

On the sample complexity of parameter estimation in logistic regression with normal design

Daniel Hsu (Columbia)

Arya Mazumdar (UCSD)

COLT 2024

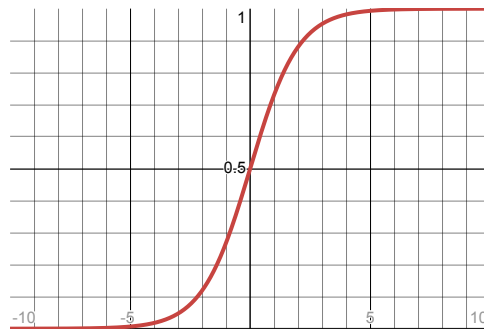


Data model for noisy binary classification

- $(X, Y) \sim P_w$: covariate vector $X \sim \mathbf{N}(0, I_d)$; binary response Y

$$\Pr_w(Y = y \mid X = x) = \frac{1}{1 + e^{-y\langle x, w \rangle}}, \quad \forall x \in \mathbb{R}^d, y \in \{-1, 1\}$$

"logistic regression"



- Parameter $w \in \mathbb{R}^d$
- **Estimation goal:** Given i.i.d. sample from P_{w^*} (w^* unknown), construct estimate \hat{w} such that

$$\|\hat{w} - w^*\| \leq \epsilon$$

How large should the sample size be?

Clues from classical asymptotic theory?

- Maximum likelihood estimator given data $(x_i, y_i)_{i=1}^n$

$$\hat{w}_{\text{mle}} = \arg \min_w \sum_{i=1}^n \ln(1 + e^{-y_i \langle x_i, w \rangle})$$

MLE may not exist!

- Asymptotically (as $n \rightarrow \infty$),

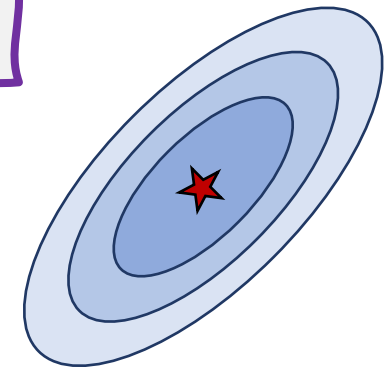
$$\sqrt{n}(\hat{w}_{\text{mle}} - w^*) \xrightarrow{\text{dist.}} \text{N}(0, \mathcal{J}(w^*)^{-1})$$

- Very roughly:

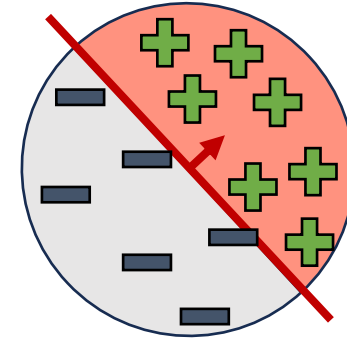
$$\mathbb{E} \|\hat{w}_{\text{mle}} - w^*\| \rightarrow \sqrt{d/n}$$

Dependence on $\|w^*\|$?

- "Conclusion": sample complexity is d/ϵ^2 ???



Learning half-spaces



- As $\|w^*\| \rightarrow \infty$, response Y is determined by X :

$$Y = \text{sign}(\langle X, \theta^* \rangle)$$

where $\theta^* = w^* / \|w^*\|$

- PAC learning homogeneous half-spaces under $X \sim \text{Uniform}(S^{d-1})$
- **Long (1995, 2003):** **sample complexity** is d/ϵ (cf. d/ϵ^2)
 - ... to guarantee **classification error rate** $\leq \epsilon$
 - ... which is proportional to **parameter error** $\|\hat{\theta} - \theta^*\|$
- AKA "1-bit compressed sensing"

Question

- What is the role of $\|w^*\|$?

- Fix $\|w^*\| = \beta$, and only consider estimating $\theta^* = w^*/\|w^*\| \in S^{d-1}$
- β is akin to **signal-to-noise ratio**, also called **inverse temperature**

$$\Pr_{\beta\theta^*}(Y = 1 \mid X = x) = \frac{1}{1 + \exp(-\beta\langle x, \theta^* \rangle)}$$

- No signal ($\beta = 0$): hopeless
- No noise ($\beta = \infty$): PAC learning half-spaces

- **Revised goal:** Given i.i.d. sample from $P_{\beta\theta^*}$ for some unknown θ^* , construct estimate $\hat{\theta}$ such that

$$\|\hat{\theta} - \theta^*\| \leq \epsilon$$

Main result

- **Sample complexity*** to ensure $\|\hat{\theta} - \theta^*\| \leq \epsilon$ (in expectation or w.h.p.):

$$n^*(d, \epsilon, \beta) \asymp \begin{cases} \frac{d}{\beta^2 \epsilon^2} & \text{if } \beta \lesssim 1 & \text{"high temperature"} \\ \frac{d}{\beta \epsilon^2} & \text{if } 1 \lesssim \beta \lesssim 1/\epsilon & \text{"moderate temperature"} \\ \frac{d}{\epsilon} & \text{if } 1/\epsilon \lesssim \beta & \text{"low temperature"} \end{cases}$$

*up to logarithmic factors in d and $1/\epsilon$

Logistic loss

- Logistic loss (i.e., negative log-likelihood of $\text{Ber}(\beta \langle x, \theta \rangle)$ on (x, y)):

$$\ell(\theta; x, y) = \ln(1 + e^{-\beta y \langle x, \theta \rangle})$$

- Excess risk with logistic loss:

$$\mathbb{E}[\ell(\theta; X, Y) - \ell(\theta^*; X, Y)] = \mathbb{E}[\text{KL}(\text{Ber}(\beta \langle X, \theta^* \rangle) \parallel \text{Ber}(\beta \langle X, \theta \rangle))]$$

- **Normal design** \rightarrow very good estimates of expected KL divergence

Sample complexity lower bound

- To use Fano's inequality, suffices to prove good upper bound on $\mathbb{E}[\text{KL}(\text{Ber}(\beta \langle X, \theta^* \rangle) \parallel \text{Ber}(\beta \langle X, \theta \rangle))]$

- High temp ($\beta \lesssim 1$): textbook exercise

$$n^*(d, \epsilon, \beta) \gtrsim \frac{d}{\beta^2 \epsilon^2}$$

- Moderate temp ($1 \lesssim \beta \lesssim 1/\epsilon$): not well-known?

$$n^*(d, \epsilon, \beta) \gtrsim \frac{d}{\beta \epsilon^2}$$

- Low temp ($1/\epsilon \lesssim \beta$): unclear how to get tight bound with Fano

Instead, extend Long's 1995 lower bound for $\beta = \infty$ to all $\beta \gtrsim 1/\epsilon$

Sample complexity upper bound

- Three different estimators, depending on temperature
 - High temp ($\beta \lesssim 1$): minimize average **linear loss** (Servedio, 1999: "Average" algorithm)

Also: Plan & Vershynin (2012)

$$\hat{\theta} = \arg \min_{\theta \in \mathcal{S}^{d-1}} \sum_{i=1}^n -y_i \langle x_i, \theta \rangle$$

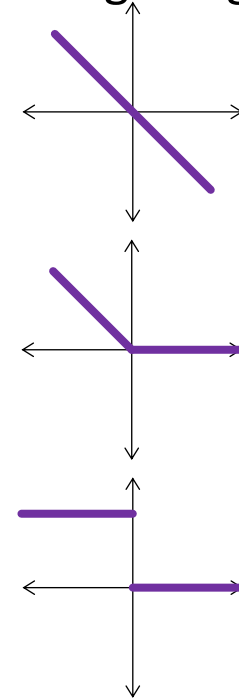
- Moderate or low temp ($1 \lesssim \beta$): minimize average **ReLU loss**

Inspired by Kuchelmeister & van de Geer (2023)

$$\hat{\theta} = \arg \min_{\theta \in \mathcal{S}^{d-1}} \sum_{i=1}^n \max\{0, -y_i \langle x_i, \theta \rangle\}$$

- Low temp ($1/\epsilon \lesssim \beta$): minimize average **0-1 loss**

$$\hat{\theta} = \arg \min_{\theta \in \mathcal{S}^{d-1}} \sum_{i=1}^n 1\{y_i \langle x_i, \theta \rangle \leq 0\}$$



What we couldn't get to work

- Minimize average **logistic loss** (i.e., MLE)
 - Taylor-expand the estimation error (Portnoy, 1988; He & Shao, 2000; ...)
 - Use self-concordance of logistic loss (Bach, 2010; Ostrovskii & Bach, 2021)
 - Our attempts gave suboptimal dependence on β



Recap and open problems

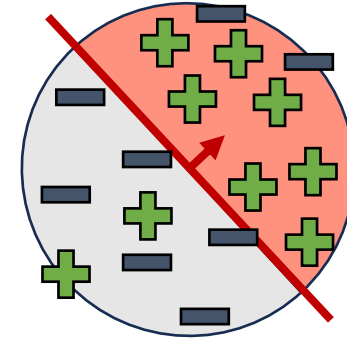
- Two "change points" in sample complexity for logistic regression

$$n^*(d, \epsilon, \beta) \asymp \begin{cases} \frac{d}{\beta^2 \epsilon^2} & \text{if } \beta \lesssim 1 & \text{"high temperature"} \\ \frac{d}{\beta \epsilon^2} & \text{if } 1 \lesssim \beta \lesssim 1/\epsilon & \text{"moderate temperature"} \\ \frac{d}{\epsilon} & \text{if } 1/\epsilon \lesssim \beta & \text{"low temperature"} \end{cases}$$

- Q: Efficient algorithms? MLE? Estimation of $\|w^*\|$?

Thank you!

Learning noisy half-spaces



- (Distribution-free) agnostic PAC learning half-spaces

- **VC theory:** To ensure $\leq \epsilon$ **excess classification error rate**

$$\text{err}(\hat{\theta}) - \text{err}(\theta^*) \leq \epsilon$$

sample complexity is at most $d(1/\epsilon + \text{err}(\theta^*)/\epsilon^2)$ (up to logs)

- But we want guarantee about **parameter error** $\|\hat{\theta} - \theta^*\|$
 - Can relate in low temp ($1/\epsilon \lesssim \beta$) regime, but unclear otherwise
 - **Useful fact:** $\text{err}(\theta^*) \approx 1/\beta$ when $\beta \gtrsim 1$

Bregman divergence

- Bernoulli distribution $\text{Ber}(\eta)$ has "mean parameter" $g'(\eta) = \frac{1}{1+e^{-\eta}}$;
 - $g(\eta) = \ln(1 + e^\eta)$ is log partition function; g' is its derivative

- KL between Bernoulli distributions as Bregman divergence:

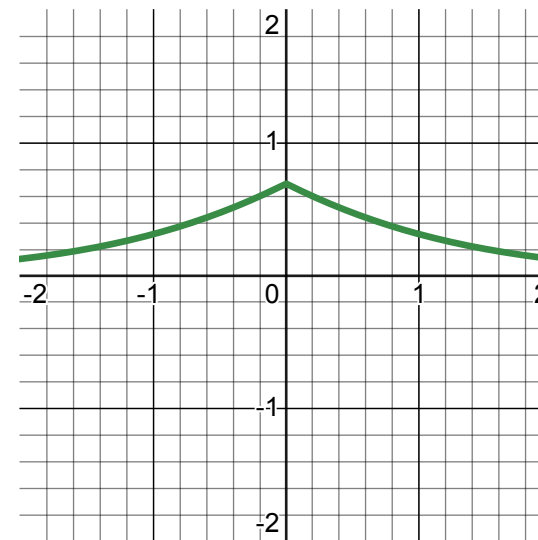
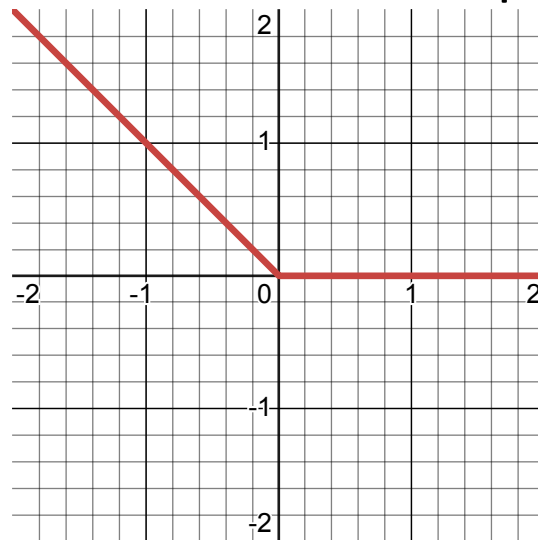
$$\text{KL}(\text{Ber}(\eta^*) \parallel \text{Ber}(\eta)) = g(\eta) - g(\eta^*) - g'(\eta^*)(\eta - \eta^*)$$

- When $\eta^* = \beta \langle X, \theta^* \rangle$ and $\eta = \beta \langle X, \theta \rangle$ and $X \sim \mathbf{N}(0, I_d)$:

$$\begin{aligned} & \mathbb{E}[\text{KL}(\text{Ber}(\beta \langle X, \theta^* \rangle) \parallel \text{Ber}(\beta \langle X, \theta \rangle))] \\ &= \beta \mathbb{E}[g'(\beta \langle X, \theta^* \rangle) \langle X, \theta - \theta^* \rangle] \end{aligned}$$

ReLU loss

- Nice observation of Kuchelmeister and van de Geer (2023):
$$\ln(1 + e^{-\beta y \langle x, \theta \rangle}) = \text{ReLU}(-\beta y \langle x, \theta \rangle) + \ln(1 + e^{-\beta |\langle x, \theta \rangle|})$$
- (Scaled) excess risk with ReLU loss = excess risk with logistic loss
 - Uses spherical symmetry of $\mathbf{N}(0, I_d)$
 - **Caveat:** optimization over the sphere



Adaptivity

- **If β is unknown:** suffices to coarsely distinguish "high temp" ($\beta \lesssim 1$) and "medium or low temp" ($1 \lesssim \beta$) regimes
 - Estimate classification error rate of θ^*
 - Can use training error rate of ERM (with zero-one loss) on dataset of size d/ϵ
 - Based on outcome, decide whether to use linear loss or ReLU loss on full data

