# Taming Heavy-tailed Features by Shrinkage

Ziwei Zhu, Wenjing Zhou

University of Michigan, Department of Statistics

April 2021

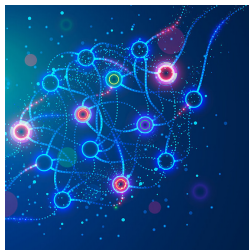# Motivation

► Heavy-tailed data abound in modern data analytics.



Figure 1: Real-life scenarios with heavy-tailed data

# Motivation

▶ Heavy-tailed features can aggravate the corruption on the response and jeopardize standard statistical approaches.



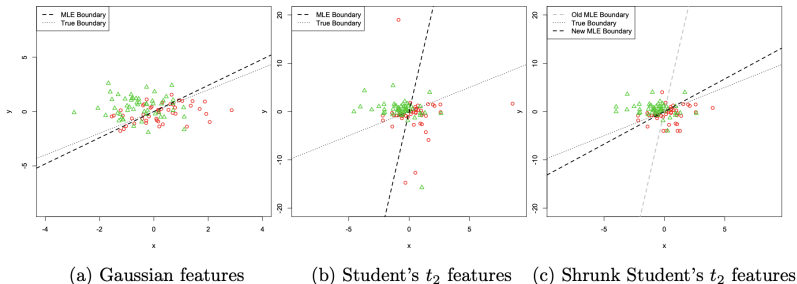(a) Gaussian features    (b) Student's $t_2$ features    (c) Shrunk Student's $t_2$ features

Figure 2: Logistic regression with 10% mislabeled data on different features

# Generalized linear model (GLM)

▶ Suppose we have $n$ observations $\{(y_i, x_i)\}_{i=1}^n$, where $y_i$ is the response and $x_i$ is the feature vector valued in $\mathbb{R}^d$. Under the GLM with the canonical link, the probability density function of the response $y_i$ is defined as

$$f_n(y; X, \beta^*) = \prod_{i=1}^n f(y_i; \eta_i^*) = \prod_{i=1}^n \left\{ c(y_i) \exp\left( \frac{y_i \eta_i^* - b(\eta_i^*)}{\phi} \right) \right\}, \tag{1}$$

where $y = (y_1, \cdots, y_n)^\top$, $X = (x_1, \cdots, x_n)^\top$, $\beta^* \in \mathbb{R}^d$ is the regression coefficient vector, $\eta_i^* := x_i^\top \beta^*$, $b(\cdot)$ is a known function that is twice differentiable with a positive second derivative and $\phi > 0$ is the dispersion parameter.

▶ The response $y$ is assumed to be generated from a particular distribution in an exponential family.

# Corrupted generalized linear model (CGLM)

- For the $i$th observation, only corrupted response $z_i = y_i + \epsilon_i$ can be observed, where $\epsilon_i$ is random noise.
- The response is not limited within the exponential family.
- More real-world problems with complex structures:
  - the linear regression model with heavy-tailed noise
  - the logistic regression with mislabeled samples.
- The flexibility of the original GLM is significantly improved.

# Shrinkage

- Low-dimensional regime:

  $\ell_4$-norm shrinkage on features, clipping on response

$$\widetilde{x}_i := \frac{\min(\|x_i\|_4, \tau_1)}{\|x_i\|_4} x_i$$

$$\widetilde{z}_i := \min(|z_i|, \tau_2) z_i / |z_i|$$

- High-dimensional regime:

  elementwise shrinkage on features, clipping on response

$$\widetilde{x}_{ij} := \min(|x_{ij}|, \tau_1) x_{ij} / |x_{ij}|.$$

$$\widetilde{z}_i := \min(|z_i|, \tau_2) z_i / |z_i|$$

- $\tau_1$ and $\tau_2$ are predetermined thresholds

# $\ell_4$-norm ball

- ▶ The norm determines the strength of the constraint.
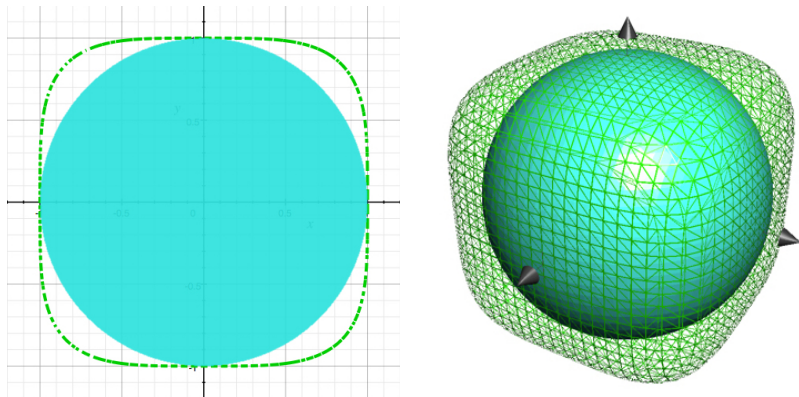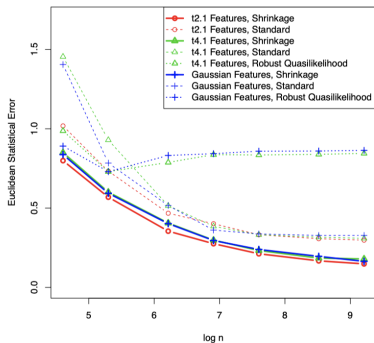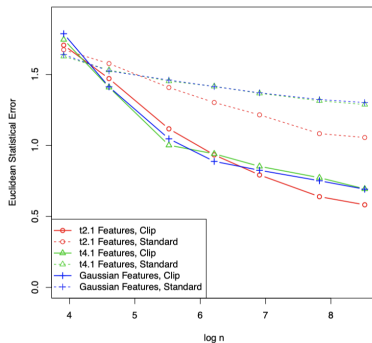- ▶ $\ell_4$-norm shrinkage balances the bias and the variance.



Figure 3: Unit Euclidean ball and $\ell_4$-norm ball in 2D and 3D

# Simulation: logistic regression with mislabeled data



Figure 4: Logistic regression with 10% mislabeled data

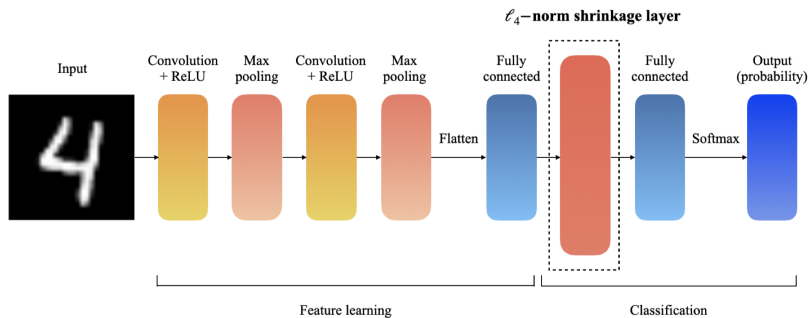# Simulation: CNN on MNIST dataset



Figure 5: Architecture of the shrinkage convolutional neural network

# Simulation: CNN on MNIST dataset

Table 1: Average testing misclassification rate (with standard error in the parentheses) on noisy MNIST images under mislabeling probability 40%

| Noisy Pixel Ratio | Original CNN | Shrinkage CNN |
|:---:|:---:|:---:|
| 0 | $3.64\%_{(0.20\%)}$ | $2.93\%_{(0.09\%)}$ |
| 0.1 | $6.88\%_{(0.22\%)}$ | $4.18\%_{(0.17\%)}$ |
| 0.2 | $6.90\%_{(0.21\%)}$ | $4.37\%_{(0.16\%)}$ |
| 0.4 | $10.69\%_{(0.29\%)}$ | $6.65\%_{(0.24\%)}$ |
| 0.6 | $18.82\%_{(0.88\%)}$ | $12.80\%_{(0.65\%)}$ |

# Taming Heavy-tailed Features by Shrinkage