# Sparse Learning for Noisy Data/Labels:
## A Simple yet Effective Framework for Vision Applications

**Yikai Wang**
yikai-wang.github.io
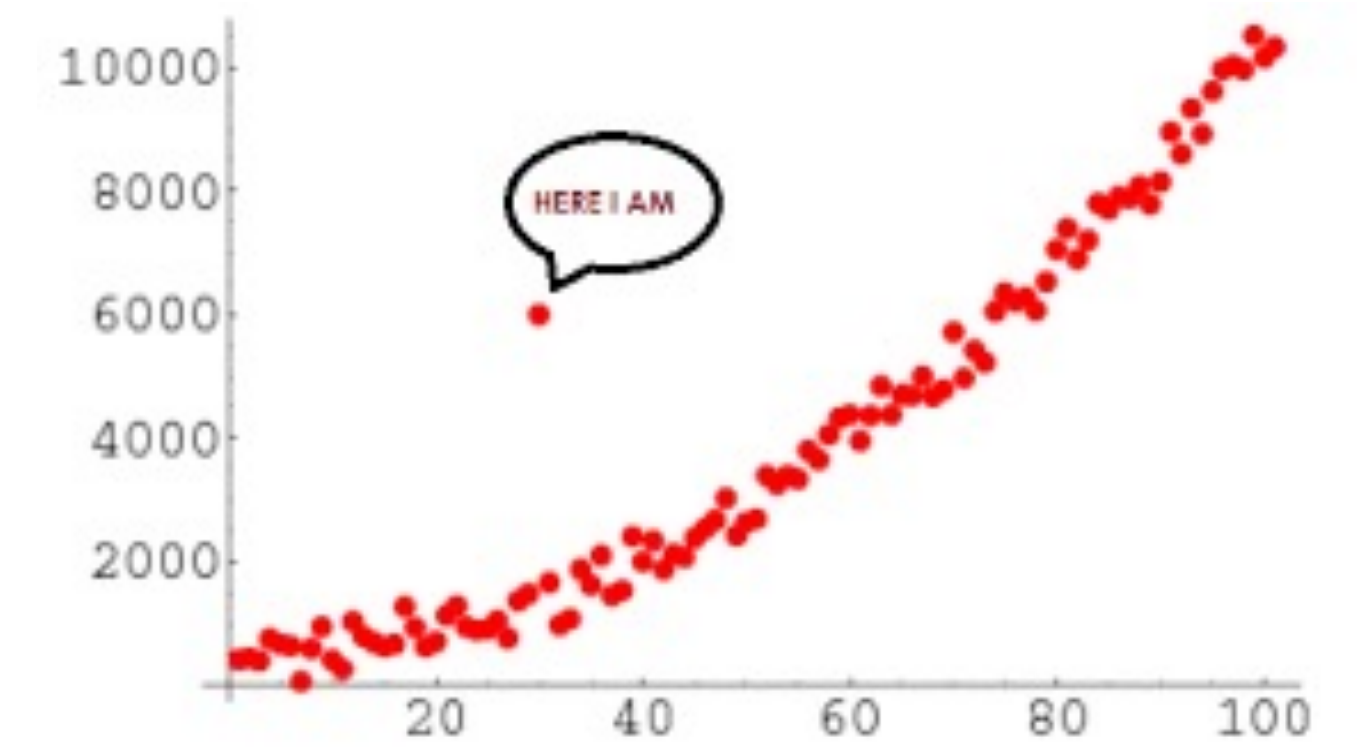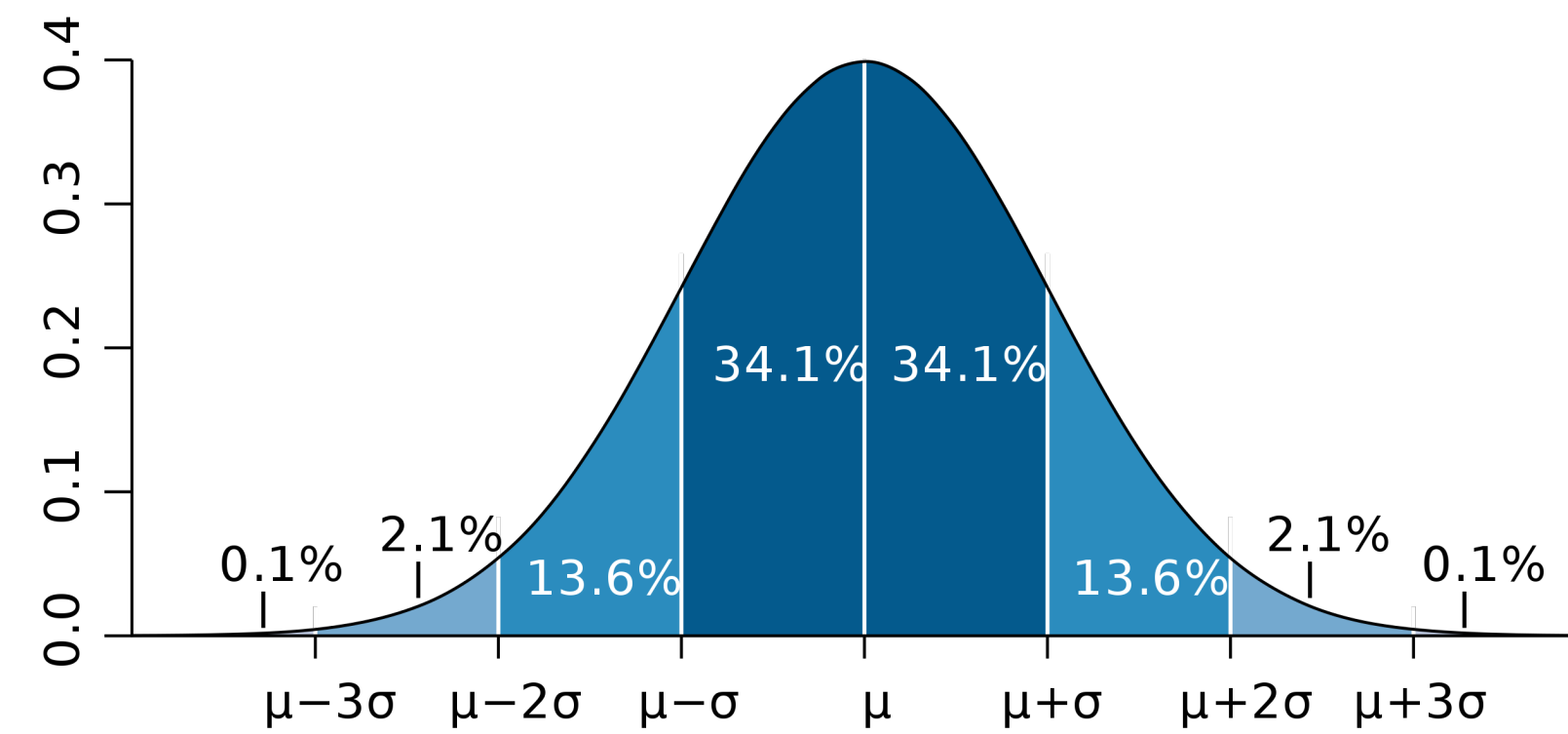
**Yanwei Fu**
http://yanweifu.github.io

School of Data Science
Fudan University

大数据学院
School of Data Science

# Sparse Learning

## for Noise Data Detection

# Examples of Noisy Data/Outliers



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

Figures from
[1] towardsdatascience.com/this-article-is-about-identifying-outliers-through-funnel-plots-using-the-microsoft-power-bi-d7ad16ac9ccc
[2] en.wikipedia.org/wiki/Outlier#/media/File:Standard_deviation_diagram_micro.svg
[3] medium.com/analytics-vidhya/its-all-about-outliers-cbe172aa1309

# Noisy Data in Label Space

- Random Corruptions



synthetic dataset: outlies(red); good observations(blue)

- 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

- Annotator mistakes



- Complex/Confusing items identified

# Identify Noisy Data in Label Space

*Linear system* $$Y = X\beta$$



**Noisy One-hot Labels**     **Deep Features**     **Fitted Coef.**

$$Y \in \mathbb{R}^{n \times c} \quad X \in \mathbb{R}^{n \times d} \quad \beta \in \mathbb{R}^{d \times c}$$

$\beta$ is sensitive to noisy data!

# Approximated Linear Assumption in Networks



$$y_i = \mathrm{SoftMax}(\boldsymbol{x}_i^\top \beta)$$

$$y_i = \boldsymbol{x}_i^\top \beta + \varepsilon$$

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

# Identify Noisy Data in Label Space: The Indicator

$$Y = X\beta + \gamma$$

*Linear system with Noisy Data/Labels*



**Noisy One-hot Labels**   **Deep Features**   **Fitted Coef.**   **Noisy Data Indicator**

$$Y \in \mathbb{R}^{n \times c} \qquad X \in \mathbb{R}^{n \times d} \qquad \beta \in \mathbb{R}^{d \times c} \qquad \gamma \in \mathbb{R}^{n \times c}$$

[Wright et al. TPAMI 09] [She et al. JASA 11] [Fu et al. ECCV 14, TPAMI 16.] [Fan et al. Statistical Sinica 18] [Wang et al. CVPR 20, TPAMI 21, CVPR 22]

大数据学院
School of Data Science

# Understanding γ in Statistics

$$y = x^\top \beta + \varepsilon + \textcolor{red}{\gamma}$$



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

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

[1] Yiyuan She and Art B Owen. Outlier detection using nonconvex penalized regression. Journal of the American Statistical Association, 2011.
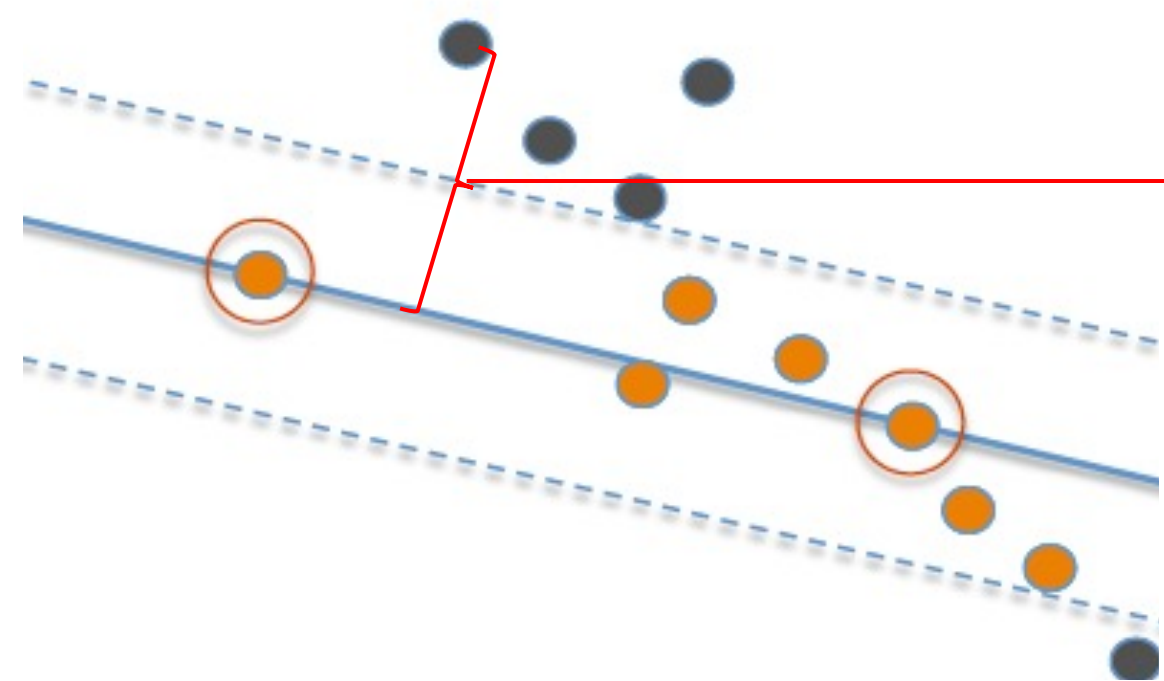
# Understanding γ in Statistics

$$y = x^\top \beta + \varepsilon + \textcolor{red}{\gamma}$$



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

Leave-one-out externally studentized residual:

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

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

[1] Yiyuan She and Art B Owen. Outlier detection using nonconvex penalized regression. Journal of the American Statistical Association, 2011.

# Understanding γ in Statistics

$$y = x^\top \beta + \varepsilon + {\color{red}\gamma}$$



${\color{red}\gamma_i}$ equals to the residual predict error ${\color{red}\gamma_i = y_i - x_i^\top \hat{\beta}}$
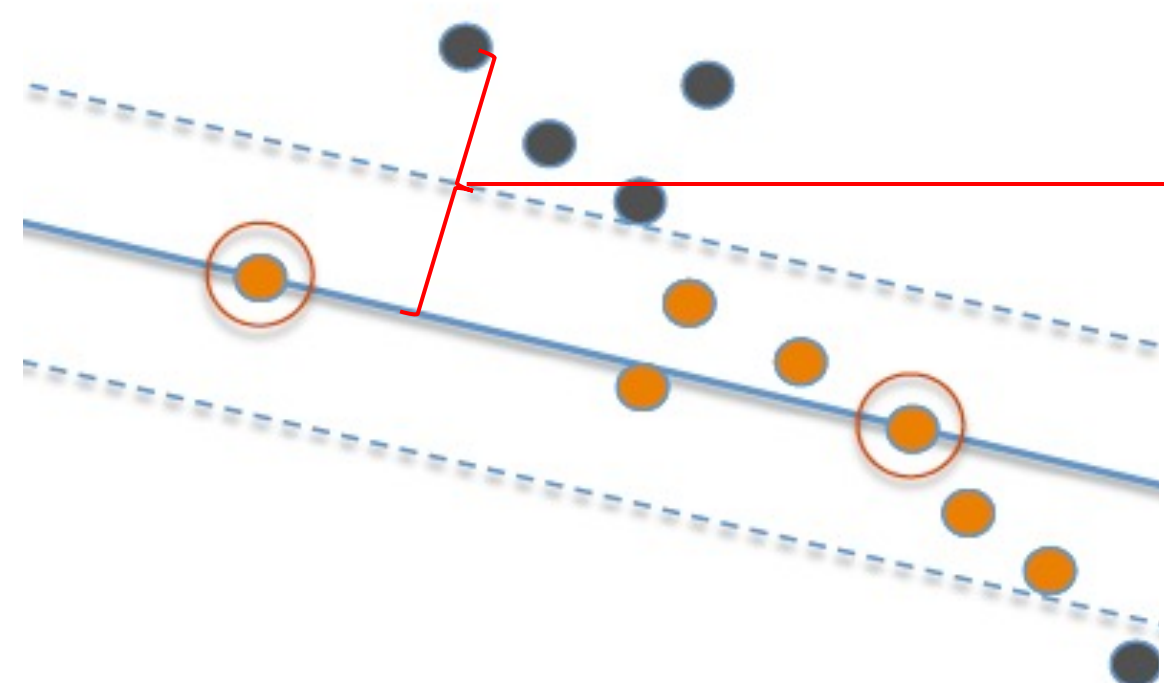
Leave-one-out externally studentized residual:

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

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

When there are multiple outliers:

**1. masking**: multiple outliers may mask each other and being ${\color{red}\text{undetected}}$;

**2. swamping**: multiple outliers may lead the ${\color{red}\text{large } t_i \text{ for clean data}}$.

[1] Yiyuan She and Art B Owen. Outlier detection using nonconvex penalized regression. Journal of the American Statistical Association, 2011.

大数据学院
School of Data Science

# Understanding γ in Statistics

$$y = x^\top \beta + \varepsilon + {\color{red}\gamma}$$



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

Leave-one-out externally studentized residual:

$$t_i = \frac{y_i - \boldsymbol{x}_i^\top \hat\beta_{(i)}}{\hat\sigma_{(i)}(1 + \boldsymbol{x}_i(\boldsymbol{X}_{(i)}^\top \boldsymbol{X}_{(i)})^{-1}\boldsymbol{x}_i)^{1/2}}$$

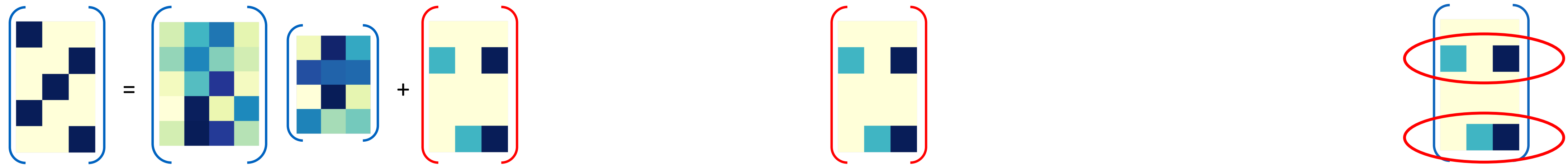$\Leftrightarrow$ test whether $\gamma = 0$ in $\boldsymbol{y} = \boldsymbol{X}\boldsymbol{\beta} + \gamma 1_i + \boldsymbol{\varepsilon}.$

$$\boldsymbol{y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} + \boldsymbol{\gamma}$$

[1] Yiyuan She and Art B Owen. Outlier detection using nonconvex penalized regression. Journal of the American Statistical Association, 2011.

# Identify Noisy Data in the Dataset

$$y_i = x_i^\top \beta + \varepsilon + {\color{red}\gamma_i} \qquad \longrightarrow \qquad {\color{red}\hat{\gamma}_i} \qquad \longrightarrow \qquad O = \{i : \hat{\gamma}_i \neq 0\}$$



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

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.

大数据学院
School of Data Science

# Simplification

$$\underset{\boldsymbol{\beta},\boldsymbol{\gamma}}{\arg\min} L\left(\boldsymbol{\beta},\boldsymbol{\gamma}\right) := \left\|\boldsymbol{Y} - \boldsymbol{X}\boldsymbol{\beta} - \boldsymbol{\gamma}\right\|_{\mathrm{F}}^{2} + \lambda R\left(\boldsymbol{\gamma}\right)$$

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

$$\underset{\boldsymbol{\gamma}}{\arg\min} \left\|\boldsymbol{Y} - \boldsymbol{X}\left(\boldsymbol{X}^{\top}\boldsymbol{X}\right)^{\dagger}\boldsymbol{X}^{\top}\left(\boldsymbol{Y} - \boldsymbol{\gamma}\right) - \boldsymbol{\gamma}\right\|_{\mathrm{F}}^{2} + \lambda R\left(\boldsymbol{\gamma}\right)$$

$$\boldsymbol{H} = \boldsymbol{X}\left(\boldsymbol{X}^{\top}\boldsymbol{X}\right)^{\dagger}\boldsymbol{X}^{\top} \quad \Bigg\downarrow \quad \tilde{\boldsymbol{X}} = \boldsymbol{I} - \boldsymbol{H}, \tilde{\boldsymbol{Y}} = \tilde{\boldsymbol{X}}\boldsymbol{Y}$$

$$\underset{\boldsymbol{\gamma}}{\arg\min} \left\|\tilde{\boldsymbol{Y}} - \tilde{\boldsymbol{X}}\boldsymbol{\gamma}\right\|_{\mathrm{F}}^{2} + \lambda R\left(\boldsymbol{\gamma}\right)$$

<span style="color:red">A linear regression problem!</span>

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.

大数据学院
School of Data Science

# Solving Gamma in Linear Regression

$$\underset{\boldsymbol{\gamma}}{\mathrm{argmin}} \left\| \tilde{\boldsymbol{Y}} - \tilde{\boldsymbol{X}}\boldsymbol{\gamma} \right\|_{\mathrm{F}}^2 + \lambda R(\boldsymbol{\gamma})$$

How to select $\lambda$?

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

It is hard to select a proper $\lambda$.

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

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

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

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

This can be sovled by GLMnet[1].

[1] Friedman, et al. 2010. "Regularization Paths for Generalized Linear Models via Coordinate Descent." Journal of Statistical Software.

# Solving Gamma in Linear Regression

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

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

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

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

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

This can be sovled by GLMnet[1].

[1] Friedman, et al. 2010. "Regularization Paths for Generalized Linear Models via Coordinate Descent." Journal of Statistical Software.

# Instance Credibility Inference

Network Features    FC output    One-hot labels

FC Layer    Softmax Operator

$$y_i = \mathrm{SoftMax}(\boldsymbol{x}_i^\top \beta)$$

Network Features    One-hot labels

FC Layer

$$y_i = \boldsymbol{x}_i^\top \beta + \varepsilon$$

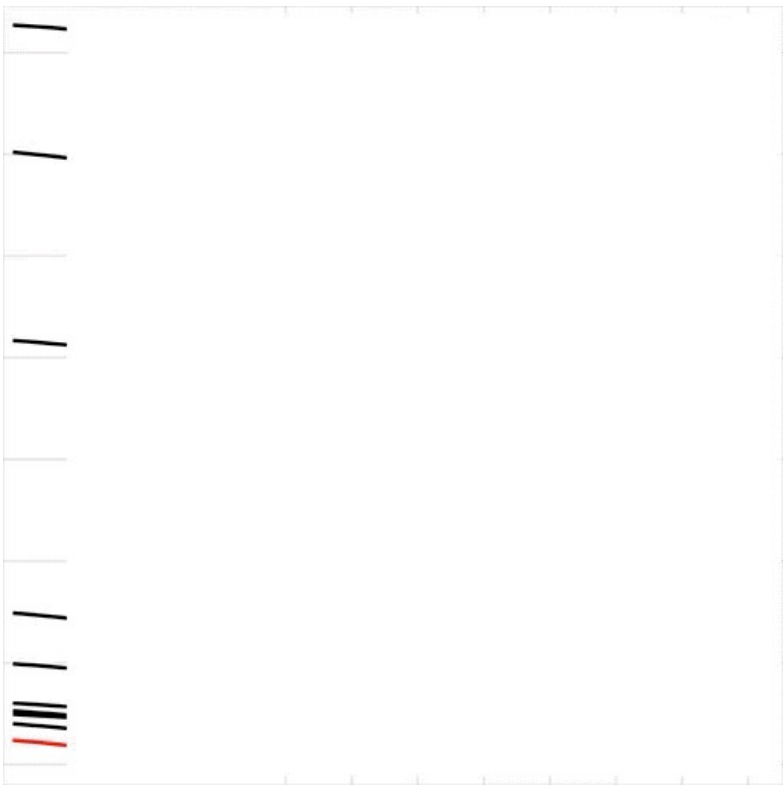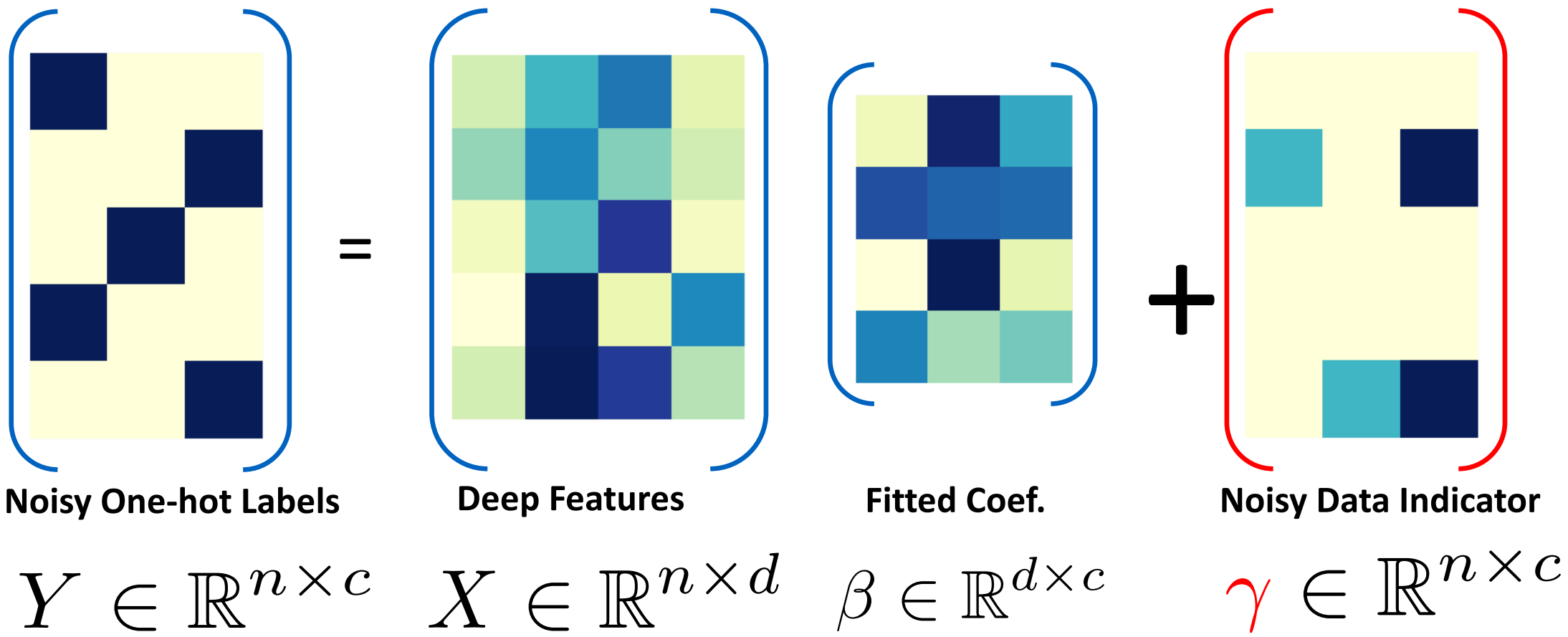$$\underset{\boldsymbol{\beta},\boldsymbol{\gamma}}{\mathrm{argmin}} L(\boldsymbol{\beta},\boldsymbol{\gamma}) := \|\boldsymbol{Y} - \boldsymbol{X}\boldsymbol{\beta} - \boldsymbol{\gamma}\|_{\mathrm{F}}^2 + \lambda R(\boldsymbol{\gamma})$$

$$\downarrow$$

$$\underset{\boldsymbol{\gamma}}{\mathrm{argmin}} \left\|\tilde{\boldsymbol{Y}} - \tilde{\boldsymbol{X}}\boldsymbol{\gamma}\right\|_{\mathrm{F}}^2 + \lambda R(\boldsymbol{\gamma})$$

$$\downarrow$$

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

**Noisy One-hot Labels**    **Deep Features**    **Fitted Coef.**    **Noisy Data Indicator**

$$Y \in \mathbb{R}^{n \times c} \quad X \in \mathbb{R}^{n \times d} \quad \beta \in \mathbb{R}^{d \times c} \quad \gamma \in \mathbb{R}^{n \times c}$$

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.

大数据学院
School of Data Science

# Noise Set Recovery

# When will the model identify all the outliers?

Assume $\boldsymbol{\varepsilon}$ is i.i.d zero-mean sub-Gaussian noise. We give three conditions:

- (C1: Restricted eigenvalue)

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

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

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

- (C3: Large error)

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

[1]: M. J. Wainwright, Sharp thresholds for high-dimensional and noisy sparsity recovery using l1-constrained quadratic programming. TIT 2009.

大数据学院
School of Data Science

# A non-asymptotic probabilistic result

Based on these conditions, we could provide the following theorem:
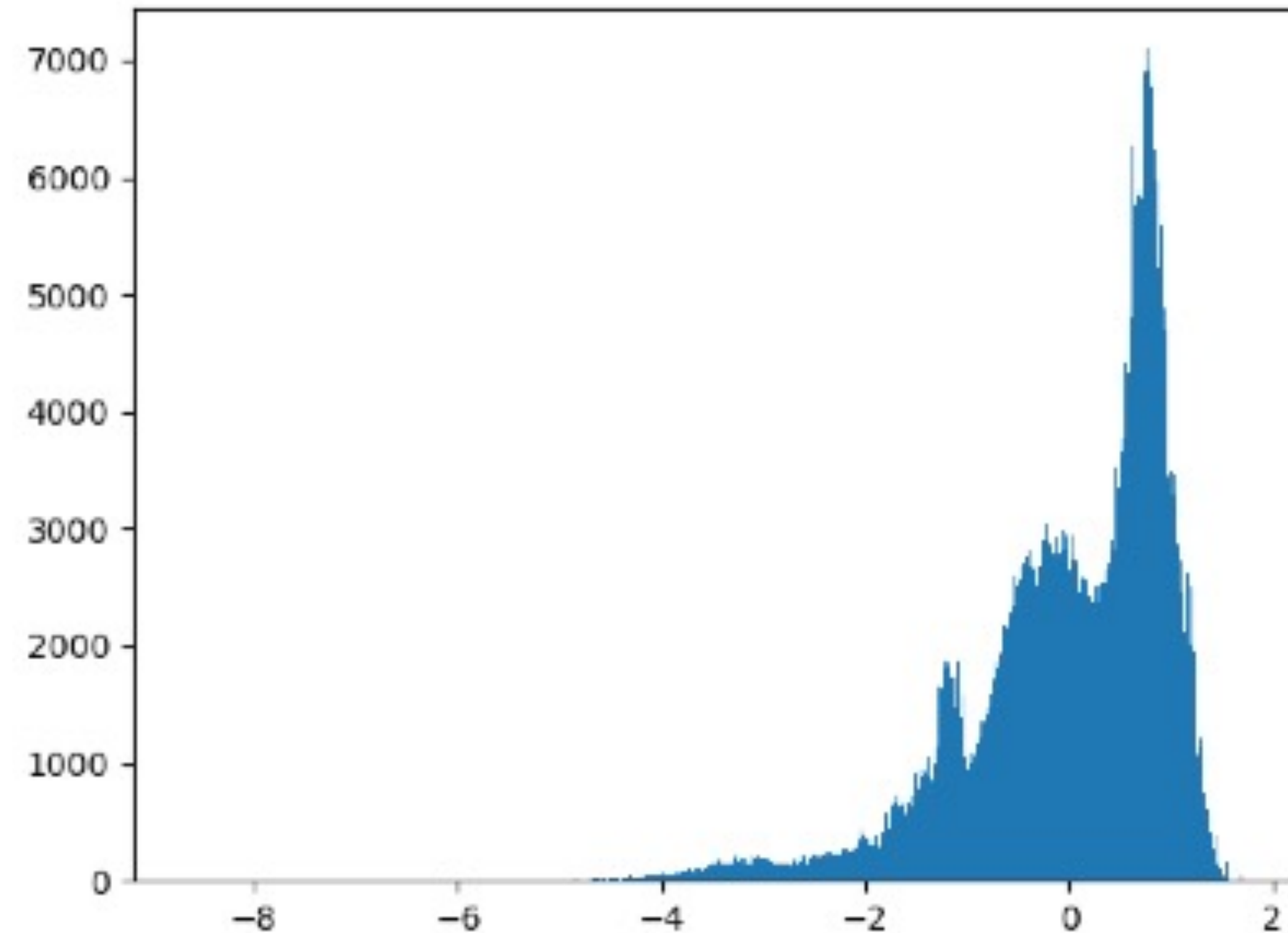
**Theorem 1** (Identifiability of ICI). *Let $\lambda \geq \frac{2\sigma\sqrt{\mu_{\tilde{U}}}}{\eta}\sqrt{\log cn}$. Then with probability greater than $1 - 2\left(cn\right)^{-1}$, the problem has a unique solution $\hat{\boldsymbol{\gamma}}$ satisfies the following properties:*

*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*

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

*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 $\mathrm{sign}\left(\hat{\vec{\boldsymbol{\gamma}}}\right) = \mathrm{sign}\left(\vec{\boldsymbol{\gamma}}^*\right)$).*

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

# Identifiability in reality: sub-Gaussian noise



$$\mathbb{E}\left[\hat{\varepsilon}\right] \approx 10^{-19}$$

$$\mathrm{Var}\left[\hat{\varepsilon}\right] \approx 0.99$$

# Sparse Learning

## in Few-Shot Learning

# Definition of Few-Shot Learning

Tackle machine learning problem with only limited training data provided.

**<span style="color:red">Few-Shot Learning</span>**

**Low cost** ───────────────────────────────────────────► **High cost**

**Unsupervised Learning**  **Learning From Noisy Labels**  **Semi-Supervised Learning**  **Supervised Learning**



Binary classification
with many labeled data

Few-shot binary classification

Few-shot binary classification
with <span style="color:red">unlabeled data</span>

# Motivation



Labeled Image

Labels

Train

Few-Shot Models

Unlabeled Image

Pseudo-Labels

Inference

Self-Taught Learning

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

大数据学院
School of Data Science

# Framework

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

# Sparse Learning in ICI

$$y_i = x_i^\top \beta + \varepsilon + \textcolor{red}{\gamma_i}$$



$$\underset{\boldsymbol{\beta}, \boldsymbol{\gamma}}{\operatorname{argmin}} L\left(\boldsymbol{\beta}, \boldsymbol{\gamma}\right) := \left\| \boldsymbol{Y} - \boldsymbol{X}\boldsymbol{\beta} - \boldsymbol{\gamma} \right\|_{\mathrm{F}}^2 + \lambda R\left(\boldsymbol{\gamma}\right)$$

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



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

# Sparse Learning: Extend to Logistic Regression

$$\underset{\boldsymbol{\beta},\boldsymbol{\gamma}}{\mathrm{argmin}} L\left(\boldsymbol{\beta},\boldsymbol{\gamma}\right) := \|\boldsymbol{Y} - \boldsymbol{X}\boldsymbol{\beta} - \boldsymbol{\gamma}\|_{\mathrm{F}}^2 + \lambda R\left(\boldsymbol{\gamma}\right)$$

$$\underset{\boldsymbol{\gamma}}{\mathrm{argmin}} \left\| \boldsymbol{Y} - \boldsymbol{X}\left(\boldsymbol{X}^\top \boldsymbol{X}\right)^\dagger \boldsymbol{X}^\top \left(\boldsymbol{Y} - \boldsymbol{\gamma}\right) - \boldsymbol{\gamma} \right\|_{\mathrm{F}}^2 + \lambda R\left(\boldsymbol{\gamma}\right)$$

$$\underset{\boldsymbol{\gamma}}{\mathrm{argmin}} \left\| \tilde{\boldsymbol{Y}} - \tilde{\boldsymbol{X}}\boldsymbol{\gamma} \right\|_{\mathrm{F}}^2 + \lambda R\left(\boldsymbol{\gamma}\right)$$

$$\boldsymbol{Y}_{i,c} = \frac{\exp\left(\boldsymbol{X}_{i,\cdot}\boldsymbol{\beta}_{\cdot,c} + \textcolor{red}{\boldsymbol{\gamma}_{i,c}}\right)}{\sum_{l=1}^{C} \exp\left(\boldsymbol{X}_{i,\cdot}\boldsymbol{\beta}_{\cdot,l} + \textcolor{red}{\boldsymbol{\gamma}_{i,l}}\right)} + \boldsymbol{\varepsilon}_{i,c}$$

$$\bar{\boldsymbol{X}} = \left(\boldsymbol{X}, \boldsymbol{I}\right) \qquad \bar{\boldsymbol{\beta}} = \left(\boldsymbol{\beta}, \boldsymbol{\gamma}\right)^\top$$

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

大数据学院
School of Data Science

# 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.

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

# 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%).

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

大数据学院
School of Data Science

# 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% |

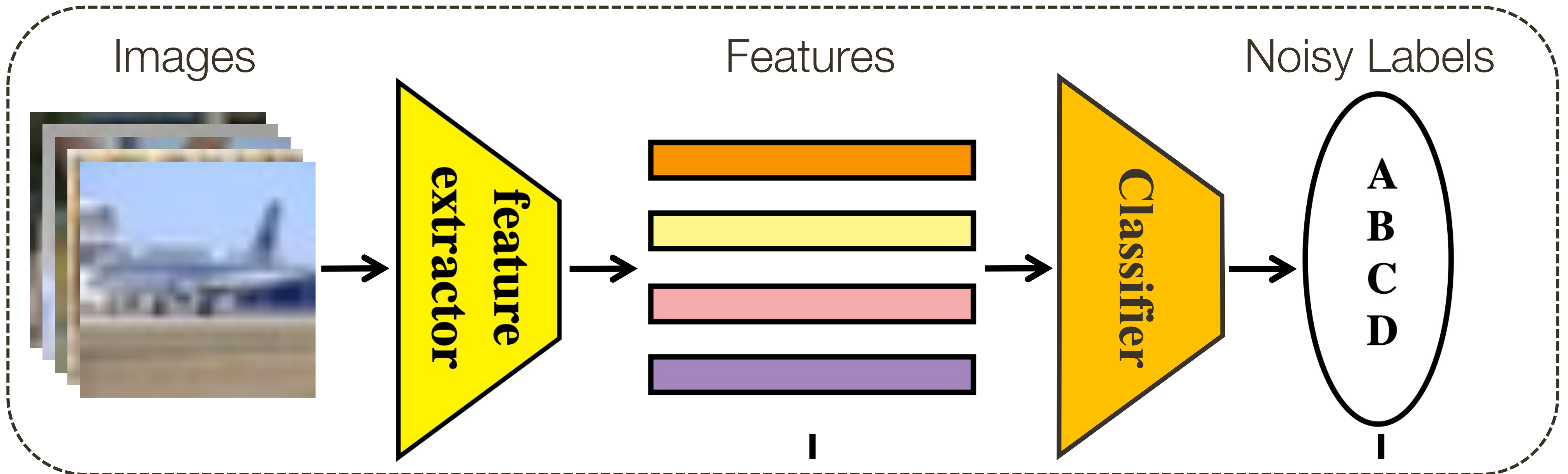3) Even if C2-C3 are not satisfied, we still have the chance of improving the performance (53.5%).

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

大数据学院
School of Data Science

# Sparse Learning

## in Learning with Noisy Labels

# Definition of learning with noisy labels



Deep Models → Robustly Trained Models

synthetic dataset: outlies(red); good observations(blue)

Shogun: Total War - IGN
ign.com

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

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

# Framework



Stage 1: Feature Learning

Images → feature extractor → Features → Classifier → Noisy Labels (A B C D)

Clean or Noisy?

$$y_i = \boldsymbol{x}_i^\top \beta + \gamma_i + \varepsilon$$

Stage 2: Sample Selection

Solution Path

$$\underset{\boldsymbol{\gamma}}{\mathrm{argmin}} \left\| \tilde{\boldsymbol{Y}} - \tilde{\boldsymbol{X}}\boldsymbol{\gamma} \right\|_{\mathrm{F}}^2 + \sum_{i=1}^{n} P\left(\gamma_i; \lambda_i\right)$$

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

# Make it scalable to large datasets



$$y_i = \boldsymbol{x}_i^\top \beta + \textcolor{red}{\gamma_i} + \varepsilon$$

Compute Class Similarity

Group Dissimilar classes

Split into pieces

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

# Strategies to help train the network

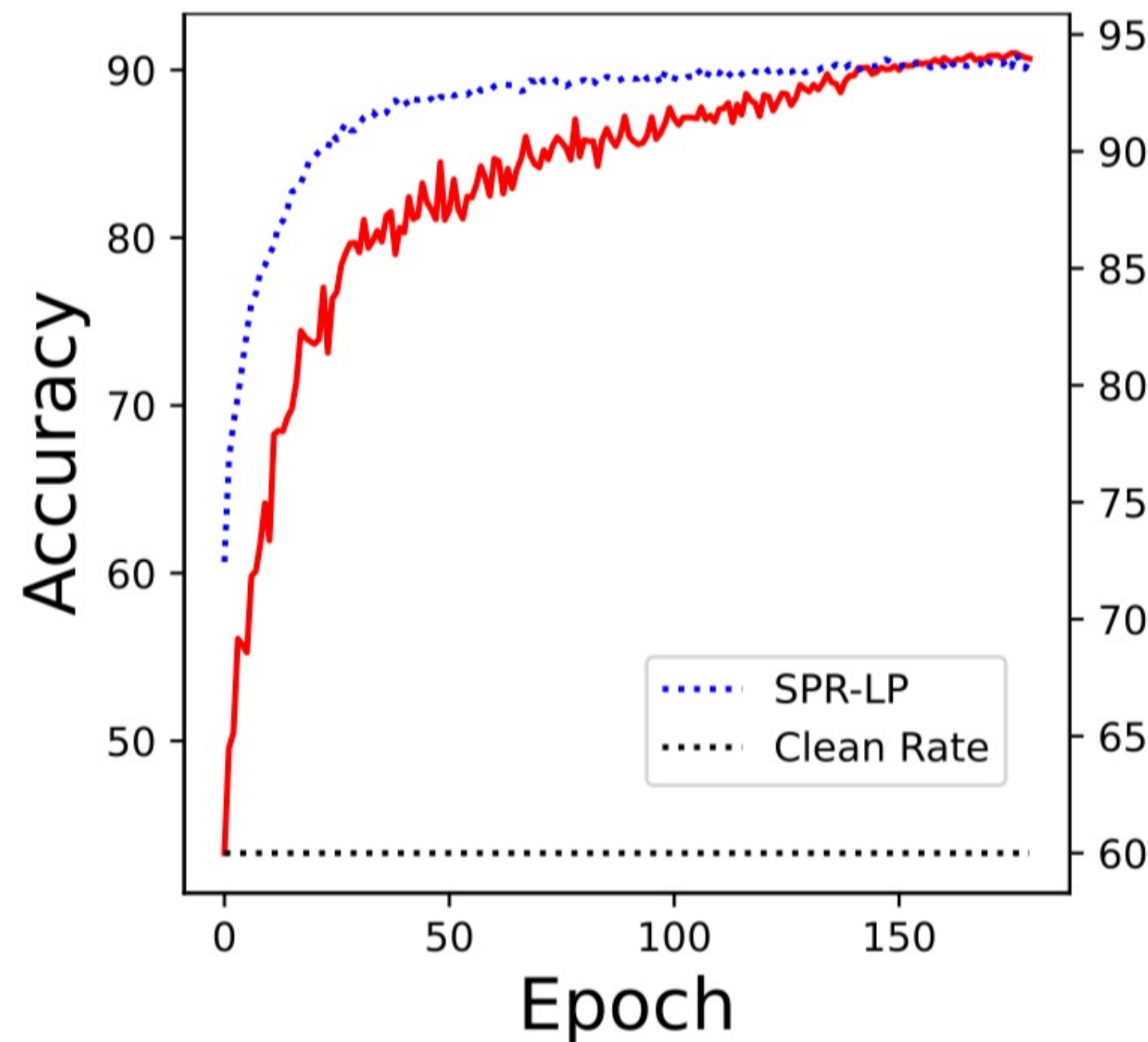- Append a $\ell_q(q < 1)$ penalty to encourage the linear relation between feature and one-hot encoded vector:

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

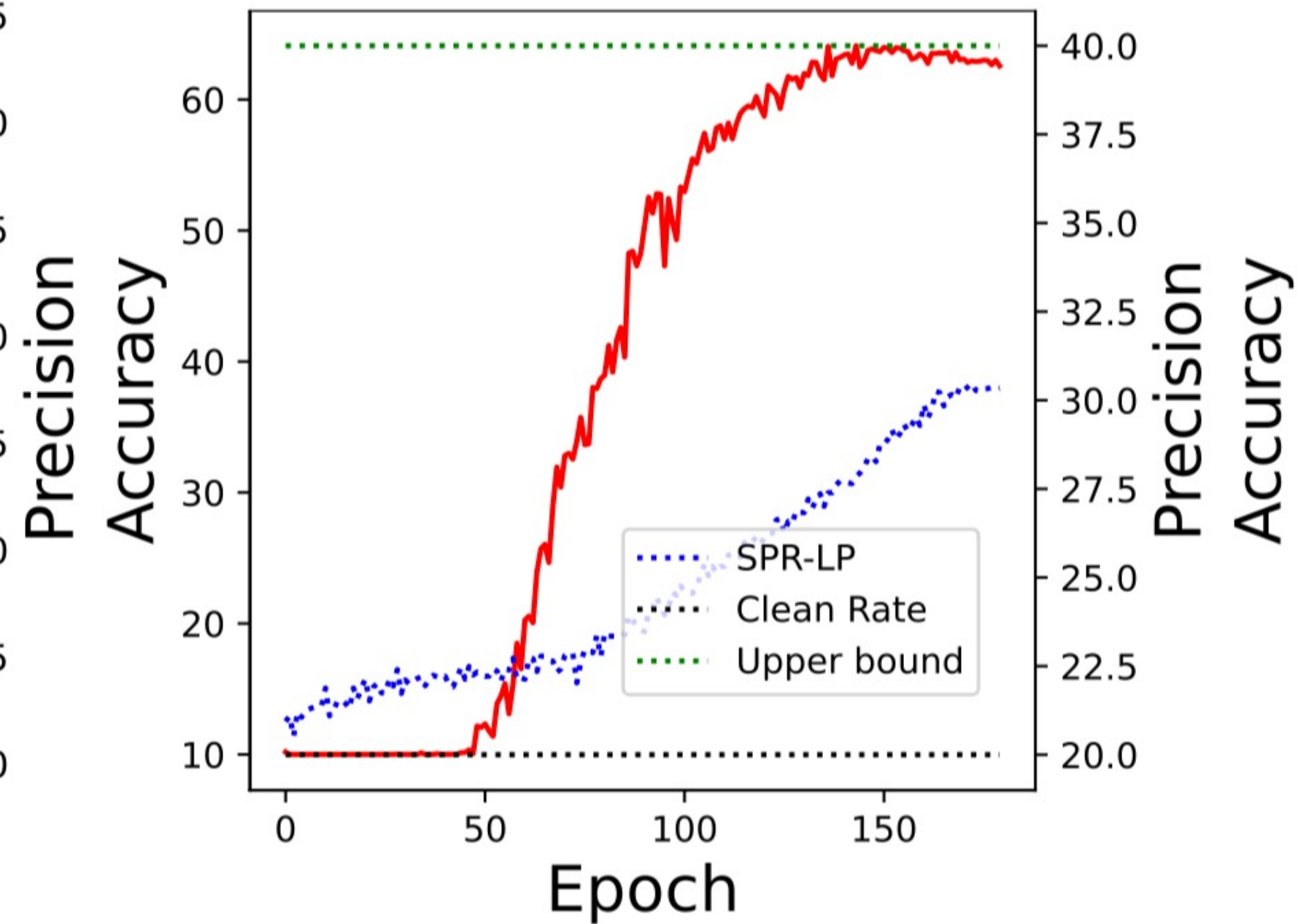- Use CutMix to further exploit the support of noisy data

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

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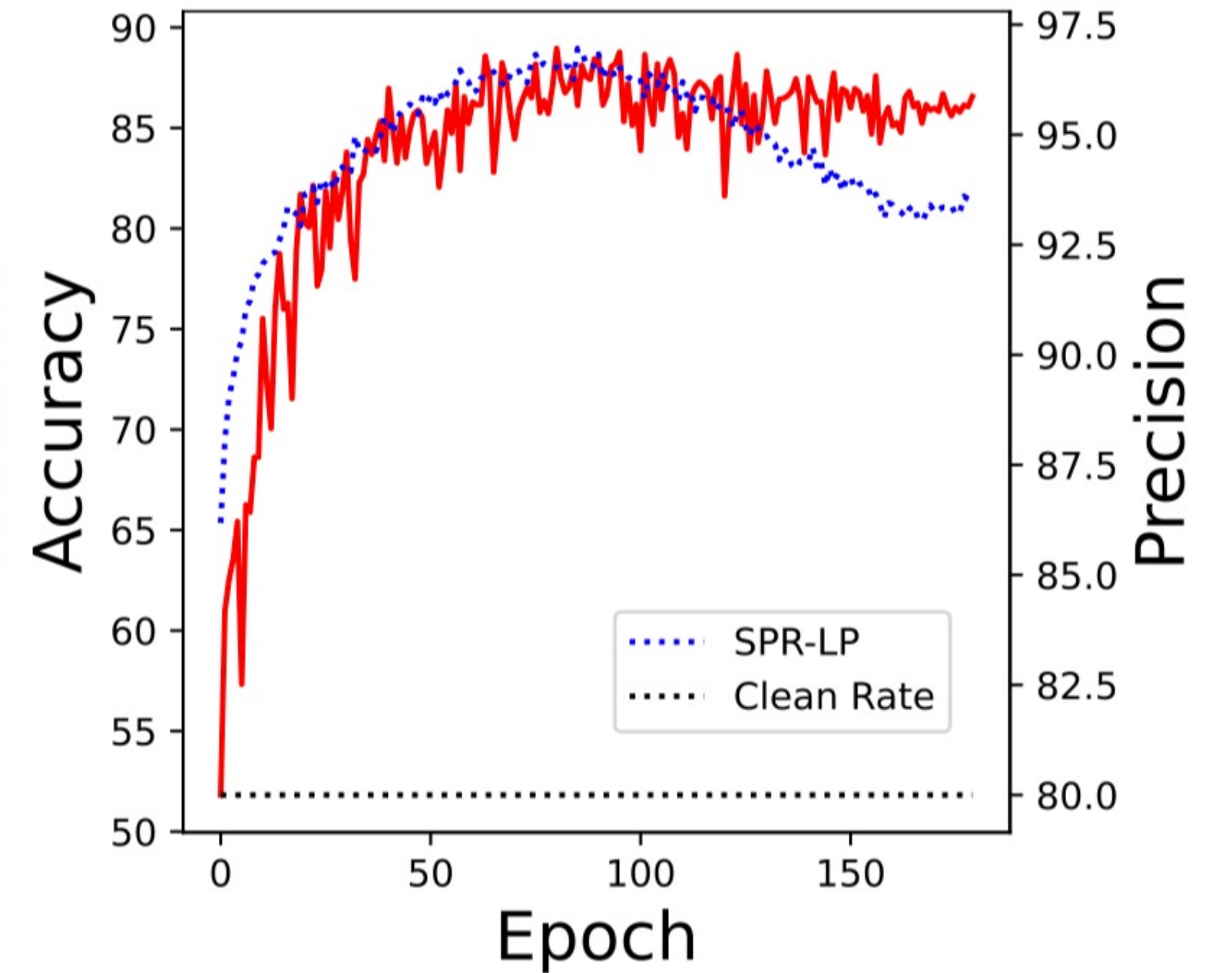
# Label precision performance



(a) Symmetric-40%   (b) Symmetric-80%   (c) Asymmetric-40%

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