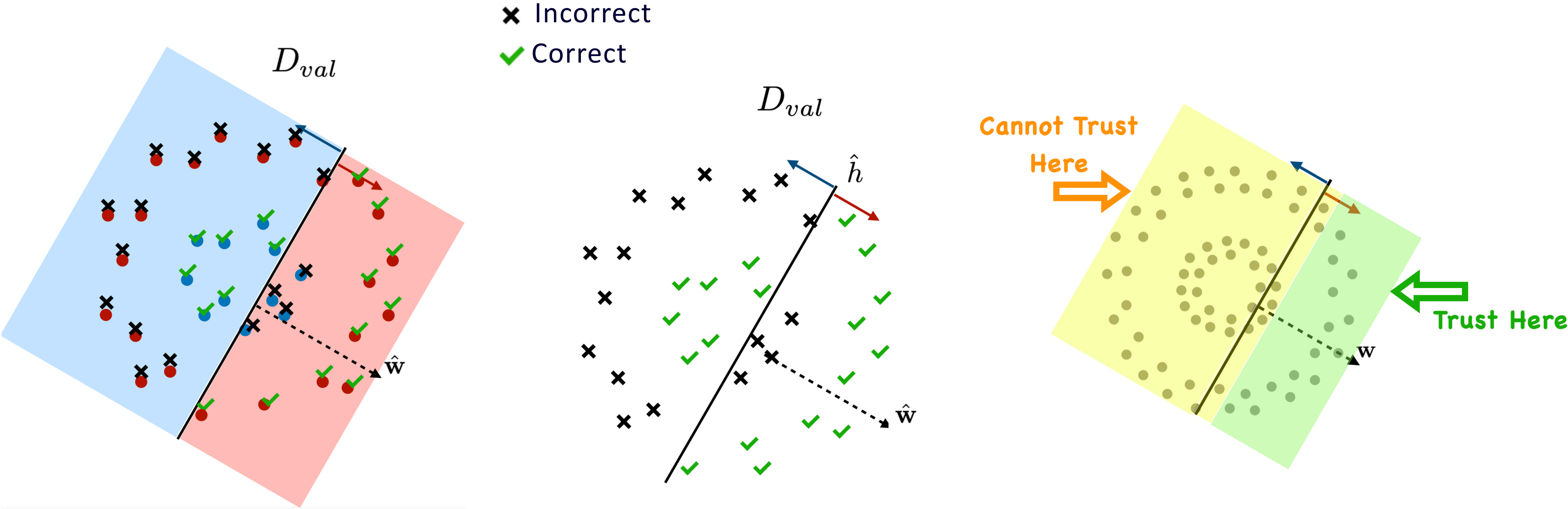


How to find the yellow and green regions?

TBAL Workflow: Step 2

Find the Auto-labeling region

Use the **validation data** to find the region where the classifier can be trusted



TBAL Workflow: Step 2

Find the Auto-labeling region

$$g(\mathbf{x}; \hat{\mathbf{w}}) = |\hat{\mathbf{w}}^T \mathbf{x}|$$

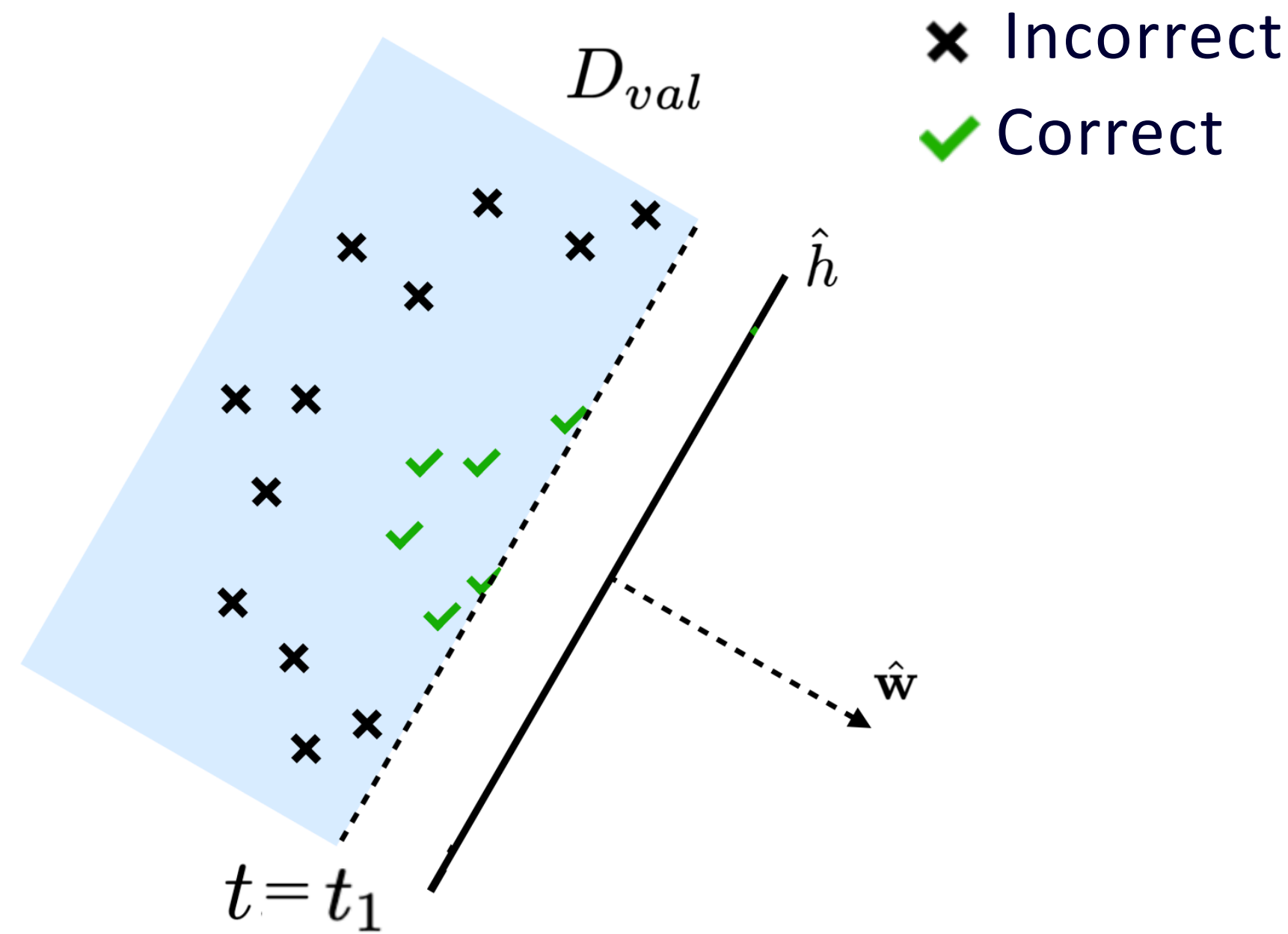
Regions defined by the confidence function

$$A_v(\hat{\mathbf{w}}, t, y) = \{\mathbf{x} \in X_v : g(\mathbf{x}; \hat{\mathbf{w}}) \geq t, \hat{h}(\mathbf{x}, \hat{\mathbf{w}}) = y\}$$

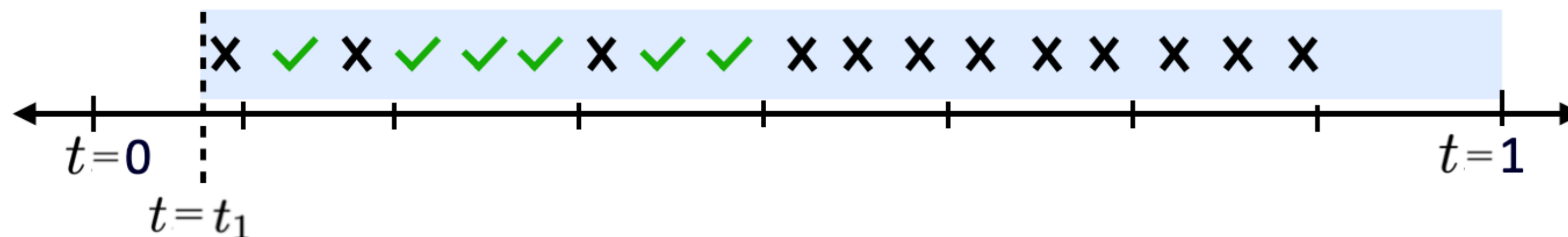
Auto-labeling Error estimation in these regions

$$\hat{\mathcal{E}}_v(\hat{\mathbf{w}}|t, y) = \frac{1}{|A_v(\hat{\mathbf{w}}, t, y)|} \sum_{\mathbf{x} \in A_v(\hat{\mathbf{w}}, t, y)} \mathbb{1}\{\hat{h}(\mathbf{x}; \hat{\mathbf{w}}) \neq f^*(\mathbf{x})\}$$

$$\blacktriangle = \frac{\#\times}{\#\checkmark + \#\times}$$



Predictions sorted by confidence scores

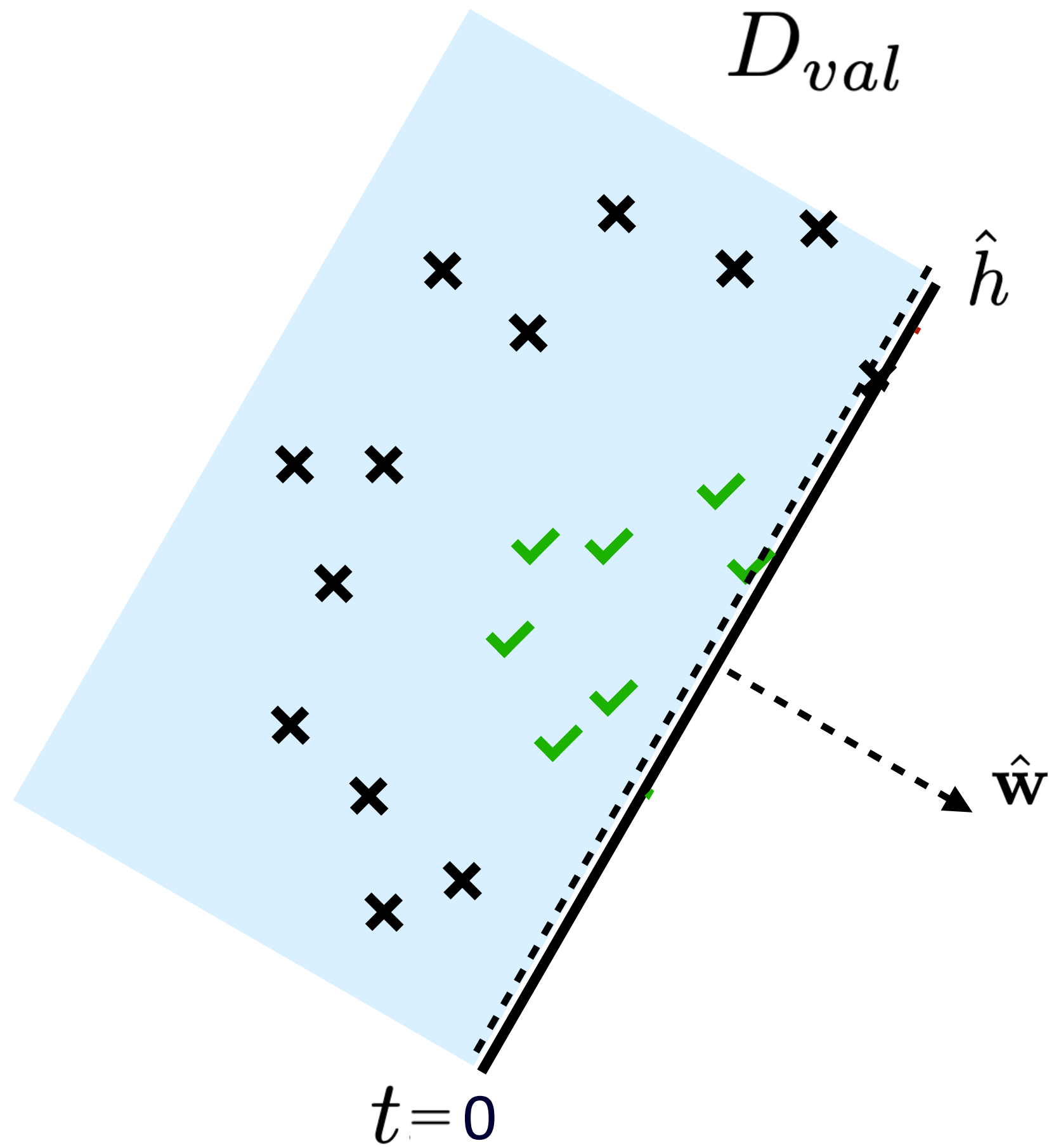


TBAL Workflow: Step 2

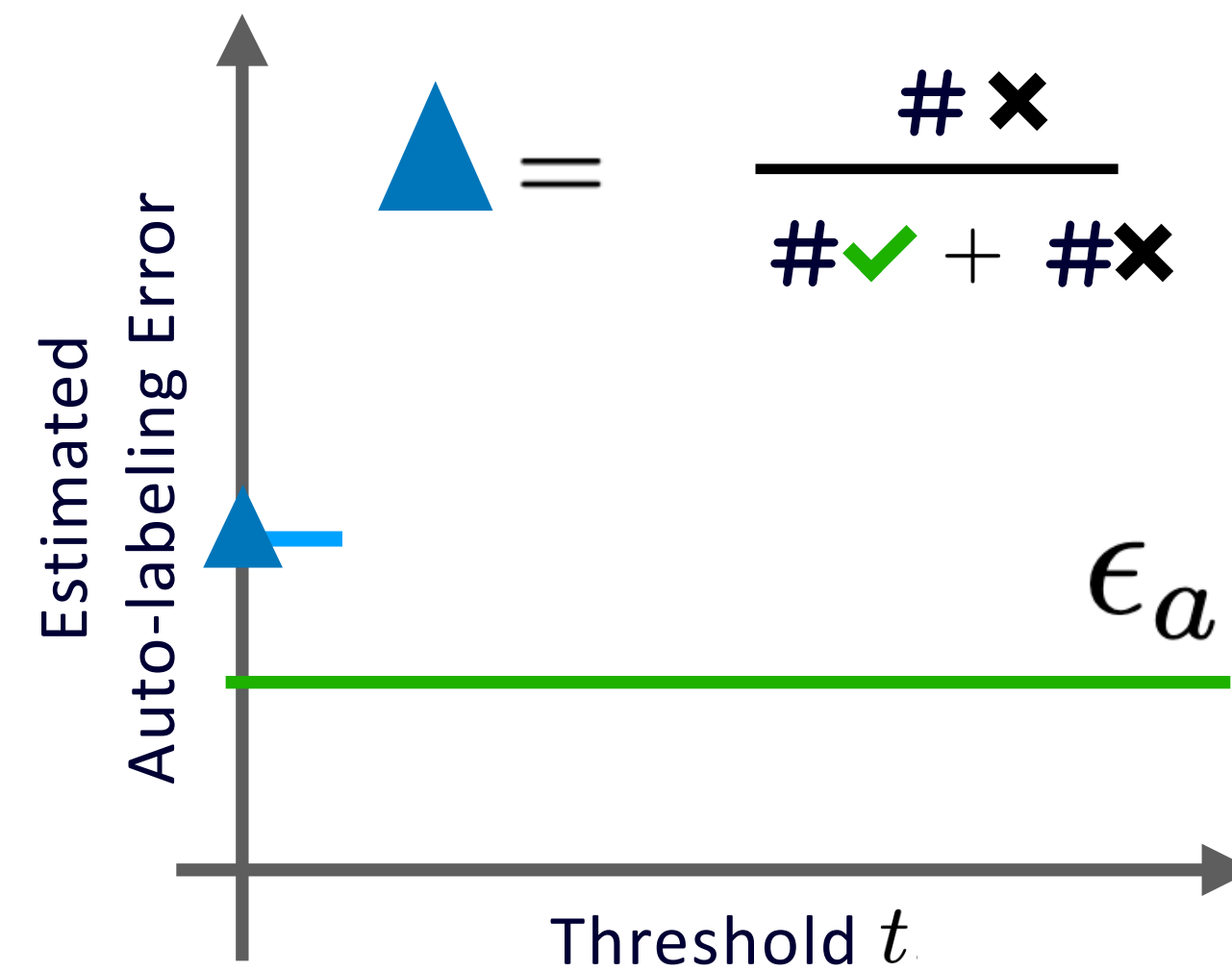
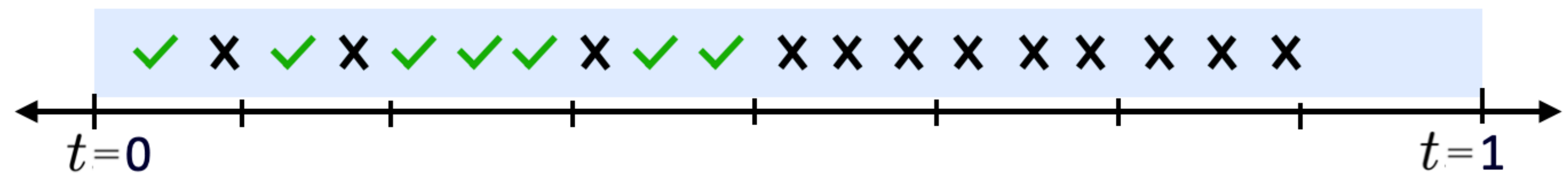
Find the Auto-labeling region

$$g(\mathbf{x}; \hat{\mathbf{w}}) = |\hat{\mathbf{w}}^T \mathbf{x}|$$

$$A_v(\hat{\mathbf{w}}, t, y) = \{\mathbf{x} \in X_v : g(\mathbf{x}; \hat{\mathbf{w}}) \geq t, \hat{h}(\mathbf{x}, \hat{\mathbf{w}}) = y\}$$



Predictions sorted by confidence scores

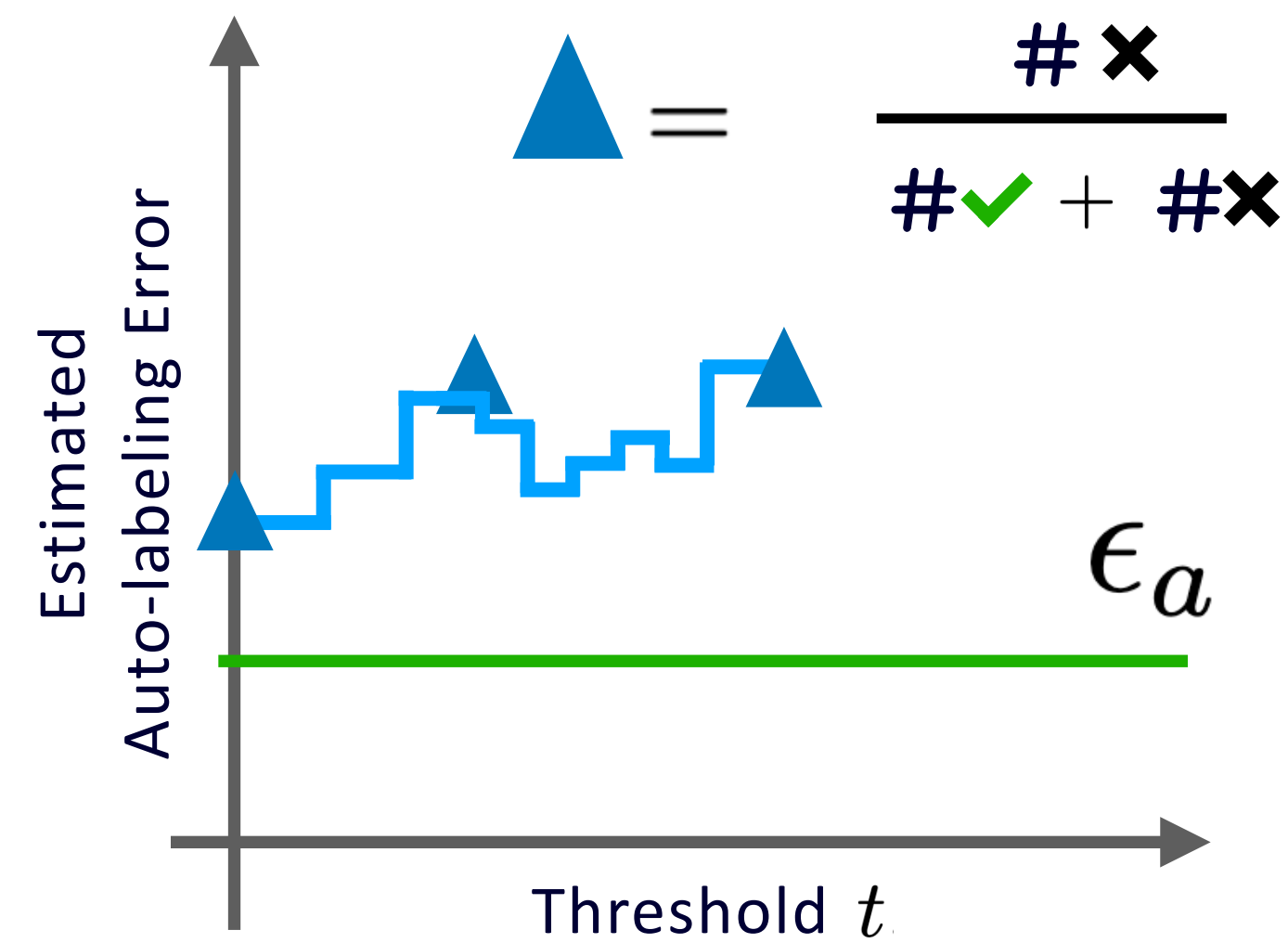
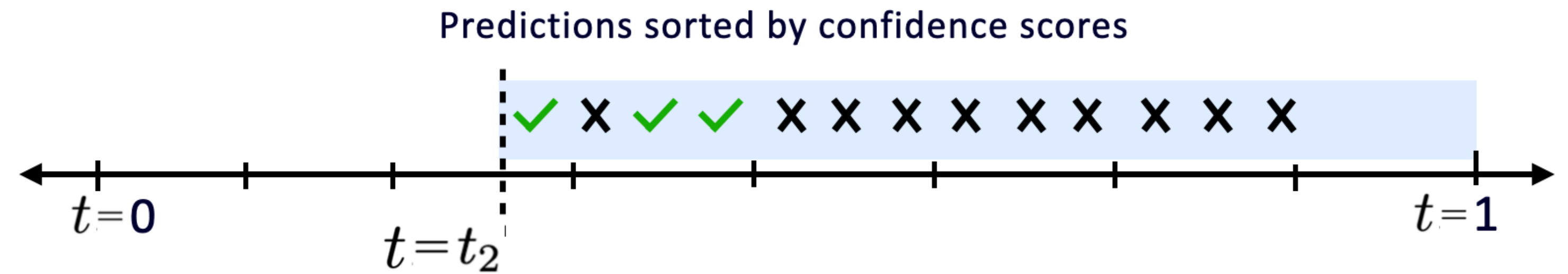
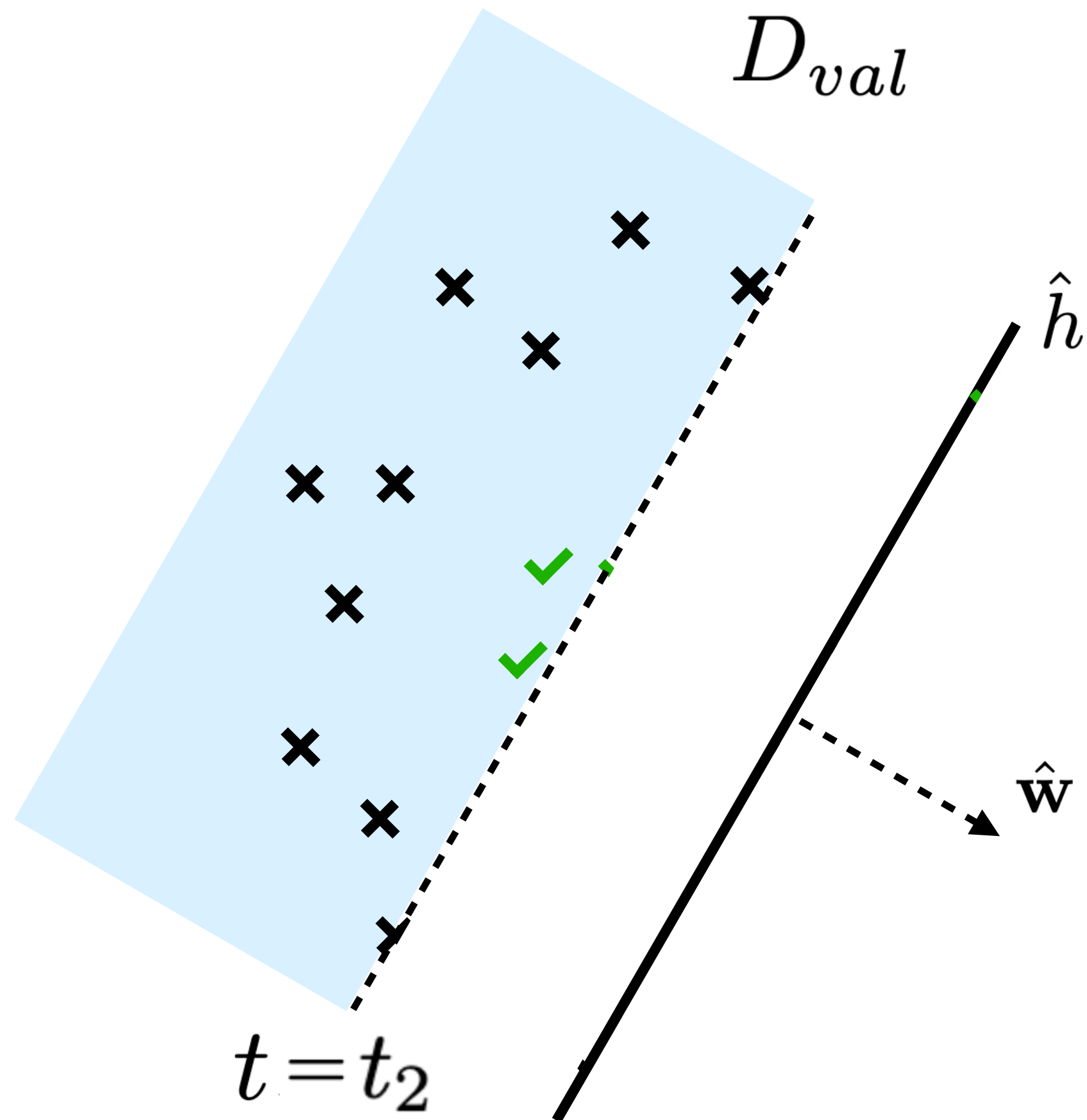


TBAL Workflow: Step 2

Find the Auto-labeling region

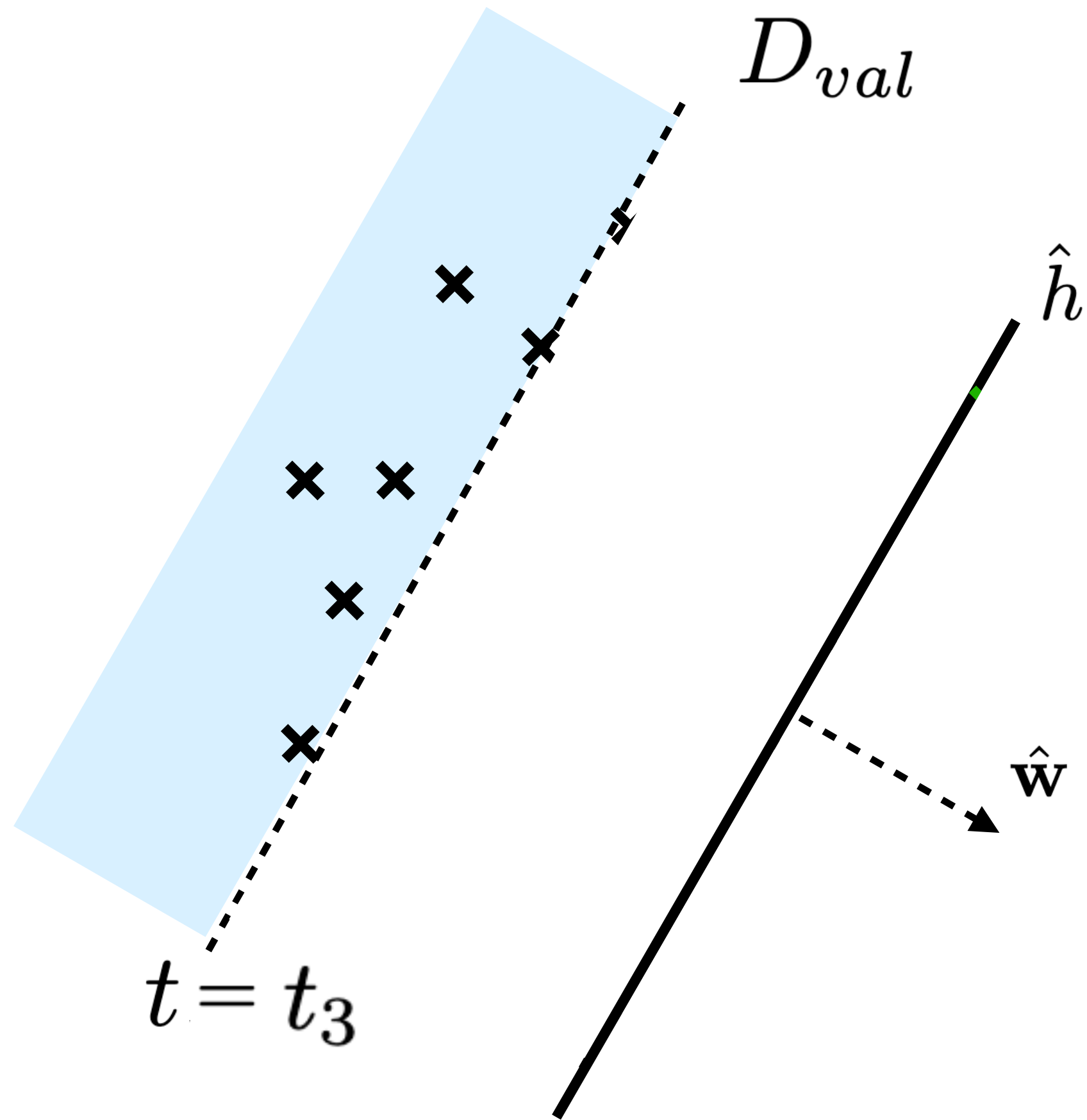
$$g(\mathbf{x}; \hat{\mathbf{w}}) = |\hat{\mathbf{w}}^T \mathbf{x}|$$

$$A_v(\hat{\mathbf{w}}, t, y) = \{\mathbf{x} \in X_v : g(\mathbf{x}; \hat{\mathbf{w}}) \geq t, \hat{h}(\mathbf{x}, \hat{\mathbf{w}}) = y\}$$



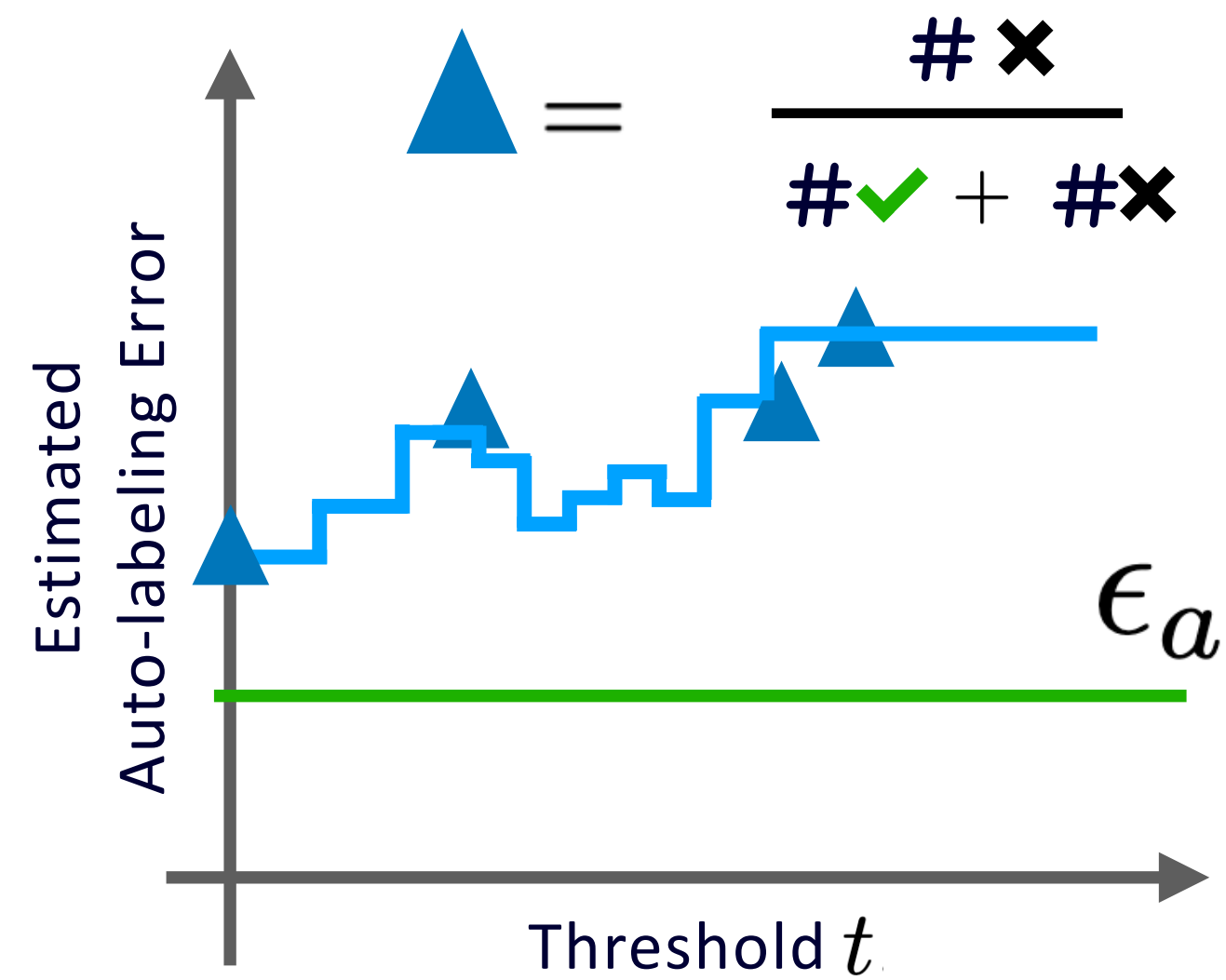
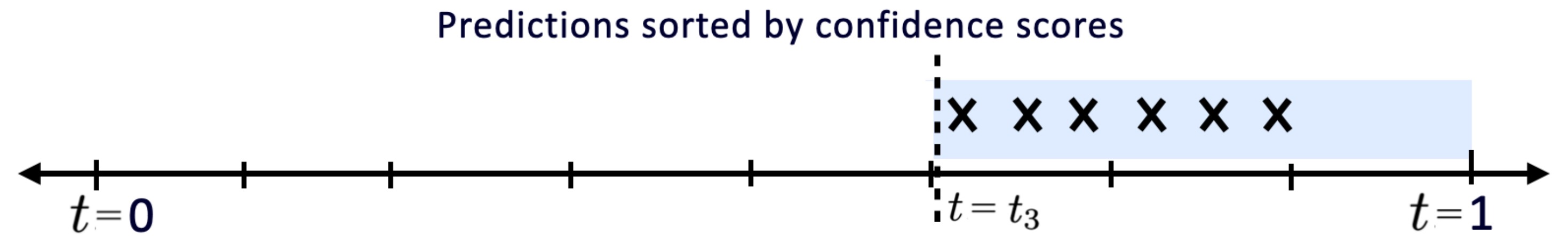
TBAL Workflow: Step 2

Find the Auto-labeling region



$$g(\mathbf{x}; \hat{\mathbf{w}}) = |\hat{\mathbf{w}}^T \mathbf{x}|$$

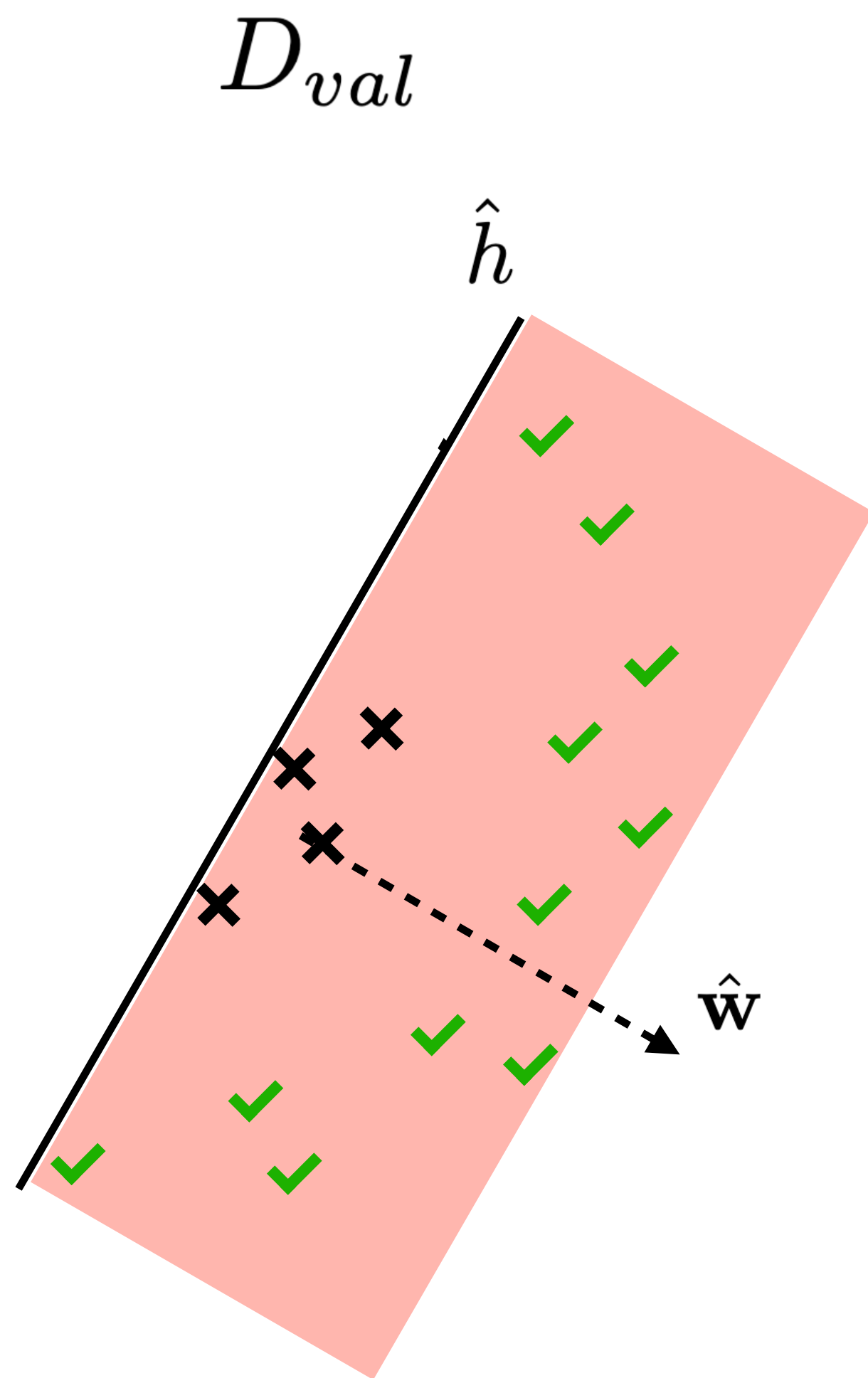
$$A_v(\hat{\mathbf{w}}, t, y) = \{\mathbf{x} \in X_v : g(\mathbf{x}; \hat{\mathbf{w}}) \geq t, \hat{h}(\mathbf{x}, \hat{\mathbf{w}}) = y\}$$



Cannot find a threshold on this side.

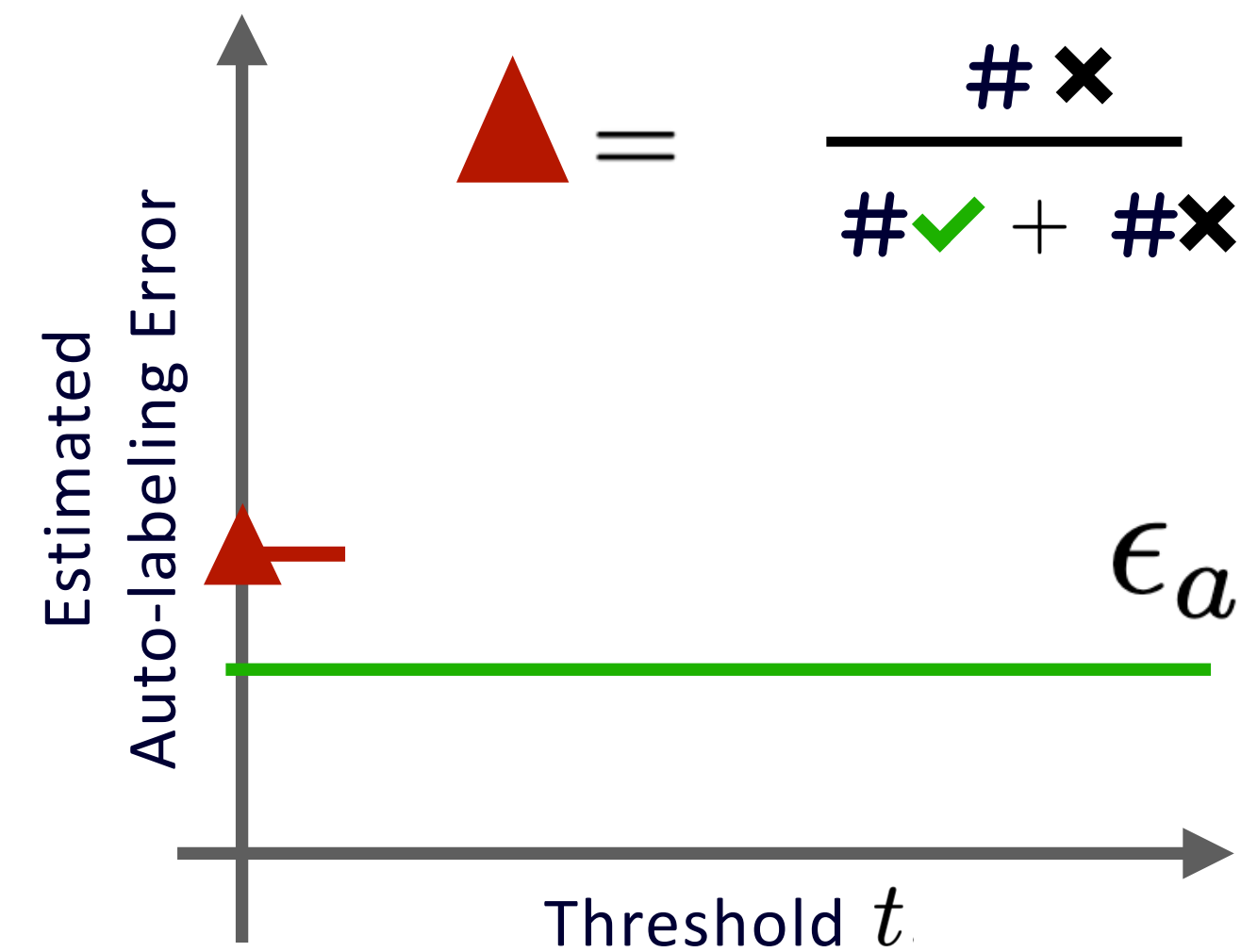
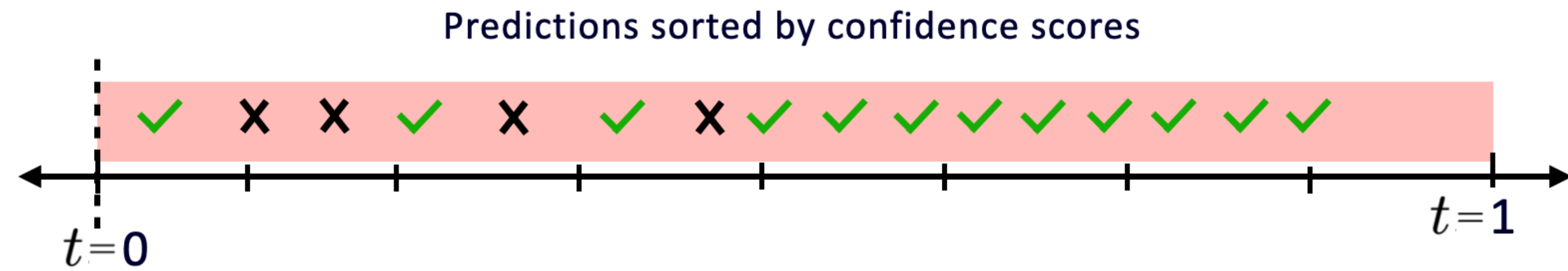
TBAL Workflow: Step 2

Find the Auto-labeling region



$$g(\mathbf{x}; \hat{\mathbf{w}}) = |\hat{\mathbf{w}}^T \mathbf{x}|$$

$$A_v(\hat{\mathbf{w}}, t, y) = \{\mathbf{x} \in X_v : g(\mathbf{x}; \hat{\mathbf{w}}) \geq t, \hat{h}(\mathbf{x}, \hat{\mathbf{w}}) = y\}$$

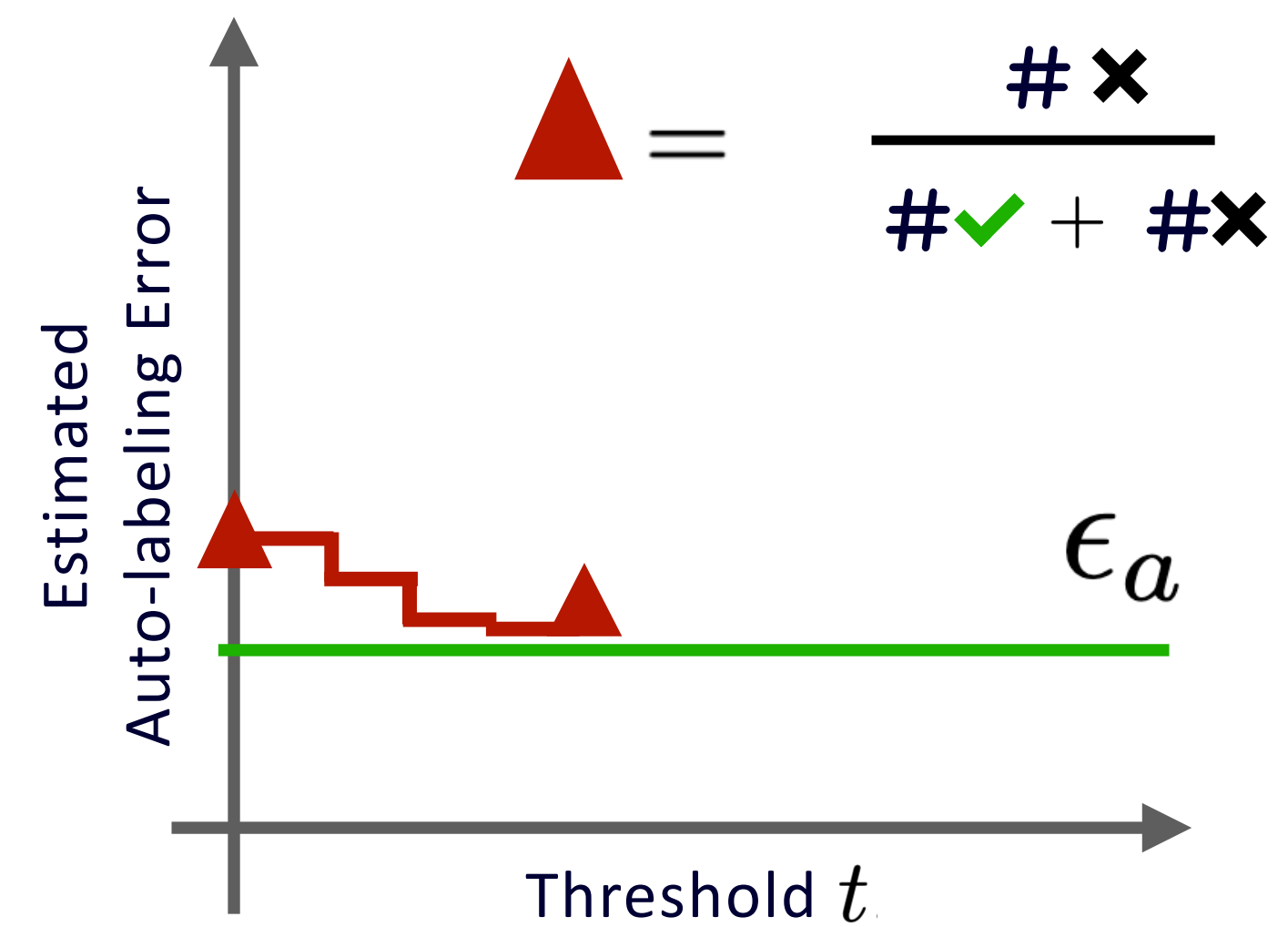
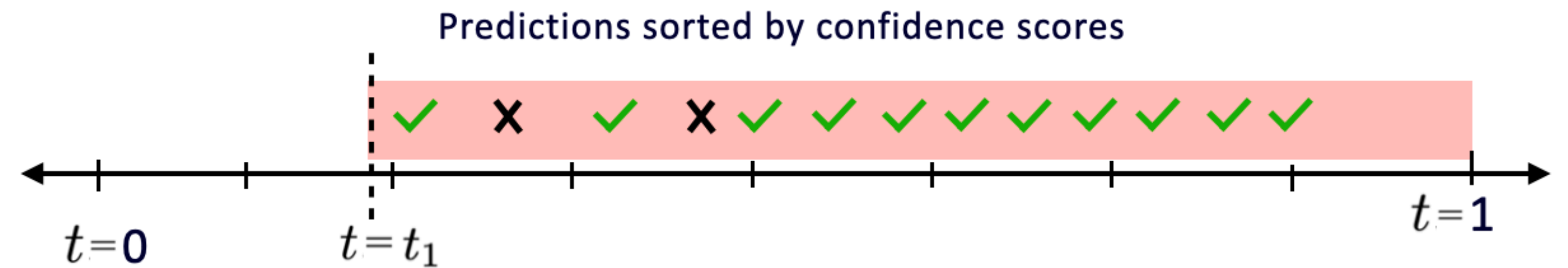
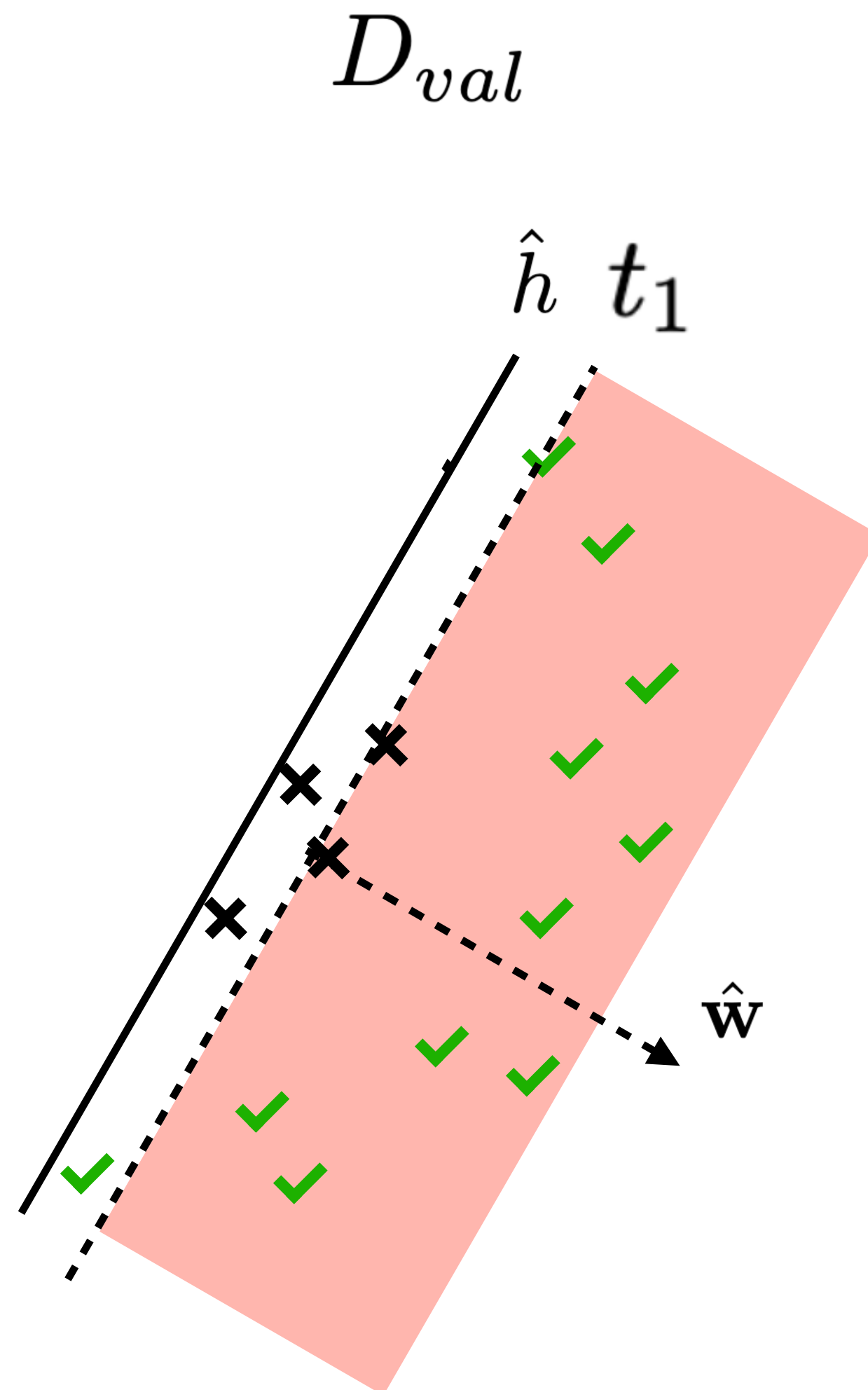


TBAL Workflow: Step 2

Find the Auto-labeling region

$$g(\mathbf{x}; \hat{\mathbf{w}}) = |\hat{\mathbf{w}}^T \mathbf{x}|$$

$$A_v(\hat{\mathbf{w}}, t, y) = \{\mathbf{x} \in X_v : g(\mathbf{x}; \hat{\mathbf{w}}) \geq t, \hat{h}(\mathbf{x}, \hat{\mathbf{w}}) = y\}$$

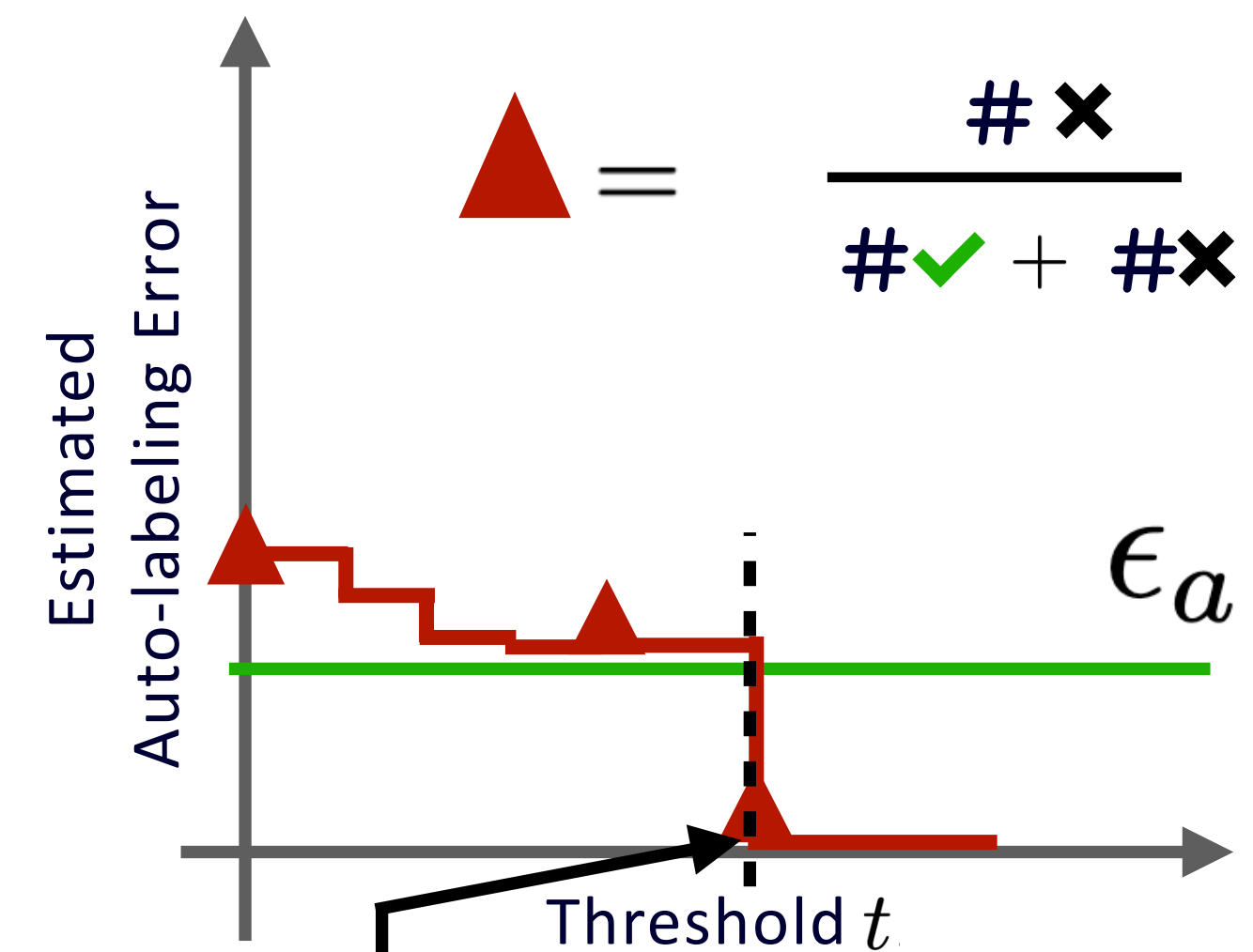
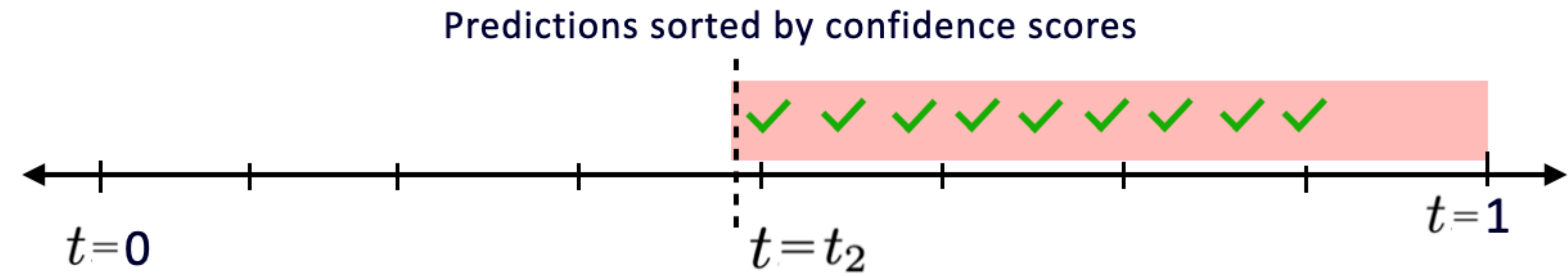
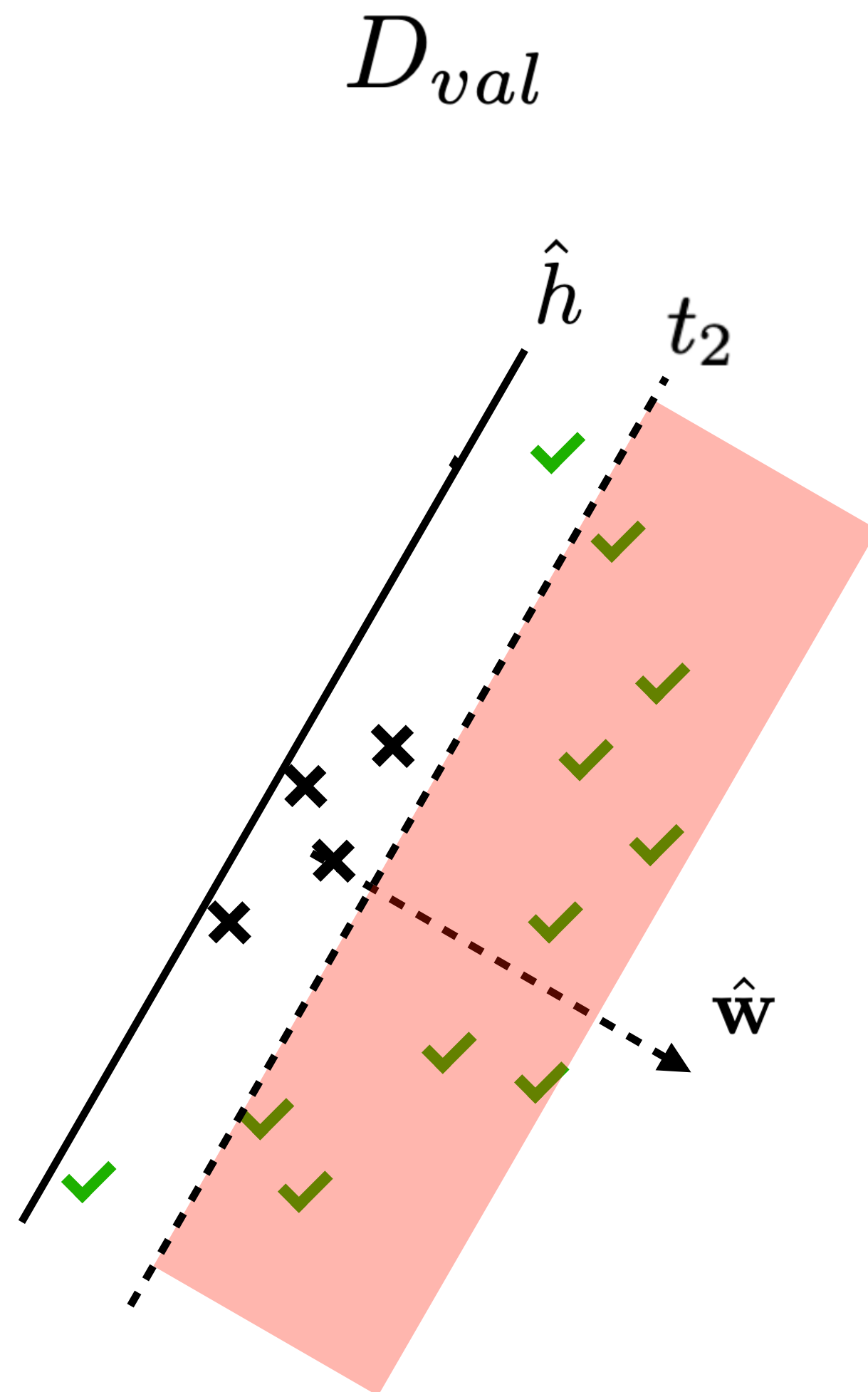


TBAL Workflow: Step 2

Find the Auto-labeling region

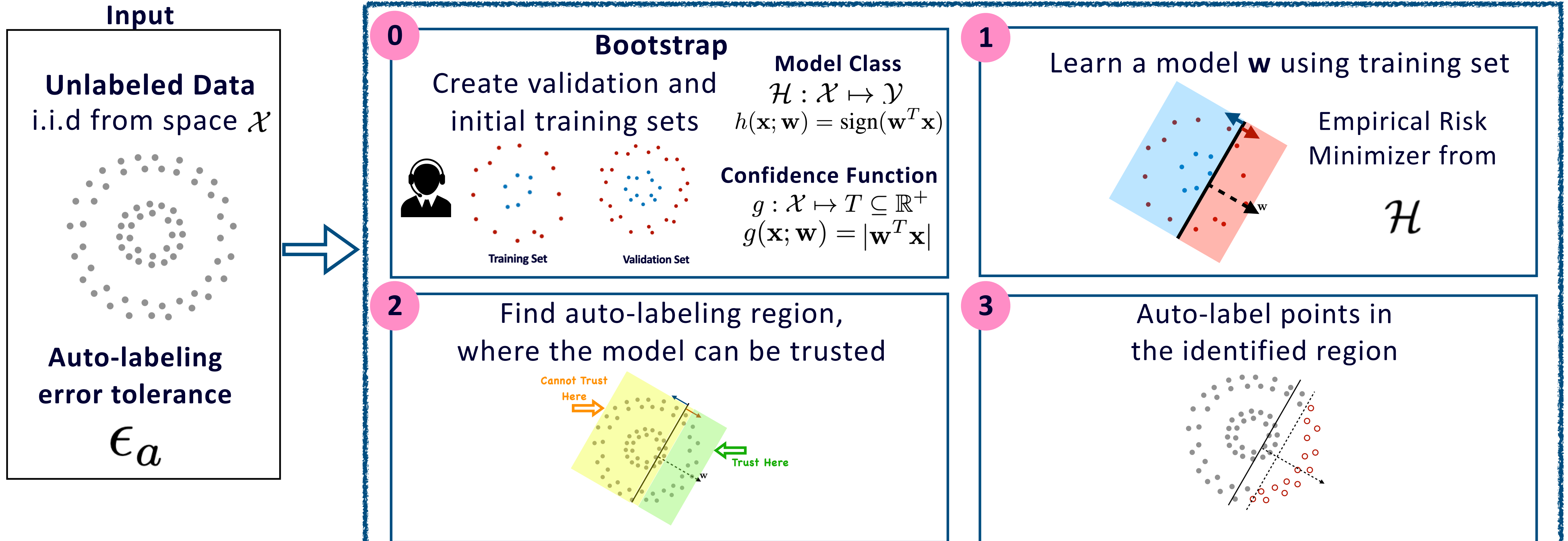
$$g(\mathbf{x}; \hat{\mathbf{w}}) = |\hat{\mathbf{w}}^T \mathbf{x}|$$

$$A_v(\hat{\mathbf{w}}, t, y) = \{\mathbf{x} \in X_v : g(\mathbf{x}; \hat{\mathbf{w}}) \geq t, \hat{h}(\mathbf{x}, \hat{\mathbf{w}}) = y\}$$



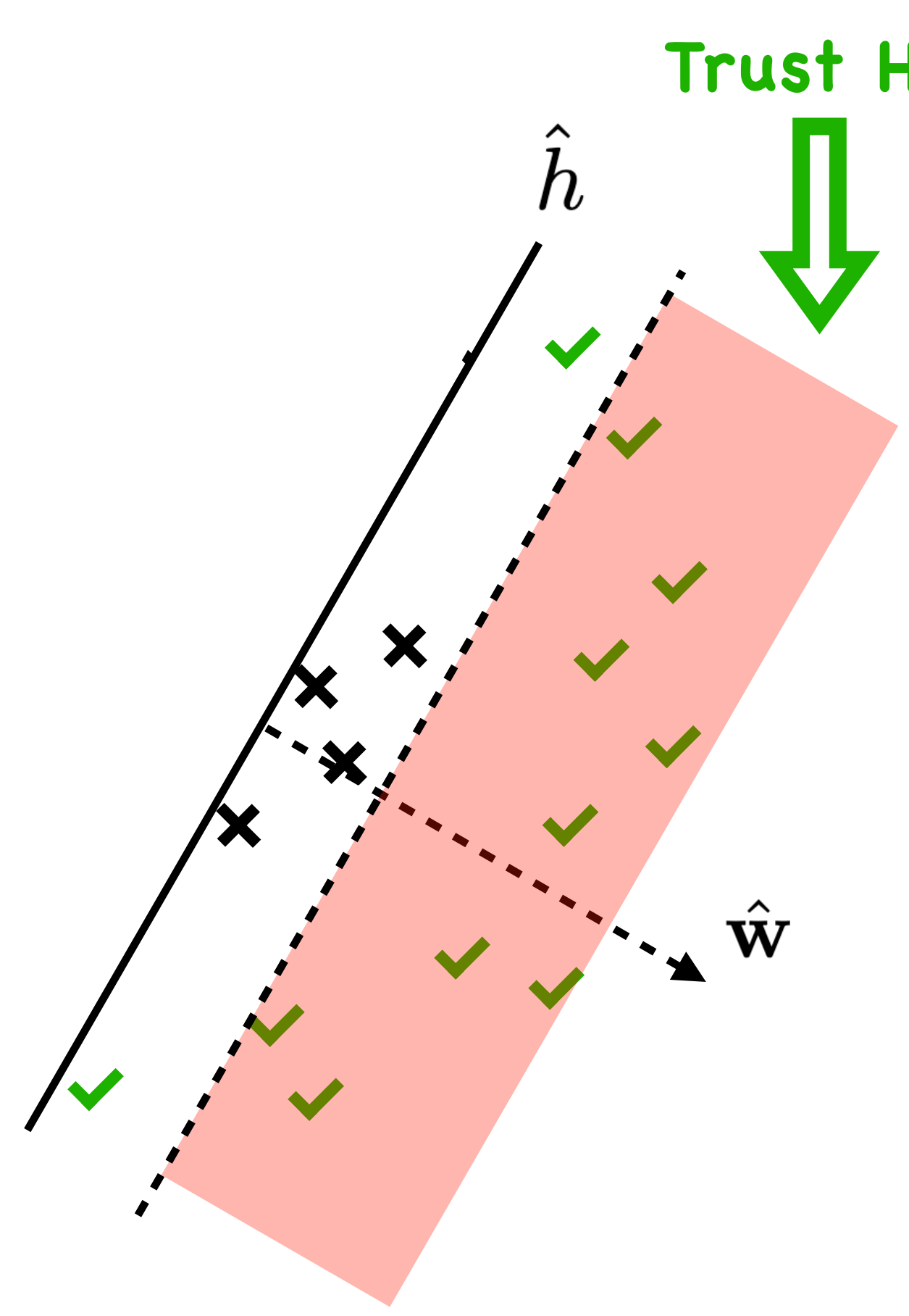
We found a threshold that has error $< \epsilon_a$

Threshold-based Auto-labeling Workflow(TBAL)

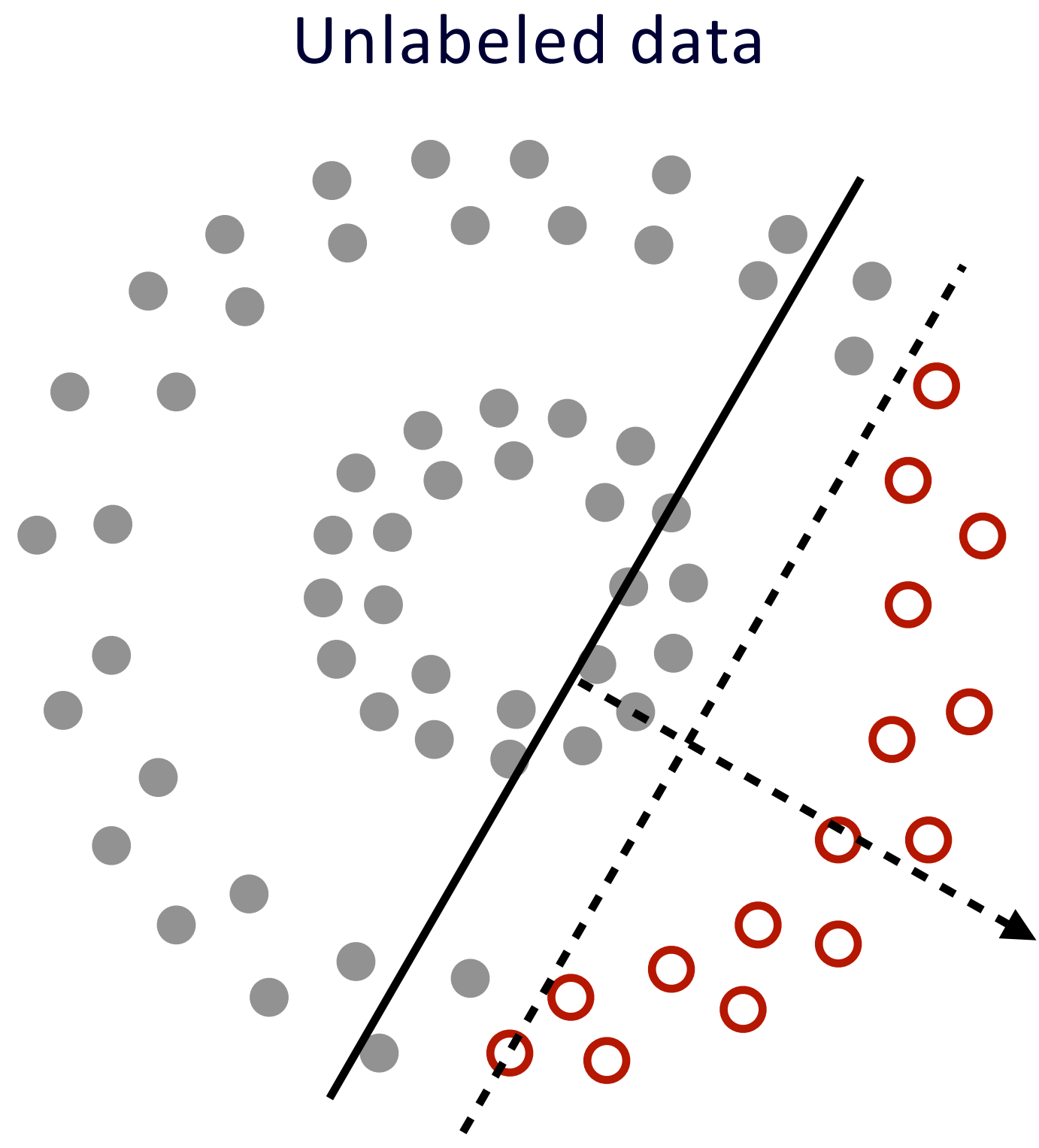


TBAL Workflow: Step 3 Auto-label points in the identified region

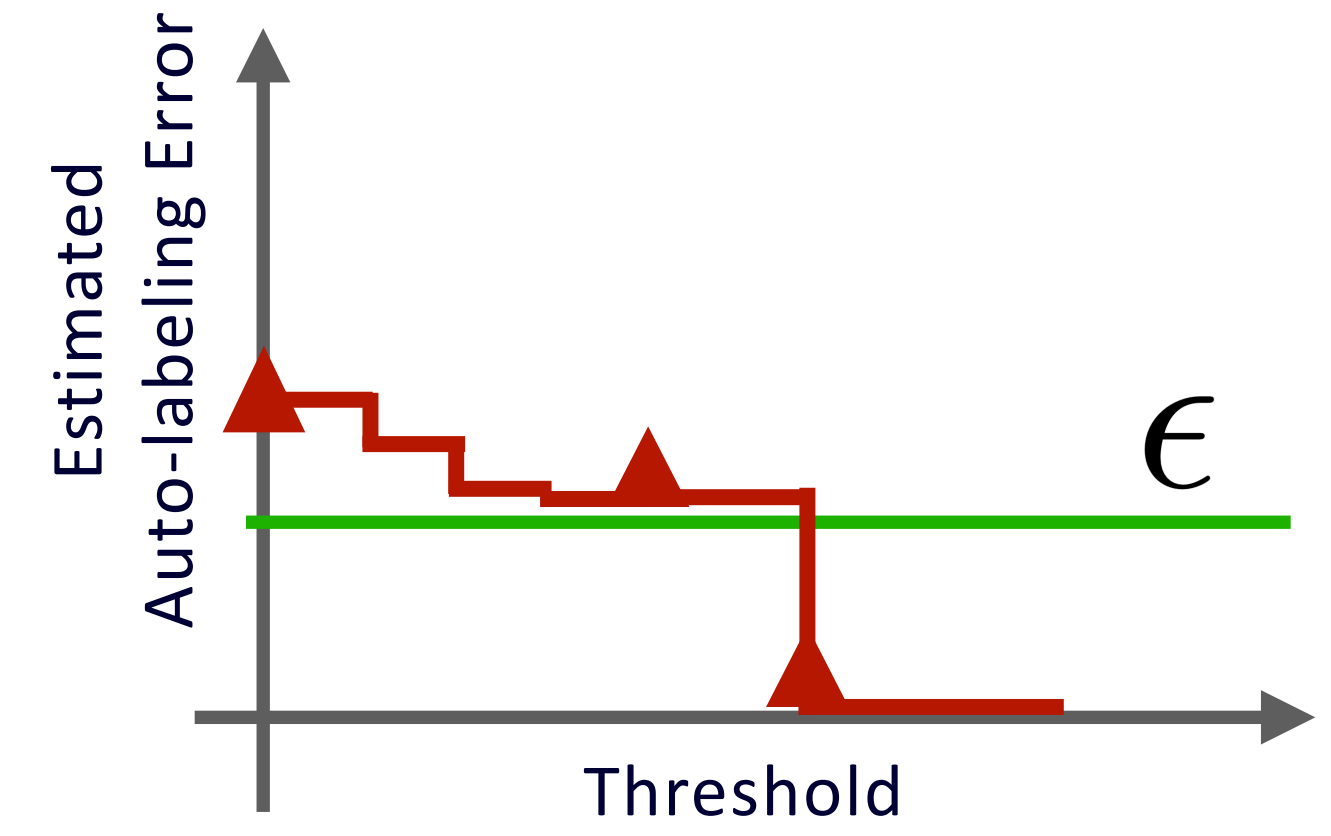
We found a threshold that has error $< \epsilon_a$



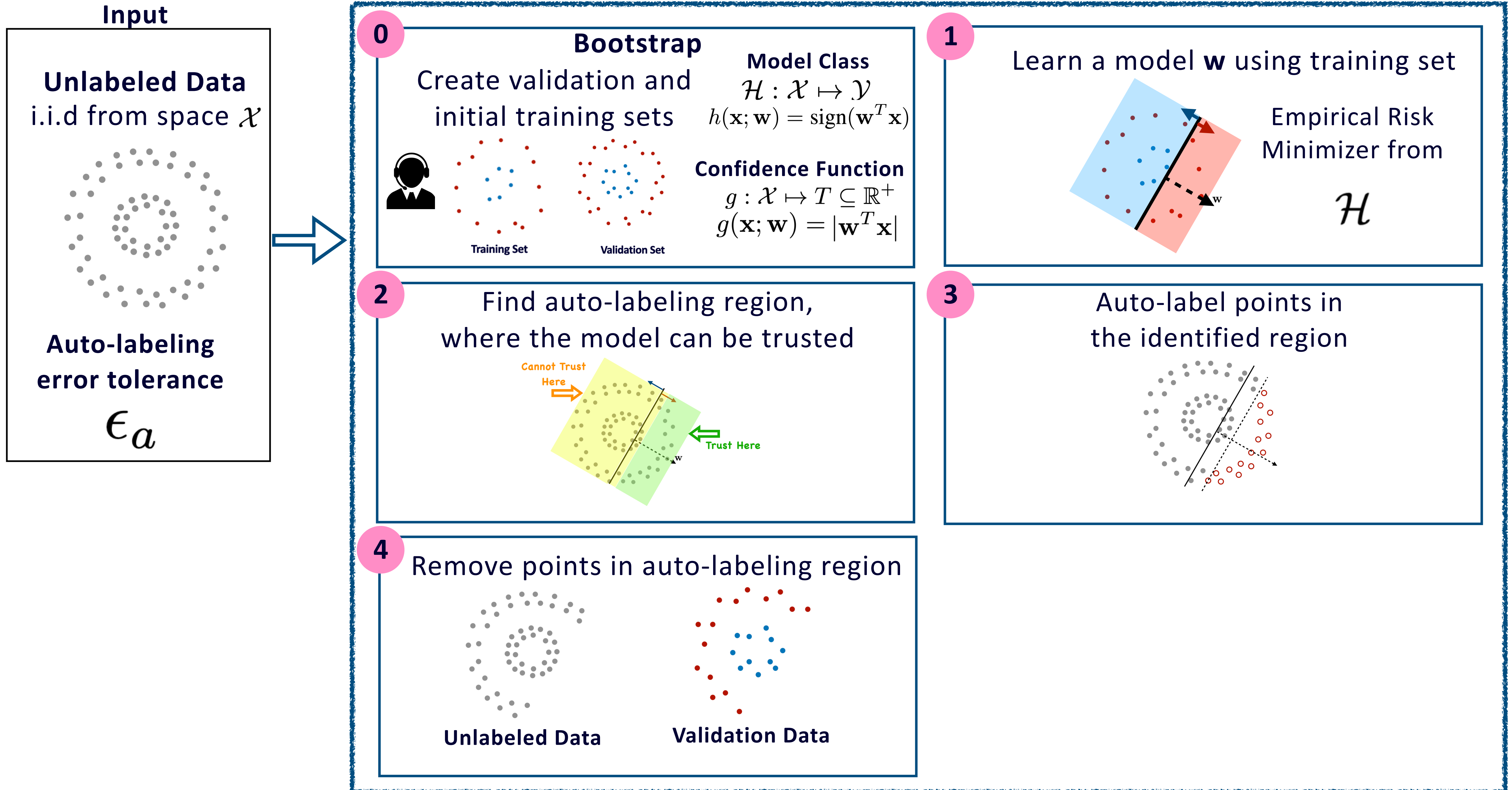
Trust Here



Auto-labeled

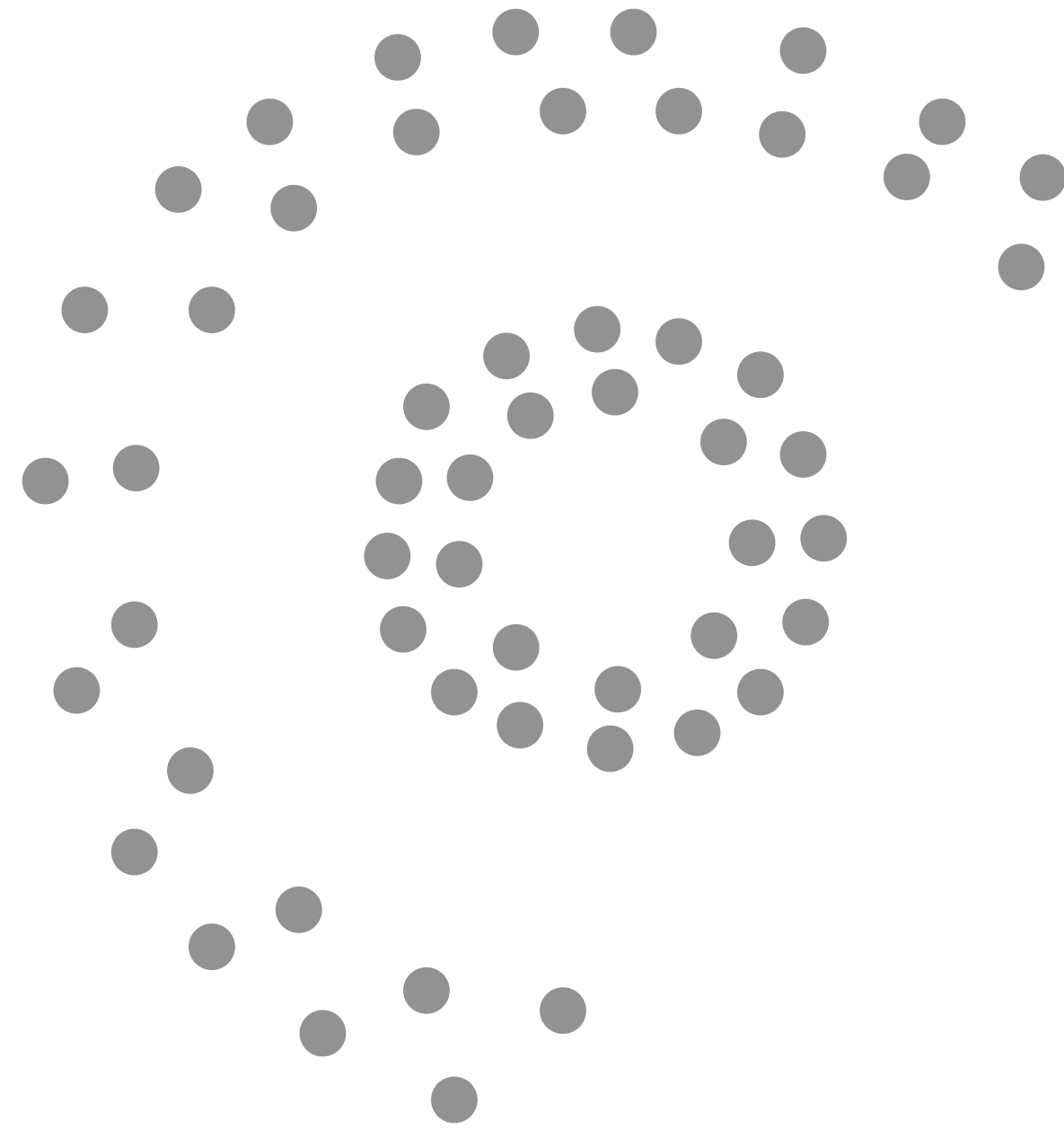


Threshold-based Auto-labeling Workflow(TBAL)



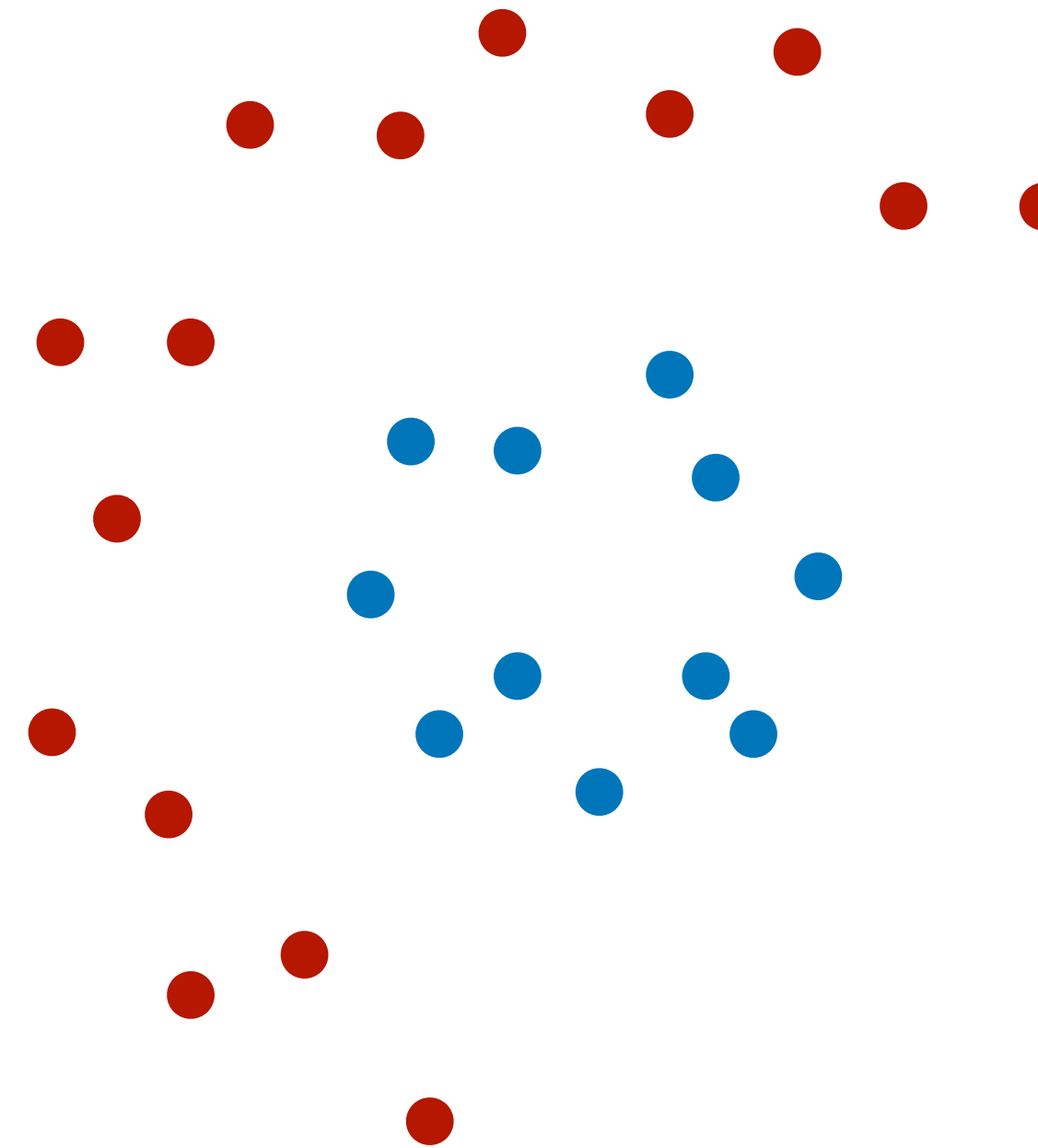
TBAL Workflow: Step 4 Prepare for the next round

Remove auto-labeled points from the pool.



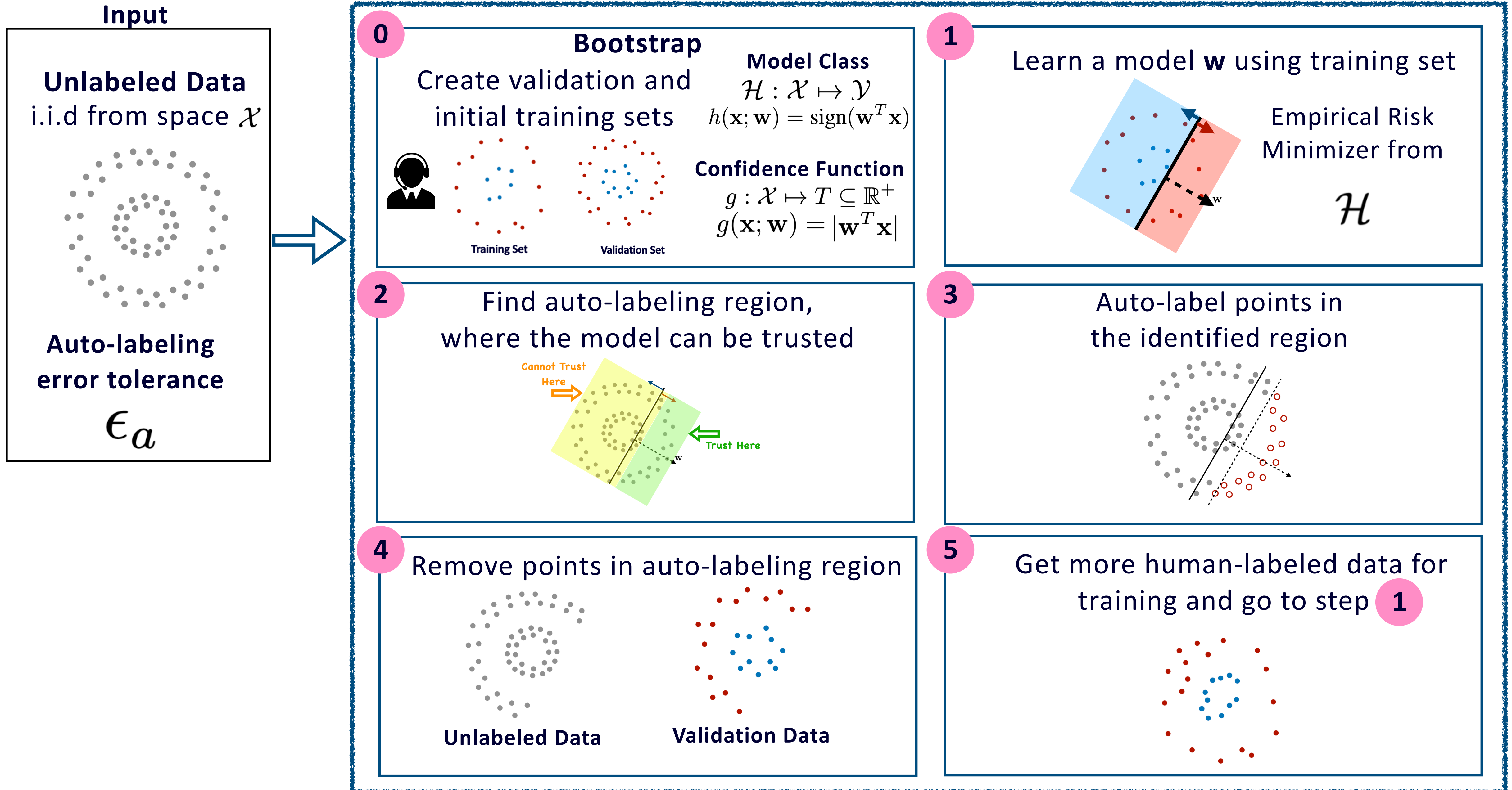
Remaining unlabeled data

Remove points from the validation set falling in the auto-labeling region.



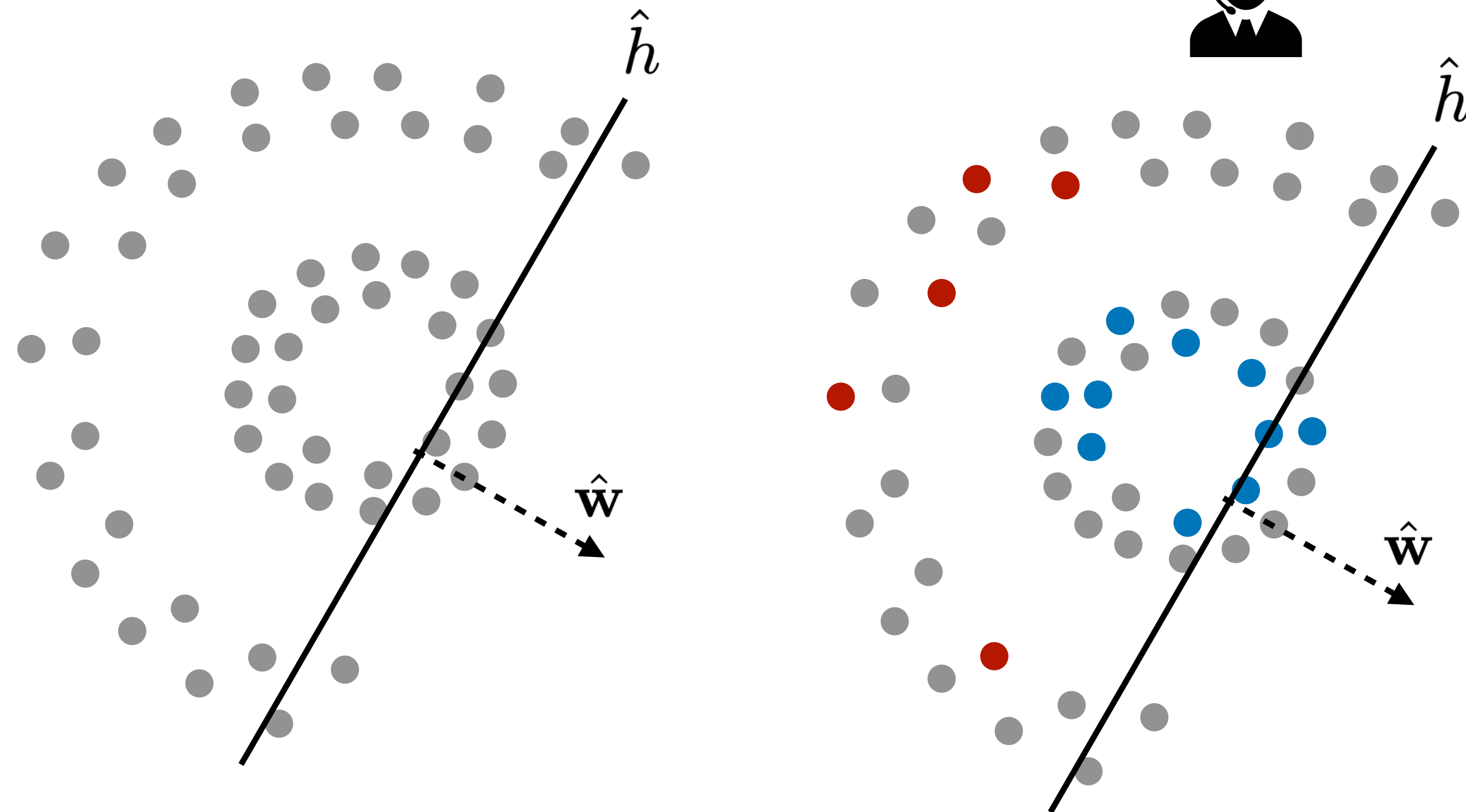
Remaining validation data

Threshold-based Auto-labeling Workflow(TBAL)

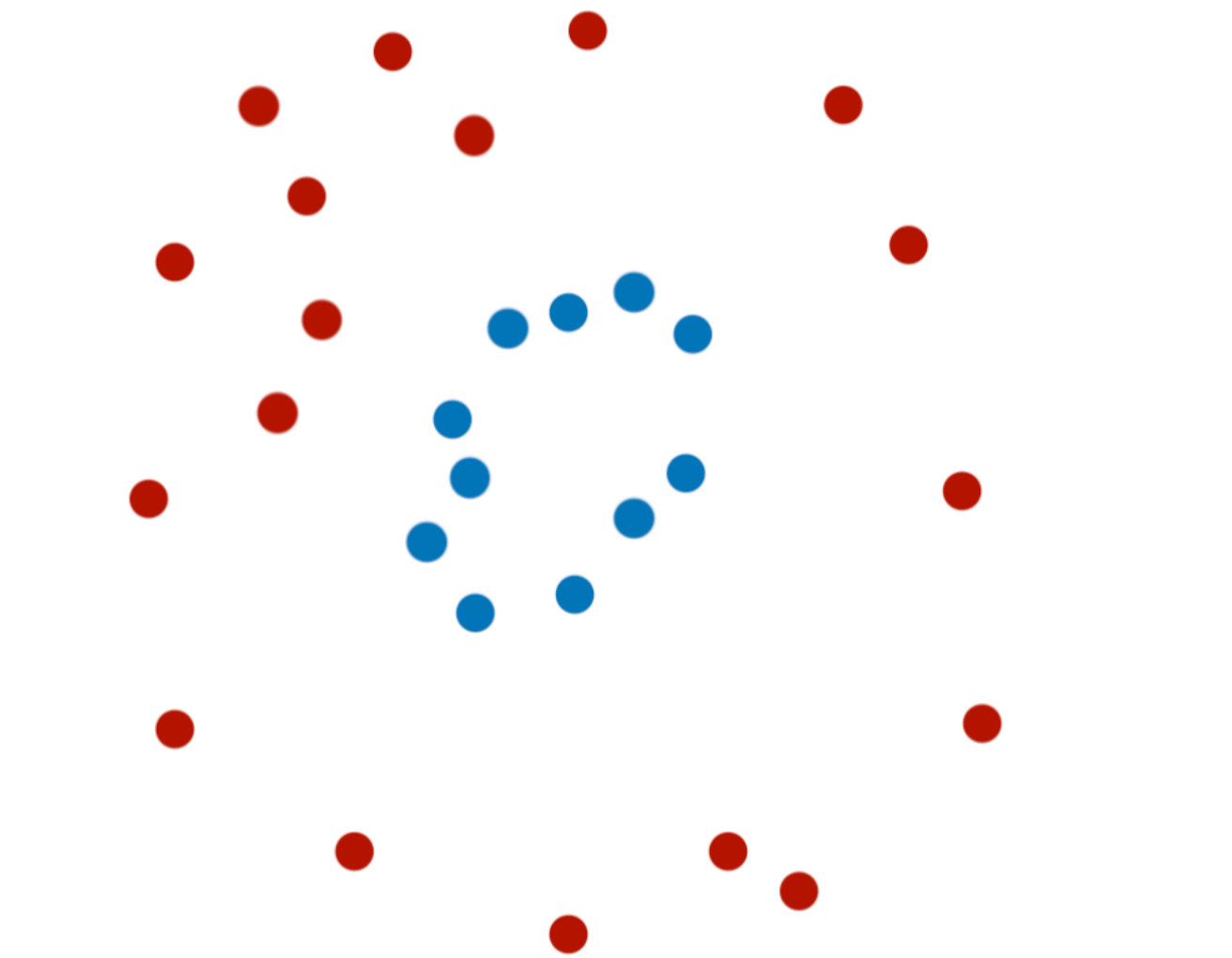


Step 5: Query next batch of human-labeled data for training

Use some active querying strategy
example: uncertainty sampling



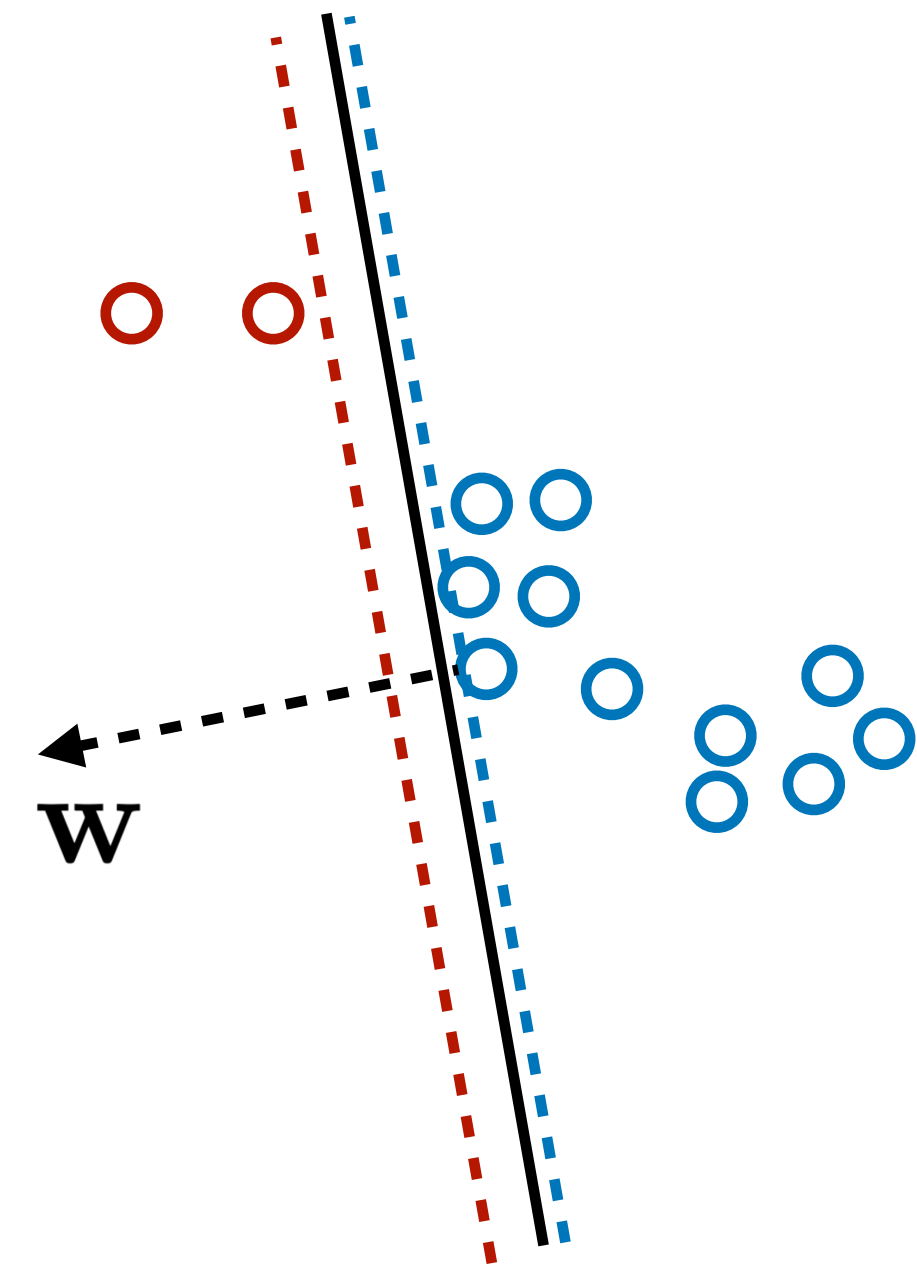
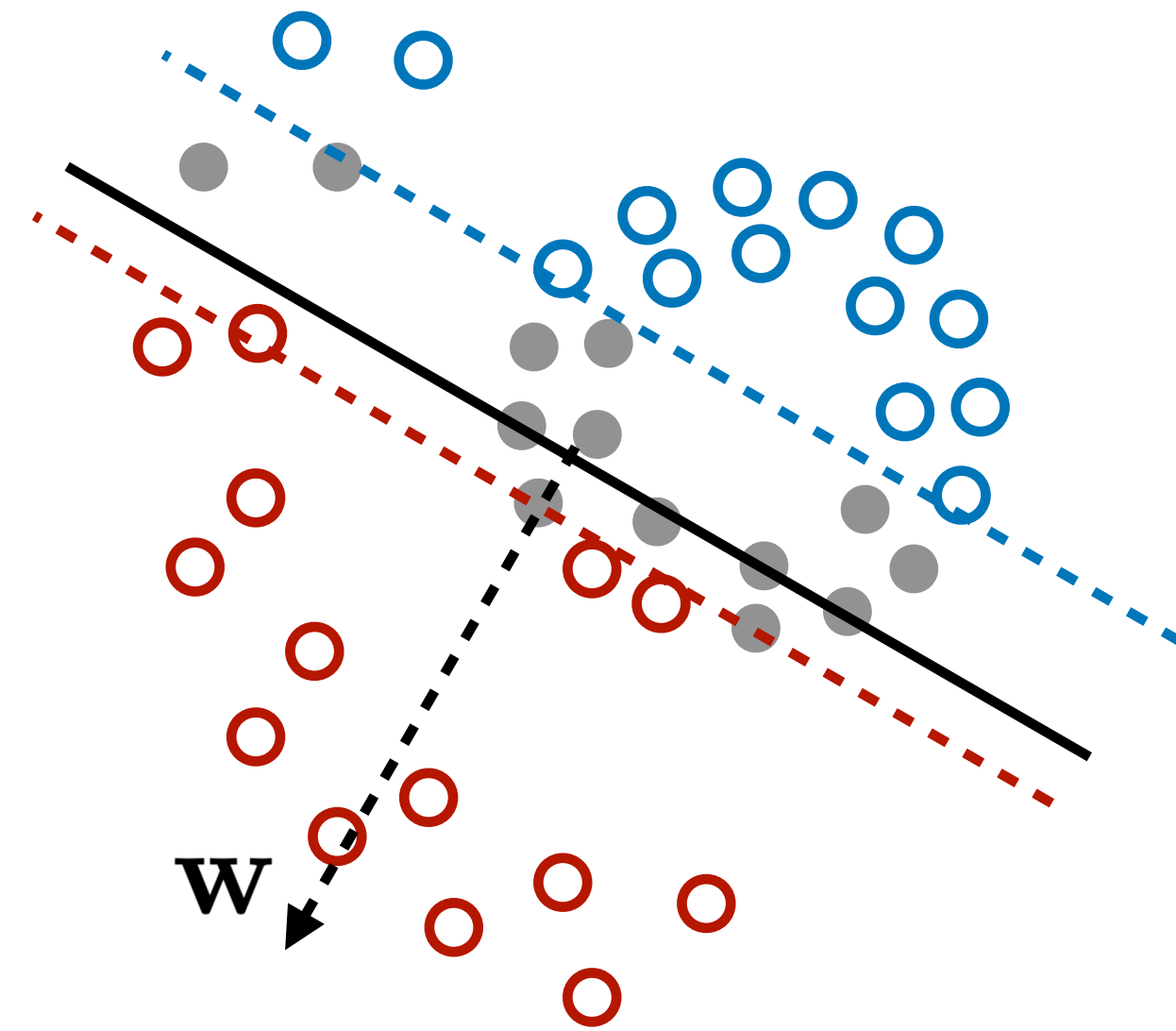
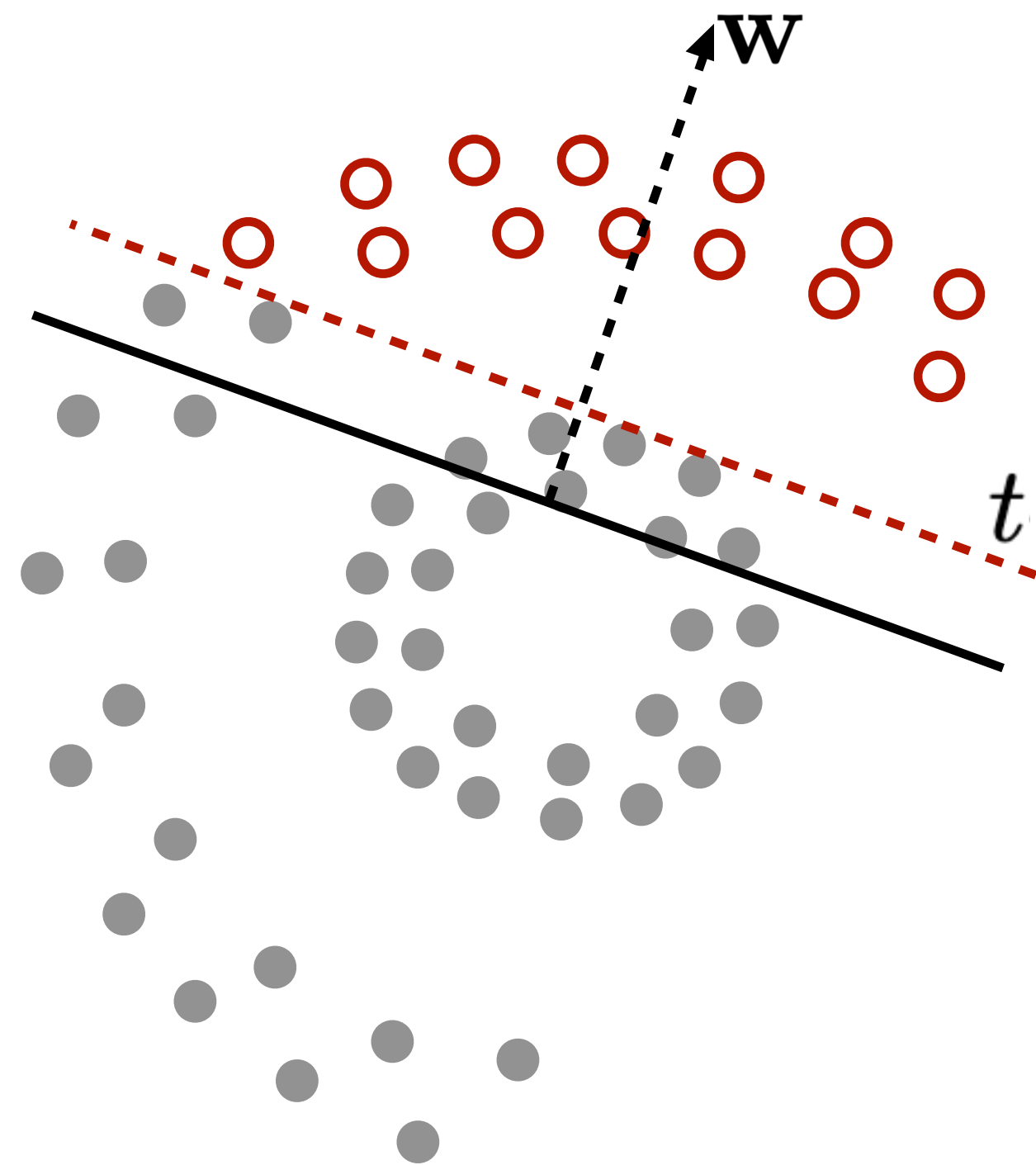
Next round's training data



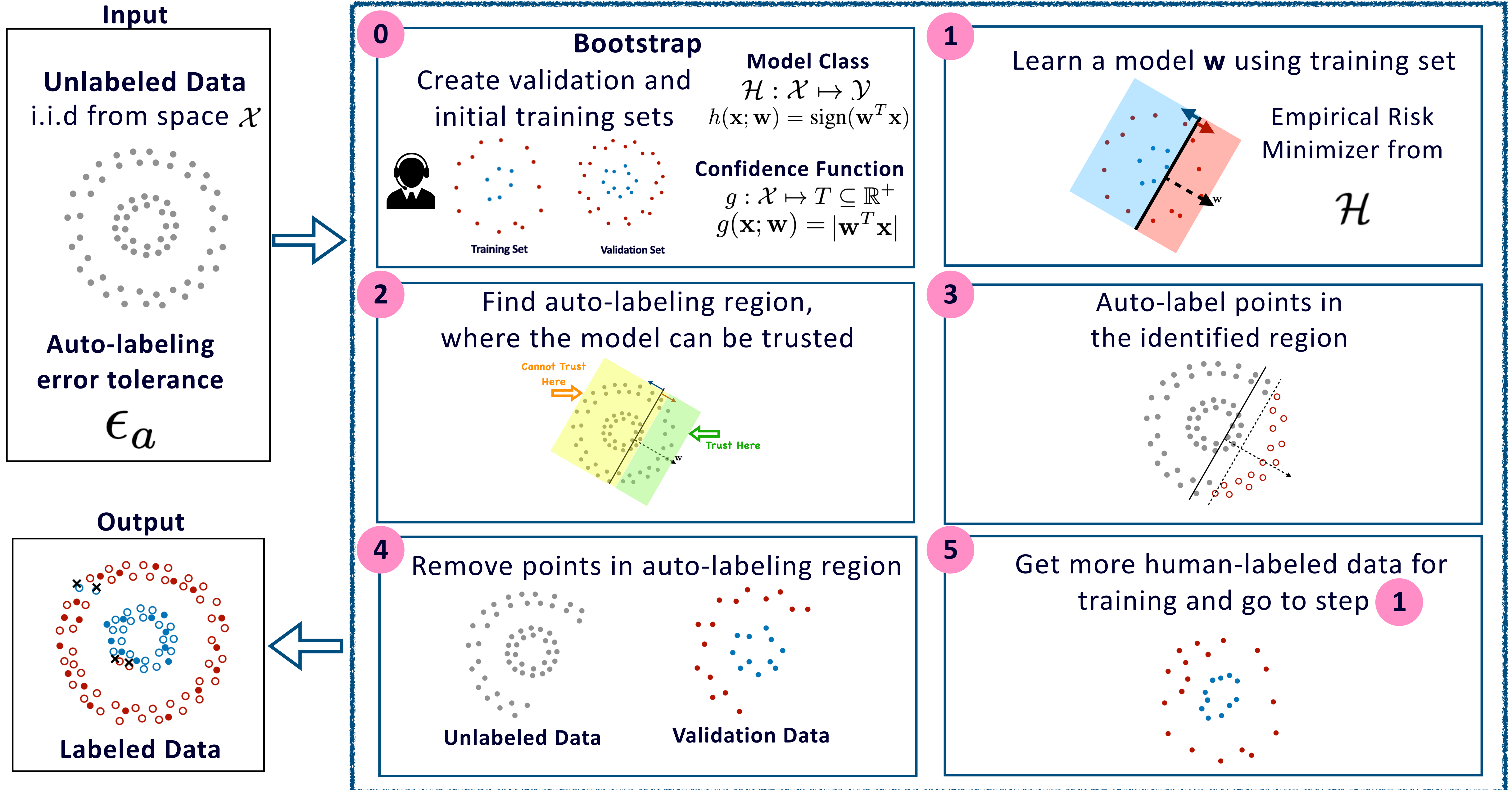
If there is unlabeled data left

Go to Step 1

Intermediate Rounds Output



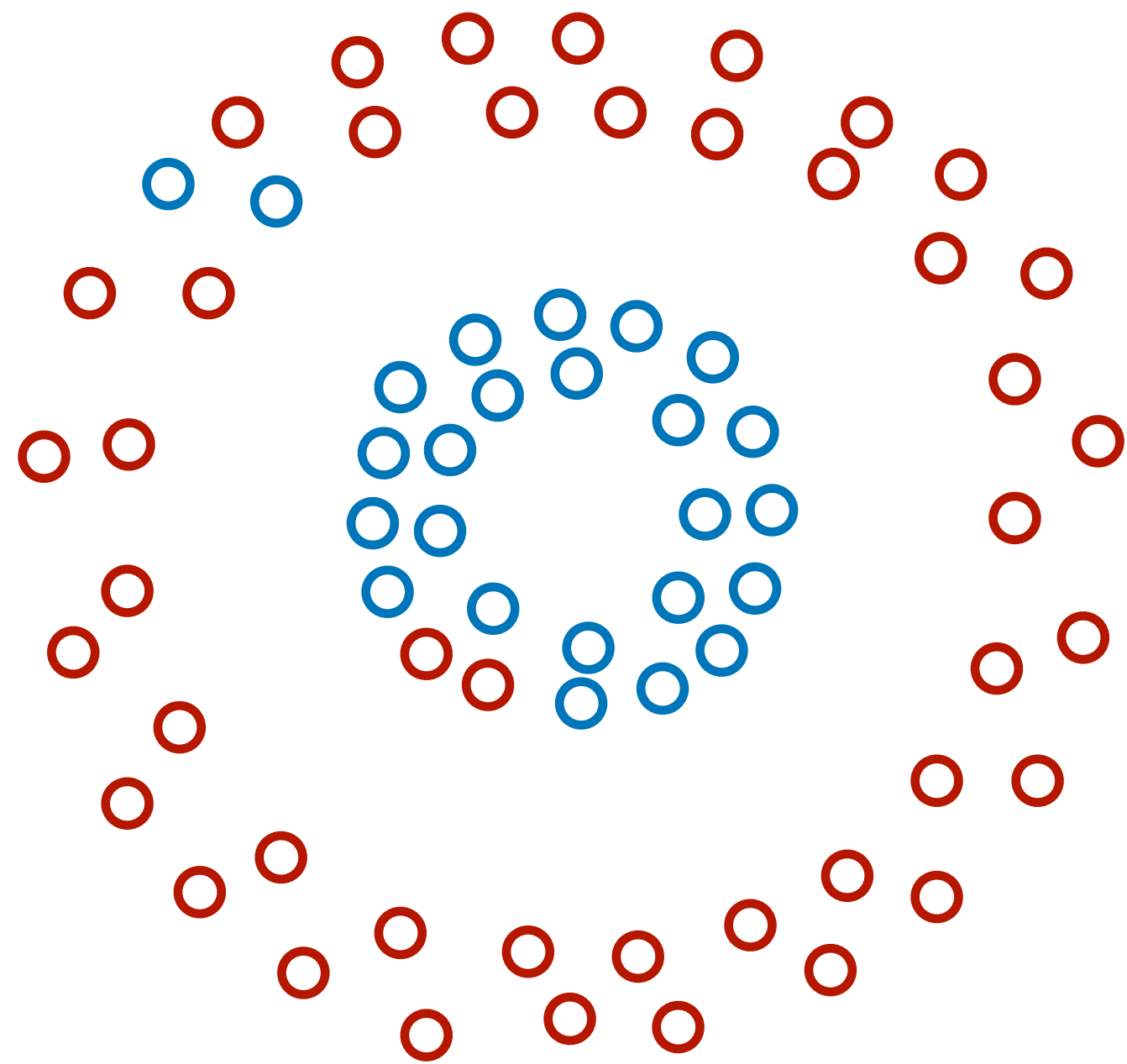
Threshold-based Auto-labeling Workflow(TBAL)



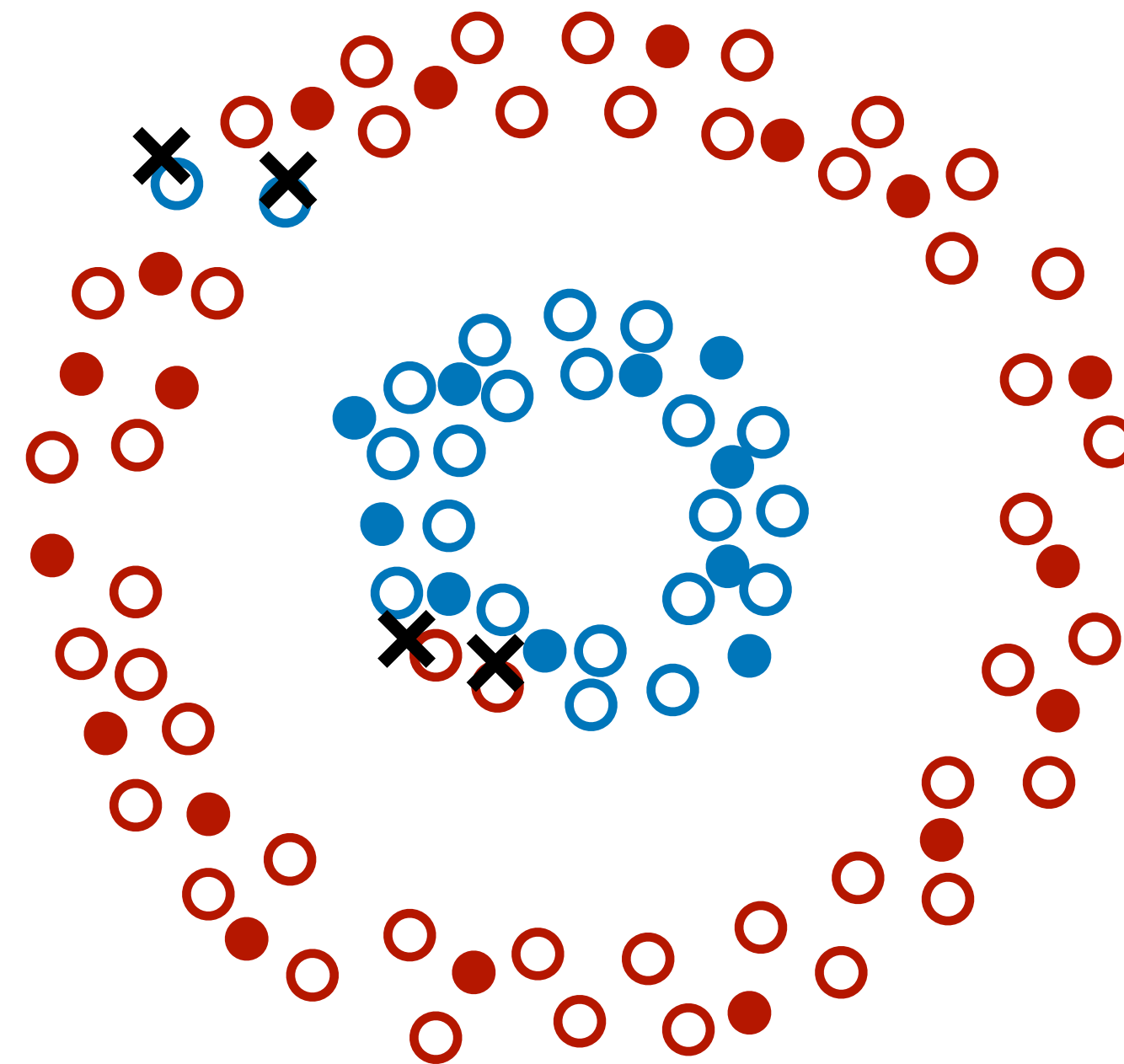
Final Output

- ● Human-labeled
- ○ Auto-labeled
- ✕ Labeling mistake

Auto-labeled data in the end



Output Labeled Dataset



Error and Coverage

Auto-labeling Error < 1%

Coverage > 95%

Roadmap

What & Why auto-labeling?

Data labeling problem

Adoption of auto-labeling

How does it work?

Workflow of TBAL

Finding the
auto-labeling region

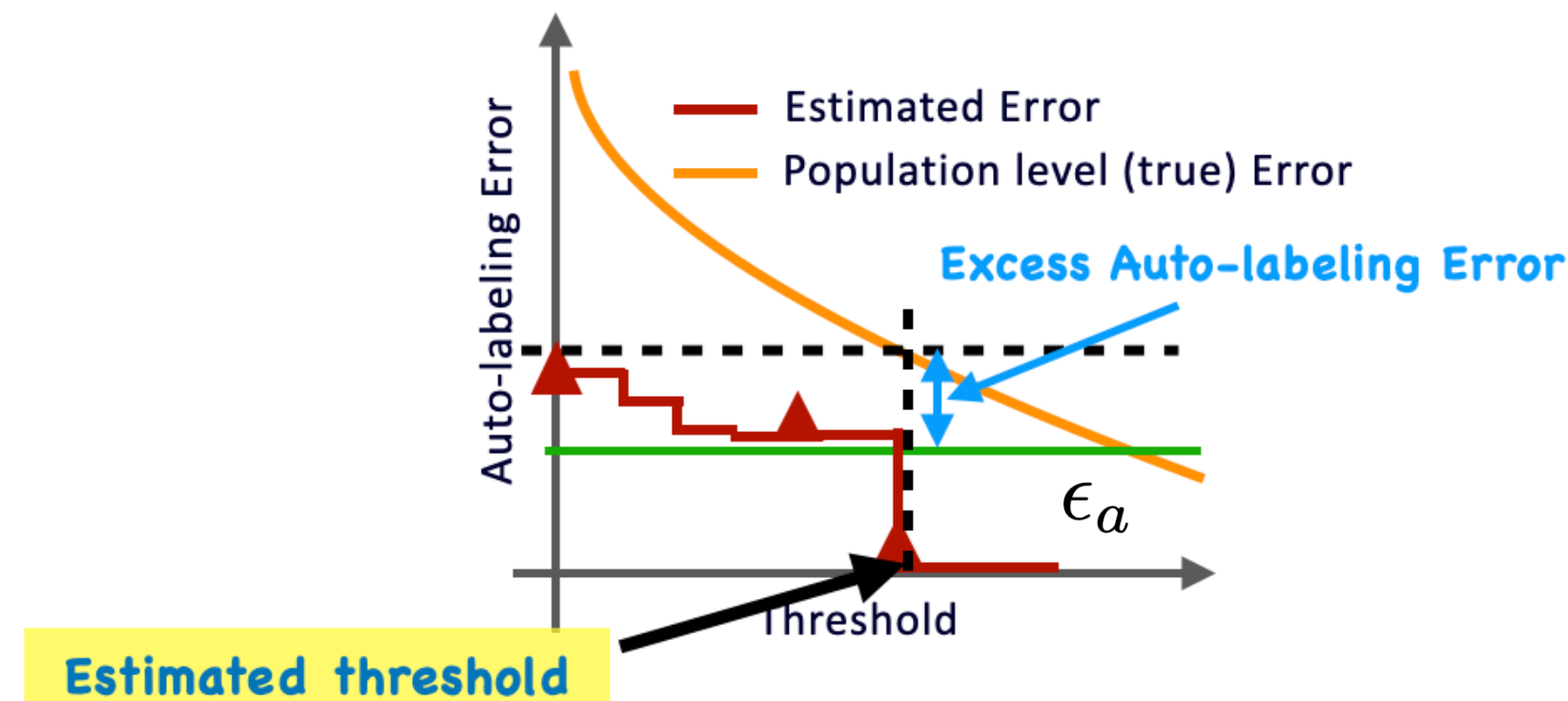
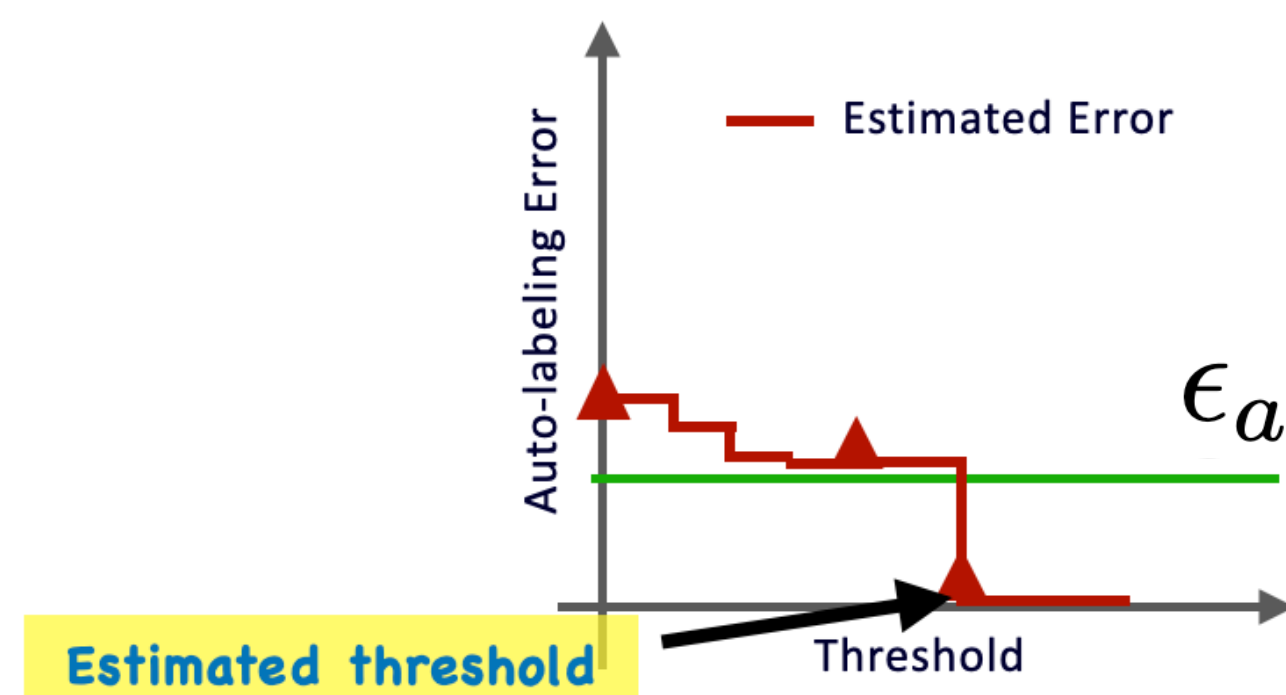
Analysis & Results

Conditions when TBAL works.

Comparison with
Active Learning, Selective
Classification

Theoretical Results

Conditions on the **validation data** for accurate auto-labeling



In the general setup: **No assumptions on data distribution and function classes**

Upper bound on excess auto-labeling error

$$\mathcal{O} \left(\frac{1}{\sqrt{N_v}} + \mathfrak{R}_{N_v}(\mathcal{H}^{T,g}) \right)$$

N_v
Validation points

$$\mathcal{H}^{T,g} := \mathcal{H} \times \mathcal{T}, (h, t) \in \mathcal{H}^{T,g}$$

$$(h, t)(\mathbf{x}) := \begin{cases} h(\mathbf{x}) & \text{if } g(h, \mathbf{x}) \geq t \\ \text{abstain} & \text{o.w.} \end{cases}$$

Lower bound on number of validation samples to ensure auto-labeling error is below ϵ_a

$$\Omega \left(\frac{1}{\epsilon_a^2} \right)$$

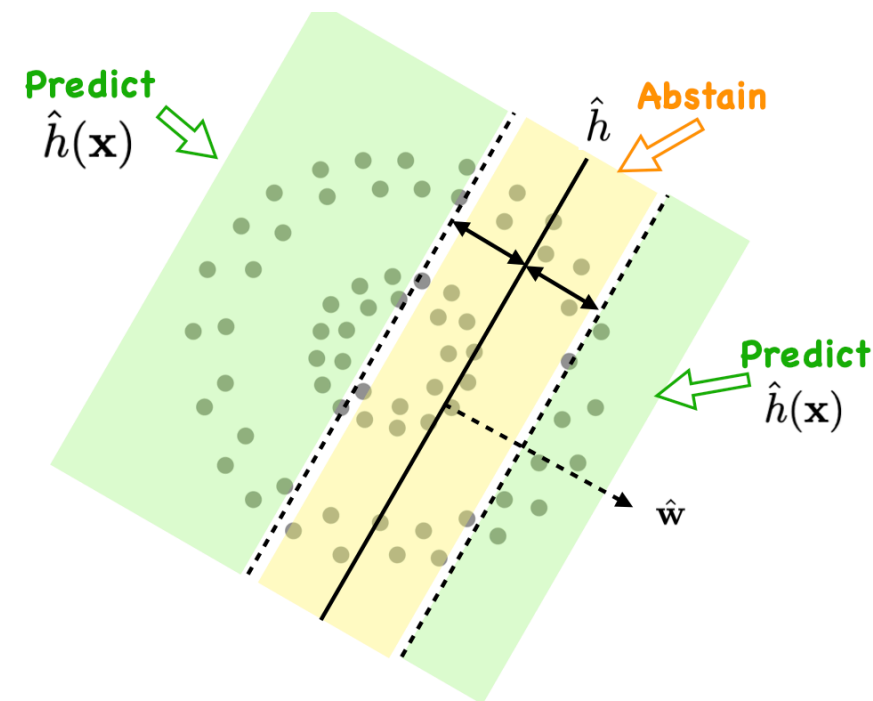
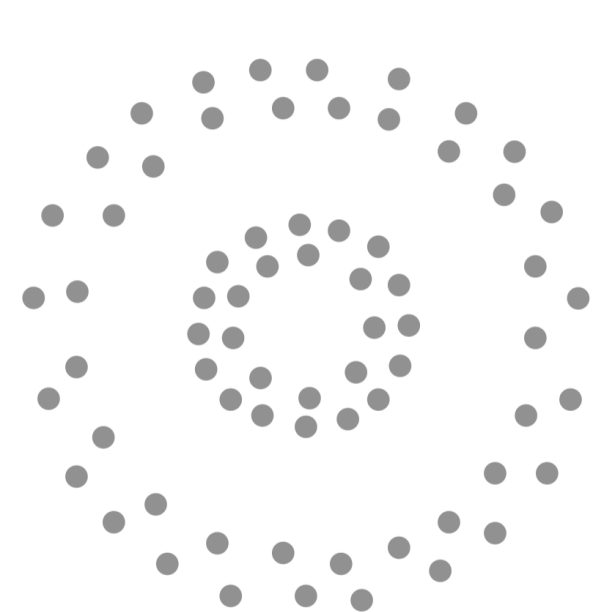
Instantiate the upper bound for uniform distribution on unit-ball in \mathbb{R}^d with homogeneous linear separators

Proof Sketch

With Finite Samples

$$A_v(h, t) = \{\mathbf{x} \in X_v : g(\mathbf{x}; h) \geq t\}$$

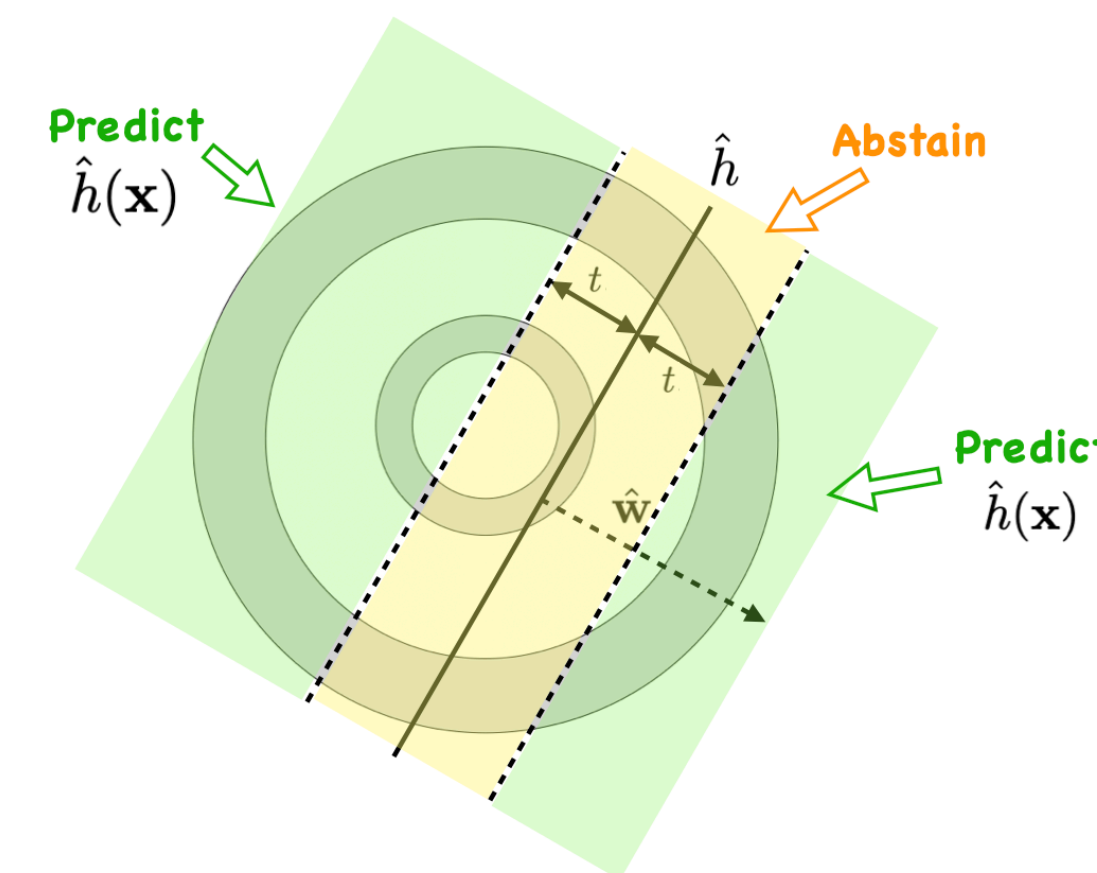
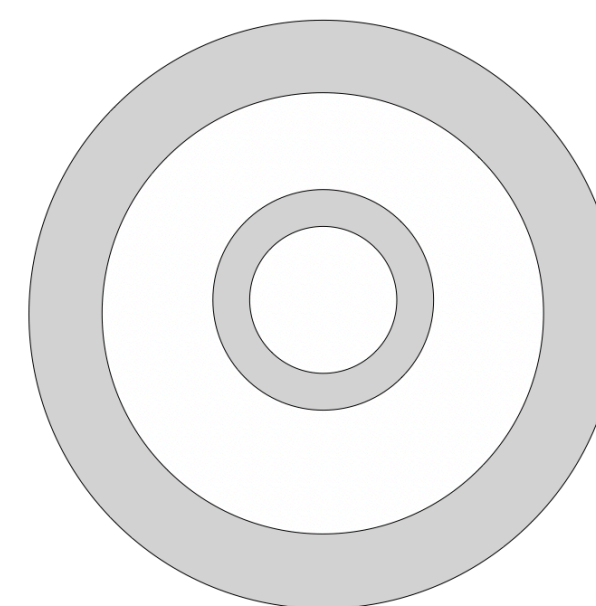
$$\hat{\mathcal{E}}_v(h|t) = \frac{1}{|A_v(h, t)|} \sum_{\mathbf{x} \in A_v(h, t)} \mathbb{1}\{h(\mathbf{x}) \neq f^*(\mathbf{x})\}$$



Population Level

$$\mathcal{A}(h, t) = \{\mathbf{x} \in \mathcal{X} : g(\mathbf{x}; h) \geq t\}$$

$$\mathcal{E}(h|t) = \mathbb{E}_{x|\mathcal{A}(h,t)}[\mathbb{1}\{h(\mathbf{x}) \neq f^*(\mathbf{x})\}]$$



Want this

$$\text{w.p. } 1 - \delta$$

$$\mathcal{E}(h|t) \leq \hat{\mathcal{E}}_v(h|t) + \psi(N_v, \delta, \mathcal{H}, g, T) \quad \forall h \in \mathcal{H}, \forall t \in T$$

Experiments

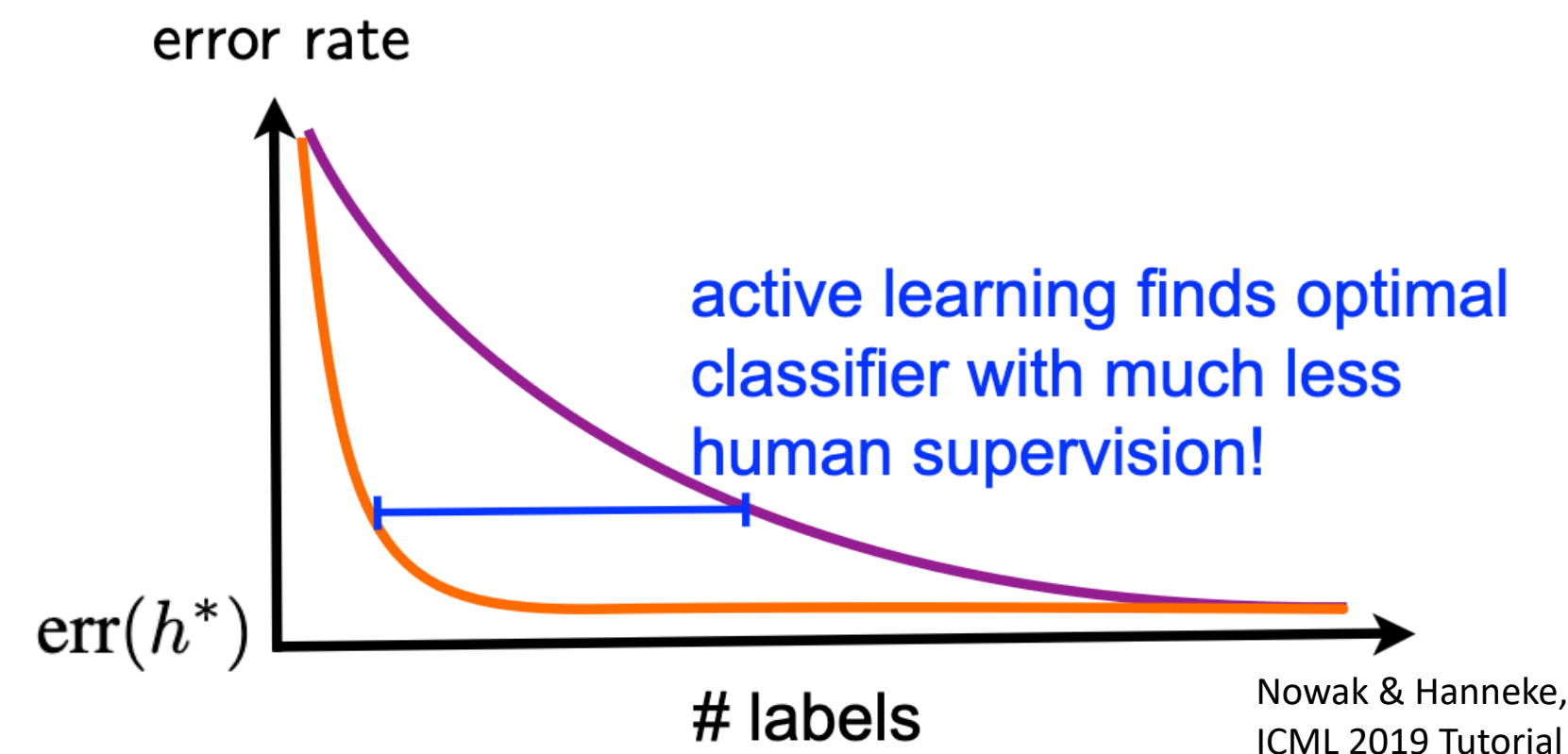
Active Learning and Selective Classification

Active Learning (AL)

$$\text{err}(h) = \mathbb{E}_{\mathbf{x}}[\mathbb{1}\{h(\mathbf{x}) \neq y\}]$$

$$h^* \in \arg \min_{h \in \mathcal{H}} \mathbb{E}_{\mathbf{x}}[\mathbb{1}\{h(\mathbf{x}) \neq y\}]$$

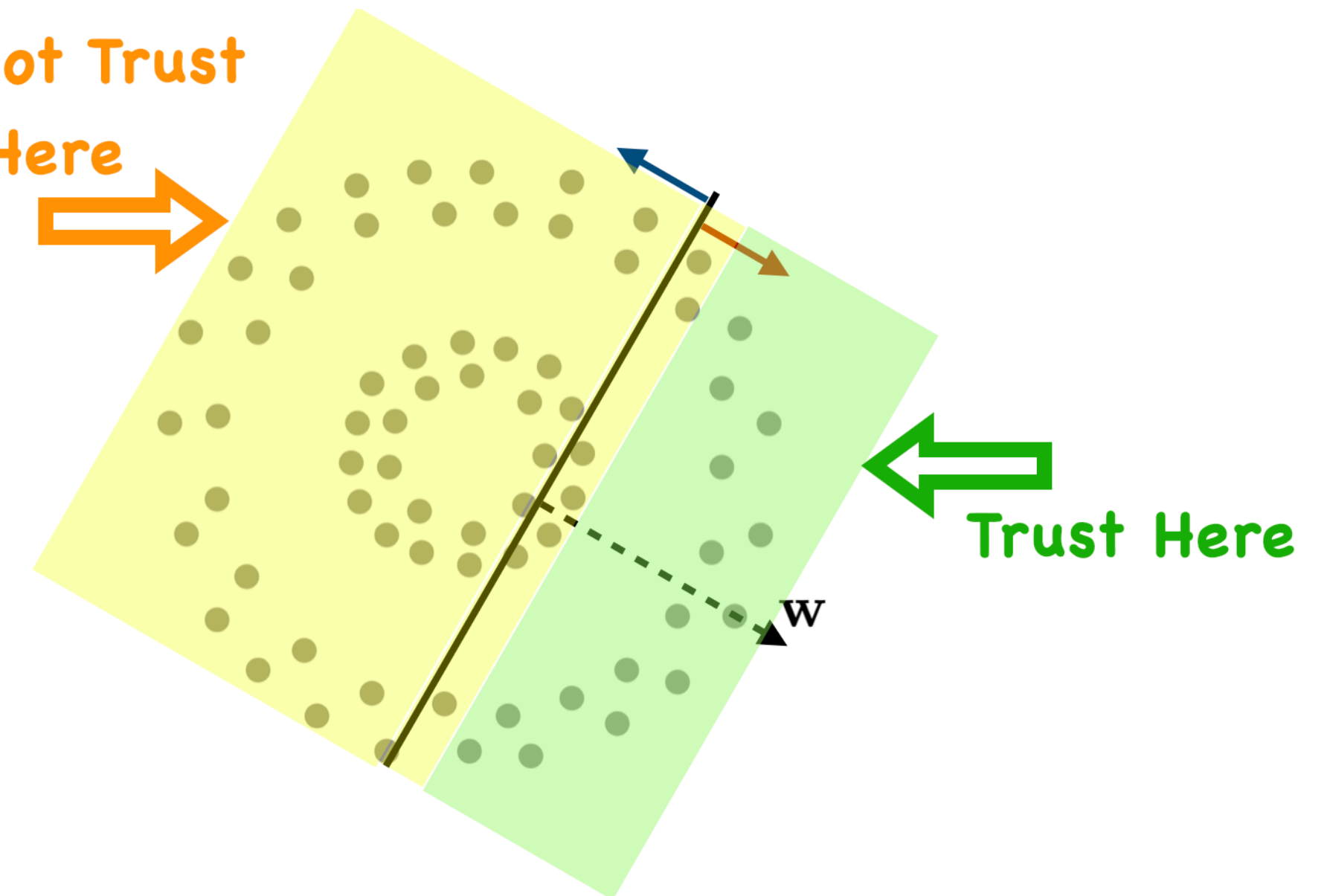
$$\text{err}(\hat{h}) - \text{err}(h^*) \rightarrow 0$$



Cohn et al. 1994;
Balcan, Dasgupta, Nowak, Zhu, Hanneke, Jamieson,
Chaudhury.... (Over the last 3 decades)

Selective Classification (SC)

Cannot Trust
Here

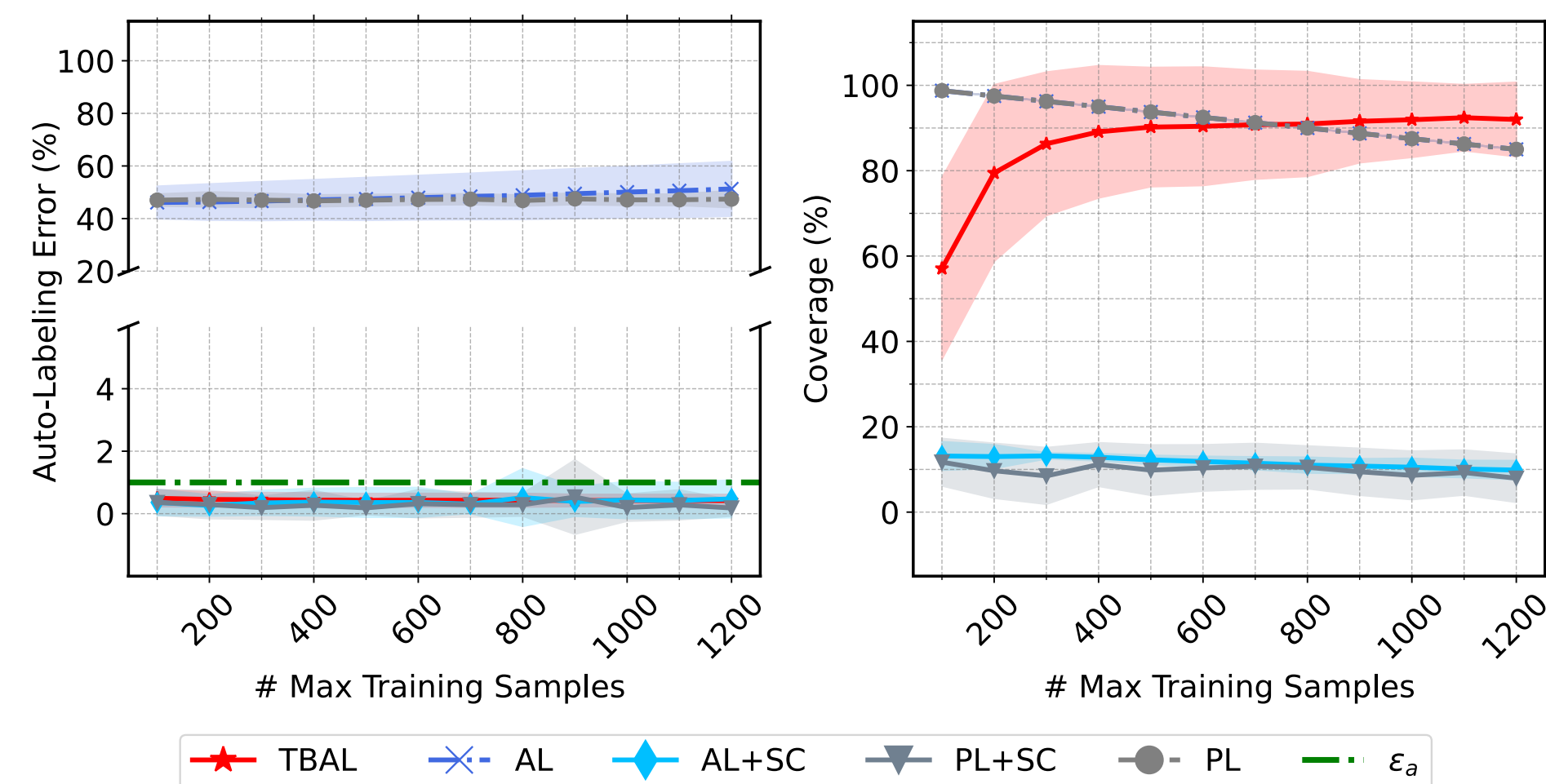
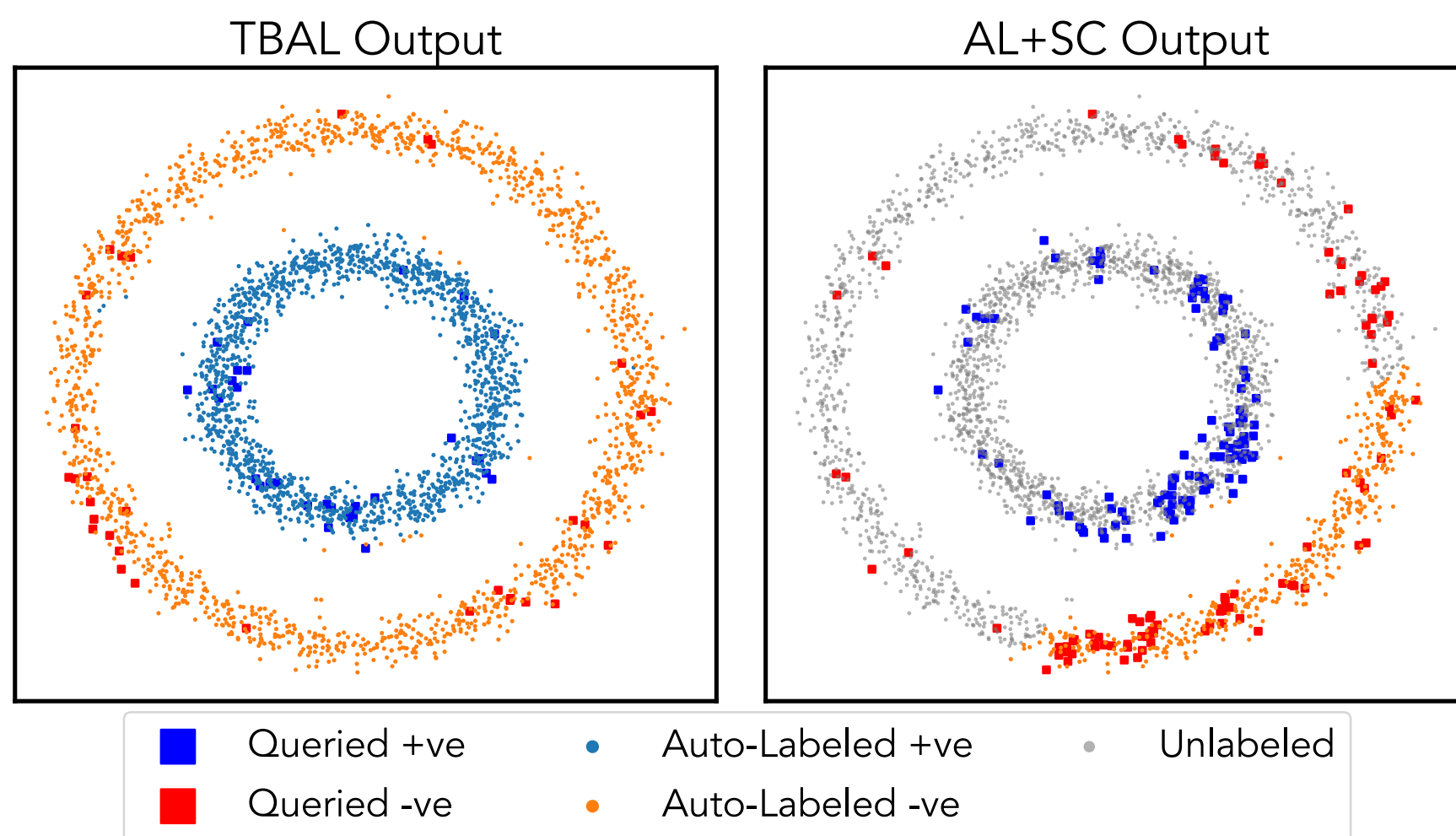
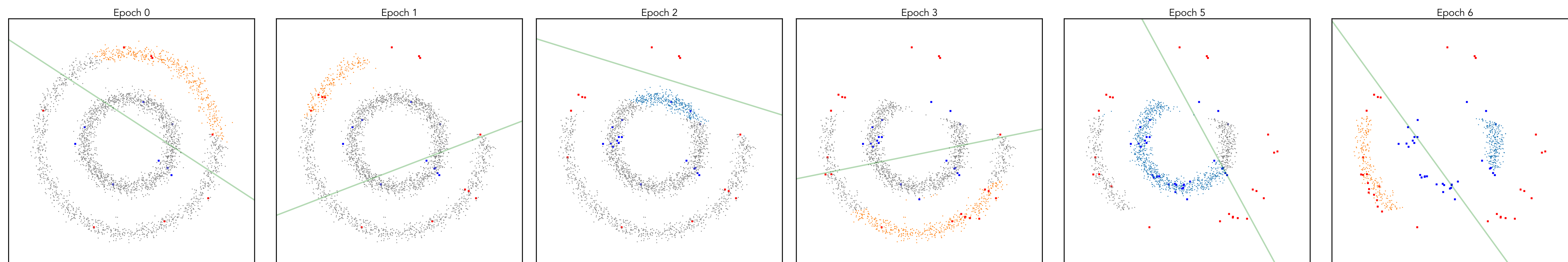


El-Yaniv & Weiner, 2010; Cortes, Desalvo, Mohri 2016;
Gelbhart & El-Yaniv 2019; Fisch, Jakkola et al. 2022;

**A natural auto-labeling strategy (AL+SC):
First learn the best classifier using Active Learning,
then auto-label using selective classification.**

The methods work as expected on the circles example

Misspecified setting: Using incorrect model class, (in practice the correct class is not known)



We validate the results empirically

Fix the auto-labeling error tolerance and the max number of training points algorithm can use.

Vary the number of validation points

Unit ball (Synthetic)

Increasing Validation data ↓

N_v	Error (%)		Coverage (%)	
	TBAL	AL+SC	TBAL	AL+SC
100	3.10 ±1.80	0.68 ±0.81	71.43 ±8.86	96.95 ±1.01
400	1.65 ±0.65	0.32 ±0.15	93.27 ±2.50	96.91 ±0.99
800	1.08 ±0.47	0.24 ±0.16	96.01 ±1.16	96.31 ±1.36
1200	0.78 ±0.27	0.17 ±0.11	96.82 ±0.84	95.96 ±1.40
1600	0.65 ±0.20	0.13 ±0.08	96.93 ±0.57	95.70 ±1.38
2000	0.54 ±0.16	0.21 ±0.11	97.23 ±0.42	96.36 ±1.13

Classes = 2 $\epsilon_a = 1\%$

Max # training points = 500

IMDB

N_v	Error (%)		Coverage (%)	
	TBAL	AL+SC	TBAL	AL+SC
200	2.28 ±0.21	3.11 ±0.86	68.24 ±6.20	57.77 ±13.09
400	1.29 ±0.10	1.98 ±0.40	63.81 ±4.86	63.06 ±10.70
600	1.41 ±0.20	1.81 ±0.22	69.64 ±3.98	62.92 ±9.20
800	1.62 ±0.30	2.04 ±0.35	67.45 ±3.72	63.22 ±7.89
1000	1.64 ±0.23	1.97 ±0.26	70.28 ±2.82	66.11 ±8.00

Classes = 2 $\epsilon_a = 5\%$

Max # training points = 500

Tiny Imagenet

N_v	Error (%)		Coverage (%)	
	TBAL	AL+SC	TBAL	AL+SC
2000	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0	0.0 ±0.0
4000	10.50 ±6.01	7.37 ±4.57	0.47 ±0.05	0.48 ±0.06
6000	10.61 ±0.62	7.71 ±1.03	10.16 ±1.10	4.31 ±1.10
8000	9.90 ±0.63	6.80 ±0.77	25.84 ±1.57	14.43 ±2.01
10000	8.97 ±0.36	6.87 ±0.48	32.19 ±1.34	21.96 ±1.35

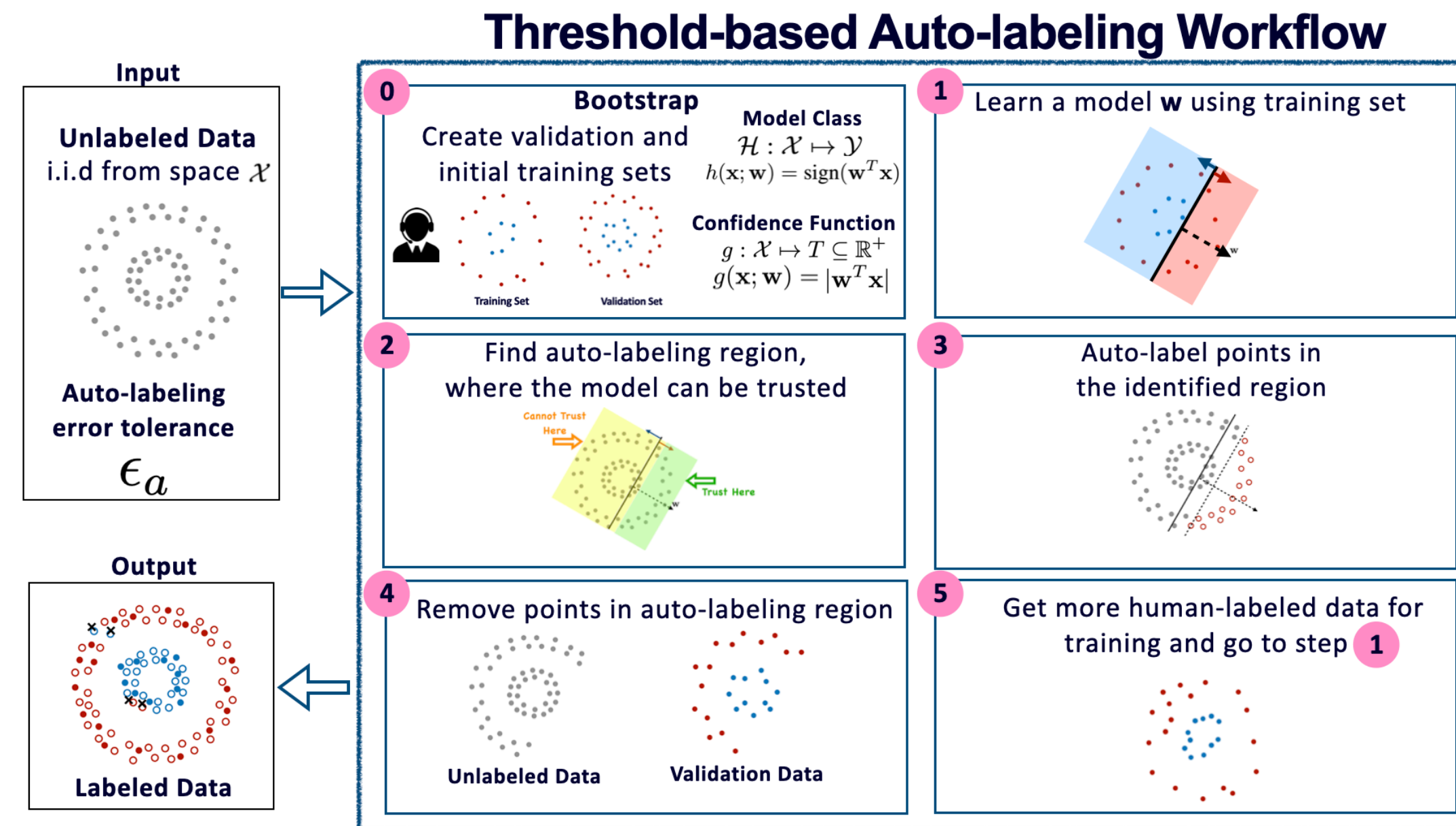
Classes = 200 $\epsilon_a = 10\%$

Max # training points = 10000

As expected, we observe

Less validation data \Rightarrow high auto-labeling errors and high variance in coverage
 Suff. Large validation data \Rightarrow less auto-labeling errors and less variance in coverage

Summary and Takeaways



1. Auto labeling is a promising solution to obtain labeled data.

2. Our work develops a theoretical understanding of auto-labeling systems.

3. **The promise** — Seemingly bad models can auto-label significant portion of data with good accuracy.

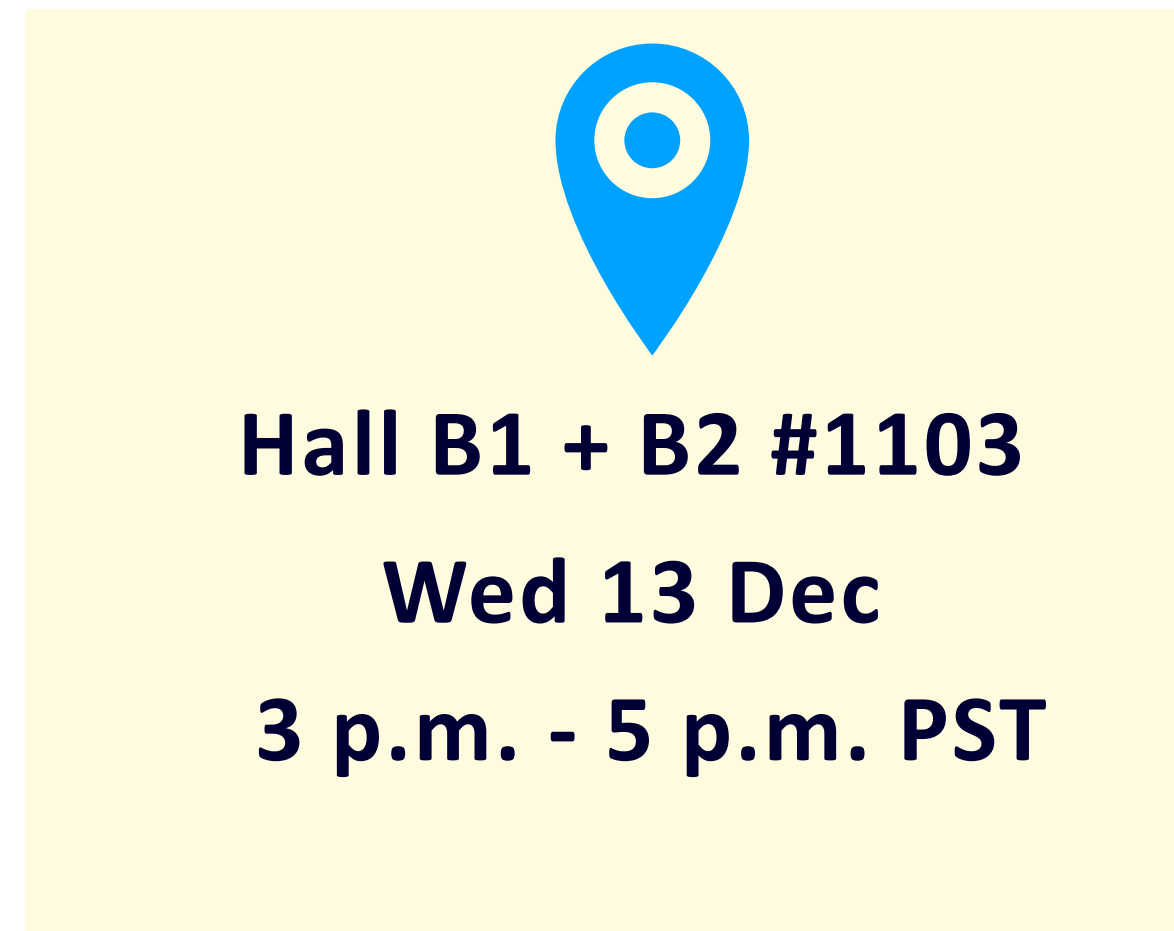
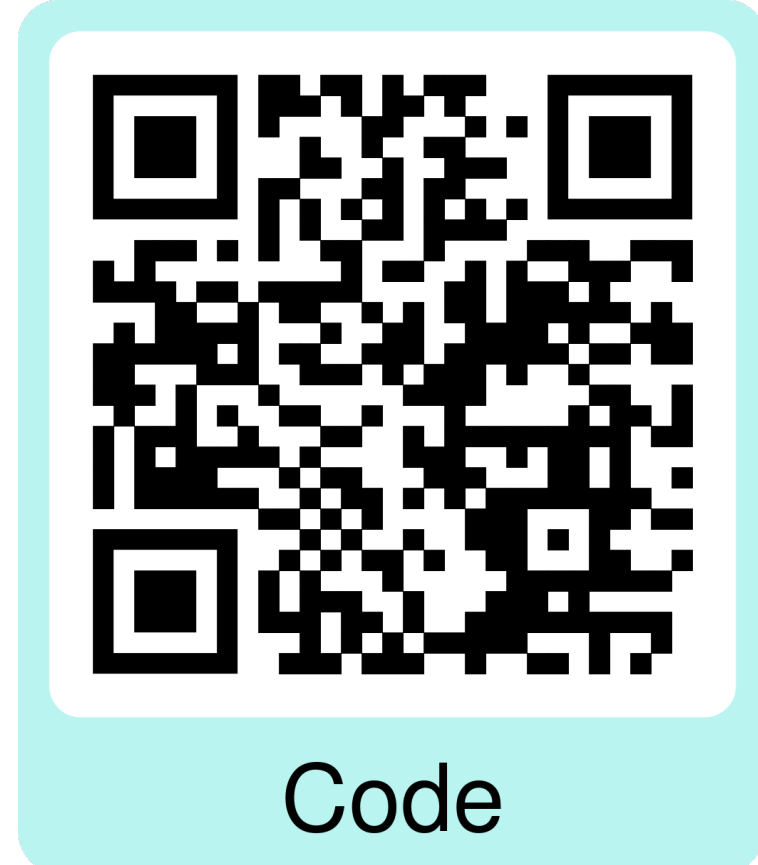
4. **The pitfall** — Hidden downside is it may need large amount validation data to ensure good accuracy.

Thank You

Checkout our paper and code!

Come to our poster @ NeurIPS

Contact us



Harit Vishwakarma
hvishwakarma@cs.wisc.edu



Huguang Lin
hglin@seas.upenn.edu



Frederic Sala
fredsala@cs.wisc.edu



Ramya Korlakai Vinayak
ramya@ece.wisc.edu

Paper <https://openreview.net/pdf?id=RUCFAKNDb2>

Code <https://github.com/harit7/TBAL-NeurIPS-23>