

机器学习

09 半监督学习

李祎

liyi@dlut.edu.cn

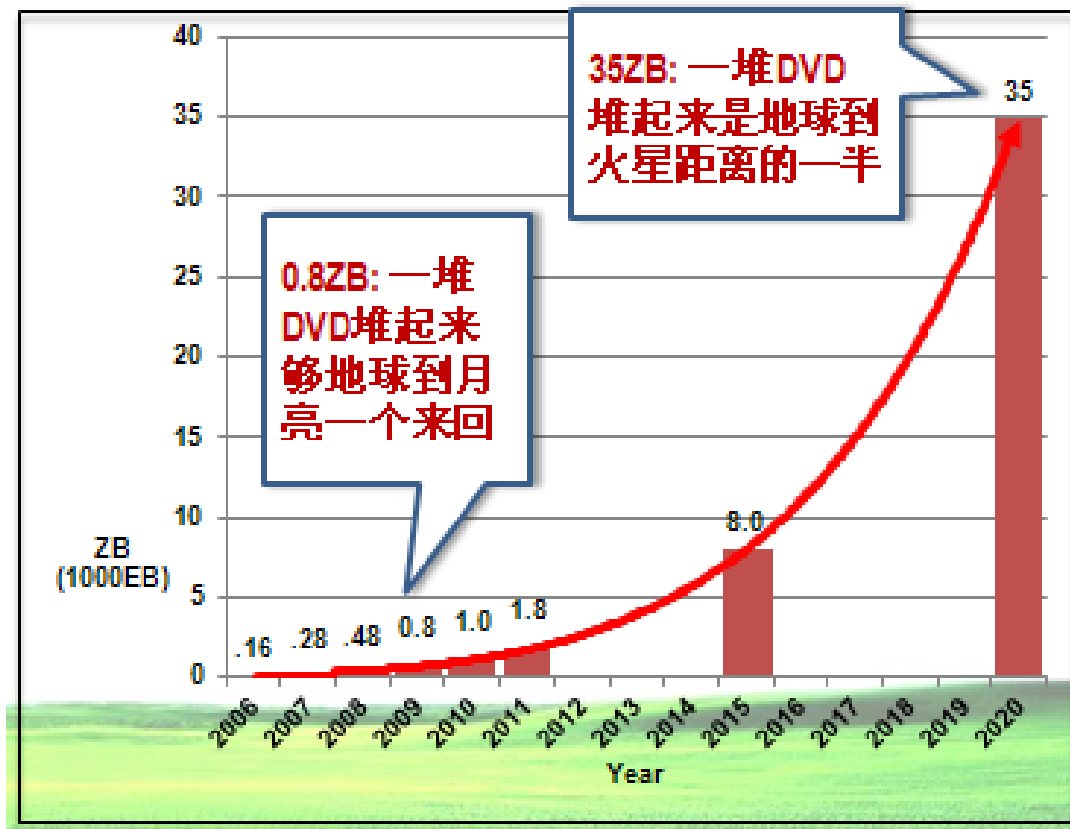
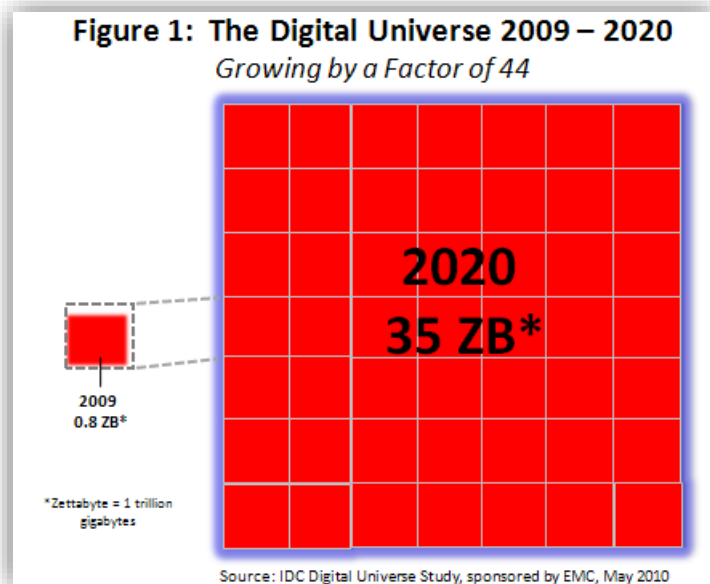


大连理工大学 人工智能学院

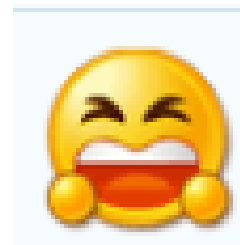
School of Artificial Intelligence, Dalian University of Technology

- 未标记样本与半监督学习
- 基于单学习器的方法：半监督SVM
- 基于分歧的方法：协同训练
- 半监督聚类

□ Big data era, obtaining data is getting **easier and easier**.



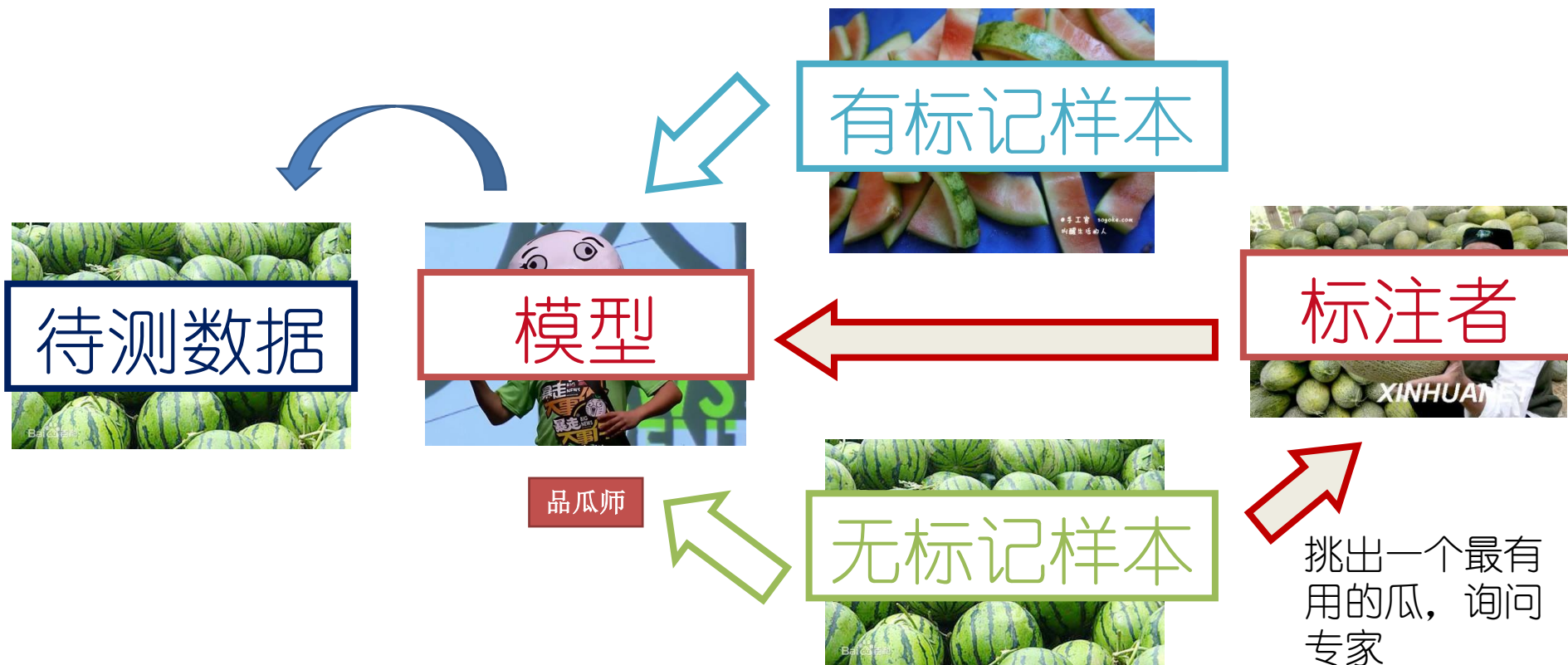
- Labeling the data is difficult! **expensive, time consuming,** sometimes need **experts.**



- For example, medical image analysis, webpage recommendation.

- How to use **few** labeled data and **large amount** of unlabeled data ?

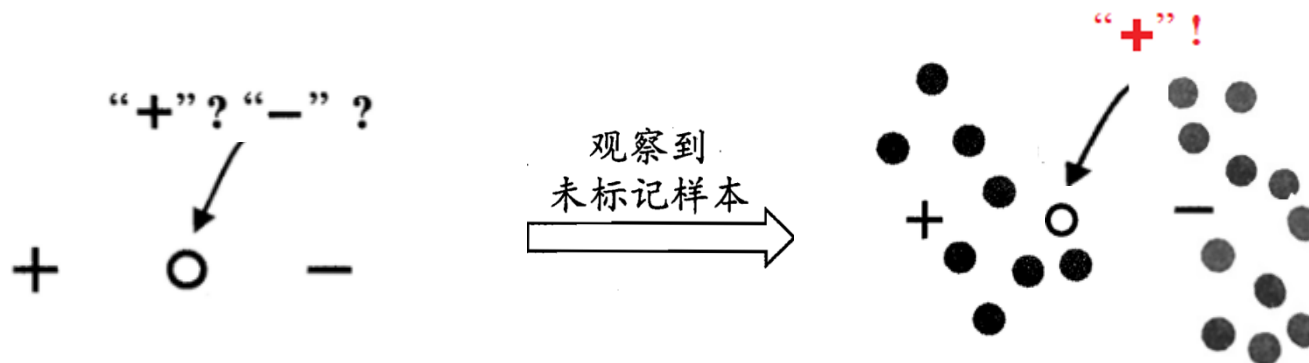




- 主动学习引入了额外的专家知识，通过与外界交互将部分未标记样本变为有标记样本。
- 若不与专家交互，还能利用未标记样本吗？

- **Usefulness of unlabeled data**

➤ Can we **use unlabeled data** to improve the learning ability?



未标记样本效用的例示. 右边的灰色点表示未标记样本

The distribution of the unlabeled data tell us ***something***.

□ 要利用未标记样本，必然要做一些将未标记样本所揭示的数据分布信息与类别标记相联系的假设，其中有两种常见的假设。

- 聚类假设 (clustering assumption) :

假设数据存在簇结构，同一簇的样本属于同一类别。

- 流形假设 (manifold assumption) :

假设数据分布在一个流形结构上，邻近的样本具有相似的输出值。

流形假设可看做聚类假设的推广

半监督学习



大连理工大学 人工智能学院
School of Artificial Intelligence, Dalian University of Technology

(纯) 半监督学习 基于“开放世界”假设

待测数据

模型

品瓜师

有标记样本

无标记样本

直推学习 基于“封闭世界”假设

Transductive learning

■ Definition

Give a training dataset

$D_l = \{(x_1, y_1), (x_2, y_2), \dots, (x_l, y_l)\}$, $y_i (i=1 \dots l)$ is the label ,

and dataset $D_u = \{x_{l+1}, x_{l+2}, \dots, x_{l+u}\}$, $l \ll u$.

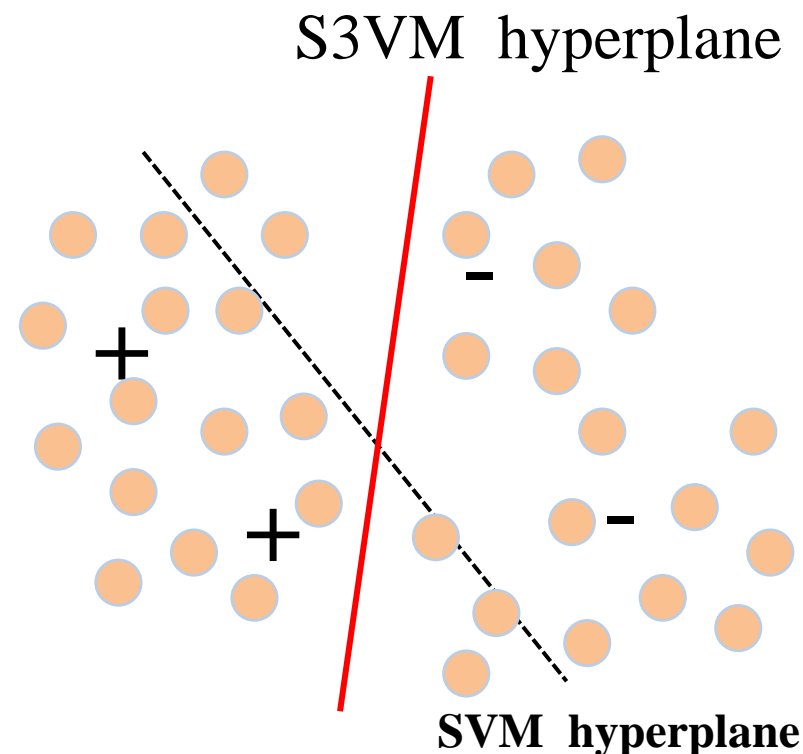
Semi-supervised learning (SSL) is learning a function $f(x)$ by **using both D_l and D_u** , and then **improve** (at least not reduce) the classification performance.

Semi-supervised learning is **halfway** between supervised and unsupervised learning.

半监督SVM

- Semi-supervised SVM tries to find a hyper-plane that can divide the different class data and goes through the lowest density area.

- Transductive SVM



● Unlabeled data

□ TSVM(Transductive Support Vector Machine)

□ 针对二分类问题，TSVM试图考虑对未标记样本进行各种可能的标记指派，然后在所有结果中，寻求一个在所有样本上间隔最大化的划分超平面。

□ 数学模型

$$\min_{\mathbf{w}, b, \hat{\mathbf{y}}, \boldsymbol{\xi}} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 + C_l \sum_{i=1}^l \xi_i + C_u \sum_{i=l+1}^m \xi_i$$

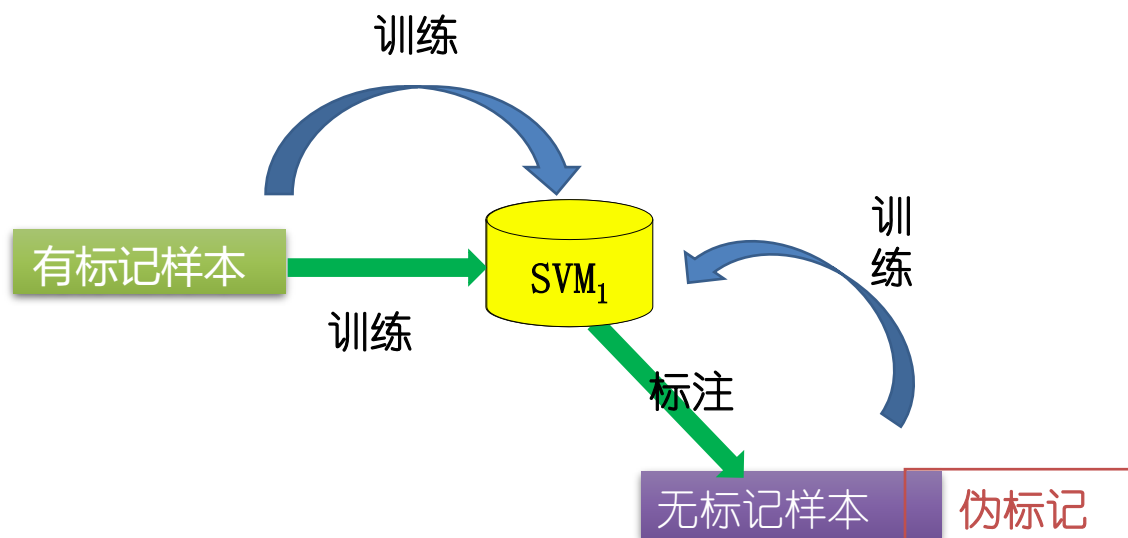
loss item of
Unlabeled
data

$$\text{s.t.} \quad y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i, \quad i = 1, \dots, l,$$

$$\hat{y}_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i, \quad i = l + 1, \dots, m,$$

$$\xi_i \geq 0, \quad i = 1, \dots, m,$$

- 显然，尝试未标记样本的各种标记指派是一个穷举过程。
- 在一般情形下，TSVM采用局部搜索来迭代地寻找近似解。



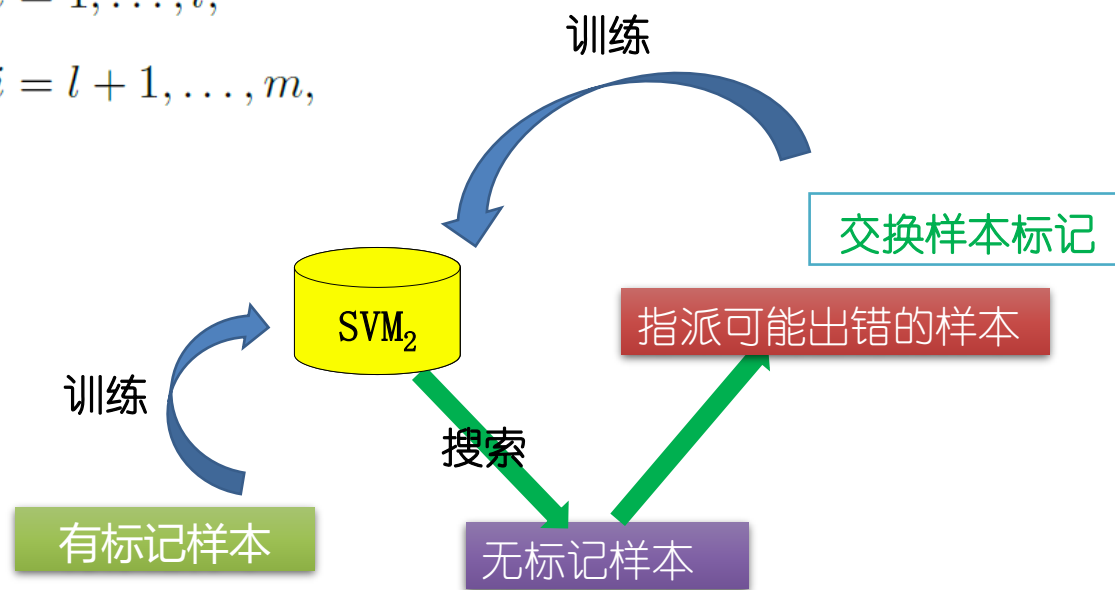
- 注意，开始时未标记样本的伪标记很可能不准确，因此设置 $C_u \ll C_l$ ，迭代过程中，逐渐增大 C_u 以提高未标记样本的作用。

$$\min_{w, b, \hat{y}, \xi} \frac{1}{2} \|w\|_2^2 + C_l \sum_{i=1}^l \xi_i + C_u \sum_{i=l+1}^m \xi_i$$

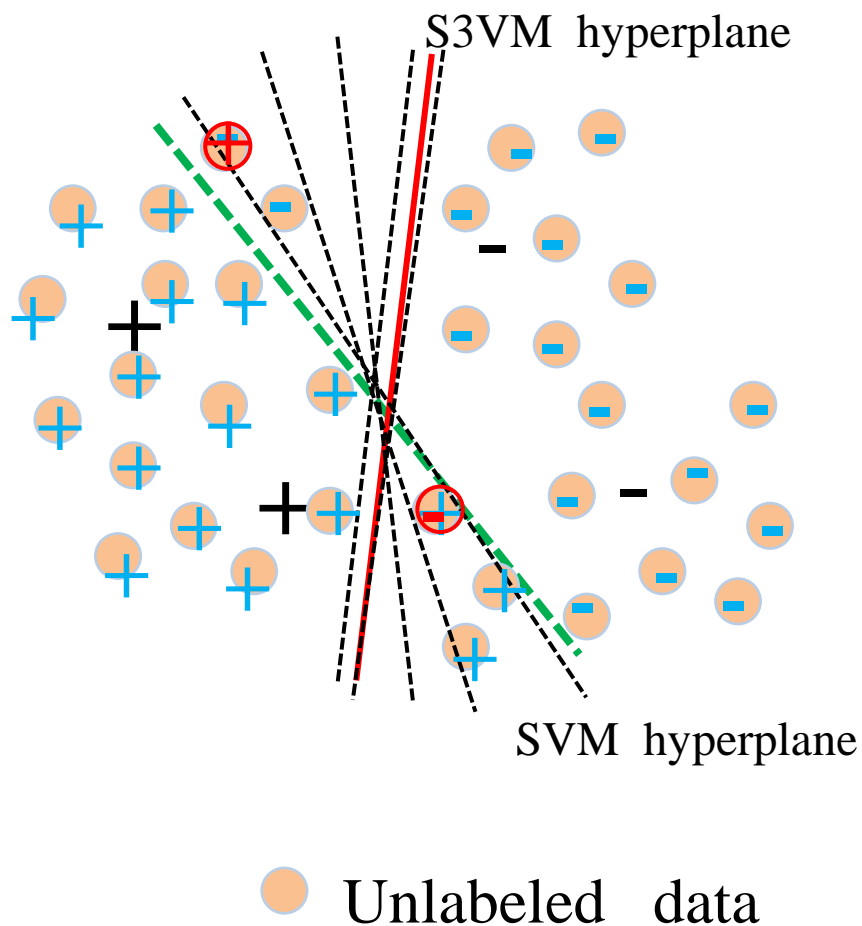
$$\text{s.t. } y_i(w^\top x_i + b) \geq 1 - \xi_i, \quad i = 1, \dots, l,$$

$$\hat{y}_i(w^\top x_i + b) \geq 1 - \xi_i, \quad i = l+1, \dots, m,$$

$$\xi_i \geq 0, \quad i = 1, \dots, m,$$



TSVM algorithm



Step 1 : train SVM on labeled data

→ Step 2 : predict unlabeled data by using SVM

Step 3 : train SVM on label data and **pseudo-label** data

Step 4 : find a pair of data which **most wrongly** predicted and **swap** their labels

Step 5 : goto step2

输入: 有标记样本集 $D_l = \{(x_1, y_1), (x_2, y_2), \dots, (x_l, y_l)\}$;
未标记样本集 $D_u = \{x_{l+1}, x_{l+2}, \dots, x_{l+u}\}$;
折中参数 C_l, C_u .

过程:

1: 用 D_l 训练一个 SVM_l ;

2: 用 SVM_l 对 D_u 中样本进行预测, 得到 $\hat{y} = (\hat{y}_{l+1}, \hat{y}_{l+2}, \dots, \hat{y}_{l+u})$;

未标记样本的
伪标记不准确

3: 初始化 $C_u \ll C_l$;

4: while $C_u < C_l$ do

5: 基于 $D_l, D_u, \hat{y}, C_l, C_u$ 求解式(13.9), 得到 $(w, b), \xi$;

6: while $\exists \{i, j \mid (\hat{y}_i \hat{y}_j < 0) \wedge (\xi_i > 0) \wedge (\xi_j > 0) \wedge (\xi_i + \xi_j > 2)\}$ do

7: $\hat{y}_i = -\hat{y}_i$;

8: $\hat{y}_j = -\hat{y}_j$;

9: 基于 $D_l, D_u, \hat{y}, C_l, C_u$ 重新求解式(13.9), 得到 $(w, b), \xi$

10: end while

11: $C_u = \min\{2C_u, C_l\}$

12: end while

输出: 未标记样本的预测结果: $\hat{y} = (\hat{y}_{l+1}, \hat{y}_{l+2}, \dots, \hat{y}_{l+u})$

图 13.4 TSVM 算法

- 显然, 搜寻标记指派可能出错的每一对未标记样本进行调整, 仍是一个涉及巨大计算开销的大规模优化问题。
- 因此, 半监督SVM研究的一个重点是如何设计出高效的优化求解策略。
- 例如基于图核(graph kernel)函数梯度下降的Laplacian SVM[Chapelle and Zien, 2005]、基于标记均值估计的meanS3VM[Li et al., 2009]等。

基于分歧的方法

- 基于分歧的方法(disagreement-based methods)使用多学习器，而学习器之间的“分歧”对未标记数据的利用至关重要。
- 协同训练(co-training)[Blum and Mitchell, 1998]是基于分歧的方法的重要代表，它最初是针对“多视图”(multi-view)数据设计的，因此也被看作“多视图学习”(multi-view learning)的代表。

Multi-view



大连理工大学 人工智能学院
School of Artificial Intelligence, Dalian University of Technology

- web pages
- driverless vehicle



Data from different sensors

滚动 | 赛程 | 赛事新闻 | 前方 | 图片 | 我要上全运 | 体育

View 1: text



1 2 3 4 5

江苏3-0上海夺冠 张常宁拥抱教练

8月28日, 2017年全运会女排赛在天津市人民体育馆战至收官日, 由蔡斌挂帅的联赛冠军江苏女排, 直落三局以3-0击败上海最终折桂, 这是自1959年以来创历史首夺全运会冠军……

头条新闻

全运女排-江苏3-0夺历史首冠 季军赛看台爆冲突(图)

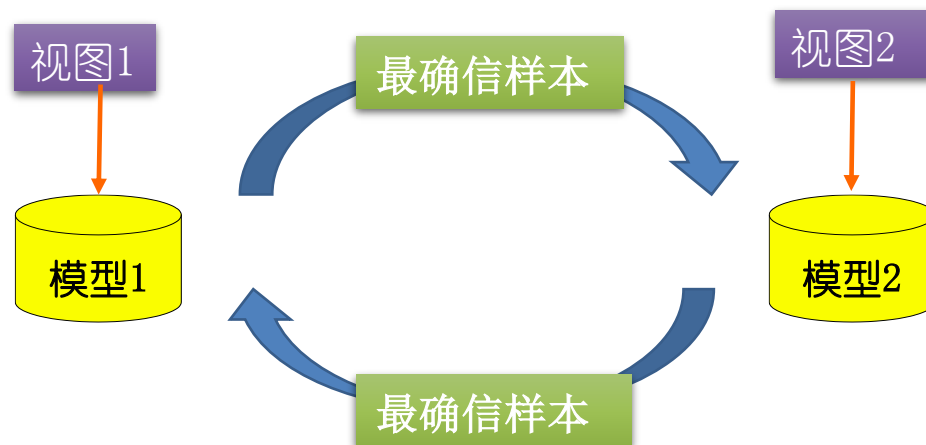
【回顾】-全运综合:七夕“狗粮”满天飞“联姻”赛艇夺首金 | 全运女子曲棍球四川险胜夺冠
【跳水】-广东队包揽全运男子跳水五金 陈艾森超600分登顶 | 三米板施廷懋摘金 成就个人大满贯
【女排】-北京女排时隔34年重返全运前3 | 北京女排赛后喜悦流泪 | 名将魏秋月宣布退役
【乒乓】-马龙时隔34年摘牌 | 独揽两冠马龙霸气 | 马龙自创“马龙式”离开
【总结】-名将观察:徐云丽重伤 | 五大老将告别全运 | 最大黑马京沪两队 | 4大副攻谁能入郎平法眼?
【伤情】-郎平爱将全运会恐怖扭伤 女排大冠军怀妊容生变 | 组图:徐云丽表情痛苦送医治疗
【举重】-女子举重48公斤级 湖南小将摘金 | 廖秋云夺女子举重53KG冠军 世界冠军黎雅君摘银

View 2: hyperlinks

调查: 您最关注全运会哪个比赛项目?

- 足球
- 篮球
- 排球
- 乒乓球
- 羽毛球
- 游泳
- 网球
- 体操
- 其他项目

- 协同训练正是很好地利用了多视图的“相容互补性”，假设数据拥有两个“充分” (sufficient) 且“条件独立”视图，每个属性集都足以描述该类样本，且相互独立。



Co-training trains **two classifiers** separately on **two views**, then uses the **predictions of each classifier** on unlabeled examples to augment the training set of **the other**.

x_i 的上标仅用于指代两个视图, 不表示序关系, 即 $\langle x_i^1, x_i^2 \rangle$ 与 $\langle x_i^2, x_i^1 \rangle$ 表示的是同一个样本.

令 $p, n \ll s$.

初始化每个视图上的有标记训练集.

在视图 j 上用有标记样本训练 h_j .

扩充有标记数据集.

输入: 有标记样本集 $D_l = \{(\langle x_1^1, x_1^2 \rangle, y_1), \dots, (\langle x_l^1, x_l^2 \rangle, y_l)\}$;
未标记样本集 $D_u = \{(\langle x_{l+1}^1, x_{l+1}^2 \rangle), \dots, (\langle x_{l+u}^1, x_{l+u}^2 \rangle)\}$;
缓冲池大小 s ;
每轮挑选的正例数 p ;
每轮挑选的反例数 n ;
基学习算法 \mathcal{L} ;
学习轮数 T .

过程:

```
1: 从  $D_u$  中随机抽取  $s$  个样本构成缓冲池  $D_s$ ;  
2:  $D_u = D_u \setminus D_s$ ;  
3: for  $j = 1, 2$  do  
4:    $D_l^j = \{(\langle x_i^j, y_i \rangle) \mid (\langle x_i^j, x_i^{3-j} \rangle, y_i) \in D_l\}$ ;  
5: end for  
6: for  $t = 1, 2, \dots, T$  do
```

```
7:   for  $j = 1, 2$  do  
8:      $h_j \leftarrow \mathcal{L}(D_l^j)$ ;  
9:     考察  $h_j$  在  $D_s^j = \{(\langle x_i^j, x_i^{3-j} \rangle) \mid (\langle x_i^j, x_i^{3-j} \rangle, y_i) \in D_s\}$  上的分类置信度, 挑选  $p$  个正例  
       置信度最高的样本  $D_p \subset D_s$ 、 $n$  个反例置信度最高的样本  $D_n \subset D_s$ ;  
10:    由  $D_p^j$  生成伪标记正例  $\tilde{D}_p^{3-j} = \{(\langle x_i^{3-j}, +1 \rangle) \mid x_i^j \in D_p^j\}$ ;  
11:    由  $D_n^j$  生成伪标记反例  $\tilde{D}_n^{3-j} = \{(\langle x_i^{3-j}, -1 \rangle) \mid x_i^j \in D_n^j\}$ ;  
12:     $D_s = D_s \setminus (D_p \cup D_n)$ ;  
13:   end for
```

```
14:   if  $h_1, h_2$  均未发生改变 then  
15:     break  
16:   else
```

```
17:     for  $j = 1, 2$  do  
18:        $D_l^j = D_l^j \cup (\tilde{D}_p^j \cup \tilde{D}_n^j)$ ;  
19:     end for
```

```
20:     从  $D_u$  中随机抽取  $2p + 2n$  个样本加入  $D_s$   
21:   end if  
22: end for
```

输出: 分类器 h_1, h_2

图 13.6 协同训练算法

- 协同训练过程虽简单，但令人惊讶的是，理论证明显示出，若两个视图充分且条件独立，则可利用未标记样本通过协同训练将弱分类器的泛化性能提升到任意高[Blum and Mitchell, 1998].
- 不过，视图的条件独立性在现实任务中通常很难满足，因此性能提升幅度不会那么大，但研究表明，即使在更弱的条件下，协同训练仍可有效地提升弱分类器的性能。

- 协同训练算法本身是为多视图数据而设计的，性集合的常见数据但此后出现了一些能在单视图数据上使用的变体算法。
- 它们或是使用不同的学习算法[Goldman and Zhou,2000]、或使用不同的数据采样[Zhou and Li, 2005b]、甚至使用不同的参数设置[Zhou and Li, 2005a]来产生不同的学习器，也能有效地利用未标记数据来提升性能。
- 后续理论研究发现，此类算法事实上无需数据拥有多视图，仅需弱学习器之间具有显著的分歧(或差异)，即可通过相互提供伪标记样本的方式来提高泛化性能[周志华, 2013]。

半监督聚类

- 聚类是一种典型的无监督学习任务，然而在现实聚类任务中我们往往能获得一些额外的监督信息，于是可通过“半监督聚类” (semi-supervised clustering) 来利用监督信息以获得更好的聚类效果。

- 聚类任务中获得的监督信息大致有两种类型：
 - 第一种类型是“必连” (must-link) 与“勿连” (cannot-link) 约束，前者是指样本必属于同一个簇，后者则是指样本必不属于同一个簇；
 - 第二种类型的监督信息则是少量的有标记样本。

- 约束 k 均值(Constrained k -means)算法[Wagstaff et al., 2001]是利用第一类监督信息的代表。
- 该算法是 k 均值算法的扩展,它在聚类过程中要确保“必连”关系集合与“勿连”关系集合中的约束得以满足,否则将返回错误提示。

输入: 样本集 $D = \{x_1, x_2, \dots, x_m\}$;
必连约束集合 \mathcal{M} ;
勿连约束集合 \mathcal{C} ;
聚类簇数 k .

过程:

```
1: 从  $D$  中随机选取  $k$  个样本作为初始均值向量  $\{\mu_1, \mu_2, \dots, \mu_k\}$ ;  
2: repeat  
3:    $C_j = \emptyset$  ( $1 \leq j \leq k$ );  
4:   for  $i = 1, 2, \dots, m$  do  
5:     计算样本  $x_i$  与各均值向量  $\mu_j$  ( $1 \leq j \leq k$ ) 的距离:  $d_{ij} = \|x_i - \mu_j\|_2$ ;  
6:      $\mathcal{K} = \{1, 2, \dots, k\}$ ;  
7:     is_merged=false;  
8:     while  $\neg$  is_merged do  
9:       基于  $\mathcal{K}$  找出与样本  $x_i$  距离最近的簇:  $r = \arg \min_{j \in \mathcal{K}} d_{ij}$ ;  
10:      检测将  $x_i$  划入聚类簇  $C_r$  是否会违背  $\mathcal{M}$  与  $\mathcal{C}$  中的约束;  
11:      if  $\neg$  is_violated then  
12:         $C_r = C_r \cup \{x_i\}$ ;  
13:        is_merged=true  
14:      else  
15:         $\mathcal{K} = \mathcal{K} \setminus \{r\}$ ;  
16:        if  $\mathcal{K} = \emptyset$  then  
17:          break并返回错误提示  
18:        end if  
19:      end if  
20:    end while  
21:  end for  
22:  for  $j = 1, 2, \dots, k$  do  
23:     $\mu_j = \frac{1}{|C_j|} \sum_{x \in C_j} x$ ;  
24:  end for  
25: until 均值向量均未更新  
输出: 簇划分  $\{C_1, C_2, \dots, C_k\}$ 
```

初始化 k 个空簇.

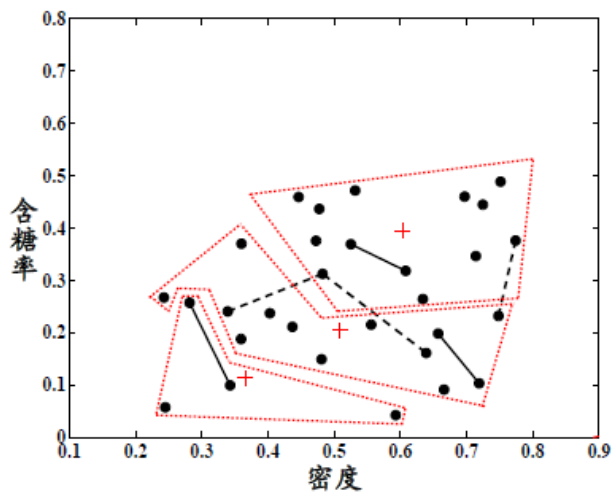
更新均值向量.

不冲突, 选择最近的簇

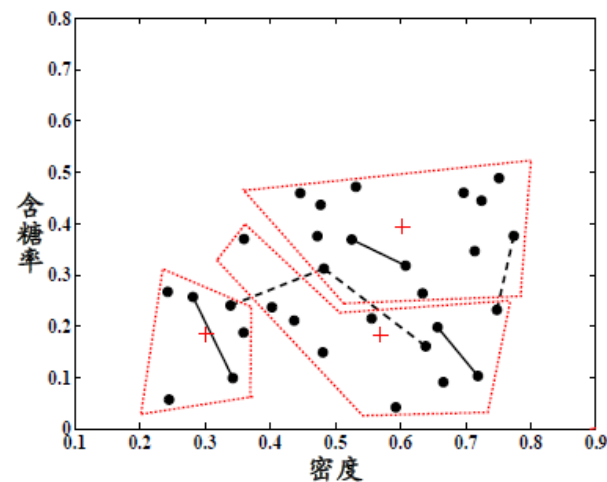
冲突, 尝试次近的簇

图 13.7 约束 k 均值算法

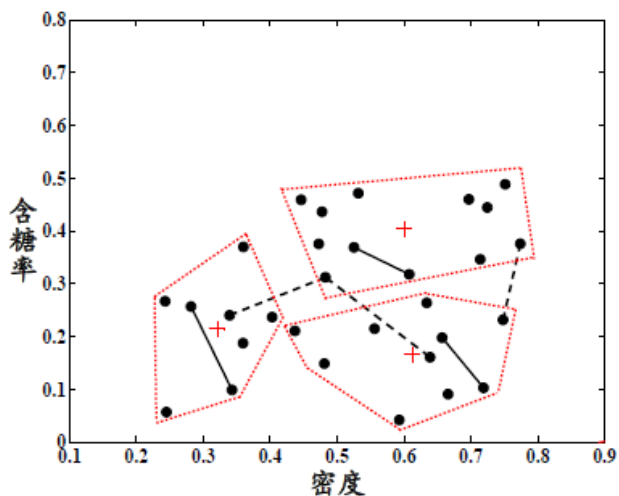
半监督聚类



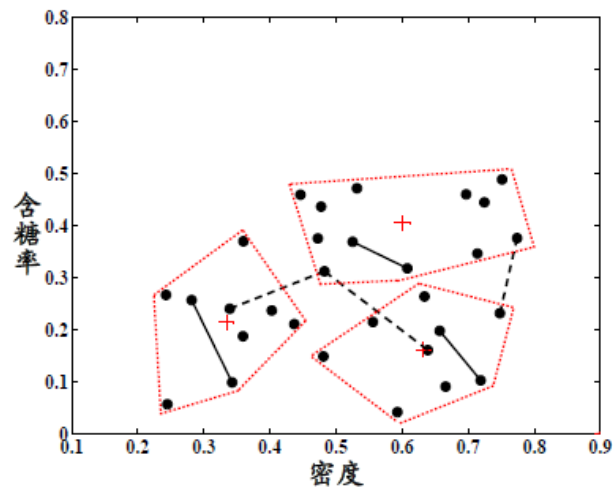
(a) 第 1 轮迭代后



(b) 第 2 轮迭代后



(c) 第 3 轮迭代后



(d) 第 4 轮迭代后

- 第二种监督信息是少量有标记样本，即假设少量有标记样本属于 k 个聚类簇。
- 这样的监督信息利用起来很容易：直接将它们作为“种子”，用它们初始化 k 均值算法的 k 个聚类中心，并且在聚类簇迭代更新过程中不改变种子样本的簇隶属关系。这样就得到了约束种子 k 均值 (Constrained Seed k -means) 算法[Basu et al., 2002]。

$S \subset D, |S| \ll |D|$.

输入: 样本集 $D = \{x_1, x_2, \dots, x_m\}$;
少量有标记样本 $S = \bigcup_{j=1}^k S_j$;
聚类簇数 k .

过程:

用有标记样本初始化簇
中心.

```
1: for  $j = 1, 2, \dots, k$  do  
2:    $\mu_j = \frac{1}{|S_j|} \sum_{x \in S_j} x$   
3: end for
```

用有标记样本初始化 k
个簇.

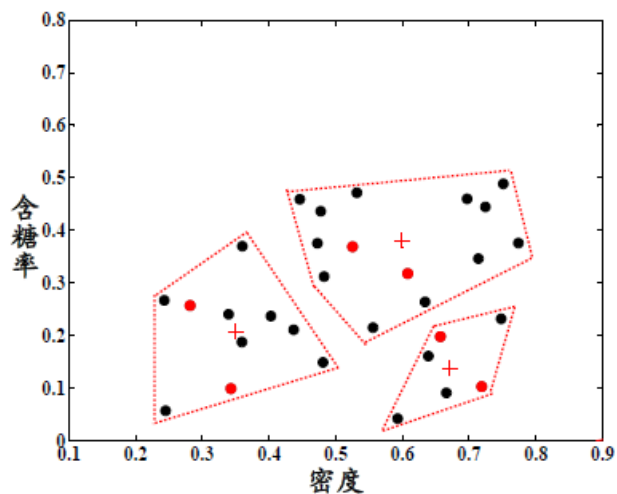
```
4: repeat  
5:    $C_j = \emptyset$  ( $1 \leq j \leq k$ );  
6:   for  $j = 1, 2, \dots, k$  do  
7:     for all  $x \in S_j$  do  
8:        $C_j = C_j \cup \{x\}$   
9:     end for
```

```
10:  end for  
11:  for all  $x_i \in D \setminus S$  do  
12:    计算样本  $x_i$  与各均值向量  $\mu_j$  ( $1 \leq j \leq k$ ) 的距离:  $d_{ij} = \|x_i - \mu_j\|_2$ ;  
13:    找出与样本  $x_i$  距离最近的簇:  $r = \arg \min_{j \in \{1, 2, \dots, k\}} d_{ij}$ ;  
14:    将样本  $x_i$  划入相应的簇:  $C_r = C_r \cup \{x_i\}$   
15:  end for  
16:  for  $j = 1, 2, \dots, k$  do  
17:     $\mu_j = \frac{1}{|C_j|} \sum_{x \in C_j} x$ ;  
18:  end for  
19: until 均值向量均未更新  
输出: 簇划分  $\{C_1, C_2, \dots, C_k\}$ 
```

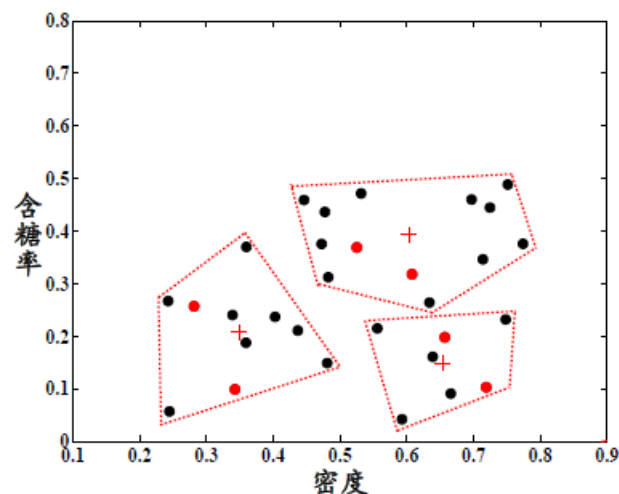
更新均值向量.

图 13.9 约束种子 k 均值算法

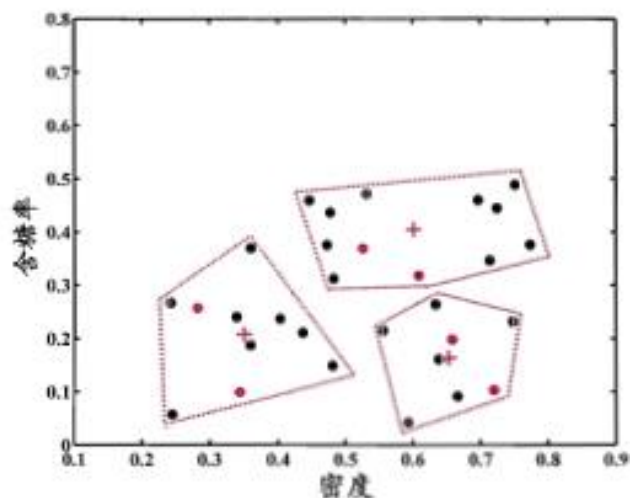
半监督聚类



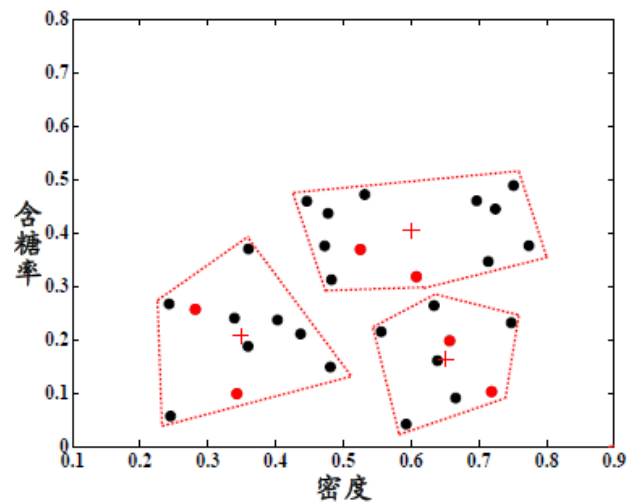
(a) 第 1 轮迭代后



(b) 第 2 轮迭代后



(c) 第 3 轮迭代后



(d) 第 4 轮迭代后

- 未标记样本与半监督学习
- 基于单学习器的方法：半监督SVM
- 基于分歧的方法：协同训练
- 半监督聚类：两种监督信息