

Sparse Learning for Noisy Data/Labels:

A Simple yet Effective Framework for Vision Applications



Yikai Wang
yikai-wang.github.io



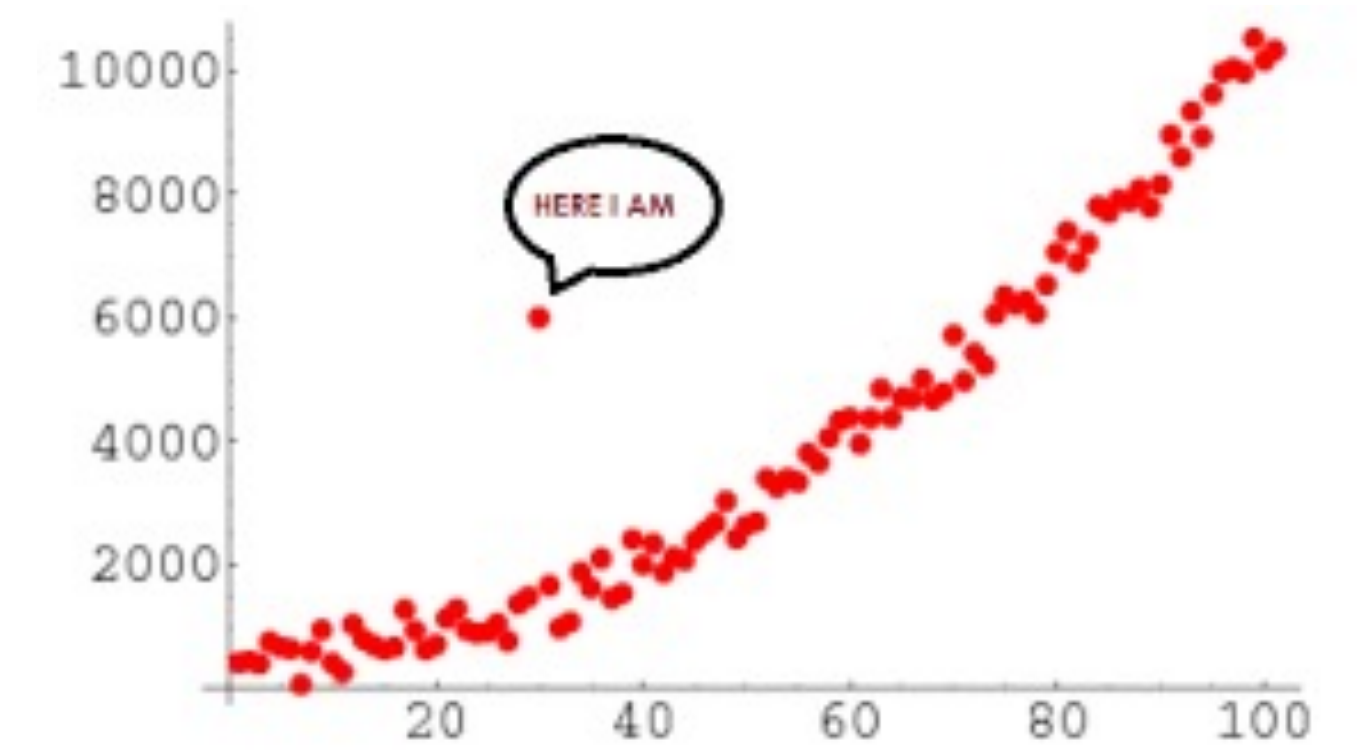
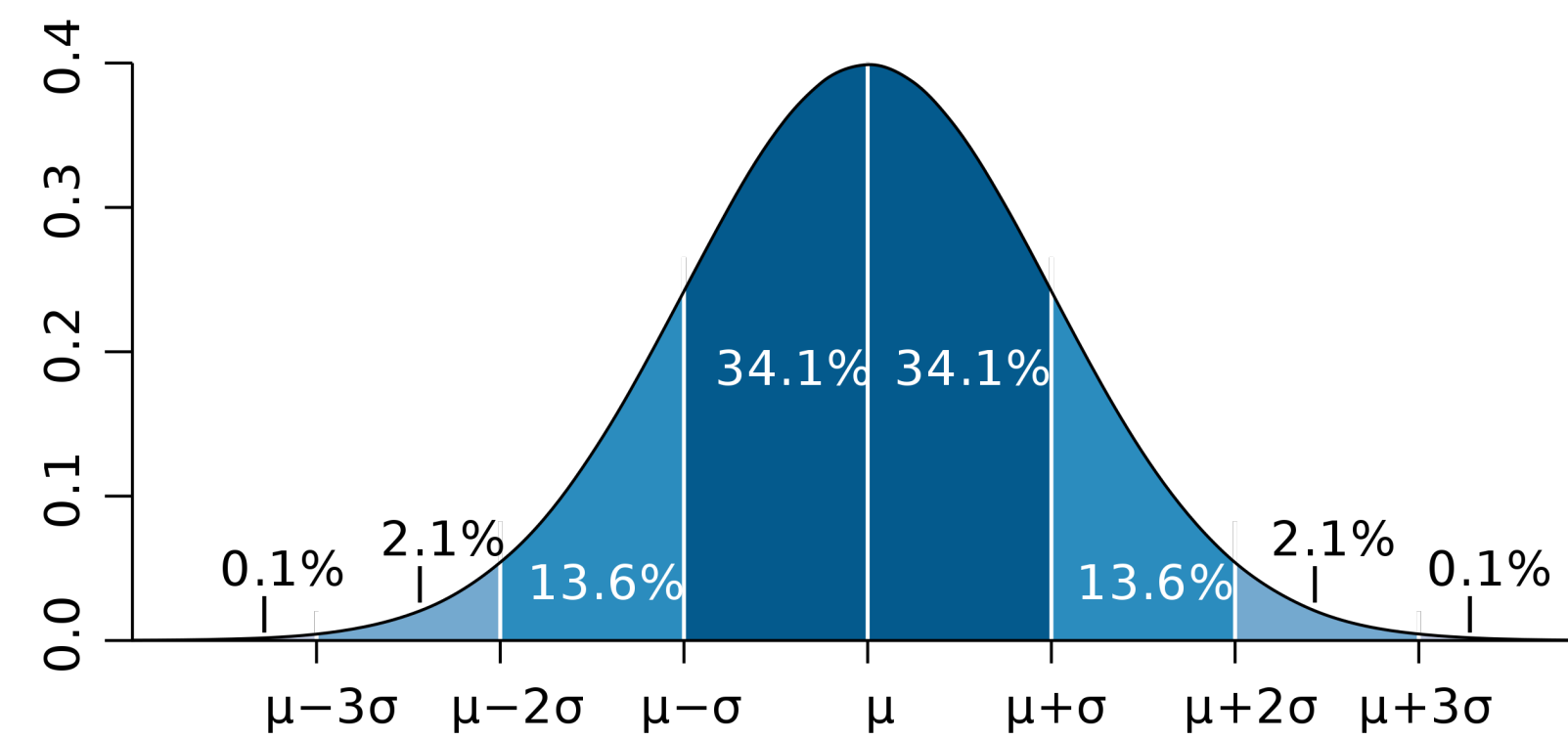
Yanwei Fu
<http://yanweifu.github.io>

School of Data Science
Fudan University

Sparse Learning

for Noisy Data Detection

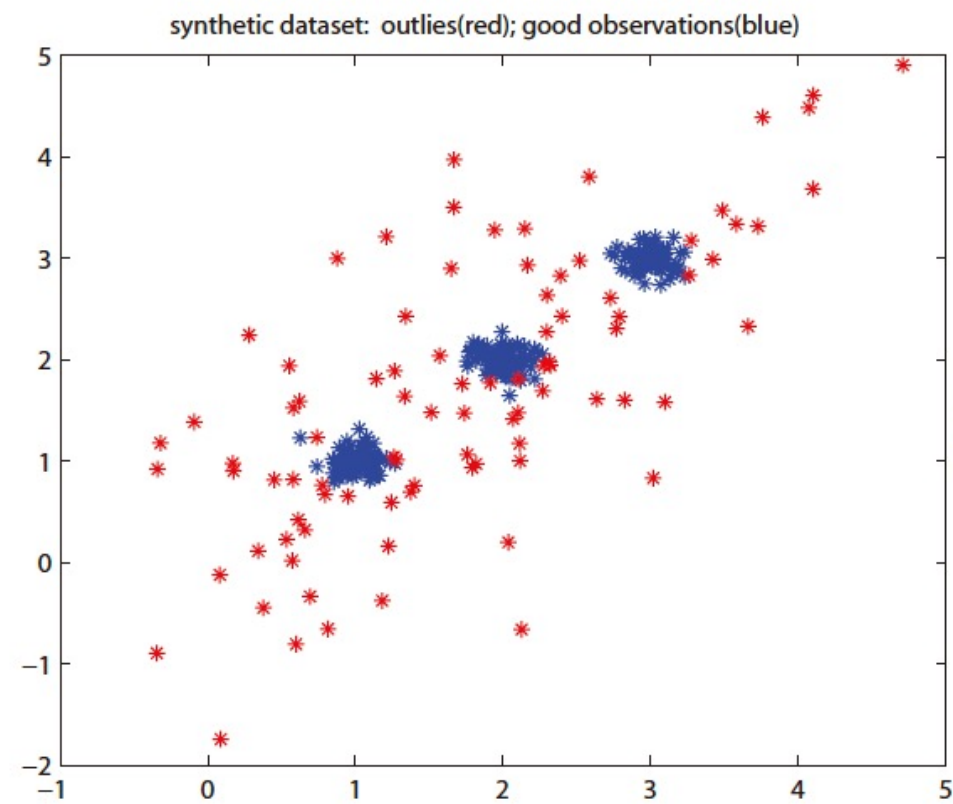
Examples of Noisy Data/Outliers



Outliers are the **irregular** data compared with the majority of the dataset.

Noisy Data in Label Space

- Random Corruptions



- Annotator mistakes



- Noisy search engine results



Shogun: Total War - IGN
ign.com



Aug 21, 1192 CE: First S...
nationalgeographic.org



Baal Ascension Materials: What To Farm For G...
forbes.com

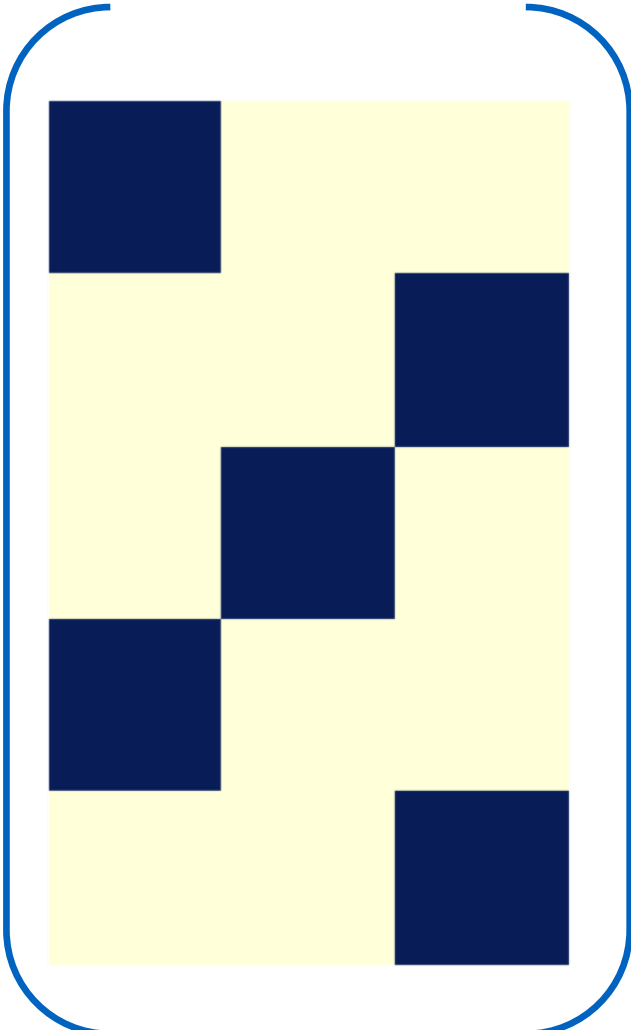
- Complex/Confusing items identified



Identify Noisy Data in Label Space

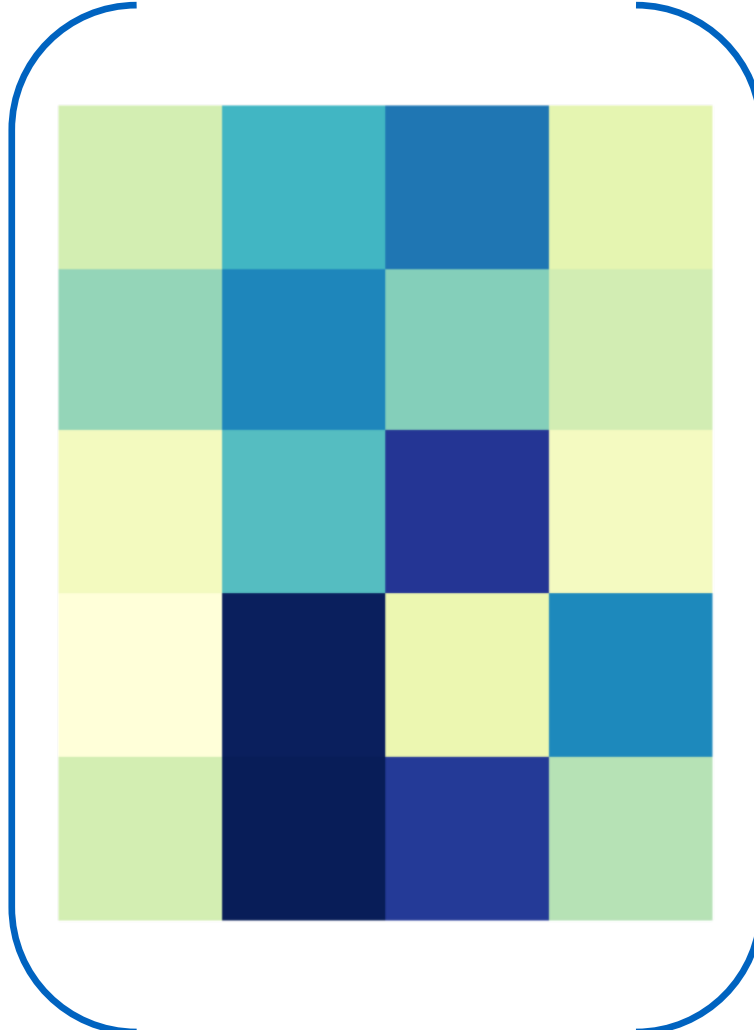
Linear system

$$Y = X\beta$$



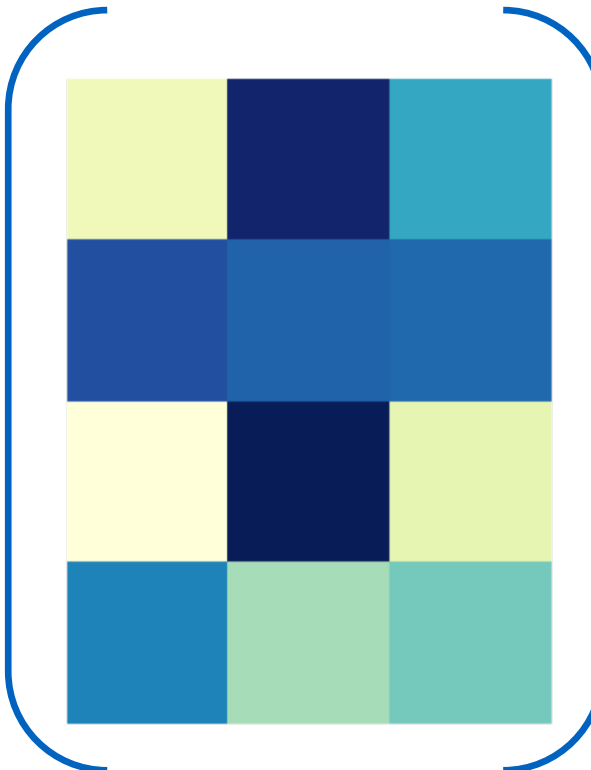
Noisy One-hot Labels

$$Y \in \mathbb{R}^{n \times c}$$



Deep Features

$$X \in \mathbb{R}^{n \times d}$$

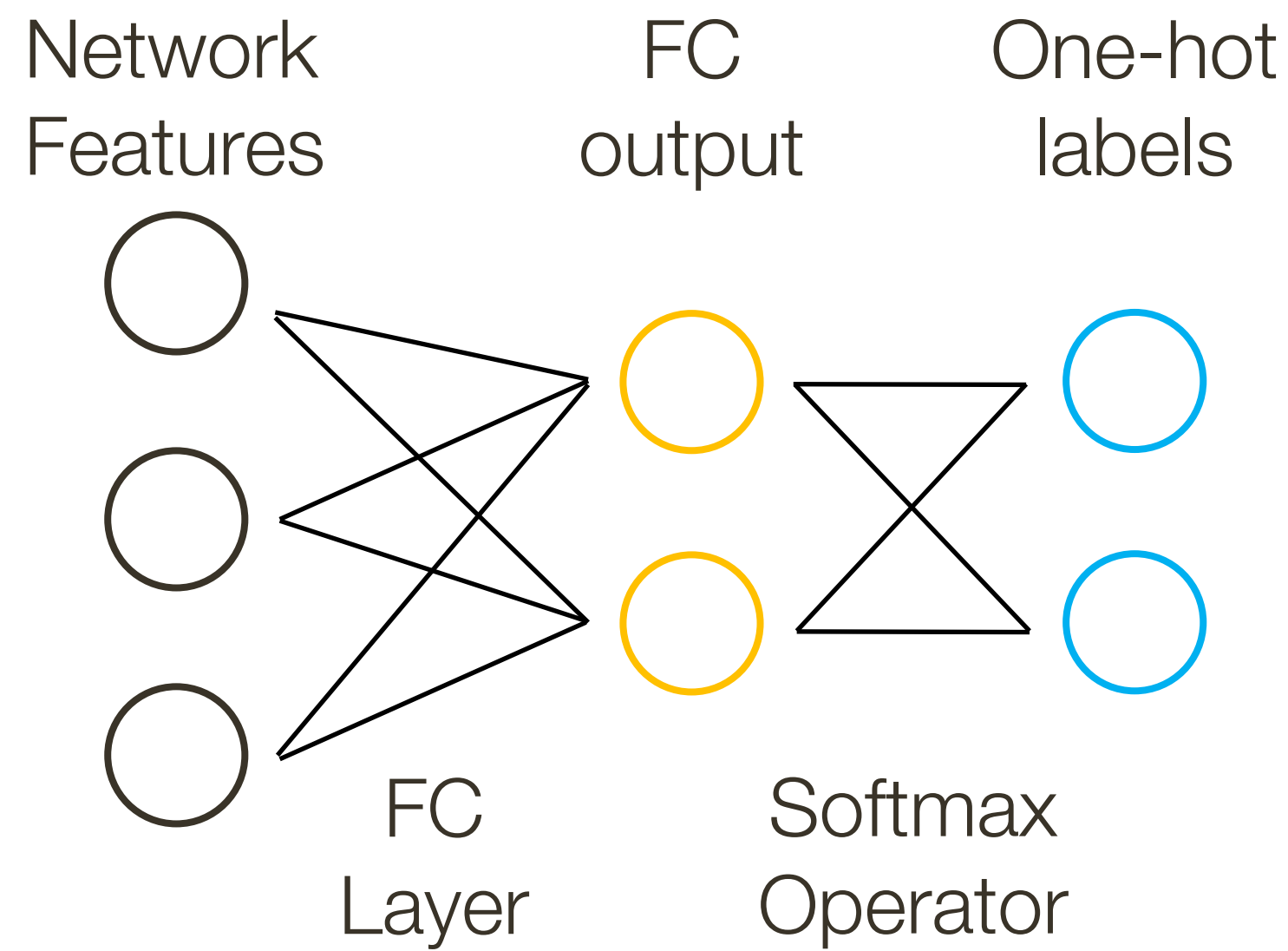


Fitted Coef.

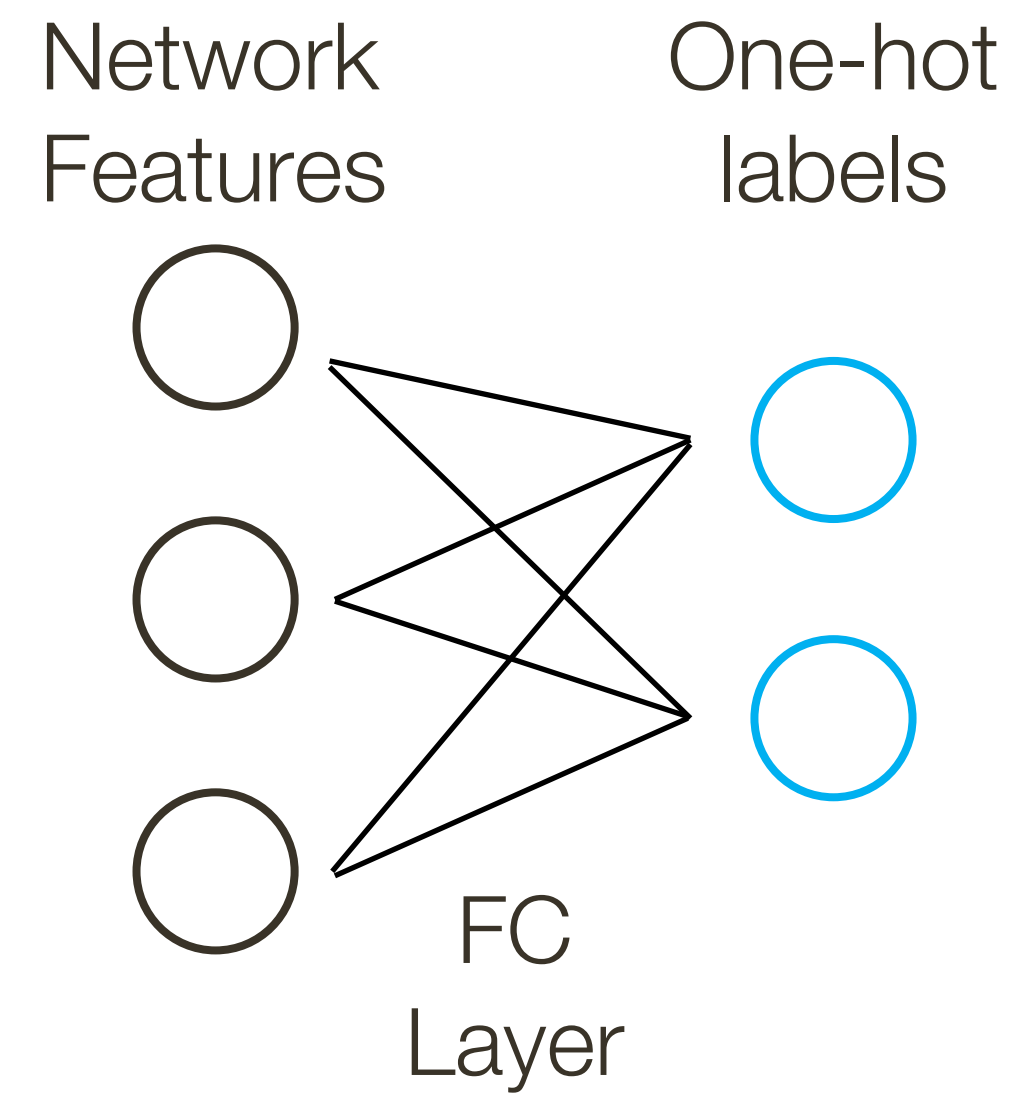
$$\beta \in \mathbb{R}^{d \times c}$$

β is sensitive to noisy data!

Approximated Linear Assumption in Networks



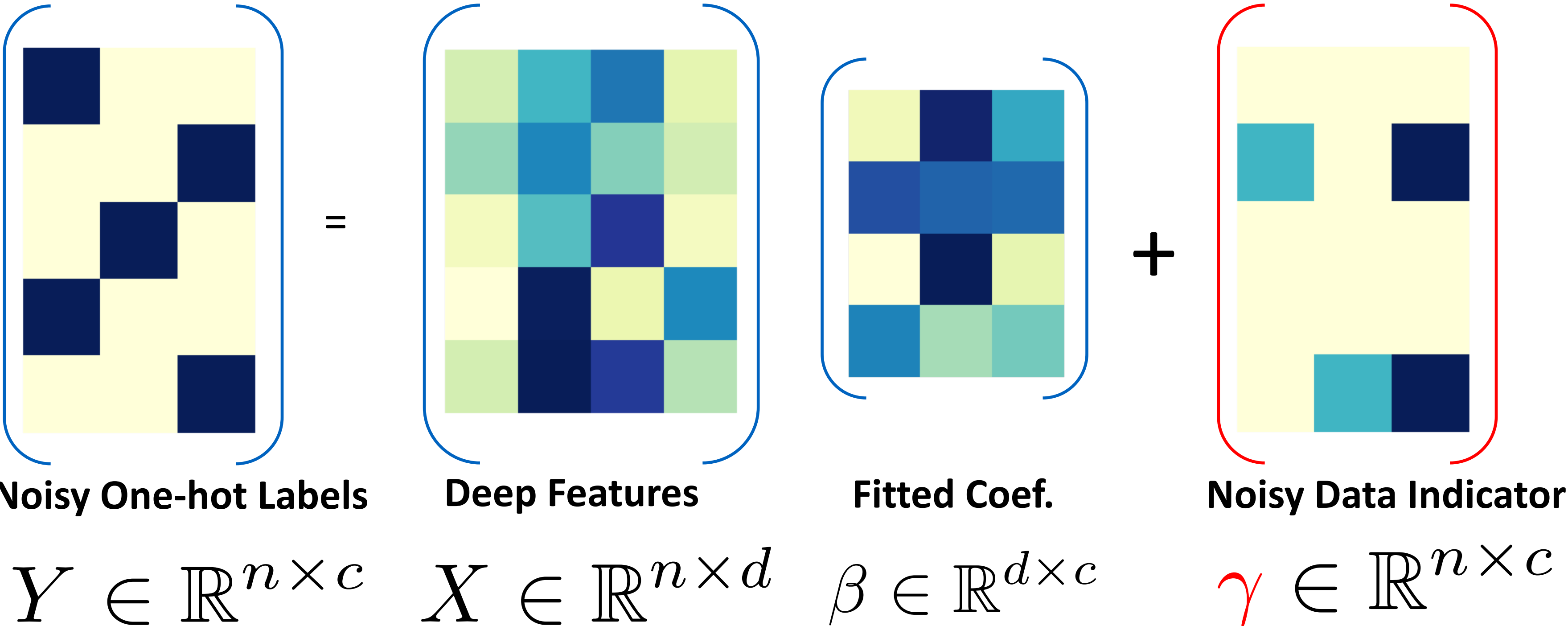
$$y_i = \text{SoftMax}(\mathbf{x}_i^T \beta)$$



$$y_i = \mathbf{x}_i^T \beta + \varepsilon$$

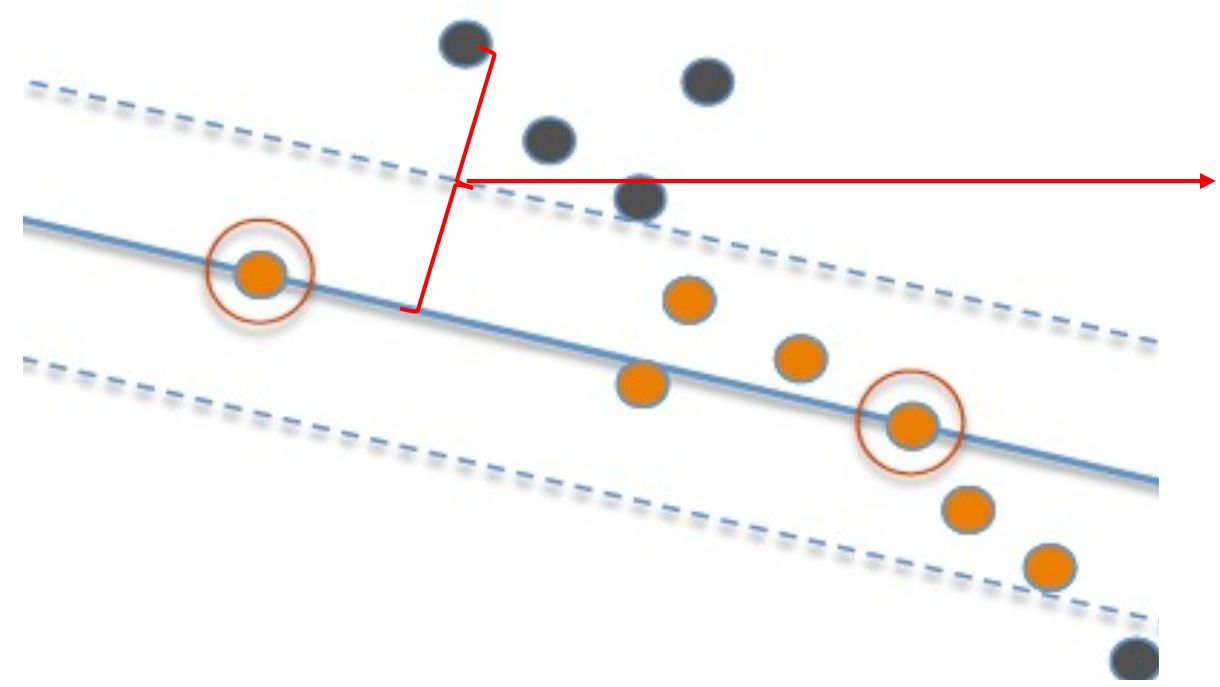
Identify Noisy Data in Label Space: The Indicator

Linear system with **Noisy** Data/Labels $Y = X\beta + \gamma$



Understanding γ in Statistics

$$y = x^\top \beta + \varepsilon + \gamma$$

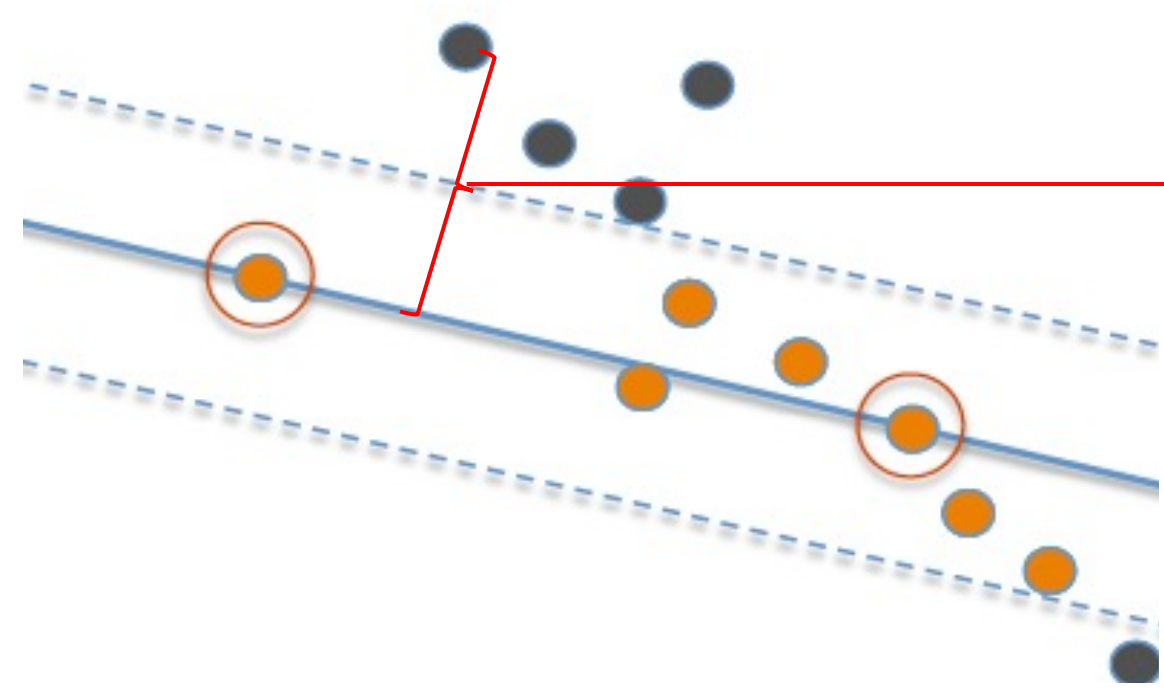


γ_i equals to the residual predict error $\gamma_i = y_i - x_i^\top \hat{\beta}$

Row residuals fail to detect outliers at *leverage points*.

Understanding γ in Statistics

$$y = x^\top \beta + \varepsilon + \gamma$$



γ_i equals to the residual predict error $\gamma_i = y_i - x_i^\top \hat{\beta}$

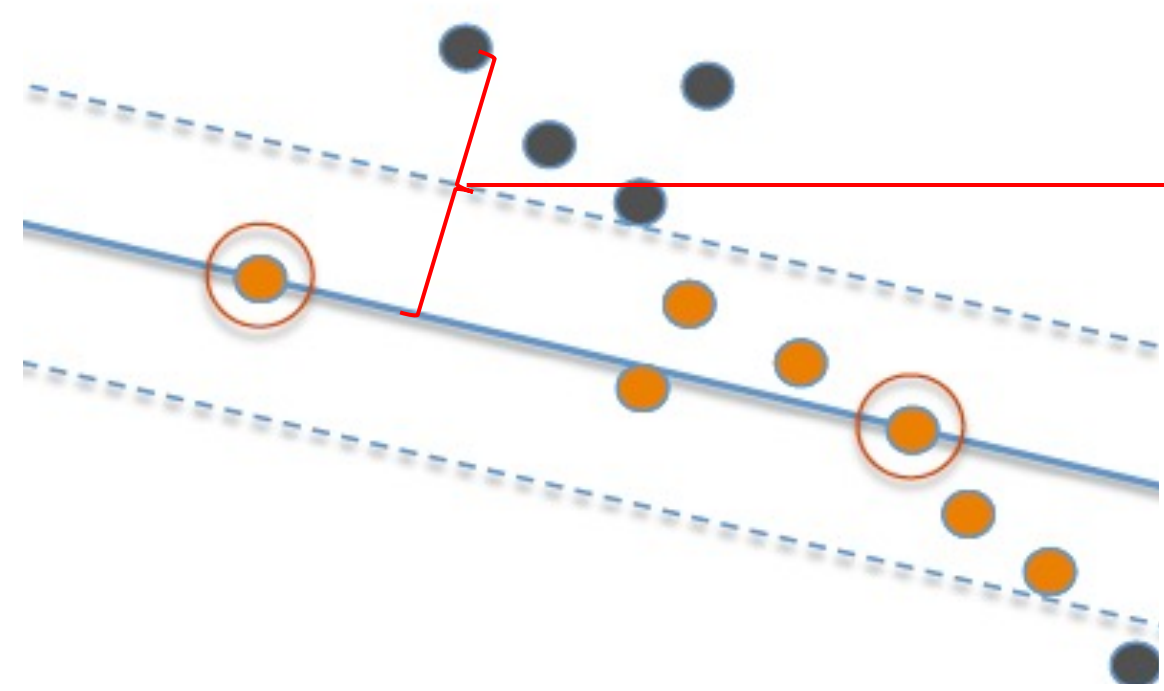
Leave-one-out externally studentized residual:

$$t_i = \frac{y_i - \mathbf{x}_i^\top \hat{\beta}_{(i)}}{\hat{\sigma}_{(i)} (1 + \mathbf{x}_i (\mathbf{X}_{(i)}^\top \mathbf{X}_{(i)})^{-1} \mathbf{x}_i)^{1/2}}$$

\Leftrightarrow test whether $\gamma = 0$ in $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \gamma\mathbf{1}_i + \boldsymbol{\varepsilon}$.

Understanding γ in Statistics

$$y = x^\top \beta + \varepsilon + \gamma$$



γ_i equals to the residual predict error $\gamma_i = y_i - x_i^\top \hat{\beta}$

Leave-one-out externally studentized residual:

$$t_i = \frac{y_i - \mathbf{x}_i^\top \hat{\beta}_{(i)}}{\hat{\sigma}_{(i)} (1 + \mathbf{x}_i (\mathbf{X}_{(i)}^\top \mathbf{X}_{(i)})^{-1} \mathbf{x}_i)^{1/2}}$$

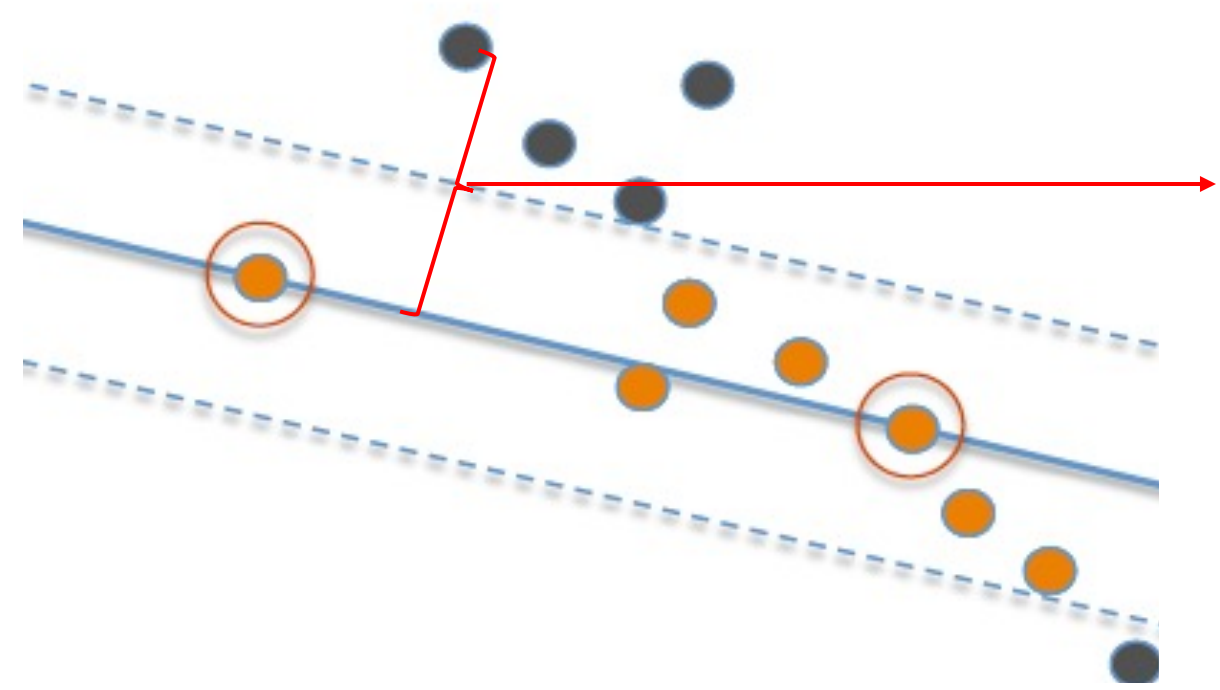
\Leftrightarrow test whether $\gamma = 0$ in $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \gamma\mathbf{1}_i + \boldsymbol{\varepsilon}$.

When there are multiple outliers:

- 1. masking:** multiple outliers may mask each other and being **undetected**;
- 2. swamping:** multiple outliers may lead the **large t_i for clean data**.

Understanding γ in Statistics

$$y = x^\top \beta + \varepsilon + \gamma$$



γ_i equals to the residual predict error $\gamma_i = y_i - x_i^\top \hat{\beta}$

Leave-one-out externally studentized residual:

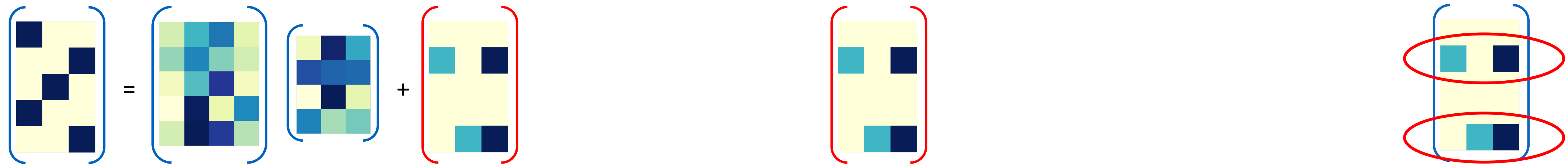
$$t_i = \frac{y_i - \mathbf{x}_i^\top \hat{\beta}_{(i)}}{\hat{\sigma}_{(i)} (1 + \mathbf{x}_i (\mathbf{X}_{(i)}^\top \mathbf{X}_{(i)})^{-1} \mathbf{x}_i)^{1/2}}$$

\Leftrightarrow test whether $\gamma = 0$ in $\mathbf{y} = \mathbf{X}\beta + \gamma 1_i + \varepsilon$.

$$\mathbf{y} = \mathbf{X}\beta + \varepsilon + \gamma$$

Identify Noisy Data in the Dataset

$$y_i = x_i^T \beta + \varepsilon + \gamma_i \quad \longrightarrow \quad \hat{\gamma}_i \quad \longrightarrow \quad O = \{i : \hat{\gamma}_i \neq 0\}$$



$$\operatorname{argmin}_{\beta, \gamma} L(\beta, \gamma) := \|\mathbf{Y} - \mathbf{X}\beta - \gamma\|_F^2 + \lambda R(\gamma)$$

Simplification

$$\operatorname{argmin}_{\beta, \gamma} L(\beta, \gamma) := \|\mathbf{Y} - \mathbf{X}\beta - \gamma\|_{\mathbb{F}}^2 + \lambda R(\gamma)$$

$$\frac{\partial L}{\partial \beta} = 0 \quad \downarrow \quad \hat{\beta} = (\mathbf{X}^{\top} \mathbf{X})^{\dagger} \mathbf{X}^{\top} (\mathbf{Y} - \gamma)$$

$$\operatorname{argmin}_{\gamma} \left\| \mathbf{Y} - \mathbf{X} (\mathbf{X}^{\top} \mathbf{X})^{\dagger} \mathbf{X}^{\top} (\mathbf{Y} - \gamma) - \gamma \right\|_{\mathbb{F}}^2 + \lambda R(\gamma)$$

$$\mathbf{H} = \mathbf{X} (\mathbf{X}^{\top} \mathbf{X})^{\dagger} \mathbf{X}^{\top} \quad \downarrow \quad \tilde{\mathbf{X}} = \mathbf{I} - \mathbf{H}, \tilde{\mathbf{Y}} = \tilde{\mathbf{X}} \mathbf{Y}$$

$$\operatorname{argmin}_{\gamma} \left\| \tilde{\mathbf{Y}} - \tilde{\mathbf{X}} \gamma \right\|_{\mathbb{F}}^2 + \lambda R(\gamma)$$

A linear regression problem!

Solving Gamma in Linear Regression

$$\operatorname{argmin}_{\gamma} \left\| \tilde{\mathbf{Y}} - \tilde{\mathbf{X}} \gamma \right\|_{\text{F}}^2 + \lambda R(\gamma)$$

How to select λ ?

- heuristics rules $\lambda = 2.5\hat{\sigma}$?
- Cross-validation?
- Data adaptive techniques?
- AIC, BIC?

It is hard to select a proper λ .

We regard $\hat{\gamma} = f(\lambda)$.

When $\lambda \rightarrow \infty$, $\hat{\gamma} \rightarrow 0$.

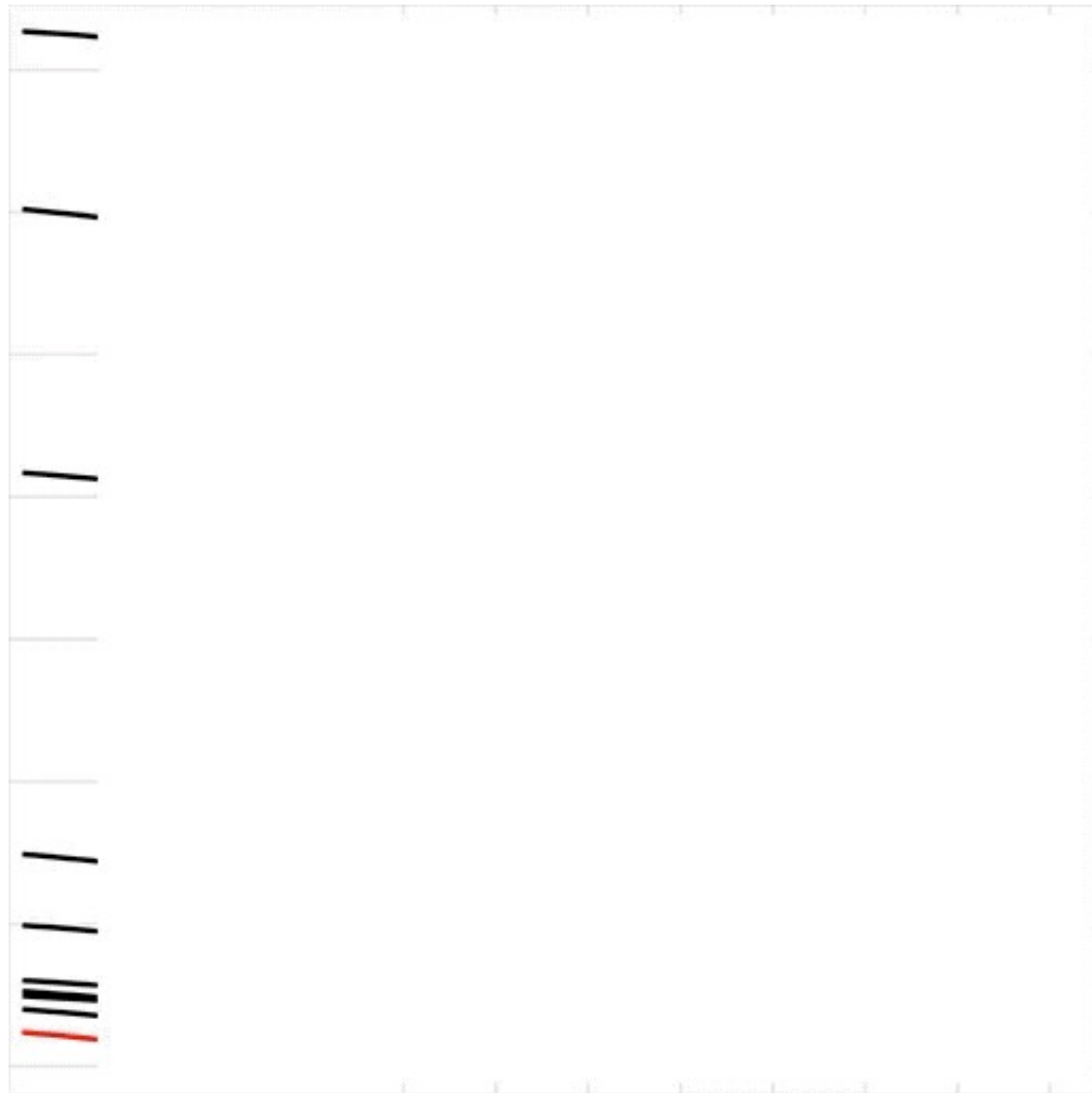
With $R(\gamma) = \sum_{i=1}^n \|\gamma_i\|_2$,
 γ vanishes instance by instance.

$$C_i = \sup\{\lambda : \|\hat{\gamma}_i(\lambda)\| \neq 0\}$$

This can be solved by GLMnet[1].

Solving Gamma in Linear Regression

$$\operatorname{argmin}_{\gamma} \left\| \tilde{\mathbf{Y}} - \tilde{\mathbf{X}} \gamma \right\|_{\text{F}}^2 + \lambda R(\gamma)$$



We regard $\hat{\gamma} = f(\lambda)$.

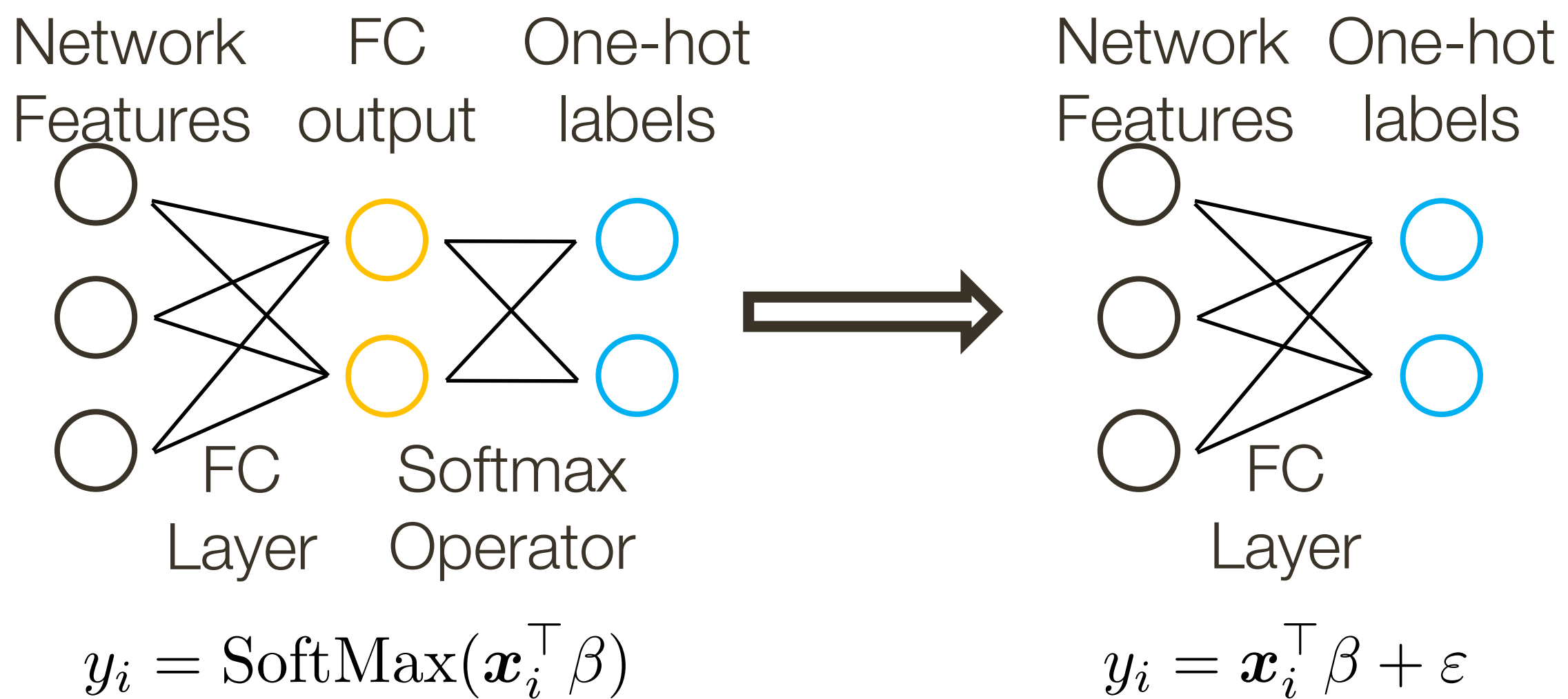
When $\lambda \rightarrow \infty$, $\hat{\gamma} \rightarrow 0$.

With $R(\gamma) = \sum_{i=1}^n \|\gamma_i\|_2$,
 γ vanishes instance by instance.

$$C_i = \sup\{\lambda : \|\hat{\gamma}_i(\lambda)\| \neq 0\}$$

This can be solved by GLMnet[1].

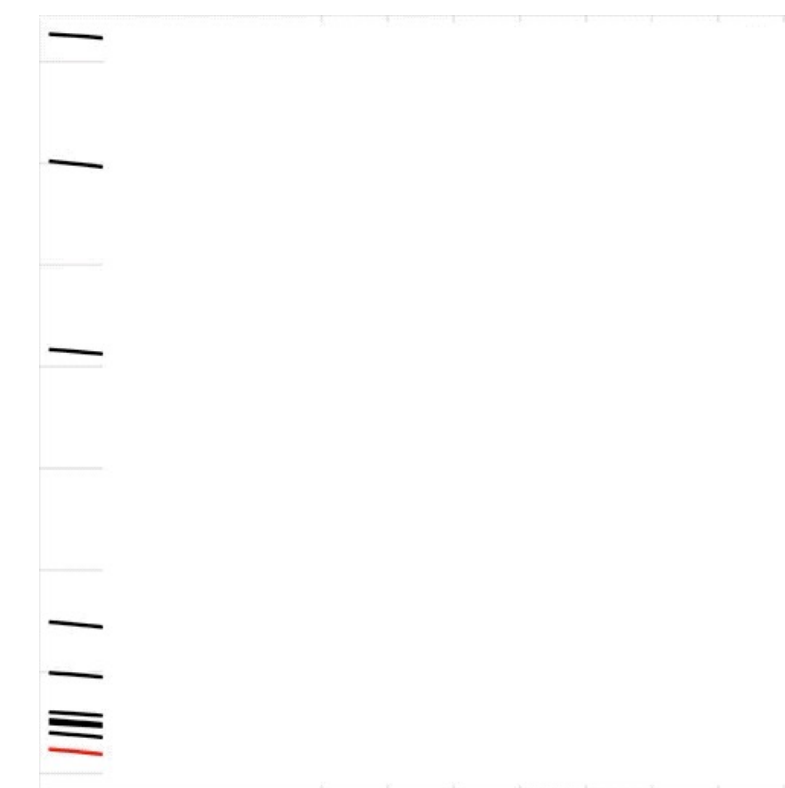
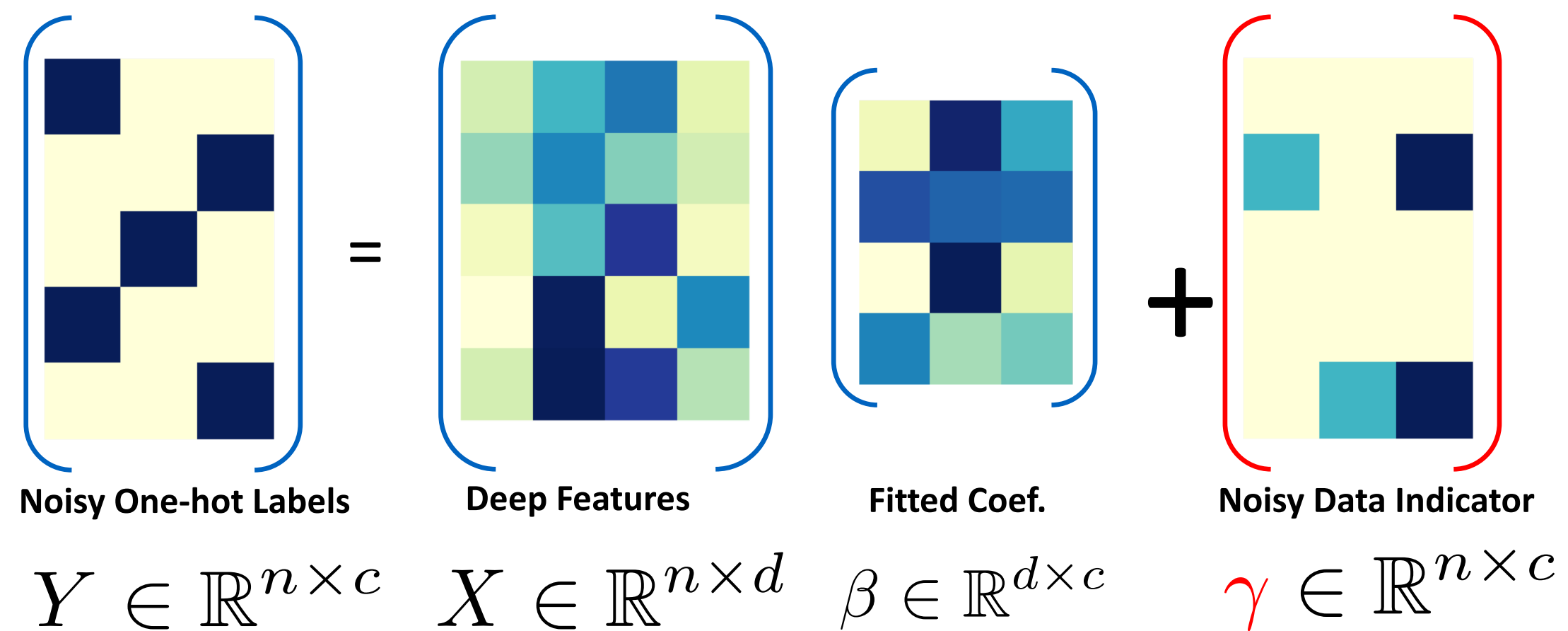
Instance Credibility Inference



$$\underset{\beta, \gamma}{\operatorname{argmin}} L(\beta, \gamma) := \|\mathbf{Y} - \mathbf{X}\beta - \gamma\|_F^2 + \lambda R(\gamma)$$

$$\underset{\gamma}{\operatorname{argmin}} \left\| \tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\gamma \right\|_F^2 + \lambda R(\gamma)$$

$$C_i = \sup\{\lambda : \|\hat{\gamma}_i(\lambda)\| \neq 0\}$$



Wang et al. Instance Credibility Inference for Few-Shot Learning. CVPR 2020

Wang et al. How to Trust Unlabeled Data? Instance Credibility Inference for Few-Shot Learning. IEEE TPAMI 2021.

Wang et al. Scalable Penalized Regression for Noise Detection in Learning with Noisy Labels. CVPR 2022.

Noise Set Recovery

When will the model identify all the outliers?

Assume ε is i.i.d zero-mean sub-Gaussian noise. We give three conditions:

- (C1: Restricted eigenvalue)

$$\lambda_{\min} \left(\tilde{\mathbf{U}}_S^\top \tilde{\mathbf{U}}_S \right) = C_{\min} > 0.$$

- (C2: Irrepresentability) $\exists \eta \in (0, 1]$,

$$\left\| \tilde{\mathbf{U}}_{S^c}^\top \tilde{\mathbf{U}}_S \left(\tilde{\mathbf{U}}_S^\top \tilde{\mathbf{U}}_S \right)^{-1} \right\|_{\infty} \leq 1 - \eta.$$

- (C3: Large error)

$$\vec{\gamma}_{\min} := \min_{i \in S} |\vec{\gamma}^*| > h \left(\lambda, \eta, \tilde{\mathbf{U}}, \vec{\gamma}^* \right).$$

A non-asymptotic probabilistic result

Based on these conditions, we could provide the following theorem:

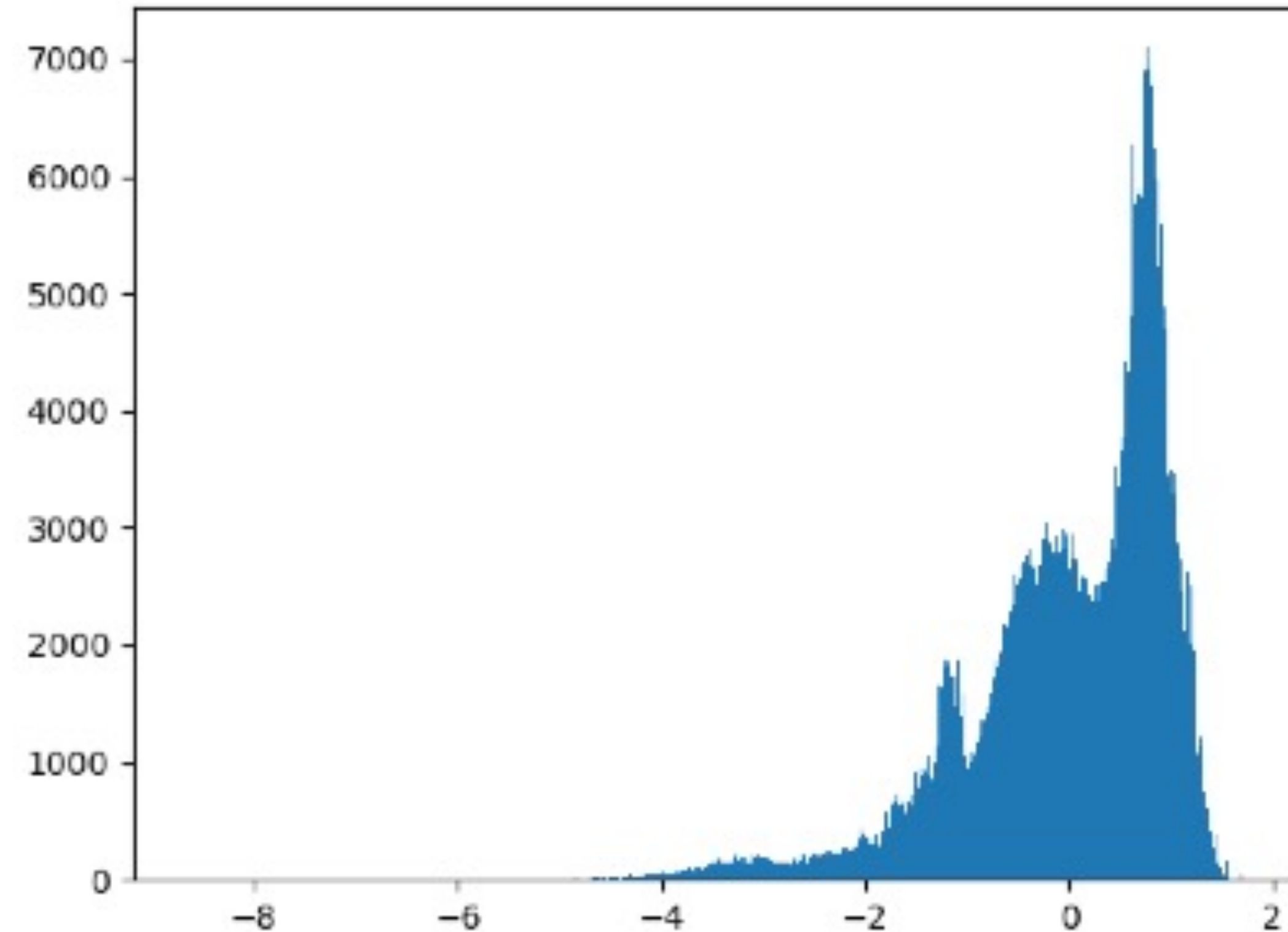
Theorem 1 (Identifiability of ICI). *Let $\lambda \geq \frac{2\sigma\sqrt{\mu\tilde{U}}}{\eta}\sqrt{\log cn}$. Then with probability greater than $1 - 2(cn)^{-1}$, the problem has a unique solution $\hat{\gamma}$ satisfies the following properties:*

1) *If C1 and C2 hold, the wrong-predicted instances indicated by ICI has no false positive error, i.e., $\hat{S} \subseteq S$ and hence $\hat{O} \subseteq O$, and*

$$\left\| \hat{\vec{\gamma}}_S - \vec{\gamma}_S^* \right\|_{\infty} \leq h\left(\lambda, \eta, \tilde{U}, \vec{\gamma}^*\right);$$

2) *If C1, C2, and C3 hold, ICI will identify all the correctly-predicted instance, i.e., $\hat{S} = S$ and hence $\hat{O} = O$ (in fact $\text{sign}\left(\hat{\vec{\gamma}}\right) = \text{sign}\left(\vec{\gamma}^*\right)$).*

Identifiability in reality: sub-Gaussian noise



$$\mathbb{E} [\hat{\varepsilon}] \approx 10^{-19}$$
$$\text{Var} [\hat{\varepsilon}] \approx 0.99$$

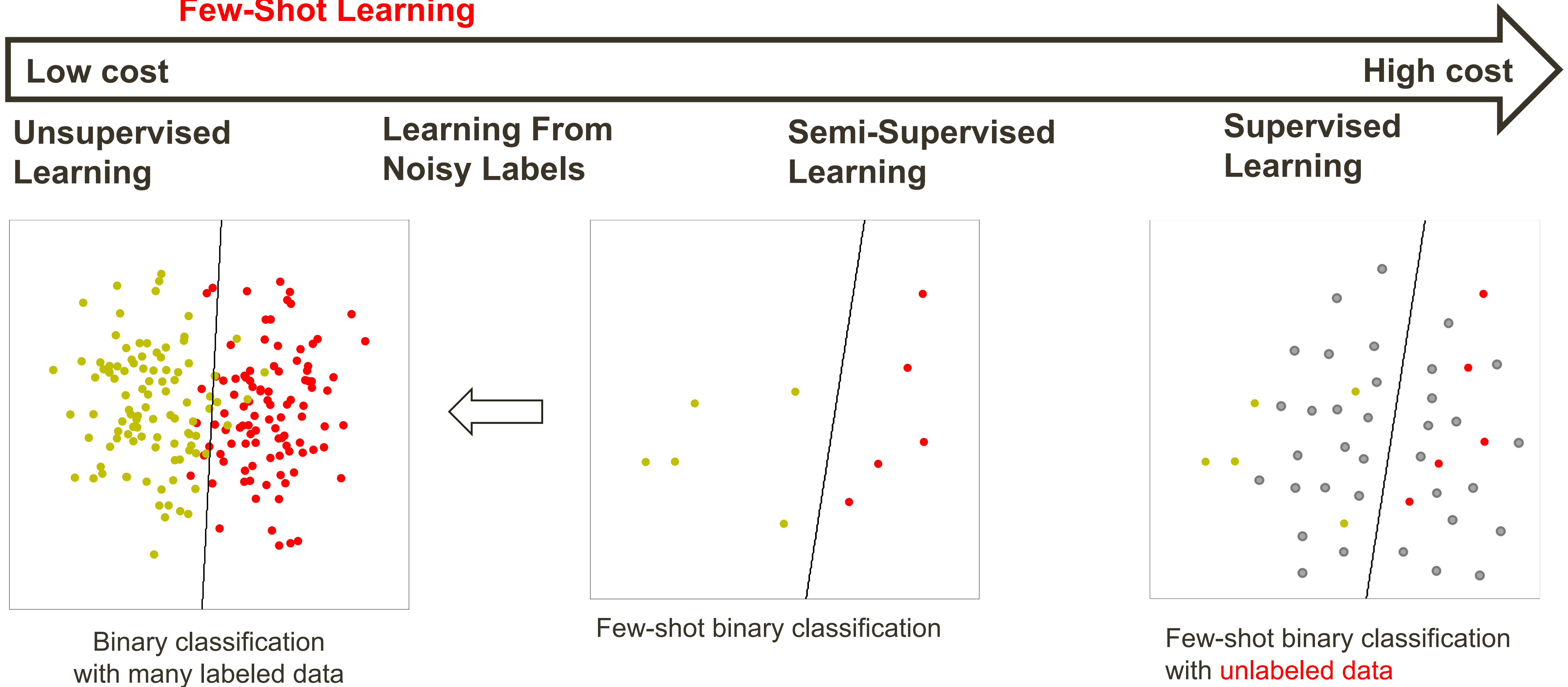
Sparse Learning

in Few-Shot Learning

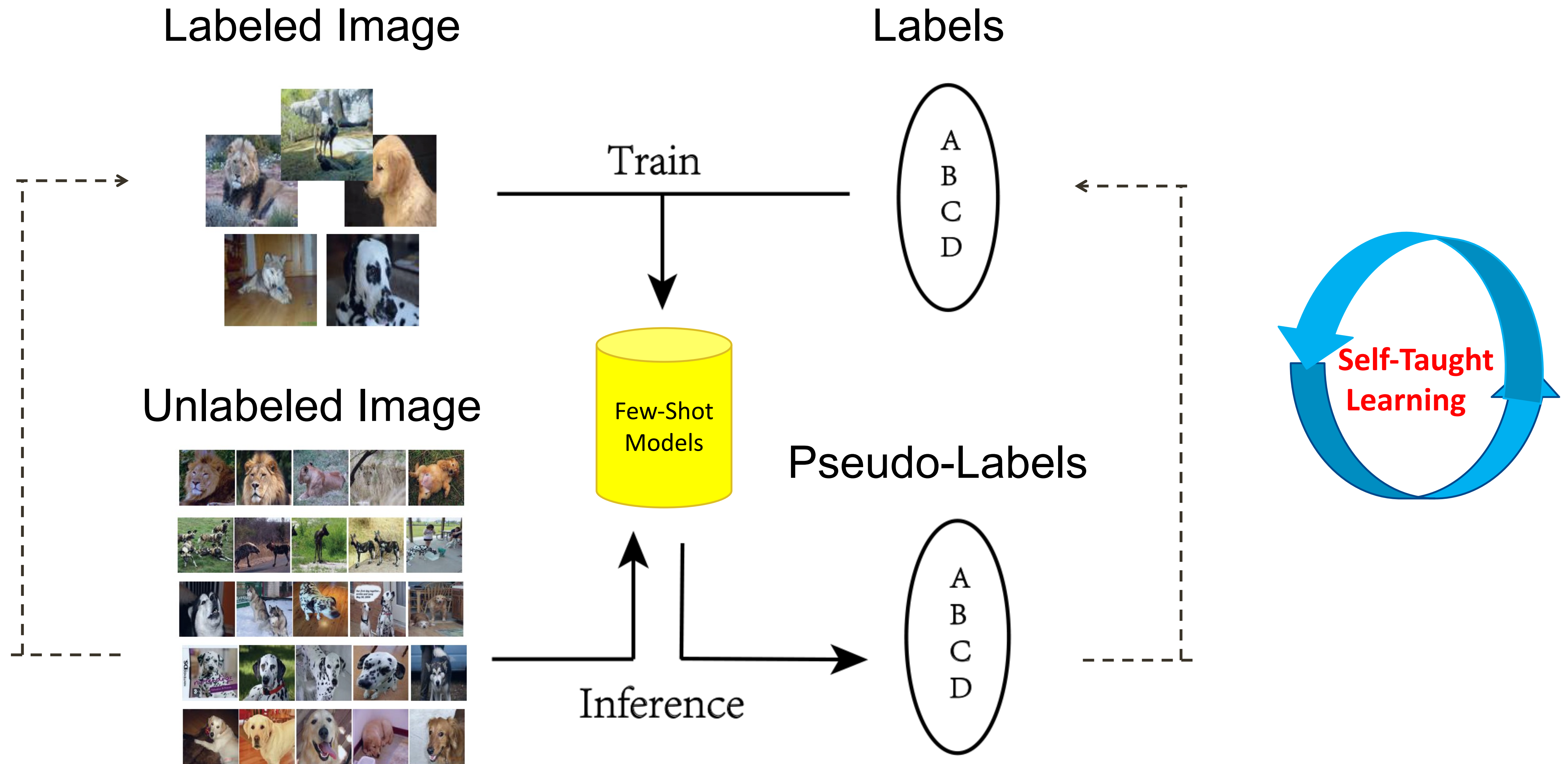
Definition of Few-Shot Learning

Tackle machine learning problem with only limited training data provided.

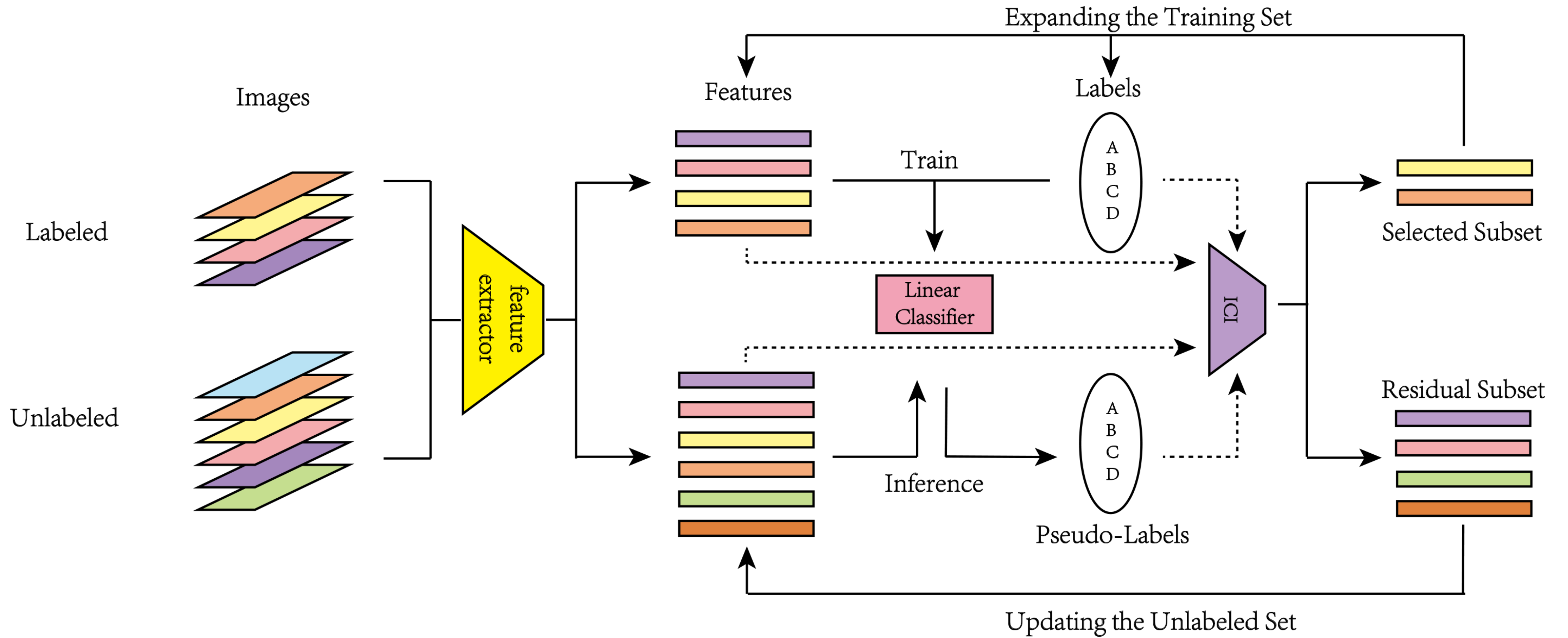
Few-Shot Learning



Motivation

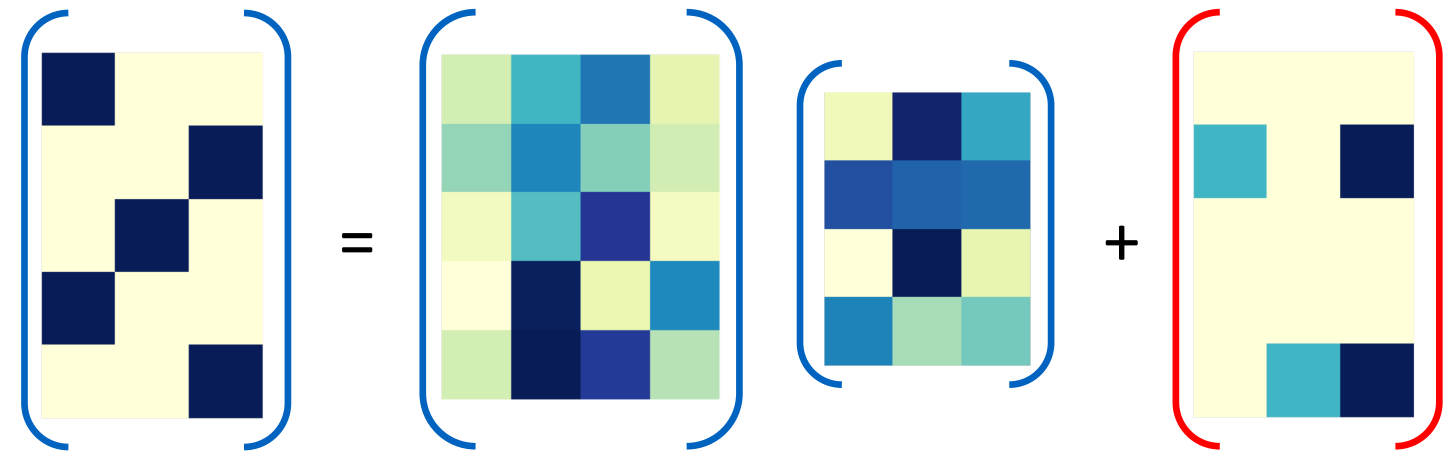


Framework



Sparse Learning in ICI

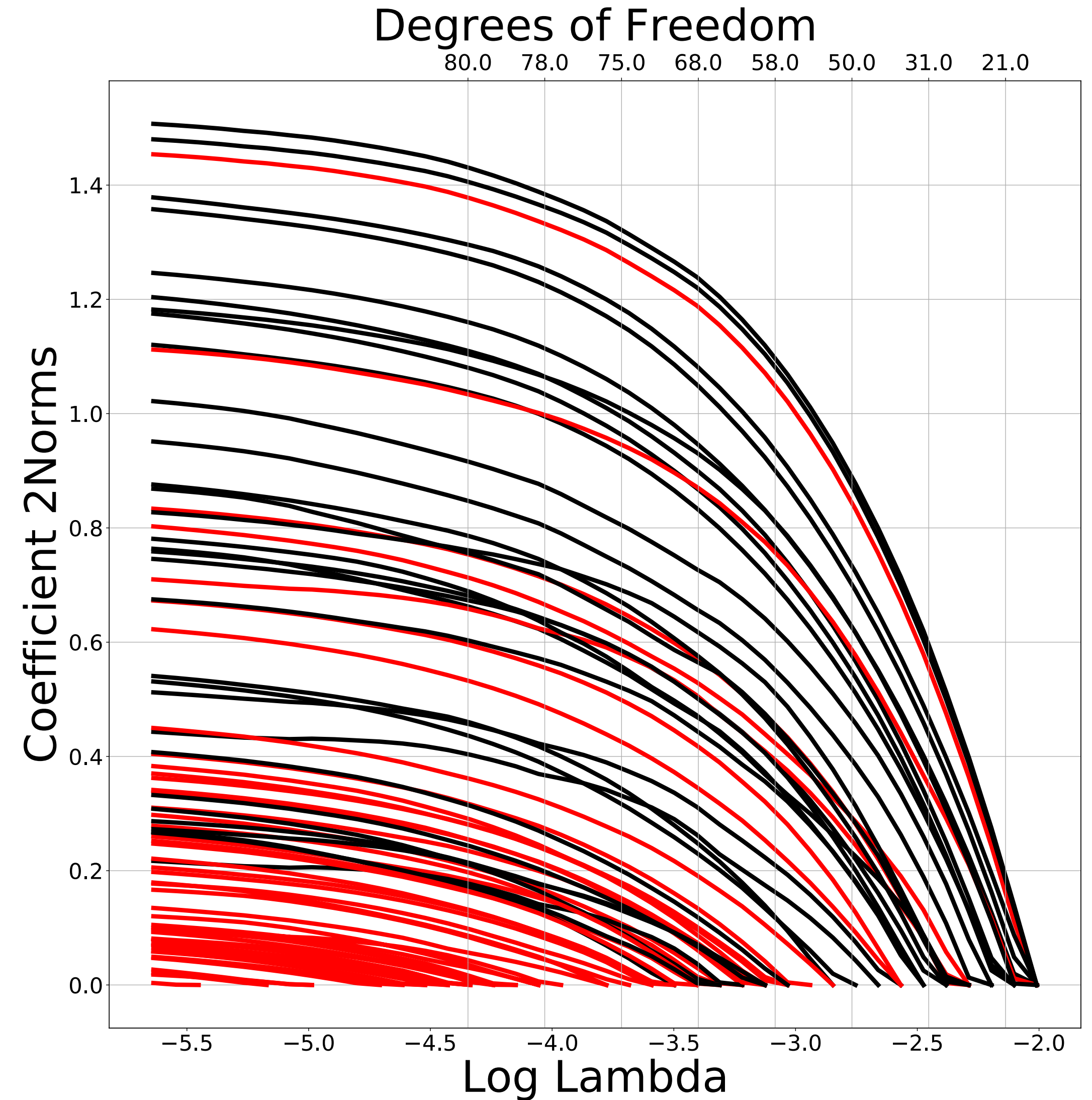
$$y_i = x_i^\top \beta + \varepsilon + \gamma_i$$



$$\operatorname{argmin}_{\beta, \gamma} L(\beta, \gamma) := \|Y - X\beta - \gamma\|_F^2 + \lambda R(\gamma)$$



$$\operatorname{argmin}_{\gamma} \left\| \tilde{Y} - \tilde{X}\gamma \right\|_F^2 + \lambda R(\gamma)$$



Sparse Learning: Extend to Logistic Regression

$$\operatorname{argmin}_{\beta, \gamma} L(\beta, \gamma) := \|\mathbf{Y} - \mathbf{X}\beta - \gamma\|_F^2 + \lambda R(\gamma)$$

$$\operatorname{argmin}_{\gamma} \left\| \mathbf{Y} - \mathbf{X} (\mathbf{X}^\top \mathbf{X})^\dagger \mathbf{X}^\top (\mathbf{Y} - \gamma) - \gamma \right\|_F^2 + \lambda R(\gamma)$$

$$\operatorname{argmin}_{\gamma} \left\| \tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\gamma \right\|_F^2 + \lambda R(\gamma)$$

$$\mathbf{Y}_{i,c} = \frac{\exp(\mathbf{X}_{i,\cdot} \beta_{\cdot,c} + \gamma_{i,c})}{\sum_{l=1}^C \exp(\mathbf{X}_{i,\cdot} \beta_{\cdot,l} + \gamma_{i,l})} + \epsilon_{i,c}$$

$$\bar{\mathbf{X}} = (\mathbf{X}, \mathbf{I}) \quad \bar{\beta} = (\beta, \gamma)^\top$$

$$\mathbf{Y}_{i,c} = \frac{\exp(\bar{\mathbf{X}}_{i,\cdot} \bar{\beta}_{\cdot,c})}{\sum_{l=1}^C \exp(\bar{\mathbf{X}}_{i,\cdot} \bar{\beta}_{\cdot,l})} + \epsilon_{i,c}$$

Identifiability in Reality: Conditions and Accuracy

Satisfied Assumptions	None	C1	C1 and C2	All
Improved Episodes	0	424	1035	40
Total Episodes	0	793	1164	43
I/T	—	53.5%	88.9%	93.0%

1) In more than half of the experiments the assumptions C1-C2 are satisfied. Most of them (89.0%) will achieve better performance after self-taught with ICI.

Identifiability in Reality: Conditions and Accuracy

Satisfied Assumptions	None	C1	C1 and C2	All
Improved Episodes	0	424	1035	40
Total Episodes	0	793	1164	43
I/T	—	53.5%	88.9%	93.0%

2) When all the assumptions are satisfied, we will get better performance in a high ratio (93.0%).

Identifiability in Reality: Conditions and Accuracy

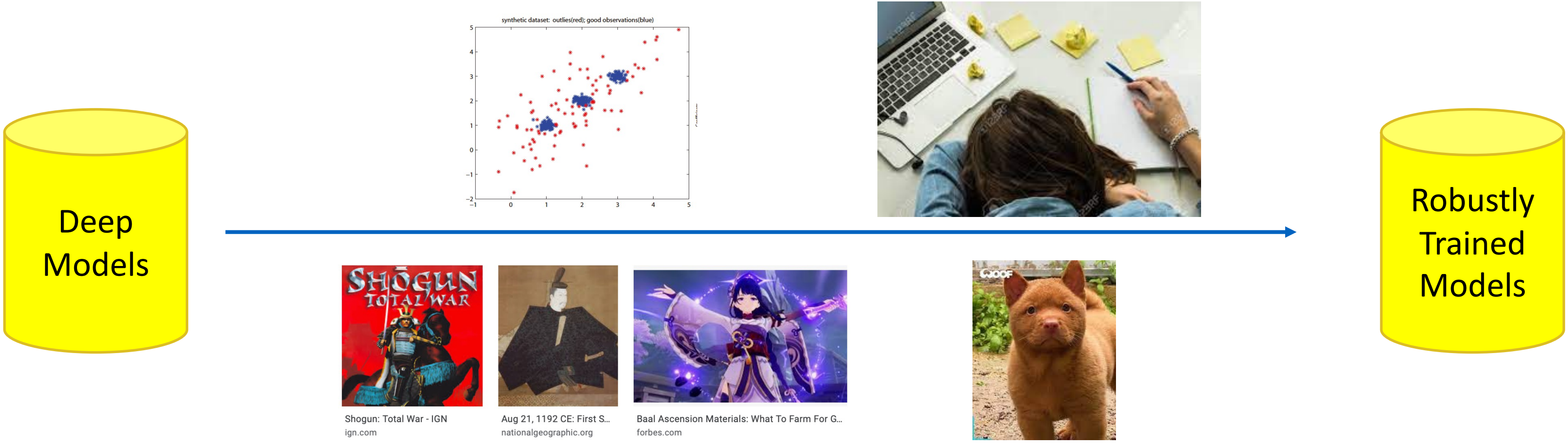
Satisfied Assumptions	None	C1	C1 and C2	All
Improved Episodes	0	424	1035	40
Total Episodes	0	793	1164	43
I/T	—	53.5%	88.9%	93.0%

3) Even if C2-C3 are not satisfied, we still have the chance of improving the performance (53.5%).

Sparse Learning

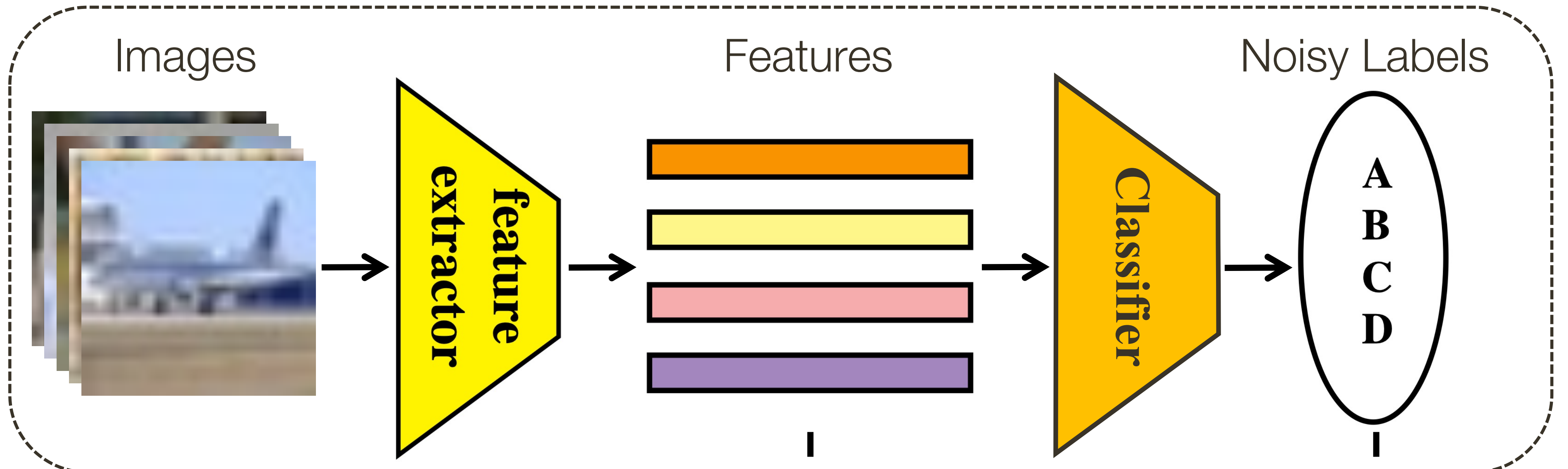
in Learning with Noisy Labels

Definition of learning with noisy labels

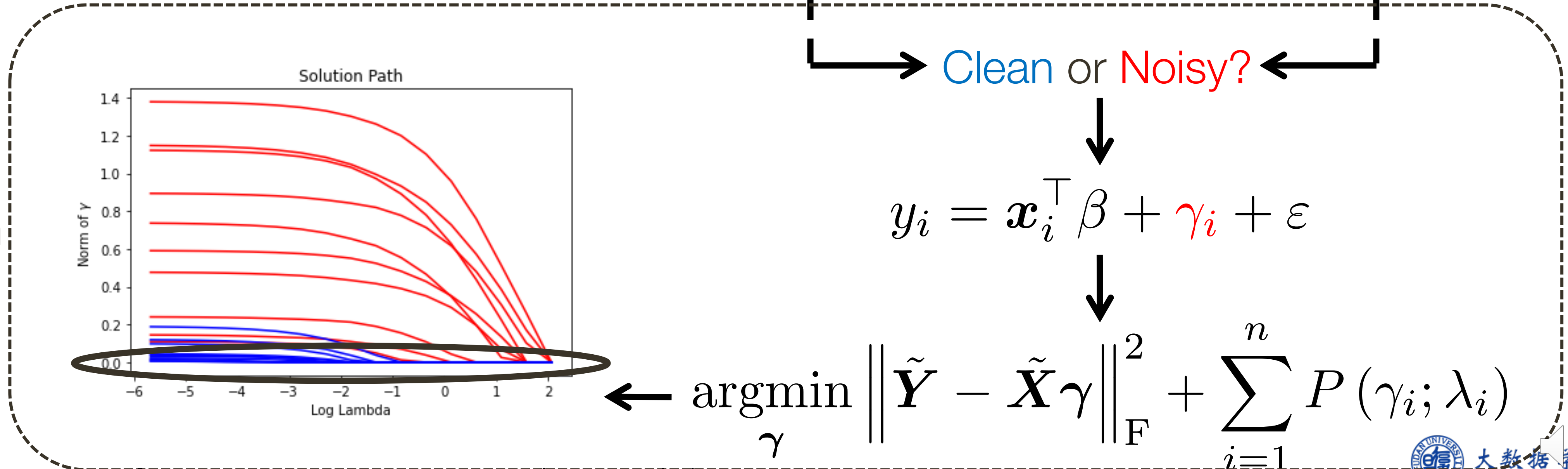


Framework

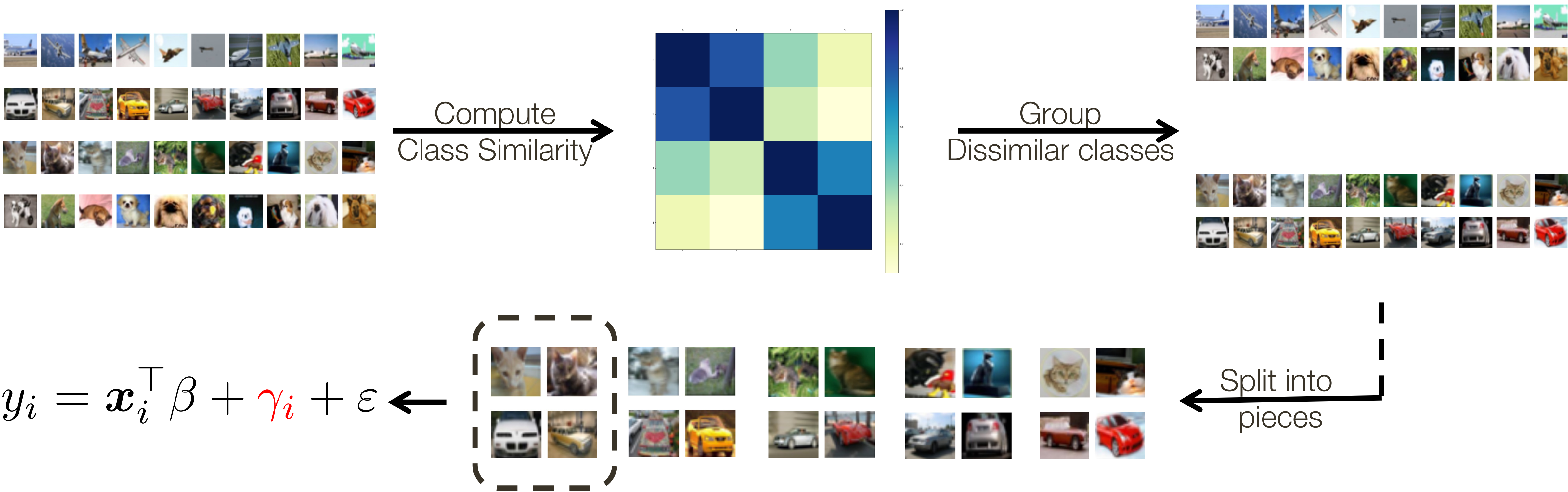
Stage 1:
Feature Learning



Stage 2:
Sample Selection



Make it scalable to large datasets



Strategies to help train the network

- Append a $\ell_q (q < 1)$ penalty to encourage the linear relation between feature and one-hot encoded vector:

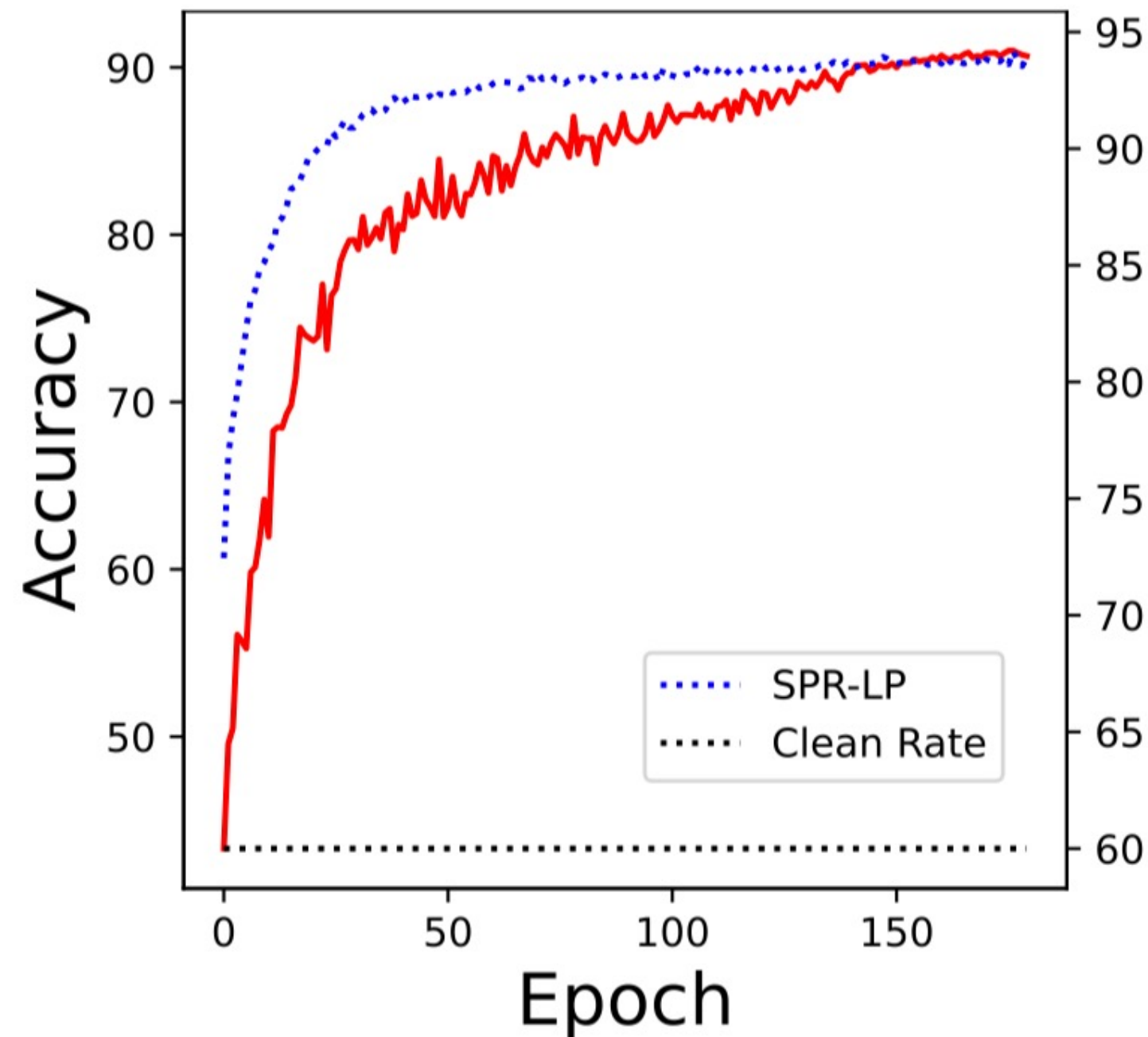
$$\mathcal{L}(\mathbf{x}_i, \mathbf{y}_i) = 1_{i \notin O} \left(\mathcal{L}_{\text{CE}}(\mathbf{x}_i, \mathbf{y}_i) + \lambda \|\mathbf{x}_i^\top W_{\text{fc}}\|_q \right)$$

- Use CutMix to further exploit the support of noisy data

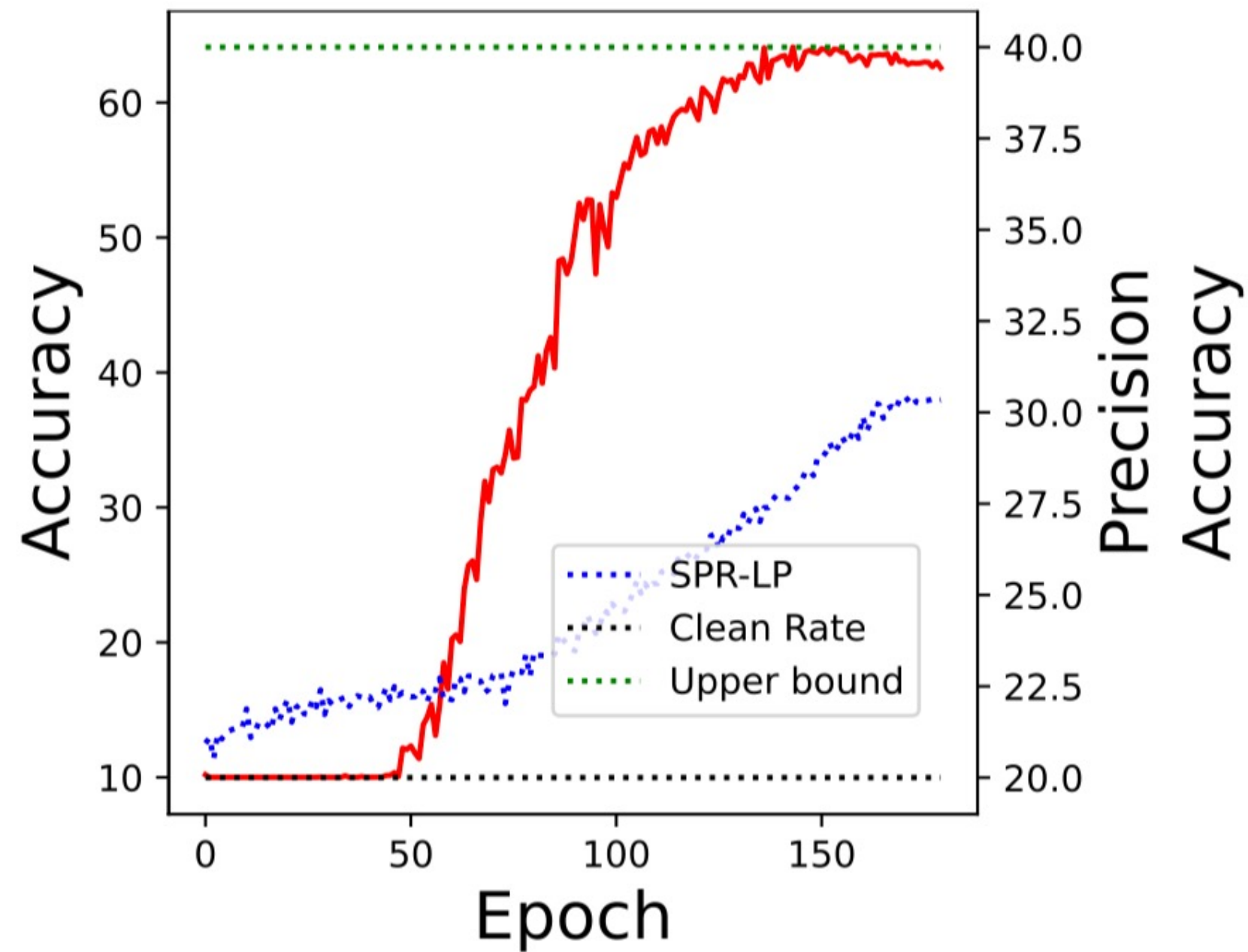
$$\tilde{\mathbf{x}} = \mathbf{M} \odot \mathbf{x}_{\text{clean}} + (1 - \mathbf{M}) \odot \mathbf{x}_{\text{noisy}}$$

$$\tilde{\mathbf{y}} = \lambda \mathbf{y}_{\text{clean}} + (1 - \lambda) \mathbf{y}_{\text{noisy}}$$

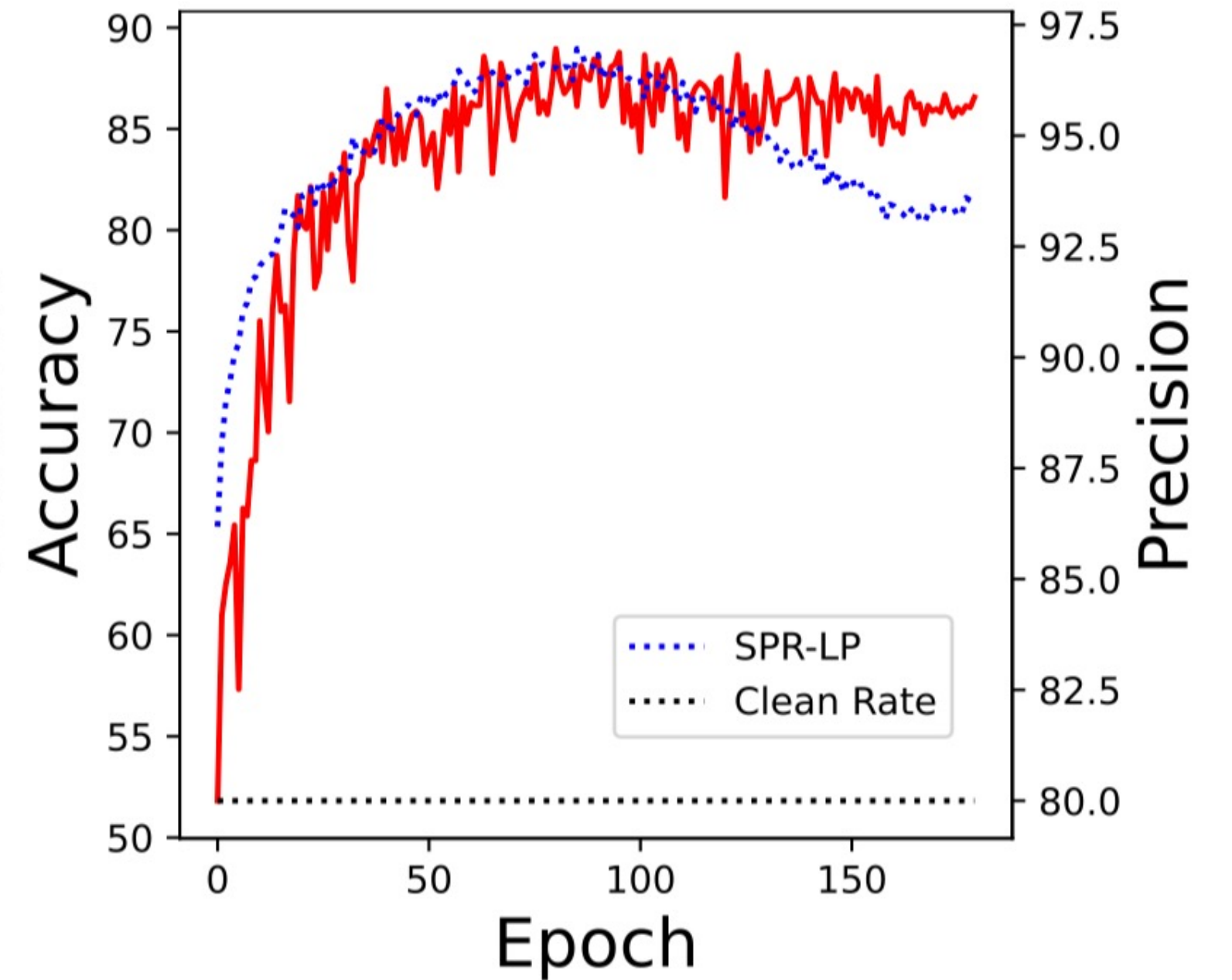
Label precision performance



(a) Symmetric-40%



(b) Symmetric-80%



(c) Asymmetric-40%

THANKS