

# CPSC 340: Machine Learning and Data Mining

Semi-Supervised Learning

Fall 2019

# Today: Semi-Supervised Learning

- Our usual supervised learning framework:

Egg	Milk	Fish	Wheat	Shellfish	Peanuts	...		Sick?
0	0.7	0	0.3	0	0		→	1
0.3	0.7	0	0.6	0	0.01		→	1
0	0	0	0.8	0	0		→	0
0.3	0.7	1.2	0	0.10	0.01		→	1

- In semi-supervised learning, we also have unlabeled examples:

Egg	Milk	Fish	Wheat	Shellfish	Peanuts	...
0.3	0	1.2	0.3	0.10	0.01	
0.6	0.7	0	0.3	0	0.01	
0	0.7	0	0.6	0	0	
0.3	0.7	0	0	0.20	0.01	

# Semi-Supervised Learning

- The semi-supervised learning (SSL) framework:

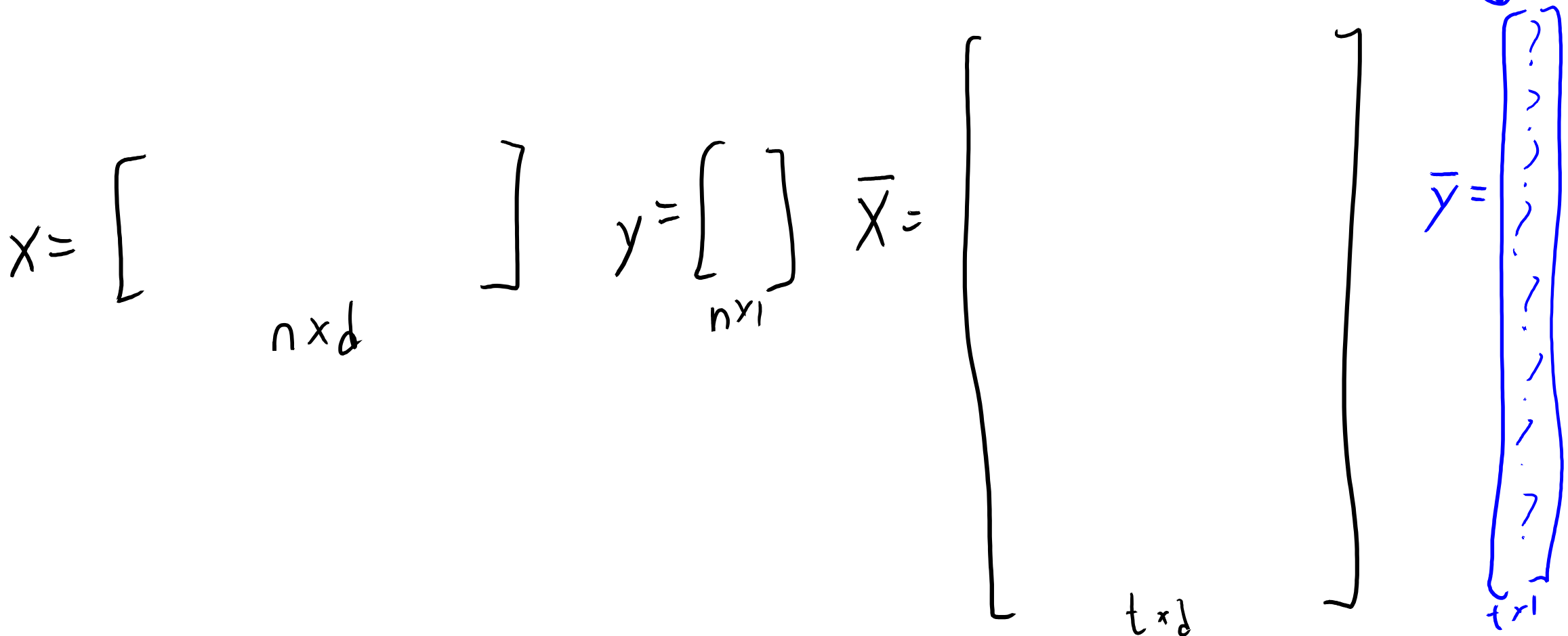
$$X = \begin{bmatrix} & \\ & \\ & \\ & \end{bmatrix}_{n \times d} \quad y = \begin{bmatrix} \\ \\ \\ \end{bmatrix}_{n \times 1} \quad \bar{X} = \begin{bmatrix} & \\ & \\ & \\ & \end{bmatrix}_{t \times d}$$

- This arises a lot:
  - Usually getting unlabeled data is easy but getting labeled data is hard.
  - Why build a classifier if getting labels is easy?
- Common situation:
  - A small number of labeled examples.
  - A huge number of unlabeled examples:  $t \gg n$ .

# Transductive vs. Inductive SSL

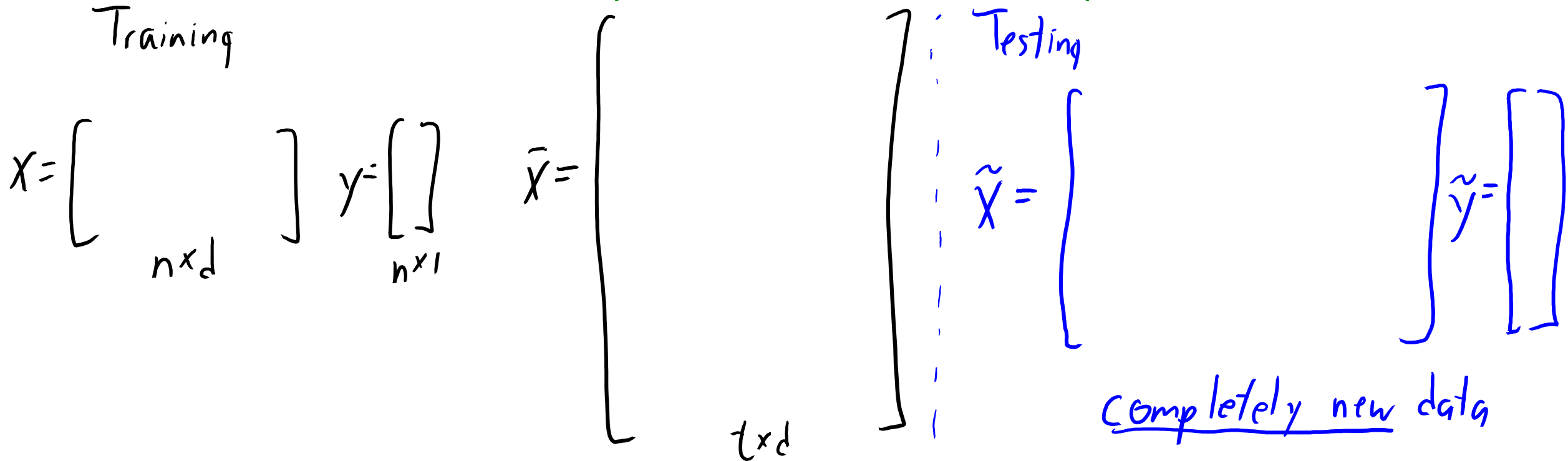
- Transductive SSL:

- Only interested in labels of the **given** unlabeled examples.



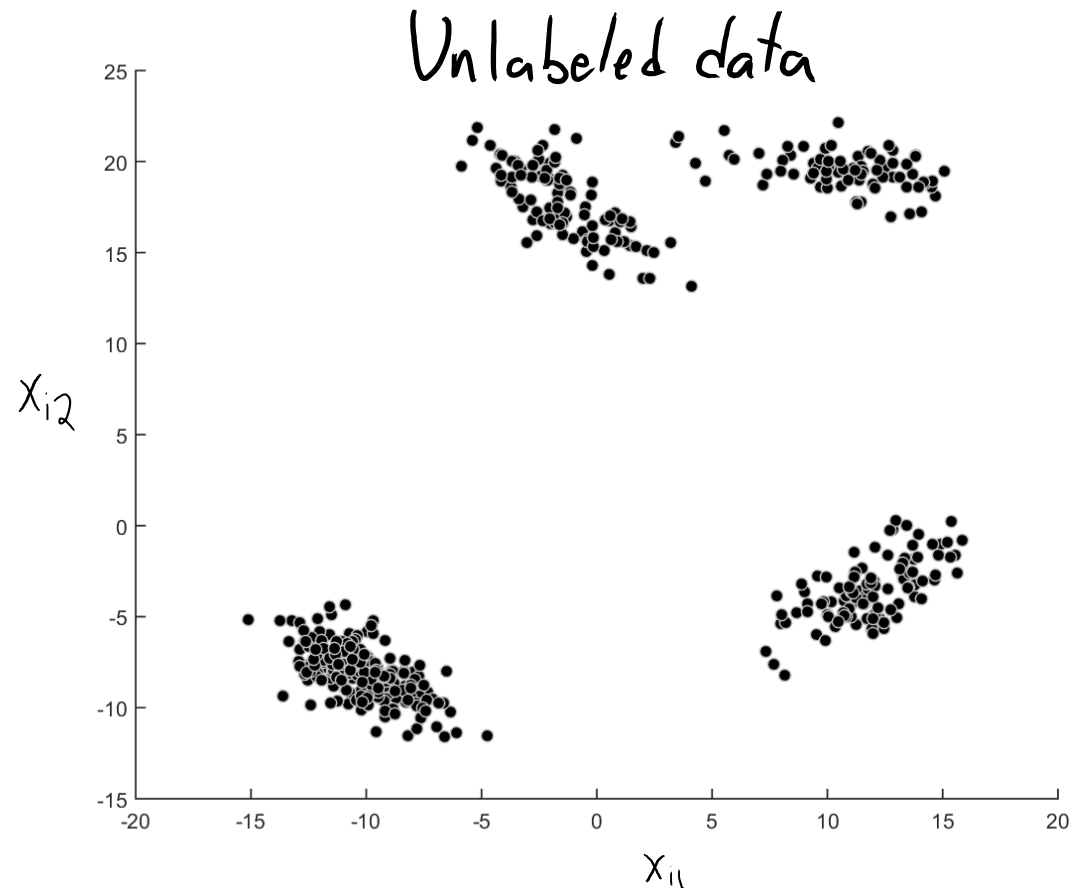
# Transductive vs. Inductive SSL

- **Transductive SSL:**
  - Only interested in **labels of the given** unlabeled examples.
- **Inductive SSL:**
  - Interested in the **test set performance on new examples.**



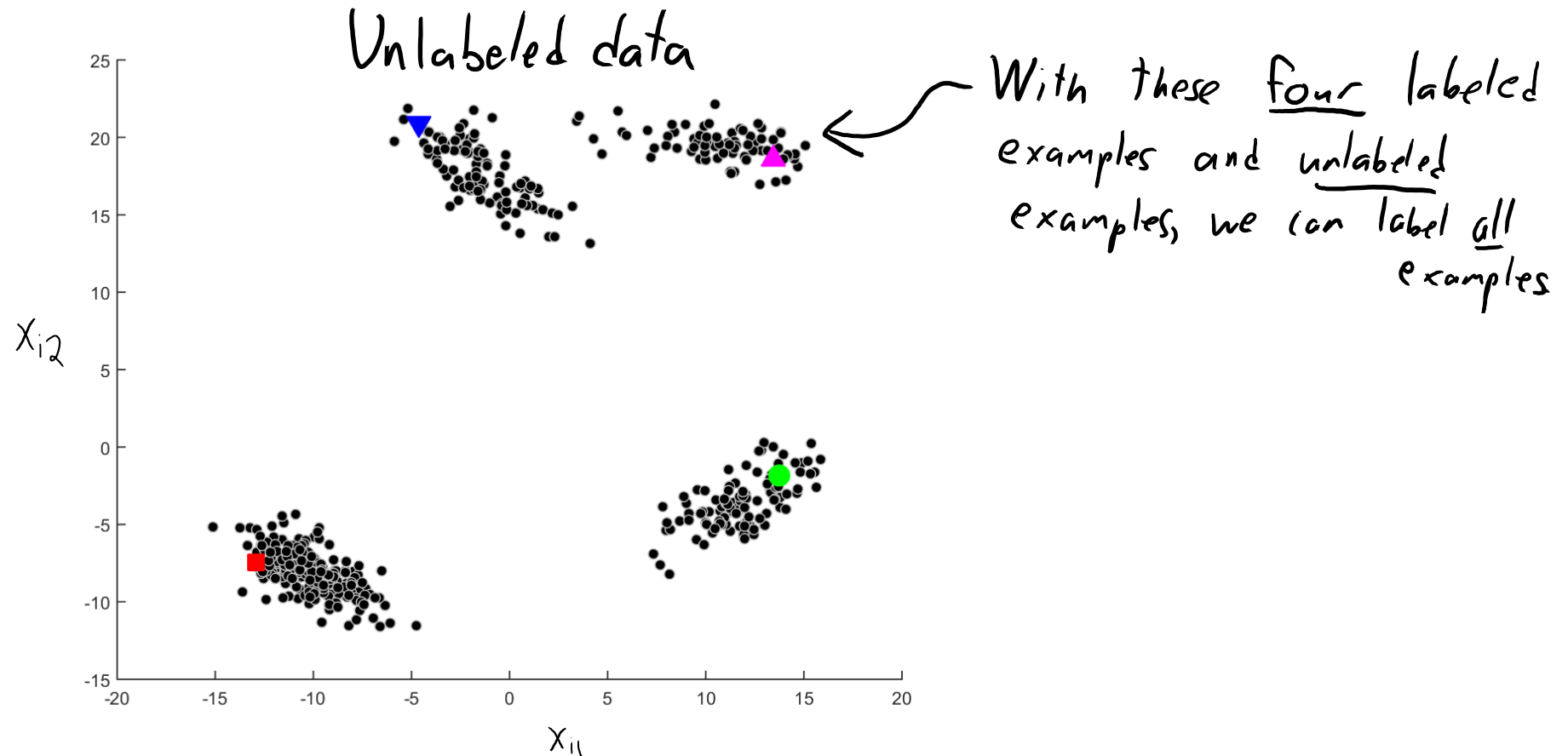
# Semi-Supervised Learning

- Why should unlabeled data tell us anything about the labels?
  - Usually, we assume that: (similar features  $\Leftrightarrow$  similar labels).



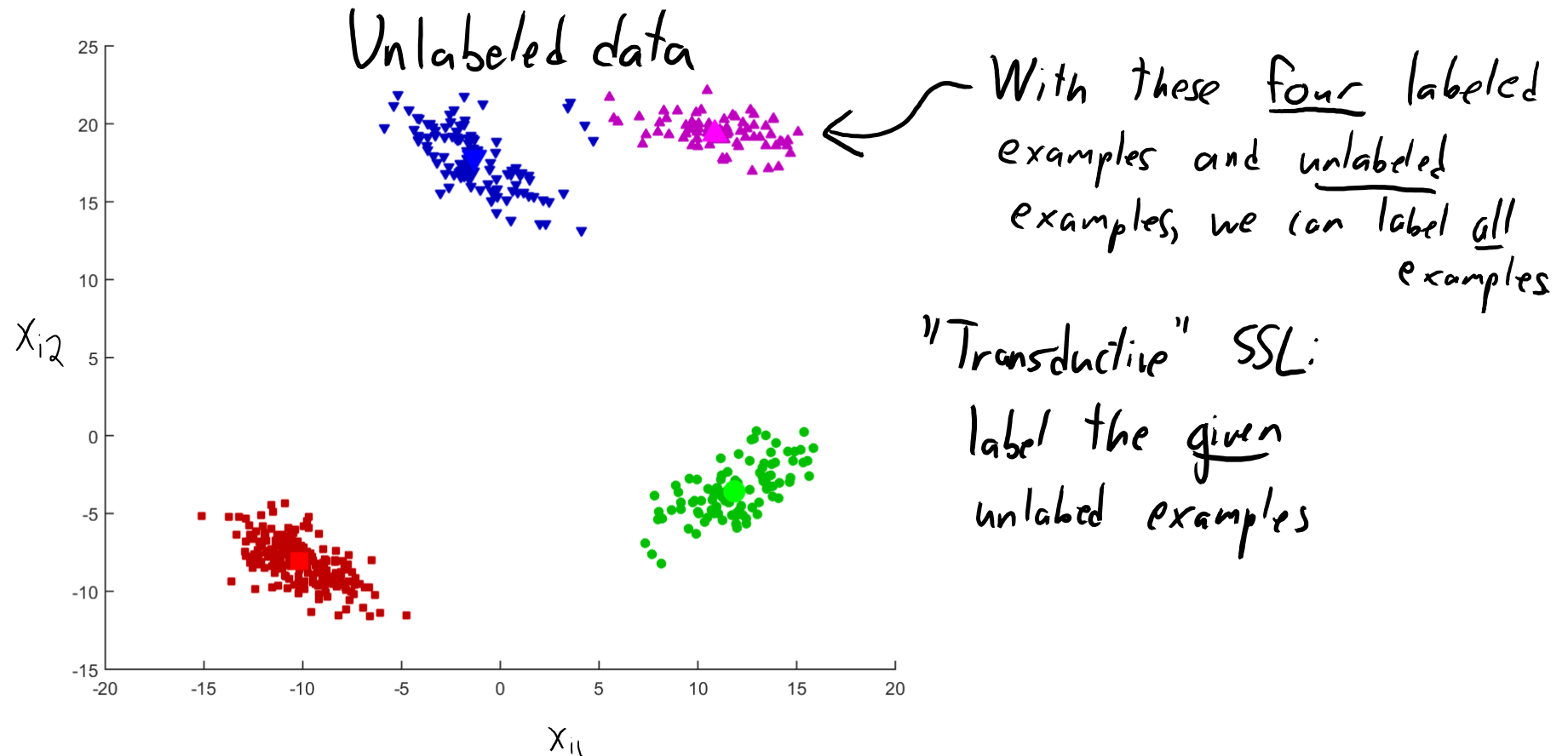
# Semi-Supervised Learning

- Why should unlabeled data tell us anything about the labels?
  - Usually, we assume that: (similar features  $\Leftrightarrow$  similar labels).



# Semi-Supervised Learning

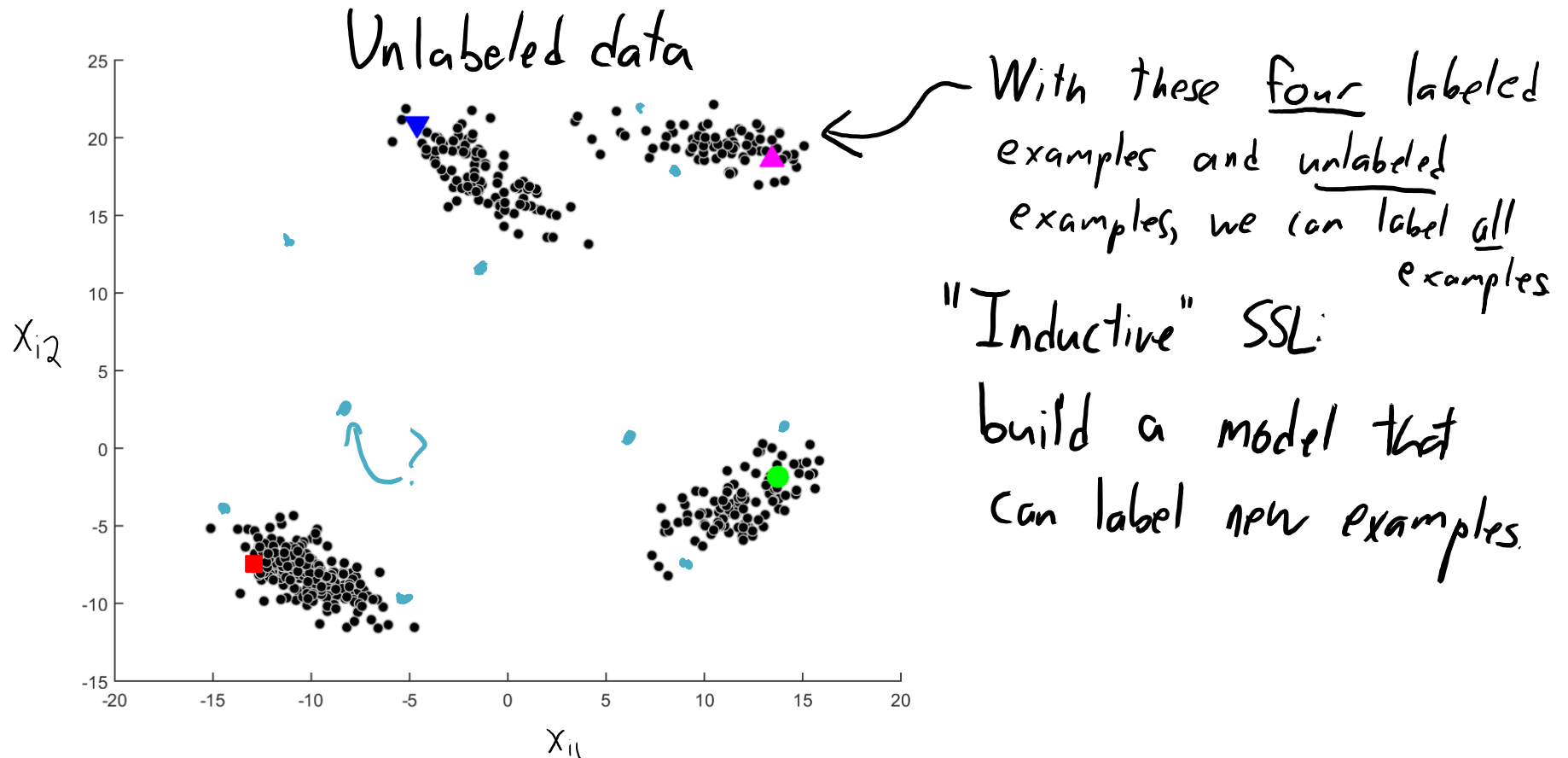
- Why should unlabeled data tell us anything about the labels?
  - Usually, we assume that: (similar features  $\Leftrightarrow$  similar labels).





# Semi-Supervised Learning

- Why should unlabeled data tell us anything about the labels?
  - Usually, we assume that: (similar features  $\Leftrightarrow$  similar labels).



# Philosophical Digression: Can we do SSL?

- Will unlabeled examples help in general?
  - No!
- Consider choosing random  $x_i$  values, then computing  $y_i$ .
  - Unlabeled examples collected in this way will not help.
  - By construction, distribution of  $x_i$  says nothing about  $y_i$ .

# Philosophical Digression: Can we do SSL?

- Example where SSL is not possible:
  - Try to detect food allergy by trying 'random' combinations of food.
    - The actual 'random' process isn't important, as long it doesn't depend on 'y<sub>i</sub>'.
  - Unlabeled data would be more random combinations:

$$X = \begin{bmatrix} \text{"random"} \\ \text{values} \end{bmatrix} \quad y = \begin{bmatrix} \text{labels} \\ \text{of} \\ \text{random} \\ \text{features} \end{bmatrix} \quad \tilde{X} = \begin{bmatrix} \text{more} \\ \text{"random"} \\ \text{values} \end{bmatrix}$$

- You can generate all possible unlabeled data, but it says nothing about labels.

# Philosophical Digression: Can we do SSL?

- Example where SSL is possible:
  - Trying to classify images as ‘cat’ vs. ‘dog’:
    - Unlabeled data would be images of cats or dogs: **not random images**.
      - Unlabeled data contains information about what images of cats *and* dogs look like.
      - E.g., clusters or manifolds in unlabeled images.
- Contrast this with ‘cat’ vs. ‘not cat’:
  - If we generate random images then label them, unlabeled data won’t help.
  - If we know that half our unlabeled images are cats, unlabeled could help.



# Philosophical Digression: Can we do SSL?

- When can unlabeled examples help?
- Consider ' $y_i$ ' somehow influencing data we collect:
  - Now there is information about labels contained in unlabeled examples.
  - Example 1: we try to have an even number of  $y_i = +1$  and  $y_i = -1$ .
  - Example 2: we need to choose non-random ' $x_i$ ' to correspond to a valid ' $y_i$ '
  - We are almost always in this case.

(pause)

# SSL Approach 1: Self-Taught Learning

- **Self-taught** learning is similar to k-means:
  1. Fit a model based on the labeled data.
  2. Use the model to label the unlabeled data.
  3. Use estimated labels to fit model based on labeled and unlabeled data.
  4. Go back to 2.
- Obvious problem: it can **reinforce errors** and even diverge.

- Possible fixes:

- Only use labels are you very confident about.
- Regularize the loss from the unlabeled examples:

$$f(w) = \frac{1}{2} \|Xw - y\|^2 + \frac{\lambda}{2} \|\bar{X}w - \hat{y}\|^2$$

prediction from step 2

$\lambda$  controls how much we trust guesses on unlabeled data

# SSL Approach 1: Self-Taught Learning

Input:

- Labeled examples  $\{X, y\}$

- Unlabeled examples  $\bar{X}$

1. Train on  $\{X, y\}$ :

$$\text{model} = \text{fit}(X, y)$$

2. Guess labels:

$$\hat{y} = \text{model.predict}(\bar{X})$$

3. Train on bigger data set:

$$\text{model} = \text{fit}\left(\begin{bmatrix} X \\ \bar{X} \end{bmatrix}, \begin{bmatrix} y \\ \hat{y} \end{bmatrix}, \lambda\right)$$

repeat



Popular variants:

1. "Expectation maximization"

2. "Yarowsky" algorithm (language)



# SSL Approach 2: Co-Training

- Assumes that we have **2 sets of features**:
  - Both sets are sufficient to give high accuracy.
  - The sets are conditionally independent given the label.
  - E.g., image features (set 1) and caption features (set 2).
- **Co-training**:
  1. Using labeled set, fit model 1 based on set 1, fit model 2 based on set 2.
  2. **Label a random subset** of unlabeled examples based on both models.
  3. Move examples where each classifier is **most confident** to labeled set.
  4. Go back to 1.
- Hope is that models “teach” each other to achieve consensus.
  - Theoretically works if **assumptions above** are satisfied.



# SSL Approach 2: Co-Training

0. Split features into  $X_1$  and  $X_2$

$$X = \begin{bmatrix} X_1 & \vdots & X_2 \end{bmatrix}$$

1. Train models on  $X_1$  and  $X_2$ :

$$\text{model1} = \text{fit}(X_1, y)$$

$$\text{model2} = \text{fit}(X_2, y)$$

We choose subset to avoid a bias from having similar "most confident" examples.

2. Choose random subset of unlabeled examples and predict label:

$$\hat{y}_1 = \text{model.predict}(\bar{X}_1)$$

$$\hat{y}_2 = \text{model.predict}(\bar{X}_2)$$

3. For each model find the  $\bar{x}_i$  values in the sample that model is most

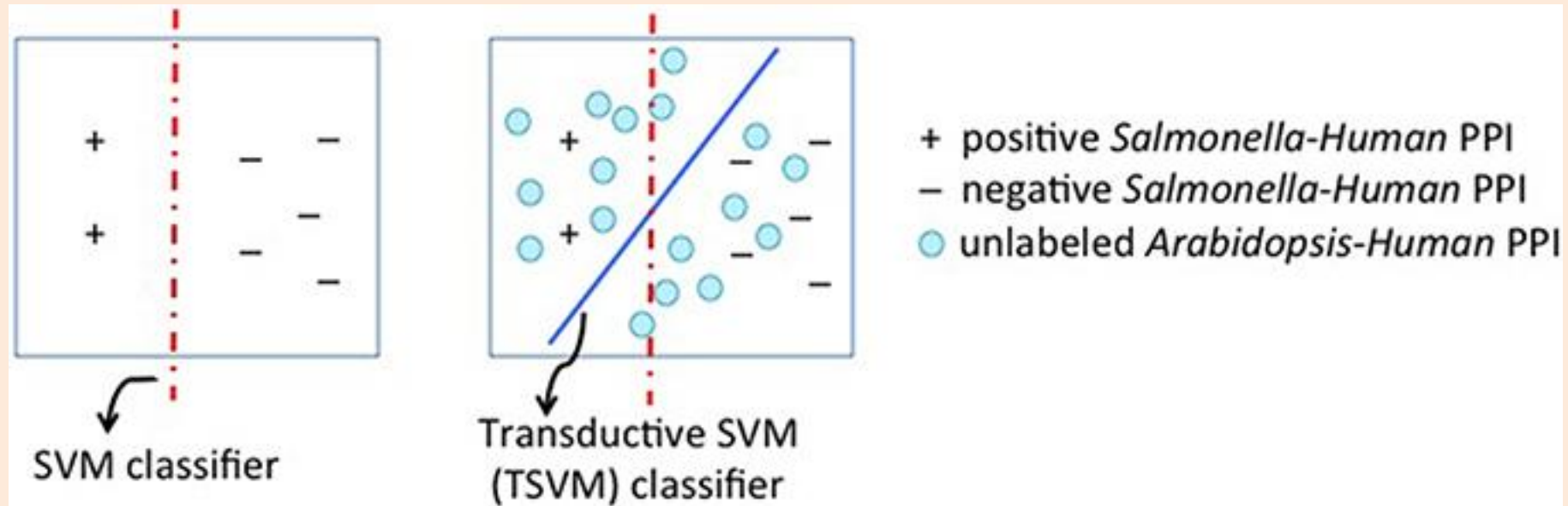
4. Add these examples to labeled set.

confident in for each class.

repeat

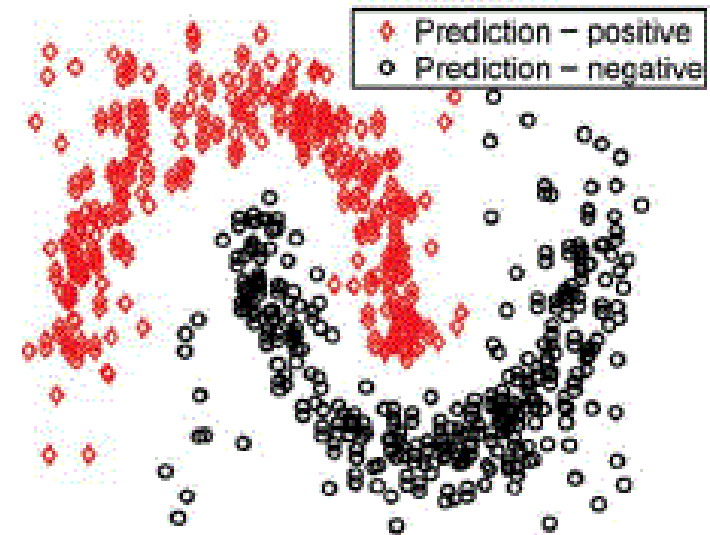
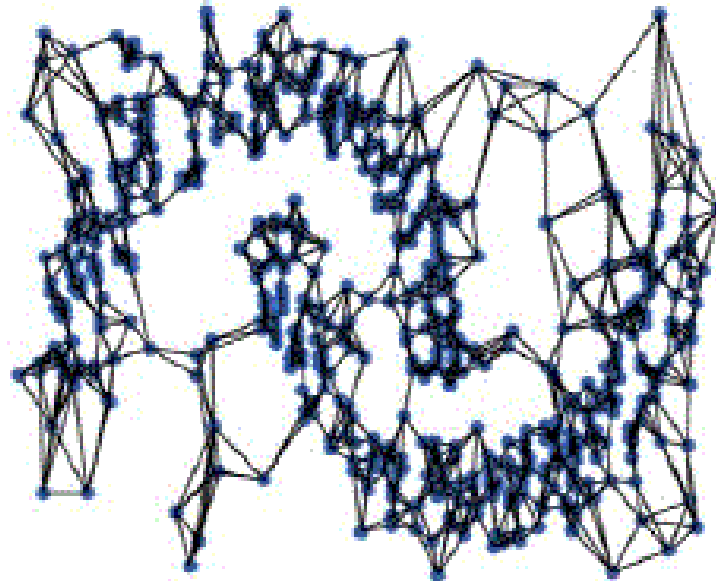
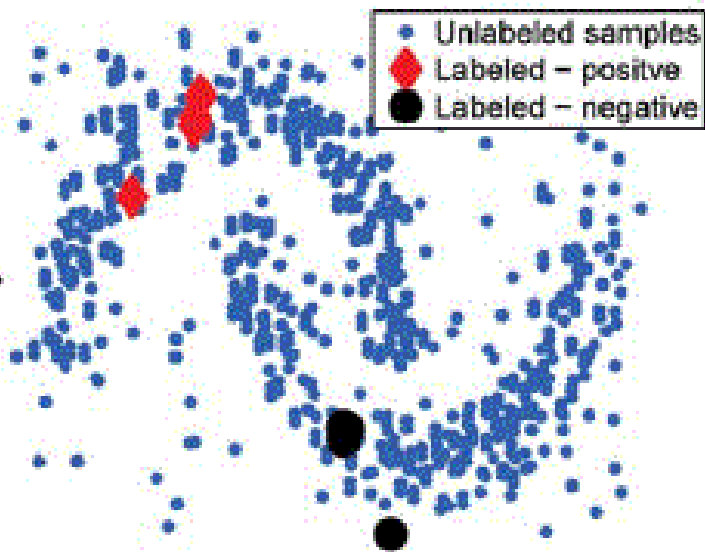
# SSL Approach 3: Entropy Regularization

- Self-taught and co-training predictions may **propagate errors**.
- Instead of making predictions, encourage “predictability”:
  - **Entropy regularization**: penalize “randomness” of labels on unlabeled.
  - **Transductive SVMs**: avoid decision boundaries in dense regions.

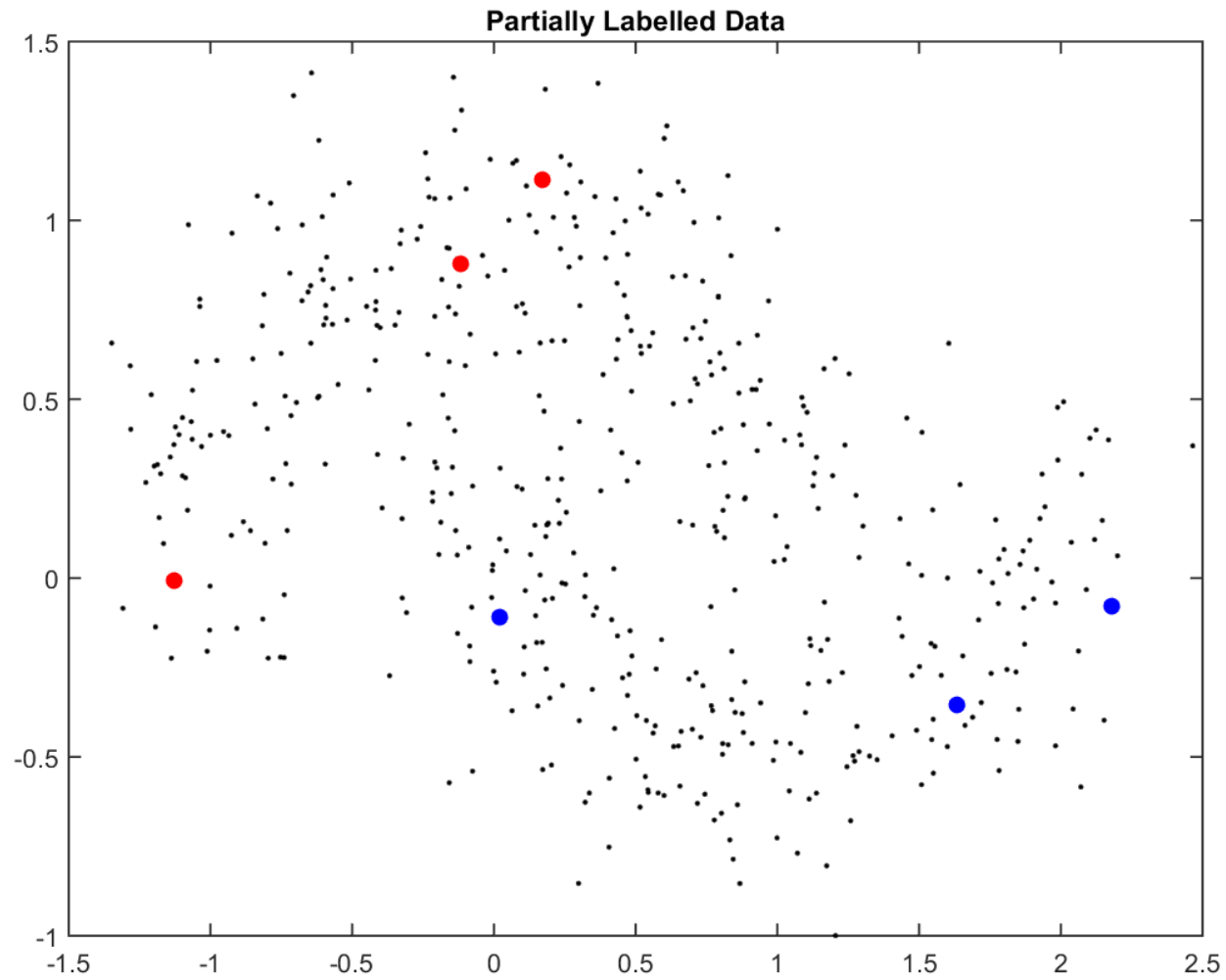


# Graph-Based Methods (Label Propagation)

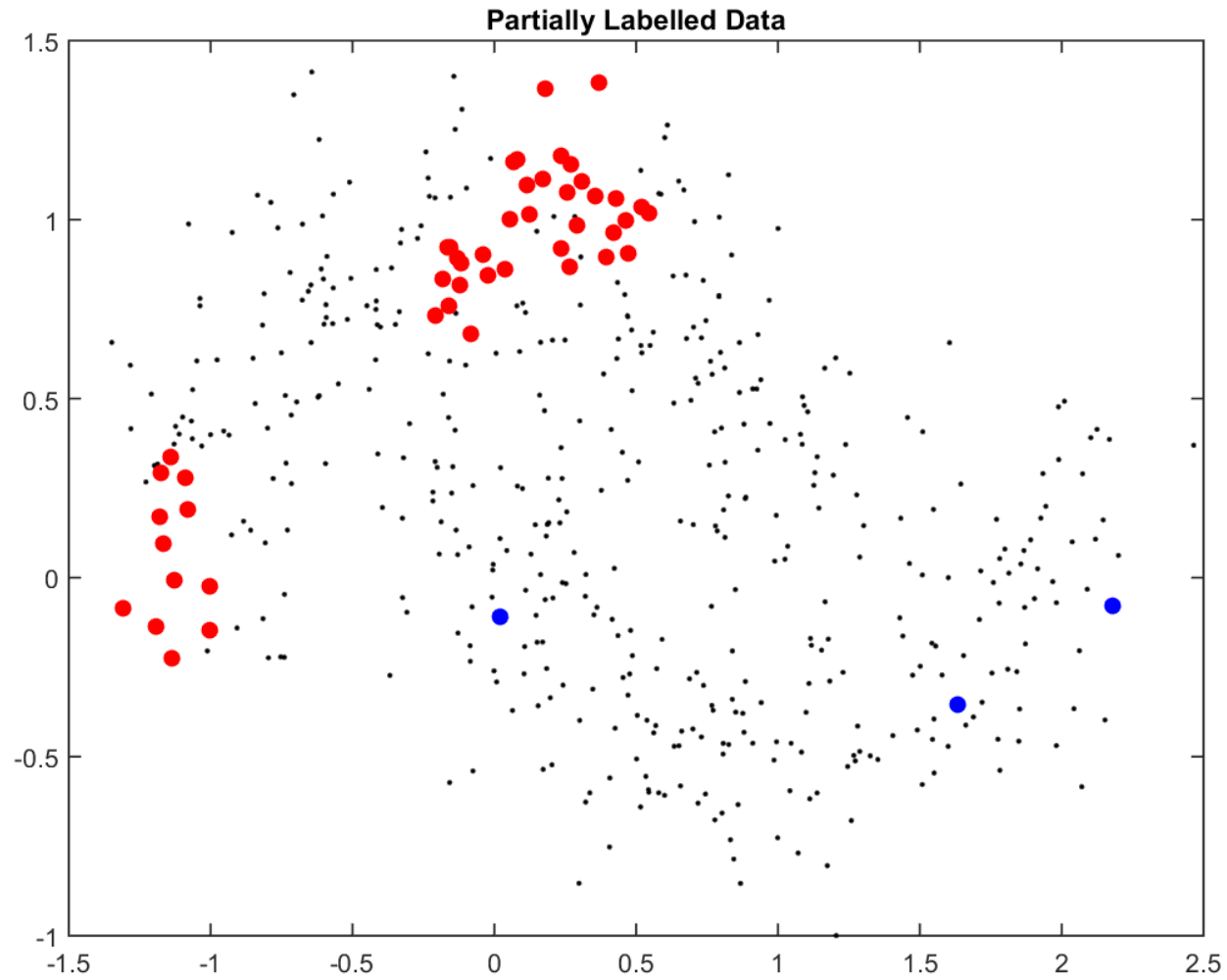
- We can only do SSL because (similar features  $\Leftrightarrow$  similar labels).
- Graph-based SSL uses this directly.
  - Define weighted graph on training examples:
    - For example, use KNN graph or points within radius ' $\epsilon$ '.
    - Weight is how 'important' it is for nodes to share label.



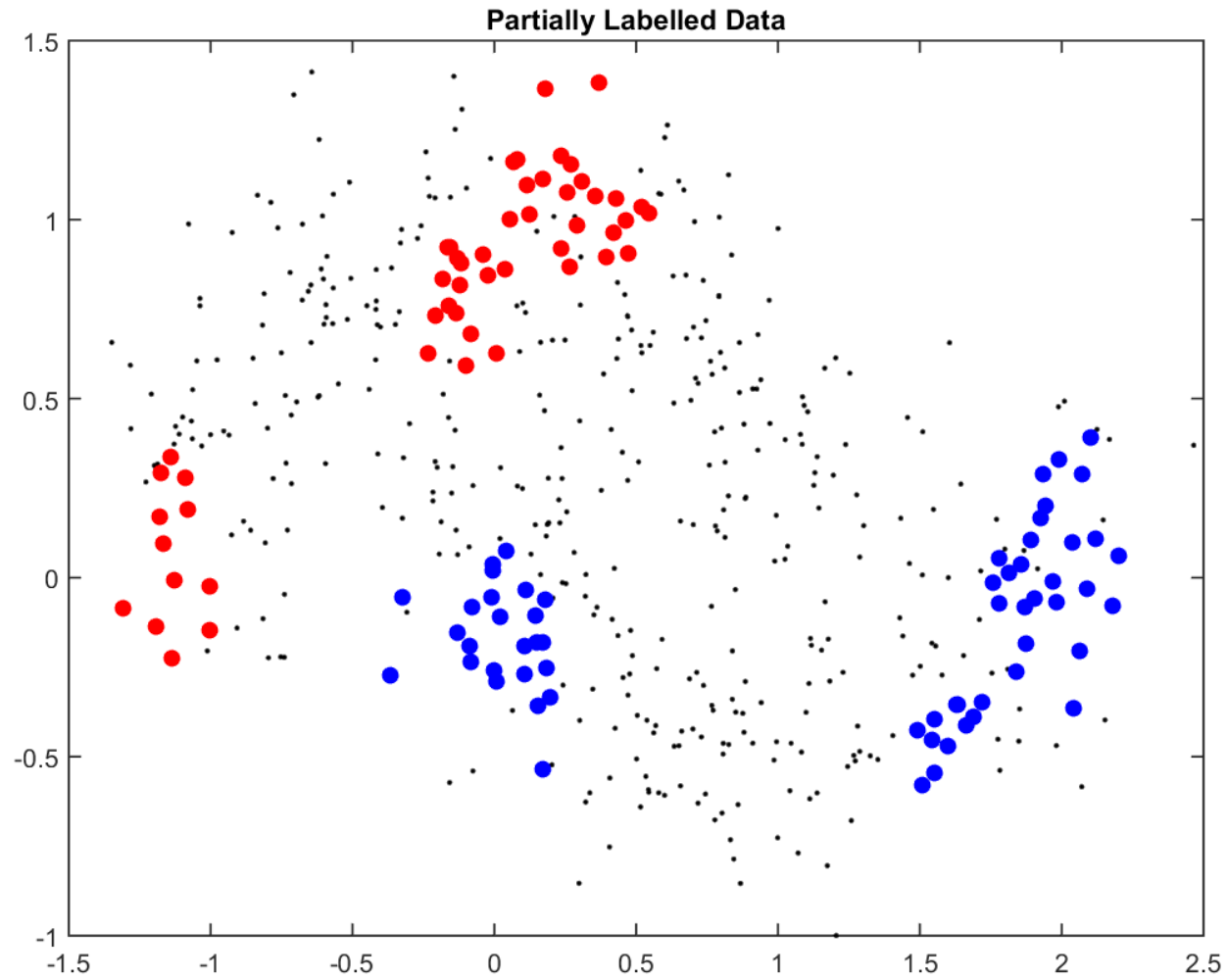
# Label Propagation in Action



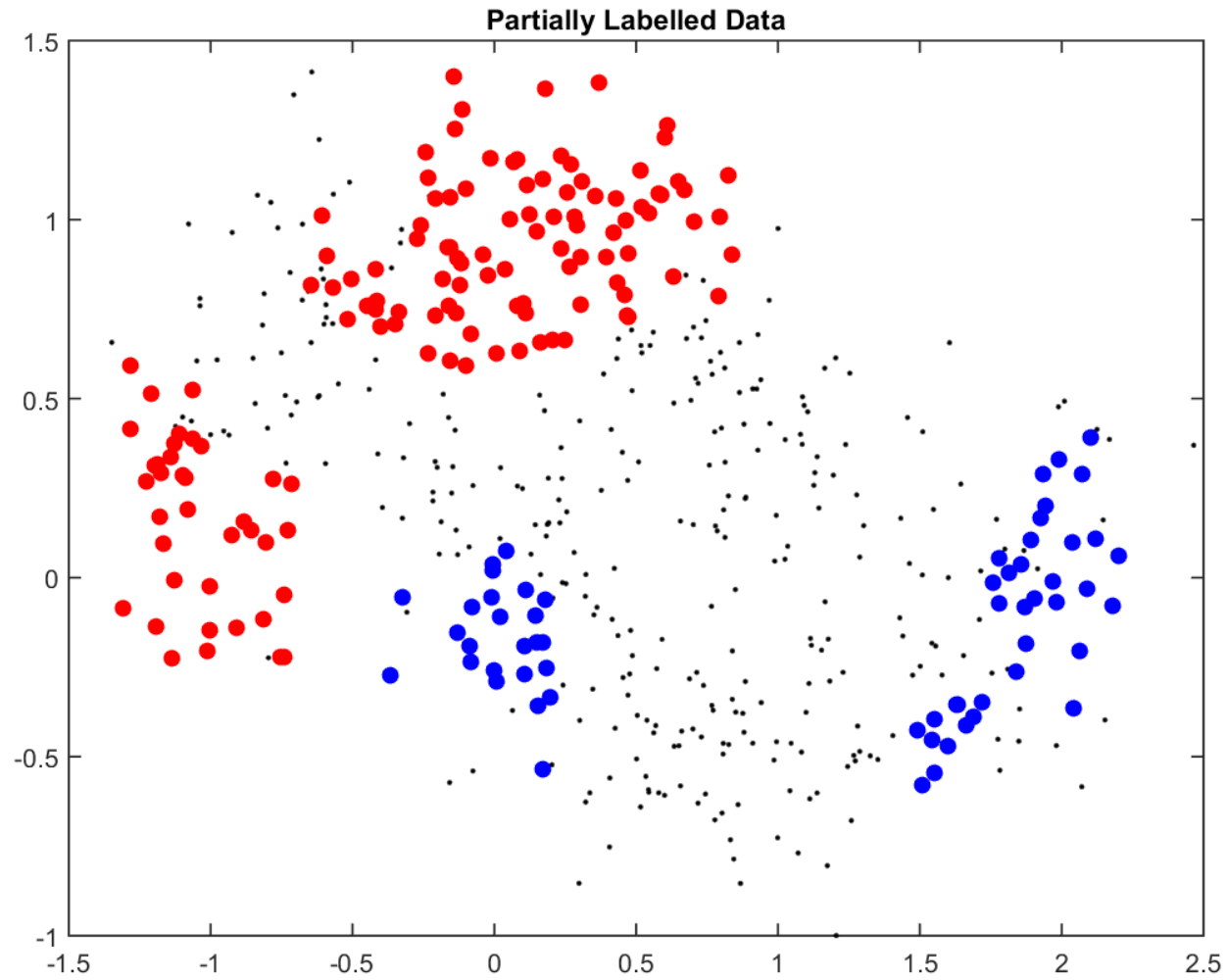
# Label Propagation in Action



# Label Propagation in Action

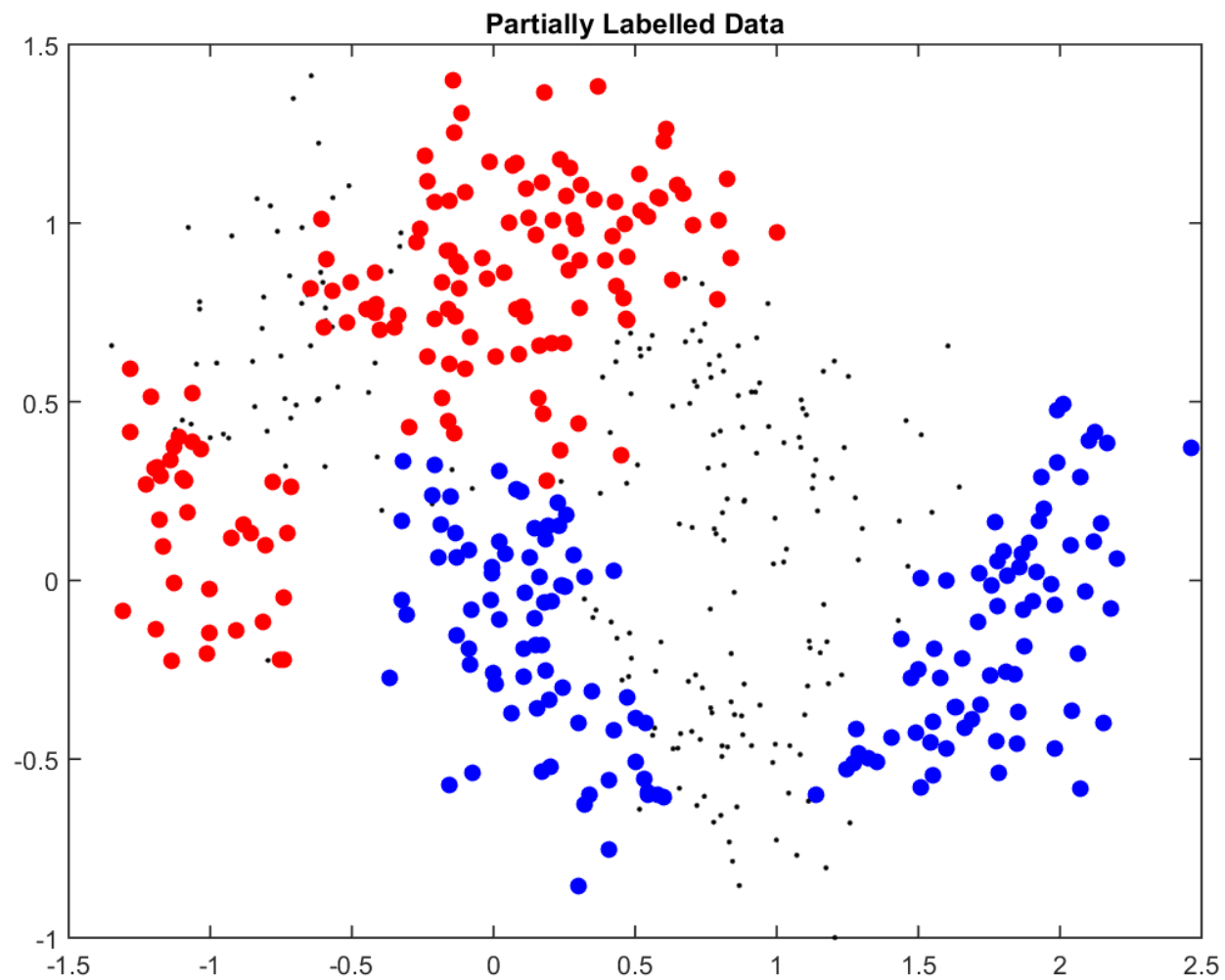


# Label Propagation in Action

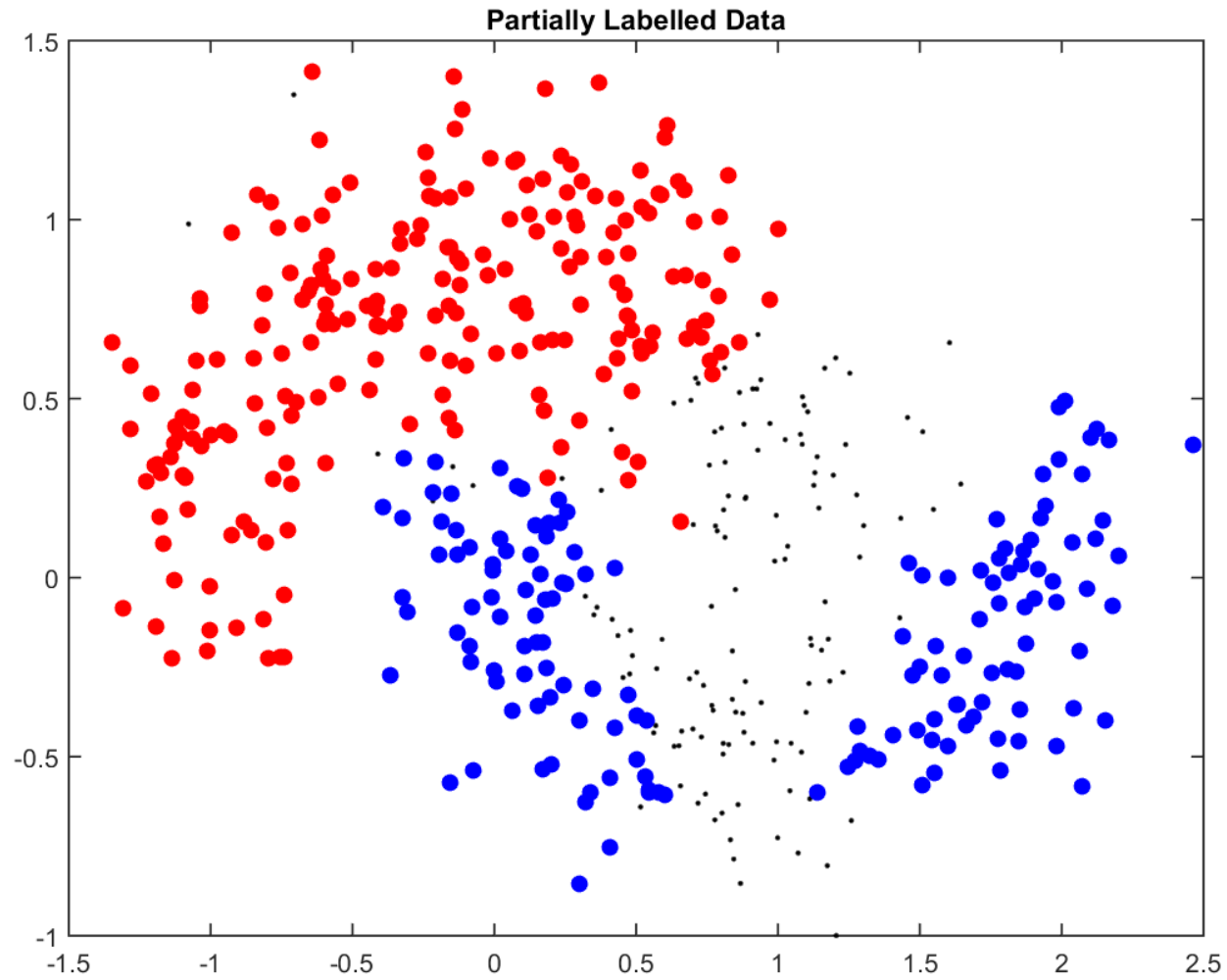




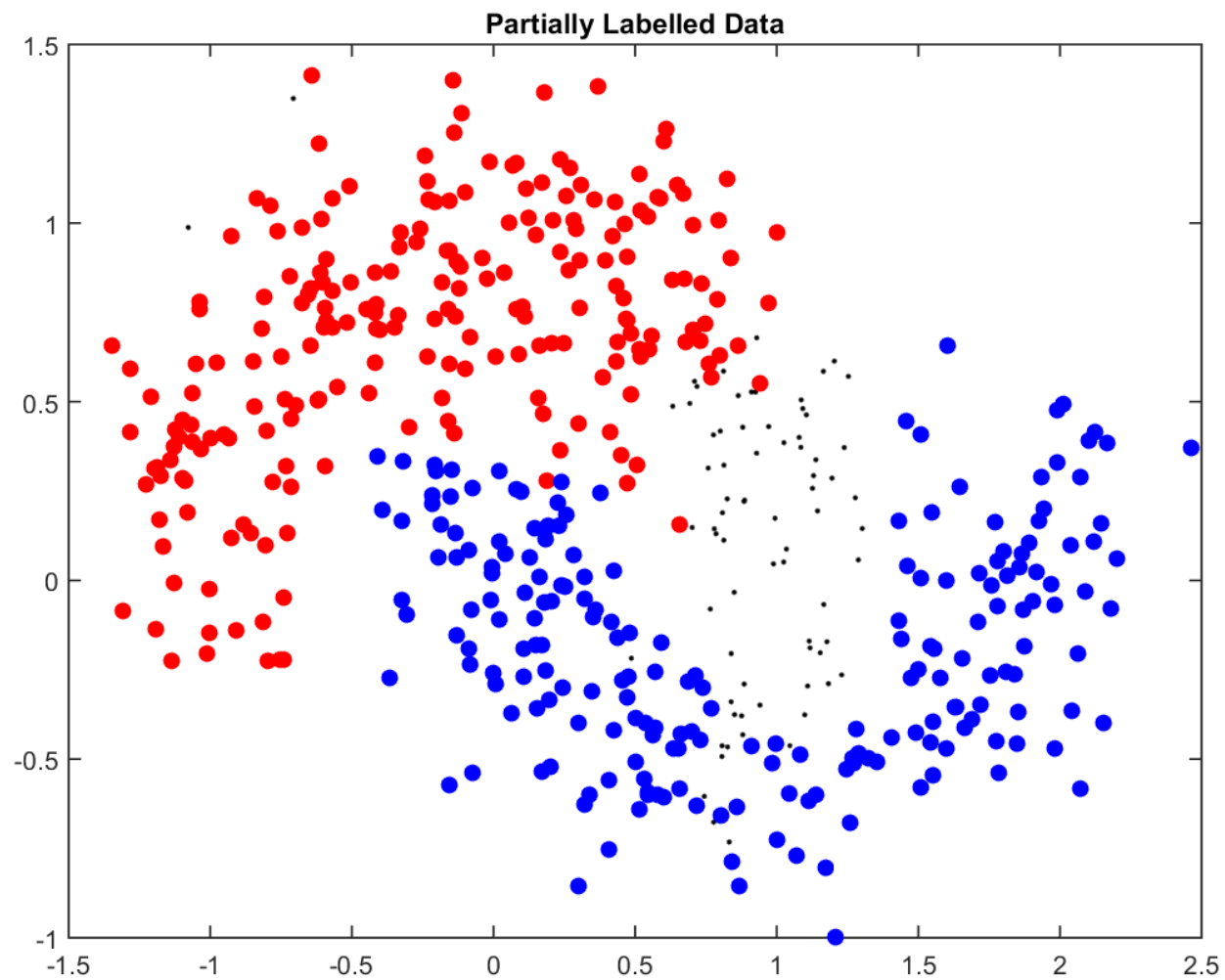
# Label Propagation in Action



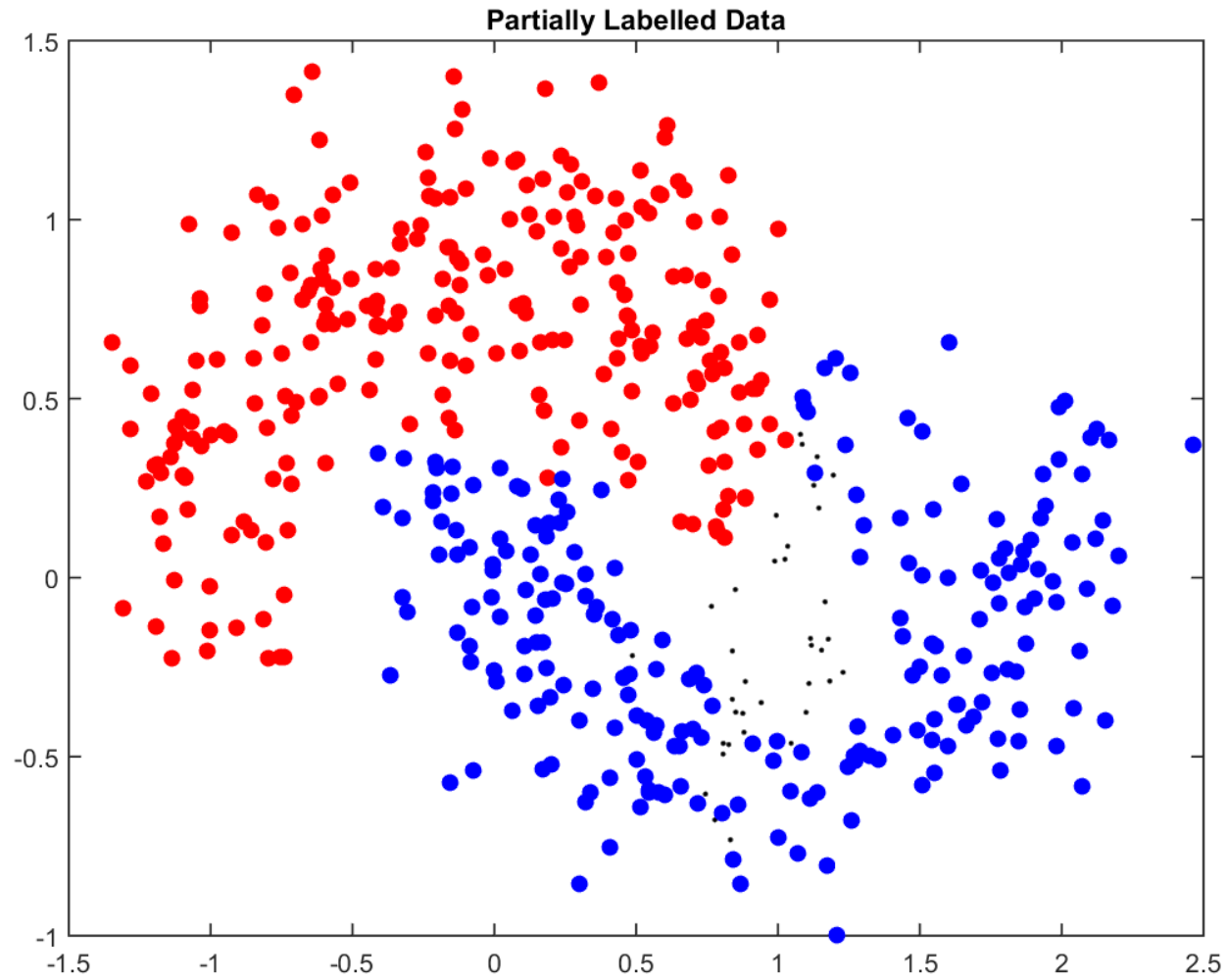
# Label Propagation in Action



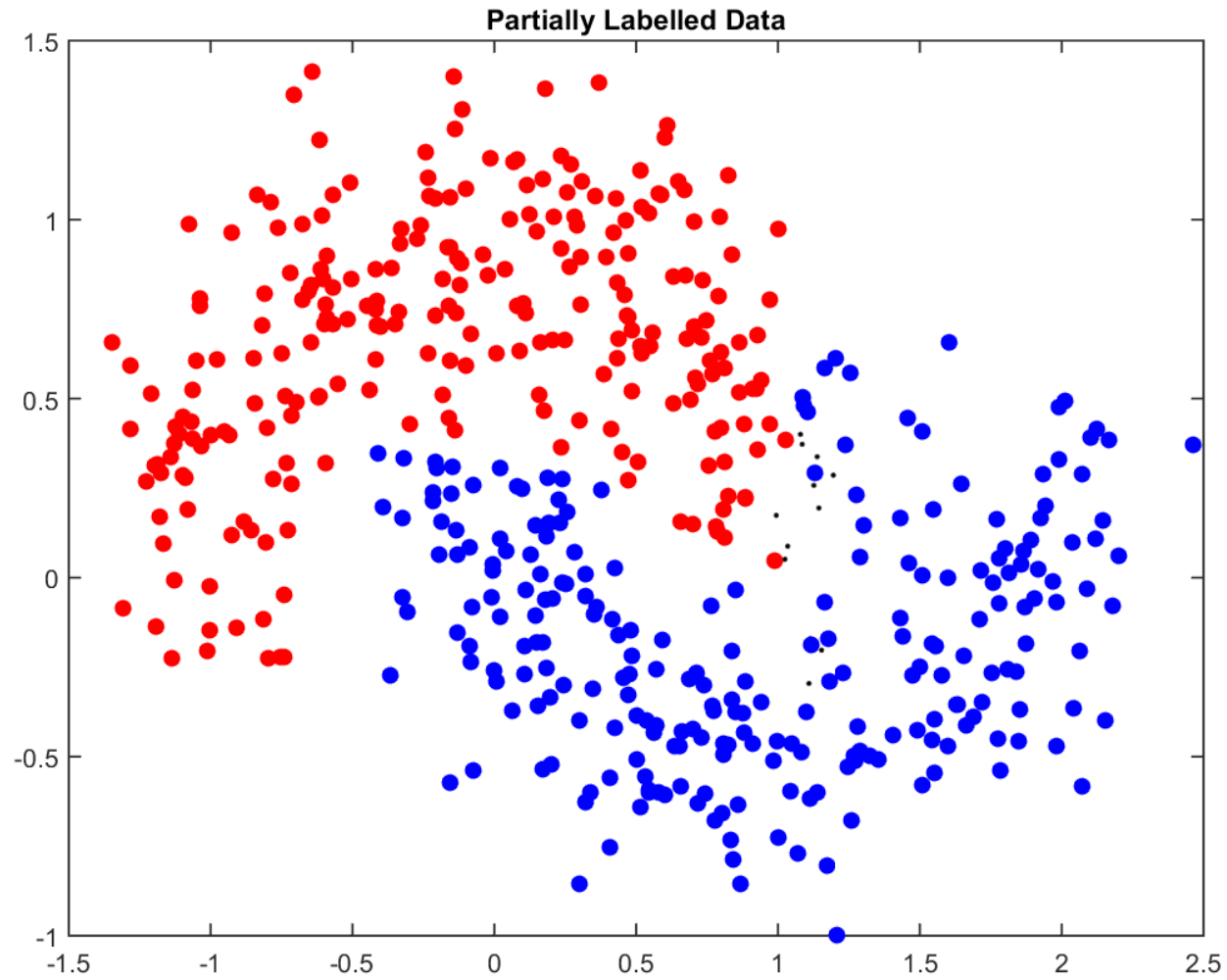
# Label Propagation in Action



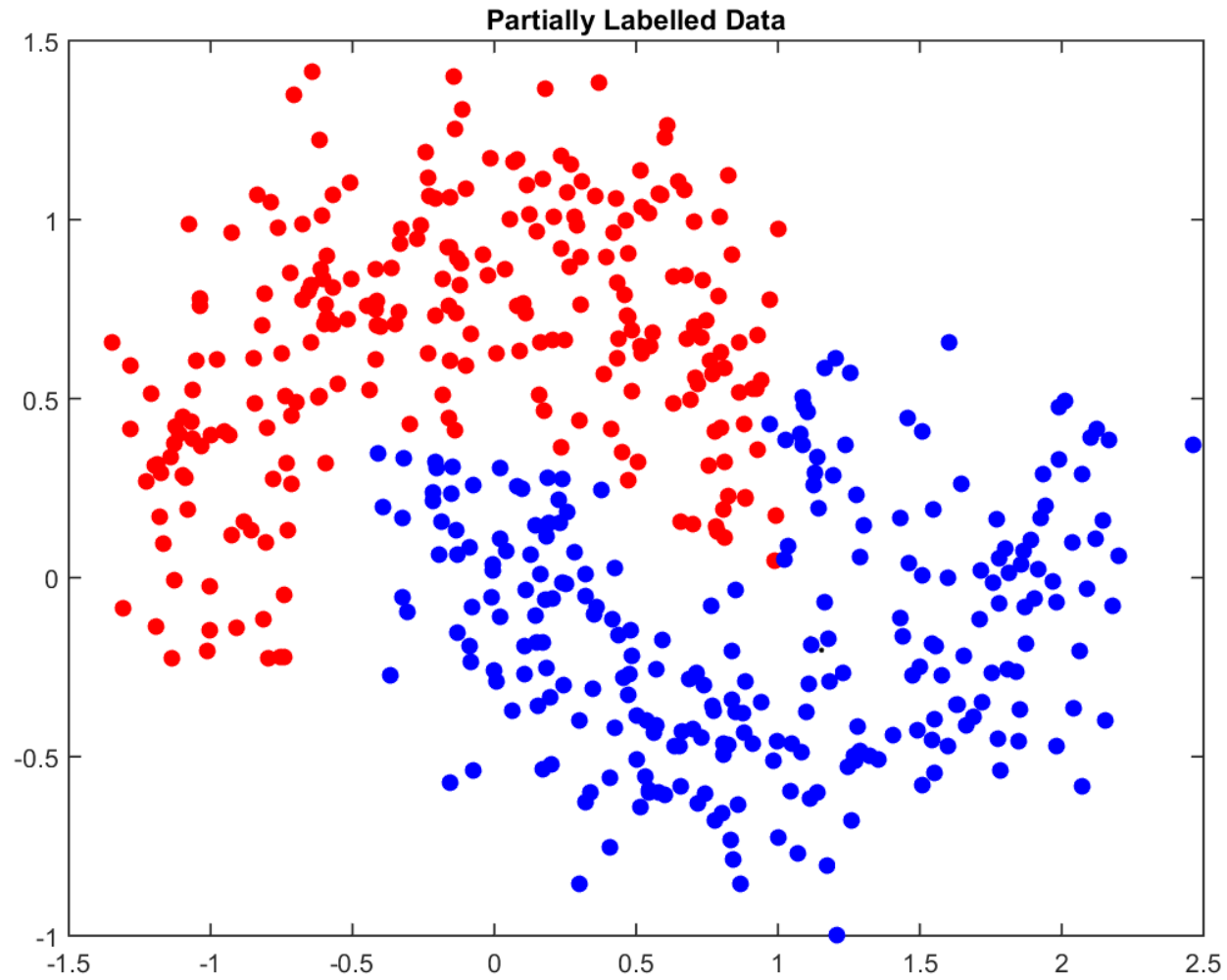
# Label Propagation in Action



# Label Propagation in Action



# Label Propagation in Action



# Graph-Based SSL (Label Propagation)

- Treat **unknown labels as variables**, minimize **cost of disagreement**:

$$f(\bar{y}) = \sum_{i=1}^n \sum_{j=1}^t w_{ij} (y_i - \bar{y}_j)^2 + \frac{1}{2} \sum_{i=1}^t \sum_{j=1}^t \bar{w}_{ij} (\bar{y}_i - \bar{y}_j)^2$$

Do gradient descent on labels of unlabeled examples

graph weight between labeled and unlabeled

make  $\tilde{y}_j$  similar to labeled neighbour

make unlabeled neighbours similar to each other

- Common variations:

- Treat labels  $y_i$  as variables (they might be wrong).
  - Weight how much you trust original labels.
- Regularize the unlabeled  $\bar{y}_i$  towards a default value.
  - Can reflect that example is really far from any labeled example.

Leads to "label propagation" through graph.

# Example: Tagging YouTube Videos

- Example:
  - Consider assigning ‘tags’ to YouTube videos (e.g., ‘cat’).
  - Construct a graph based on sequences of videos that people watch.
    - Give high weight if video A is often followed/preceded by video B.
  - Use label propagation to tag all videos.



- Becoming popular in **bioinformatics**:
  - Label a subset of genes using manual experiments.
  - Find out which genes interact using more manual experiments.
  - Predict function/location/etc of genes using label propagation.
- Comments on **graph-based SSL**:
  - **Transductive** method: only estimates the unknown labels.
  - **Often surprisingly effective** even if you only have a few labels.
  - **Does not need features** if you have the weighted graph.



# Summary

- **Semi-supervised learning** uses unlabeled data in supervised task.
  - **Transductive learning** only focuses on labeling this data.
  - **SSL may or may not help**, depending on structure of data.
- **Self-taught/co-training** alternate labeling/fitting.
- **Graph-based SSL** propagates labels in graph (no features needed).