



# Random Matrix Analysis to Balance between Supervised and Unsupervised Learning

**Vasilii Feofanov\***, Malik Tiomoko\*, Aladin Virmaux\*

Huawei Paris Noah's Ark Lab  
firstname.lastname@huawei.com  
\*equal contribution

CAp RFIAP 2024, July 2  
Published in ICML 2023



In some applications data acquisition is cheaper than labeling,



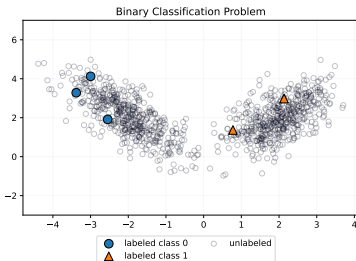
And supervised learning is inefficient.



**Semi-supervised learning:** learn with both few labeled and many unlabeled training examples.

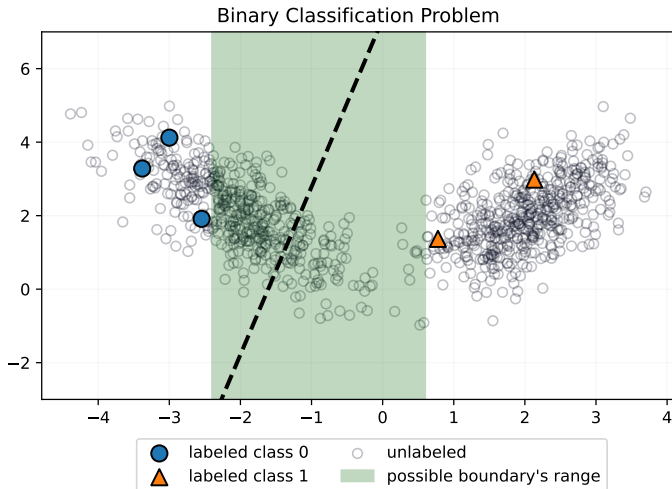
Family of SSL Methods:

- **Pseudo-labeling,**
- Graph-based algorithms,
- Cluster-then-label,
- Unsupervised feature learning.



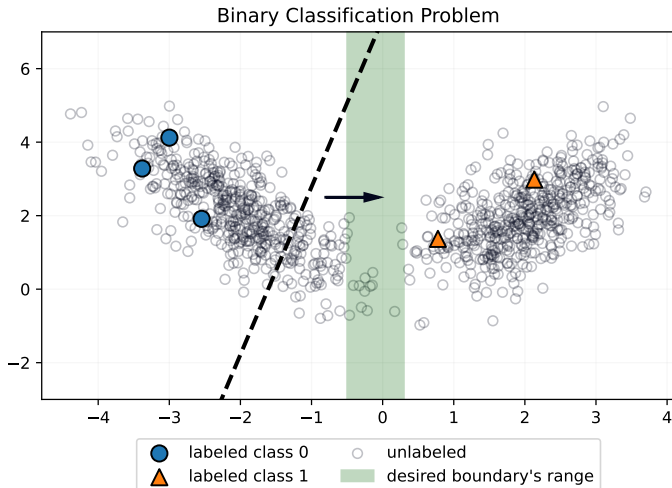


Range of possible supervised classifiers is vast: we need to make assumptions.



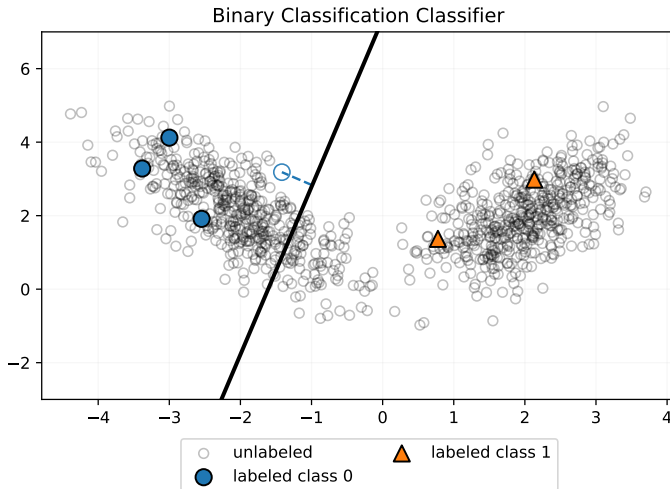


Low Density Separation (LDS) assumption: decision boundary is far away from dense regions of unlabeled data.





Implementation of LDS: push the boundary away from unlabeled data with high confident scores.

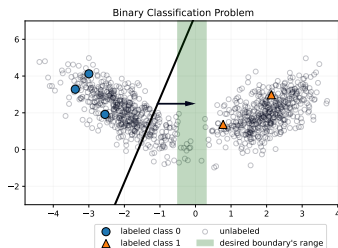




We introduce the QLDS classifier defined by solving:

$$\operatorname{argmin}_{\omega} \underbrace{\frac{\alpha_{\ell}}{2} \sum_{i=1}^{n_{\ell}} \left( y_i - \frac{\omega^{\top} \mathbf{x}_i}{\sqrt{n}} \right)^2}_{\text{labeled loss}} - \underbrace{\frac{\alpha_u}{2} \sum_{i=n_{\ell}+1}^{n_{\ell}+n_u} \left( \frac{\omega^{\top} \mathbf{x}_i}{\sqrt{n}} \right)^2}_{\text{unlabeled LDS loss}} + \underbrace{\frac{\lambda}{2} \|\omega\|^2}_{\text{regularization}}$$

- **Linear classification:**  
for  $\mathbf{x}$ , output  $\operatorname{sign}(\omega^{\top} \mathbf{x})$ ;



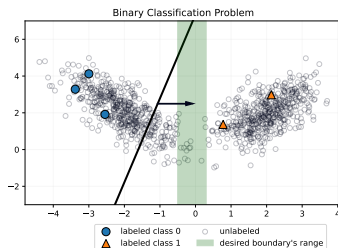




We introduce the QLDS classifier defined by solving:

$$\operatorname{argmin}_{\omega} \underbrace{\frac{\alpha_{\ell}}{2} \sum_{i=1}^{n_{\ell}} \left( y_i - \frac{\omega^{\top} \mathbf{x}_i}{\sqrt{n}} \right)^2}_{\text{labeled loss}} - \underbrace{\frac{\alpha_u}{2} \sum_{i=n_{\ell}+1}^{n_{\ell}+n_u} \left( \frac{\omega^{\top} \mathbf{x}_i}{\sqrt{n}} \right)^2}_{\text{unlabeled LDS loss}} + \underbrace{\frac{\lambda}{2} \|\omega\|^2}_{\text{regularization}}$$

- Linear classification:  
for  $\mathbf{x}$ , output  $\operatorname{sign}(\omega^{\top} \mathbf{x})$ ;
- Square margin maximization;

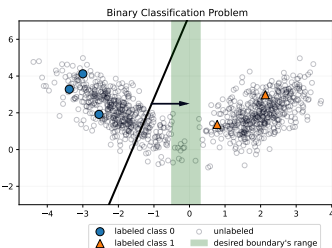




We introduce the QLDS classifier defined by solving:

$$\operatorname{argmin}_{\omega} \underbrace{\frac{\alpha_{\ell}}{2} \sum_{i=1}^{n_{\ell}} \left( y_i - \frac{\omega^{\top} \mathbf{x}_i}{\sqrt{n}} \right)^2}_{\text{labeled loss}} - \underbrace{\frac{\alpha_u}{2} \sum_{i=n_{\ell}+1}^{n_{\ell}+n_u} \left( \frac{\omega^{\top} \mathbf{x}_i}{\sqrt{n}} \right)^2}_{\text{unlabeled LDS loss}} + \underbrace{\frac{\lambda}{2} \|\omega\|^2}_{\text{regularization}}$$

- Linear classification:  
for  $\mathbf{x}$ , output  $\operatorname{sign}(\omega^{\top} \mathbf{x})$ ;
- Square margin maximization;
- Hyperparameters  $\alpha_{\ell}$ ,  $\alpha_u$ ,  
 $\lambda$  to balance the components.





$$\operatorname{argmin}_{\omega} \underbrace{\frac{\alpha_l}{2} \sum_{i=1}^{n_l} \left( y_i - \frac{\omega^\top \mathbf{x}_i}{\sqrt{n}} \right)^2}_{\text{labeled loss}} - \underbrace{\frac{\alpha_u}{2} \sum_{i=n_l+1}^{n_l+n_u} \left( \frac{\omega^\top \mathbf{x}_i}{\sqrt{n}} \right)^2}_{\text{unlabeled LDS loss}} + \underbrace{\frac{\lambda}{2} \|\omega\|^2}_{\text{regularization}}$$

Difference between QLDS and Transductive SVM:

- Quadratic loss instead of hinge loss;



$$\operatorname{argmin}_{\omega} \underbrace{\frac{\alpha_{\ell}}{2} \sum_{i=1}^{n_{\ell}} \left( y_i - \frac{\omega^{\top} \mathbf{x}_i}{\sqrt{n}} \right)^2}_{\text{labeled loss}} - \underbrace{\frac{\alpha_u}{2} \sum_{i=n_{\ell}+1}^{n_{\ell}+n_u} \left( \frac{\omega^{\top} \mathbf{x}_i}{\sqrt{n}} \right)^2}_{\text{unlabeled LDS loss}} + \underbrace{\frac{\lambda}{2} \|\omega\|^2}_{\text{regularization}}$$

Difference between QLDS and Transductive SVM:

- Quadratic loss instead of hinge loss;
- Margin squared  $(\omega^{\top} \mathbf{x})^2$  instead of  $|\omega^{\top} \mathbf{x}|$ , regularization instead of pseudo-labeling loss;



$$\operatorname{argmin}_{\omega} \underbrace{\frac{\alpha_l}{2} \sum_{i=1}^{n_l} \left( y_i - \frac{\omega^\top \mathbf{x}_i}{\sqrt{n}} \right)^2}_{\text{labeled loss}} - \underbrace{\frac{\alpha_u}{2} \sum_{i=n_l+1}^{n_l+n_u} \left( \frac{\omega^\top \mathbf{x}_i}{\sqrt{n}} \right)^2}_{\text{unlabeled LDS loss}} + \underbrace{\frac{\lambda}{2} \|\omega\|^2}_{\text{regularization}}$$

Difference between QLDS and Transductive SVM:

- Quadratic loss instead of hinge loss;
- Margin squared  $(\omega^\top \mathbf{x})^2$  instead of  $|\omega^\top \mathbf{x}|$ , regularization instead of pseudo-labeling loss;
- TSVM loss is non-convex and difficult to optimize, QLDS loss is convex and has a **closed form solution**;



$$\operatorname{argmin}_{\omega} \underbrace{\frac{\alpha_{\ell}}{2} \sum_{i=1}^{n_{\ell}} \left( y_i - \frac{\omega^{\top} \mathbf{x}_i}{\sqrt{n}} \right)^2}_{\text{labeled loss}} - \underbrace{\frac{\alpha_u}{2} \sum_{i=n_{\ell}+1}^{n_{\ell}+n_u} \left( \frac{\omega^{\top} \mathbf{x}_i}{\sqrt{n}} \right)^2}_{\text{unlabeled LDS loss}} + \underbrace{\frac{\lambda}{2} \|\omega\|^2}_{\text{regularization}}$$

Particular cases of QLDS are:

- $\alpha_u = 0 \Rightarrow$  Least-Square SVM (supervised regime);



$$\operatorname{argmin}_{\omega} \underbrace{\frac{\alpha_l}{2} \sum_{i=1}^{n_l} \left( y_i - \frac{\omega^\top \mathbf{x}_i}{\sqrt{n}} \right)^2}_{\text{labeled loss}} - \underbrace{\frac{\alpha_u}{2} \sum_{i=n_l+1}^{n_l+n_u} \left( \frac{\omega^\top \mathbf{x}_i}{\sqrt{n}} \right)^2}_{\text{unlabeled LDS loss}} + \underbrace{\frac{\lambda}{2} \|\omega\|^2}_{\text{regularization}}$$

Particular cases of QLDS are:

- $\alpha_u = 0 \Rightarrow$  Least-Square SVM (supervised regime);
- $\alpha_l \rightarrow 0 \Rightarrow$  Graph-based SSL (Mai, X. and Couillet, R., 2018);



$$\operatorname{argmin}_{\omega} \underbrace{\frac{\alpha_l}{2} \sum_{i=1}^{n_l} \left( y_i - \frac{\omega^\top \mathbf{x}_i}{\sqrt{n}} \right)^2}_{\text{labeled loss}} - \underbrace{\frac{\alpha_u}{2} \sum_{i=n_l+1}^{n_l+n_u} \left( \frac{\omega^\top \mathbf{x}_i}{\sqrt{n}} \right)^2}_{\text{unlabeled LDS loss}} + \underbrace{\frac{\lambda}{2} \|\omega\|^2}_{\text{regularization}}$$

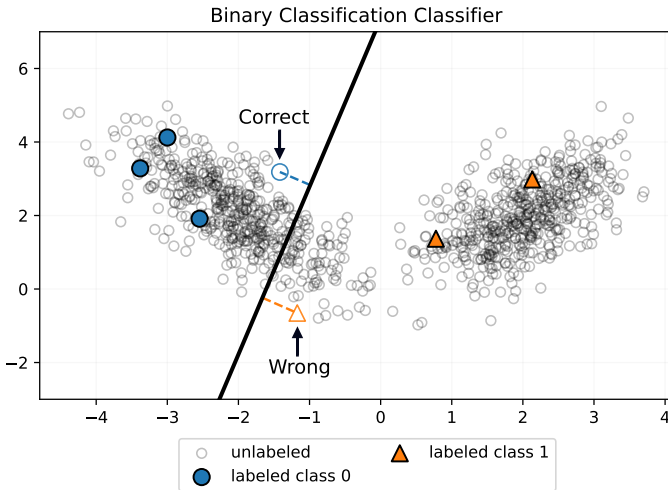
Particular cases of QLDS are:

- $\alpha_u = 0 \Rightarrow$  Least-Square SVM (supervised regime);
- $\alpha_l \rightarrow 0 \Rightarrow$  Graph-based SSL (Mai, X. and Couillet, R., 2018);
- $\alpha_l \rightarrow 0$  and  $\lambda \rightarrow \lambda_{\max}(\mathbf{X}_u) \Rightarrow$  Linear spectral clustering.





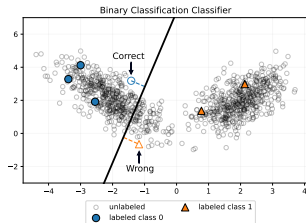
Problem: It is not safe since prediction can be wrong.





Problem: It is not safe since prediction can be wrong.

$$\frac{\alpha_\ell}{2} \sum_{i=1}^{n_\ell} \left( y_i - \frac{\omega^\top \mathbf{x}_i}{\sqrt{n}} \right)^2 - \frac{\alpha_u}{2} \sum_{i=n_\ell+1}^{n_\ell+n_u} \left( \frac{\omega^\top \mathbf{x}_i}{\sqrt{n}} \right)^2 + \frac{\lambda}{2} \|\omega\|$$



Questions:

- 1 What are generalization guarantees of the classifier?
- 2 How to properly choose the hyperparameters?



Three types of asymptotic analysis:

- Traditional:  $d$  is fixed,  $n$  is large ( $n \gg d, n \rightarrow \infty$ );



Three types of asymptotic analysis:

- Traditional:  $d$  is fixed,  $n$  is large ( $n \gg d, n \rightarrow \infty$ );
- Small-data:  $n$  is small,  $d$  is large, ( $n \ll d, d \rightarrow \infty$ );



Three types of asymptotic analysis:

- Traditional:  $d$  is fixed,  $n$  is large ( $n \gg d, n \rightarrow \infty$ );
- Small-data:  $n$  is small,  $d$  is large, ( $n \ll d, d \rightarrow \infty$ );
- Large-dimensional:  $n, d$  are both large ( $d = \mathcal{O}(n), (n, d) \rightarrow \infty$ ).



Three types of asymptotic analysis:

- Traditional:  $d$  is fixed,  $n$  is large ( $n \gg d, n \rightarrow \infty$ );
- Small-data:  $n$  is small,  $d$  is large, ( $n \ll d, d \rightarrow \infty$ );
- **Large-dimensional**:  $n, d$  are both large ( $d = \mathcal{O}(n), (n, d) \rightarrow \infty$ ).

Assumptions on data distribution:

- Gaussian Mixture Model (GMM);



Three types of asymptotic analysis:

- Traditional:  $d$  is fixed,  $n$  is large ( $n \gg d, n \rightarrow \infty$ );
- Small-data:  $n$  is small,  $d$  is large, ( $n \ll d, d \rightarrow \infty$ );
- **Large-dimensional**:  $n, d$  are both large ( $d = \mathcal{O}(n), (n, d) \rightarrow \infty$ ).

Assumptions on data distribution:

- Gaussian Mixture Model (GMM);
- **Concentrated Data** (Louart, C. and Couillet, R., 2018): variance of  $\omega^\top \mathbf{x}$  does not grow with dimension  $d$ ,



Three types of asymptotic analysis:

- Traditional:  $d$  is fixed,  $n$  is large ( $n \gg d, n \rightarrow \infty$ );
- Small-data:  $n$  is small,  $d$  is large, ( $n \ll d, d \rightarrow \infty$ );
- **Large-dimensional**:  $n, d$  are both large ( $d = \mathcal{O}(n), (n, d) \rightarrow \infty$ ).

Assumptions on data distribution:

- Gaussian Mixture Model (GMM);
- **Concentrated Data** (Louart, C. and Couillet, R., 2018): variance of  $\omega^\top \mathbf{x}$  does not grow with dimension  $d$ ,

Particular cases:

- Standard Gaussian distribution,
- Lipschitz transformation of Gaussian (e.g., GAN images),
- Open question: learned features by DNN?





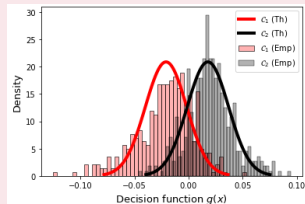
Assumptions:

- 1 Large-dimensional regime  $d = \mathcal{O}(n)$ ;
- 2 Concentrated data distribution: variance of  $\omega^\top \mathbf{x}$  does not grow with dimension  $d$ .

## Theorem

*Under the assumptions, we have:*

- $\omega^\top \mathbf{x}|y=-1$  and  $\omega^\top \mathbf{x}|y=+1$  are asymptotically normally distributed with known parameters;
- Classification error is explicitly evaluated;
- The classification problem concentrates into two-dimensional sufficient statistics.





## ■ Sketch of Proof.

- CLT: if  $\mathbf{x}$  is a concentrated random vector, then  $\boldsymbol{\omega}^\top \mathbf{x}$  is asymptotically Gaussian,
- Following Marchenko, V. A. and Pastur, L. A. (1967), for each class  $\mathcal{C}_j$ , we compute:

$$\mathbb{E}_{\mathbf{X}_\ell, \mathbf{X}_u} \left[ (\boldsymbol{\omega}^*(\mathbf{X}_\ell, \mathbf{X}_u))^\top \mathbf{x} \mid \mathbf{x} \in \mathcal{C}_j \right],$$
$$\text{Var}_{\mathbf{X}_\ell, \mathbf{X}_u} \left[ (\boldsymbol{\omega}^*(\mathbf{X}_\ell, \mathbf{X}_u))^\top \mathbf{x} \mid \mathbf{x} \in \mathcal{C}_j \right].$$

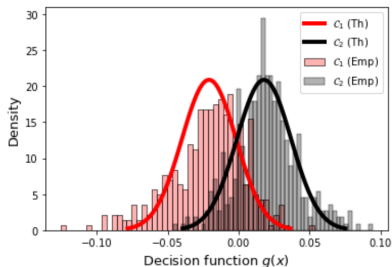


## ■ Sketch of Proof.

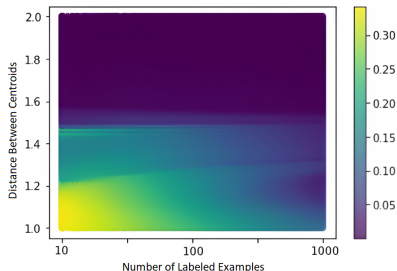
- CLT: if  $\mathbf{x}$  is a concentrated random vector, then  $\boldsymbol{\omega}^\top \mathbf{x}$  is asymptotically Gaussian,
- Following Marchenko, V. A. and Pastur, L. A. (1967), for each class  $\mathcal{C}_j$ , we compute:

$$\mathbb{E}_{\mathbf{X}_\ell, \mathbf{X}_u} \left[ (\boldsymbol{\omega}^*(\mathbf{X}_\ell, \mathbf{X}_u))^\top \mathbf{x} \mid \mathbf{x} \in \mathcal{C}_j \right],$$
$$\text{Var}_{\mathbf{X}_\ell, \mathbf{X}_u} \left[ (\boldsymbol{\omega}^*(\mathbf{X}_\ell, \mathbf{X}_u))^\top \mathbf{x} \mid \mathbf{x} \in \mathcal{C}_j \right].$$

- In the final expression, some quantities depend on class mean and covariance and need to be estimated from data.
  - For example,  $[\mu_{-1}, \mu_{+1}]^\top [\mu_{-1}, \mu_{+1}]$  is better to estimate directly rather than estimating  $\mu_{-1}, \mu_{+1}$  separately.



(a) Theoretical vs Empirical Distribution.



(b) Utility of Unlabeled Data.

- Theory can fit the empirical distribution of  $\omega^\top \mathbf{x}$ .
- The theoretical expression can be viewed as a function of different variables:  $n_\ell, \alpha_\ell, \alpha_u$ , etc.



- Find  $\alpha_\ell$  and  $\alpha_u$  automatically based on asymptotic error given by Theorem.



- Find  $\alpha_\ell$  and  $\alpha_u$  automatically based on asymptotic error given by Theorem.
- Experimental results showed that

Data set	Baselines		Model Selection	
	QLDS (1,0) (LS-SVM)	QLDS (0,1) (Graph SSL)	QLDS (cv)	QLDS (th)
books	37.47 <sup>↓</sup> ± 2.25	26.47 ± 0.72	27.91 ± 3.32	<b>26.03</b> ± 0.79
dvd	38.33 <sup>↓</sup> ± 1.72	29.12 ± 1.35	29.53 ± 3.48	<b>28.53</b> ± 1.33
electronics	34.15 <sup>↓</sup> ± 3.25	19.4 ± 0.29	20.1 <sup>↓</sup> ± 1.03	<b>19.41</b> ± 0.46
kitchen	32.39 <sup>↓</sup> ± 3.02	19.31 ± 0.16	19.98 <sup>↓</sup> ± 2.28	<b>19.11</b> ± 0.32
splice	39.81 <sup>↓</sup> ± 2.93	35.48 ± 0.86	37.02 ± 3.04	<b>35.35</b> ± 1.26
adult	33.35 ± 0.68	36.28 <sup>↓</sup> ± 0.06	<b>32.25</b> ± 1.92	32.88 ± 2.46
mushrooms	6.55 <sup>↓</sup> ± 2.07	11.33 <sup>↓</sup> ± 0.04	<b>2.57</b> ± 1.86	8.49 <sup>↓</sup> ± 3.63



- Find  $\alpha_\ell$  and  $\alpha_u$  automatically based on asymptotic error given by Theorem.
- Experimental results showed that

Data set	Baselines		Model Selection	
	QLDS (1,0) (LS-SVM)	QLDS (0,1) (Graph SSL)	QLDS (cv)	QLDS (th)
books	37.47 <sup>↓</sup> ± 2.25	26.47 ± 0.72	27.91 ± 3.32	<b>26.03</b> ± 0.79
dvd	38.33 <sup>↓</sup> ± 1.72	29.12 ± 1.35	29.53 ± 3.48	<b>28.53</b> ± 1.33
electronics	34.15 <sup>↓</sup> ± 3.25	19.4 ± 0.29	20.1 <sup>↓</sup> ± 1.03	<b>19.41</b> ± 0.46
kitchen	32.39 <sup>↓</sup> ± 3.02	19.31 ± 0.16	19.98 <sup>↓</sup> ± 2.28	<b>19.11</b> ± 0.32
splice	39.81 <sup>↓</sup> ± 2.93	35.48 ± 0.86	37.02 ± 3.04	<b>35.35</b> ± 1.26
adult	33.35 ± 0.68	36.28 <sup>↓</sup> ± 0.06	<b>32.25</b> ± 1.92	32.88 ± 2.46
mushrooms	6.55 <sup>↓</sup> ± 2.07	11.33 <sup>↓</sup> ± 0.04	<b>2.57</b> ± 1.86	8.49 <sup>↓</sup> ± 3.63

- Model selection outperforms both LS-SVM and Graph-based SSL;



- Find  $\alpha_\ell$  and  $\alpha_u$  automatically based on asymptotic error given by Theorem.
- Experimental results showed that

Data set	Baselines		Model Selection	
	QLDS (1,0) (LS-SVM)	QLDS (0,1) (Graph SSL)	QLDS (cv)	QLDS (th)
books	37.47 <sup>↓</sup> ± 2.25	26.47 ± 0.72	27.91 ± 3.32	<b>26.03</b> ± 0.79
dvd	38.33 <sup>↓</sup> ± 1.72	29.12 ± 1.35	29.53 ± 3.48	<b>28.53</b> ± 1.33
electronics	34.15 <sup>↓</sup> ± 3.25	19.4 ± 0.29	20.1 <sup>↓</sup> ± 1.03	<b>19.41</b> ± 0.46
kitchen	32.39 <sup>↓</sup> ± 3.02	19.31 ± 0.16	19.98 <sup>↓</sup> ± 2.28	<b>19.11</b> ± 0.32
splice	39.81 <sup>↓</sup> ± 2.93	35.48 ± 0.86	37.02 ± 3.04	<b>35.35</b> ± 1.26
adult	33.35 ± 0.68	36.28 <sup>↓</sup> ± 0.06	<b>32.25</b> ± 1.92	32.88 ± 2.46
mushrooms	6.55 <sup>↓</sup> ± 2.07	11.33 <sup>↓</sup> ± 0.04	<b>2.57</b> ± 1.86	8.49 <sup>↓</sup> ± 3.63

- Model selection outperforms both LS-SVM and Graph-based SSL;
- Selecting  $\alpha_\ell$  and  $\alpha_u$  by cross-validation is more costly and can lead to over-confidence towards labeled data.



Thanks for your attention !