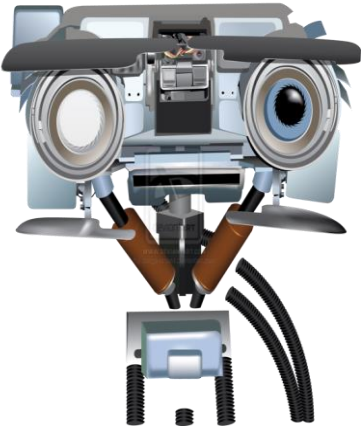


# Using Zoom for Lectures

- **Sign in using:**
  - your name
- **Please mute both:**
  - your video cameras for the entire lecture
  - your audio/mics unless asking or answering a question
- **Asking/answering a question, option 1:**
  - click on Participants
  - use the hand icon to raise your hand
  - I will call on you and ask you to unmute yourself
- **Asking/answering a question, option 2:**
  - click on Chat
  - type your question, and I will answer it

# Today: Outline

- **Semi-supervised Learning**
- **Reminders:** Class Challenge is posted, due Apr 24  
(3-week challenge)  
Midterm Exam, Apr 15 during class time  
(covering material up to and including Apr 3)  
Practice Problems are posted  
Trial Exam Submission in class today



# Semi-Supervised Learning

---

Slides credit: Jerry Zhu, Aarti Singh

# Supervised Learning

**Feature Space**  $\mathcal{X}$

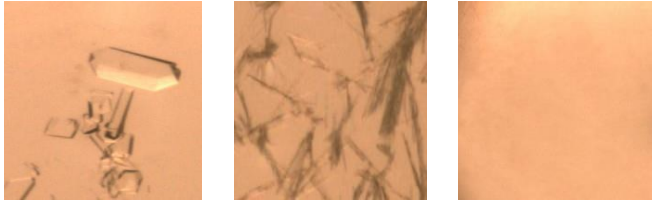
**Label Space**  $\mathcal{Y}$

**Goal:** Construct a **predictor**  $f : \mathcal{X} \rightarrow \mathcal{Y}$  to minimize



$\text{loss}(Y, f(X))$

# Labeled and Unlabeled data



0 1 2 3 4 5 6 7 8 9  
8 9 0 1 2 3 4 5 6 7



Unlabeled data,  $X_i$

**Cheap and abundant !**



Human expert/  
Special equipment/  
Experiment

“Crystal” “Needle” “Empty”

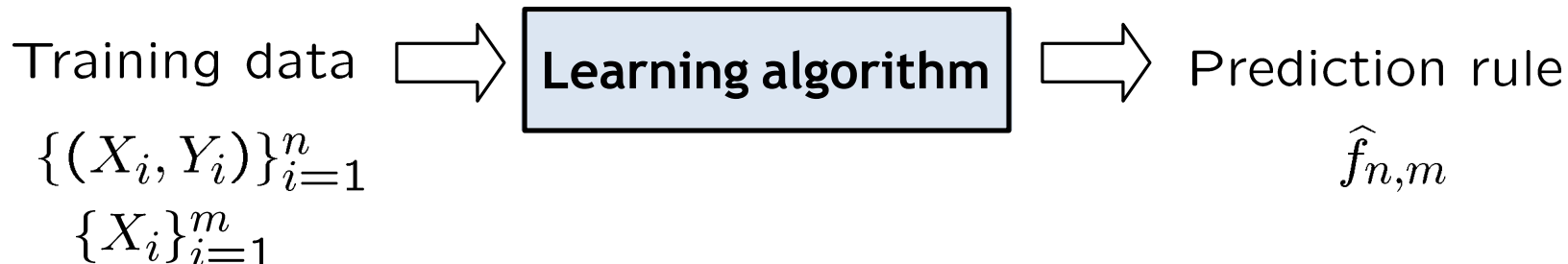
“0” “1” “2” ...

“Sports”  
“News”  
“Science”  
...

Labeled data,  $Y_i$

**Expensive and scarce !**

# Semi-Supervised learning



## Supervised learning (SL)

Labeled data  $\{X_i, Y_i\}_{i=1}^n$



$X_i$

“Crystal”

$Y_i$

## Semi-Supervised learning (SSL)

Labeled data  $\{X_i, Y_i\}_{i=1}^n$  **and** Unlabeled data  $\{X_i\}_{i=1}^m$

$m \gg n$

**Goal: Learn a better prediction rule than based on labeled data alone.**

# Semi-Supervised learning in Humans

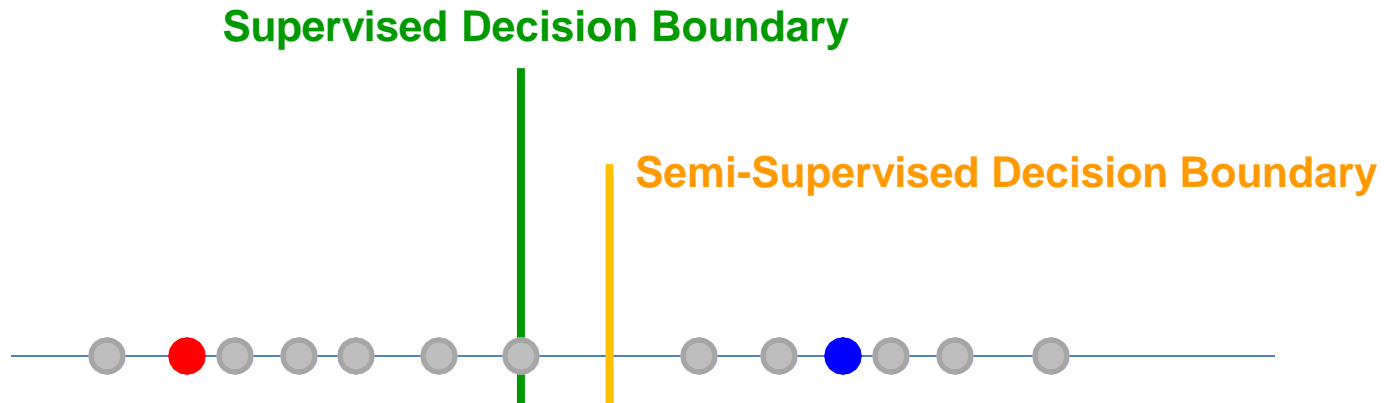
## Cognitive science

Computational model of how humans learn from labeled and unlabeled data.

- concept learning in children:  $x$ =animal,  $y$ =concept (e.g., dog)
- Daddy points to a brown animal and says “dog!”
- Children also observe animals by themselves

# Can unlabeled data help?

- Positive labeled data
- Negative labeled data
- Unlabeled data



Assume each class is a coherent group (e.g. Gaussian)

Then unlabeled data can help identify the boundary more accurately.

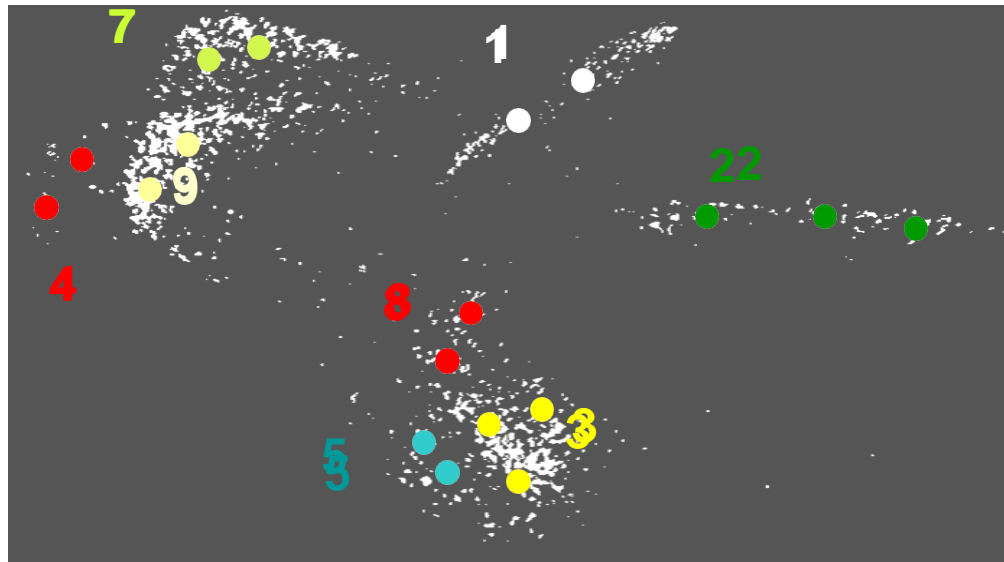


# Can unlabeled data help?

Unlabeled Images

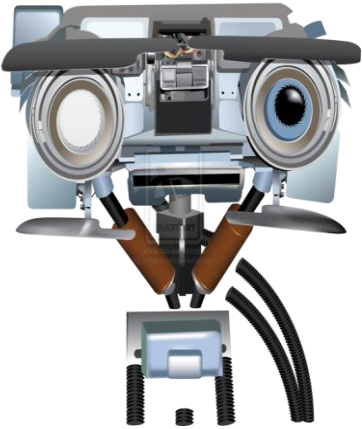
0 1 2 3 4 5 6 7 8 9  
8 9 0 1 2 3 4 5 6 7  
6 7 8 9 0 1 2 3 4 5

Labels “0” “1” “2” ...



This embedding can be done by manifold learning algorithms, e.g. t-SNE

**“Similar” data points have “similar” labels**



# Algorithms

---

## Semi-Supervised Learning

Slides credit: Jerry Zhu, Aarti Singh

# Some SSL Algorithms

- Self-Training
- Generative methods, mixture models
- Graph-based methods
- Co-Training
- Semi-supervised SVM
- Many others

# Notation

- instance  $\mathbf{x}$ , label  $y$
- learner  $f : \mathcal{X} \mapsto \mathcal{Y}$
- labeled data  $(X_l, Y_l) = \{(x_{1:l}, y_{1:l})\}$
- unlabeled data  $X_u = \{\mathbf{x}_{l+1:l+u}\}$ , **available** during training. Usually  $l \ll u$ . Let  $n = l + u$
- test data  $\{(x_{n+1...}, y_{n+1...})\}$ , **not available** during training

# Self-training

Our first SSL algorithm:

Input: labeled data  $\{(\mathbf{x}_i, y_i)\}_{i=1}^l$ , unlabeled data  $\{\mathbf{x}_j\}_{j=l+1}^{l+u}$ .

1. Initially, let  $L = \{(\mathbf{x}_i, y_i)\}_{i=1}^l$  and  $U = \{\mathbf{x}_j\}_{j=l+1}^{l+u}$ .
2. Repeat:
3.       Train  $f$  from  $L$  using supervised learning.
4.       Apply  $f$  to the unlabeled instances in  $U$ .
5.       Remove a subset  $S$  from  $U$ ; add  $\{(\mathbf{x}, f(\mathbf{x})) | \mathbf{x} \in S\}$  to  $L$ .

Self-training is a *wrapper* method

- the choice of learner for  $f$  in step 3 is left completely open
- good for many real world tasks like natural language processing
- but mistake by  $f$  can reinforce itself

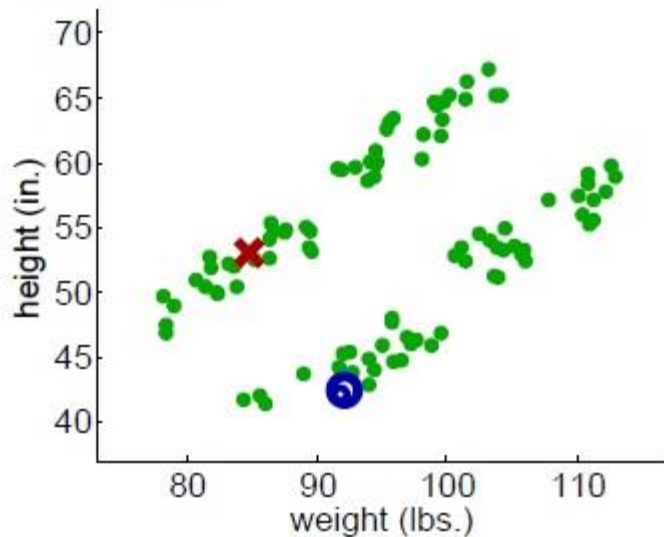
# Self-training Example

## Propagating 1-NN

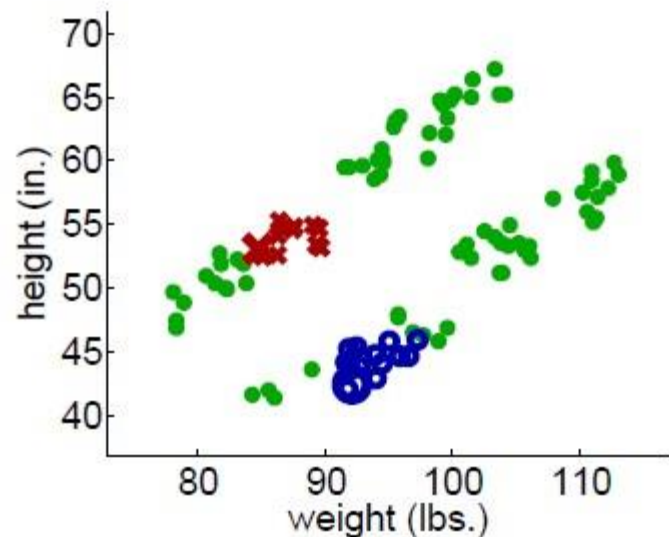
Input: labeled data  $\{(\mathbf{x}_i, y_i)\}_{i=1}^l$ , unlabeled data  $\{\mathbf{x}_j\}_{j=l+1}^{l+u}$ , distance function  $d()$ .

1. Initially, let  $L = \{(\mathbf{x}_i, y_i)\}_{i=1}^l$  and  $U = \{\mathbf{x}_j\}_{j=l+1}^{l+u}$ .
2. Repeat until  $U$  is empty:
3.     Select  $\mathbf{x} = \operatorname{argmin}_{\mathbf{x} \in U} \min_{\mathbf{x}' \in L} d(\mathbf{x}, \mathbf{x}')$ .
4.     Set  $f(\mathbf{x})$  to the label of  $\mathbf{x}$ 's nearest instance in  $L$ .  
      Break ties randomly.
5.     Remove  $\mathbf{x}$  from  $U$ ; add  $(\mathbf{x}, f(\mathbf{x}))$  to  $L$ .

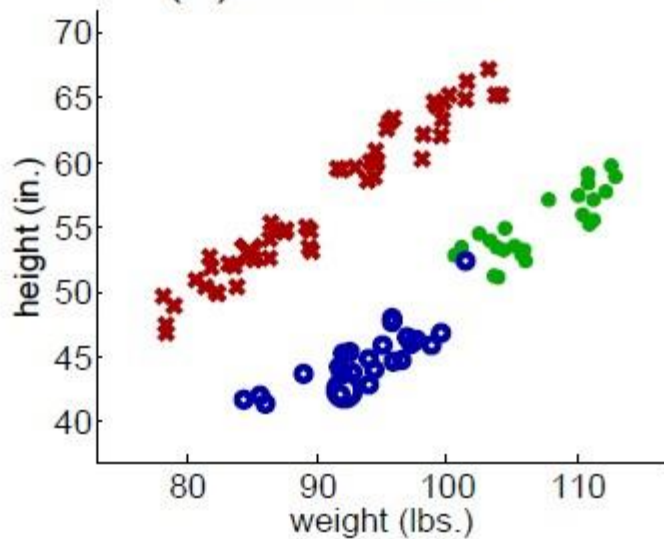
# Propagating 1-Nearest-Neighbor: now it works



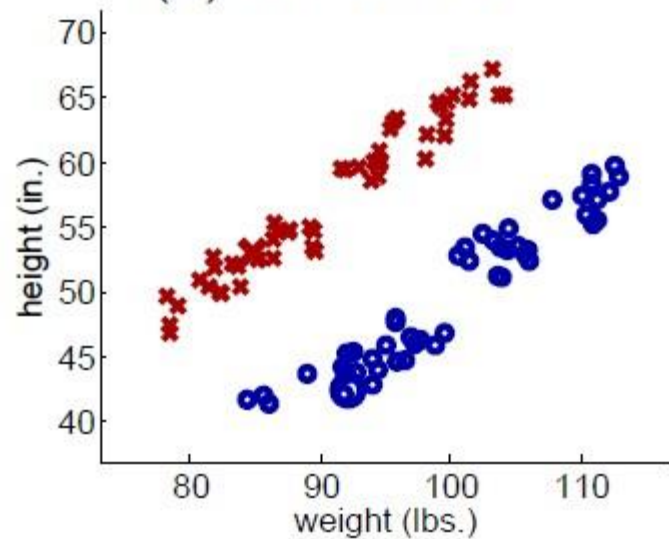
(a) Iteration 1



(b) Iteration 25



(c) Iteration 74

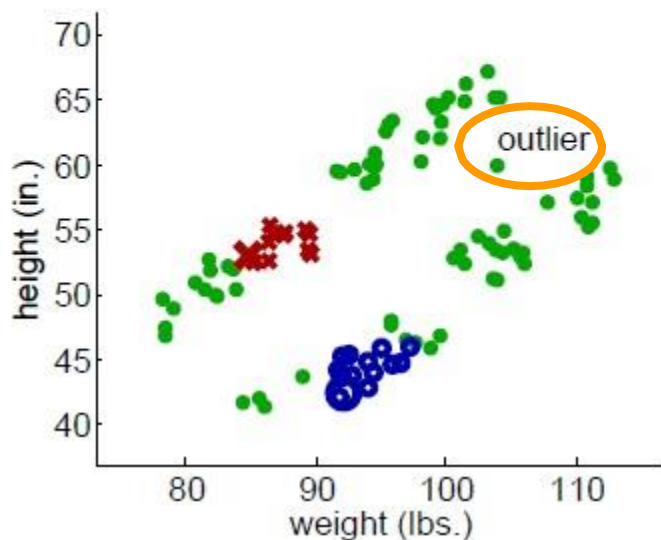


(d) Final labeling of all instances

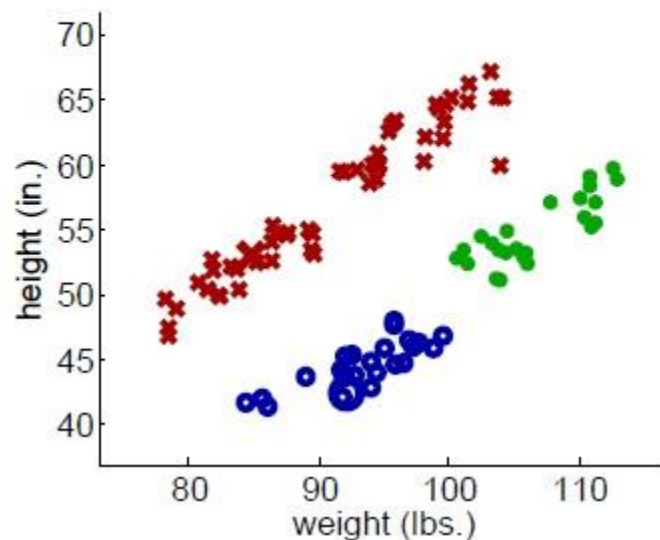


# Propagating 1-Nearest-Neighbor: now it doesn't

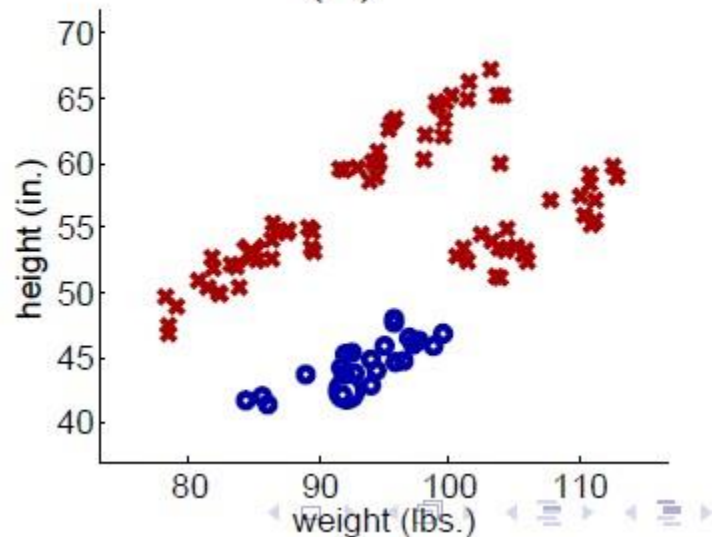
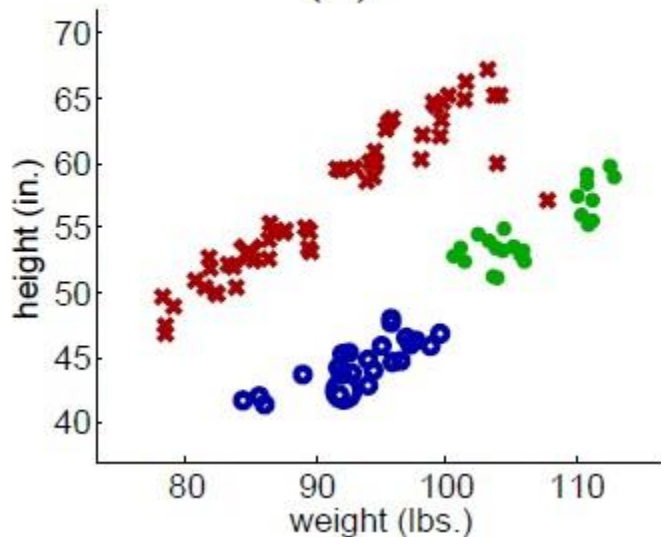
But with a single outlier...



(a)



(b)

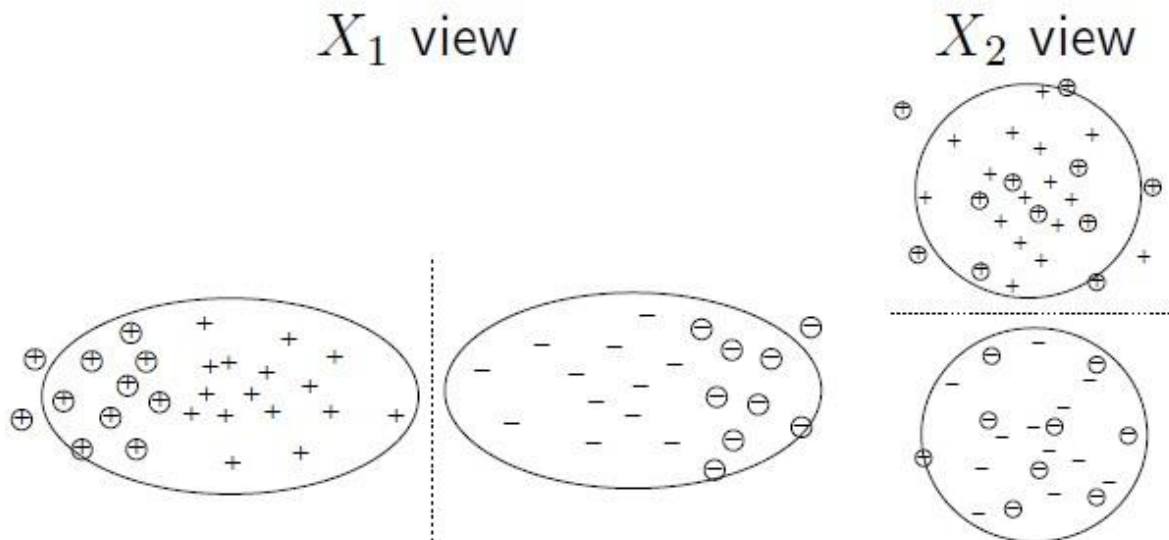




# Co-training

## Assumptions

- feature split  $x = [x^{(1)}; x^{(2)}]$  exists
- $x^{(1)}$  or  $x^{(2)}$  alone is sufficient to train a good classifier



# Co-training Algorithm

Co-training (Blum & Mitchell, 1998) (Mitchell, 1999) assumes that

- (i) features can be split into two sets;
  - (ii) each sub-feature set is sufficient to train a good classifier.
- 
- Initially two separate classifiers are trained with the labeled data, on the two sub-feature sets respectively.
  - Each classifier then classifies the unlabeled data, and ‘teaches’ the other classifier with the few unlabeled examples (and the predicted labels) they feel most confident.
  - Each classifier is retrained with the additional training examples given by the other classifier, and the process repeats.

# Co-training Algorithm

Blum & Mitchell'98

**Input:** labeled data  $\{(\mathbf{x}_i, y_i)\}_{i=1}^l$ , unlabeled data  $\{\mathbf{x}_j\}_{j=l+1}^{l+u}$   
each instance has two views  $\mathbf{x}_i = [\mathbf{x}_i^{(1)}, \mathbf{x}_i^{(2)}]$ ,  
and a learning speed  $k$ .

1. let  $L_1 = L_2 = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_l, y_l)\}$ .
2. Repeat until unlabeled data is used up:
3.     Train view-1  $f^{(1)}$  from  $L_1$ , view-2  $f^{(2)}$  from  $L_2$ .
4.     Classify unlabeled data with  $f^{(1)}$  and  $f^{(2)}$  separately.
5.     Add  $f^{(1)}$ 's top  $k$  most-confident predictions  $(\mathbf{x}, f^{(1)}(\mathbf{x}))$  to  $L_2$ .  
      Add  $f^{(2)}$ 's top  $k$  most-confident predictions  $(\mathbf{x}, f^{(2)}(\mathbf{x}))$  to  $L_1$ .  
      Remove these from the unlabeled data.

# Trial Exam Submission

[Submission Link](#)