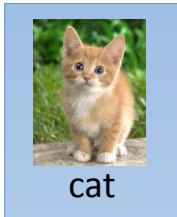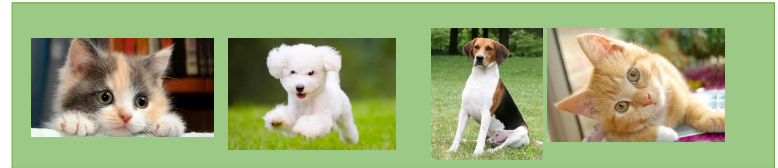# Semi-supervised Learning

# Introduction

cat

dog

Unlabeled data

(Image of cats and dogs without labeling)

- Supervised learning: $\{(x^r, \hat{y}^r)\}_{r=1}^R$
  - E.g. $x^r$: image, $\hat{y}^r$: class labels
- Semi-supervised learning: $\{(x^r, \hat{y}^r)\}_{r=1}^R, \{x^u\}_{u=R}^{R+U}$
  - A set of unlabeled data, usually U >> R
  - Transductive learning: unlabeled data is the testing data
  - Inductive learning: unlabeled data is not the testing data
- Why semi-supervised learning?
  - Collecting data is easy, but collecting "labelled" data is expensive
  - We do semi-supervised learning in our lives

# Outline

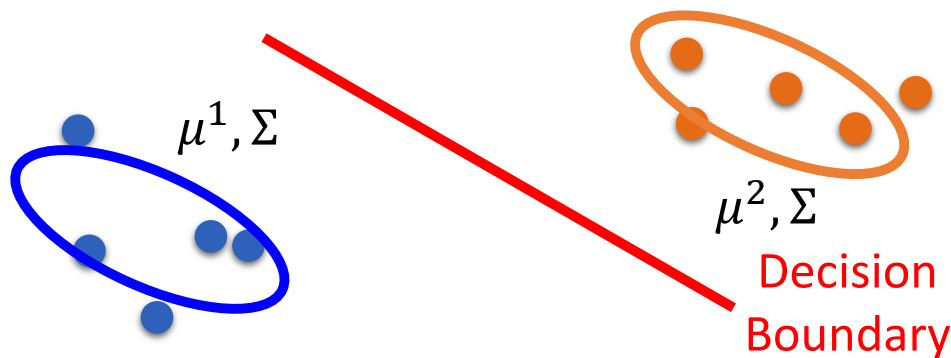Semi-supervised Learning for Generative Model

Low-density Separation Assumption

Smoothness Assumption

Better Representation

# Supervised Generative Model

- Given labelled training examples $x^r \in C_1, C_2$
  - looking for most likely prior probability P(C$_i$) and class-dependent probability P(x|C$_i$)
  - P(x|C$_i$) is a Gaussian parameterized by $\mu^i$ and $\Sigma$



$\mu^1, \Sigma$

$\mu^2, \Sigma$

Decision Boundary

With $P(C_1), P(C_2), \mu^1, \mu^2, \Sigma$    $P(C_1|x) = \dfrac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)}$

The unlabeled data $x^u$ help re-estimate $P(C_1), P(C_2), \mu^1, \mu^2, \Sigma$

# Semi-supervised Generative Model

- Initialization: $\theta = \{P(C_1), P(C_2), \mu^1, \mu^2, \Sigma\}$

**E**

- Step 1: compute the posterior probability of unlabeled data

$$P_\theta(C_1|x^u)$$

Depending on model $\theta$

↑ Back to step 1

**M**

- Step 2: update model

$$P(C_1) = \frac{N_1 + \sum_{x^u} P(C_1|x^u)}{N}$$

$N$: total number of examples

$N_1$: number of examples belonging to $C_1$

$$\mu^1 = \frac{1}{N_1} \sum_{x^r \in C_1} x^r + \frac{1}{\sum_{x^u} P(C_1|x^u)} \sum_{x^u} P(C_1|x^u) x^u \quad ......$$

# Why?

$$\theta = \{P(C_1), P(C_2), \mu^1, \mu^2, \Sigma\}$$

- Maximum likelihood with labelled data

$$logL(\theta) = \sum_{(x^r, \hat{y}^r)} logP_\theta(x^r | \hat{y}^r)$$

- Maximum likelihood with labelled + unlabeled data

$$logL(\theta) = \sum_{(x^r, \hat{y}^r)} logP_\theta(x^r | \hat{y}^r) + \sum_{x^u} logP_\theta(x^u)$$

Solved iteratively

$$\boxed{P_\theta(x^u) = P_\theta(x^u | C_1)P(C_1) + P_\theta(x^u | C_2)P(C_2)}$$

($x^u$ can come from either C$_1$ and C$_2$)

# Low-density Separation 非黑即白
*"Black-or-white"*

# Self-training

- Given: labelled data set = $\{(x^r, \hat{y}^r)\}_{r=1}^R$, unlabeled data set = $\{x^u\}_{u=1}^U$
- Repeat:
  - Train model $f^*$ from labelled data set

    You can use any model here. | Regression?
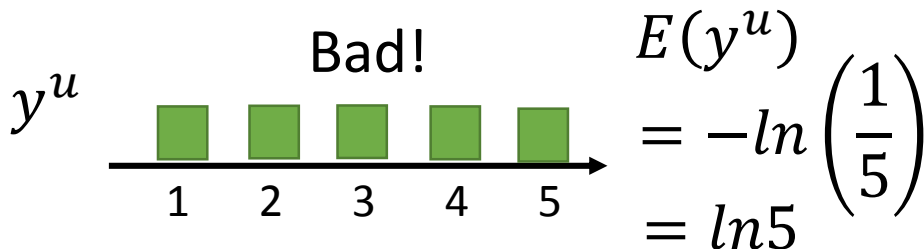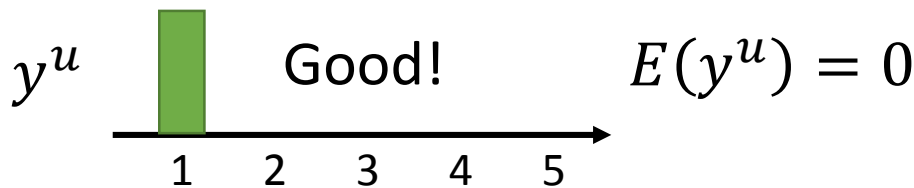
  - Apply $f^*$ to the unlabeled data set
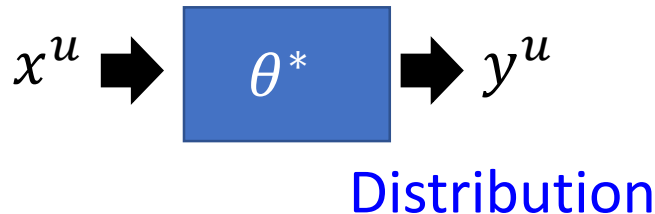    - Obtain $\{(x^u, y^u)\}_{u=1}^U$    Pseudo-label
  - Remove <u>a set of data</u> from unlabeled data set, and add them into the labeled data set

    How to choose the data set remains open | You can also provide a weight to each data.

# Entropy-based Regularization



$x^u \rightarrow \boxed{\theta^*} \rightarrow y^u$

Distribution

Entropy of $y^u$ :
Evaluate how concentrate the distribution $y^u$ is

$y^u$   Good!   $E(y^u) = 0$

1   2   3   4   5

$y^u$   Good!   $E(y^u) = 0$

1   2   3   4   5

$y^u$   Bad!

1   2   3   4   5

$E(y^u)$
$= -ln\left(\frac{1}{5}\right)$
$= ln5$

$$E(y^u) = -\sum_{m=1}^{5} y_m^u ln(y_m^u)$$

As small as possible

$$L = \sum_{x^r} C(y^r, \hat{y}^r) \quad \boxed{\text{labelled data}}$$

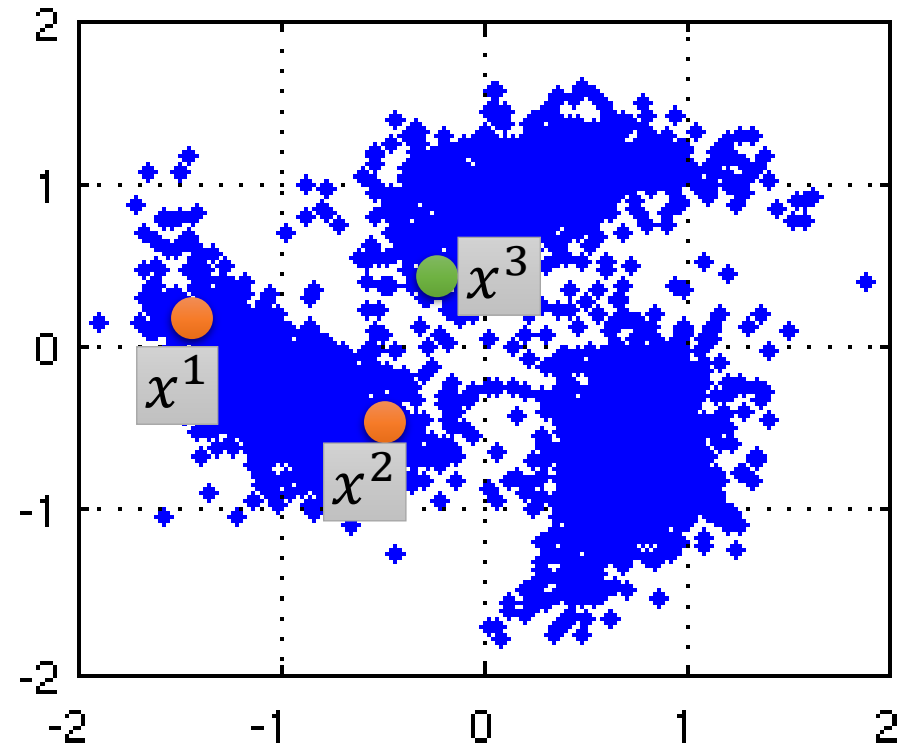$$+\lambda \sum_{x^u} E(y^u) \quad \boxed{\text{unlabeled data}}$$

# Smoothness Assumption

- Assumption: "similar" $x$ has the same $\hat{y}$

- More precisely:
  - x is not uniform.
  - If $x^1$ and $x^2$ are close in <u>a high density region,</u> $\hat{y}^1$ and $\hat{y}^2$ are the same.
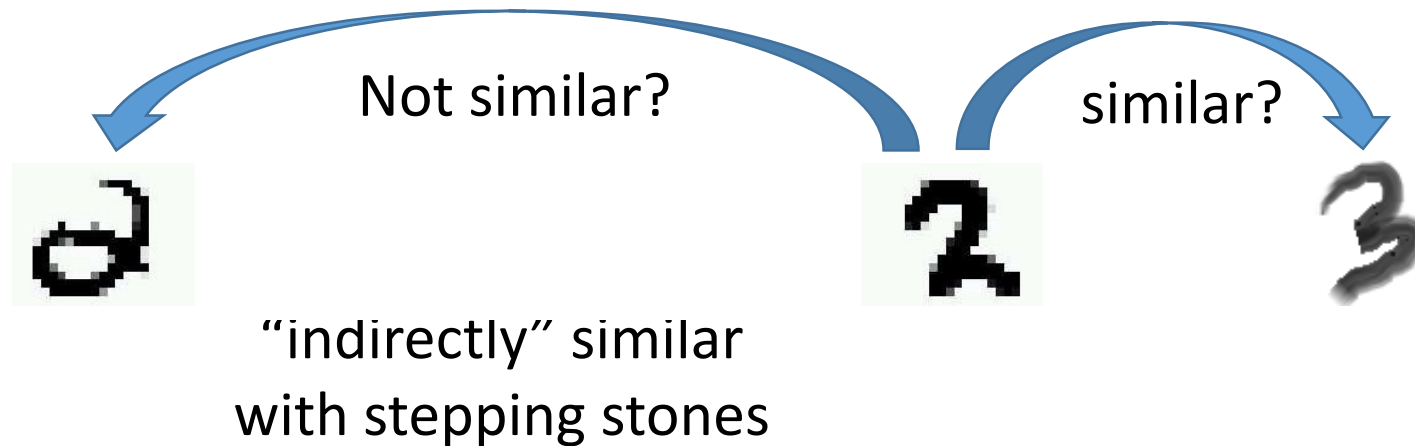
> connected by a
> high density path

Source of image:
http://hips.seas.harvard.edu/files
/pinwheel.png



$x^1$ and $x^2$ have the same label

$x^2$ and $x^3$ have different labels

# Smoothness Assumption

Not similar?          similar?

"indirectly" similar
with stepping stones

(The example is from the tutorial slides of Xiaojin Zhu.)

正侧面          正侧面

Source of image: http://www.moehui.com/5833.html/5/
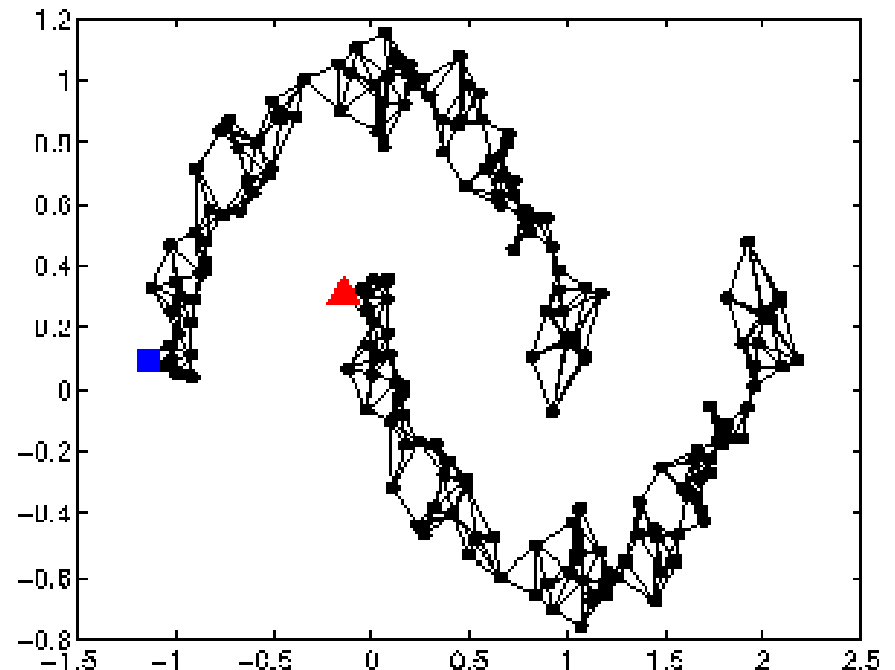
# Graph-based Approach

- How to know $x^1$ and $x^2$ are connected by a high density path

Represented the data points as a ***graph***

Graph representation is nature sometimes.

E.g. Hyperlink of webpages, citation of papers

Sometimes you have to construct the graph yourself.
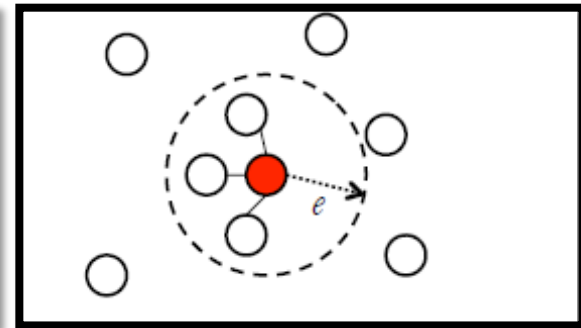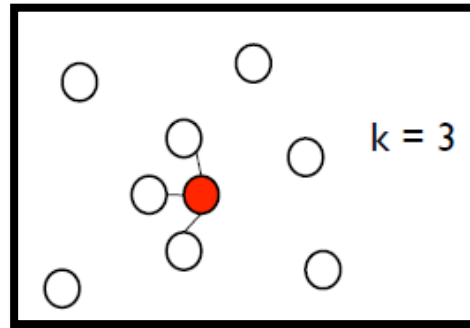
# Graph-based Approach - Graph Construction

- Define the similarity $s(x^i, x^j)$ between $x^i$ and $x^j$
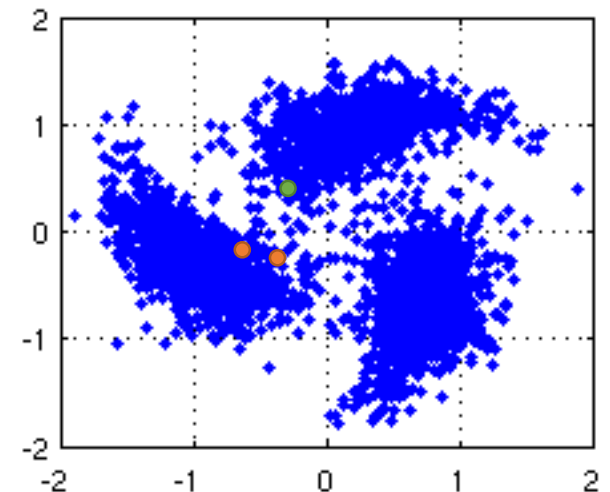
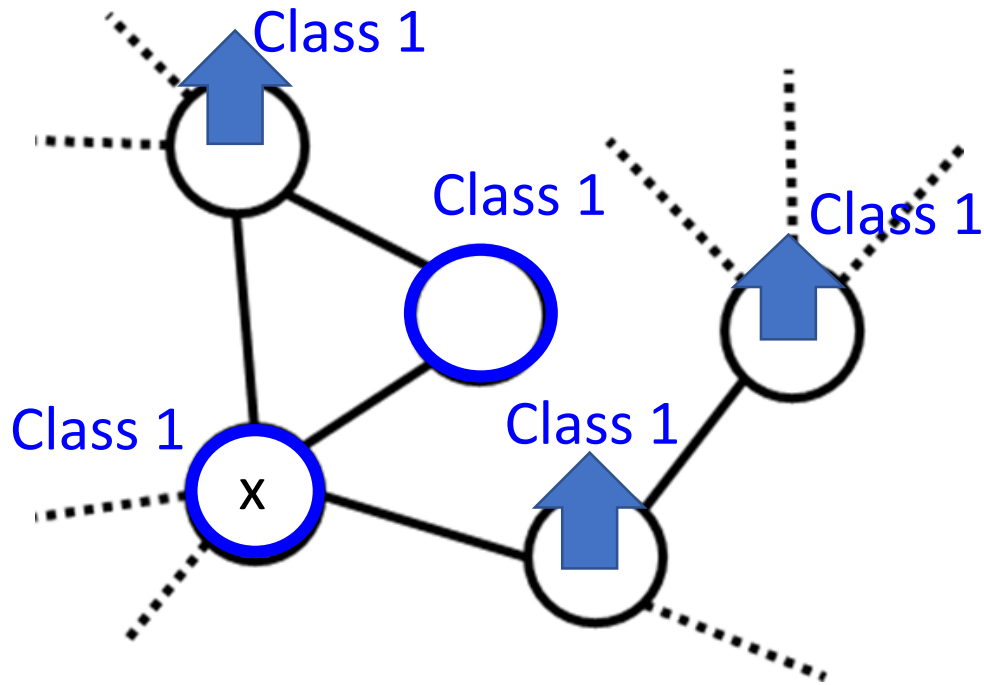- Add edge:
  - K Nearest Neighbor

  - $e$-Neighborhood



- Edge weight is proportional to $s(x^i, x^j)$

Gaussian Radial Basis Function:
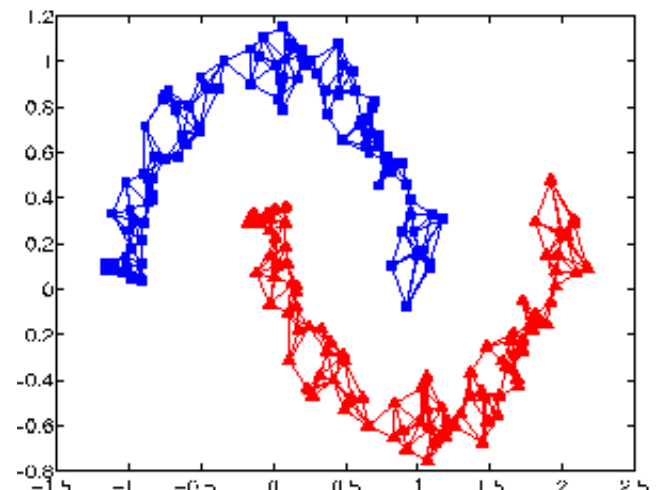$$s(x^i, x^j) = exp\left(-\gamma \|x^i - x^j\|^2\right)$$
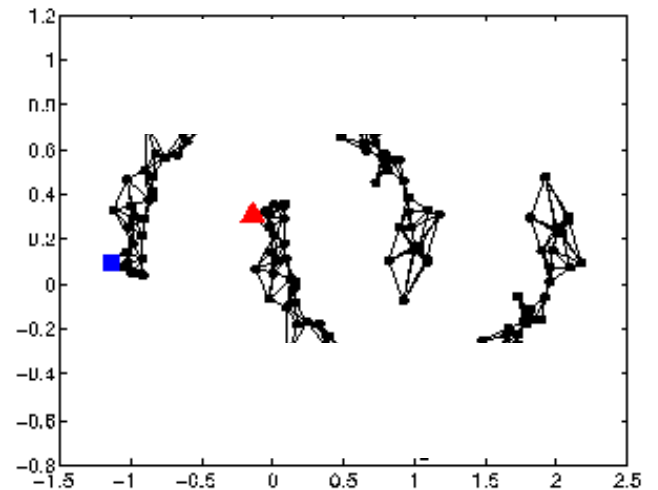
# Graph-based Approach



The labelled data influence their neighbors.
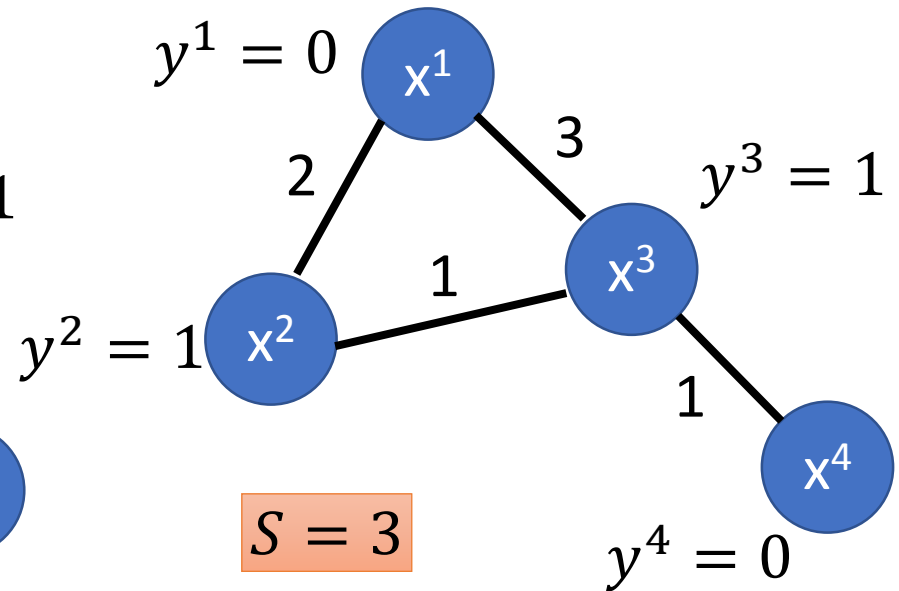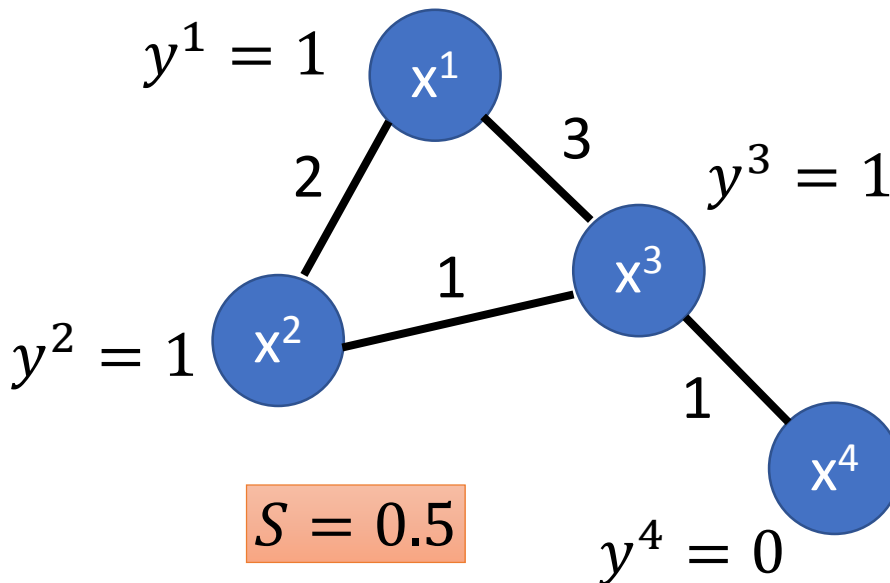
Propagate through the graph

# Graph-based Approach

- Define the smoothness of the labels on the graph

$$S = \frac{1}{2} \sum_{i,j} w_{i,j} (y^i - y^j)^2$$

Smaller means smoother

For all data (no matter labelled or not)



$y^1 = 1$

$x^1$

$2$

$3$

$y^3 = 1$

$x^3$

$1$

$x^2$

$y^2 = 1$

$1$

$x^4$

$S = 0.5$

$y^4 = 0$

$y^1 = 0$

$x^1$

$2$

$3$

$y^3 = 1$

$x^3$

$1$

$x^2$

$y^2 = 1$

$1$

$x^4$

$S = 3$

$y^4 = 0$

# Graph-based Approach

- Define the smoothness of the labels on the graph

$$S = \frac{1}{2}\sum_{i,j} w_{i,j}\left(y^i - y^j\right)^2 = \boldsymbol{y}^T L \boldsymbol{y}$$

**y**: (R+U)-dim vector

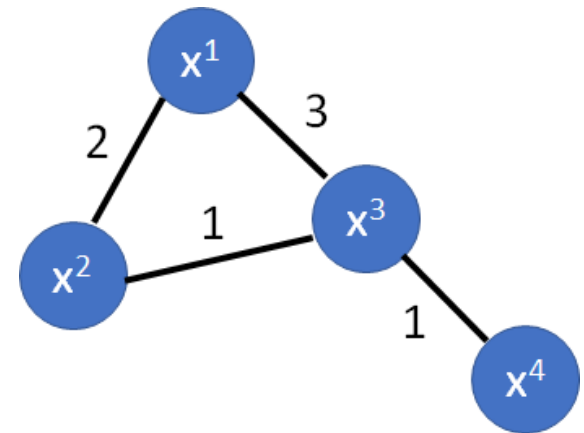$$\boldsymbol{y} = \left[\cdots y^i \cdots y^j \cdots\right]^T$$

L: (R+U) x (R+U) matrix

Graph Laplacian

$$L = \underline{D} - \underline{W}$$

$$W = \begin{vmatrix} 0 & 2 & 3 & 0 \\ 2 & 0 & 1 & 0 \\ 3 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{vmatrix}
\quad D = \begin{vmatrix} 5 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 5 & 0 \\ 0 & 0 & 0 & 1 \end{vmatrix}$$

# Graph-based Approach

- Define the smoothness of the labels on the graph

$$S = \frac{1}{2} \sum_{i,j} w_{i,j} (y^i - y^j)^2 = \boldsymbol{y}^T L \boldsymbol{y}$$

Depending on model parameters

$$L = \sum_{x^r} C(y^r, \hat{y}^r) \boxed{+\lambda S}$$

As a regularization term

J. Weston, F. Ratle, and R. Collobert, "Deep learning via semi-supervised embedding," ICML, 2008