

Exponential tail bounds and Large Deviation Principle for Heavy-Tailed U-Statistics

Milad Bakhshizadeh

Stanford University

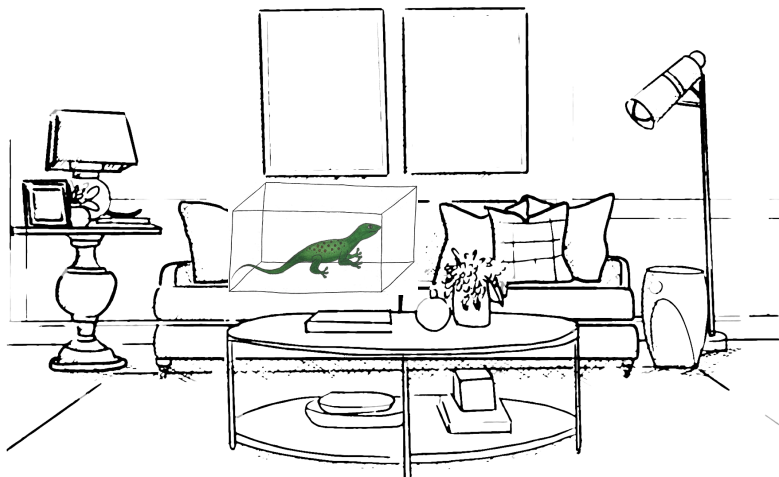
May 24, 2023

Uncertainty

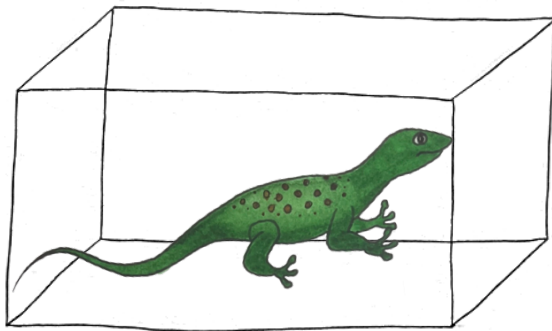
Uncertainty



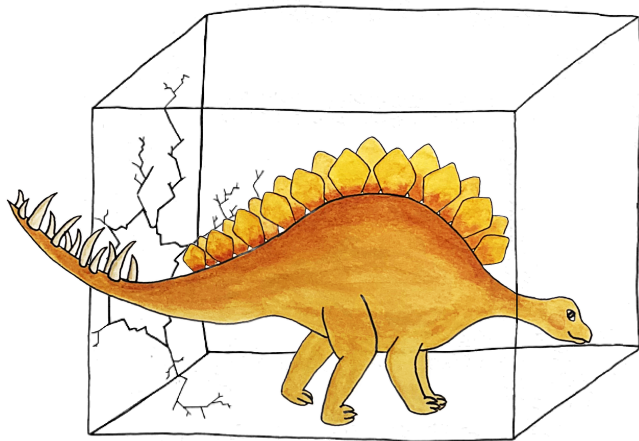
Uncertainty



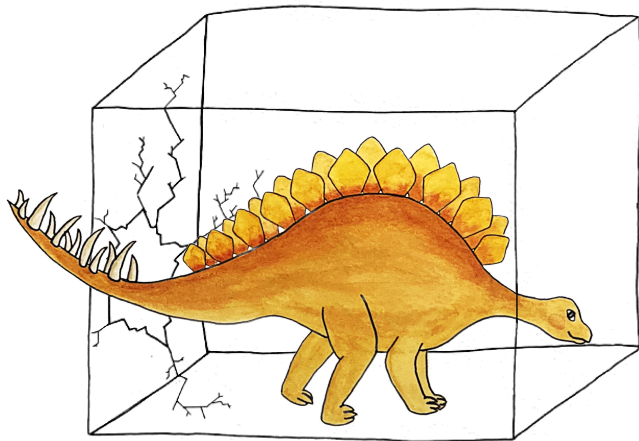
Uncertainty



Uncertainty



Uncertainty : How to make the Dinosaur a lizard?!

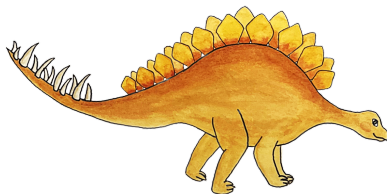


What is this talk about?

1. Heavy-tail makes it challenging to bound uncertainty
2. We can control heavy-tail by truncation
3. Truncation gives optimal bound

Part 1: problem setup

- ▶ Heavy-tail distributions: the dinosaur



Heavy-tails don't have finite MGF

▶ SubGaussian distributions

▶ $\mathbb{P}(|X| > t) \leq \exp(-ct^2)$

▶ $\|X\|_p = O(\sqrt{p})$

▶ $\mathbb{E}[\exp(\lambda X)] \leq \exp(c\lambda^2)$



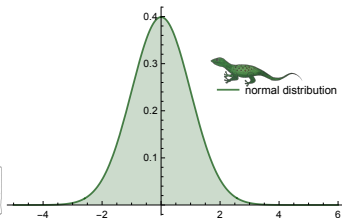
▶ Heavy-tailed distributions

▶ $\mathbb{E}[\exp(\lambda X)] = \infty, \forall \lambda > 0$



▶ e.g. Weibull with $k < 1$,

$\mathcal{N}(0, 1)^\alpha$ $\alpha > 2$, Log-Normal, ...



Heavy-tails don't have finite MGF

▶ SubGaussian distributions

▶ $\mathbb{P}(|X| > t) \leq \exp(-ct^2)$

▶ $\|X\|_p = O(\sqrt{p})$

▶ $\mathbb{E}[\exp(\lambda X)] \leq \exp(c\lambda^2)$



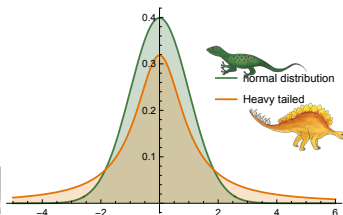
▶ Heavy-tailed distributions

▶ $\mathbb{E}[\exp(\lambda X)] = \infty, \forall \lambda > 0$



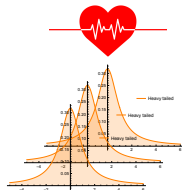
▶ e.g. Weibull with $k < 1$,

$\mathcal{N}(0, 1)^\alpha \alpha > 2$, Log-Normal, ...



Rise of the heavy-tails: A new era dawns

- ▶ Data is heavy-tailed
- ▶ Multiplication makes tail heavier
 - ▶ XY, X^n
 - ▶ $\mathcal{N}(0, 1), \mathcal{N}(0, 1)^2, \mathcal{N}(0, 1)^3$
- ▶ Real applications
 - ▶ Neural nets
 - ▶ Phase retrieval



Rise of the heavy-tails: A new era dawns

- ▶ Data is heavy-tailed
- ▶ Multiplication makes tail heavier
 - ▶ XY, X^n
 - ▶ $\mathcal{N}(0, 1), \mathcal{N}(0, 1)^2, \mathcal{N}(0, 1)^3$
- ▶ Real applications
 - ▶ Neural nets
 - ▶ Phase retrieval

Preprints

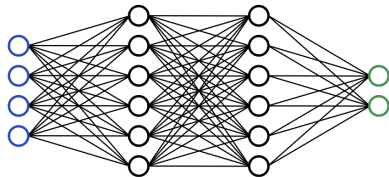
LAST UPDATE	AUTHORS	TITLE
Apr. 2023	M. Bakhshizadeh	Algebra of Sub-Weibull Random Variables (download)

Rise of the heavy-tails: A new era dawns

- ▶ Data is heavy-tailed
- ▶ Multiplication makes tail heavier
 - ▶ XY, X^n
 - ▶ $\mathcal{N}(0, 1), \mathcal{N}(0, 1)^2, \mathcal{N}(0, 1)^3$
- ▶ Real applications
 - ▶ Neural nets
 - ▶ Phase retrieval

Preprints

LAST UPDATE	AUTHORS	TITLE
Apr. 2023	M. Bakhshizadeh	Algebra of Sub-Weibull Random Variables (download)

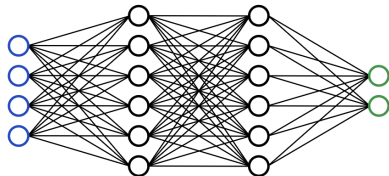


Rise of the heavy-tails: A new era dawns

- ▶ Data is heavy-tailed
- ▶ Multiplication makes tail heavier
 - ▶ XY, X^n
 - ▶ $\mathcal{N}(0, 1), \mathcal{N}(0, 1)^2, \mathcal{N}(0, 1)^3$
- ▶ Real applications
 - ▶ Neural nets
 - ▶ Phase retrieval

Preprints

LAST UPDATE	AUTHORS	TITLE
Apr. 2023	M. Bakhshizadeh	Algebra of Sub-Weibull Random Variables (download)



U-statistics, a low risk unbiased estimator

- ▶ $U_n = \frac{1}{\binom{n}{m}} \sum_{1 \leq i_1 < i_2 < \dots < i_m \leq n} h(X_{i_1}, \dots, X_{i_m}), \quad h(\cdot)$ symmetric
 - ▶ X_i iid
 - ▶ $U_n \rightarrow \mathbb{E}[h]$
 - ▶ $h = X_1 \rightarrow \hat{\mu} = \bar{X}_n$
 - ▶ $h = \frac{(X_1 - X_2)^2}{2} \rightarrow \hat{\sigma}^2$
- ▶ How fast $\mathbb{P}(|U_n - \mathbb{E}[h]| > \epsilon) \rightarrow 0$?
 - ▶ Sample size n
 - ▶ High-dimensional statistics

U-statistics, a low risk unbiased estimator

- ▶ $U_n = \frac{1}{\binom{n}{m}} \sum_{1 \leq i_1 < i_2 < \dots < i_m \leq n} h(X_{i_1}, \dots, X_{i_m}), \quad h(\cdot)$ symmetric
 - ▶ X_i iid
 - ▶ $U_n \rightarrow \mathbb{E}[h]$
 - ▶ $h = X_1 \rightarrow \hat{\mu} = \bar{X}_n$
 - ▶ $h = \frac{(X_1 - X_2)^2}{2} \rightarrow \hat{\sigma}^2$
- ▶ How fast $\mathbb{P}(|U_n - \mathbb{E}[h]| > \epsilon) \rightarrow 0$?
 - ▶ Sample size n
 - ▶ High-dimensional statistics

U-statistics, a low risk unbiased estimator

- ▶ $U_n = \frac{1}{\binom{n}{m}} \sum_{1 \leq i_1 < i_2 < \dots < i_m \leq n} h(X_{i_1}, \dots, X_{i_m}), \quad h(\cdot)$ symmetric
 - ▶ X_i iid
 - ▶ $U_n \rightarrow \mathbb{E}[h]$
 - ▶ $h = X_1 \rightarrow \hat{\mu} = \bar{X}_n$
 - ▶ $h = \frac{(X_1 - X_2)^2}{2} \rightarrow \hat{\sigma}^2$
- ▶ How fast $\mathbb{P}(|U_n - \mathbb{E}[h]| > \epsilon) \rightarrow 0$?
 - ▶ Sample size n
 - ▶ High-dimensional statistics

Concentration inequality, a tool for uncertainty control

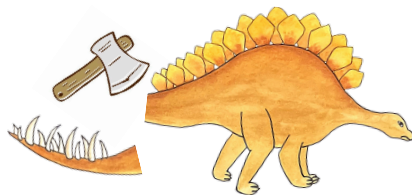
- ▶ $\mathbb{P}(|U_n| > \epsilon) \leq \exp(-L(n, \epsilon))$
 - ▶ $\mathbb{E}[h] = 0$ (no generality loss)
 - ▶ Simple
 - ▶ Tight (asymptotically)
- ▶ Recall
 - ▶ $U_n = \frac{1}{\binom{n}{m}} \sum h(X_{i_1}, \dots, X_{i_m})$
 - ▶ $\mathbb{E}[\exp(\lambda h(X_1, \dots, X_m))] = \infty, \quad \forall \lambda > 0$

Concentration inequality, a tool for uncertainty control

- ▶ $\mathbb{P}(|U_n| > \epsilon) \leq \exp(-L(n, \epsilon))$
 - ▶ $\mathbb{E}[h] = 0$ (no generality loss)
 - ▶ Simple
 - ▶ Tight (asymptotically)
- ▶ Recall
 - ▶ $U_n = \frac{1}{\binom{n}{m}} \sum h(X_{i_1}, \dots, X_{i_m})$
 - ▶ $\mathbb{E} \left[\exp(\lambda h(X_1, \dots, X_m)) \right] = \infty, \quad \forall \lambda > 0$

Part 2

► The Solution: Tail Truncation



Bound tail and body, separately

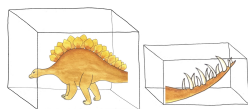
► Define:

► $k = \lfloor \frac{n}{m} \rfloor$

► $h_L = h \mathbb{1}(h \leq L)$

► $\mathbb{P}(h > t) \simeq \exp(-I(t)), I(t) \ll t$

► $\mathbb{P}(U_n > t) \leq \mathbb{P}(U_n(h_L) > t) + \mathbb{P}(\exists ij | h(X_{i_1}, \dots, X_{i_m}) > L)$



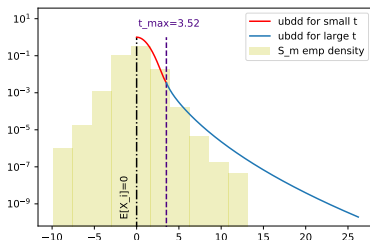
Theorem (1)

$$\mathbb{P}(U_n > t) \lesssim \exp\left(-\frac{kt^2}{2\text{Var}(h)}\right) + (1 + \binom{n}{m}) \exp(-I(kt))$$

There are two different regions of deviation

$$\mathbb{P}(U_n > t) \lesssim \exp\left(-\frac{kt^2}{2\text{Var}(h)}\right) + \left(1 + \binom{n}{m}\right) \exp(-I(kt))$$

- ▶ Regions:
 - ▶ Small t , Gaussian decay
 - ▶ Large t , like $\exp(-I(kt))$
- ▶ Change point: $kt^2 \simeq I(kt)$



Derivation

$$\begin{aligned} \mathbb{P}(U_n > t) &\leq \mathbb{P}(U_n(h_L) > t) + \mathbb{P}(\exists ij \mid h(X_{i_1}, \dots, X_{i_m}) > L) \\ &\leq e^{-\lambda t} \mathbb{E} \left[e^{\lambda U_n(h_L)} \right] + \binom{n}{m} e^{-I(L)} \end{aligned}$$

► issues

- dependent terms in $U_n(h_L)$
- $\mathbb{E} \left[e^{\lambda U_n(h_L)} \right] \xrightarrow{L \rightarrow \infty} \infty$, fixed λ
- Optimize λ, L together
- Sharpness

$$\text{First term} \leq \exp \left(-\frac{kt^2}{2v(kt, \beta \frac{I(kt)}{kt})} \right) + \exp \left(-\beta I(kt) \max\left(\frac{1}{2}, c(t, \beta, k)\right) \right)$$

- $v(L, \eta) \triangleq \mathbb{E} \left[h_L^2 \mathbf{1}(h \leq 0) + h_L^2 \exp(\eta h_L) \mathbf{1}(h > 0) \right] \rightarrow \text{Var}(h)$
- $c(t, \beta, k) \triangleq 1 - \frac{\beta}{2t} \frac{I(kt)}{kt} v(kt, \beta \frac{I(kt)}{kt}) \rightarrow 1$

Derivation

$$\begin{aligned} \mathbb{P}(U_n > t) &\leq \mathbb{P}(U_n(h_L) > t) + \mathbb{P}(\exists ij \mid h(X_{i_1}, \dots, X_{i_m}) > L) \\ &\leq e^{-\lambda t} \mathbb{E} \left[e^{\lambda U_n(h_L)} \right] + \binom{n}{m} e^{-I(L)} \end{aligned}$$

► issues

- dependent terms in $U_n(h_L)$
- $\mathbb{E} \left[e^{\lambda U_n(h_L)} \right] \xrightarrow{L \rightarrow \infty} \infty$, fixed λ
- Optimize λ, L together
- Sharpness

$$\text{First term} \leq \exp \left(-\frac{kt^2}{2v(kt, \beta \frac{I(kt)}{kt})} \right) + \exp \left(-\beta I(kt) \max\left(\frac{1}{2}, c(t, \beta, k)\right) \right)$$

$$\text{► } v(L, \eta) \triangleq \mathbb{E} \left[h_L^2 \mathbf{1}(h \leq 0) + h_L^2 \exp(\eta h_L) \mathbf{1}(h > 0) \right] \rightarrow \text{Var}(h)$$

$$\text{► } c(t, \beta, k) \triangleq 1 - \frac{\beta}{2t} \frac{I(kt)}{kt} v(kt, \beta \frac{I(kt)}{kt}) \rightarrow 1$$

Derivation

$$\begin{aligned} \mathbb{P}(U_n > t) &\leq \mathbb{P}(U_n(h_L) > t) + \mathbb{P}(\exists ij \mid h(X_{i_1}, \dots, X_{i_m}) > L) \\ &\leq e^{-\lambda t} \mathbb{E} \left[e^{\lambda U_n(h_L)} \right] + \binom{n}{m} e^{-I(L)} \end{aligned}$$

► issues

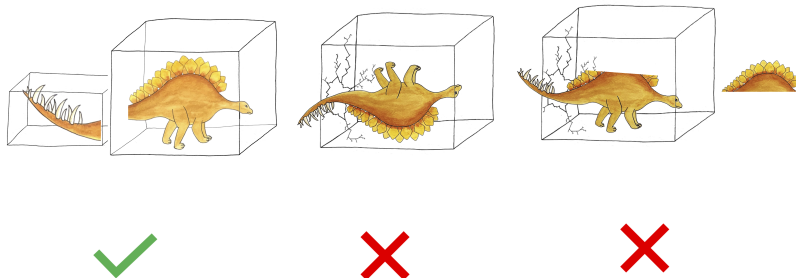
- dependent terms in $U_n(h_L)$
- $\mathbb{E} \left[e^{\lambda U_n(h_L)} \right] \xrightarrow{L \rightarrow \infty} \infty$, fixed λ
- Optimize λ, L together
- Sharpness

$$\text{First term} \leq \exp \left(-\frac{kt^2}{2v(kt, \beta \frac{I(kt)}{kt})} \right) + \exp \left(-\beta I(kt) \max\left(\frac{1}{2}, c(t, \beta, k)\right) \right)$$

- $v(L, \eta) \triangleq \mathbb{E} \left[h_L^2 \mathbf{1}(h \leq 0) + h_L^2 \exp(\eta h_L) \mathbf{1}(h > 0) \right] \rightarrow \text{Var}(h)$
- $c(t, \beta, k) \triangleq 1 - \frac{\beta}{2t} \frac{I(kt)}{kt} v(kt, \beta \frac{I(kt)}{kt}) \rightarrow 1$

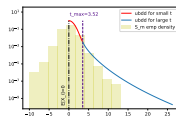
Part 3

- ▶ Tail truncation is optimal (in several cases)



For large same size, the bound is tight

$$\mathbb{P}(U_n > t) \lesssim \exp\left(-\frac{kt^2}{2\text{Var}(h)}\right) + \left(1 + \binom{n}{m}\right) \exp(-I(kt))$$



Theorem (2)

$$\lim_{n \rightarrow \infty} \frac{-\log \mathbb{P}(U_n > t)}{I(kt)} = 1$$

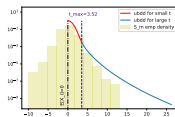
► Assumptions

- $kt^2 \gg I(kt)$
- $I(t) \geq c \sqrt[t]{t} \rightarrow$ sub-Weibull
- $(h > t) \simeq (\exists i |X_i| > f(t))$

$$\text{► } h = |X_1 - X_2|, (X_1 - X_2)^2, \max(|X_1|, |X_2|, \dots, |X_m|)$$

For large same size, the bound is tight

$$\mathbb{P}(U_n > t) \lesssim \exp\left(-\frac{kt^2}{2\text{Var}(h)}\right) + \left(1 + \binom{n}{m}\right) \exp(-I(kt))$$



Theorem (2)

$$\lim_{n \rightarrow \infty} \frac{-\log \mathbb{P}(U_n > t)}{I(kt)} = 1$$

► Assumptions

- $kt^2 \gg I(kt)$
- $I(t) \geq c \sqrt[t]{t} \rightarrow$ sub-Weibull
- $(h > t) \simeq (\exists i |X_i| > f(t))$
 - $h = |X_1 - X_2|, (X_1 - X_2)^2, \max(|X_1|, |X_2|, \dots, |X_m|)$

It yields Large Deviation Principle (LDP)

► Recall

► $\lim_{n \rightarrow \infty} \frac{-\log \mathbb{P}(U_n > t)}{I(kt)} = 1$

► $I(t) \ll t$

► LDP

► $I(kt) = c\sqrt{kt} \implies \lim_{n \rightarrow \infty} \frac{-\log \mathbb{P}(U_n > t)}{\sqrt{n}} = c\sqrt{\frac{t}{m}}$

► U_n satisfies LDP, with speed \sqrt{n} , and rate function $c\sqrt{\frac{t}{m}}$

It yields Large Deviation Principle (LDP)

► Recall

► $\lim_{n \rightarrow \infty} \frac{-\log \mathbb{P}(U_n > t)}{I(kt)} = 1$

► $I(t) \ll t$

► LDP

► $I(kt) = c \sqrt[\alpha]{kt} \implies \lim_{n \rightarrow \infty} \frac{-\log \mathbb{P}(U_n > t)}{\sqrt[\alpha]{n}} = c \sqrt[\alpha]{\frac{t}{m}}$

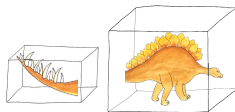
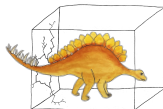
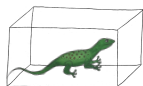
► U_n satisfies LDP, with speed $\sqrt[\alpha]{n}$, and rate function $c \sqrt[\alpha]{\frac{t}{m}}$

Future Work: This was a piece of the puzzle

- ▶ Extend Heavy-tailed Analysis toolbox
 - ▶ Why only U-statistics?
 - ▶ $F_n = F(X_1, \dots, X_n) \xrightarrow{n \rightarrow \infty} \mathbb{E}[F]$
- ▶ Applications
 - ▶ Finance
 - ▶ Differential Privacy
 - ▶ Asymptotic Hypothesis Testing
 - ▶ Bahadur Efficiency
 - ▶ ...

Thanks

Truncation can turn dinos to lizards, so heavy-tails better beware - they're next in line for a makeover!



$$\mathbb{P}(U_n > t) \lesssim \exp\left(-\frac{kt^2}{2\text{Var}(h)}\right) + \left(1 + \binom{n}{m}\right) \exp(-l(kt))$$

$$\mathbb{E}\left[\exp(\lambda h(X_1, \dots, X_m))\right] = \infty,$$

$$\lim_{n \rightarrow \infty} \frac{-\log \mathbb{P}(U_n > t)}{l(kt)} = 1$$

Cartoon characters credit: Mina Latifi

Another application

▶ Phase retrieval

▶ $\mathbf{y} = |\mathbf{X}\boldsymbol{\beta}| + \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$

▶ $\hat{\boldsymbol{\beta}} = \arg \min_{\mathbf{b}} \left\| \mathbf{y}^2 - (\mathbf{X}\mathbf{b})^2 \right\|^2$

The extreme event:

One X_i large $\implies \binom{n-1}{m-1}$ kernel terms large

► Recall:

$$\text{► } U_n = \frac{1}{\binom{n}{m}} \sum h(X_{i_1}, \dots, X_{i_m})$$

$$\text{► } (\exists i |X_i| > f(t)) \simeq (h > t)$$

$$\text{► } \frac{-\log \mathbb{P}(U_n > t)}{I(kt)} \rightarrow 1$$

► $\mathcal{E} = (\exists i \text{ s.t. } |X_i| > f(kt))$ (the event)

$$\text{► } \mathbb{P}(\mathcal{E}) \simeq \exp(-I(kt))$$

$$\text{► } \mathcal{E} \implies (U_n > t), \quad U_n > \frac{1}{\binom{n}{m}} \binom{n-1}{m-1} kt = \frac{m}{n} kt \simeq t$$