

Heavy-Tailed Large Deviations and Sharp Characterization of Global Dynamics of SGDs in Deep Learning

Chang-Han Rhee

Northwestern University

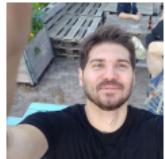
Machine Learning and Data Science Seminar, Oxford

April 8, 2024

Based on the joint works with

Mihail Bazhba, Jose Blanchet, Bohan Chen, Sewoong Oh, Zhe Su, Xingyu Wang, and Bert Zwart

Team



Mihail Bazhba
U. of Amsterdam



Jose Blanchet
Stanford



Bohan Chen
Munich Re



Sewoong Oh
U. of Washington



Xingyu Wang
Northwestern



Zhe Su
Northwestern



Bert Zwart
CWI

Generalization Mystery of Deep Learning

Empirical Success of Deep Neural Networks (DNNs)

“Deep Learning is eating the world.”

- Jorge Nocedal

Empirical Success of Deep Neural Networks (DNNs)

- Even with massive over-parametrization, DNNs generalize well.

Empirical Success of Deep Neural Networks (DNNs)

- Even with massive over-parametrization, DNNs generalize well.

0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9

Training Set



Test Set image source

Empirical Success of Deep Neural Networks (DNNs)

- Even with massive over-parametrization, DNNs generalize well.
- Training a DNN can be formulated as a non-convex loss optimization problem.

0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9

Training Set



Test Set image source

Empirical Success of Deep Neural Networks (DNNs)

- Even with massive over-parametrization, DNNs generalize well.
- Training a DNN can be formulated as a non-convex loss optimization problem.
- Algorithmic Regularazation:

0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9

Training Set



Test Set image source

Empirical Success of Deep Neural Networks (DNNs)

- Even with massive over-parametrization, DNNs generalize well.
- Training a DNN can be formulated as a non-convex loss optimization problem.
- Algorithmic Regularazation: choice of numerical optimization algorithm matters

0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9

Training Set



Test Set image source

Empirical Success of Deep Neural Networks (DNNs)

- Even with massive over-parametrization, DNNs generalize well.
- Training a DNN can be formulated as a non-convex loss optimization problem.
- Algorithmic Regularazation: choice of numerical optimization algorithm matters
- Stochastic Gradient Descent (SGD) turns out to work well.

0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9

Training Set



Test Set image source

Empirical Success of Deep Neural Networks (DNNs)

- Even with massive over-parametrization, DNNs generalize well.
- Training a DNN can be formulated as a non-convex loss optimization problem.
- Algorithmic Regularazation: choice of numerical optimization algorithm matters

Better than deterministic gradient descent.

- Stochastic Gradient Descent (SGD) turns out to work well.

0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9

Training Set



Test Set image source

Empirical Success of Deep Neural Networks (DNNs)

- Even with massive over-parametrization, DNNs generalize well.
- Training a DNN can be formulated as a non-convex loss optimization problem.
- Algorithmic Regularazation: choice of numerical optimization algorithm matters

Better than deterministic gradient descent. In fact, the more noise, the better!

- Stochastic Gradient Descent (SGD) turns out to work well.

0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9

Training Set



Test Set image source

Empirical Success of Deep Neural Networks (DNNs)

- Even with massive over-parametrization, DNNs generalize well.
- Training a DNN can be formulated as a non-convex loss optimization problem.
- Algorithmic Regularazation: choice of numerical optimization algorithm matters
 - Better than deterministic gradient descent. In fact, the more noise, the better!
- Stochastic Gradient Descent (SGD) turns out to work well.

A Central Mystery of Deep Learning

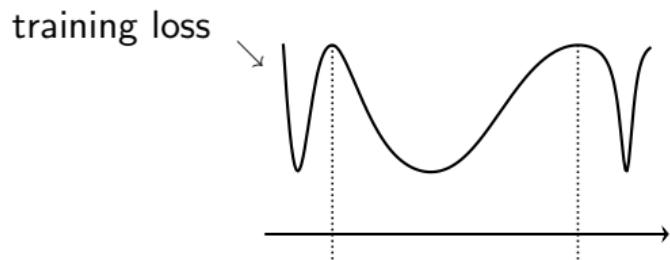
Heavy Tails may have something to do with the Mystery

Heavy Tails may have something to do with the Mystery

- Popular explanation: SGD somehow finds flat local minima.

Heavy Tails may have something to do with the Mystery

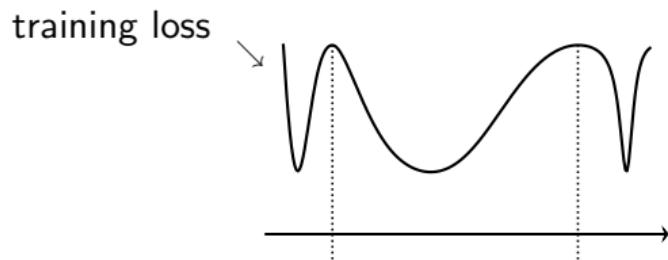
- Popular explanation: SGD somehow finds flat local minima.



Heavy Tails may have something to do with the Mystery

- Popular explanation: SGD somehow finds flat local minima.

↑ tends to generalize well

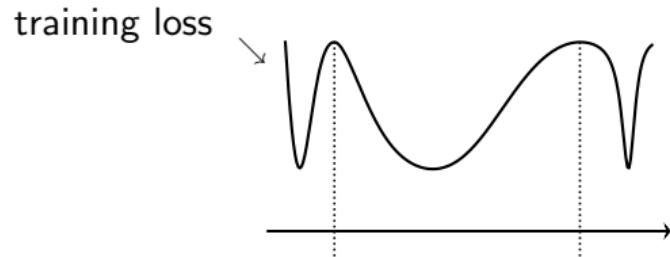


Heavy Tails may have something to do with the Mystery

- Popular explanation: SGD somehow finds flat local minima.

tends to generalize well

- But how?

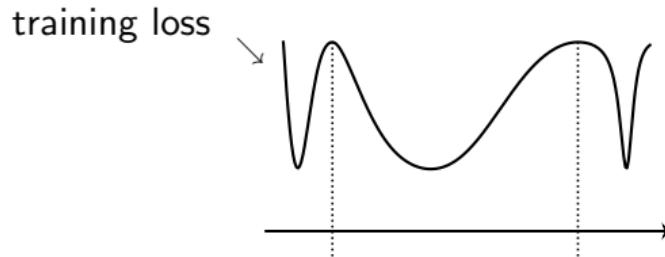


Heavy Tails may have something to do with the Mystery

- Popular explanation: SGD somehow finds flat local minima.

tends to generalize well

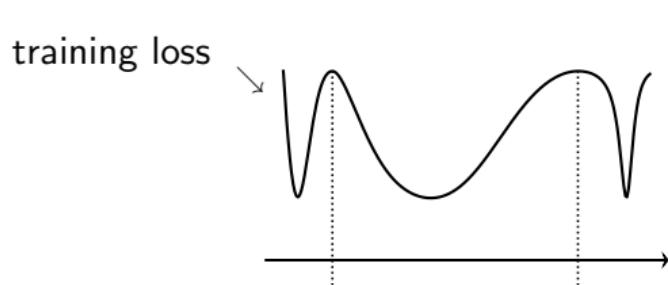
- But how?



It takes more than 10^{600} time steps
for SGD to escape from any of these.

Heavy Tails may have something to do with the Mystery

- Popular explanation: SGD somehow finds flat local minima.
 ↑ tends to generalize well
- But how? Previous attempts to explain had been unsatisfactory.



It takes more than 10^{600} time steps
for SGD to escape from any of these.

Heavy Tails may have something to do with the Mystery

- Popular explanation: SGD somehow finds flat local minima.
  tends to generalize well
- But how? Previous attempts to explain had been unsatisfactory.
- Evidence/arguments for heavy-tails in training DNNs with SGDs!
eg. Simsekli et al. (2019), Hodgkinsons & Mahoney (2020), etc

Heavy Tails may have something to do with the Mystery

- Popular explanation: SGD somehow finds flat local minima.
  tends to generalize well
- But how? Previous attempts to explain had been unsatisfactory.
- Evidence/arguments for heavy-tails in training DNNs with SGDs!
eg. Simsekli et al. (2019), Hodgkinsons & Mahoney (2020), etc

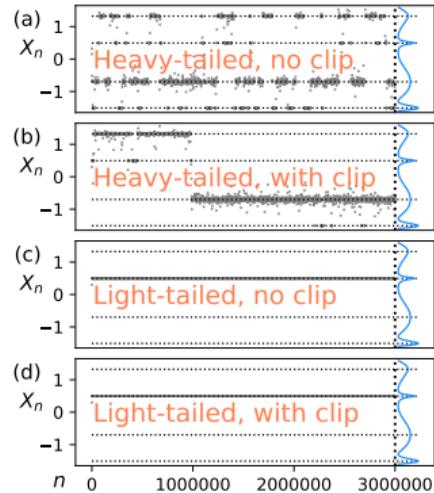
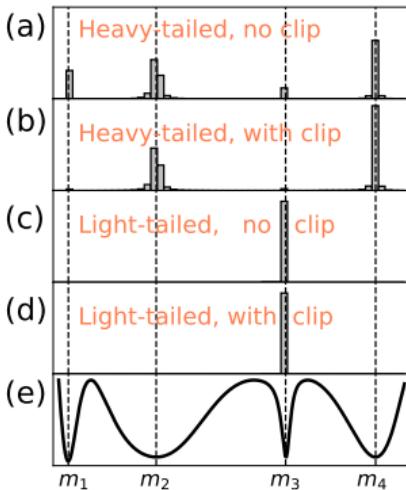
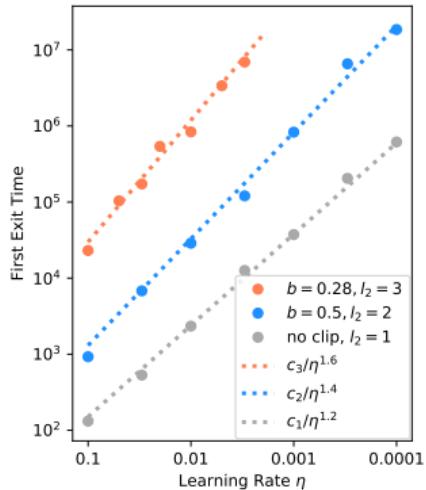
Heavy-tailed SGD escapes local minima and prefers flat local minima.

Heavy Tails may have something to do with the Mystery

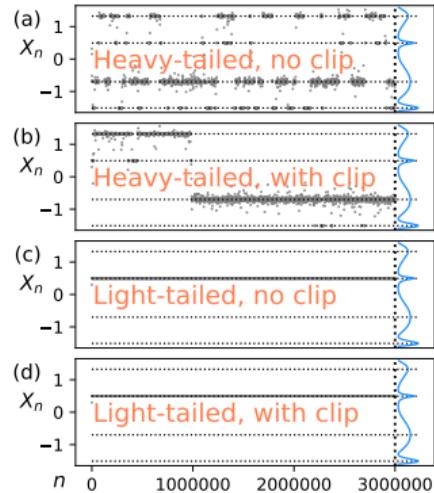
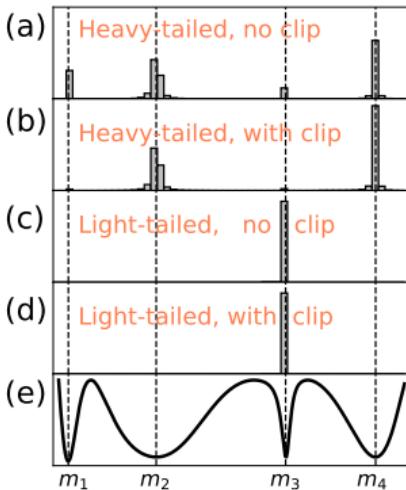
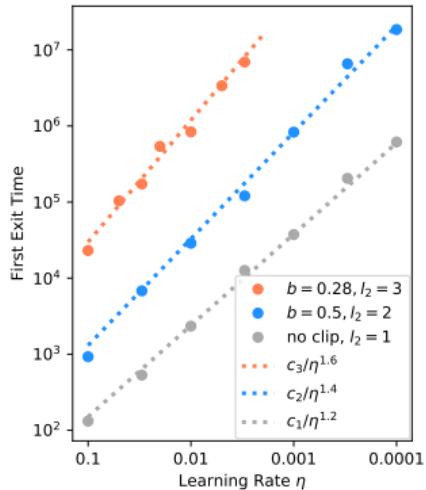
- Popular explanation: SGD somehow finds flat local minima.
  tends to generalize well
- But how? Previous attempts to explain had been unsatisfactory.
- Evidence/arguments for heavy-tails in training DNNs with SGDs!
eg. Simsekli et al. (2019), Hodgkinsons & Mahoney (2020), etc

However, when SGD is heavy-tailed, one often truncates gradients.

Entirely Different Global Dynamics Depending on Tail Behaviors

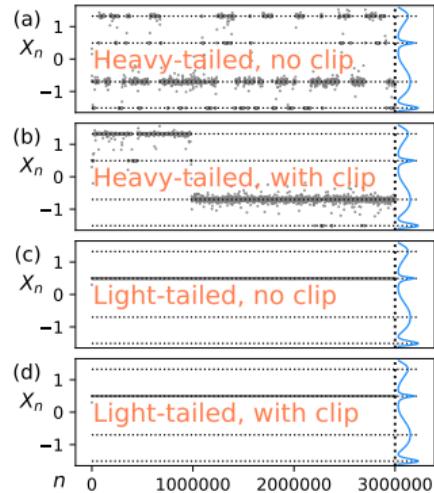
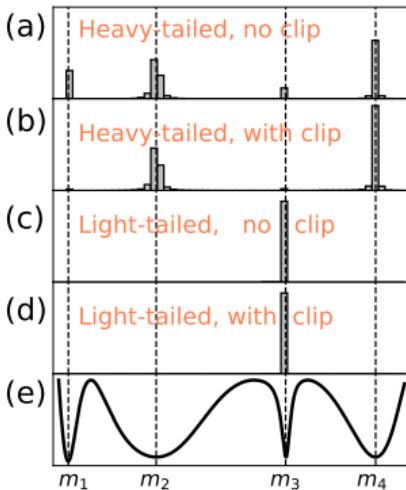
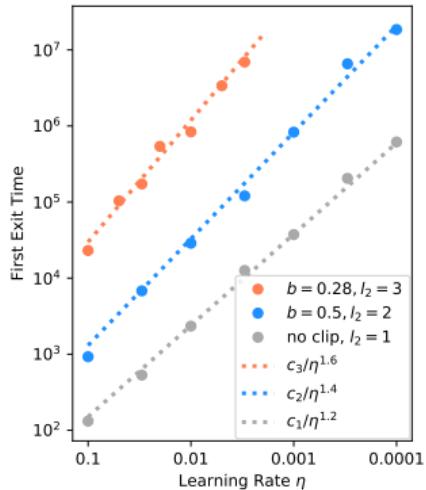


Entirely Different Global Dynamics Depending on Tail Behaviors



Explanation?

Entirely Different Global Dynamics Depending on Tail Behaviors



Explanation? Catastrophe Principle.

Heavy Tails and Catastrophe Principle

Rare Events depend on “Tail Behaviors”

Light-Tailed Distributions

- Extreme Values are Very Rare
- Normal, Exponential, etc



Heavy-Tailed Distributions

- Extreme Values are Frequent
- Power Law, Weibull, etc



Rare Events depend on “Tail Behaviors”

Light-Tailed Distributions

- Extreme Values are Very Rare
- Normal, Exponential, etc



Heavy-Tailed Distributions

- Extreme Values are Frequent
- Power Law, Weibull, etc



Structural difference in the way systemwide rare events arise.

Rare Events depend on “Tail Behaviors”

Light-Tailed Distributions

- Extreme Values are Very Rare
- Normal, Exponential, etc

Systemwide rare events

arise because

EVERYTHING goes wrong.

(Conspiracy Principle)

Heavy-Tailed Distributions

- Extreme Values are Frequent
- Power Law, Weibull, etc



Structural difference in the way systemwide rare events arise.

Rare Events depend on “Tail Behaviors”

Light-Tailed Distributions

- Extreme Values are Very Rare
- Normal, Exponential, etc

Heavy-Tailed Distributions

- Extreme Values are Frequent
- Power Law, Weibull, etc

Systemwide rare events

arise because

EVERYTHING goes wrong.

(Conspiracy Principle)

Systemwide rare events

arise because of

A FEW Catastrophes.

(Catastrophe Principle)

Structural difference in the way systemwide rare events arise.

Insurance Example: Capital Reserve

$$Y(t) = c + pt - \sum_{i=1}^{N(t)} X_i$$

Insurance Example: Capital Reserve

Initial Capital

$$Y(t) = c + pt - \sum_{i=1}^{N(t)} X_i$$

Insurance Example: Capital Reserve

Initial Capital

$$Y(t) = c + pt - \sum_{i=1}^{N(t)} X_i$$

Premium

Insurance Example: Capital Reserve

Initial Capital

$$Y(t) = c + pt - \sum_{i=1}^{N(t)} X_i$$

Premium i.i.d. Claim Size

Insurance Example: Capital Reserve

$$Y(t) = c + pt - \sum_{i=1}^{N(t)} X_i$$

Initial Capital Poisson Arrival
Premium i.i.d. Claim Size

Insurance Example: Capital Reserve

$$Y(t) = c + pt - \sum_{i=1}^{N(t)} X_i$$

Insurance Example: Capital Reserve

$$\bar{Y}_{\textcolor{red}{n}}(t) = c + pt - \sum_{i=1}^{N(\textcolor{red}{nt})} X_i / \textcolor{red}{n}$$

Insurance Example: Capital Reserve

$$\bar{Y}_{\textcolor{red}{n}}(t) = c + pt - \sum_{i=1}^{N(\textcolor{red}{n}t)} X_i / n$$

Large n : analysis of large loss over a long time period

Typical Scenario

Typical Scenario

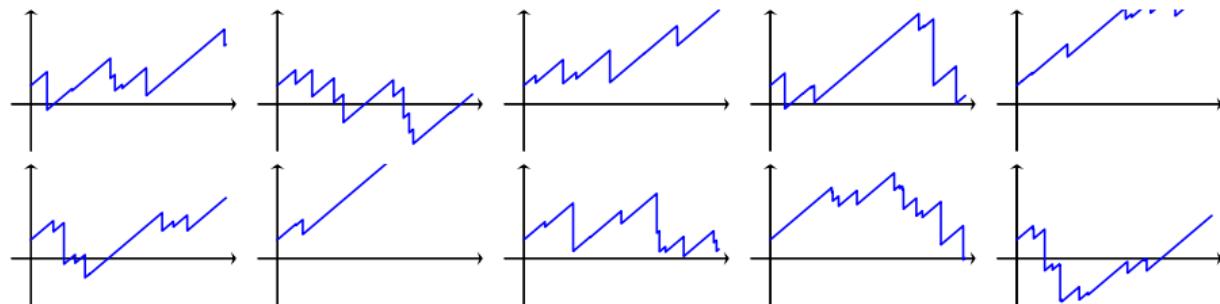
Sample paths of \bar{Y}_n :

$n=10$ & claim sizes are **light-tailed**

Typical Scenario

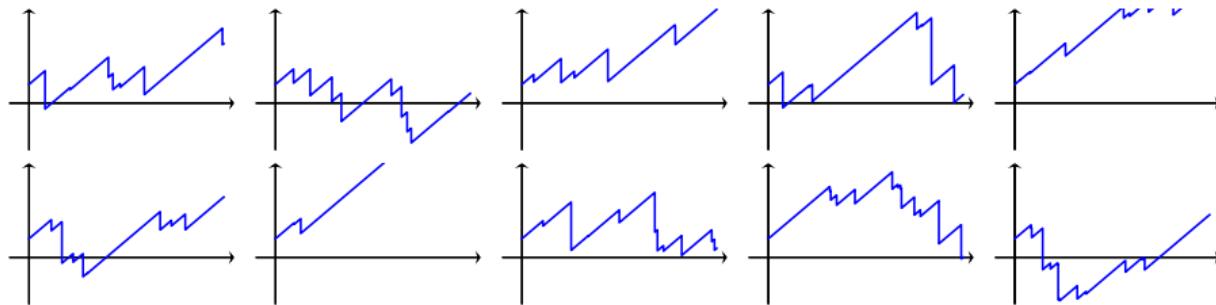
Sample paths of \bar{Y}_n :

$n=10$ & claim sizes are **light-tailed**



Typical Scenario

Sample paths of \bar{Y}_n :



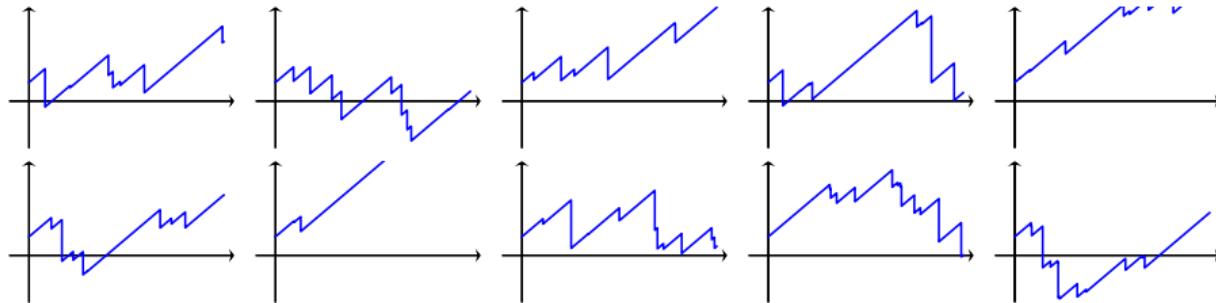
$n=10$ & claim sizes are **light-tailed**

Sample paths of \bar{Y}_n :

$n=10$ & claim sizes are **heavy-tailed**

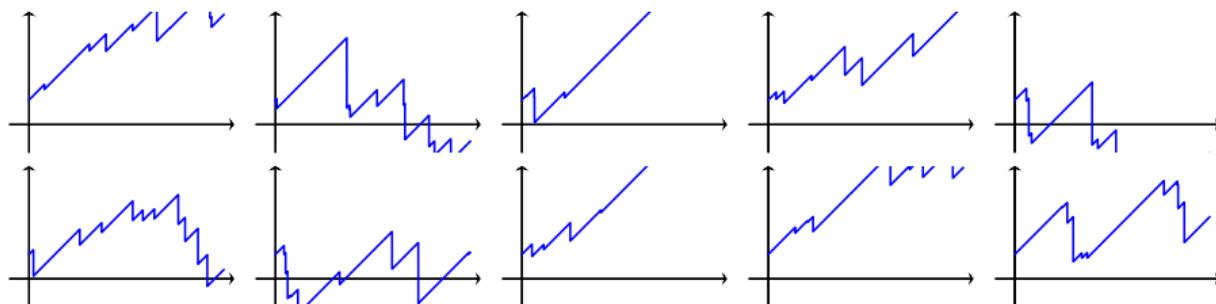
Typical Scenario

Sample paths of \bar{Y}_n :



$n=10$ & claim sizes are **light-tailed**

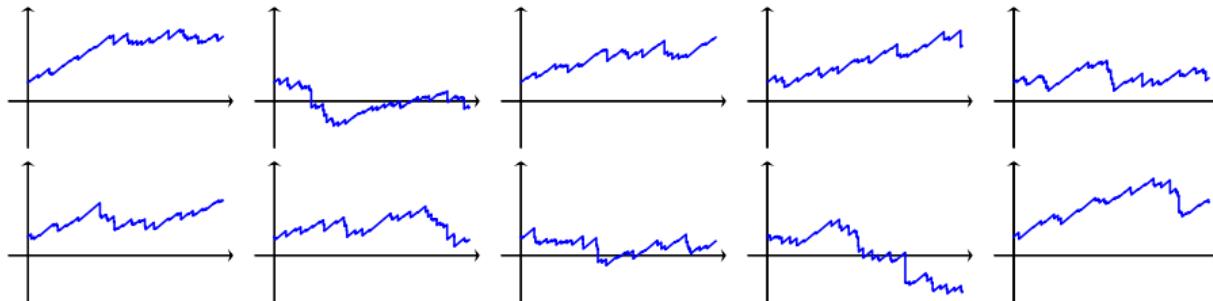
Sample paths of \bar{Y}_n :



$n=10$ & claim sizes are **heavy-tailed**

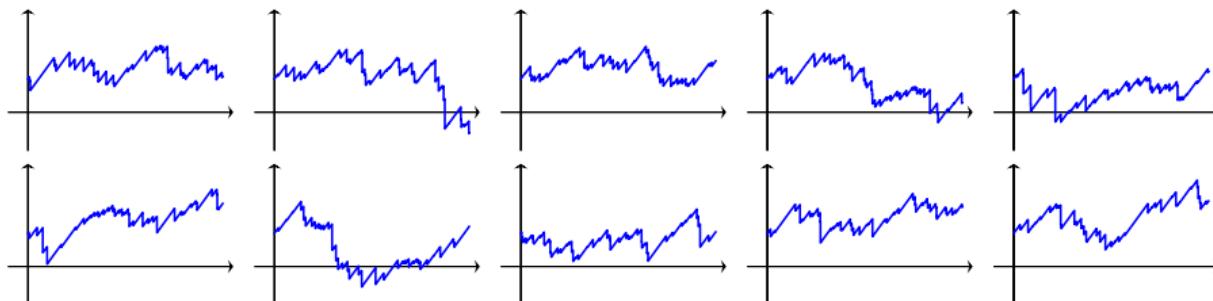
Typical Scenario

Sample paths of \bar{Y}_n :



$n=50$ & claim sizes are **light-tailed**

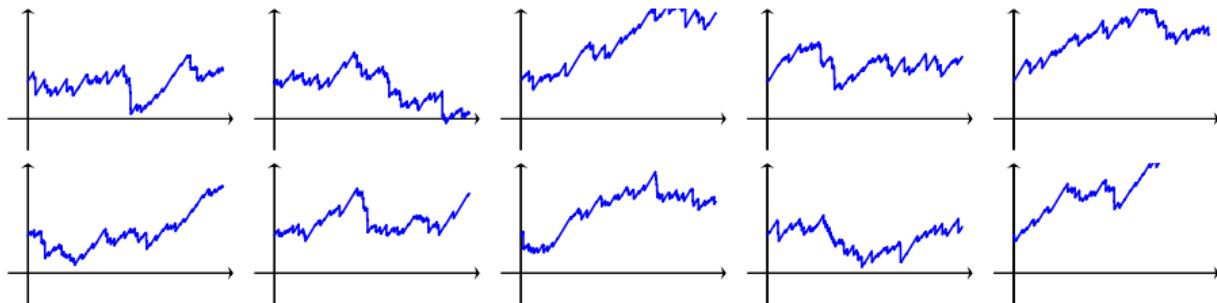
Sample paths of \bar{Y}_n :



$n=50$ & claim sizes are **heavy-tailed**

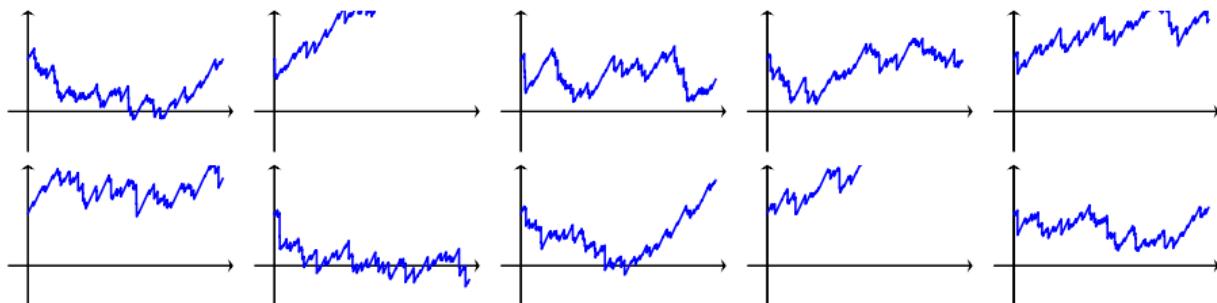
Typical Scenario

Sample paths of \bar{Y}_n :



$n=100$ & claim sizes are **light-tailed**

Sample paths of \bar{Y}_n :

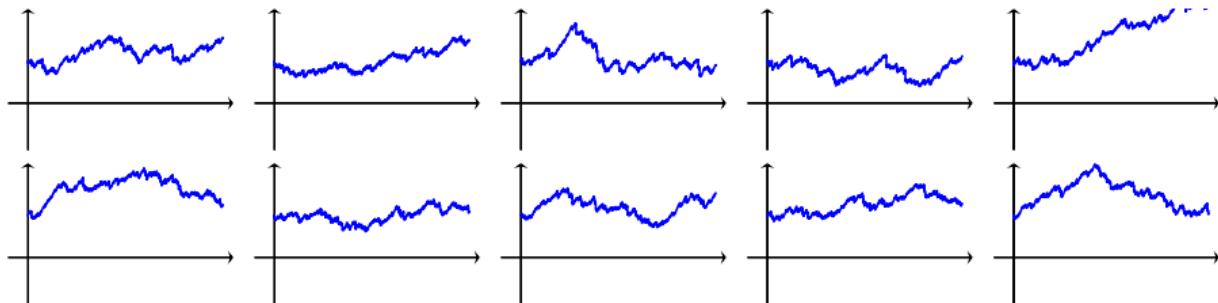


$n=100$ & claim sizes are **heavy-tailed**

Typical Scenario

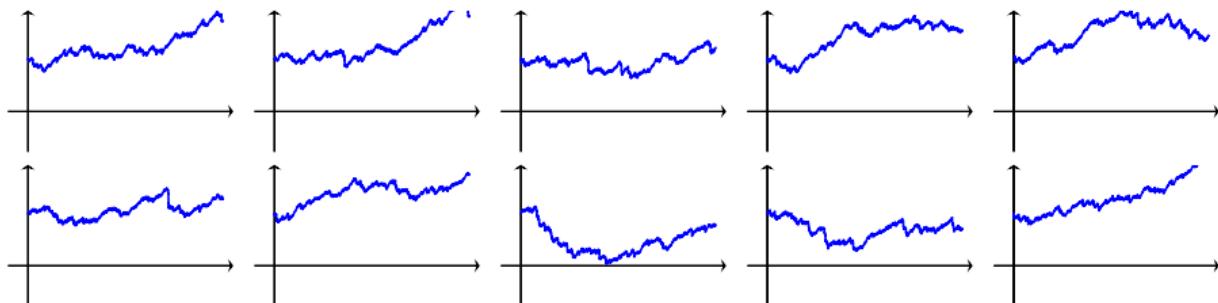
Sample paths of \bar{Y}_n :

$n=500$ & claim sizes are **light-tailed**



Sample paths of \bar{Y}_n :

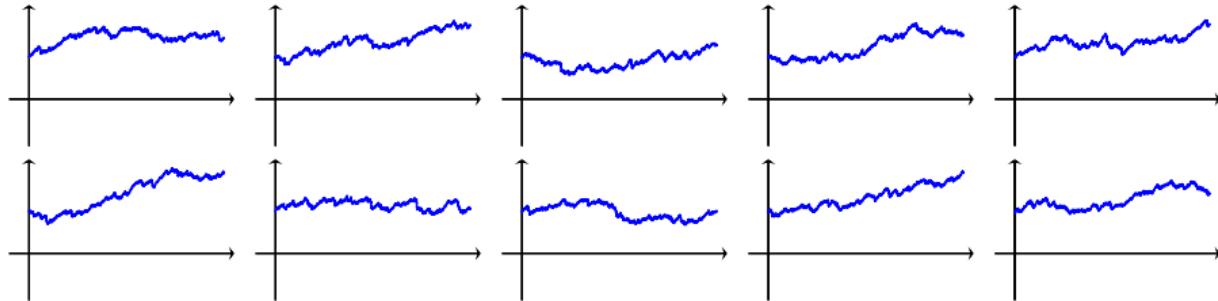
$n=500$ & claim sizes are **heavy-tailed**



Typical Scenario

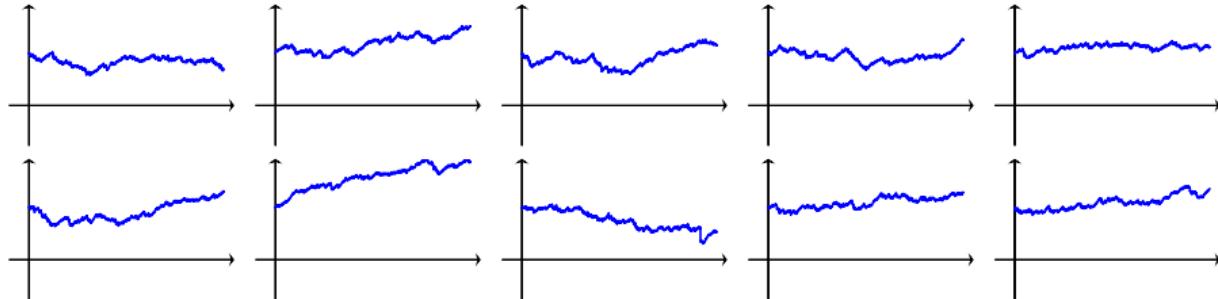
Sample paths of \bar{Y}_n :

$n=1000$ & claim sizes are **light-tailed**



Sample paths of \bar{Y}_n :

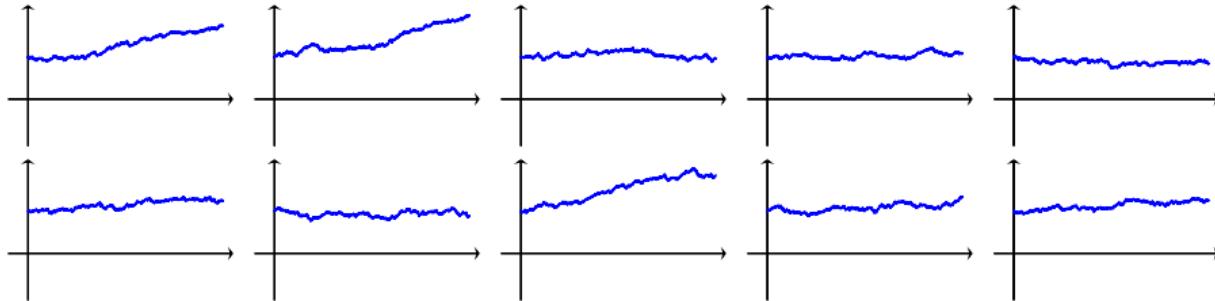
$n=1000$ & claim sizes are **heavy-tailed**



Typical Scenario

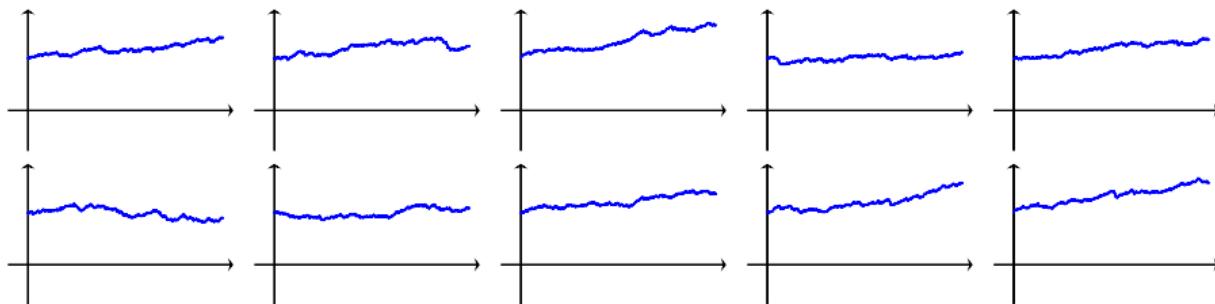
Sample paths of \bar{Y}_n :

$n=2500$ & claim sizes are **light-tailed**



Sample paths of \bar{Y}_n :

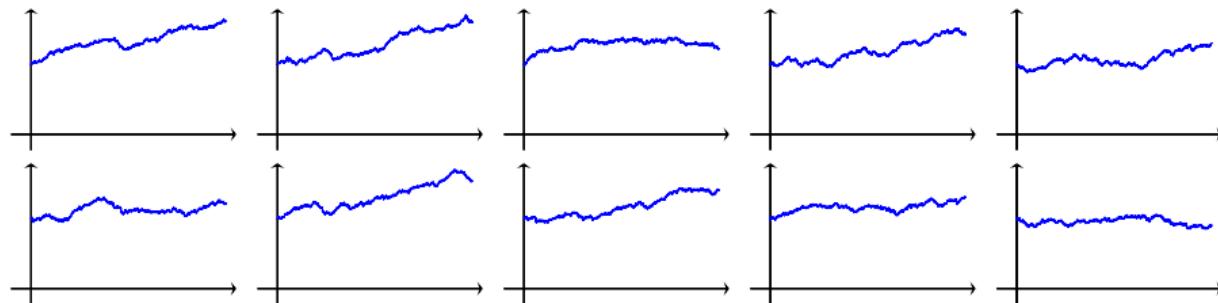
$n=2500$ & claim sizes are **heavy-tailed**



Typical Scenario

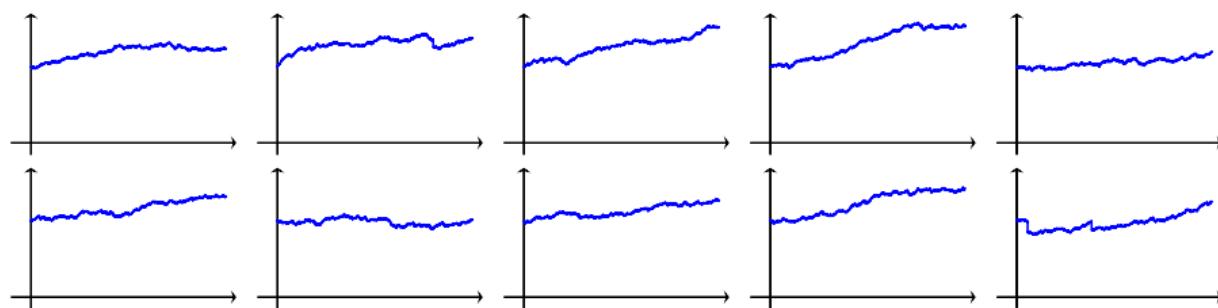
Sample paths of \bar{Y}_n :

$n=5000$ & claim sizes are **light-tailed**



Sample paths of \bar{Y}_n :

$n=5000$ & claim sizes are **heavy-tailed**

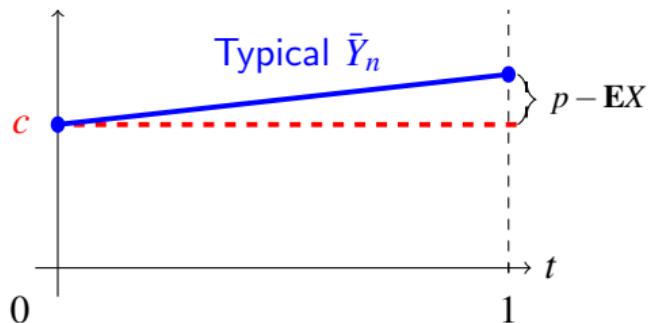


Typical Scenario

That is, $\bar{Y}_n(t) \approx c + (p - \mathbf{E}X)t$ for large n .

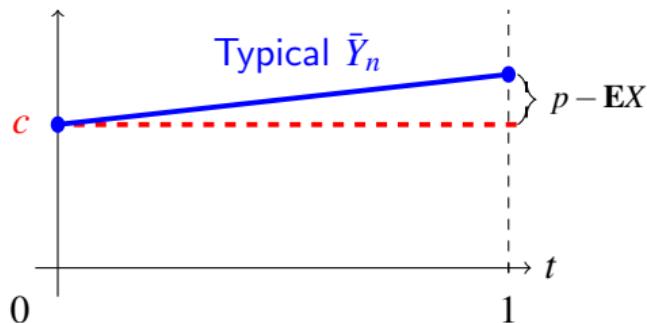
Typical Scenario

That is, $\bar{Y}_n(t) \approx c + (p - \mathbf{E}X)t$ for large n .



Typical Scenario

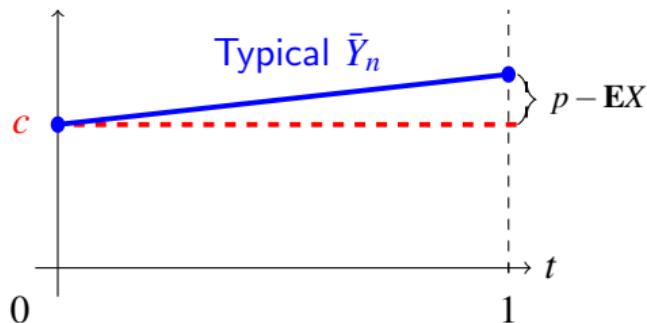
That is, $\bar{Y}_n(t) \approx c + (p - \mathbf{E}X)t$ for large n .



Typically, your business will flourish

Typical Scenario

That is, $\bar{Y}_n(t) \approx c + (p - \mathbf{E}X)t$ for large n .

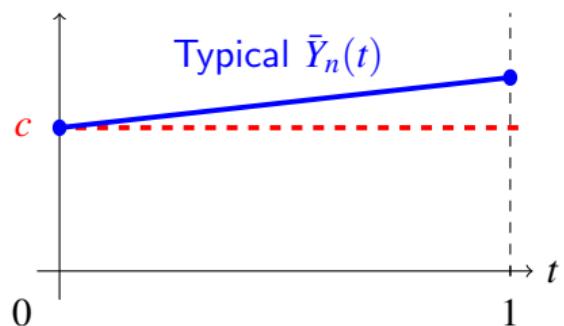


Typically, your business will flourish

regardless of the tail distributions.

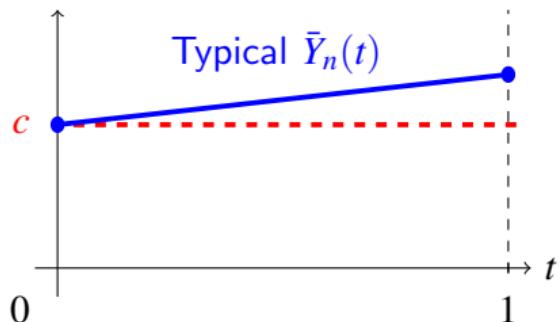
What about atypical cases?

A Rare Event: Bankruptcy



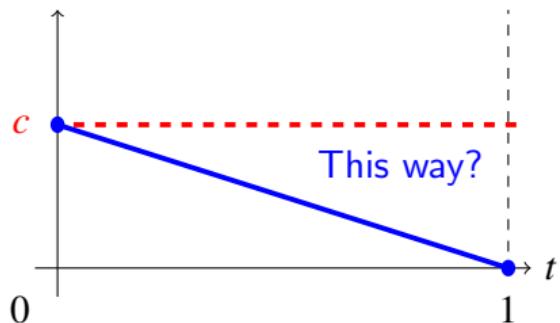
A Rare Event: Bankruptcy

Consider $B \triangleq \{ \bar{Y}_n \text{ falls below } 0 \text{ on } [0, 1] \}$. (i.e., Bankruptcy)



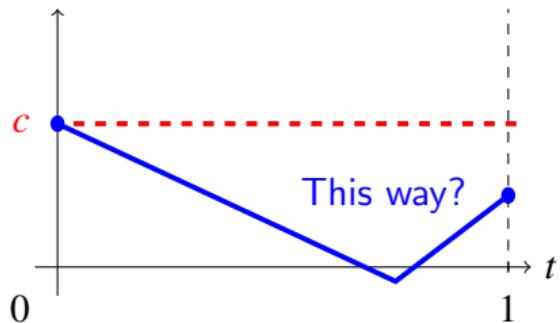
A Rare Event: Bankruptcy

Consider $B \triangleq \{ \bar{Y}_n \text{ falls below } 0 \text{ on } [0, 1] \}$. (i.e., Bankruptcy)



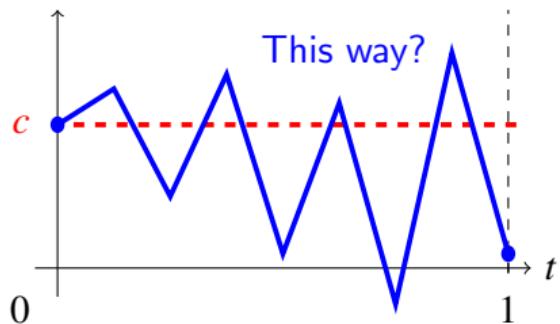
A Rare Event: Bankruptcy

Consider $B \triangleq \{ \bar{Y}_n \text{ falls below } 0 \text{ on } [0, 1] \}$. (i.e., Bankruptcy)



A Rare Event: Bankruptcy

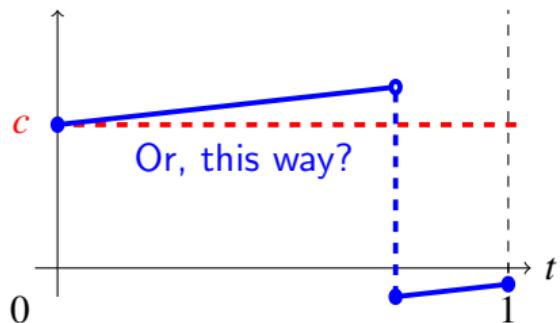
Consider $B \triangleq \{ \bar{Y}_n \text{ falls below } 0 \text{ on } [0, 1] \}$. (i.e., Bankruptcy)



A Rare Event: Bankruptcy

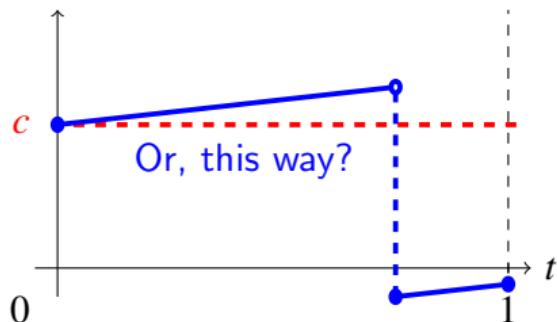
Consider $B \triangleq \{ \bar{Y}_n \text{ falls below } 0 \text{ on } [0, 1] \}$.

(i.e., Bankruptcy)



A Rare Event: Bankruptcy

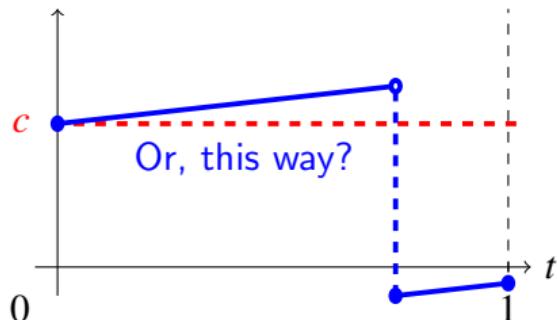
Consider $B \triangleq \{ \bar{Y}_n \text{ falls below } 0 \text{ on } [0, 1] \}$. (i.e., Bankruptcy)



Are we going to see clear patterns?

A Rare Event: Bankruptcy

Consider $B \triangleq \{ \bar{Y}_n \text{ falls below } 0 \text{ on } [0, 1] \}$. (i.e., Bankruptcy)

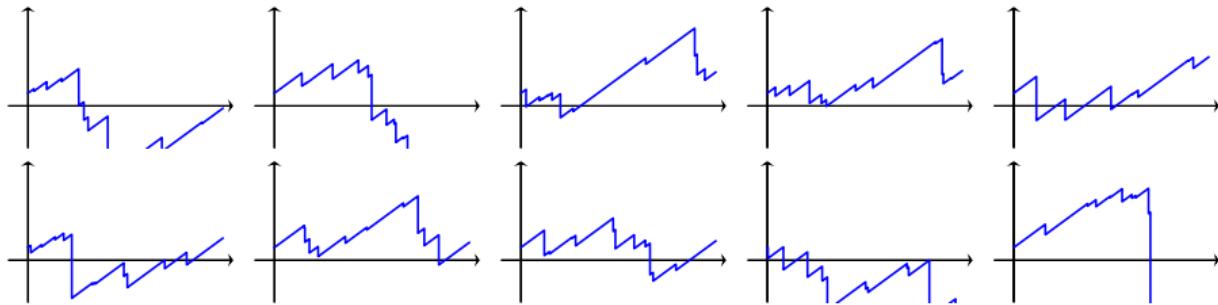


Are we going to see clear patterns?

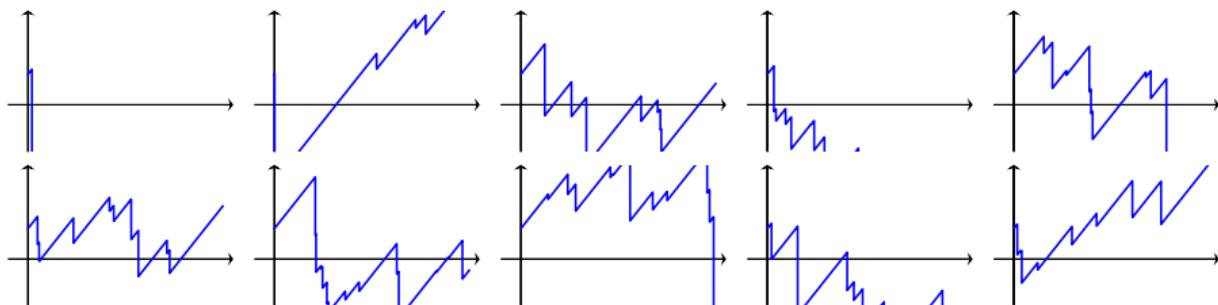
Do they depend on the tail distributions?

A Rare Event: Bankruptcy

Sample paths of \bar{Y}_{10} conditional on B for **light-tailed** claims:

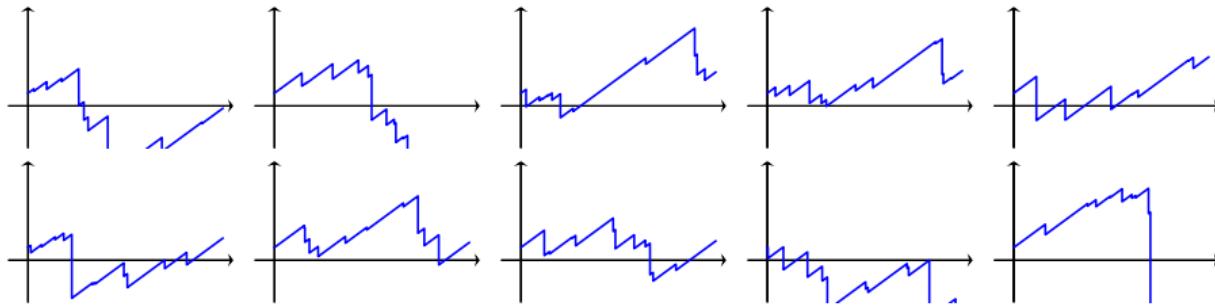


Sample paths of \bar{Y}_{10} conditional on B for **heavy-tailed** claims:

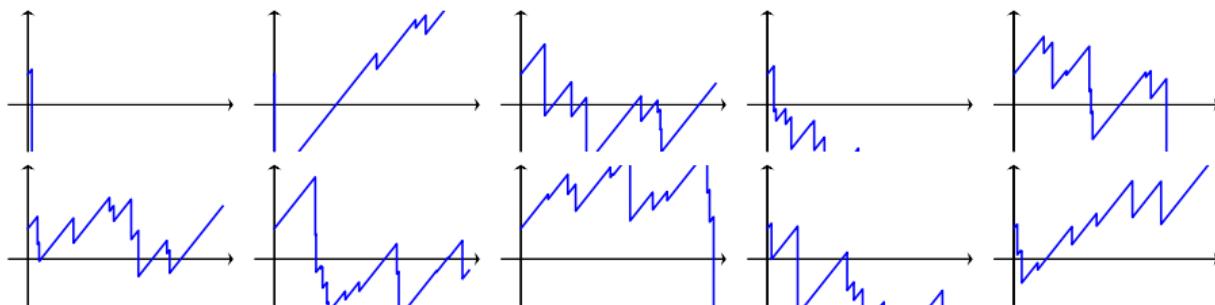


A Rare Event: Bankruptcy

Sample paths of \bar{Y}_{10} conditional on B for **light-tailed** claims:

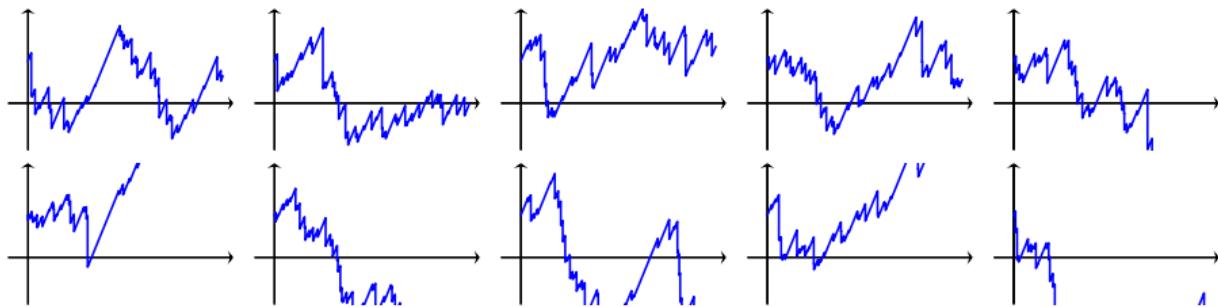


Sample paths of \bar{Y}_{10} conditional on B for **heavy-tailed** claims:

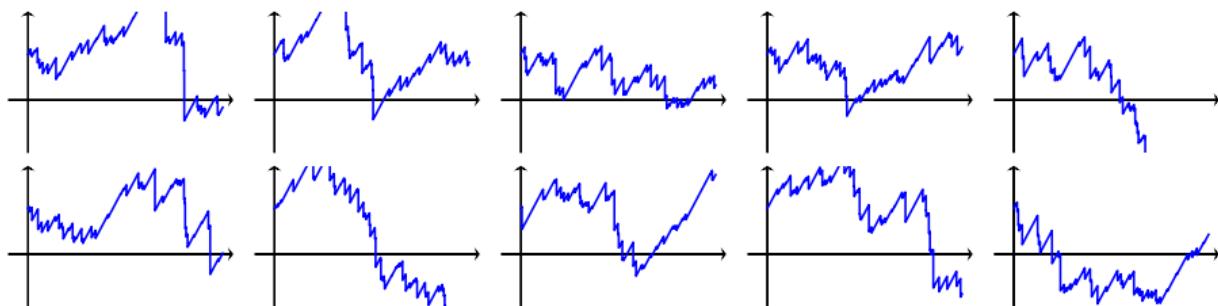


A Rare Event: Bankruptcy

Sample paths of \bar{Y}_{50} conditional on B for **light-tailed** claims:

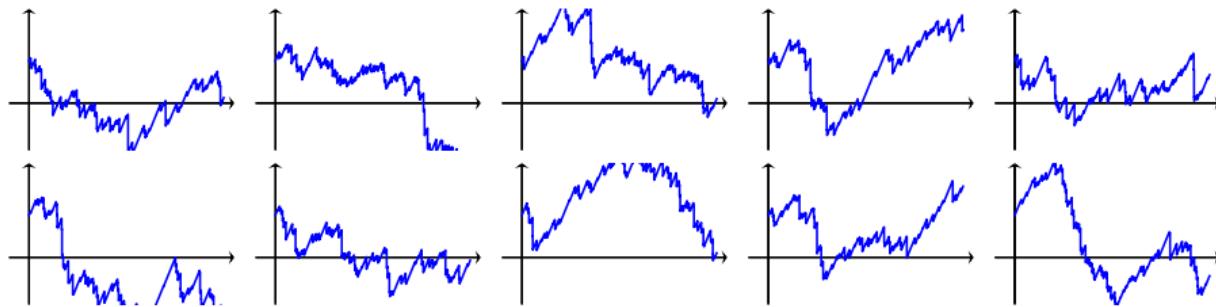


Sample paths of \bar{Y}_{50} conditional on B for **heavy-tailed** claims:

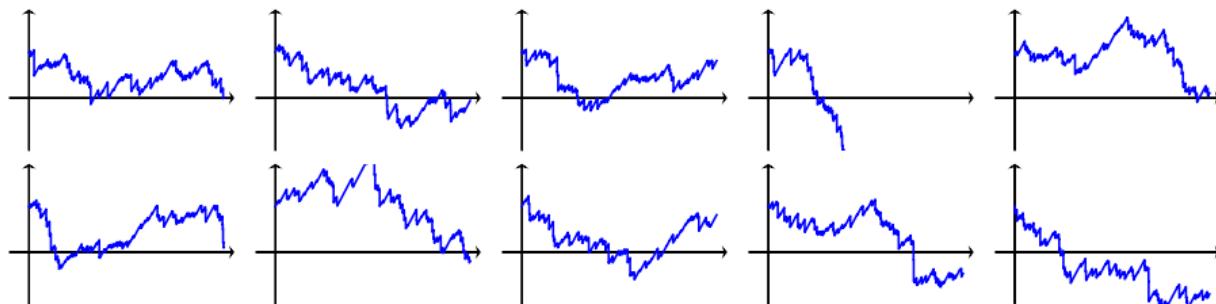


A Rare Event: Bankruptcy

Sample paths of \bar{Y}_{100} conditional on B for **light-tailed** claims:

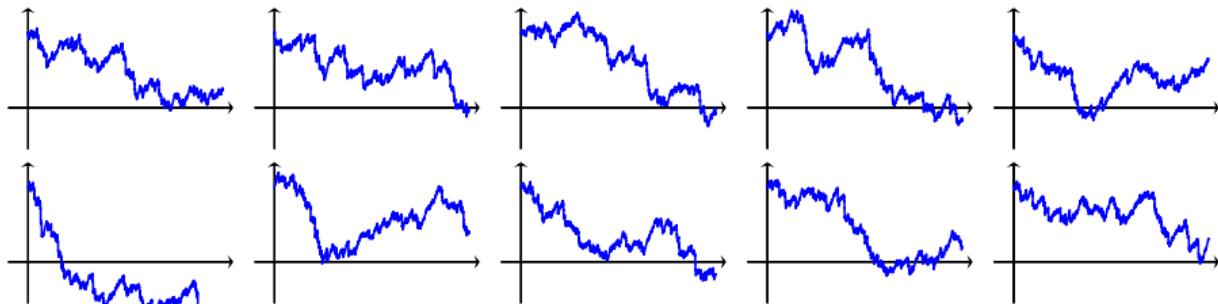


Sample paths of \bar{Y}_{100} conditional on B for **heavy-tailed** claims:

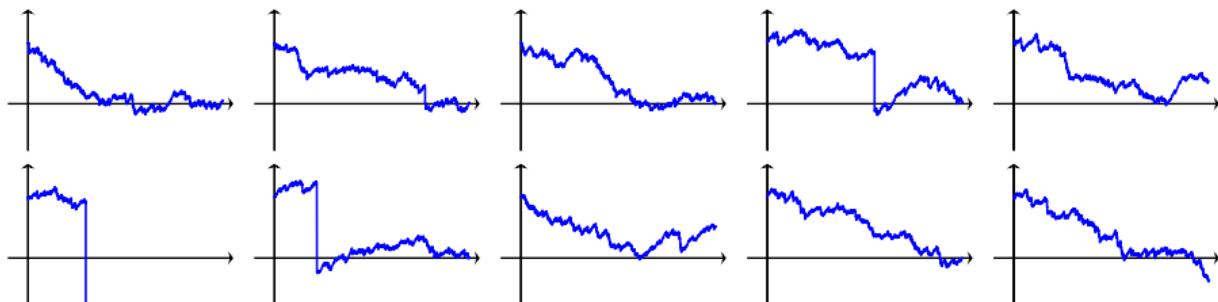


A Rare Event: Bankruptcy

Sample paths of \bar{Y}_{500} conditional on B for **light-tailed** claims:

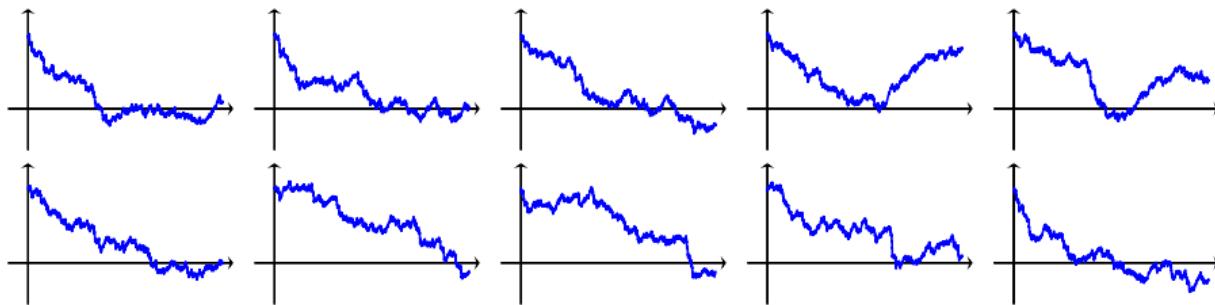


Sample paths of \bar{Y}_{500} conditional on B for **heavy-tailed** claims:

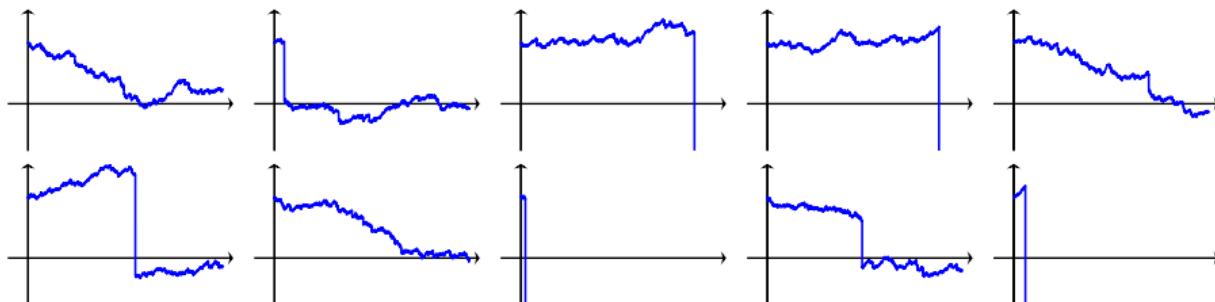


A Rare Event: Bankruptcy

Sample paths of \bar{Y}_{1000} conditional on B for **light-tailed** claims:

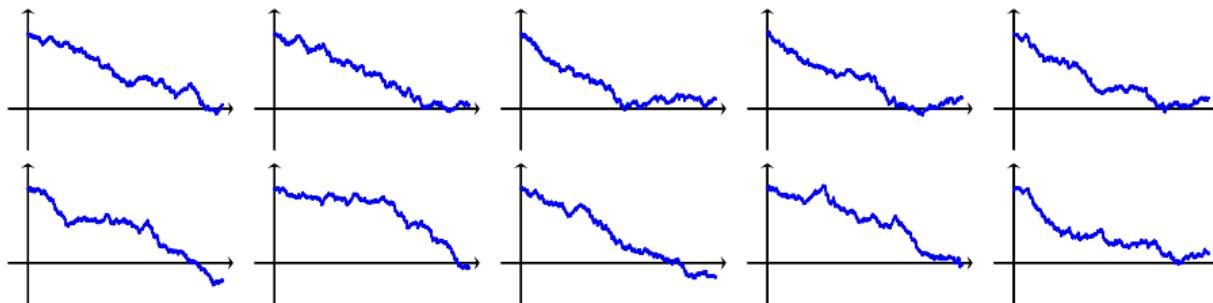


Sample paths of \bar{Y}_{1000} conditional on B for **heavy-tailed** claims:

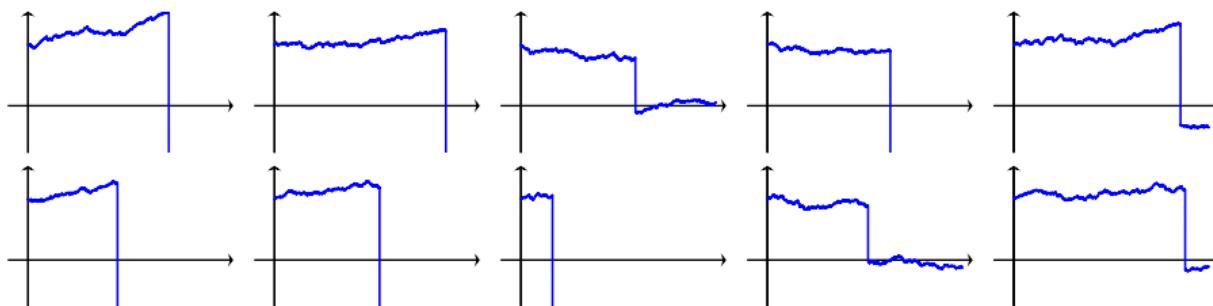


A Rare Event: Bankruptcy

Sample paths of \bar{Y}_{2500} conditional on B for **light-tailed** claims:

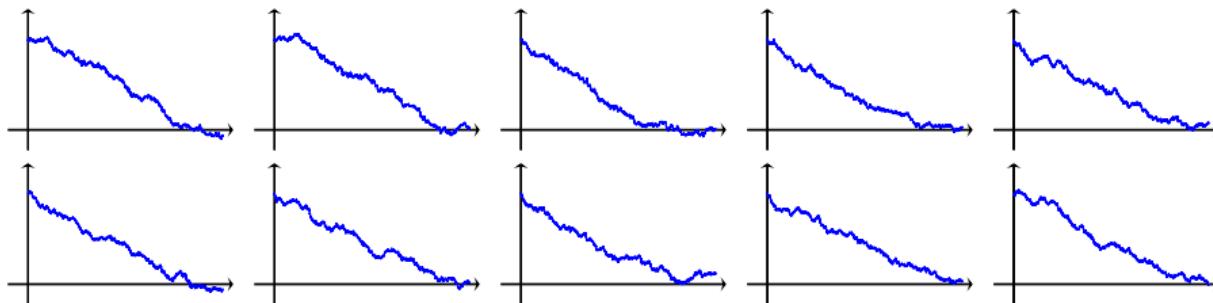


Sample paths of \bar{Y}_{2500} conditional on B for **heavy-tailed** claims:

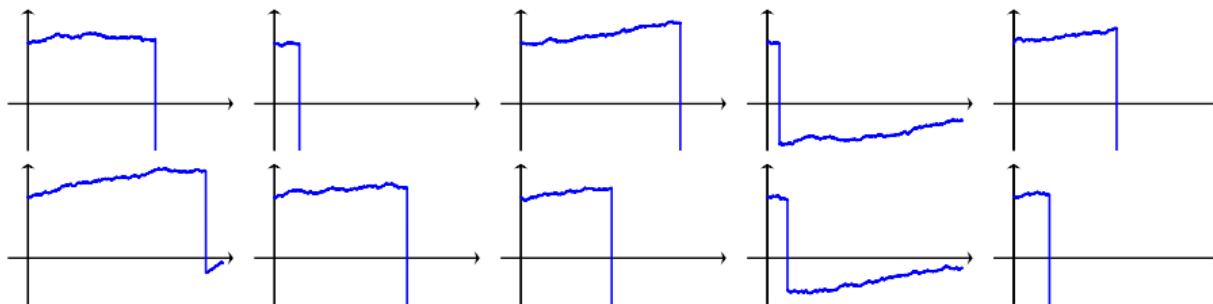


A Rare Event: Bankruptcy

Sample paths of \bar{Y}_{5000} conditional on B for **light-tailed** claims:

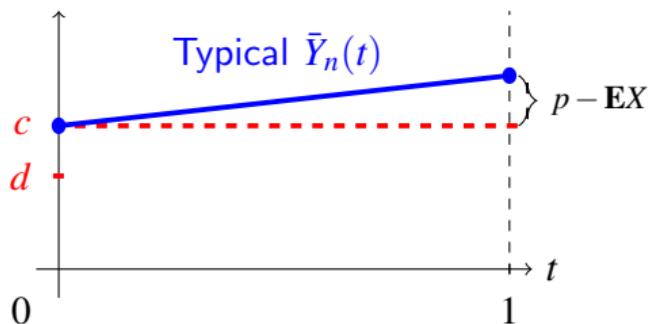


Sample paths of \bar{Y}_{5000} conditional on B for **heavy-tailed** claims:



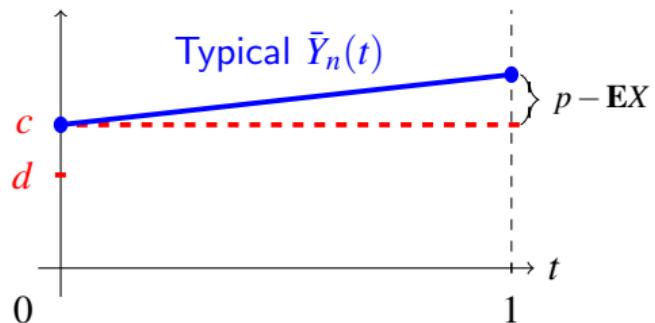
Bankruptcy

Bankruptcy Despite Reinsurance



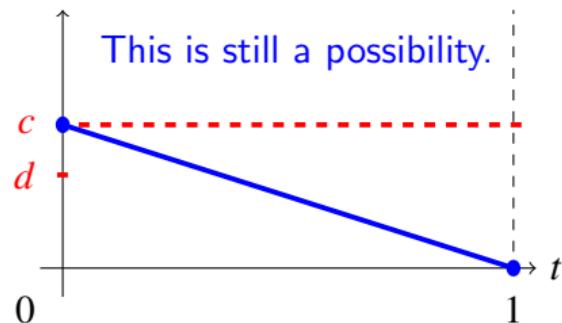
Bankruptcy Despite of Reinsurance

Consider $R \triangleq \{ \bar{Y}_n \text{ falls below } 0 \text{ on } [0,1], \text{ jump sizes } \leq d \}$



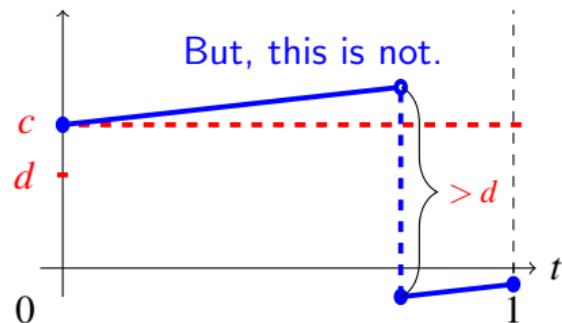
Bankruptcy Despite of Reinsurance

Consider $R \triangleq \{ \bar{Y}_n \text{ falls below } 0 \text{ on } [0,1], \text{ jump sizes } \leq d \}$



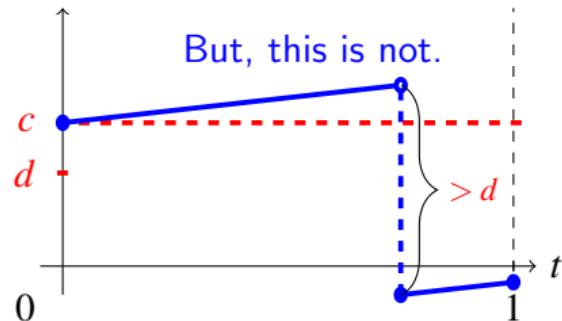
Bankruptcy Despite of Reinsurance

Consider $R \triangleq \{ \bar{Y}_n \text{ falls below } 0 \text{ on } [0,1], \text{ jump sizes } \leq d \}$



Bankruptcy Despite of Reinsurance

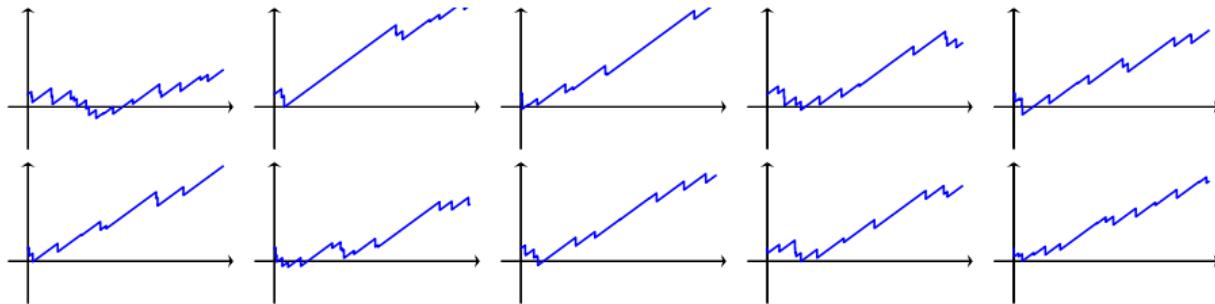
Consider $R \triangleq \{ \bar{Y}_n \text{ falls below } 0 \text{ on } [0,1], \text{ jump sizes } \leq d \}$



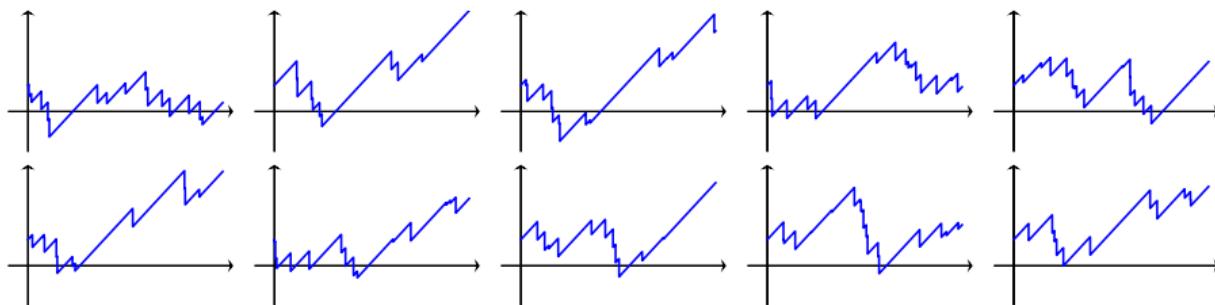
How does the pattern change in this case?

Bankruptcy Despite Reinsurance

Sample paths of \bar{Y}_{10} conditional on R for **light-tailed** claims:

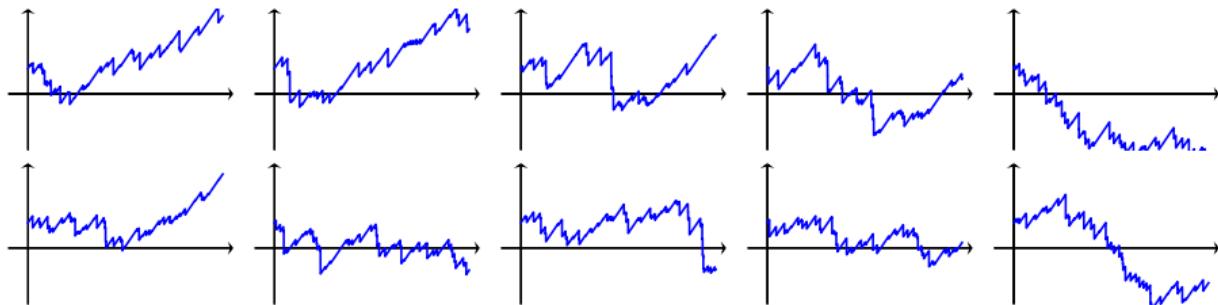


Sample paths of \bar{Y}_{10} conditional on R for **heavy-tailed** claims:

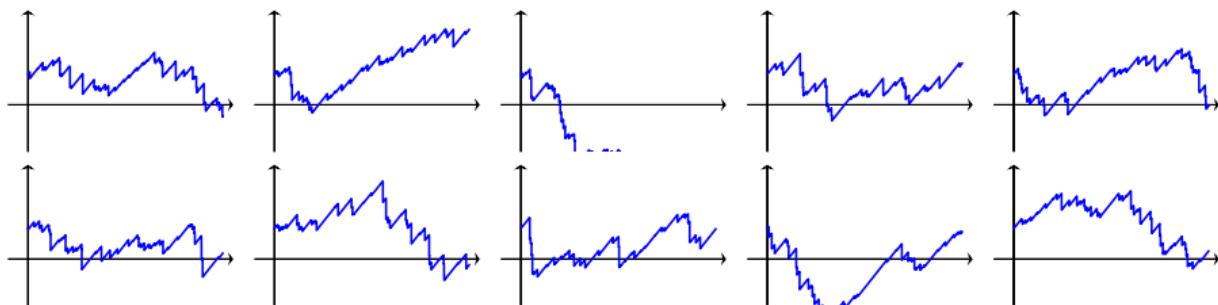


Bankruptcy Despite Reinsurance

Sample paths of \bar{Y}_{50} conditional on R for **light-tailed** claims:

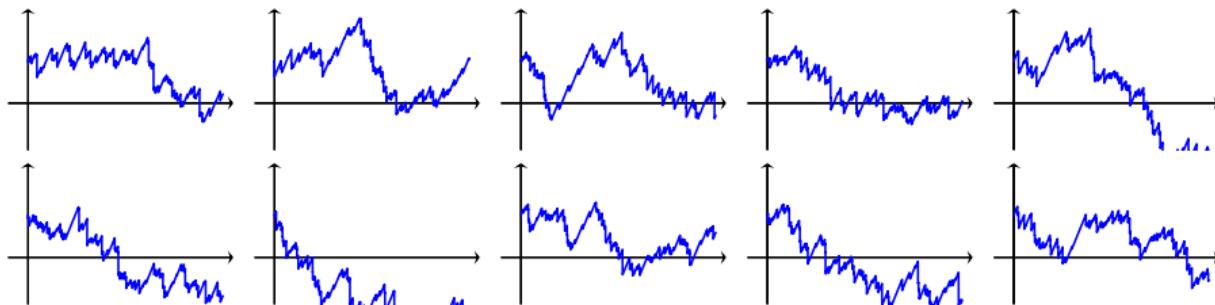


Sample paths of \bar{Y}_{50} conditional on R for **heavy-tailed** claims:

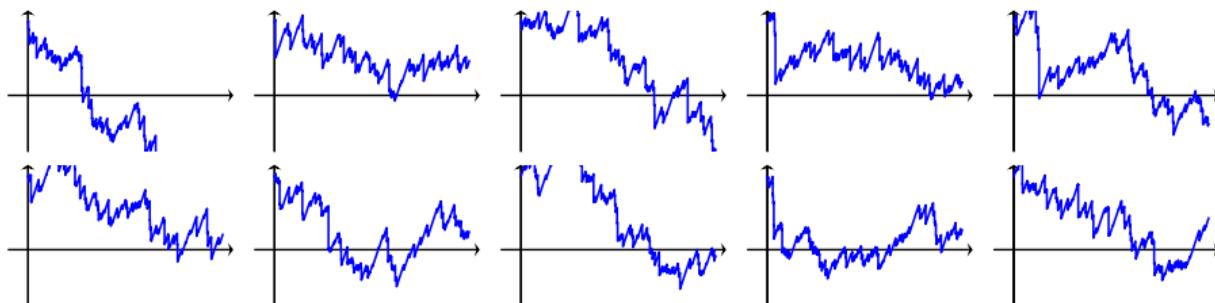


Bankruptcy Despite Reinsurance

Sample paths of \bar{Y}_{100} conditional on R for **light-tailed** claims:

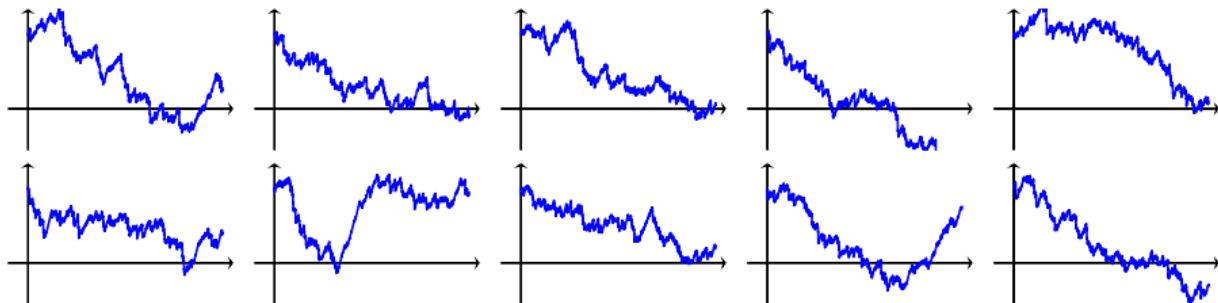


Sample paths of \bar{Y}_{100} conditional on R for **heavy-tailed** claims:

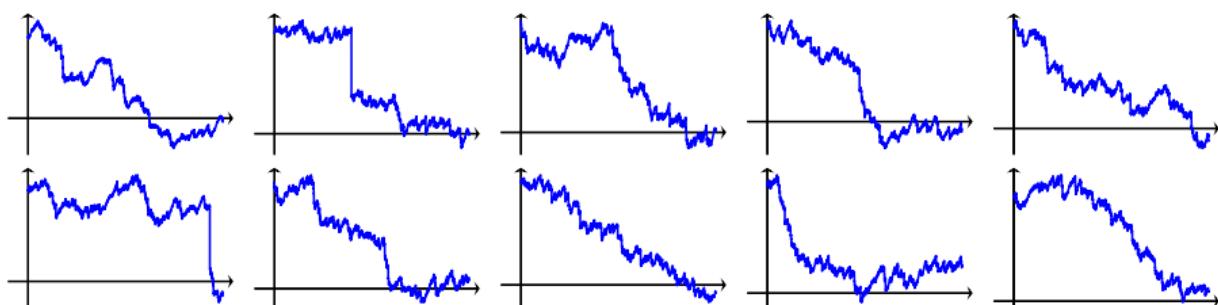


Bankruptcy Despite Reinsurance

Sample paths of \bar{Y}_{500} conditional on R for **light-tailed** claims:

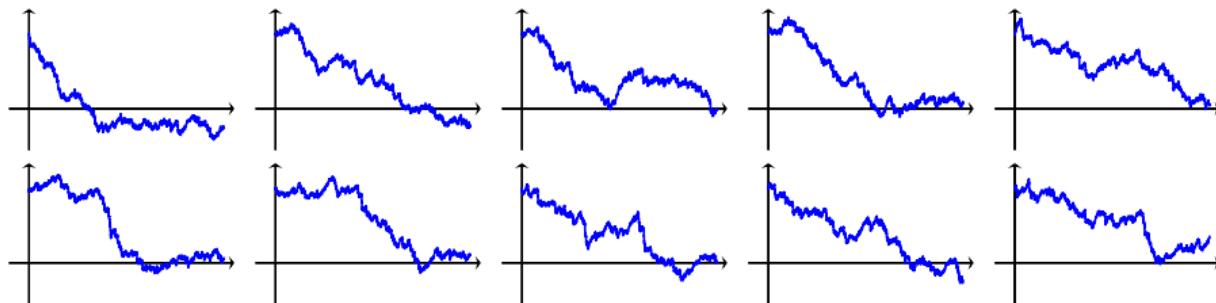


Sample paths of \bar{Y}_{500} conditional on R for **heavy-tailed** claims:

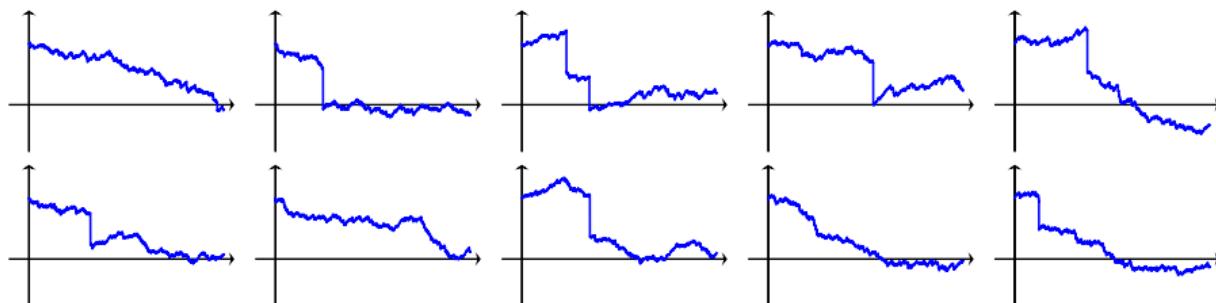


Bankruptcy Despite Reinsurance

Sample paths of \bar{Y}_{1000} conditional on R for **light-tailed** claims:

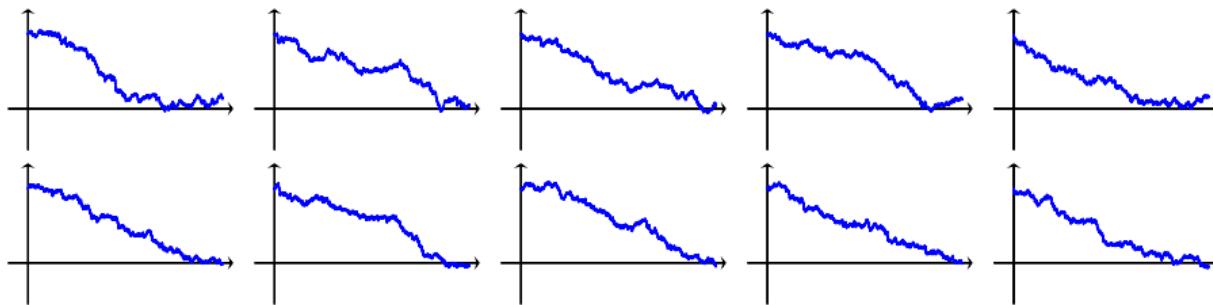


Sample paths of \bar{Y}_{1000} conditional on R for **heavy-tailed** claims:

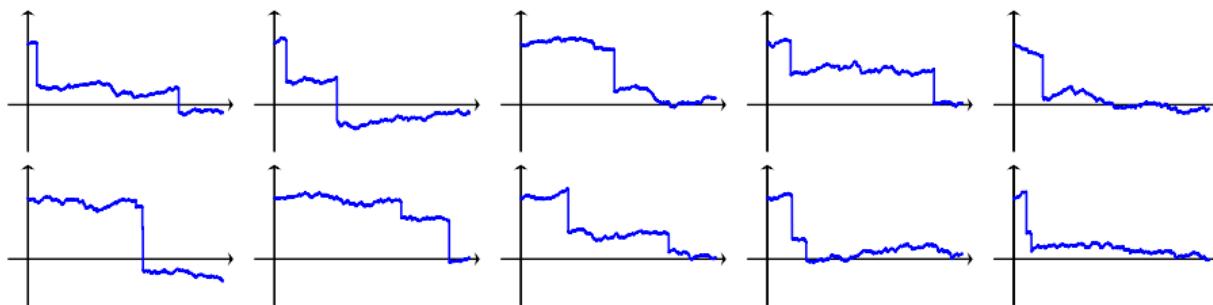


Bankruptcy Despite Reinsurance

Sample paths of \bar{Y}_{2500} conditional on R for **light-tailed** claims:

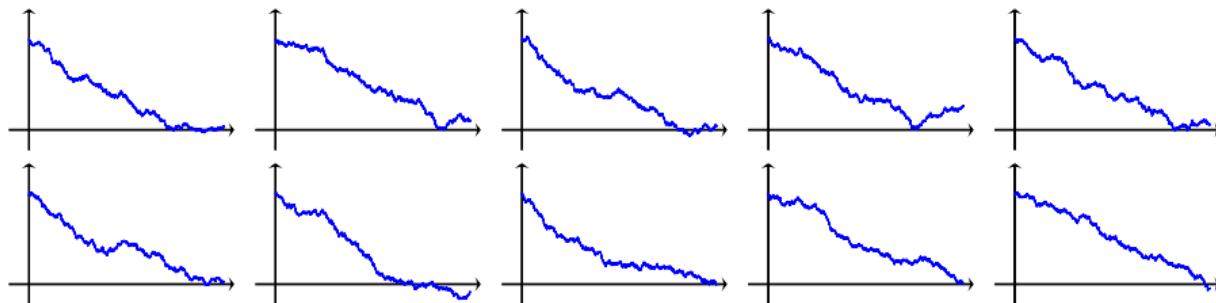


Sample paths of \bar{Y}_{2500} conditional on R for **heavy-tailed** claims:

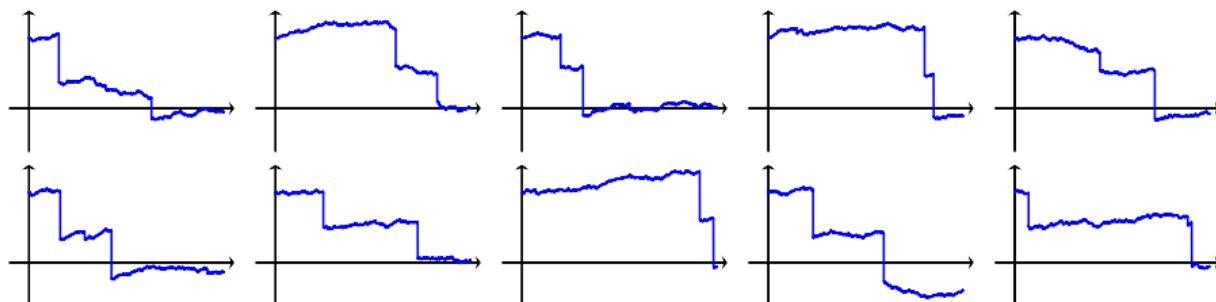


Bankruptcy Despite of Reinsurance

Sample paths of \bar{Y}_{5000} conditional on R for **light-tailed** claims:

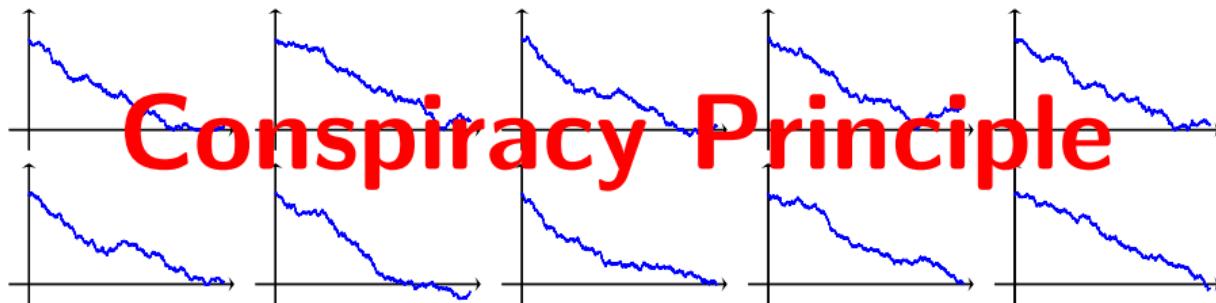


Sample paths of \bar{Y}_{5000} conditional on R for **heavy-tailed** claims:



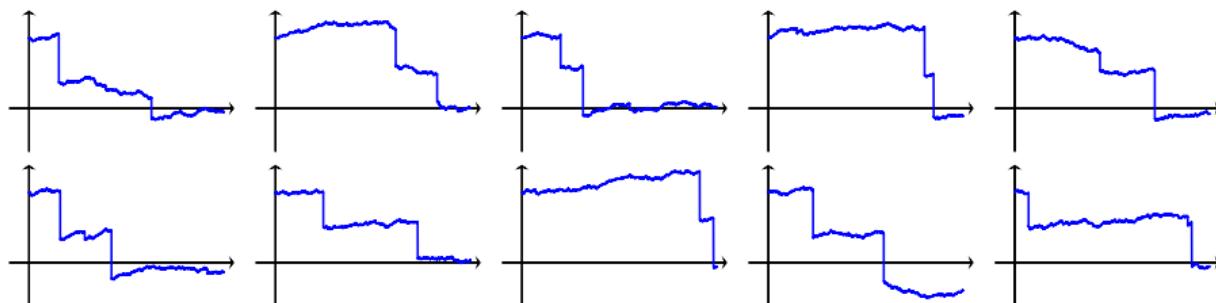
Bankruptcy Despite of Reinsurance

Sample paths of \bar{Y}_{5000} conditional on R for **light-tailed** claims:



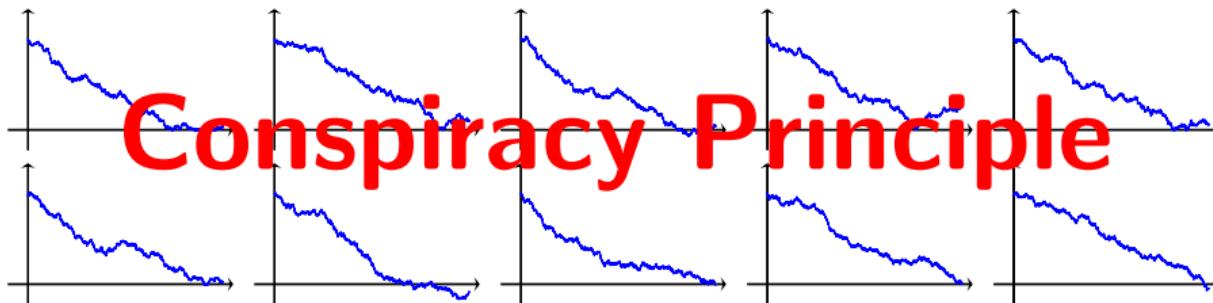
Conspiracy Principle

Sample paths of \bar{Y}_{5000} conditional on R for **heavy-tailed** claims:

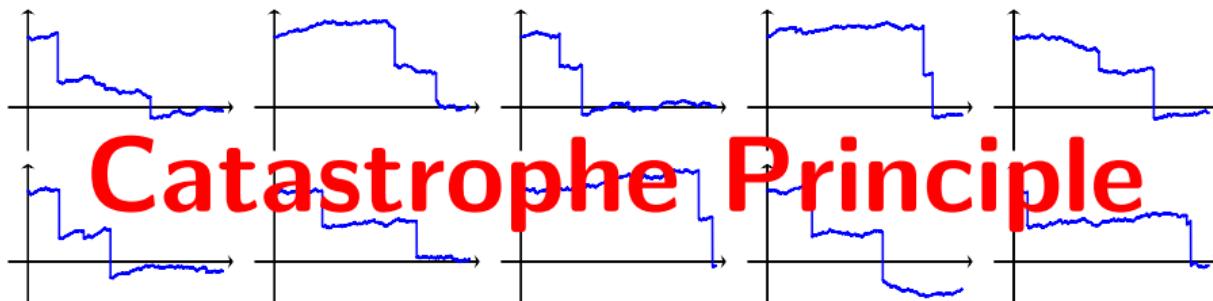


Bankruptcy Despite of Reinsurance

Sample paths of \bar{Y}_{5000} conditional on R for **light-tailed** claims:



Sample paths of \bar{Y}_{5000} conditional on R for **heavy-tailed** claims:



Sample Path Large Deviations:

Characterization of Conspiracy and Catastrophe Principles

Light-Tailed Large Deviations

$$\bar{S}_n(t) = \frac{1}{n} \sum_{i=1}^{\lfloor nt \rfloor} X_i, \quad X_i: \text{centered iid r.v. with } \Lambda(\lambda) = \log \mathbf{E} \exp(\lambda X_1) < \infty, \quad \forall \lambda > 0$$

Light-Tailed Large Deviations

$$\bar{S}_n(t) = \frac{1}{n} \sum_{i=1}^{\lfloor nt \rfloor} X_i, \quad X_i: \text{centered iid r.v. with } \Lambda(\lambda) = \log \mathbf{E} \exp(\lambda X_1) < \infty, \quad \forall \lambda > 0$$

Theorem (Mogulskii, 1993)

For any measurable $A \subseteq \mathbb{D}$

$$-\inf_{f \in A^\circ} I(f) \leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \mathbf{P}(\bar{S}_n \in A) \leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log \mathbf{P}(\bar{S}_n \in A) \leq -\inf_{f \in A^-} I(f)$$

$$\bullet \quad I(f) = \begin{cases} \int_0^1 \Lambda^*(\dot{f}(t)) dt & \text{if } f \in \mathcal{AC}, \ f(0) = 0 \\ 0 & \text{otherwise} \end{cases}$$

Light-Tailed Large Deviations

$$\bar{S}_n(t) = \frac{1}{n} \sum_{i=1}^{\lfloor nt \rfloor} X_i, \quad X_i: \text{centered iid r.v. with } \Lambda(\lambda) = \log \mathbf{E} \exp(\lambda X_1) < \infty, \quad \forall \lambda > 0$$

Theorem (Mogulskii, 1993)

For any measurable $A \subseteq \mathbb{D}$

$$\mathbf{P}(\bar{S}_n \in A) \stackrel{\log}{\sim} \exp \left(-n \cdot \inf_{f \in A} I(f) \right)$$

$$\bullet \quad I(f) = \begin{cases} \int_0^1 \Lambda^*(\dot{f}(t)) dt & \text{if } f \in \mathcal{AC}, \quad f(0) = 0 \\ 0 & \text{otherwise} \end{cases}$$

Implication: The Conspiracy Principle

Theorem (Ganesh et al. 2003; Lemma 4.2)

Suppose that X_n satisfies an LDP and

$$\inf_{f \in A} I(f) = k < \infty.$$

Under certain regularity conditions,

$$\mathbf{P}(\bar{S}_n \notin B | \bar{S}_n \in A) \rightarrow 0$$

for any neighborhood B of $\{x \in A : I(x) = k\}$.

Implication: The Conspiracy Principle

Theorem (Ganesh et al. 2003; Lemma 4.2)

Suppose that X_n satisfies an LDP and

$$\inf_{f \in A} I(f) = k < \infty.$$

Under certain regularity conditions,

$$\mathbf{P}(\bar{S}_n \notin B | \bar{S}_n \in A) \rightarrow 0$$

for any neighborhood B of $\{x \in A : I(x) = k\}$.

Rigorous Characterization of Conspiracy Principle

Heavy-Tailed Large Deviations

$$\bar{S}_n(t) = \frac{1}{n} \sum_{i=1}^{\lfloor nt \rfloor} X_i,$$

X_i : centered iid r.v. with $\mathbf{P}(X_i \geq x) = x^{-(\alpha+1)}, \quad \alpha > 0$

Heavy-Tailed Large Deviations

$$\bar{S}_n(t) = \frac{1}{n} \sum_{i=1}^{\lfloor nt \rfloor} X_i, \quad X_i: \text{centered iid r.v. with } \mathbf{P}(X_i \geq x) = x^{-(\alpha+1)}, \quad \alpha > 0$$

Theorem (R., Blanchet, Zwart, 2019)

For “general” $A \subseteq \mathbb{D}$

$$C(A^\circ) \leq \liminf_{n \rightarrow \infty} \frac{\mathbf{P}(\bar{S}_n \in A)}{n^{-\alpha \mathcal{J}(A)}} \leq \limsup_{n \rightarrow \infty} \frac{\mathbf{P}(\bar{S}_n \in A)}{n^{-\alpha \mathcal{J}(A)}} \leq C(A^-).$$

- $\mathcal{J}(A)$: min #jumps for step functions to be inside A
- $C(\cdot)$: a measure

Heavy-Tailed Large Deviations

$$\bar{S}_n(t) = \frac{1}{n} \sum_{i=1}^{\lfloor nt \rfloor} X_i, \quad X_i: \text{centered iid r.v. with } \mathbf{P}(X_i \geq x) = x^{-(\alpha+1)}, \quad \alpha > 0$$

Theorem (R., Blanchet, Zwart, 2019)

For “general” $A \subseteq \mathbb{D}$

$$\mathbf{P}(\bar{S}_n \in A) \sim n^{-\alpha \mathcal{J}(A)}$$

- $\mathcal{J}(A)$: min #jumps for step functions to be inside A

Heavy-Tailed Large Deviations

$$\bar{S}_n(t) = \frac{1}{n} \sum_{i=1}^{\lfloor nt \rfloor} X_i, \quad X_i: \text{centered iid r.v. with } \mathbf{P}(X_i \geq x) = x^{-(\alpha+1)}, \quad \alpha > 0$$

Theorem (R., Blanchet, Zwart, 2019)

For “general” $A \subseteq \mathbb{D}$

$$\mathbf{P}(\bar{S}_n \in A) \sim n^{-\alpha \mathcal{J}(A)}$$

LD power index ↗

- $\mathcal{J}(A)$: min #jumps for step functions to be inside A

Implication: The Catastrophe Principle

Theorem (R., Blanchet, Zwart. 2019)

Under certain regularity conditions on A,

$$\mathbf{P}(\bar{S}_n \in \cdot | \bar{S}_n \in A) \rightarrow \mathbf{P}(\bar{S}_{|A} \in \cdot) = \frac{C(\cdot \cap A)}{C(A)}$$

$\bar{S}_{|A}$: a (random) piecewise-constant function with $\mathcal{J}(A)$ jumps.

Implication: The Catastrophe Principle

Theorem (R., Blanchet, Zwart. 2019)

Under certain regularity conditions on A,

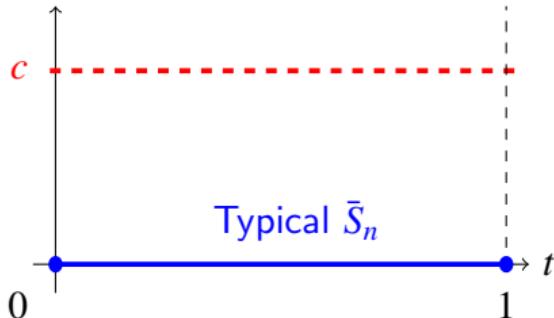
$$\mathbf{P}(\bar{S}_n \in \cdot | \bar{S}_n \in A) \rightarrow \mathbf{P}(\bar{S}_{|A} \in \cdot) = \frac{C(\cdot \cap A)}{C(A)}$$

$\bar{S}_{|A}$: a (random) piecewise-constant function with $\mathcal{J}(A)$ jumps.

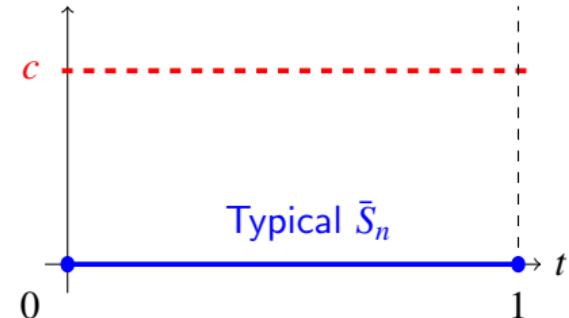
Rigorous Characterization of Catastrophe Principle

Conspiracy vs Catastrophe

$$\bar{S}_n \in \{ f \in \mathbb{D} : f \text{ crosses level } c \text{ on } [0, 1] \} = A$$



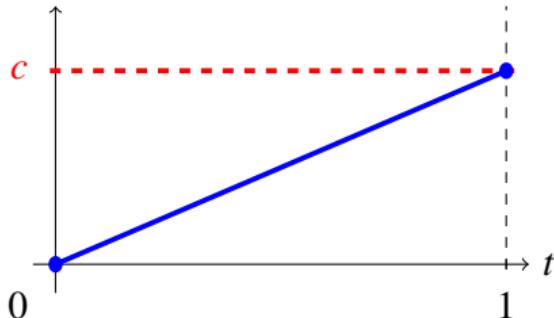
Light-Tailed Claim Size



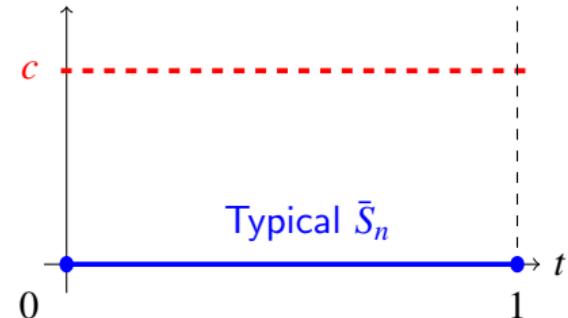
Heavy-Tailed Claim Size

Conspiracy vs Catastrophe

$$\bar{S}_n \in \{ f \in \mathbb{D} : f \text{ crosses level } c \text{ on } [0, 1] \} = A$$



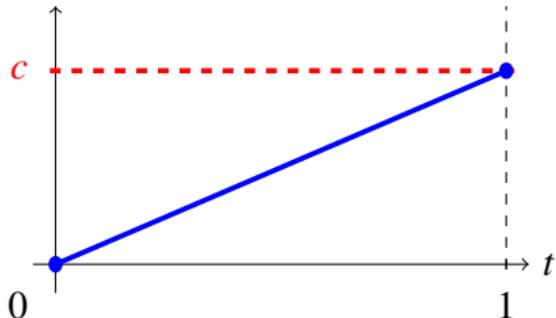
Light-Tailed Claim Size



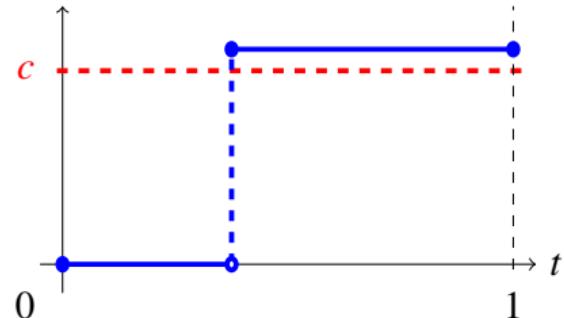
Heavy-Tailed Claim Size

Conspiracy vs Catastrophe

$$\bar{S}_n \in \{ f \in \mathbb{D} : f \text{ crosses level } c \text{ on } [0, 1] \} = A$$



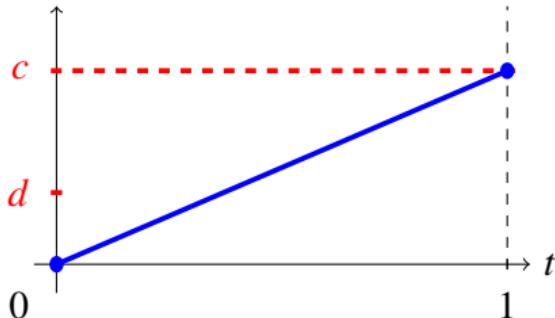
Light-Tailed Claim Size



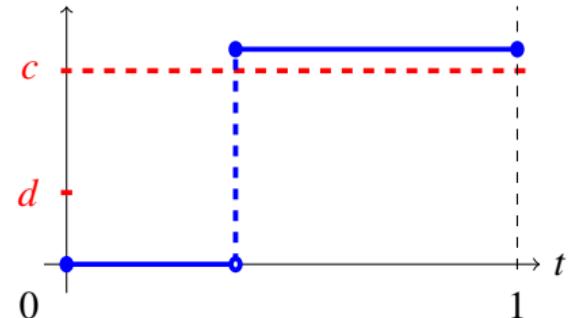
Heavy-Tailed Claim Size

Conspiracy vs Catastrophe

$$\bar{S}_n \in \{ f \in \mathbb{D} : f \text{ crosses level } c \text{ on } [0,1] \text{ & jump sizes } \leq d \} = A$$



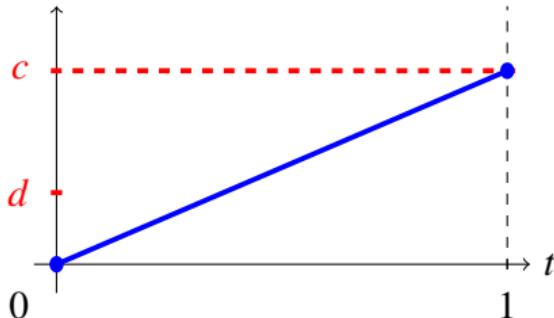
Light-Tailed Claim Size



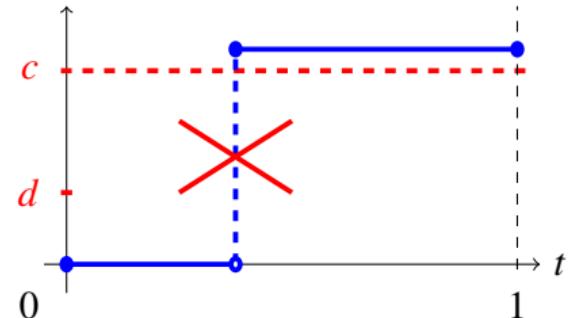
Heavy-Tailed Claim Size

Conspiracy vs Catastrophe

$$\bar{S}_n \in \{ f \in \mathbb{D} : f \text{ crosses level } c \text{ on } [0,1] \text{ & jump sizes } \leq d \} = A$$



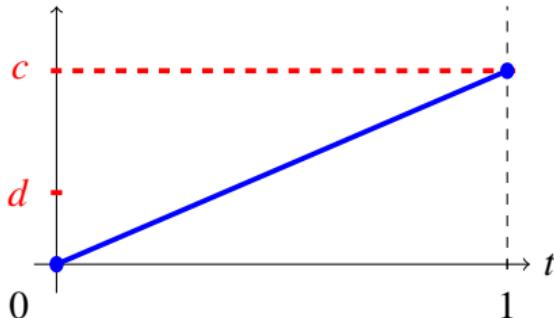
Light-Tailed Claim Size



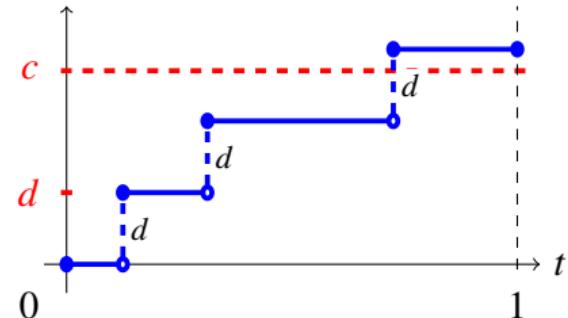
Heavy-Tailed Claim Size

Conspiracy vs Catastrophe

$$\bar{S}_n \in \{ f \in \mathbb{D} : f \text{ crosses level } c \text{ on } [0,1] \text{ & jump sizes } \leq d \} = A$$



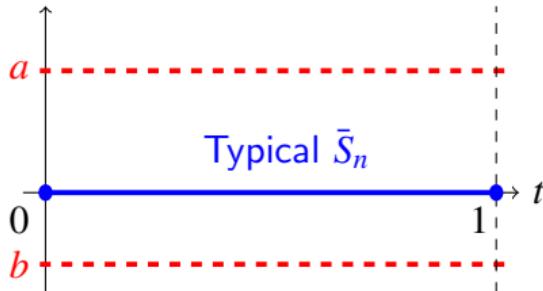
Light-Tailed Claim Size



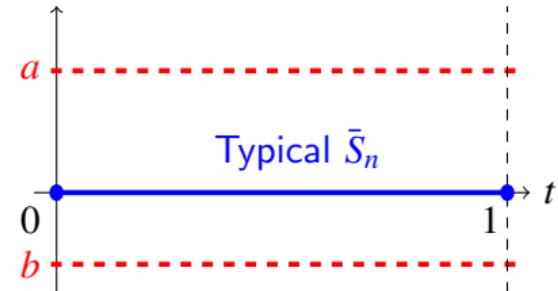
Heavy-Tailed Claim Size

Conspiracy vs Catastrophe

$\bar{S}_n \in \{ f \in \mathbb{D} : f \text{ hits below } b \text{ on } [0, 1] \text{ and ends up above } a \} = A$



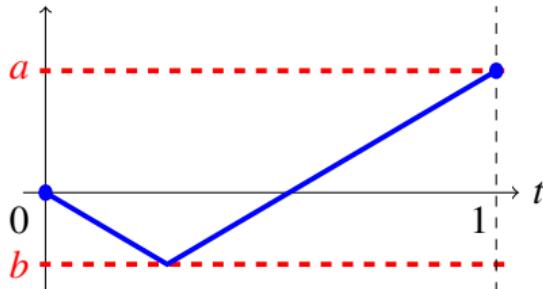
Light-Tailed Increments



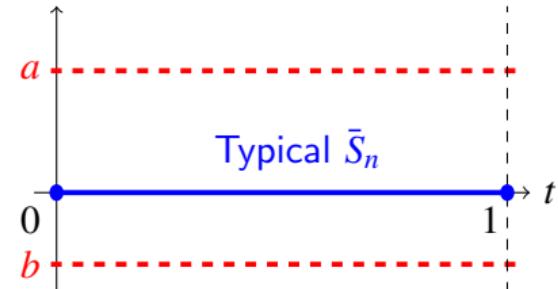
Heavy-Tailed Increments

Conspiracy vs Catastrophe

$\bar{S}_n \in \{ f \in \mathbb{D} : f \text{ hits below } b \text{ on } [0, 1] \text{ and ends up above } a \} = A$



Light-Tailed Increments



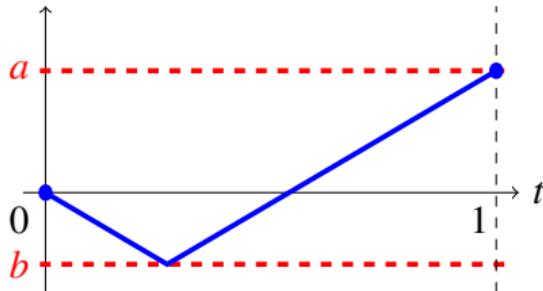
Heavy-Tailed Increments

Conspiracy

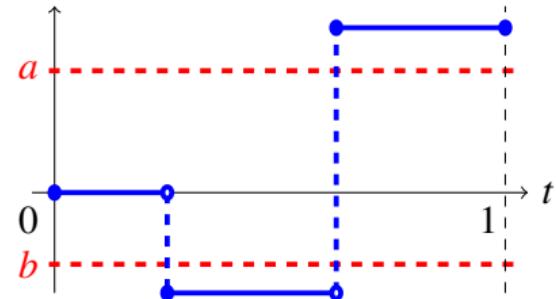
vs

Catastrophe

$$\bar{S}_n \in \{ f \in \mathbb{D} : f \text{ hits below } b \text{ on } [0, 1] \text{ and ends up above } a \} = A$$



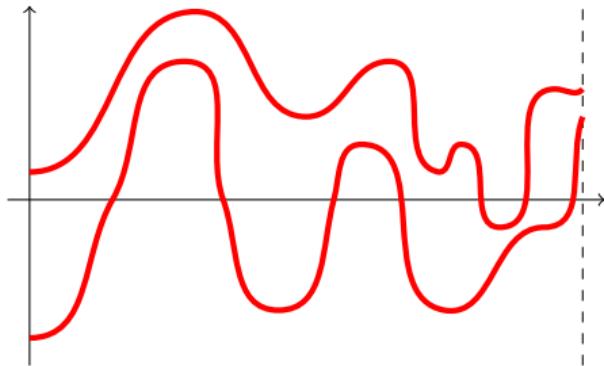
Light-Tailed Increments



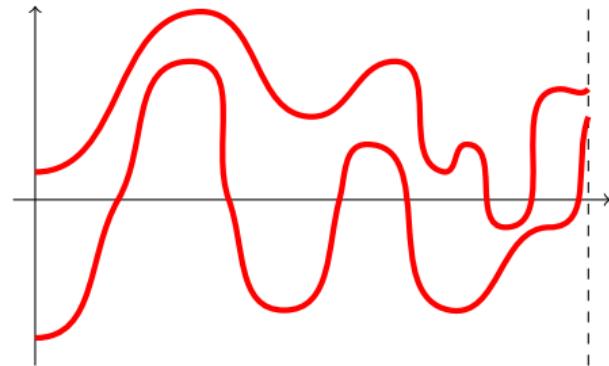
Heavy-Tailed Increments

Conspiracy vs Catastrophe

$\bar{S}_n \in \{ f \in \mathbb{D} : f \text{ lies between the two red curves} \}$



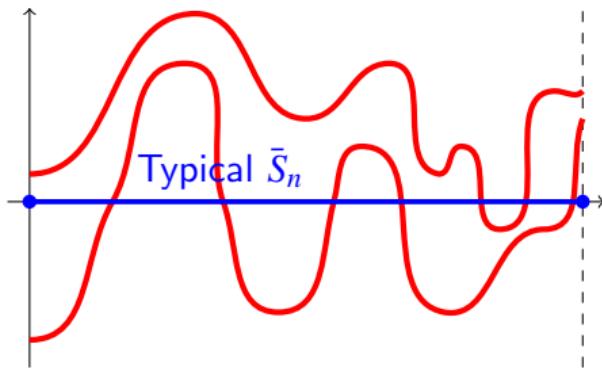
Light-Tailed Increments



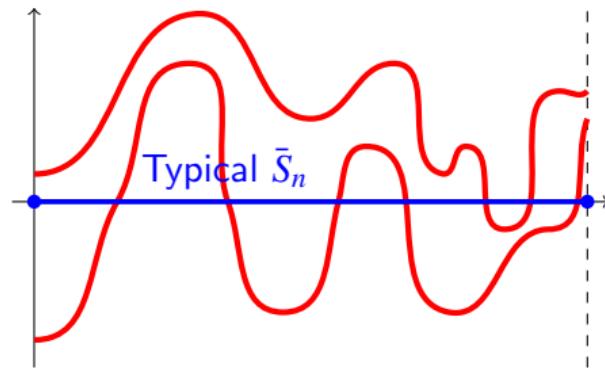
Heavy-Tailed Increments

Conspiracy vs Catastrophe

$\bar{S}_n \in \{ f \in \mathbb{D} : f \text{ lies between the two red curves} \}$



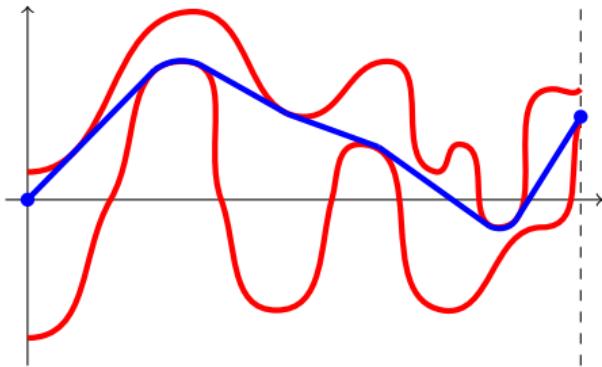
Light-Tailed Increments



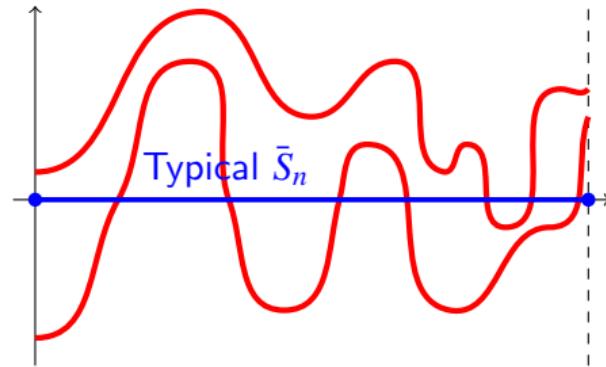
Heavy-Tailed Increments

Conspiracy vs Catastrophe

$\bar{S}_n \in \{ f \in \mathbb{D} : f \text{ lies between the two red curves} \}$



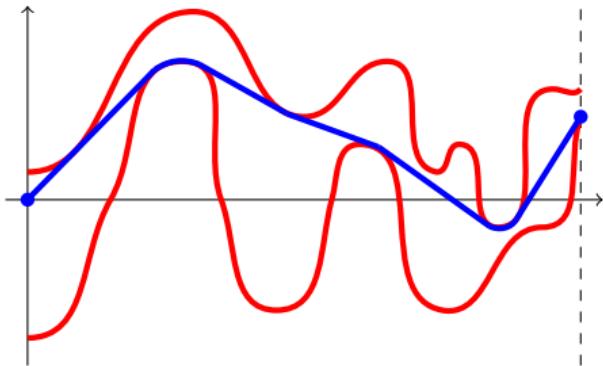
Light-Tailed Increments



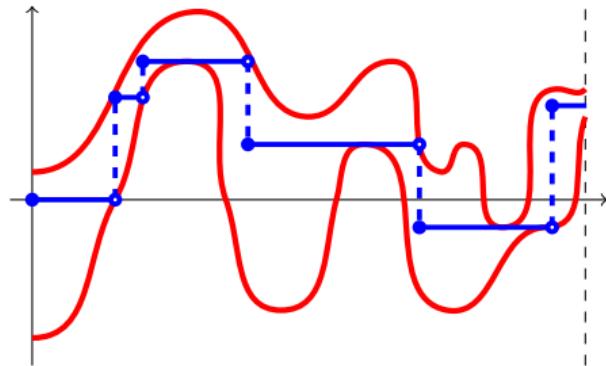
Heavy-Tailed Increments

Conspiracy vs Catastrophe

$\bar{S}_n \in \{ f \in \mathbb{D} : f \text{ lies between the two red curves} \}$



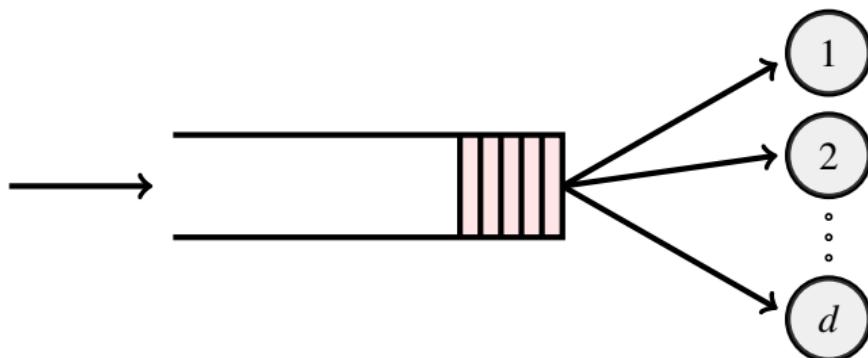
Light-Tailed Increments



Heavy-Tailed Increments

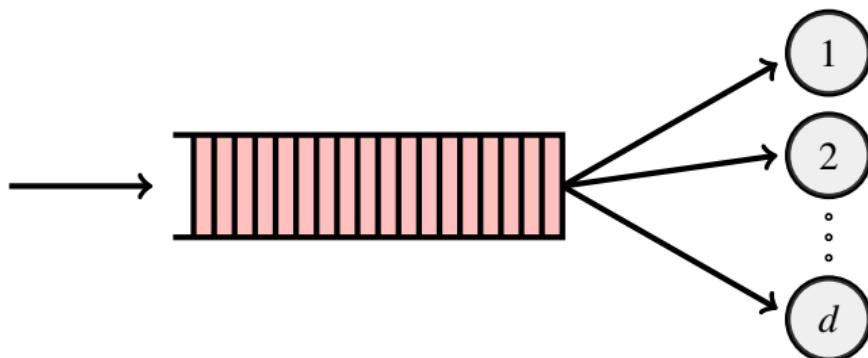
Conspiracy vs Catastrophe

Congestion of Multiple Server Queue:



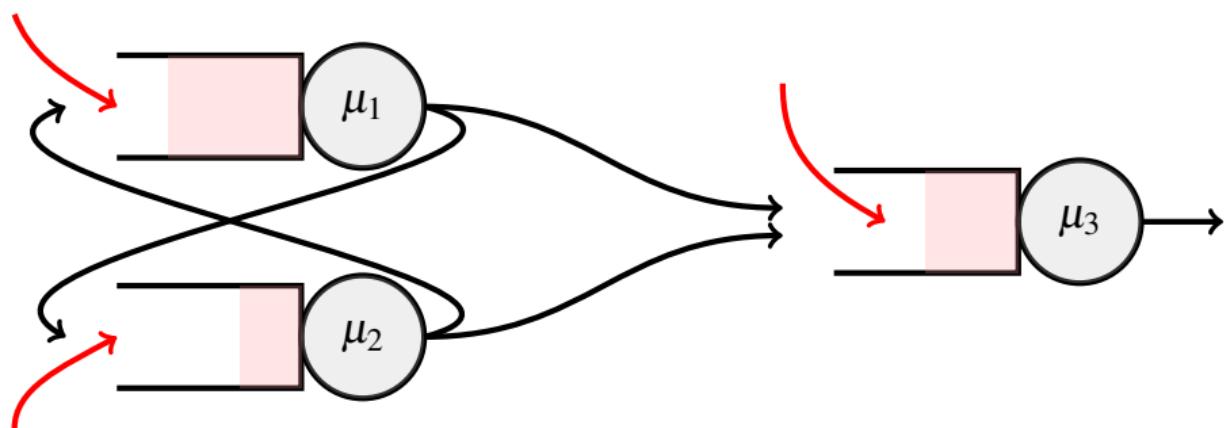
Conspiracy vs Catastrophe

Congestion of Multiple Server Queue:



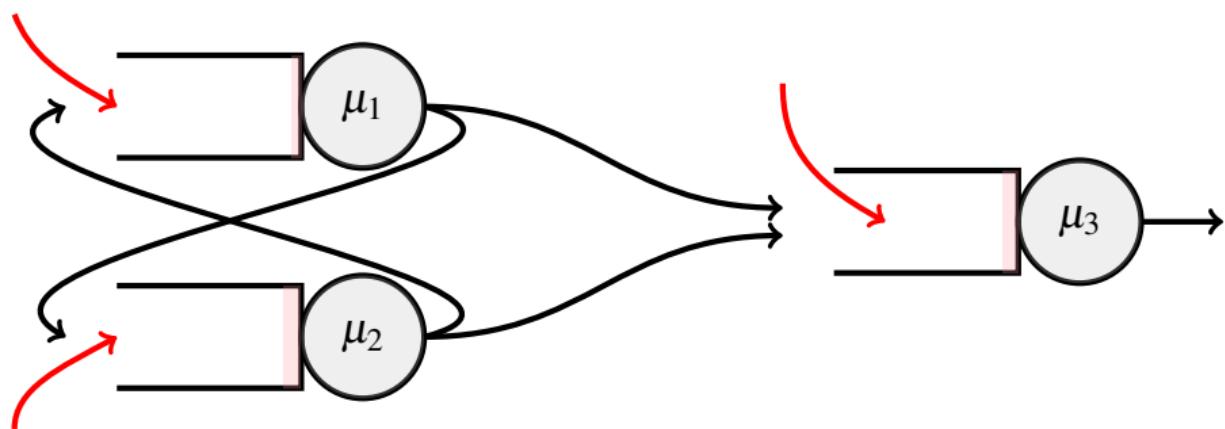
Conspiracy vs Catastrophe

Congestion of Stochastic Fluid Networks:



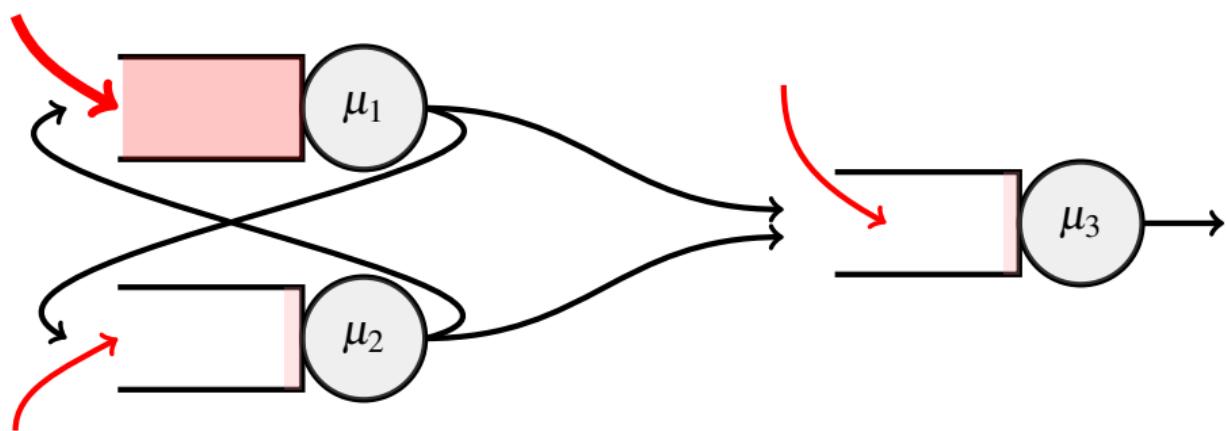
Conspiracy vs Catastrophe

Congestion of Stochastic Fluid Networks:



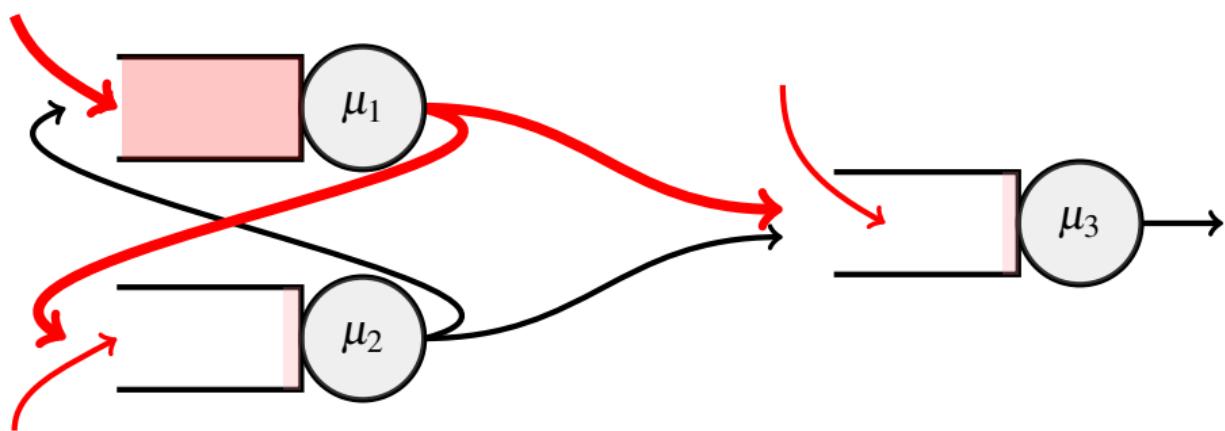
Conspiracy vs Catastrophe

Congestion of Stochastic Fluid Networks:



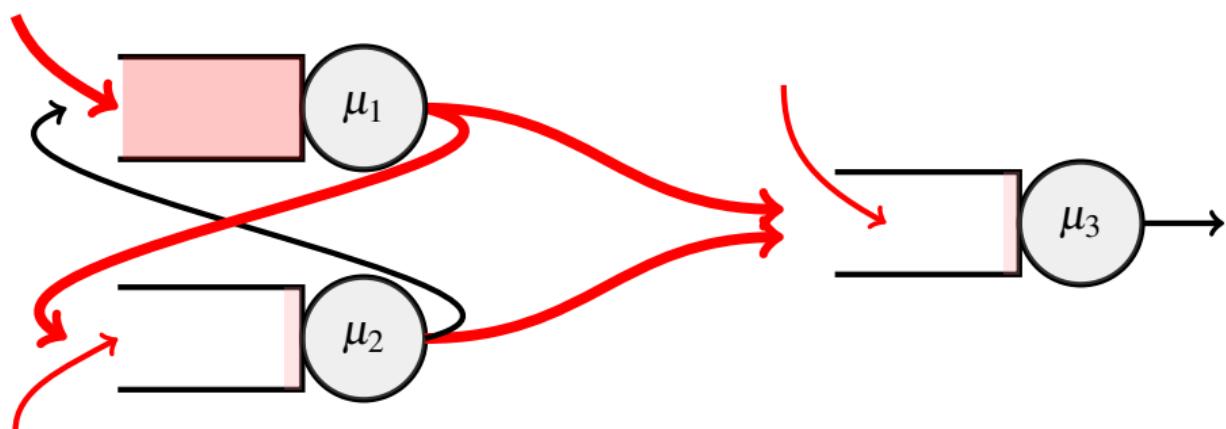
Conspiracy vs Catastrophe

Congestion of Stochastic Fluid Networks:



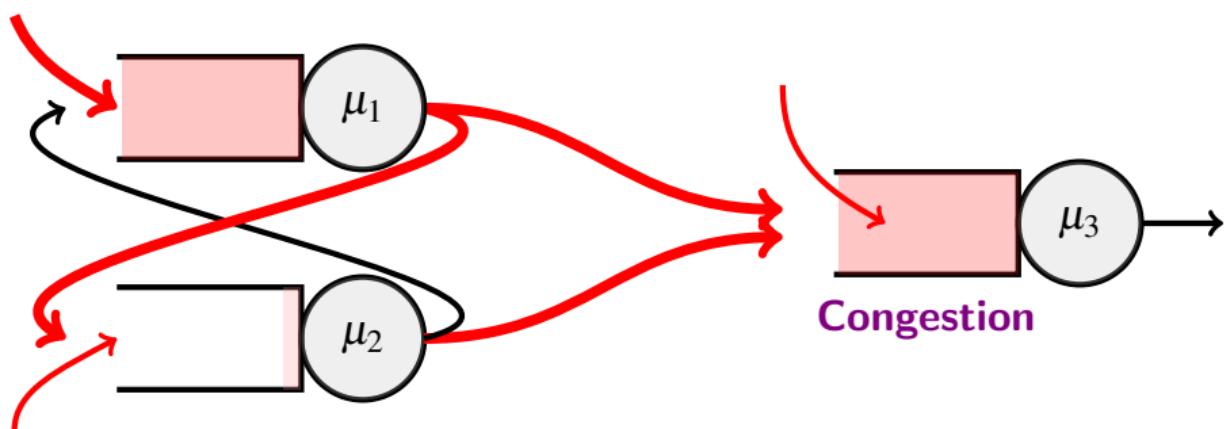
Conspiracy vs Catastrophe

Congestion of Stochastic Fluid Networks:



Conspiracy vs Catastrophe

Congestion of Stochastic Fluid Networks:



The point is

Many heavy-tailed rare events can be written as

$$\{\bar{S}_n \in A\}$$

and the decay rate is determined by $\mathcal{J}(A)$, i.e.,

$$\mathbf{P}(\bar{S}_n \in A) \sim n^{-\alpha \mathcal{J}(A)}$$

Catastrophe Principle Extends to General Heavy-Tailed Systems

Heavy-Tailed Large Deviations for

- Continuous-Time Processes

R., Blanchet, Zwart (2019), Bazhba, Blanchet, R., Zwart (2020), Wang, R. (2024+)

- Processes with Spatial and Temporal Correlations

Bazhba, R., Zwart (2022), Chen, R., Zwart (2024), Bazhba, Blanchet, R., Zwart (2024+), Su, R. (2024+), Wang, R. (2024+)

Minimal # Jumps added to Typical Paths Characterize the Catastrophe Principle

Back to Stochastic Gradient Descent

Stochastic Gradient Descent

Minimizing loss function f :

$$W_{k+1} = W_k - \eta(f'(W_k)) \quad k = 0, 1, 2, \dots$$

Stochastic Gradient Descent

Minimizing loss function f :

$$W_{k+1} = W_k - \eta (\tilde{f}'(W_k)) \quad k = 0, 1, 2, \dots$$

Stochastic Gradient Descent

Minimizing loss function f :

$$W_{k+1} = W_k - \eta (f'(W_k) + Z_k) \quad k = 0, 1, 2, \dots$$

Stochastic Gradient Descent

Minimizing loss function f :

$$W_{k+1} = W_k - \eta (f'(W_k) + Z_k) \quad k = 0, 1, 2, \dots$$

Stochastic Gradient Descent

Minimizing loss function f :

$$W^{\eta}_{k+1} = W^{\eta}_k - \eta (f'(W^{\eta}_k) + Z_k) \quad k = 0, 1, 2, \dots$$

Stochastic Gradient Descent

Minimizing loss function f :

$$W^\eta(\eta(k+1)) = W^\eta(\eta k) - \eta(f'(W^\eta(\eta k)) + Z_k) \quad k = 0, 1, 2, \dots$$

Stochastic Gradient Descent

Minimizing loss function f :

$$W^\eta(\eta(k+1)) = W^\eta(\eta k) - \eta(f'(W^\eta(\eta k)) + Z_k) \quad k = 0, 1, 2, \dots$$

Then

$$W^\eta(\cdot) \rightarrow w(\cdot) \quad \text{as} \quad \eta \rightarrow 0$$

Stochastic Gradient Descent

Minimizing loss function f :

$$W^\eta(\eta(k+1)) = W^\eta(\eta k) - \eta(f'(W^\eta(\eta k)) + Z_k) \quad k = 0, 1, 2, \dots$$

Then

$$W^\eta(\cdot) \rightarrow w(\cdot) \quad \text{as} \quad \eta \rightarrow 0$$

where

$$dw(t) = -f'(w(t))dt$$

Stochastic Gradient Descent

Minimizing loss function f :

$$W^\eta(\eta(k+1)) = W^\eta(\eta k) - \eta(f'(W^\eta(\eta k)) + Z_k) \quad k = 0, 1, 2, \dots$$

Then

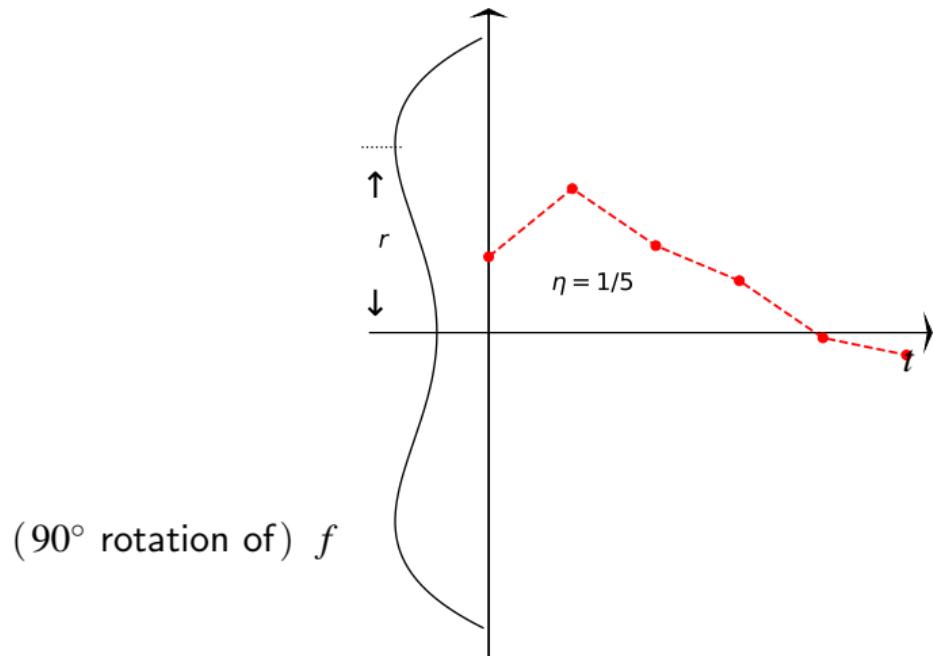
$$W^\eta(\cdot) \rightarrow w(\cdot) \quad \text{as} \quad \eta \rightarrow 0$$

↖
Gradient Flow

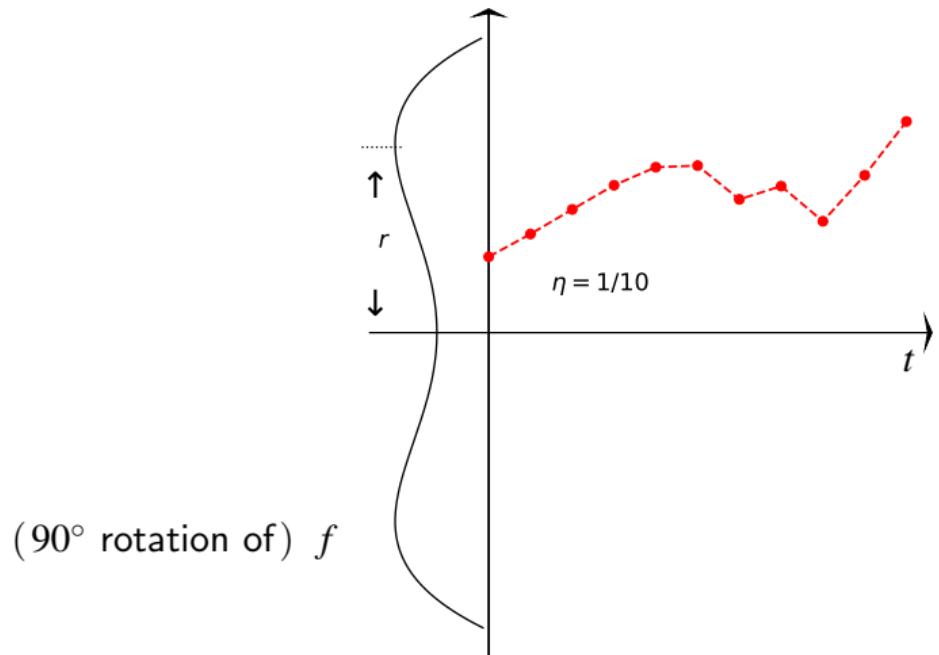
where

$$dw(t) = -f'(w(t))dt$$

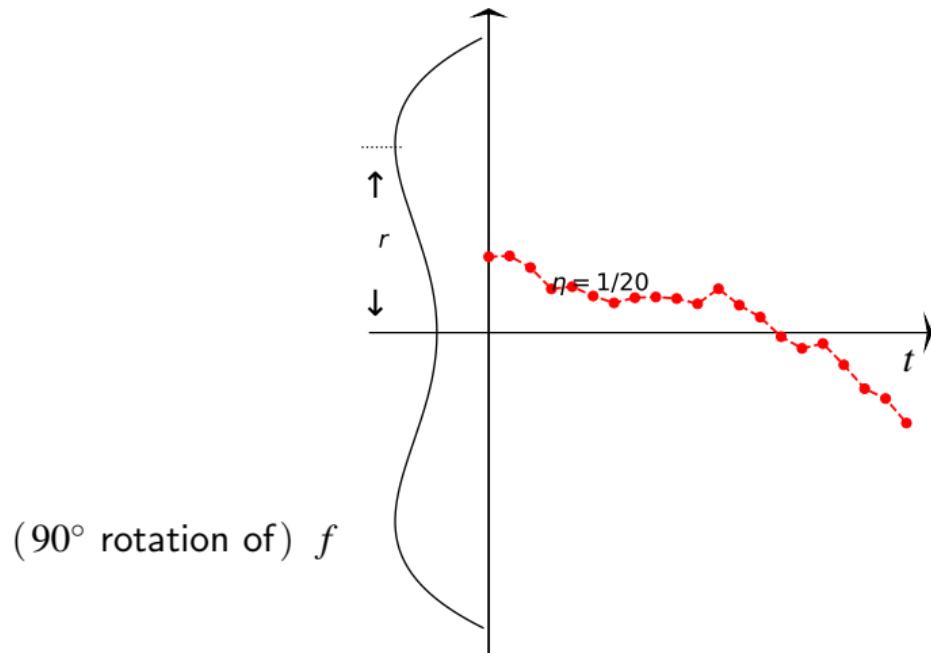
Gradient Flow: Law of Large Numbers



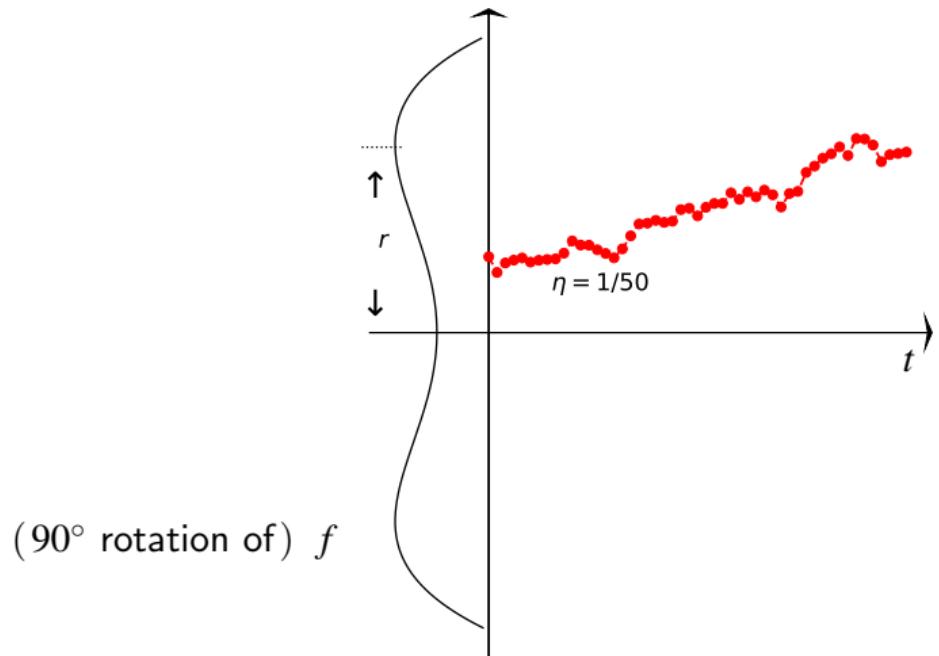
Gradient Flow: Law of Large Numbers



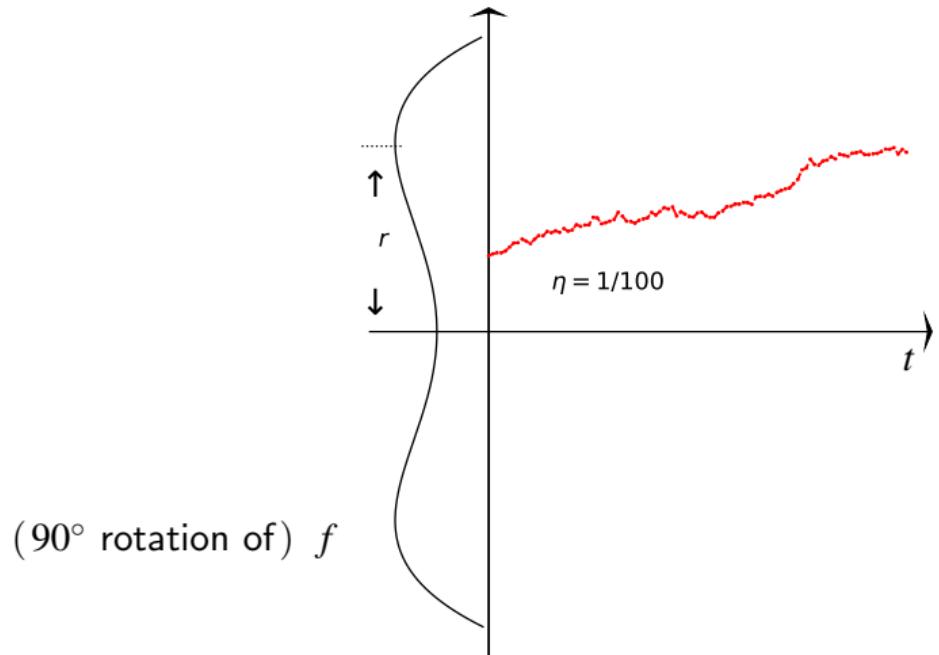
Gradient Flow: Law of Large Numbers



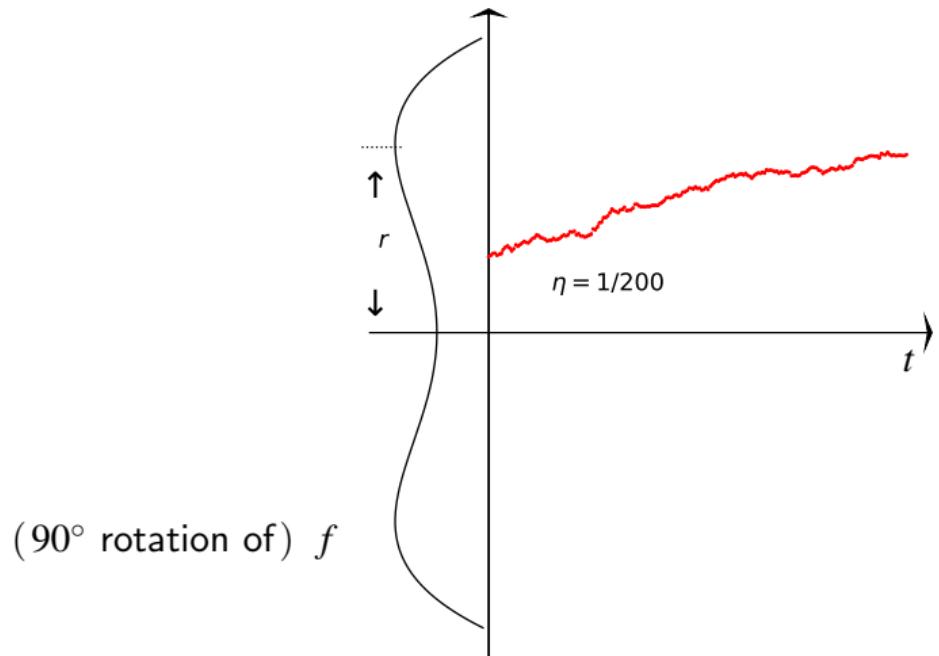
Gradient Flow: Law of Large Numbers



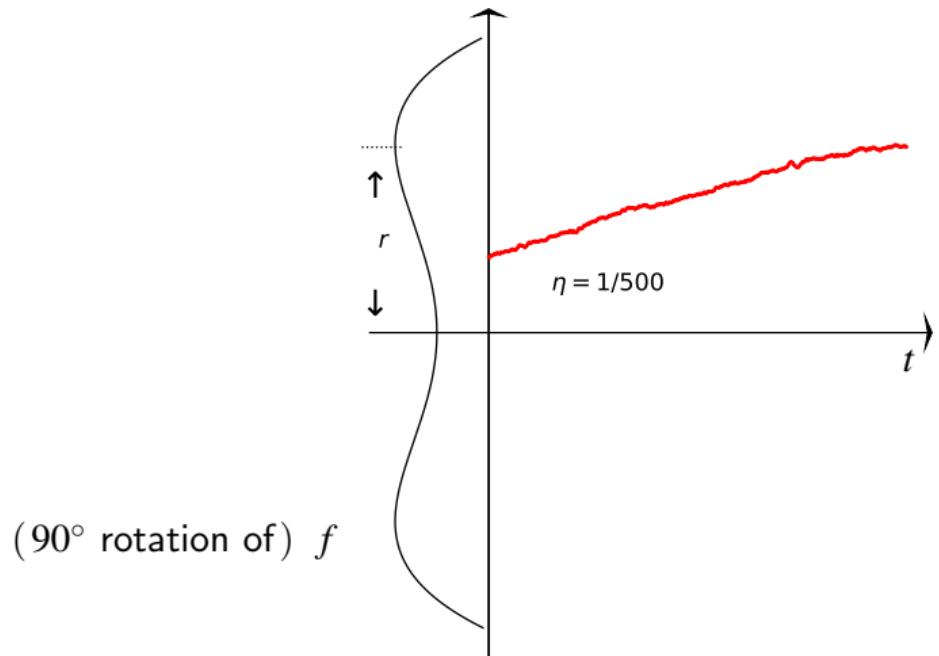
Gradient Flow: Law of Large Numbers



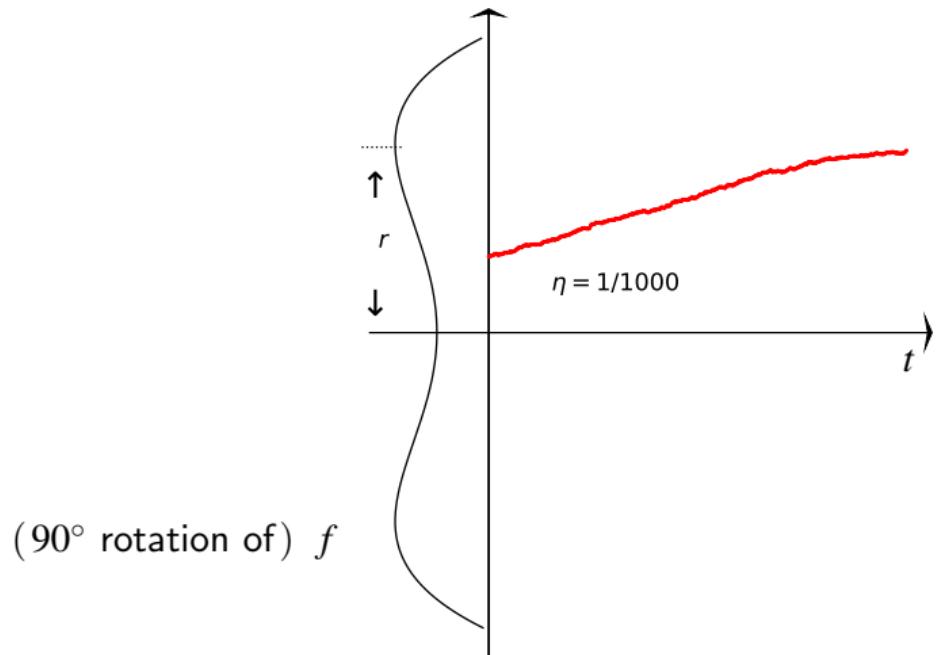
Gradient Flow: Law of Large Numbers



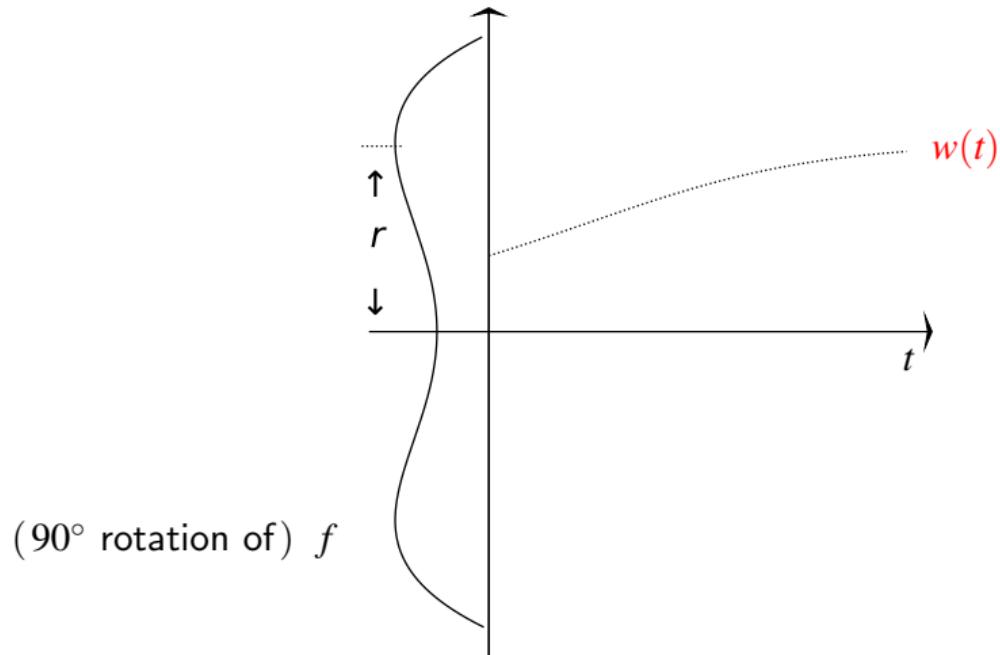
Gradient Flow: Law of Large Numbers



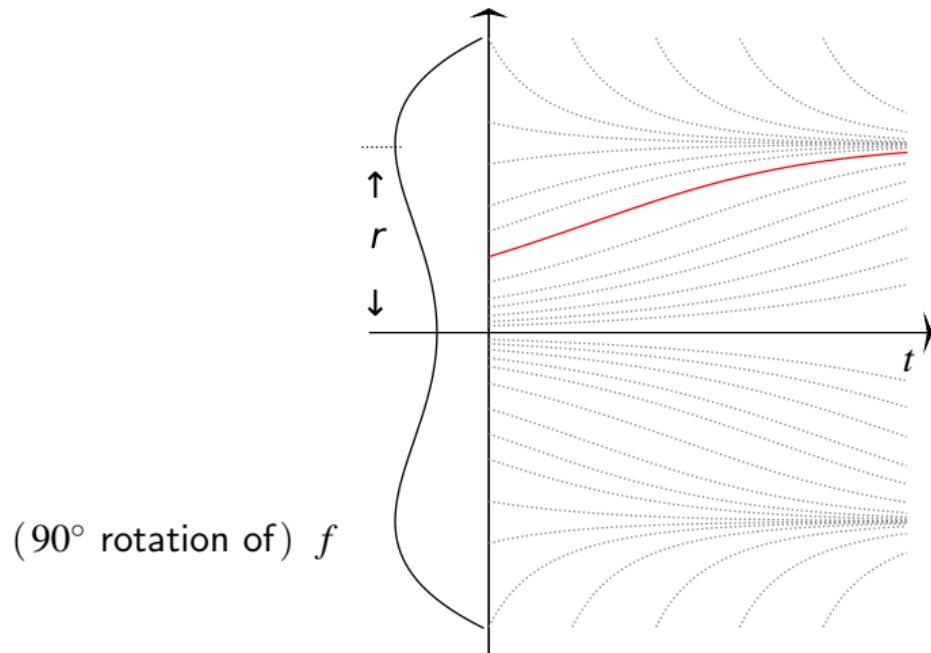
Gradient Flow: Law of Large Numbers



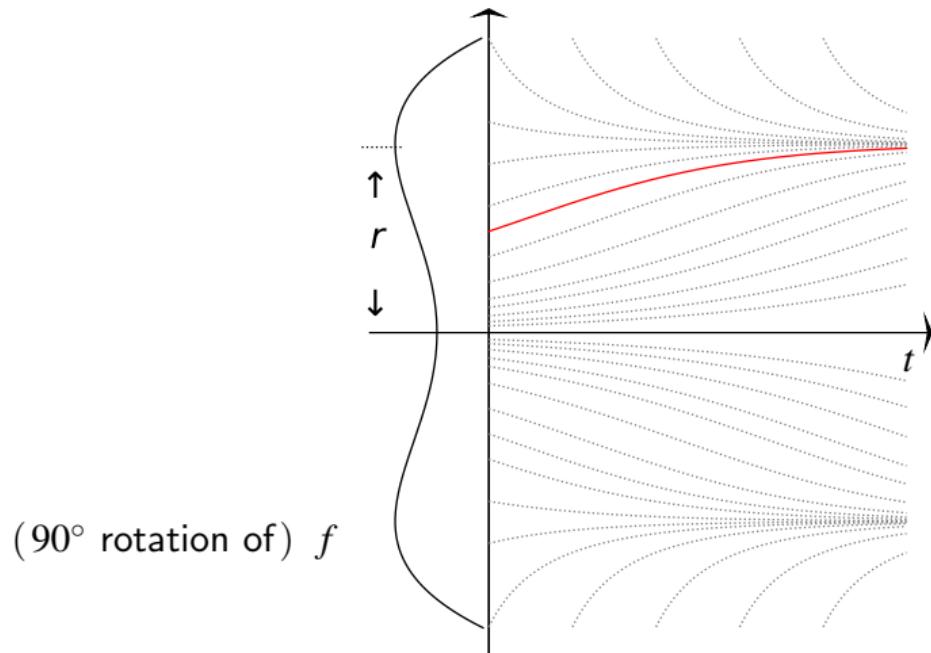
Gradient Flow: Law of Large Numbers



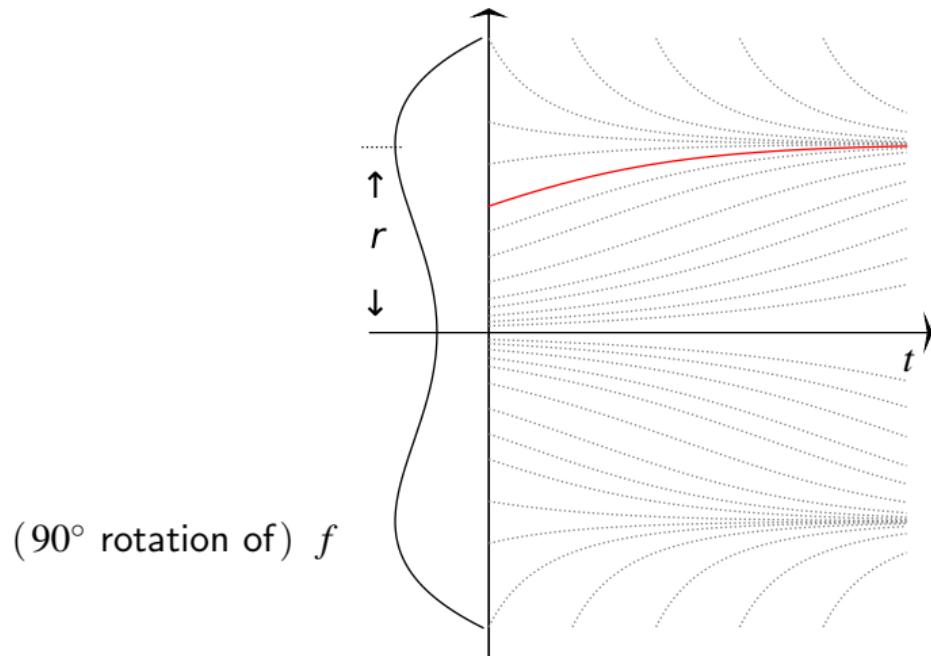
Gradient Flow: Law of Large Numbers



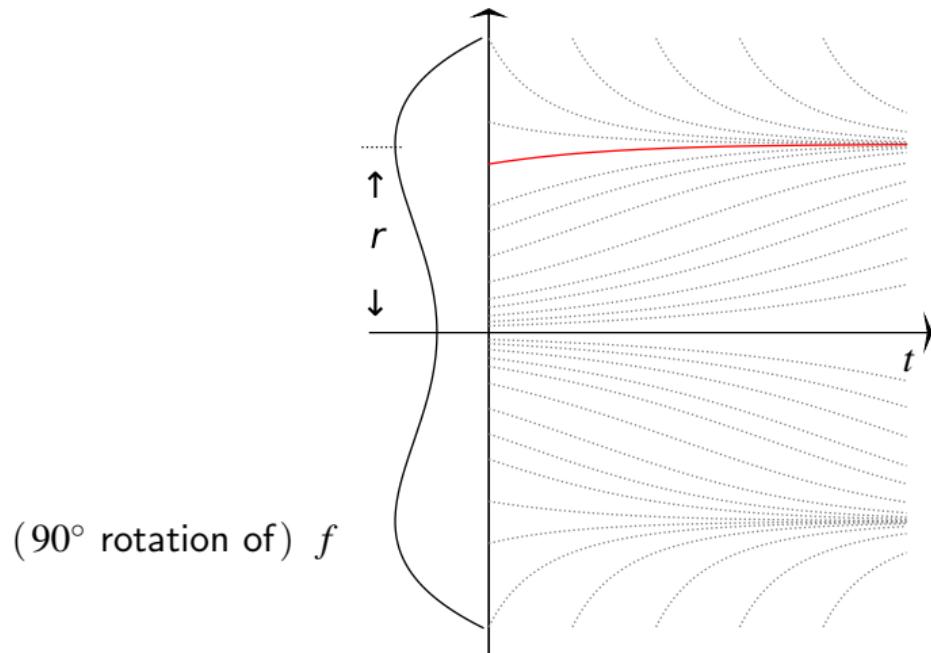
Gradient Flow: Law of Large Numbers



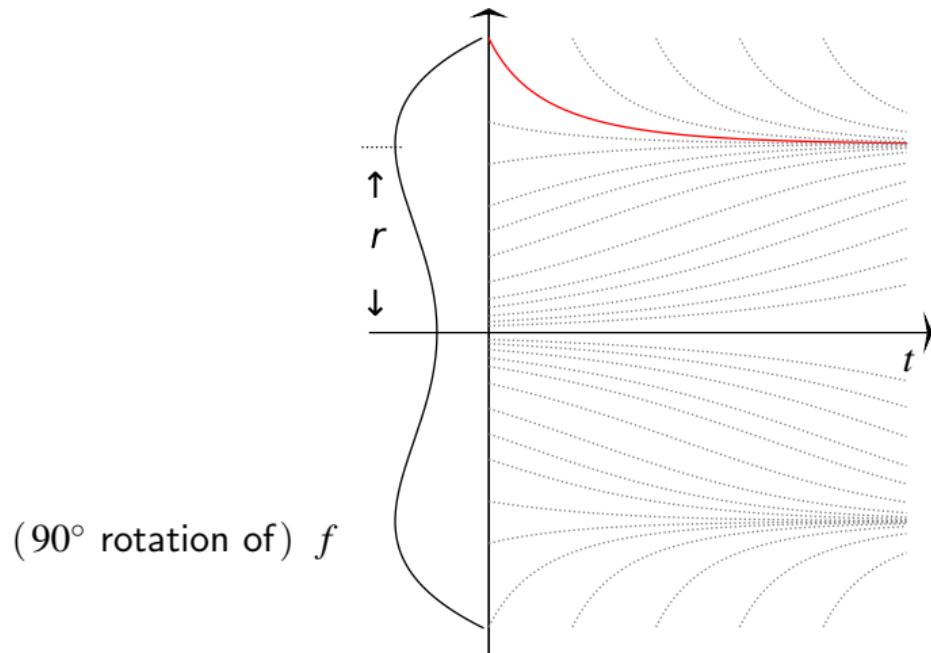
Gradient Flow: Law of Large Numbers



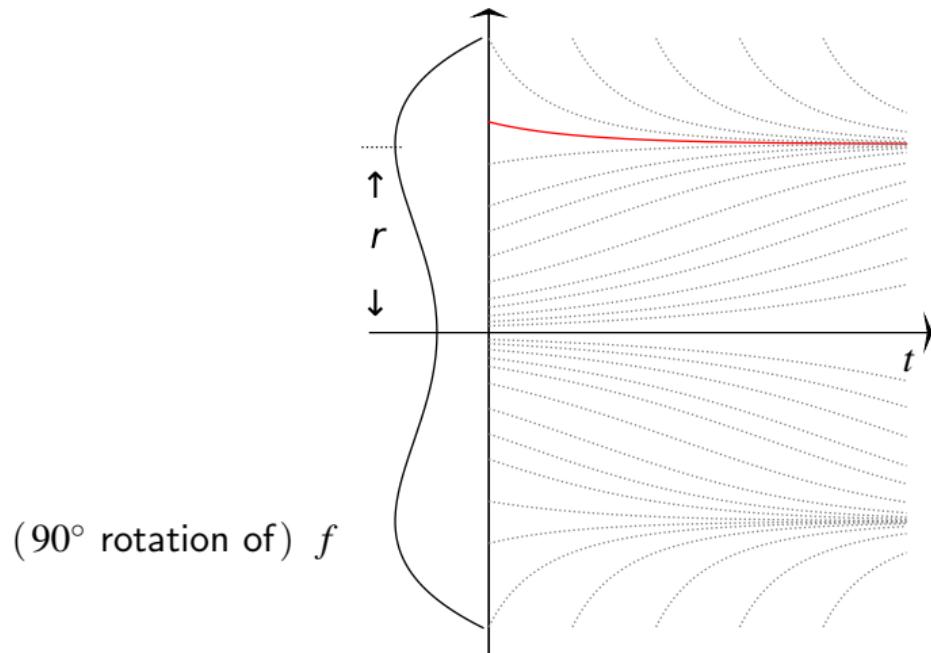
Gradient Flow: Law of Large Numbers



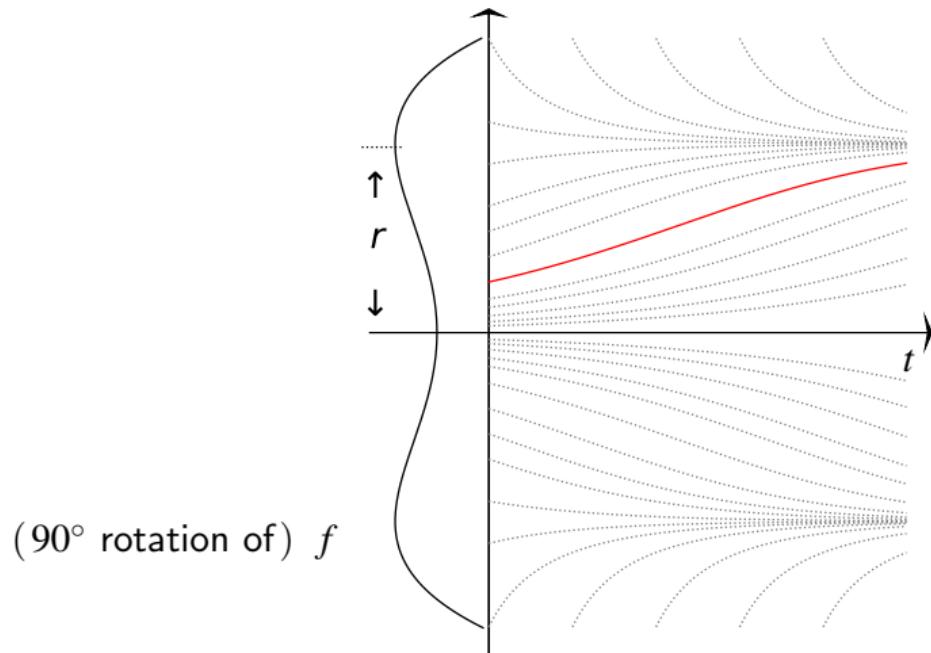
Gradient Flow: Law of Large Numbers



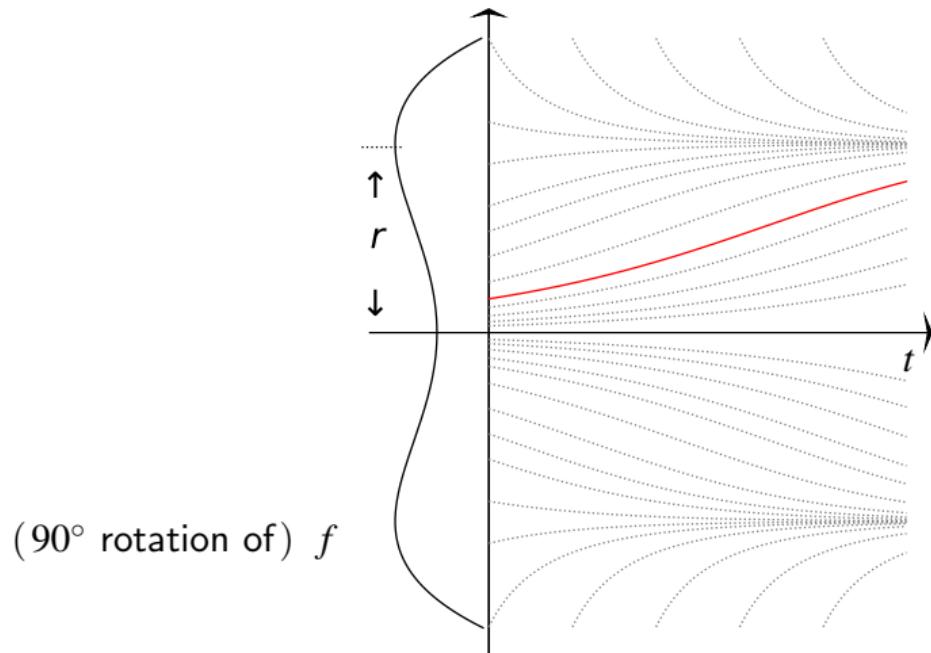
Gradient Flow: Law of Large Numbers



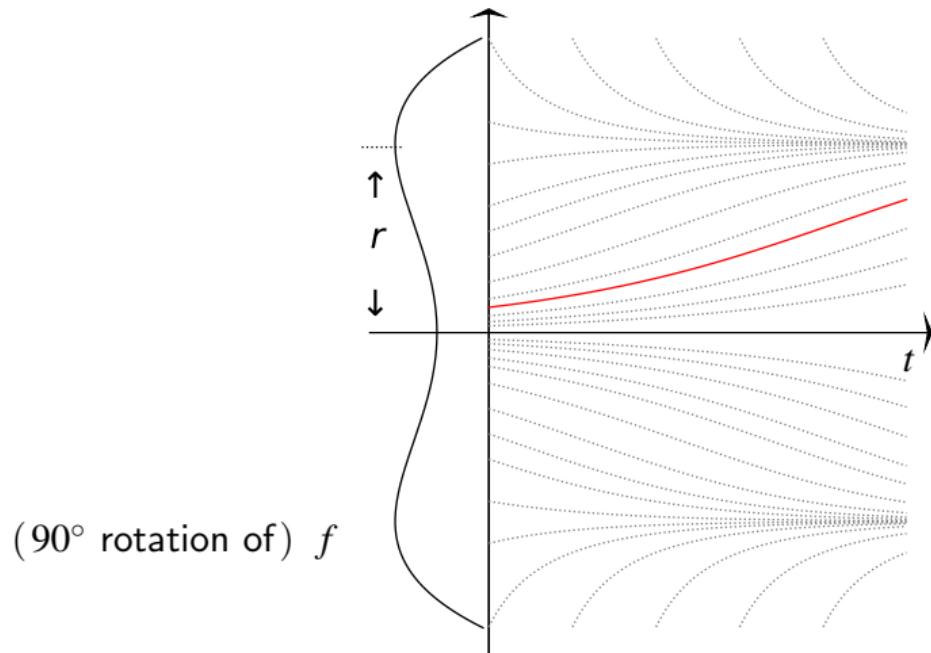
Gradient Flow: Law of Large Numbers



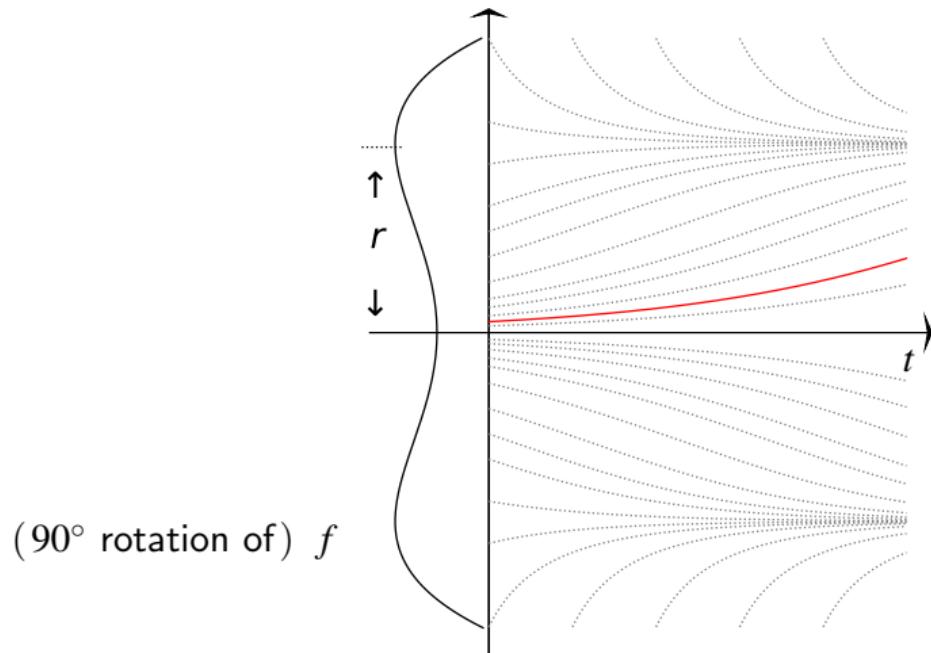
Gradient Flow: Law of Large Numbers



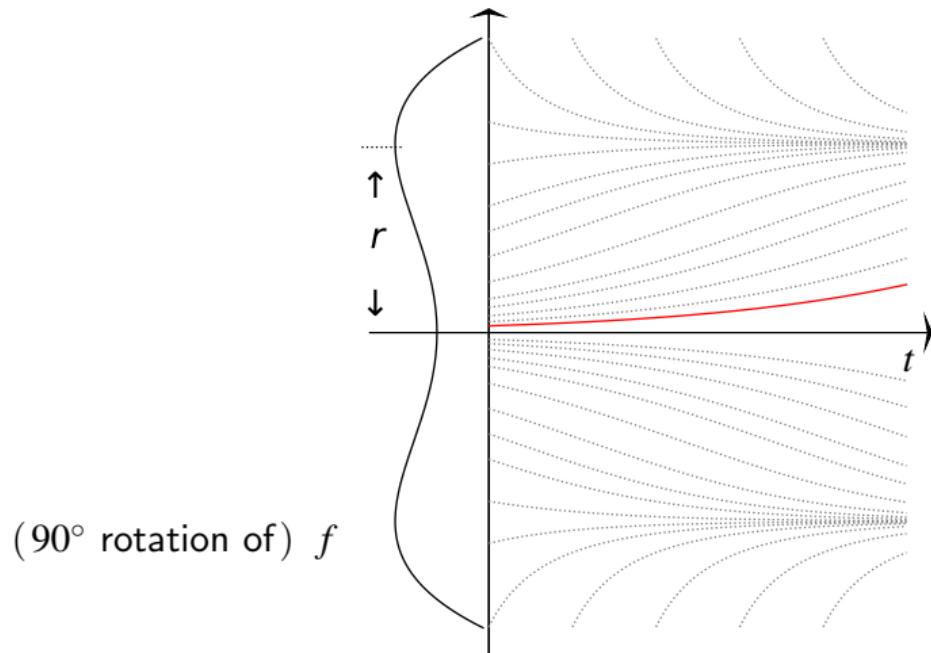
Gradient Flow: Law of Large Numbers



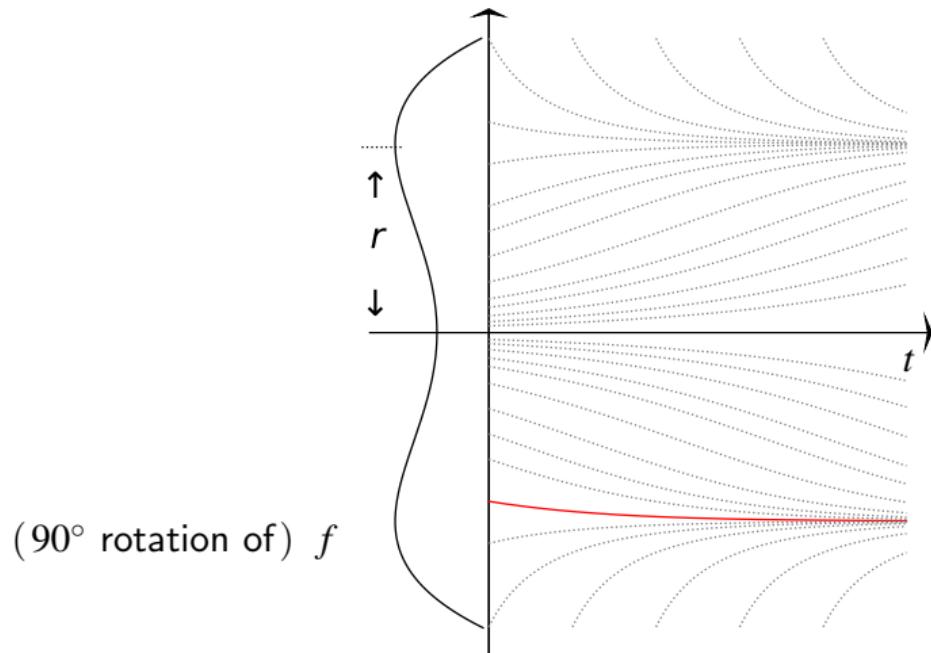
Gradient Flow: Law of Large Numbers



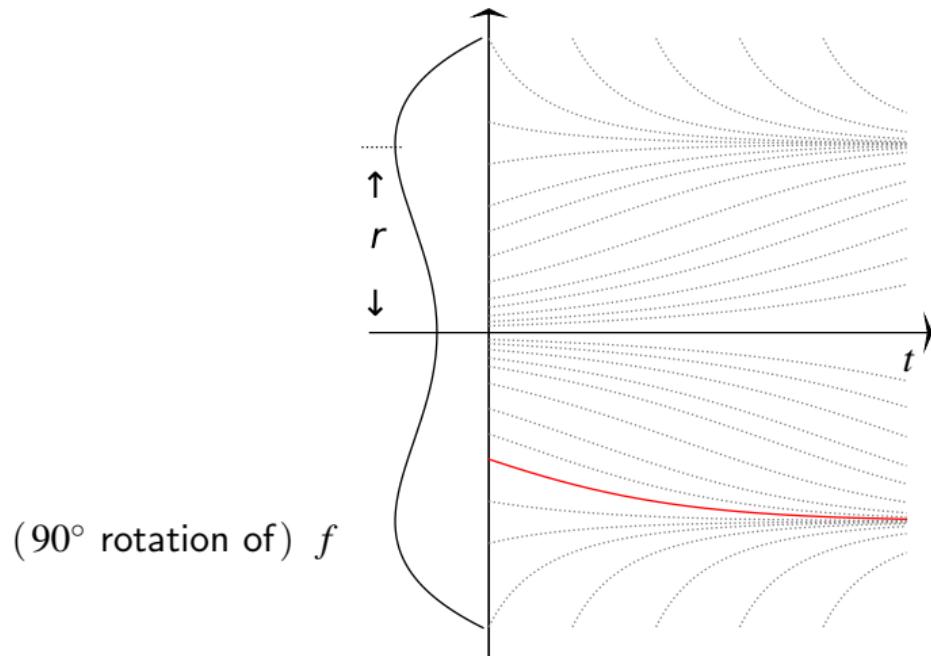
Gradient Flow: Law of Large Numbers



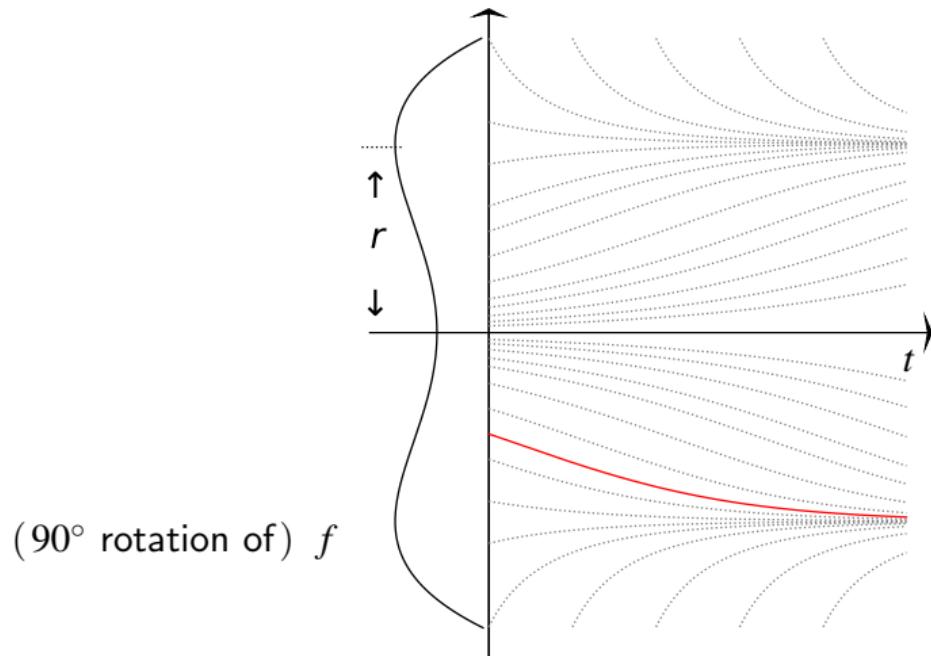
Gradient Flow: Law of Large Numbers



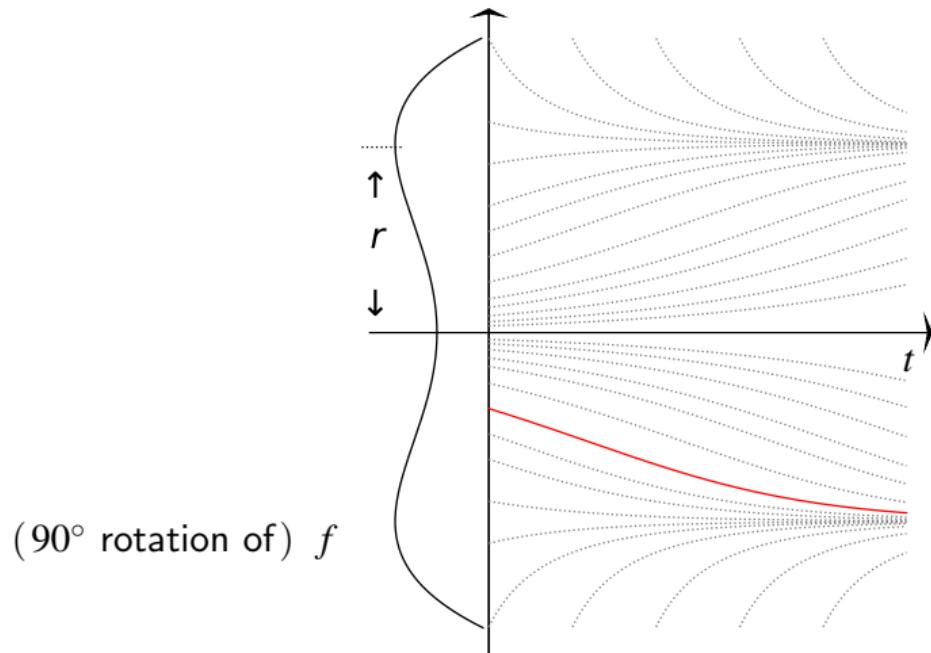
Gradient Flow: Law of Large Numbers



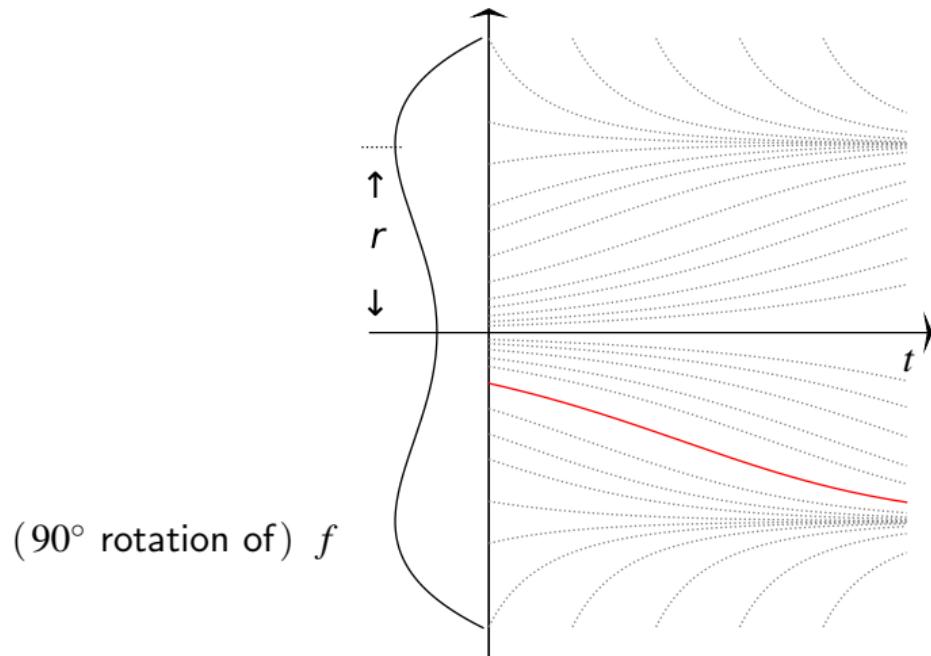
Gradient Flow: Law of Large Numbers



Gradient Flow: Law of Large Numbers



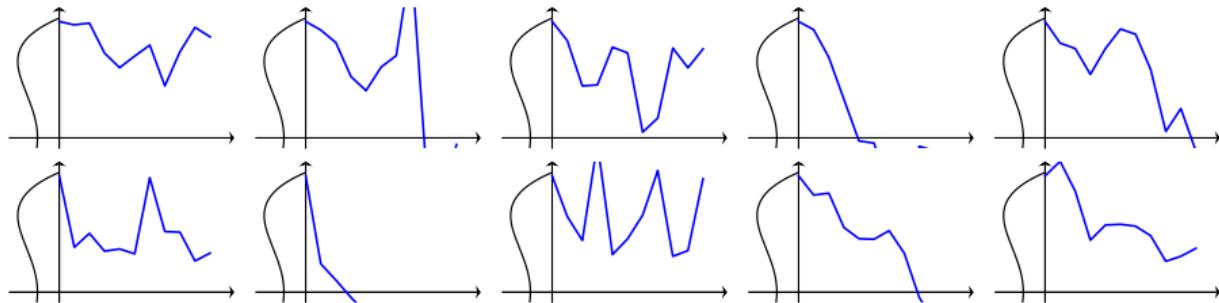
Gradient Flow: Law of Large Numbers



Typical Scenario

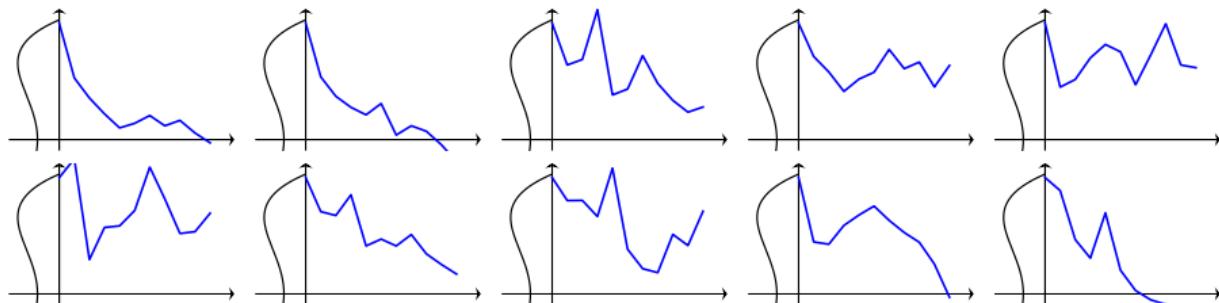
Trajectory of SGD W^η :

$\eta = 1/10$ & noises are **light-tailed**



Trajectory of SGD W^η :

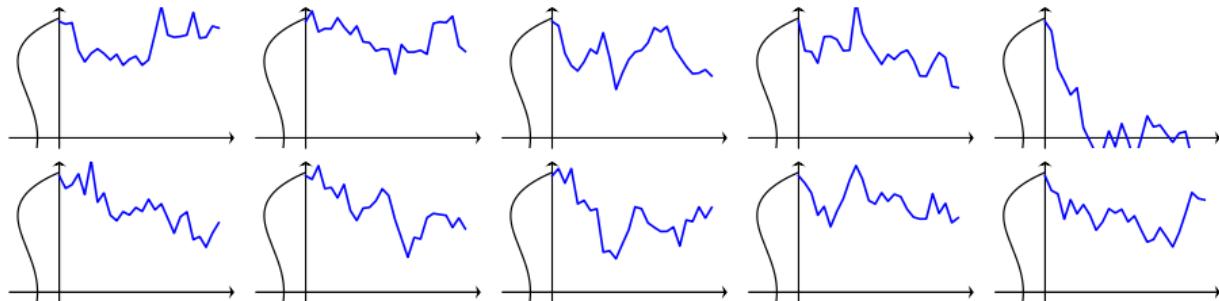
$\eta = 1/10$ & noises are **heavy-tailed**



Typical Scenario

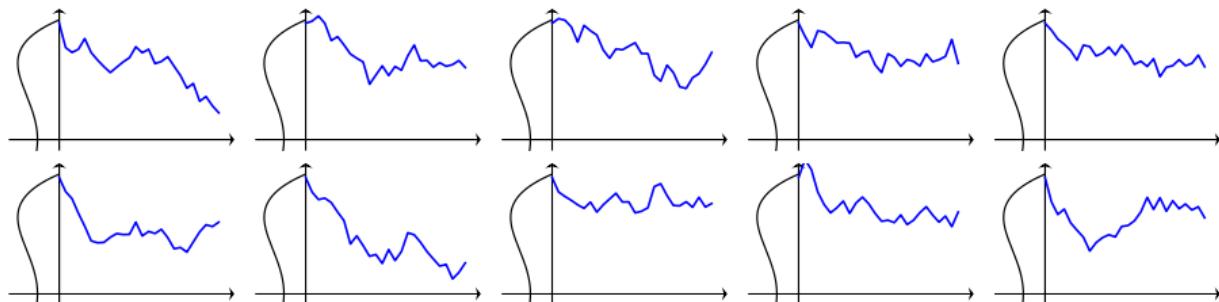
Trajectory of SGD W^η :

$\eta = 1/25$ & noises are **light-tailed**



Trajectory of SGD W^η :

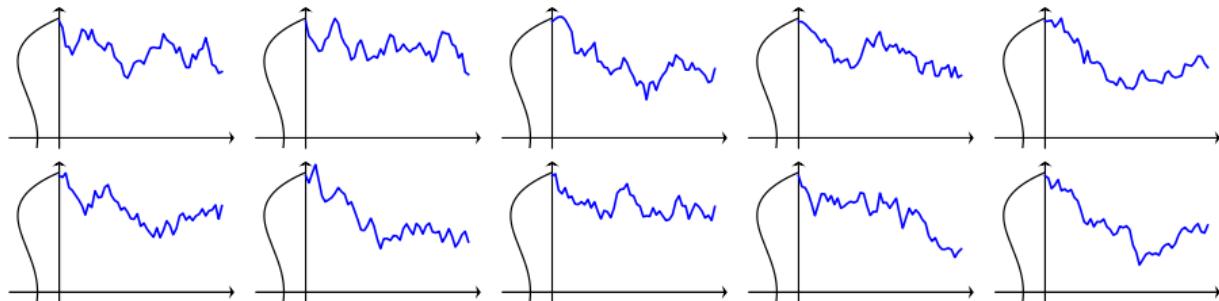
$\eta = 1/25$ & noises are **heavy-tailed**



Typical Scenario

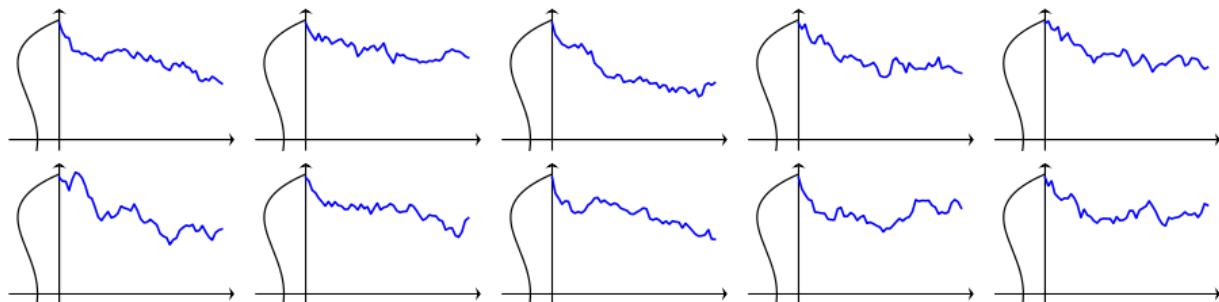
Trajectory of SGD W^η :

$\eta = 1/50$ & noises are **light-tailed**



Trajectory of SGD W^η :

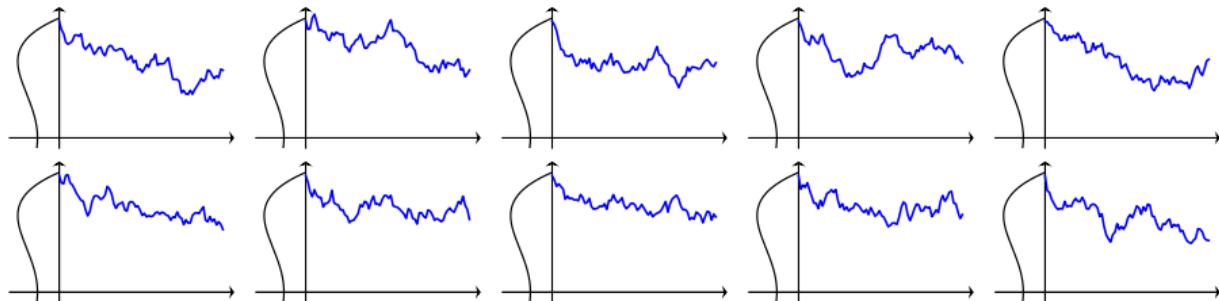
$\eta = 1/50$ & noises are **heavy-tailed**



Typical Scenario

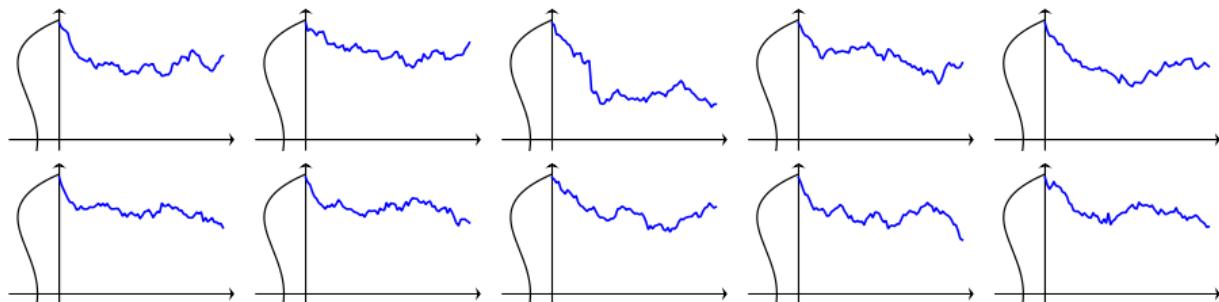
Trajectory of SGD W^η :

$\eta = 1/75$ & noises are **light-tailed**



Trajectory of SGD W^η :

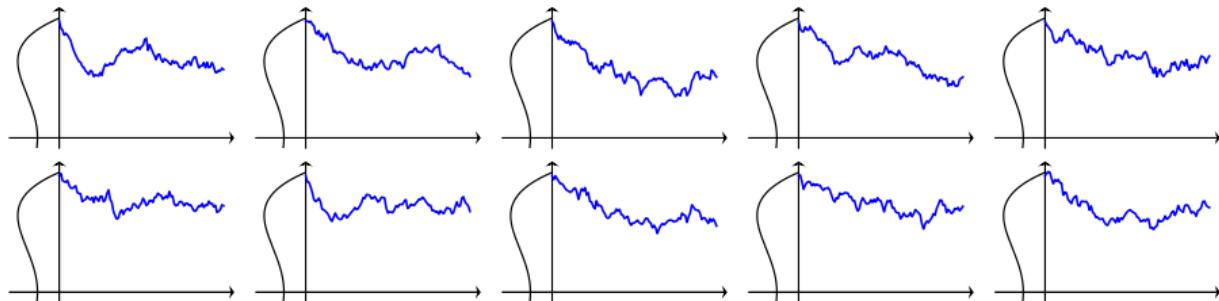
$\eta = 1/75$ & noises are **heavy-tailed**



Typical Scenario

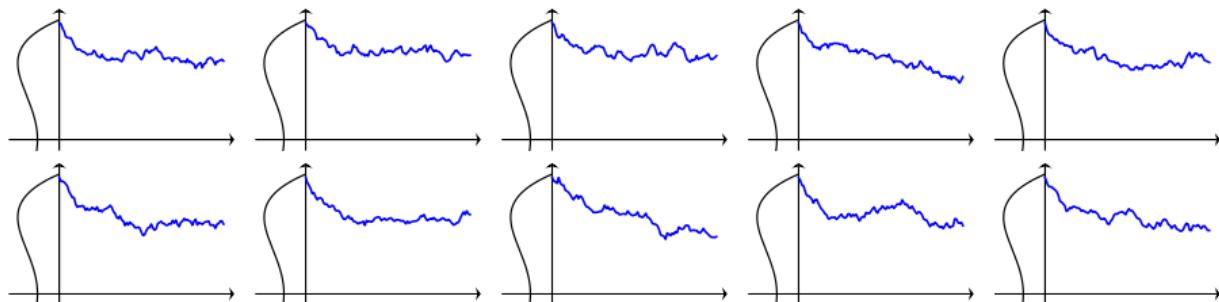
Trajectory of SGD W^η :

$\eta = 1/100$ & noises are **light-tailed**



Trajectory of SGD W^η :

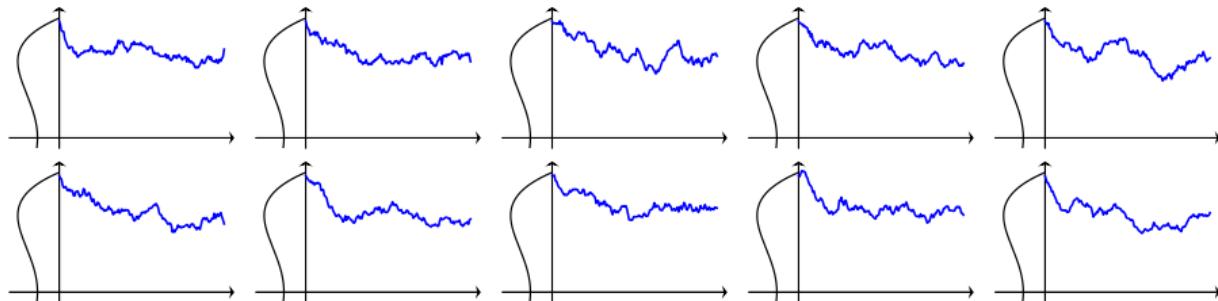
$\eta = 1/100$ & noises are **heavy-tailed**



Typical Scenario

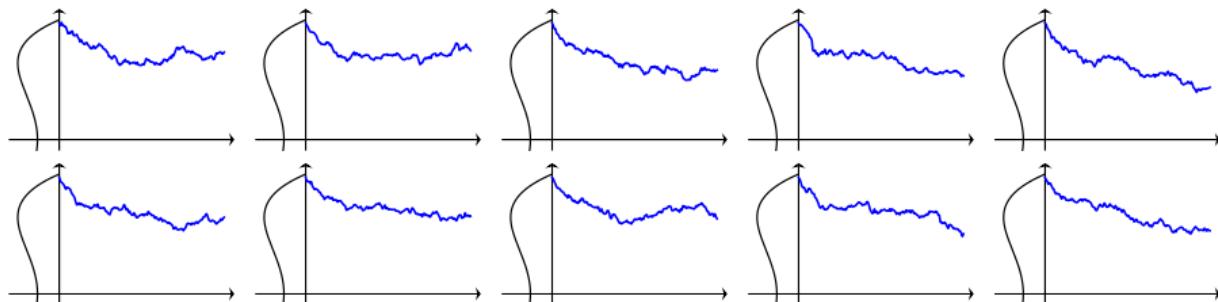
Trajectory of SGD W^η :

$\eta = 1/150$ & noises are **light-tailed**



Trajectory of SGD W^η :

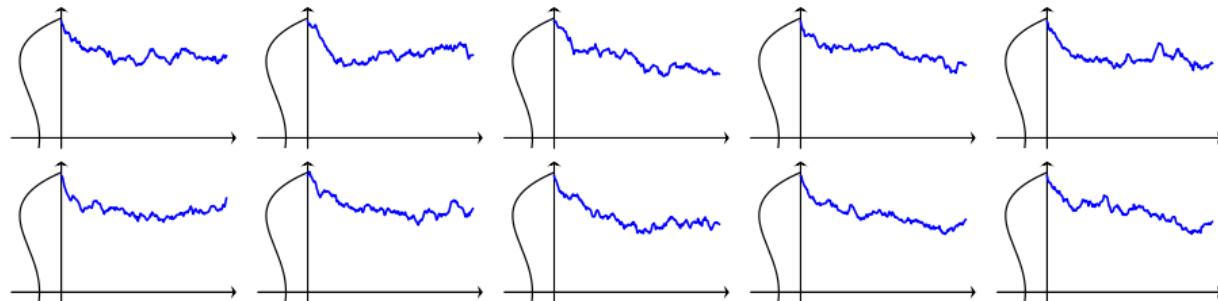
$\eta = 1/150$ & noises are **heavy-tailed**



Typical Scenario

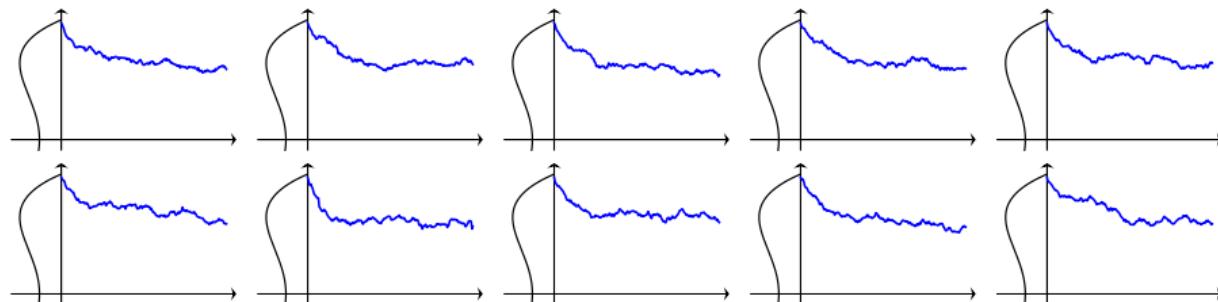
Trajectory of SGD W^η :

$\eta = 1/200$ & noises are **light-tailed**



Trajectory of SGD W^η :

$\eta = 1/200$ & noises are **heavy-tailed**



Heavy-Tailed Large Deviations for SGD

Theorem (Wang, Su, R., 2022+)

For “general” $A \subseteq \mathbb{D}$

$$C(A^\circ) \leq \liminf_{\eta \rightarrow 0} \frac{\mathbf{P}(W^\eta \in A)}{\eta^{\alpha \mathcal{J}(A)}} \leq \limsup_{\eta \rightarrow 0} \frac{\mathbf{P}(W^\eta \in A)}{\eta^{\alpha \mathcal{J}(A)}} \leq C(A^-).$$

- $\mathcal{J}(A)$: min #jumps added to $w(\cdot)$ for it to be inside A
- $C(\cdot)$: a measure

Heavy-Tailed Large Deviations for SGD

$$\mathbf{P}(|Z_i| \geq x) = x^{-(\alpha+1)}, \quad \alpha > 0$$

Theorem (Wang, Su, R., 2022+)

For “general” $A \subseteq \mathbb{D}$

$$C(A^\circ) \leq \liminf_{\eta \rightarrow 0} \frac{\mathbf{P}(W^\eta \in A)}{\eta^\alpha \mathcal{J}(A)} \leq \limsup_{\eta \rightarrow 0} \frac{\mathbf{P}(W^\eta \in A)}{\eta^\alpha \mathcal{J}(A)} \leq C(A^-).$$

- $\mathcal{J}(A)$: min #jumps added to $w(\cdot)$ for it to be inside A
- $C(\cdot)$: a measure

Heavy-Tailed Large Deviations for SGD

$$\mathbf{P}(|Z_i| \geq x) = x^{-(\alpha+1)}, \quad \alpha > 0$$

gradient noise

Theorem (Wang, Su, R., 2022+)

For “general” $A \subseteq \mathbb{D}$

$$C(A^\circ) \leq \liminf_{\eta \rightarrow 0} \frac{\mathbf{P}(W^\eta \in A)}{\eta^\alpha \mathcal{J}(A)} \leq \limsup_{\eta \rightarrow 0} \frac{\mathbf{P}(W^\eta \in A)}{\eta^\alpha \mathcal{J}(A)} \leq C(A^-).$$

- $\mathcal{J}(A)$: min #jumps added to $w(\cdot)$ for it to be inside A
- $C(\cdot)$: a measure

Heavy-Tailed Large Deviations for SGD

$$\mathbf{P}(|Z_i| \geq x) = x^{-(\alpha+1)}, \quad \alpha > 0$$

↙
gradient noise

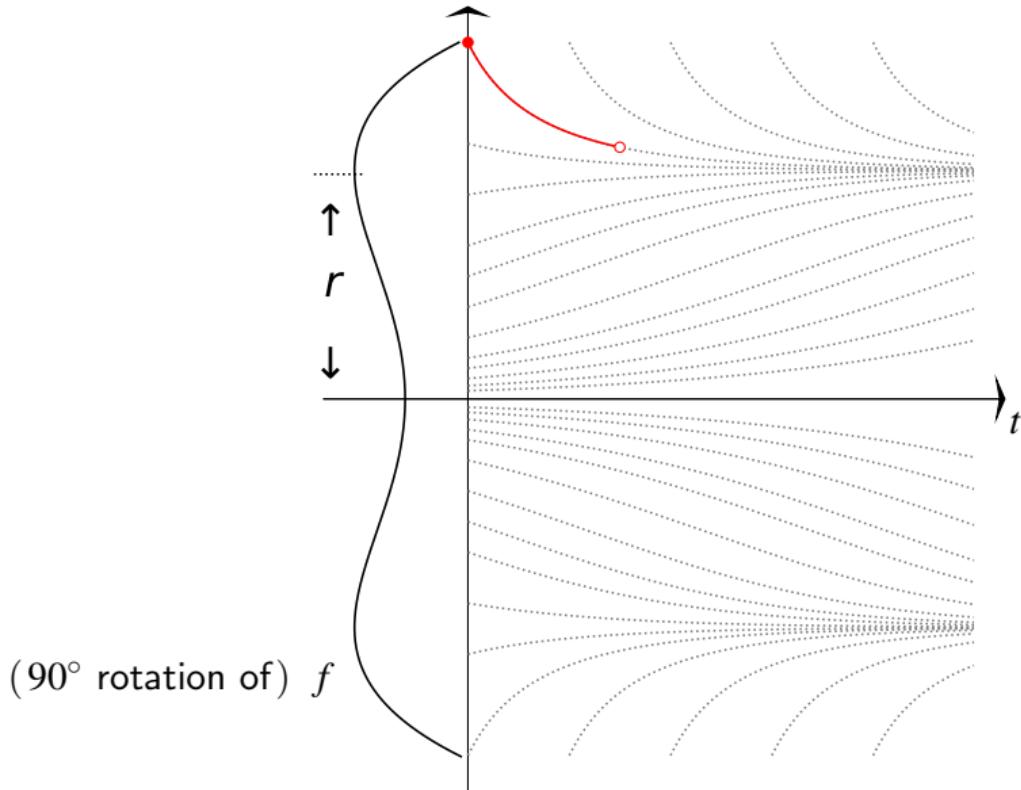
Theorem (Wang, Su, R., 2022+)

For “general” $A \subseteq \mathbb{D}$

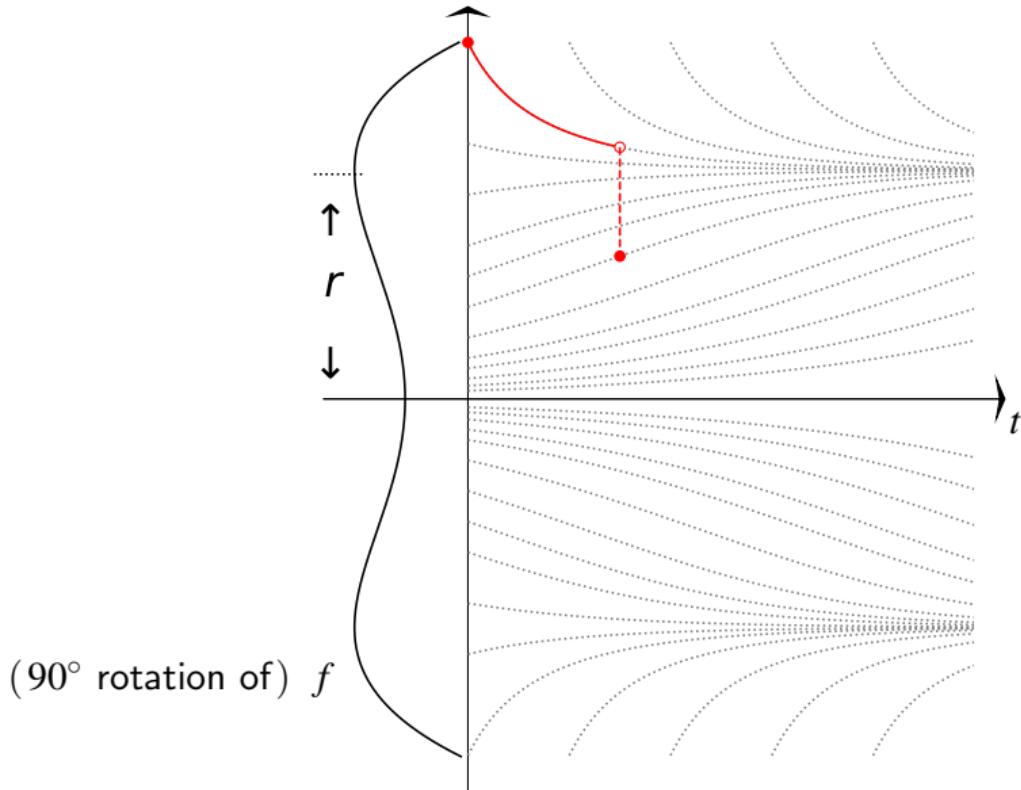
$$\mathbf{P}(W^\eta \in A) \sim n^{-\alpha \mathcal{J}(A)}$$

- $\mathcal{J}(A)$: min #jumps added to $w(\cdot)$ for it to be inside A

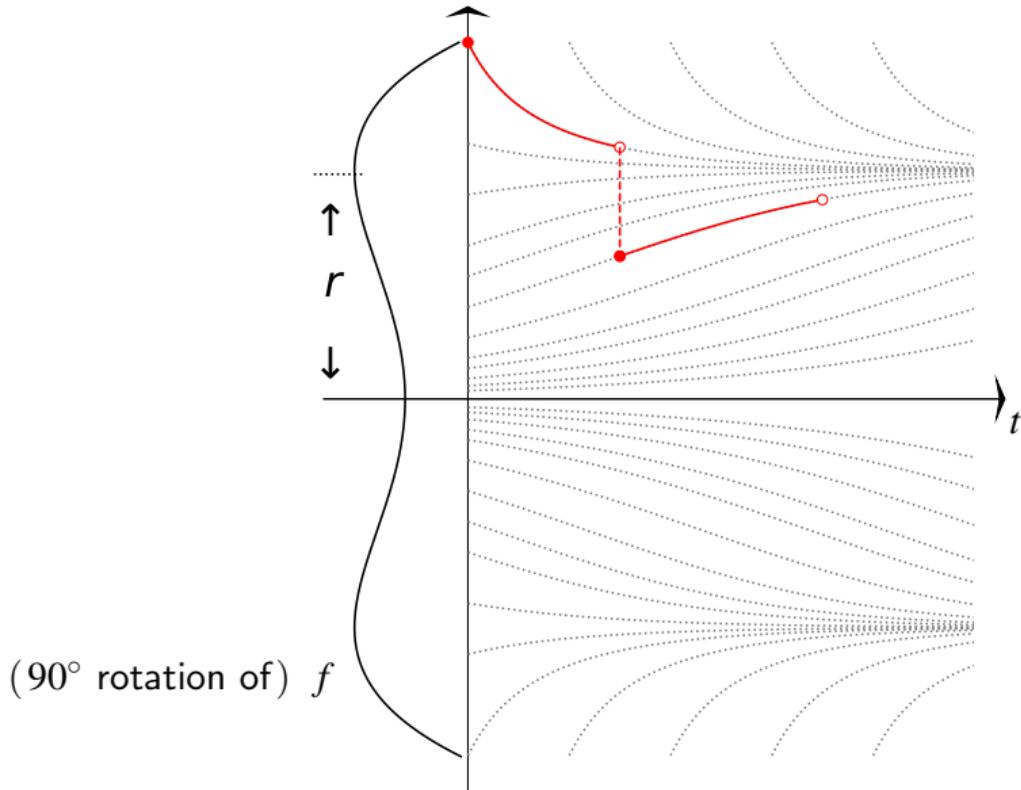
Adding Jumps to $w(\cdot)$: Piecewise Gradient Flow



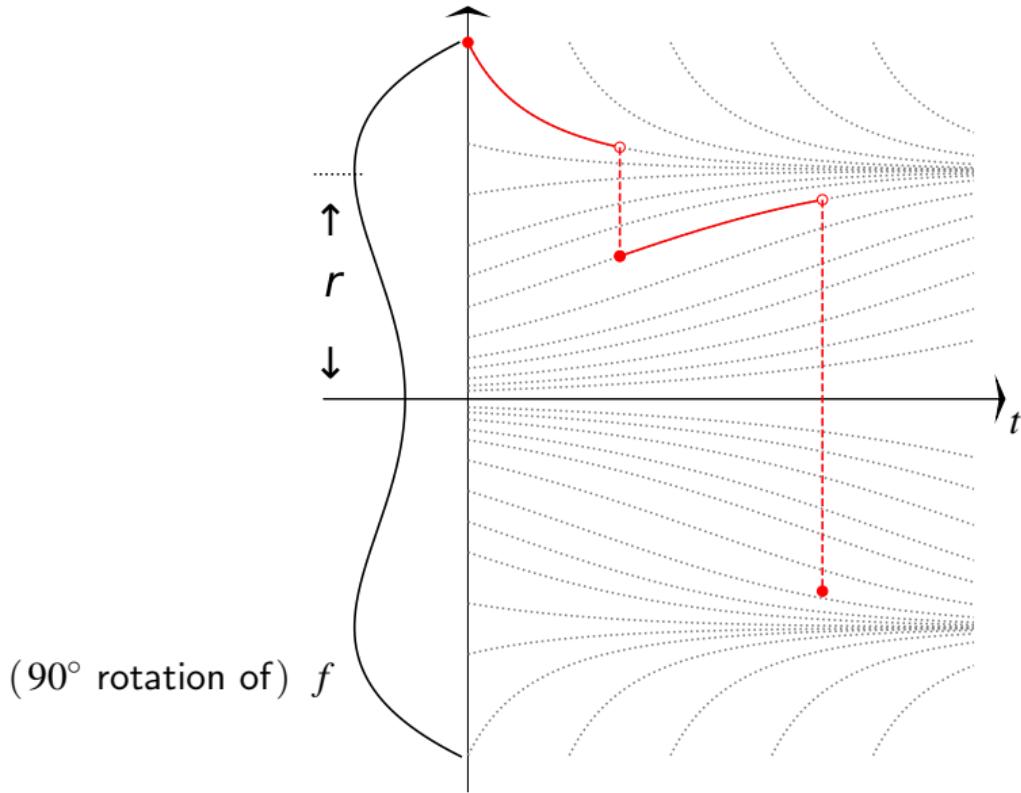
Adding Jumps to $w(\cdot)$: Piecewise Gradient Flow



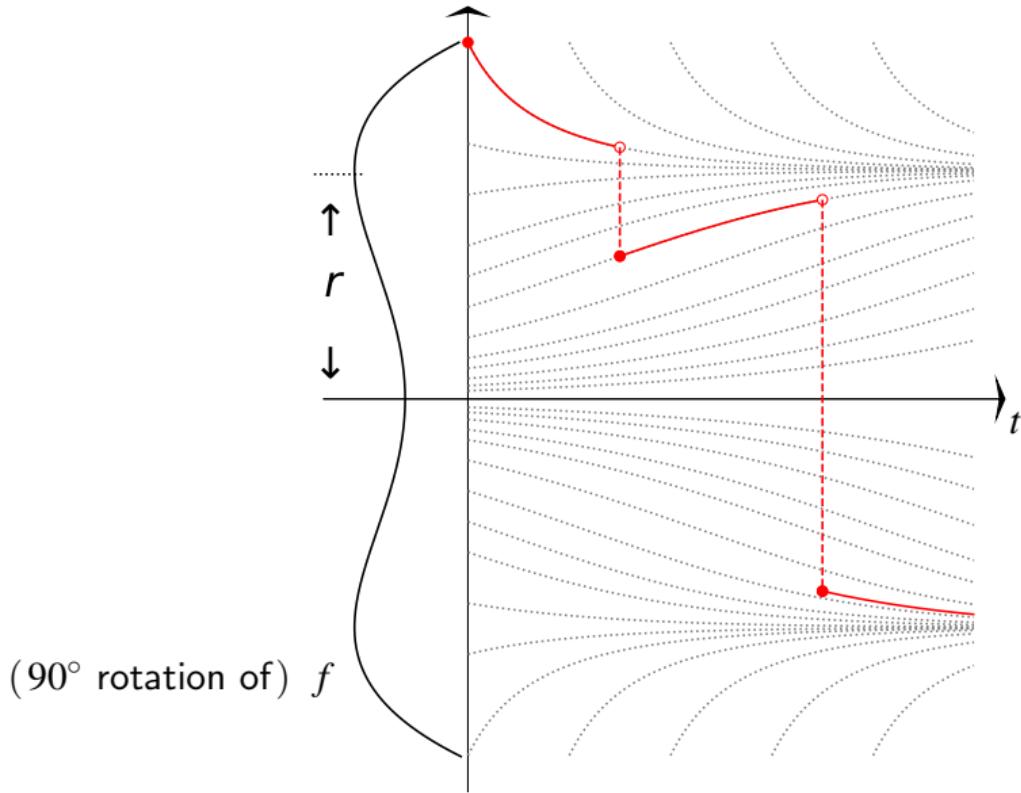
Adding Jumps to $w(\cdot)$: Piecewise Gradient Flow



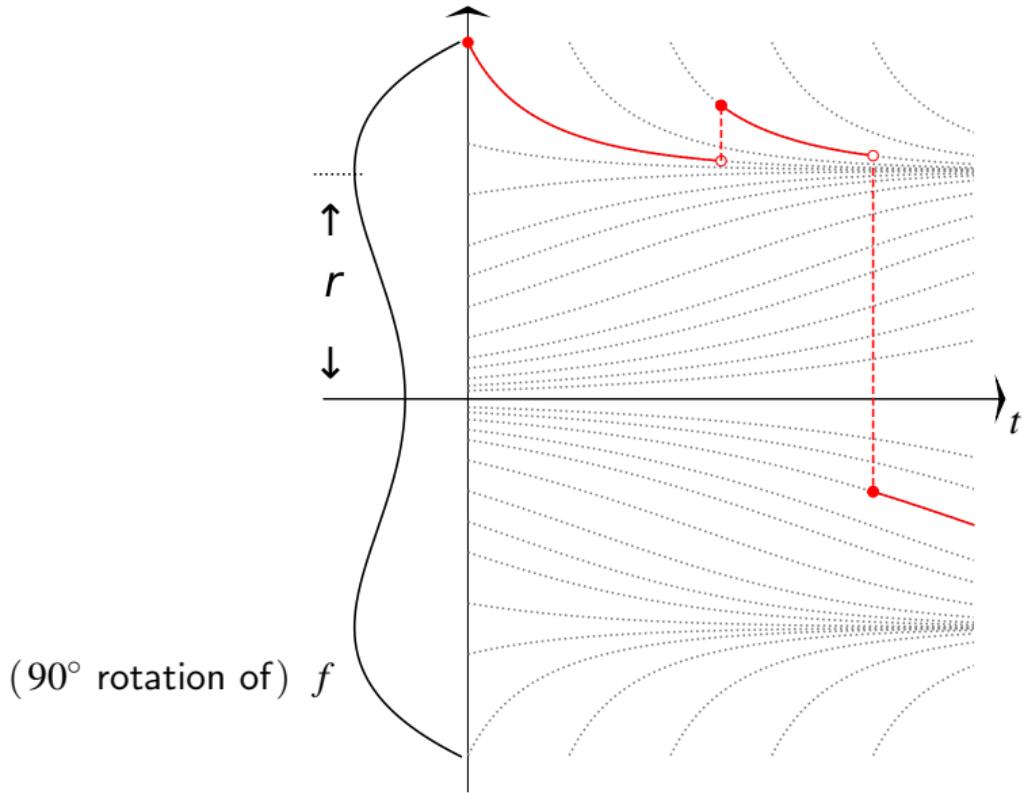
Adding Jumps to $w(\cdot)$: Piecewise Gradient Flow



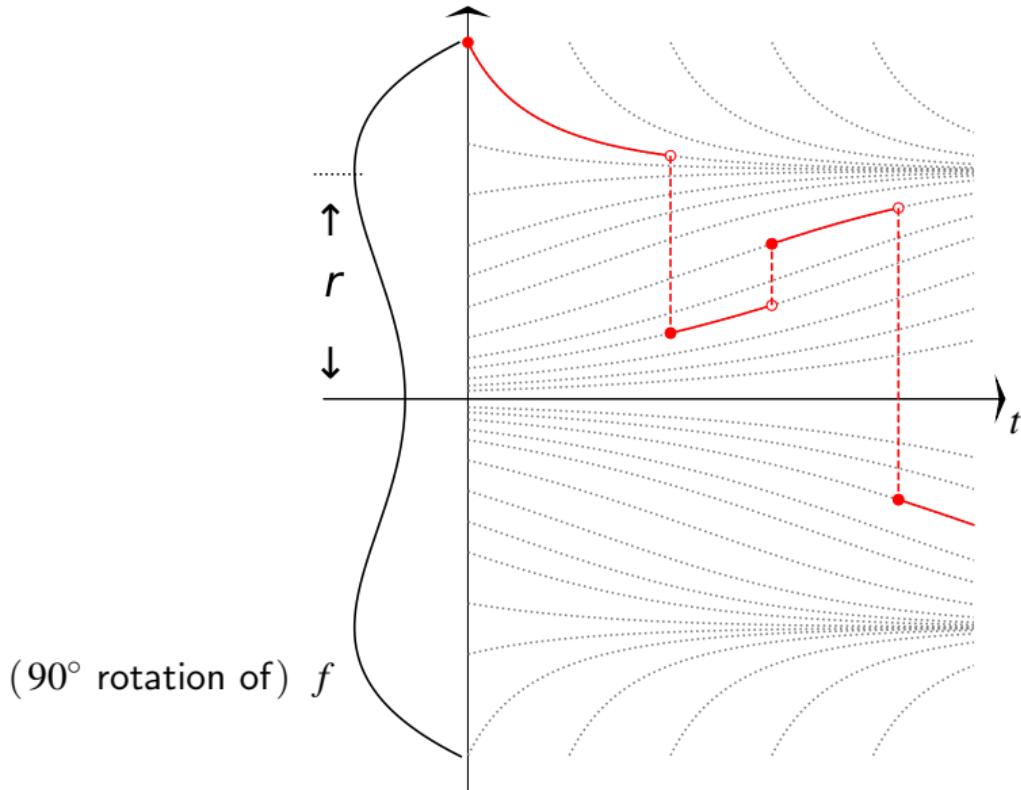
Adding Jumps to $w(\cdot)$: Piecewise Gradient Flow



Adding Jumps to $w(\cdot)$: Piecewise Gradient Flow



Adding Jumps to $w(\cdot)$: Piecewise Gradient Flow



Recall: Heavy-Tailed Large Deviations for SGD

Theorem (Wang, Su, R., 2022+)

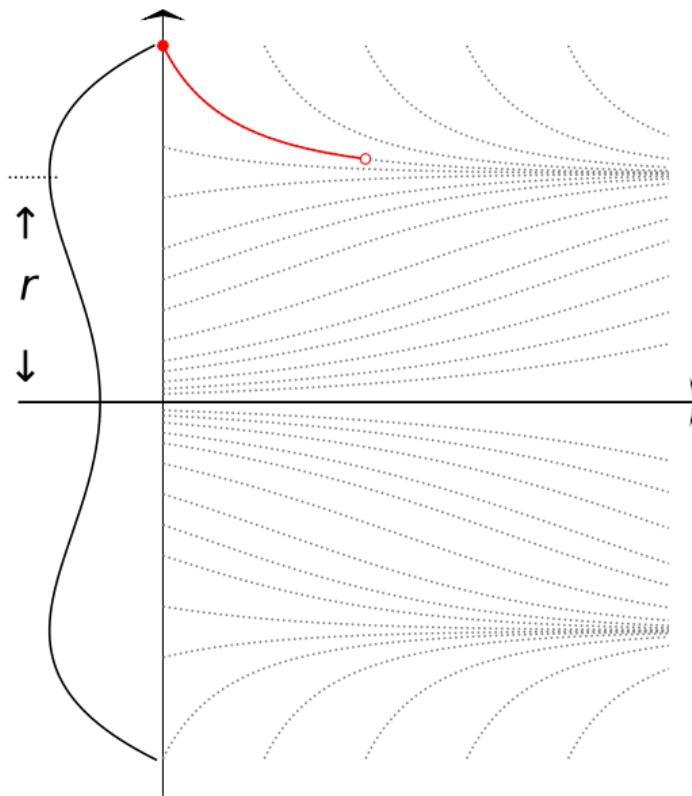
For “general” $A \subseteq \mathbb{D}$

$$\mathbf{P}(W^\eta \in A) \sim n^{-\alpha \mathcal{J}(A)}$$

- $\mathcal{J}(A)$: min #jumps added to $w(\cdot)$ for it to be inside A

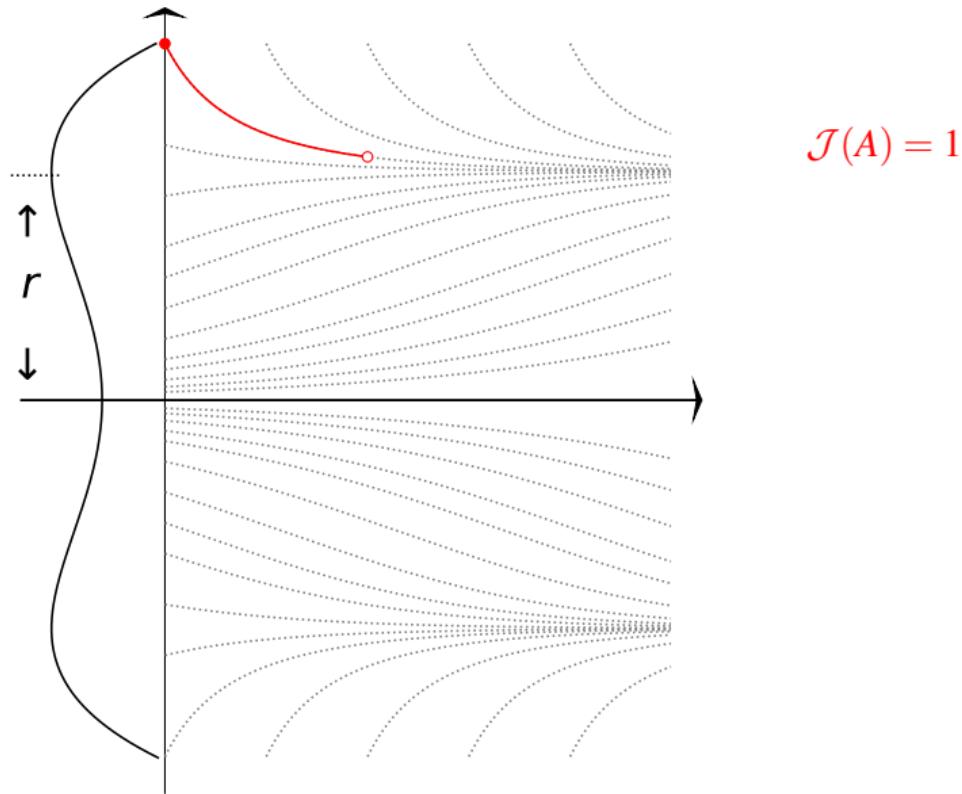
Catastrophe Principle: Most Likely Escape Route

Most Likely Paths for $\{W^\eta \text{ escapes the local minimum}\}$



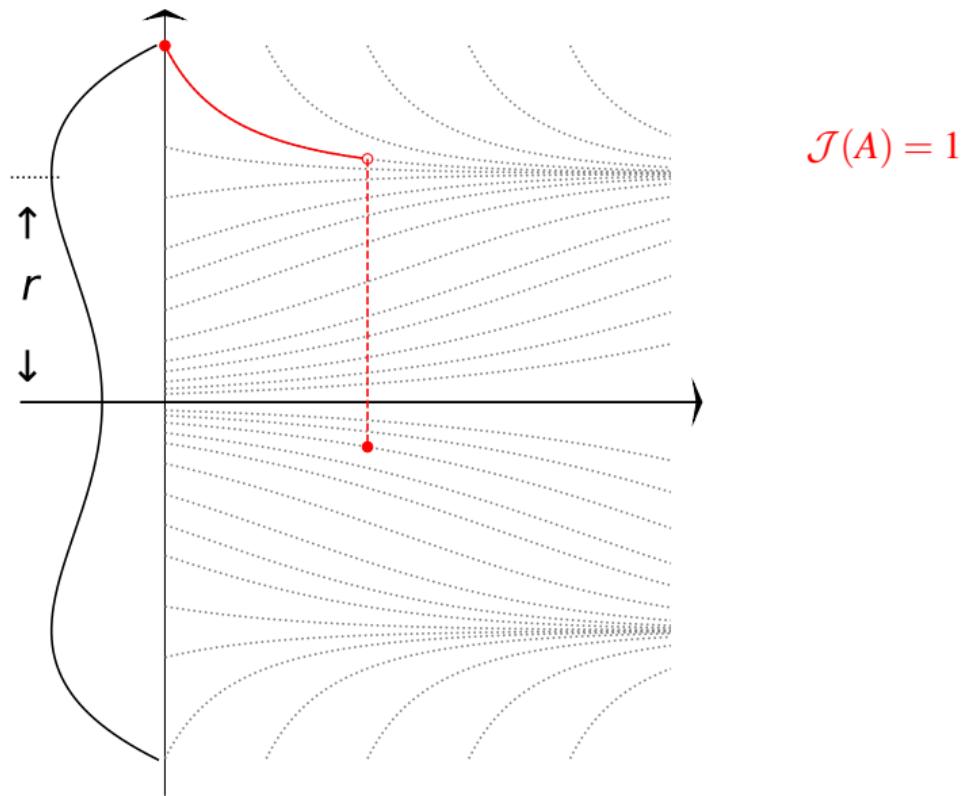
Catastrophe Principle: Most Likely Escape Route

Most Likely Paths for $\{W^\eta \text{ escapes the local minimum}\}$



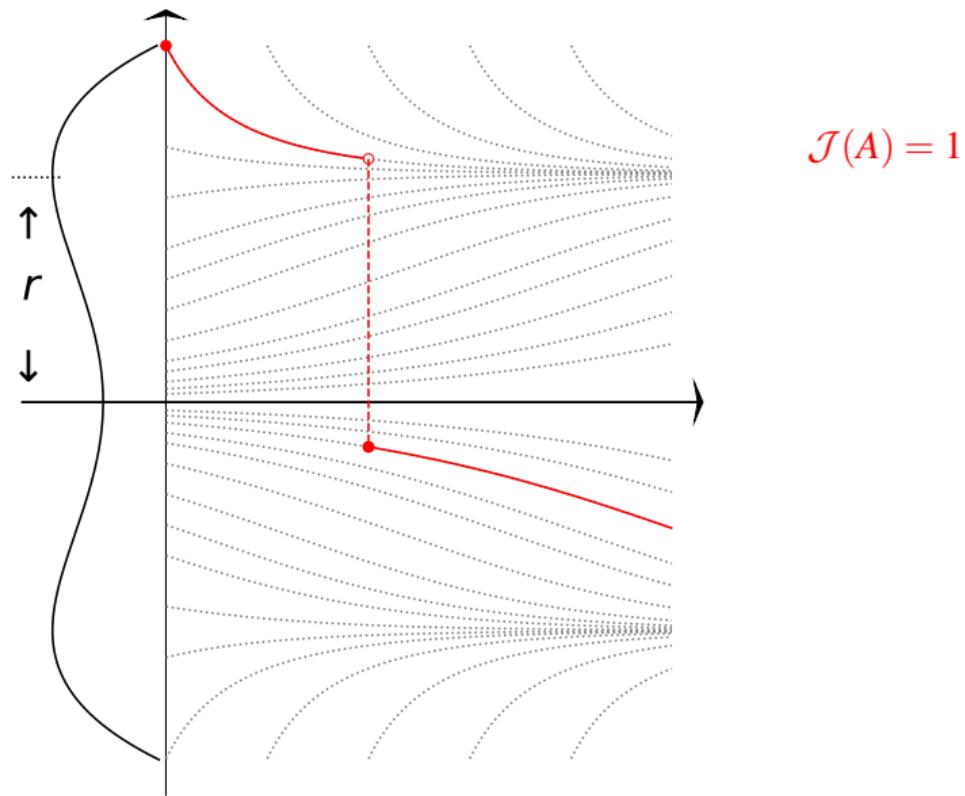
Catastrophe Principle: Most Likely Escape Route

Most Likely Paths for $\{W^\eta \text{ escapes the local minimum}\}$



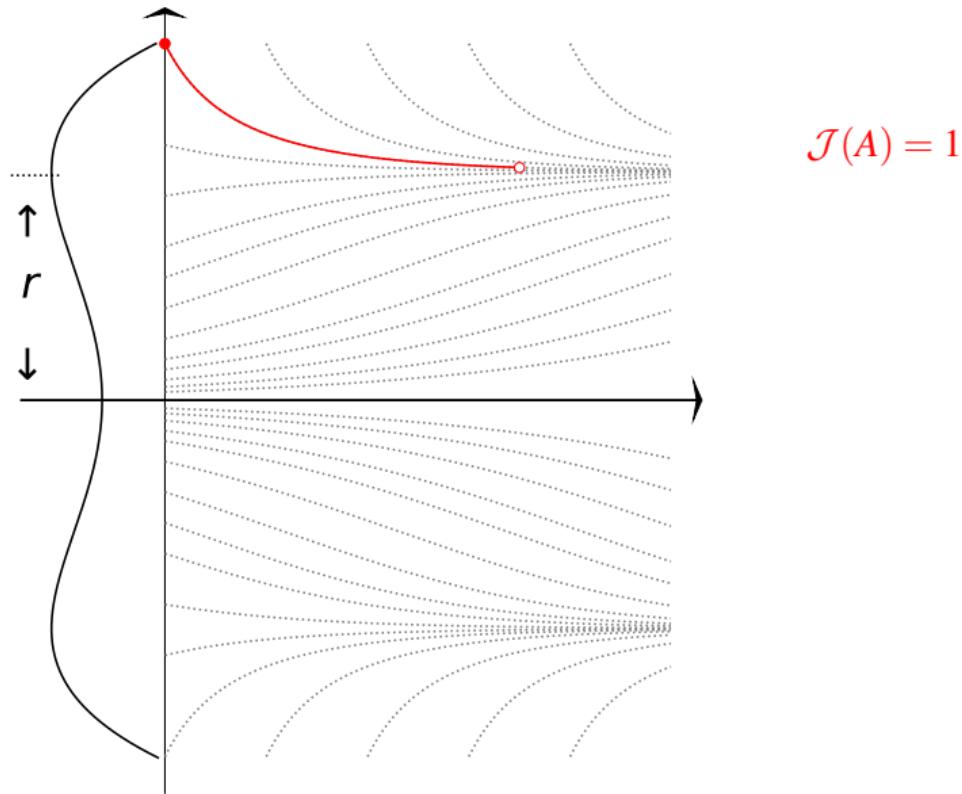
Catastrophe Principle: Most Likely Escape Route

Most Likely Paths for $\{W^\eta \text{ escapes the local minimum}\}$



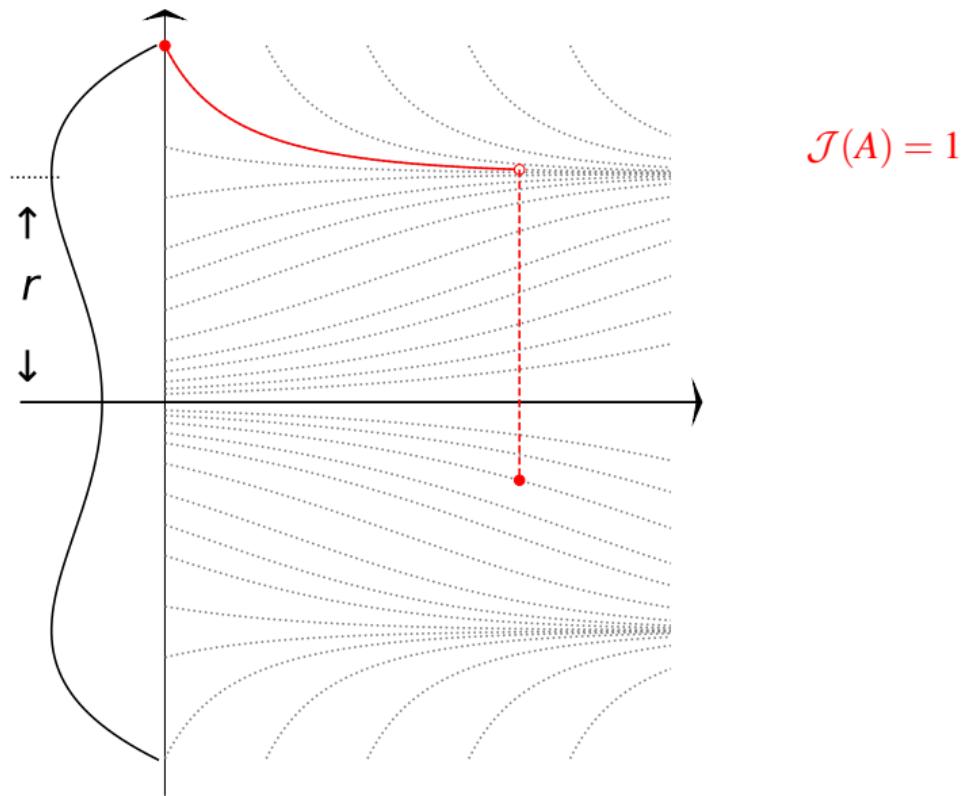
Catastrophe Principle: Most Likely Escape Route

Most Likely Paths for $\{W^\eta \text{ escapes the local minimum}\}$



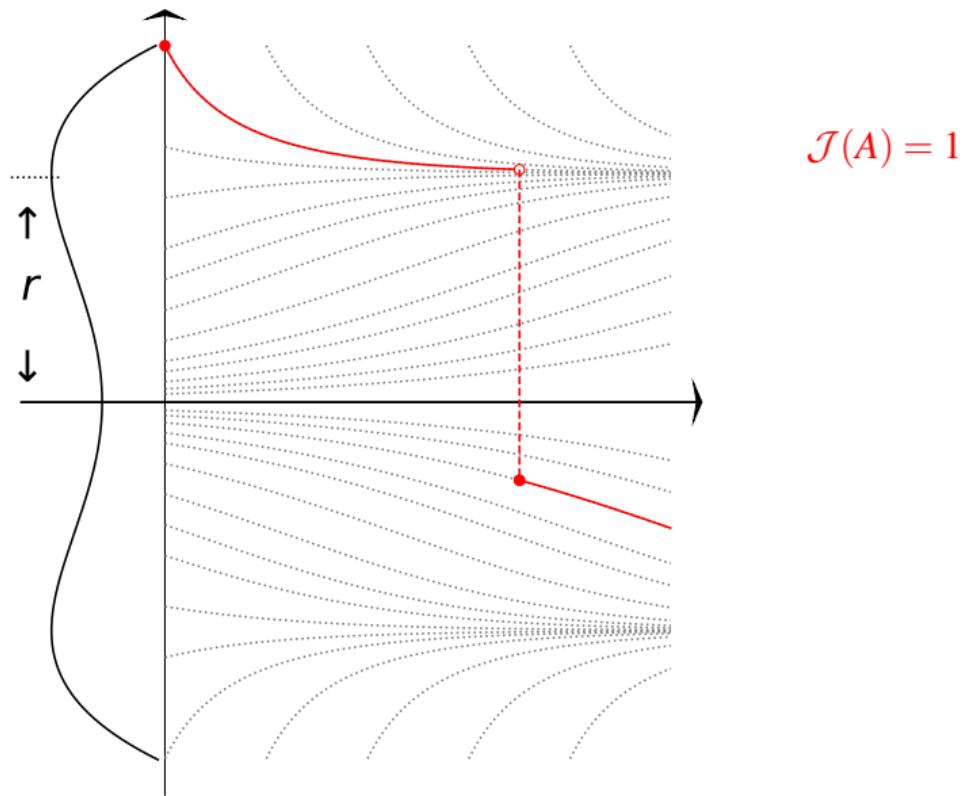
Catastrophe Principle: Most Likely Escape Route

Most Likely Paths for $\{W^\eta \text{ escapes the local minimum}\}$



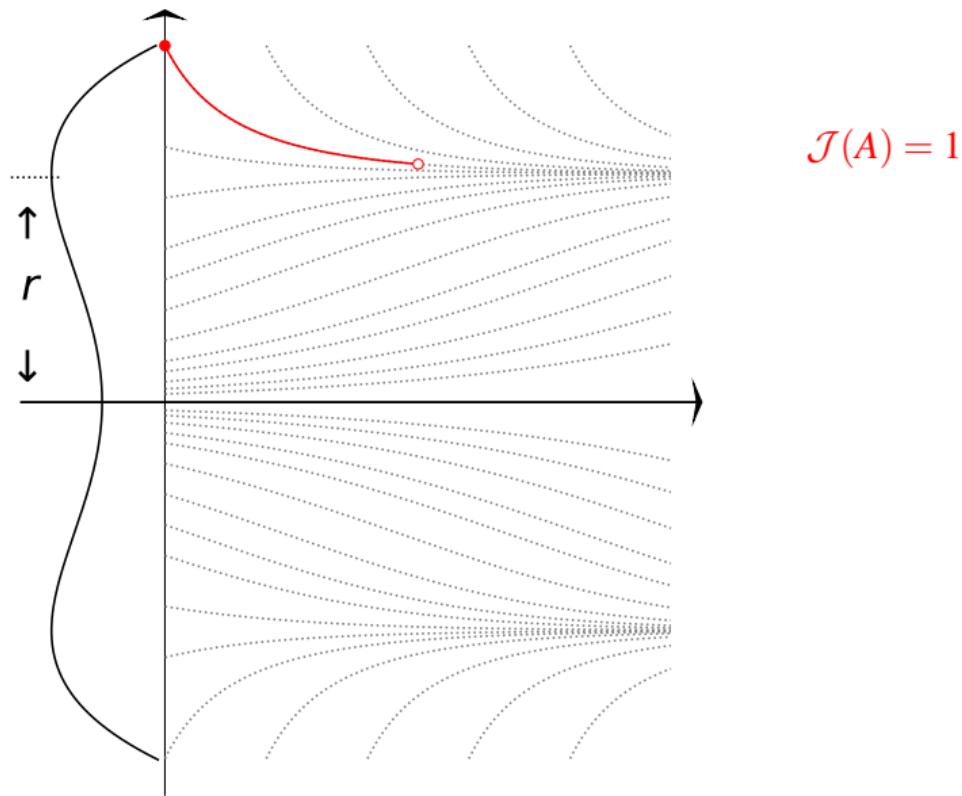
Catastrophe Principle: Most Likely Escape Route

Most Likely Paths for $\{W^\eta \text{ escapes the local minimum}\}$



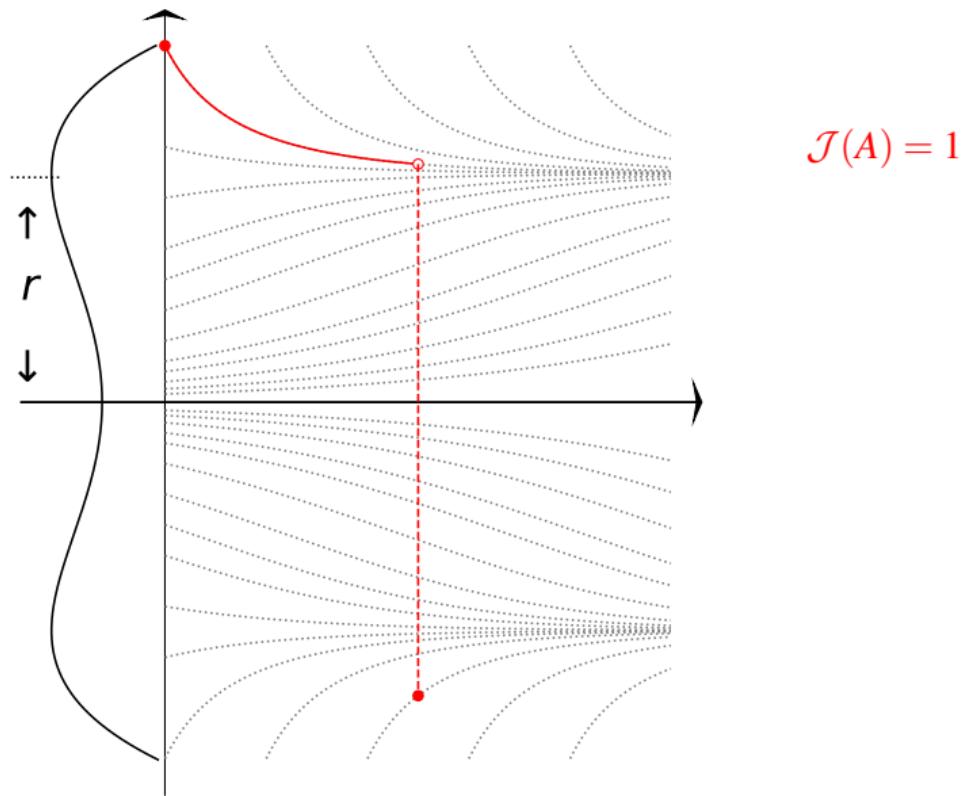
Catastrophe Principle: Most Likely Escape Route

Most Likely Paths for $\{W^\eta \text{ escapes the local minimum}\}$



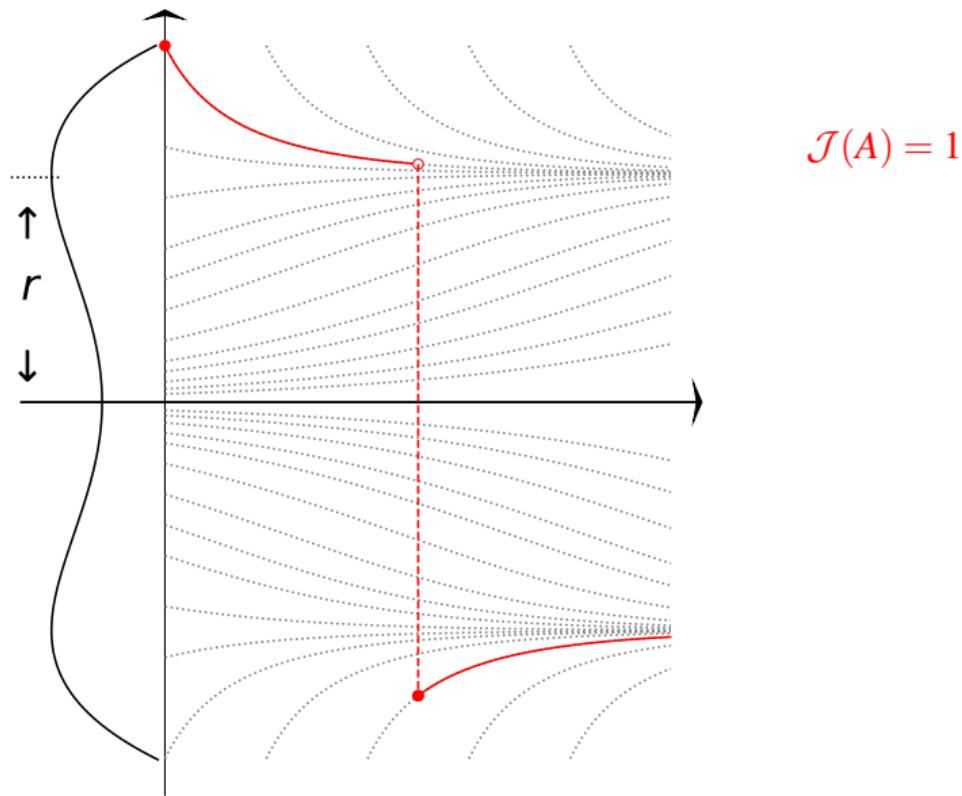
Catastrophe Principle: Most Likely Escape Route

Most Likely Paths for $\{W^\eta \text{ escapes the local minimum}\}$



Catastrophe Principle: Most Likely Escape Route

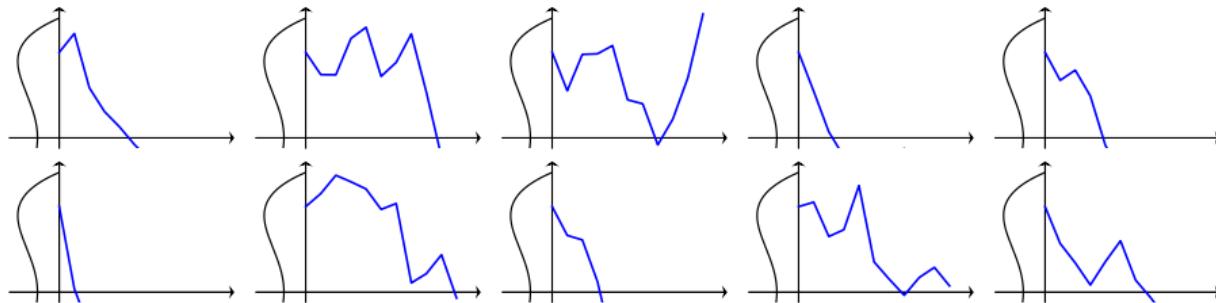
Most Likely Paths for $\{W^\eta \text{ escapes the local minimum}\}$



Catastrophe Principle Dictates SGD's Escape Route

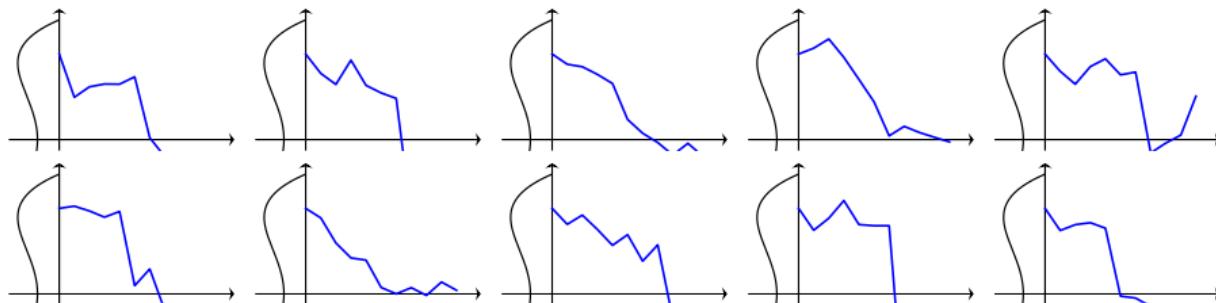
Trajectory of SGD X^η conditional on exit:

light-tailed noises with $\eta = 1/10$



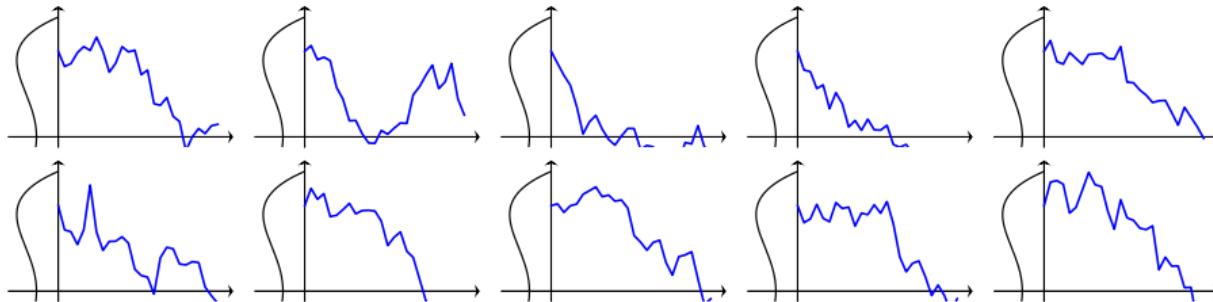
Trajectory of SGD X^η conditional on exit:

heavy-tailed noises with $\eta = 1/10$



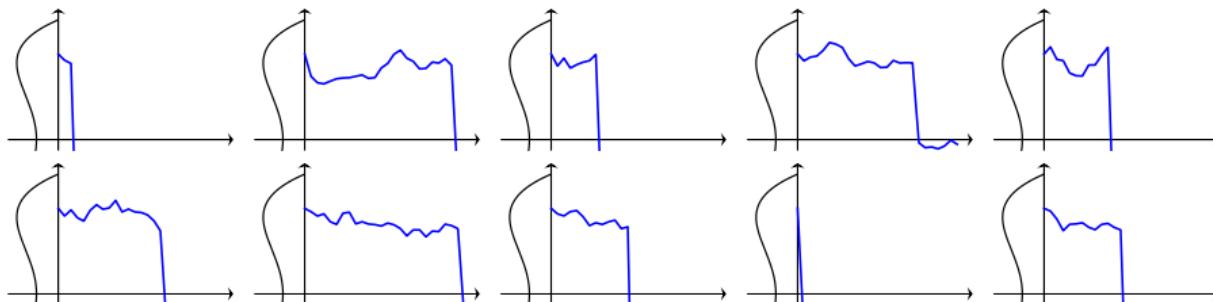
Catastrophe Principle Dictates SGD's Escape Route

Trajectory of SGD X^η conditional on exit:



light-tailed noises with $\eta = 1/25$

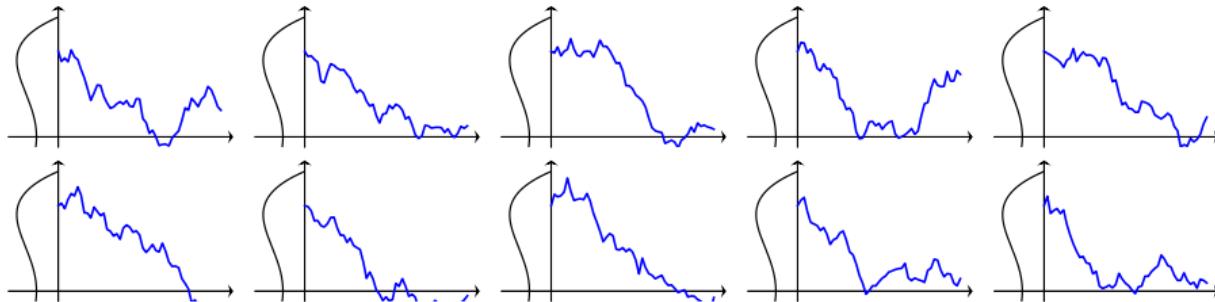
Trajectory of SGD X^η conditional on exit:



heavy-tailed noises with $\eta = 1/25$

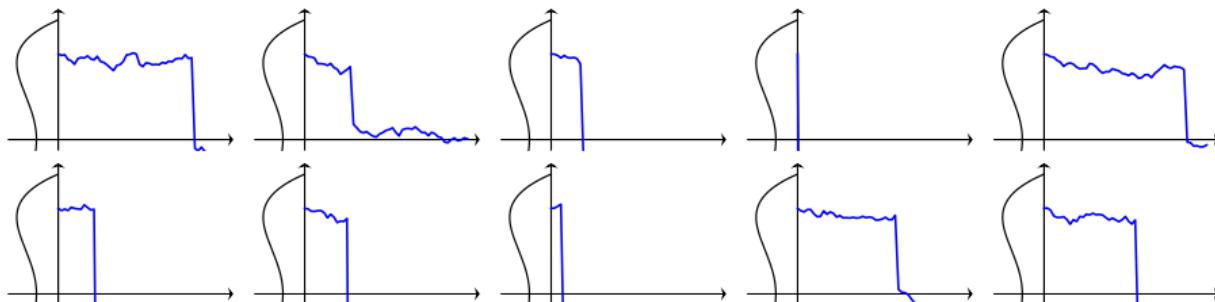
Catastrophe Principle Dictates SGD's Escape Route

Trajectory of SGD X^η conditional on exit:



light-tailed noises with $\eta = 1/50$

Trajectory of SGD X^η conditional on exit:

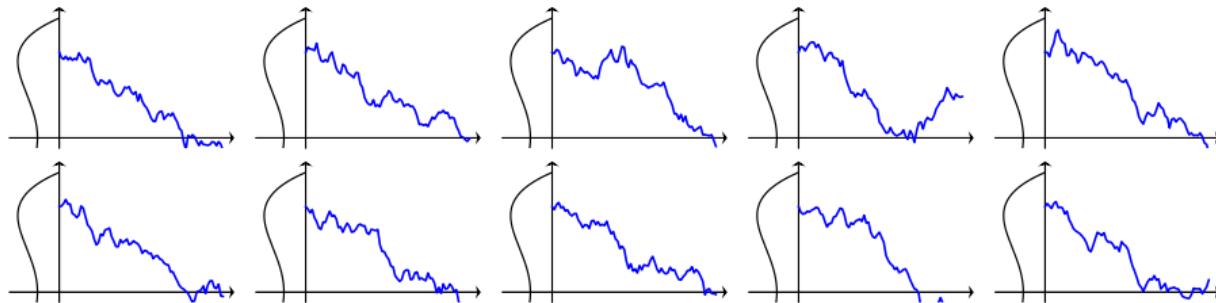


heavy-tailed noises with $\eta = 1/50$

Catastrophe Principle Dictates SGD's Escape Route

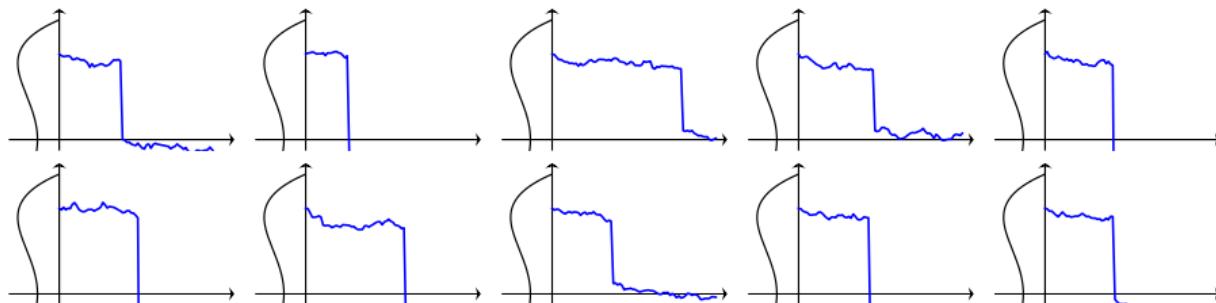
Trajectory of SGD X^η conditional on exit:

light-tailed noises with $\eta = 1/75$



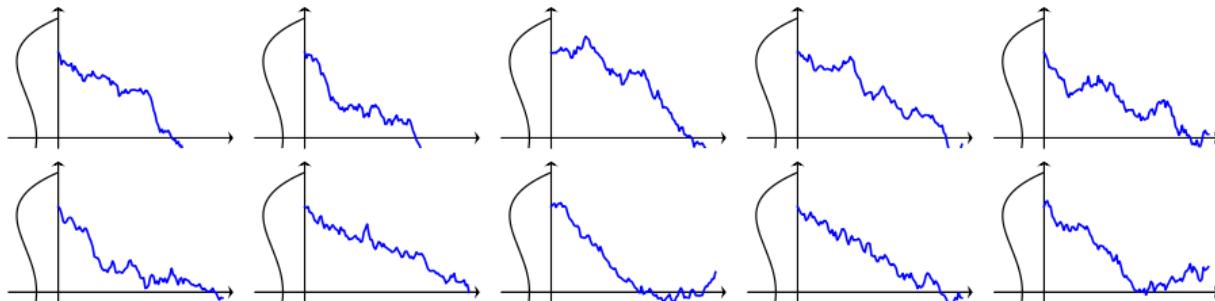
Trajectory of SGD X^η conditional on exit:

heavy-tailed noises with $\eta = 1/75$

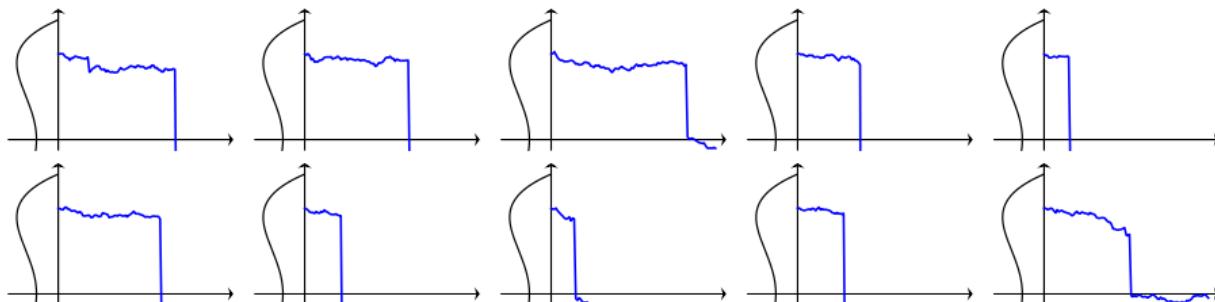


Catastrophe Principle Dictates SGD's Escape Route

Trajectory of SGD X^η conditional on exit: **light-tailed** noises with $\eta = 1/100$

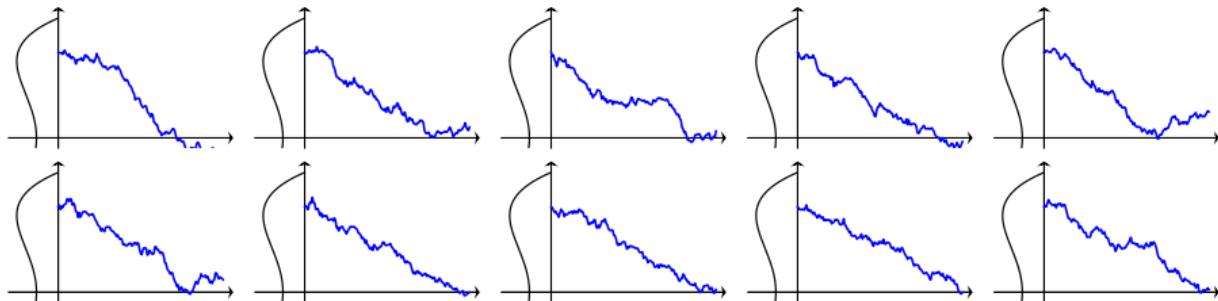


Trajectory of SGD X^η conditional on exit: **heavy-tailed** noises with $\eta = 1/100$

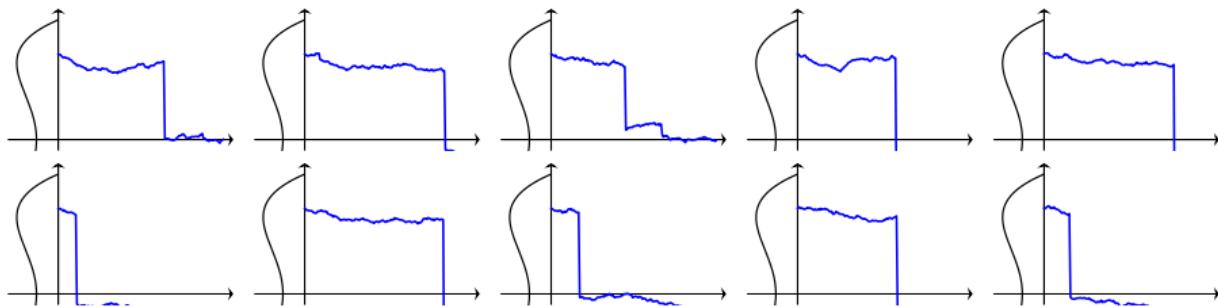


Catastrophe Principle Dictates SGD's Escape Route

Trajectory of SGD X^η conditional on exit: **light-tailed** noises with $\eta = 1/150$

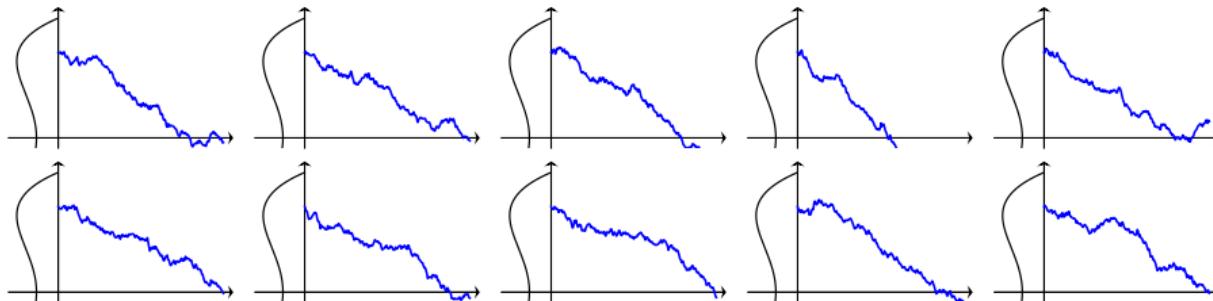


Trajectory of SGD X^η conditional on exit: **heavy-tailed** noises with $\eta = 1/150$

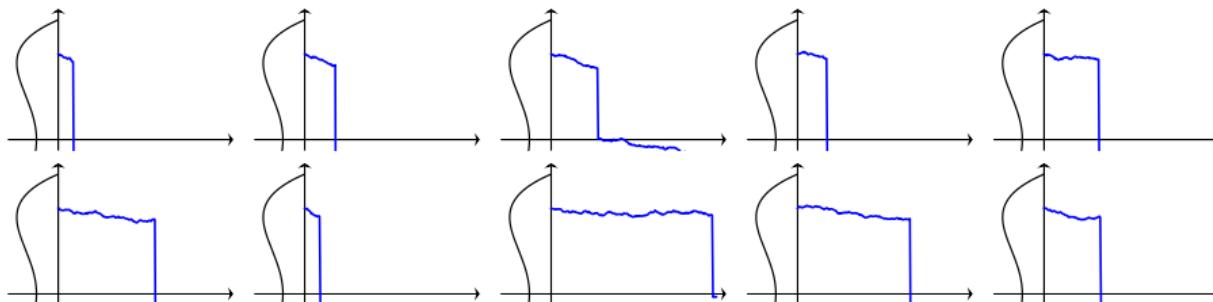


Catastrophe Principle Dictates SGD's Escape Route

Trajectory of SGD X^η conditional on exit: **light-tailed** noises with $\eta = 1/200$



Trajectory of SGD X^η conditional on exit: **heavy-tailed** noises with $\eta = 1/200$



Truncated Version of Stochastic Gradient Descent

SGD

$$W_{k+1}^{\eta} = W_k^{\eta} - \eta (f'(W_k^{\eta}) + Z_k) \quad k = 0, 1, 2, \dots$$

Truncated Version of Stochastic Gradient Descent

SGD with Gradient Clipping

$$W_{k+1}^{\eta} = W_k^{\eta} - \varphi_c(\eta(f'(W_k^{\eta}) + Z_k)) \quad k = 0, 1, 2, \dots$$

Truncated Version of Stochastic Gradient Descent

SGD with Gradient Clipping

$$W_{k+1}^\eta = W_k^\eta - \varphi_c(\eta(f'(W_k^\eta) + Z_k)) \quad k = 0, 1, 2, \dots$$

where

$$\varphi_c(x) = \frac{x}{|x|} \min\{c, |x|\}.$$

Truncated Version of Stochastic Gradient Descent

SGD with Gradient Clipping

$$W_{k+1}^\eta = W_k^\eta - \varphi_c(\eta(f'(W_k^\eta) + Z_k)) \quad k = 0, 1, 2, \dots$$

where

$$\varphi_c(x) = \frac{x}{|x|} \min\{c, |x|\}.$$

Then, again,

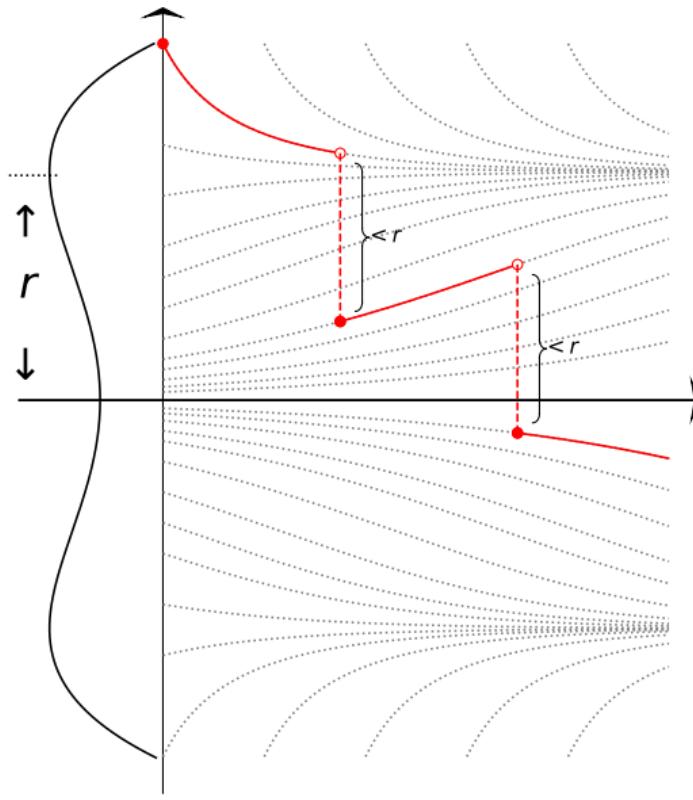
$$W^\eta(\cdot) \rightarrow w(\cdot) \quad \text{as} \quad \eta \rightarrow 0$$

where

$$dw(t) = -f'(w(t))dt.$$

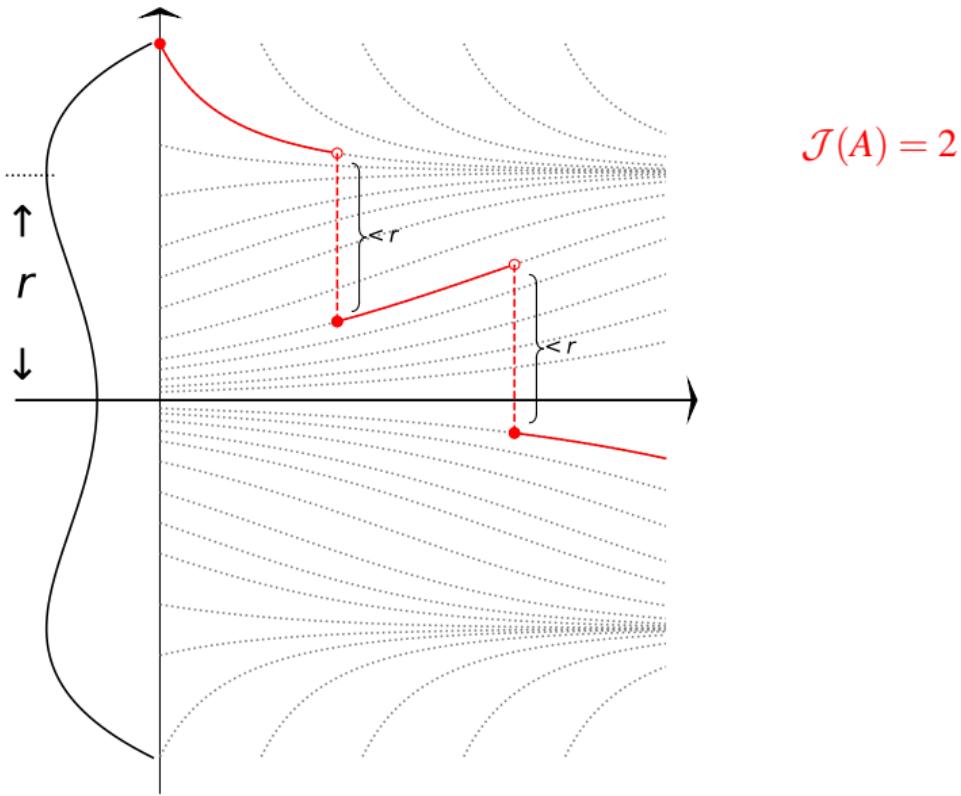
How does $\mathcal{J}(A)$ change?

If $r \in (c, 2c)$



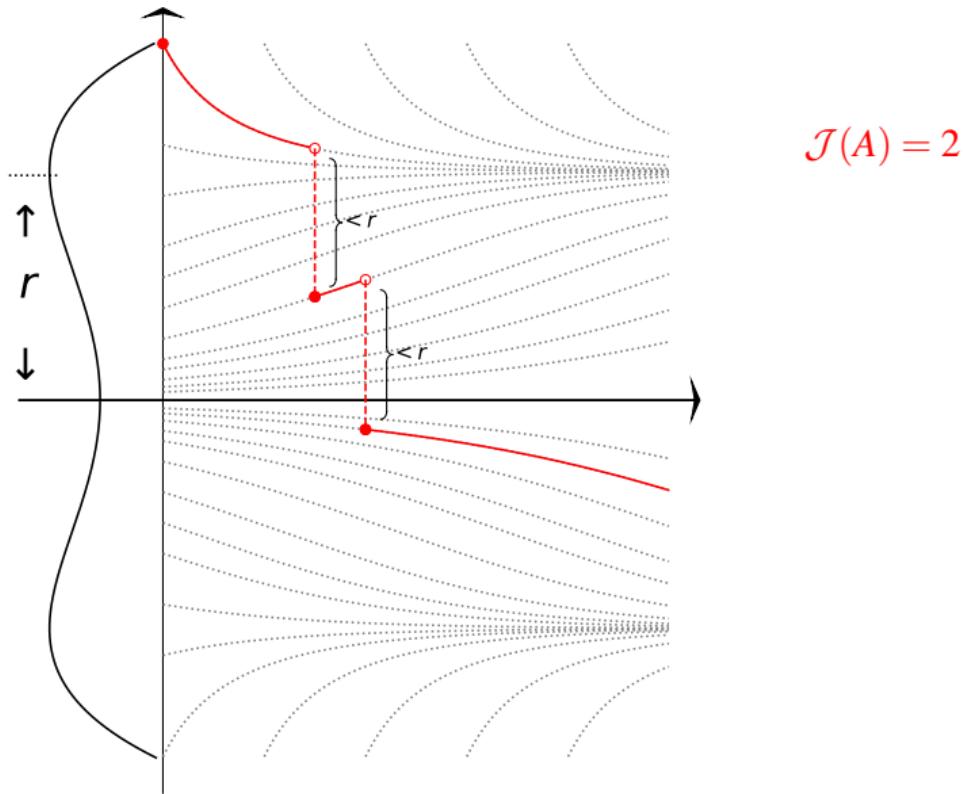
How does $\mathcal{J}(A)$ change?

If $r \in (c, 2c)$



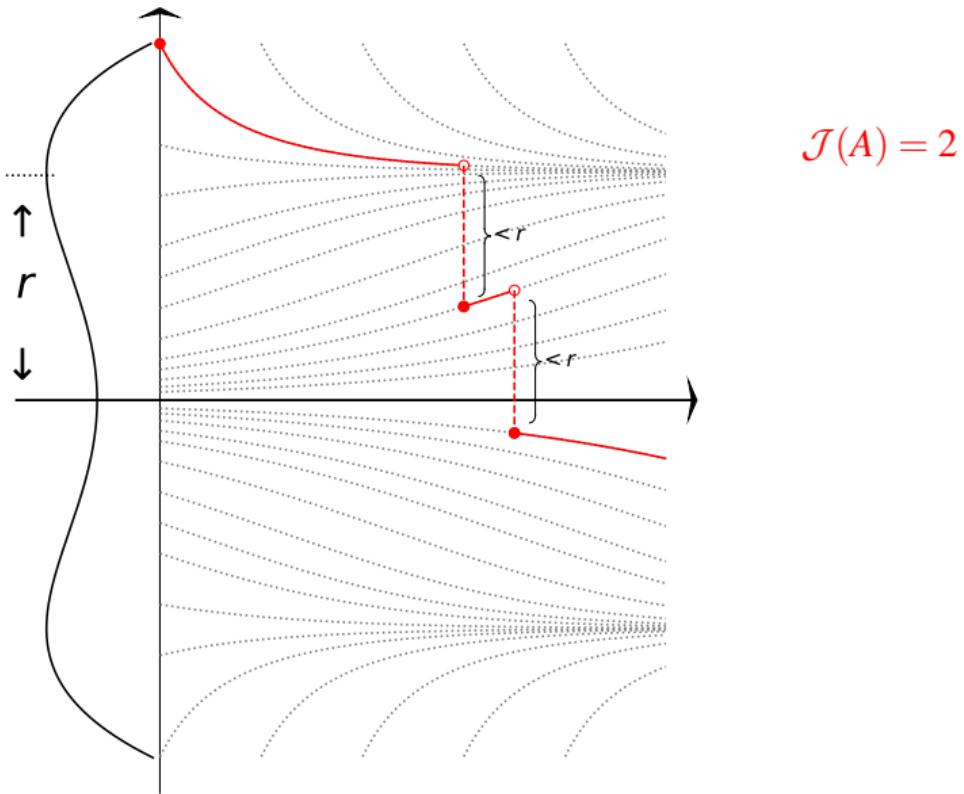
How does $\mathcal{J}(A)$ change?

If $r \in (c, 2c)$



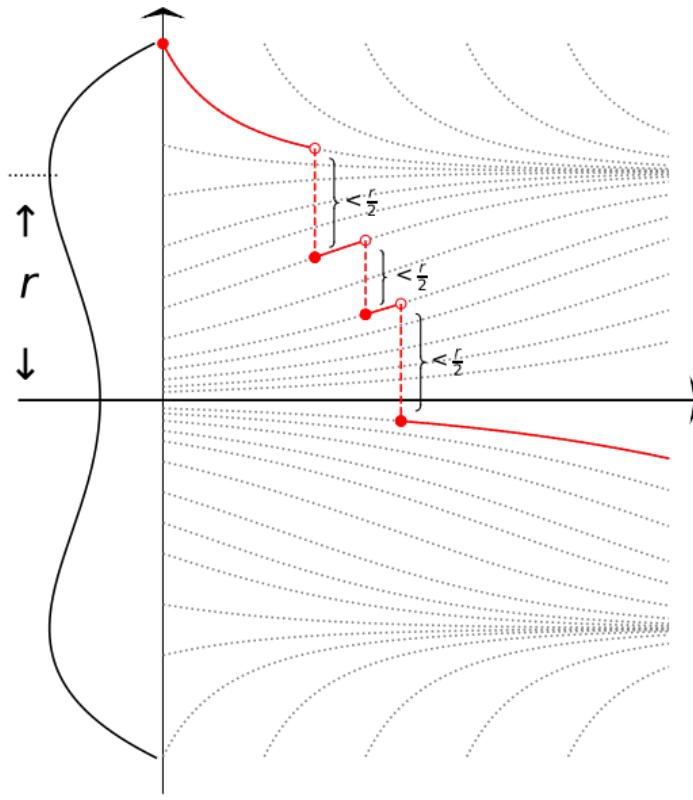
How does $\mathcal{J}(A)$ change?

If $r \in (c, 2c)$



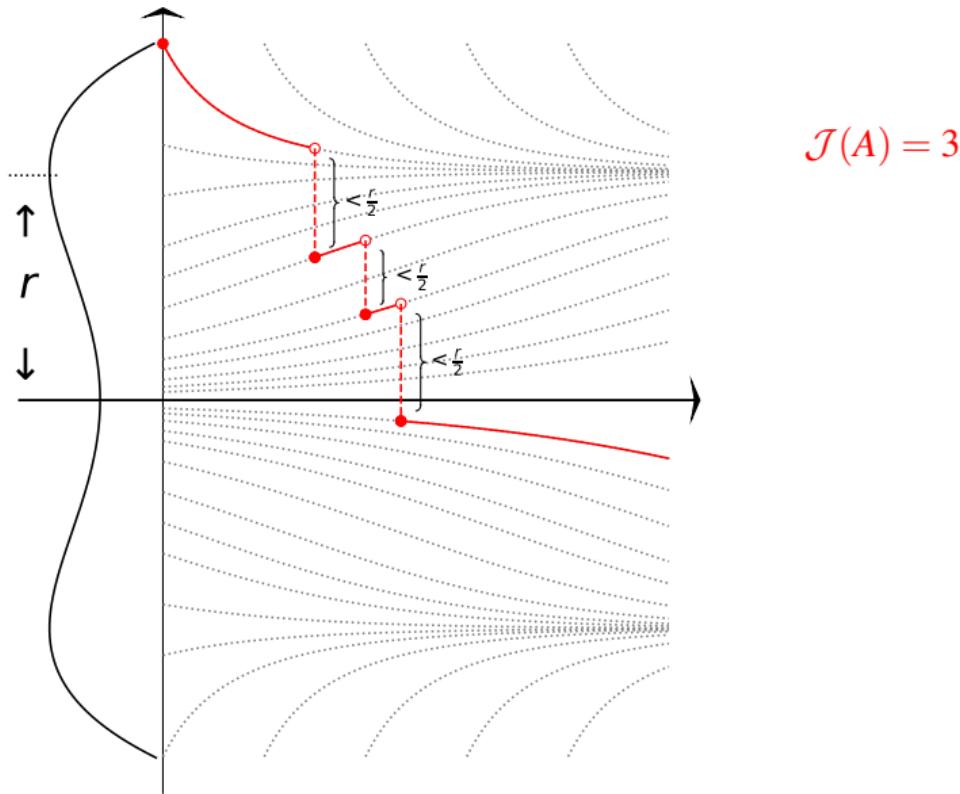
How does $\mathcal{J}(A)$ change?

If $r \in (2c, 3c)$



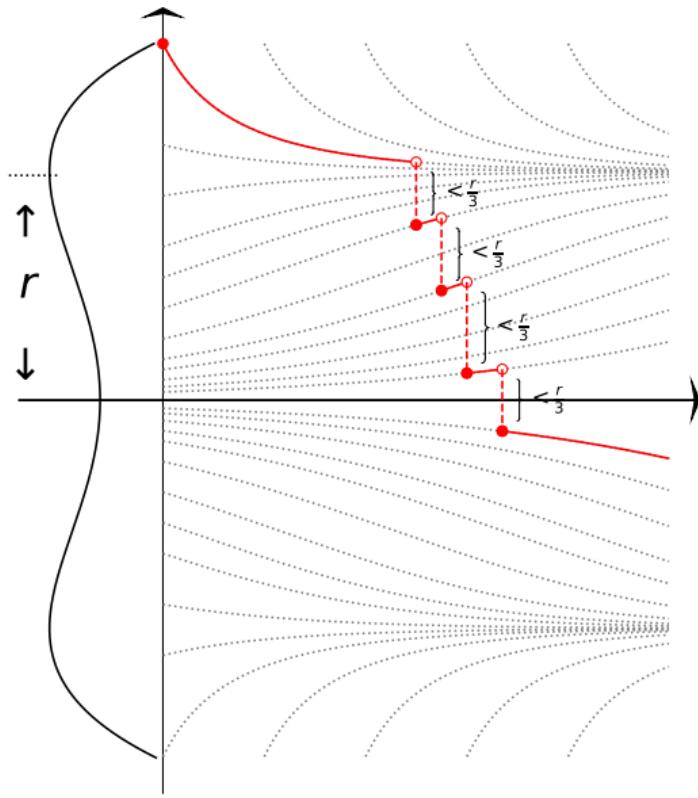
How does $\mathcal{J}(A)$ change?

If $r \in (2c, 3c)$



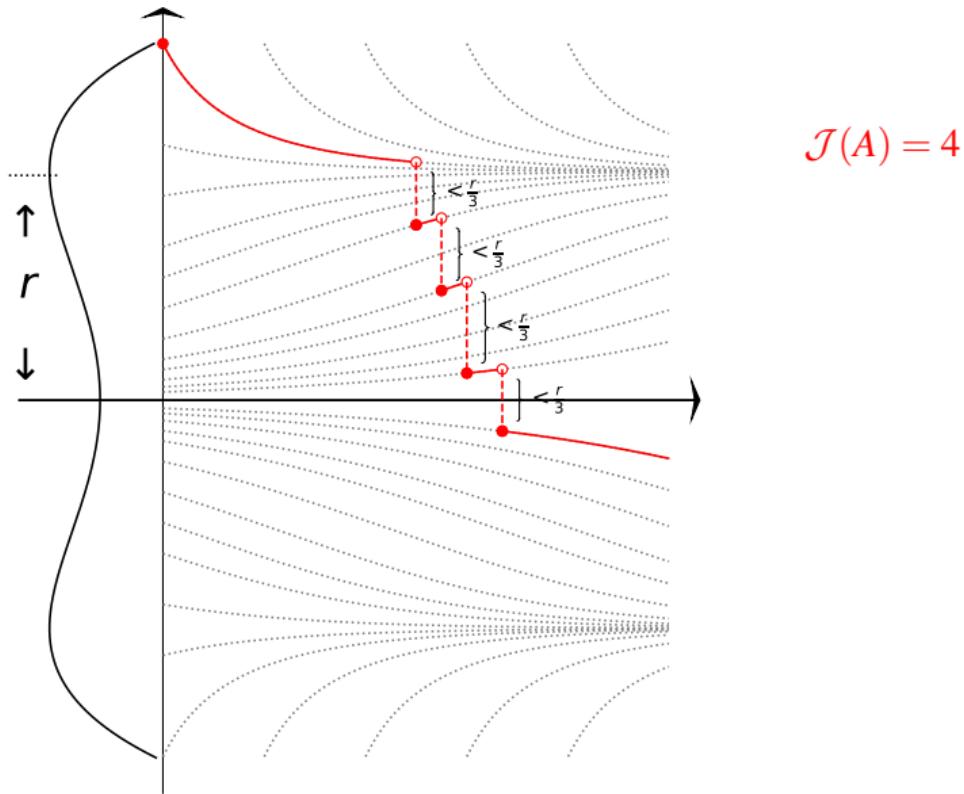
How does $\mathcal{J}(A)$ change?

If $r \in (3c, 4c)$



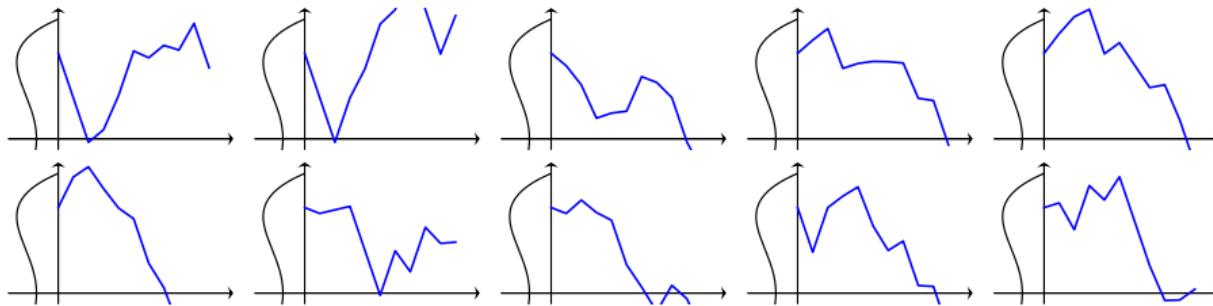
How does $\mathcal{J}(A)$ change?

If $r \in (3c, 4c)$



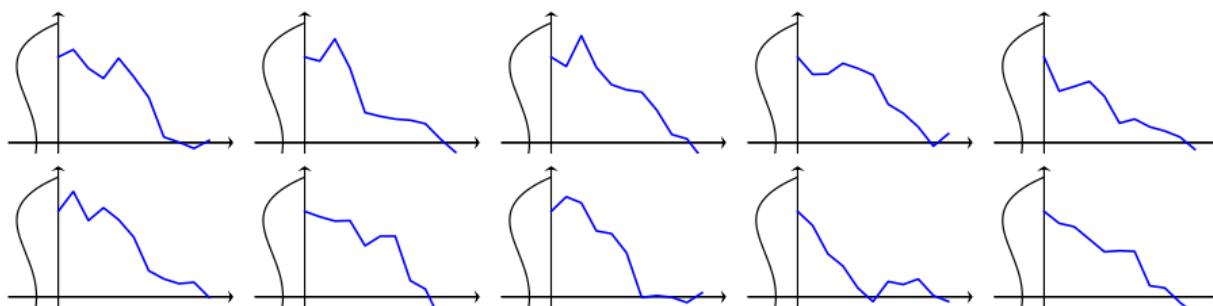
SGD's Escaping Route under Gradient Clipping with $\mathcal{J}(A) = 2$

Trajectory of SGD X^η conditional on exit:



light-tailed noises with $\eta = 1/10$

Trajectory of SGD X^η conditional on exit:

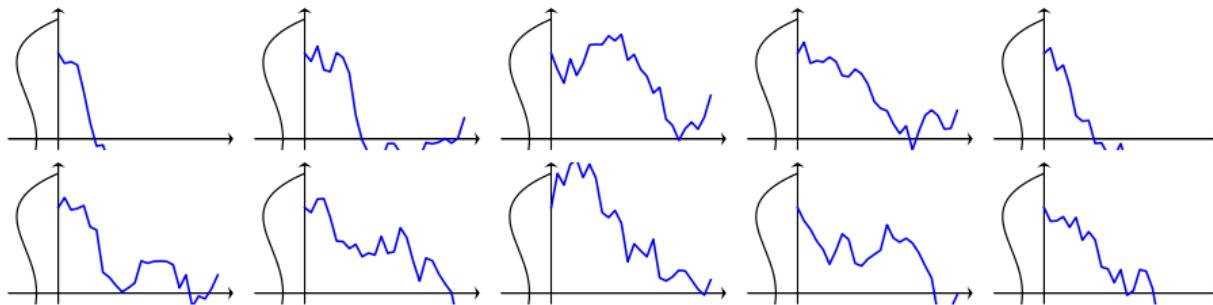


heavy-tailed noises with $\eta = 1/10$

SGD's Escaping Route under Gradient Clipping with $\mathcal{J}(A) = 2$

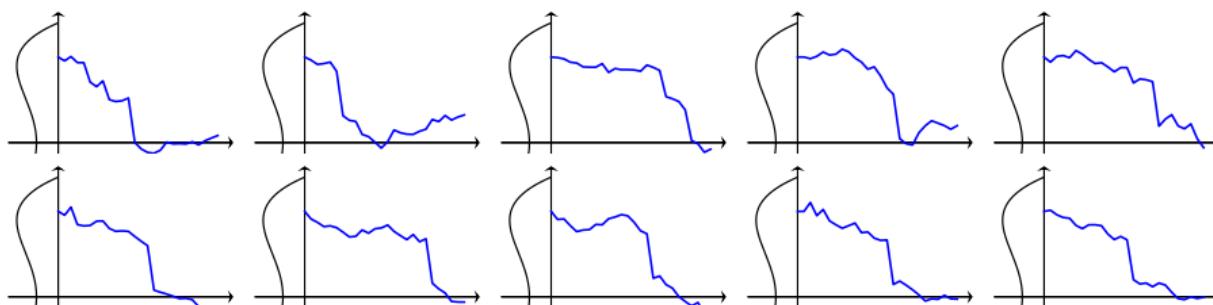
Trajectory of SGD X^η conditional on exit:

light-tailed noises with $\eta = 1/25$



Trajectory of SGD X^η conditional on exit:

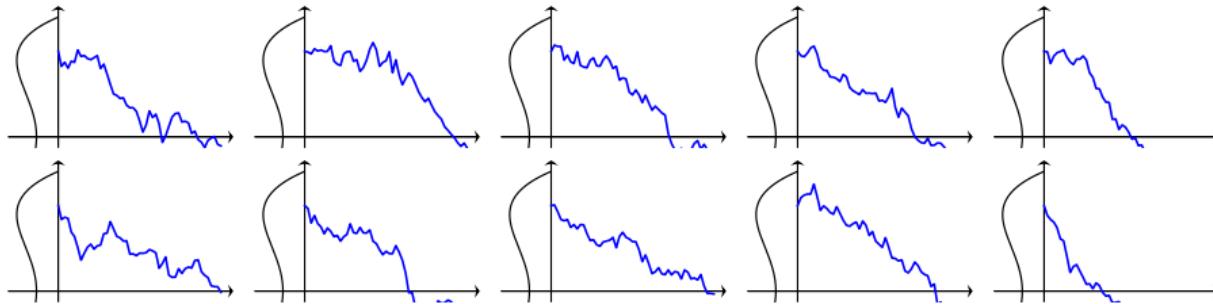
heavy-tailed noises with $\eta = 1/25$



SGD's Escaping Route under Gradient Clipping with $\mathcal{J}(A) = 2$

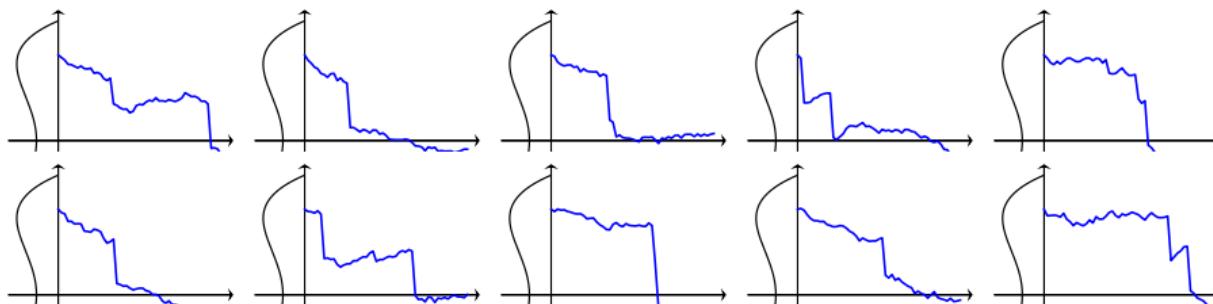
Trajectory of SGD X^η conditional on exit:

light-tailed noises with $\eta = 1/50$



Trajectory of SGD X^η conditional on exit:

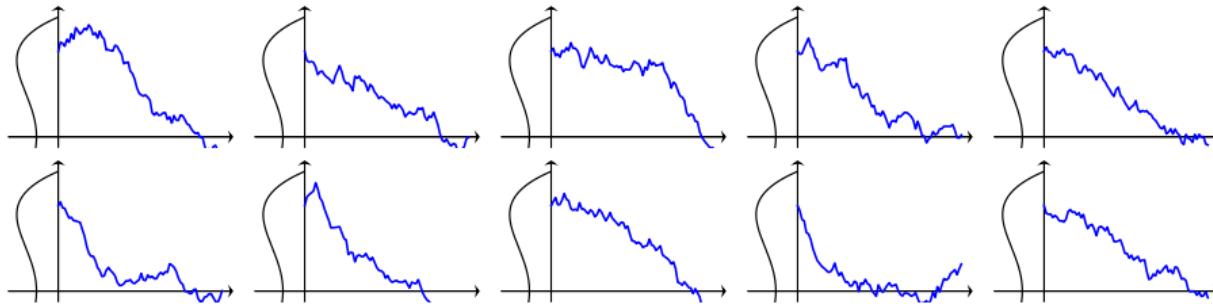
heavy-tailed noises with $\eta = 1/10$



SGD's Escaping Route under Gradient Clipping with $\mathcal{J}(A) = 2$

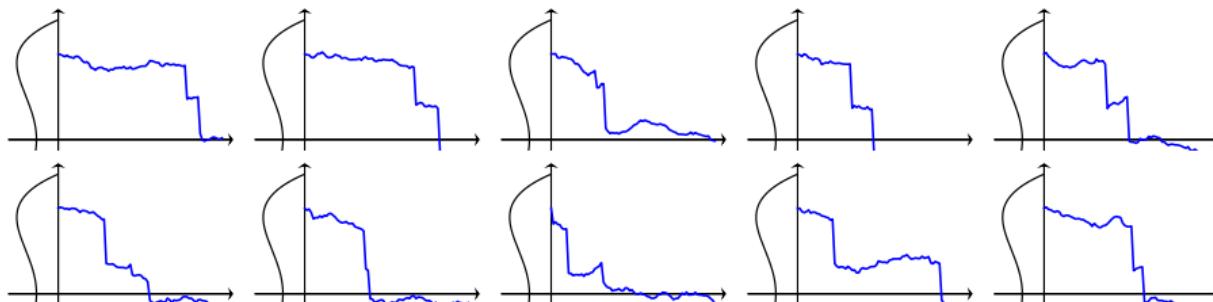
Trajectory of SGD X^η conditional on exit:

light-tailed noises with $\eta = 1/75$



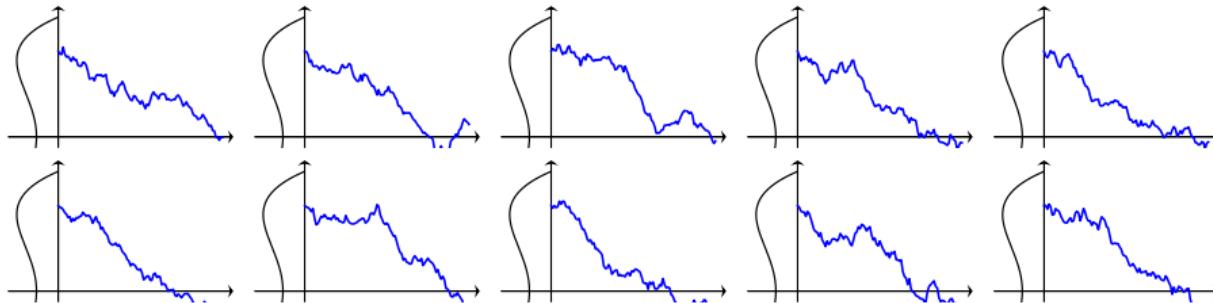
Trajectory of SGD X^η conditional on exit:

heavy-tailed noises with $\eta = 1/75$



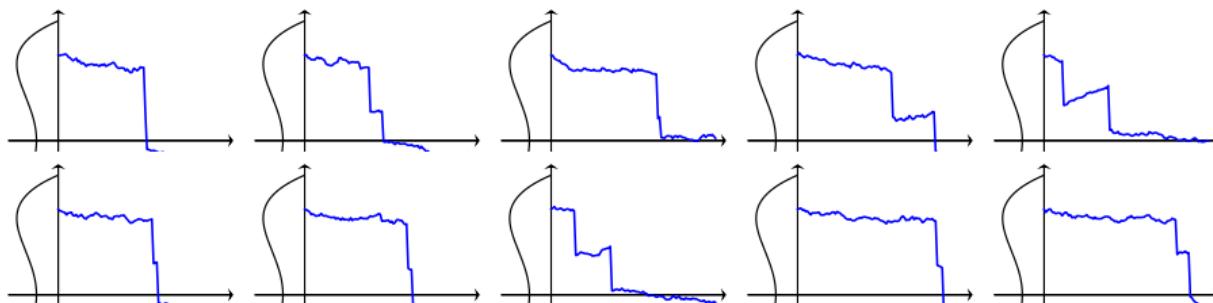
SGD's Escaping Route under Gradient Clipping with $\mathcal{J}(A) = 2$

Trajectory of SGD X^η conditional on exit:



light-tailed noises with $\eta = 1/100$

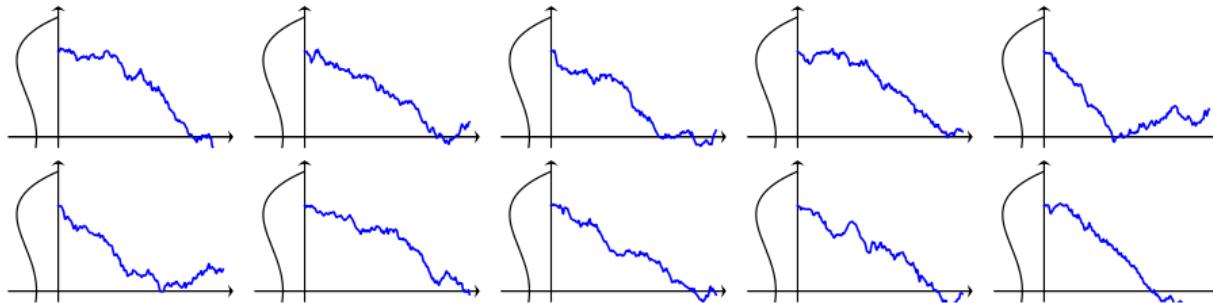
Trajectory of SGD X^η conditional on exit:



heavy-tailed noises with $\eta = 1/100$

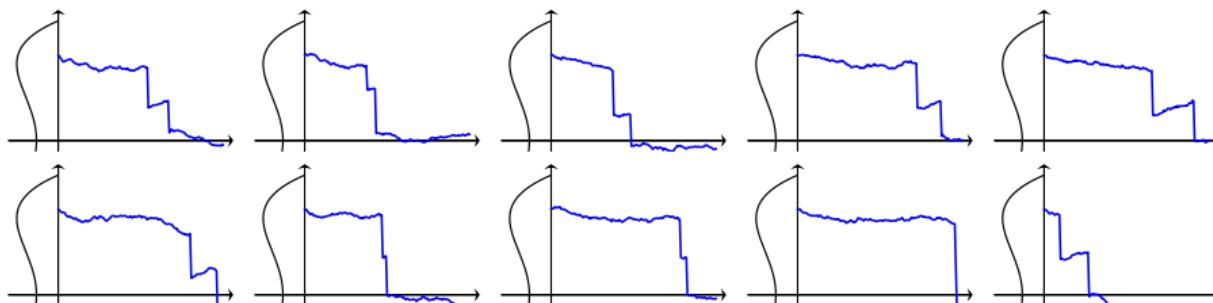
SGD's Escaping Route under Gradient Clipping with $\mathcal{J}(A) = 2$

Trajectory of SGD X^η conditional on exit:



light-tailed noises with $\eta = 1/150$

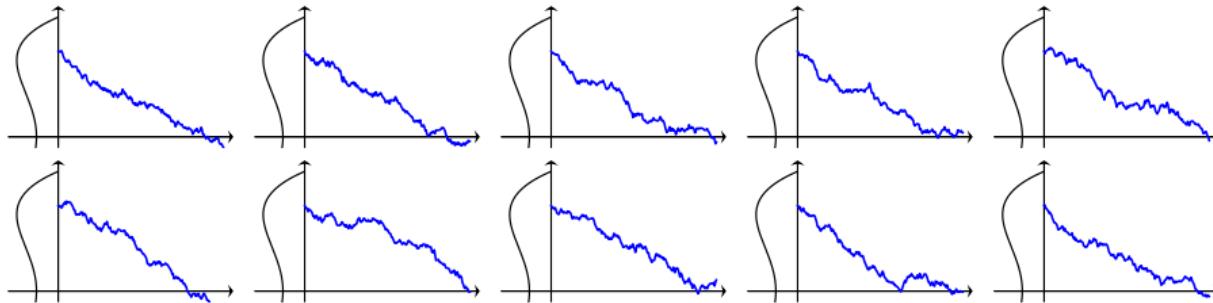
Trajectory of SGD X^η conditional on exit:



heavy-tailed noises with $\eta = 1/150$

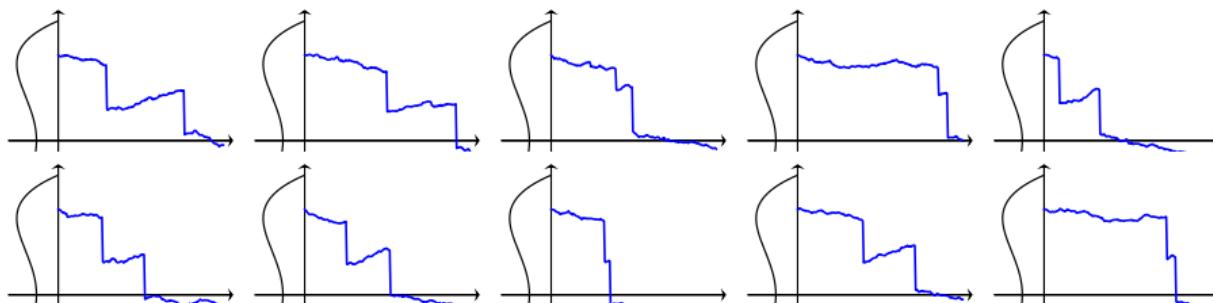
SGD's Escaping Route under Gradient Clipping with $\mathcal{J}(A) = 2$

Trajectory of SGD X^η conditional on exit:



light-tailed noises with $\eta = 1/200$

Trajectory of SGD X^η conditional on exit:

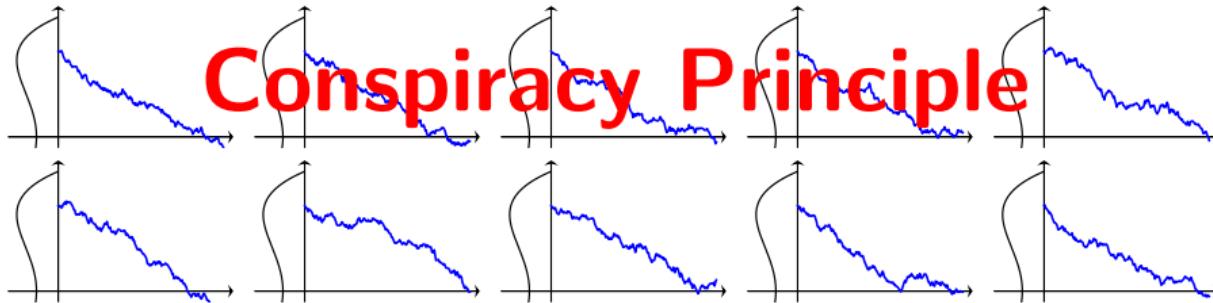


heavy-tailed noises with $\eta = 1/200$

SGD's Escaping Route under Gradient Clipping with $\mathcal{J}(A) = 2$

Trajectory of SGD X^η conditional on exit:

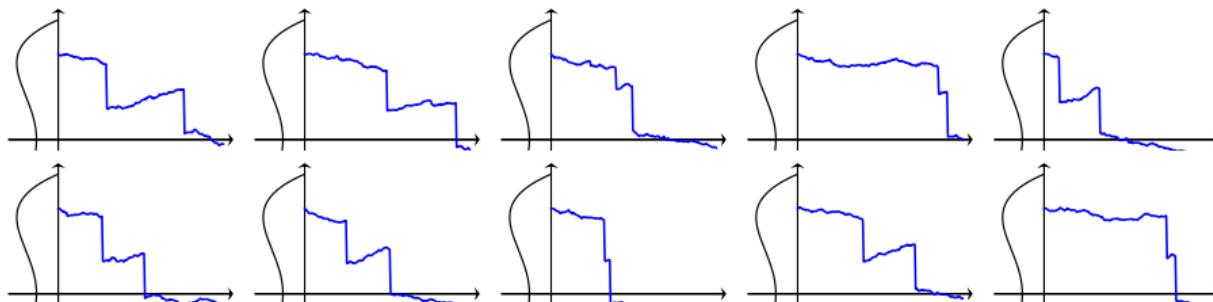
light-tailed noises with $\eta = 1/200$



Conspiracy Principle

Trajectory of SGD X^η conditional on exit:

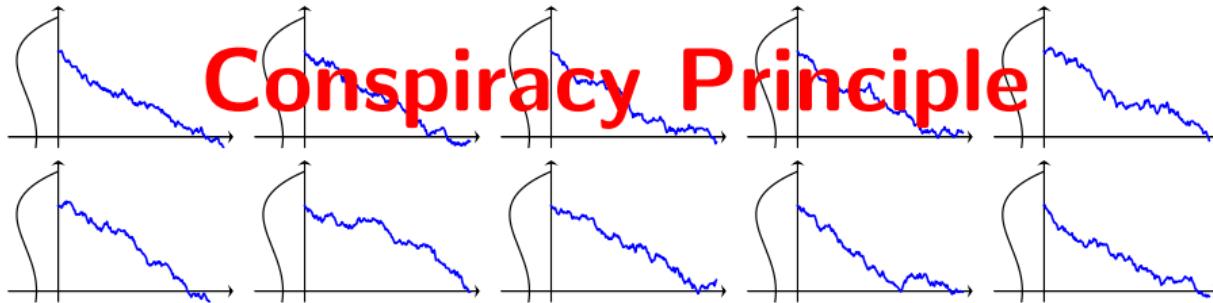
heavy-tailed noises with $\eta = 1/200$



SGD's Escaping Route under Gradient Clipping with $\mathcal{J}(A) = 2$

Trajectory of SGD X^η conditional on exit:

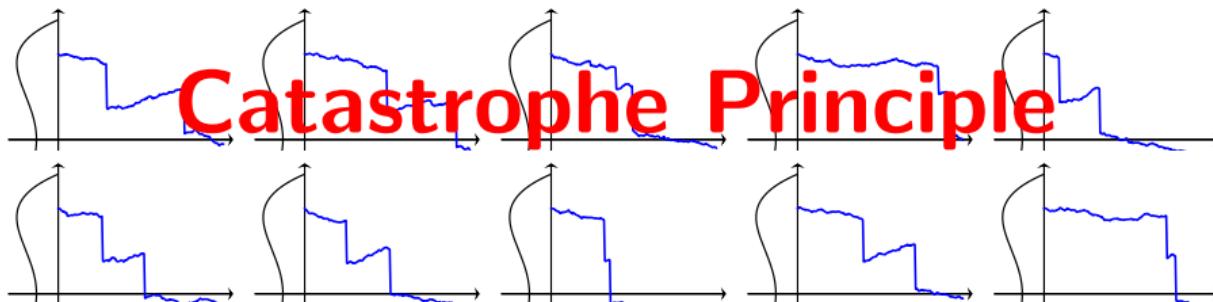
light-tailed noises with $\eta = 1/200$



Conspiracy Principle

Trajectory of SGD X^η conditional on exit:

heavy-tailed noises with $\eta = 1/200$



Catastrophe Principle

Metastability of SGD

Heavy-Tailed Large Deviations for SGD

Theorem (Wang, R., 2024+)

For “general” $B \subseteq \mathbb{D}[0, T]$

$$C(B^\circ) \leq \liminf_{\eta \rightarrow 0} \frac{\mathbf{P}(W^\eta \in B)}{\eta^{\alpha \mathcal{J}(B)}} \leq \limsup_{\eta \rightarrow 0} \frac{\mathbf{P}(W^\eta \in B)}{\eta^{\alpha \mathcal{J}(B)}} \leq C(B^-).$$

- $\mathcal{J}(B)$: min #jumps added to $w(\cdot)$ for it to be inside B
- $C(\cdot)$: a measure

Heavy-Tailed Large Deviations for SGD

$$\mathbf{P}(|Z_i| \geq x) = x^{-(\alpha+1)}, \quad \alpha > 0$$

Theorem (Wang, R., 2024+)

For “general” $B \subseteq \mathbb{D}[0, T]$

$$C(B^\circ) \leq \liminf_{\eta \rightarrow 0} \frac{\mathbf{P}(W^\eta \in B)}{\eta^{\alpha \mathcal{J}(B)}} \leq \limsup_{\eta \rightarrow 0} \frac{\mathbf{P}(W^\eta \in B)}{\eta^{\alpha \mathcal{J}(B)}} \leq C(B^-).$$

- $\mathcal{J}(B)$: min #jumps added to $w(\cdot)$ for it to be inside B
- $C(\cdot)$: a measure

Heavy-Tailed Large Deviations for SGD

$$\mathbf{P}(|Z_i| \geq x) = x^{-(\alpha+1)}, \quad \alpha > 0$$

Theorem (Wang, R., 2024+)

For “general” $B \subseteq \mathbb{D}[0, T]$

$$C(B^\circ) \leq \lim_{\varepsilon \rightarrow 0} \liminf_{\eta \rightarrow 0} \frac{\inf_{x \in A(\varepsilon)} \mathbf{P}_x(W^\eta \in B)}{\eta^{\alpha \mathcal{J}(B)}} \leq \lim_{\varepsilon \rightarrow 0} \limsup_{\eta \rightarrow 0} \frac{\sup_{x \in A(\varepsilon)} \mathbf{P}_x(W^\eta \in B)}{\eta^{\alpha \mathcal{J}(B)}} \leq C(B^-).$$

- $\mathcal{J}(B)$: min #jumps added to $w(\cdot)$ for it to be inside B
- $C(\cdot)$: a measure

Heavy-Tailed Large Deviations for SGD

$$\mathbf{P}(|Z_i| \geq x) = x^{-(\alpha+1)}, \quad \alpha > 0$$

Theorem (Wang, R., 2024+)

For “general” $B \subseteq \mathbb{D}[0, T]$

$$C(B^\circ) \leq \lim_{\varepsilon \rightarrow 0} \liminf_{\eta \rightarrow 0} \frac{\inf_{x \in A(\varepsilon)} \mathbf{P}_x(W^\eta \in B)}{\eta^{\alpha \mathcal{J}(B)}} \leq \lim_{\varepsilon \rightarrow 0} \limsup_{\eta \rightarrow 0} \frac{\sup_{x \in A(\varepsilon)} \mathbf{P}_x(W^\eta \in B)}{\eta^{\alpha \mathcal{J}(B)}} \leq C(B^-).$$

- $\mathcal{J}(B)$: min #jumps added to $w(\cdot)$ for it to be inside B
- $C(\cdot)$: a measure

Locally Uniform Large Deviations over Asymptotic Atom $\{A(\varepsilon) : \varepsilon > 0\}$

M-Convergence

$$\mathbb{M}(\mathbb{S} \setminus \mathbb{C}) = \{v(\cdot) : v(\mathbb{S} \setminus \mathbb{C}^\varepsilon) < \infty, \forall \varepsilon > 0; \quad v(\cdot) \text{ Borel measure on } \mathbb{S} \setminus \mathbb{C}\}$$

Definition (M-convergence; Lindskog, Resnick, Roy., 2014)

Let $\mu^\eta, \mu \in \mathbb{M}(\mathbb{S} \setminus \mathbb{C})$ for each $\eta > 0$. We say that μ^η converges to μ in $\mathbb{M}(\mathbb{S} \setminus \mathbb{C})$ as $\eta \rightarrow 0$ if

$$\lim_{\eta \downarrow 0} |\mu^\eta(f) - \mu(f)| = 0 \quad \forall f \in \mathcal{C}(\mathbb{S} \setminus \mathbb{C}).$$

M-Convergence

$\nwarrow \varepsilon$ -fattening of \mathbb{C}

$$\mathbb{M}(\mathbb{S} \setminus \mathbb{C}) = \{v(\cdot) : v(\mathbb{S} \setminus \mathbb{C}^\varepsilon) < \infty, \forall \varepsilon > 0; \quad v(\cdot) \text{ Borel measure on } \mathbb{S} \setminus \mathbb{C}\}$$

Definition (M-convergence; Lindskog, Resnick, Roy., 2014)

Let $\mu^\eta, \mu \in \mathbb{M}(\mathbb{S} \setminus \mathbb{C})$ for each $\eta > 0$. We say that μ^η converges to μ in $\mathbb{M}(\mathbb{S} \setminus \mathbb{C})$ as $\eta \rightarrow 0$ if

$$\lim_{\eta \downarrow 0} |\mu^\eta(f) - \mu(f)| = 0 \quad \forall f \in \mathcal{C}(\mathbb{S} \setminus \mathbb{C}).$$

M-Convergence

$$\mathbb{M}(\mathbb{S} \setminus \mathbb{C}) = \{v(\cdot) : v(\mathbb{S} \setminus \mathbb{C}^\varepsilon) < \infty, \forall \varepsilon > 0; \quad v(\cdot) \text{ Borel measure on } \mathbb{S} \setminus \mathbb{C}\}$$

$\swarrow \varepsilon\text{-fattening of } \mathbb{C}$

Definition (M-convergence; Lindskog, Resnick, Roy., 2014)

Let $\mu^\eta, \mu \in \mathbb{M}(\mathbb{S} \setminus \mathbb{C})$ for each $\eta > 0$. We say that μ^η converges to μ in $\mathbb{M}(\mathbb{S} \setminus \mathbb{C})$ as $\eta \rightarrow 0$ if

$$\lim_{\eta \downarrow 0} |\mu^\eta(f) - \mu(f)| = 0 \quad \forall f \in \mathcal{C}(\mathbb{S} \setminus \mathbb{C}).$$

↑
“bounded continuous functions supported on $\mathbb{S} \setminus \mathbb{C}$ ”

M-Convergence

$\swarrow \epsilon\text{-fattening of } \mathbb{C}$

$$\mathbb{M}(\mathbb{S} \setminus \mathbb{C}) = \{v(\cdot) : v(\mathbb{S} \setminus \mathbb{C}^\epsilon) < \infty, \forall \epsilon > 0; \quad v(\cdot) \text{ Borel measure on } \mathbb{S} \setminus \mathbb{C}\}$$

Definition (M-convergence; Lindskog, Resnick, Roy., 2014)

Let $\mu^\eta, \mu \in \mathbb{M}(\mathbb{S} \setminus \mathbb{C})$ for each $\eta > 0$. We say that μ^η converges to μ in $\mathbb{M}(\mathbb{S} \setminus \mathbb{C})$ as $\eta \rightarrow 0$ if

$$\lim_{\eta \downarrow 0} |\mu^\eta(f) - \mu(f)| = 0 \quad \forall f \in \mathcal{C}(\mathbb{S} \setminus \mathbb{C}).$$

↑
“bounded continuous functions supported on $\mathbb{S} \setminus \mathbb{C}$ ”

Definition (Uniform M-convergence; Wang, R., 2024+)

Let Θ be a set of indices. Let $\mu_\theta^\eta, \mu_\theta \in \mathbb{M}(\mathbb{S} \setminus \mathbb{C})$ for each $\eta > 0$ and $\theta \in \Theta$. We say that μ_θ^η converges to μ_θ in $\mathbb{M}(\mathbb{S} \setminus \mathbb{C})$ uniformly in θ on Θ as $\eta \rightarrow 0$ if

$$\lim_{\eta \downarrow 0} \sup_{\theta \in \Theta} |\mu_\theta^\eta(f) - \mu_\theta(f)| = 0 \quad \forall f \in \mathcal{C}(\mathbb{S} \setminus \mathbb{C}).$$

Portmanteau Theorem for Uniform $\mathbb{M}(\mathbb{S} \setminus \mathbb{C})$ -Convergence

Theorem (Wang, R., 2024+)

Let $\mu_\theta^\eta, \mu_\theta \in \mathbb{M}(\mathbb{S} \setminus \mathbb{C})$ for each $\eta > 0$ and $\theta \in \Theta$. Suppose that for any sequence $(\theta_n)_{n \geq 1}$, there exist a sub-sequence $(\theta_{n_k})_{k \geq 1}$ and $\theta^* \in \Theta$ s.t.

$$\lim_{k \rightarrow \infty} \mu_{\theta_{n_k}}(f) = \mu_{\theta^*}(f) \quad \forall f \in \mathcal{C}(\mathbb{S} \setminus \mathbb{C}).$$

Portmanteau Theorem for Uniform $\mathbb{M}(\mathbb{S} \setminus \mathbb{C})$ -Convergence

Theorem (Wang, R., 2024+)

Let $\mu_\theta^\eta, \mu_\theta \in \mathbb{M}(\mathbb{S} \setminus \mathbb{C})$ for each $\eta > 0$ and $\theta \in \Theta$. Suppose that for any sequence $(\theta_n)_{n \geq 1}$, there exist a sub-sequence $(\theta_{n_k})_{k \geq 1}$ and $\theta^* \in \Theta$ s.t.

$$\lim_{k \rightarrow \infty} \mu_{\theta_{n_k}}(f) = \mu_{\theta^*}(f) \quad \forall f \in \mathcal{C}(\mathbb{S} \setminus \mathbb{C}).$$

Then the following statements are equivalent:

- $\mu_\theta^\eta \rightarrow \mu_\theta$ in $\mathbb{M}(\mathbb{S} \setminus \mathbb{C})$ uniformly in θ on Θ as $\eta \downarrow 0$;

- For all $\varepsilon > 0$, F, G bounded away from \mathbb{C} ,
$$\liminf_{\eta \downarrow 0} \inf_{\theta \in \Theta} (\mu_\theta^\eta(G) - \mu_\theta(G_\varepsilon)) \geq 0$$
$$\limsup_{\eta \downarrow 0} \sup_{\theta \in \Theta} (\mu_\theta^\eta(F) - \mu_\theta(F^\varepsilon)) \leq 0$$

Portmanteau Theorem for Uniform $\mathbb{M}(\mathbb{S} \setminus \mathbb{C})$ -Convergence

Theorem (Wang, R., 2024+)

Let $\mu_\theta^\eta, \mu_\theta \in \mathbb{M}(\mathbb{S} \setminus \mathbb{C})$ for each $\eta > 0$ and $\theta \in \Theta$. Suppose that for any sequence $(\theta_n)_{n \geq 1}$, there exist a sub-sequence $(\theta_{n_k})_{k \geq 1}$ and $\theta^* \in \Theta$ s.t.

$$\lim_{k \rightarrow \infty} \mu_{\theta_{n_k}}(f) = \mu_{\theta^*}(f) \quad \forall f \in \mathcal{C}(\mathbb{S} \setminus \mathbb{C}).$$

Then the following statements are equivalent:

- $\mu_\theta^\eta \rightarrow \mu_\theta$ in $\mathbb{M}(\mathbb{S} \setminus \mathbb{C})$ uniformly in θ on Θ as $\eta \downarrow 0$;
- For all $\varepsilon > 0$, F, G bounded away from \mathbb{C} ,
$$\liminf_{\eta \downarrow 0} \inf_{\theta \in \Theta} (\mu_\theta^\eta(G) - \mu_\theta(G_\varepsilon)) \geq 0$$
$$\limsup_{\eta \downarrow 0} \sup_{\theta \in \Theta} (\mu_\theta^\eta(F) - \mu_\theta(F^\varepsilon)) \leq 0$$

Furthermore, they both imply

- For all open G and closed F that are bounded away from \mathbb{C} ,

$$\inf_{\theta \in \Theta} \mu_\theta(G) \leq \liminf_{\eta \downarrow 0} \inf_{\theta \in \Theta} \mu_\theta^\eta(G)$$

$$\limsup_{\eta \downarrow 0} \sup_{\theta \in \Theta} \mu_\theta^\eta(F) \leq \sup_{\theta \in \Theta} \mu_\theta(F).$$

Asymptotic Atom $\{A(\varepsilon) \subseteq \mathbb{S} : \varepsilon > 0\}$ of Markov Chain $\{V_j^\eta(x) : j \geq 0\}$

For measurable $B \subseteq \mathbb{S}$, there exist $\delta_B : (0, \infty) \times (0, \infty) \rightarrow (0, \infty)$, $\varepsilon_B > 0$, and $T_B > 0$ s.t.

$$\begin{aligned}
 C(B^\circ) - \delta_B(\varepsilon, T) &\leq \liminf_{\eta \downarrow 0} \frac{\inf_{x \in A(\varepsilon)} \mathbf{P}(\tau_{I(\varepsilon)^\complement}^\eta(x) \leq T/\eta; V_{\tau_\varepsilon}^\eta(x) \in B)}{\gamma(\eta)T/\eta} \\
 &\leq \limsup_{\eta \downarrow 0} \frac{\sup_{x \in A(\varepsilon)} \mathbf{P}(\tau_{I(\varepsilon)^\complement}^\eta(x) \leq T/\eta; V_{\tau_\varepsilon}^\eta(x) \in B)}{\gamma(\eta)T/\eta} \leq C(B^-) + \delta_B(\varepsilon, T) \\
 &\quad \limsup_{\eta \downarrow 0} \frac{\sup_{x \in I(\varepsilon)} \mathbf{P}(\tau_{(I(\varepsilon)) \setminus A(\varepsilon))^\complement}^\eta(x) > T/\eta)}{\gamma(\eta)T/\eta} = 0 \\
 &\quad \liminf_{\eta \downarrow 0} \inf_{x \in I(\varepsilon)} \mathbf{P}(\tau_{A(\varepsilon)}^\eta(x) \leq T/\eta) = 1 \quad (\{I(\varepsilon) \subseteq I : \varepsilon > 0\}: \text{covering of } I)
 \end{aligned}$$

for any $\varepsilon \leq \varepsilon_B$ and $T \geq T_B$, where $\gamma(\eta)/\eta \rightarrow 0$ as $\eta \downarrow 0$ and δ_B 's are such that

$$\lim_{\varepsilon \rightarrow 0} \lim_{T \rightarrow \infty} \delta_B(\varepsilon, T) = 0.$$

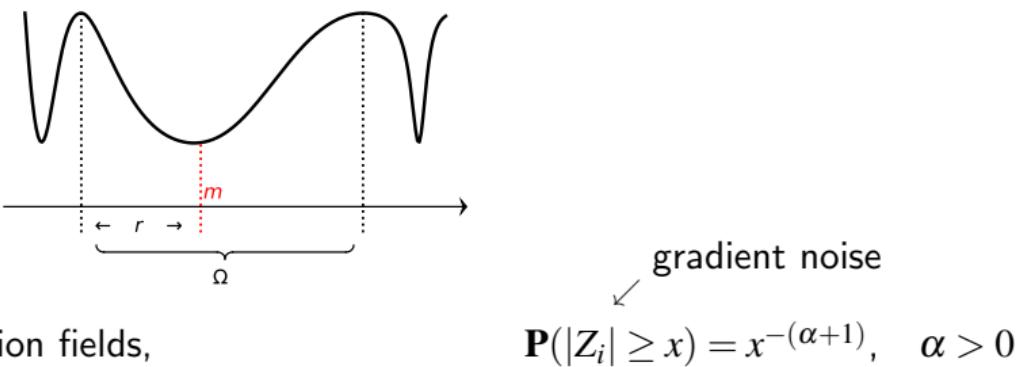
Exit Time and Location under the Presence of Asymptotic Atom

Theorem (Wang, R., 2024+)

If Markov chain $\{V_j^\eta(x) : j \geq 0\}$ possesses an asymptotic atom $\{A(\varepsilon) \subseteq \mathbb{S} : \varepsilon > 0\}$, then

$$\begin{aligned} C(B^\circ) \cdot e^{-t} &\leq \liminf_{\eta \downarrow 0} \inf_{x \in I(\varepsilon)} \mathbf{P}(\gamma(\eta) \tau_{I^c}^\eta(x) > t, V_\tau^\eta(x) \in B) \\ &\leq \limsup_{\eta \downarrow 0} \sup_{x \in I(\varepsilon)} \mathbf{P}(\gamma(\eta) \tau_{I^c}^\eta(x) > t, V_\tau^\eta(x) \in B) \leq C(B^-) \cdot e^{-t}. \end{aligned}$$

First Exit Time Analysis for SGD

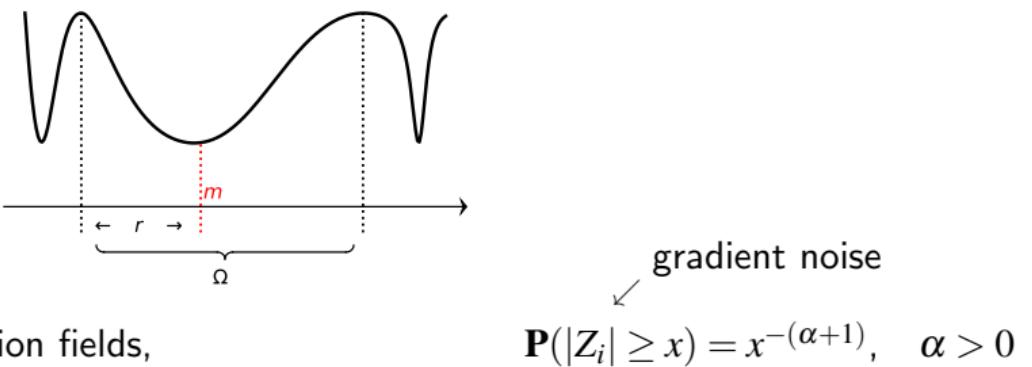


Theorem (Wang, Oh, R., 2022)

Let $\sigma(\eta) = \min\{j \geq 0 : W_j^\eta \notin \Omega\}$ and $\lambda(\eta) \sim \eta^{1+\alpha \cdot l}$

$$\sigma(n)\lambda(n) \Rightarrow \text{Exp}(1)$$

First Exit Time Analysis for SGD



$l = \lceil r/c \rceil$: “width” of the attraction fields,

$$\mathbf{P}(|Z_i| \geq x) = x^{-(\alpha+1)}, \quad \alpha > 0$$

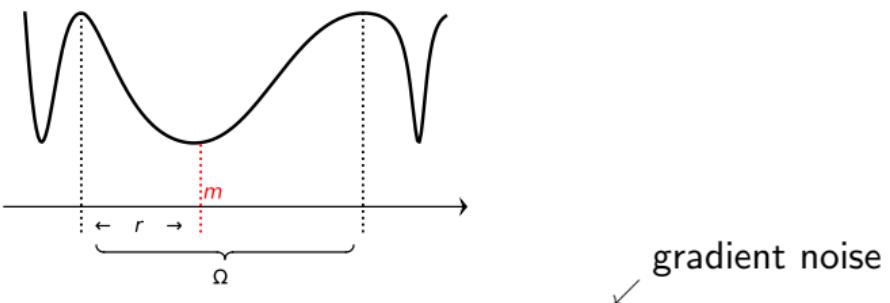
Theorem (Wang, Oh, R., 2022)

Let $\sigma(\eta) = \min\{j \geq 0 : W_j^\eta \notin \Omega\}$ and $\lambda(\eta) \sim \eta^{1+\alpha \cdot l}$

First Exit Time

$$\sigma(n)\lambda(n) \Rightarrow \text{Exp}(1)$$

First Exit Time Analysis for SGD



$l = \lceil r/c \rceil$: “width” of the attraction fields,

$$\mathbf{P}(|Z_i| \geq x) = x^{-(\alpha+1)}, \quad \alpha > 0$$

Theorem (Wang, Oh, R., 2022)

Let $\sigma(\eta) = \min\{j \geq 0 : W_j^\eta \notin \Omega\}$ and $\lambda(\eta) \sim \eta^{1+\alpha \cdot l}$

First Exit Time

$$\sigma(n)\lambda(n) \Rightarrow \text{Exp}(1)$$

$$\sim (1/\eta)^{1+\alpha \cdot l}$$

Eliminating Sharp Local Minima with Truncated Heavy-Tails

l^* : “width” of the widest attraction fields,

$$\mathbf{P}(|Z_i| \geq x) = x^{-(\alpha+1)}, \quad \alpha > 0$$

gradient noise

Theorem (Wang, Oh, R., 2022)

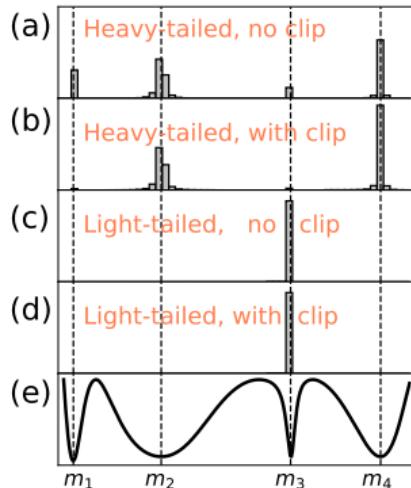
Under certain structural conditions, for any $t > 0$ and $\beta > 1 + \alpha \cdot l^*$,

$$\frac{1}{t/\eta^\beta} \int_0^{\lfloor t/\eta^\beta \rfloor} \mathbb{I}\{W_{\lfloor u \rfloor}^\eta \in \text{sharp minima}\} du \xrightarrow{p} 0$$

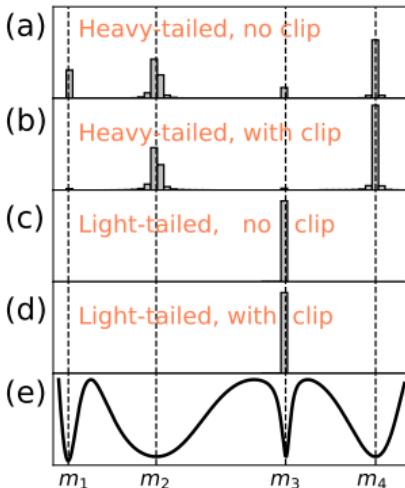
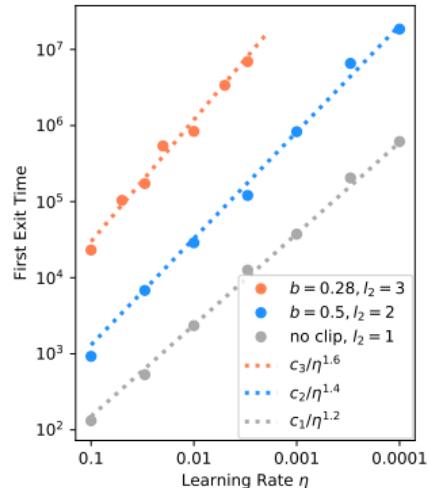
In fact, $W_{\lfloor t/\eta^{1+\alpha \cdot l^*} \rfloor}^\eta$ converges to a Markov jump processes

whose state space consists of wide local minima only.

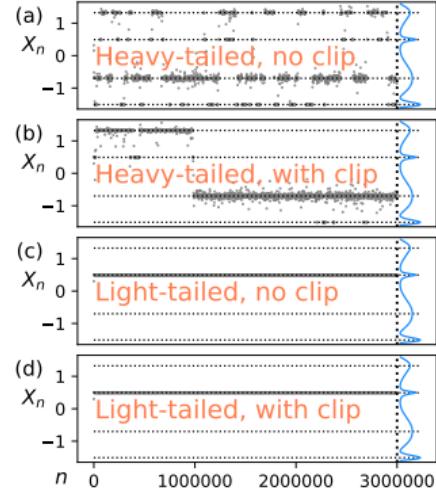
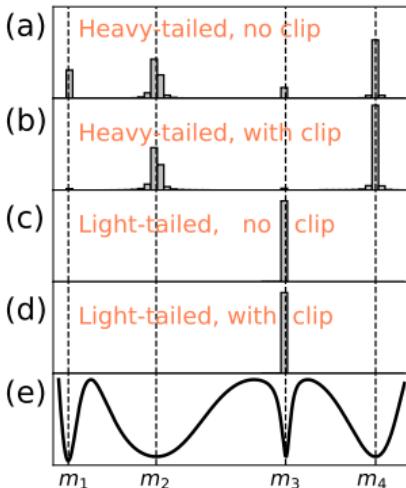
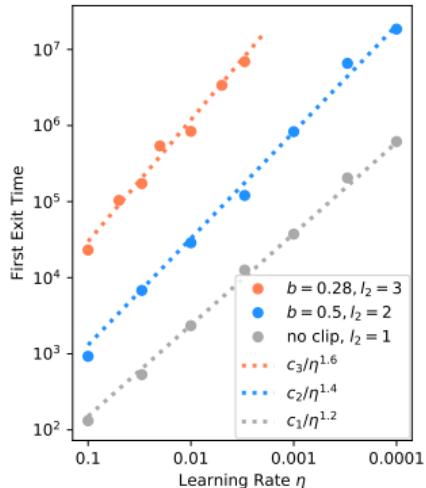
Eliminating Sharp Local Minima with Truncated Heavy-Tails



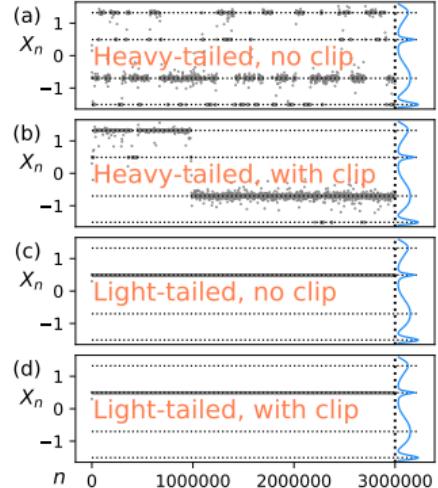
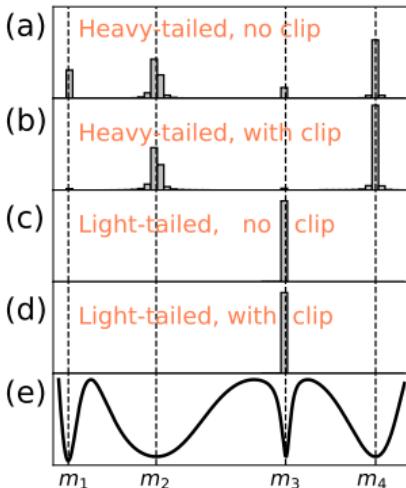
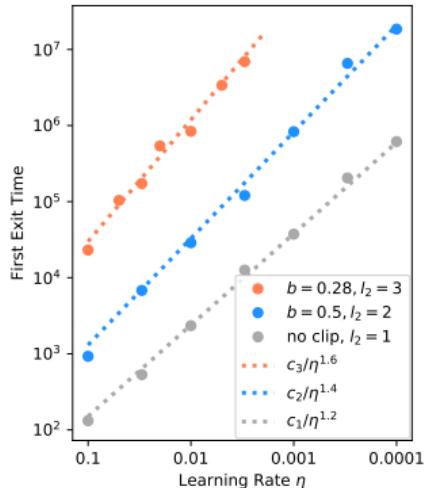
Eliminating Sharp Local Minima with Truncated Heavy-Tails



Eliminating Sharp Local Minima with Truncated Heavy-Tails

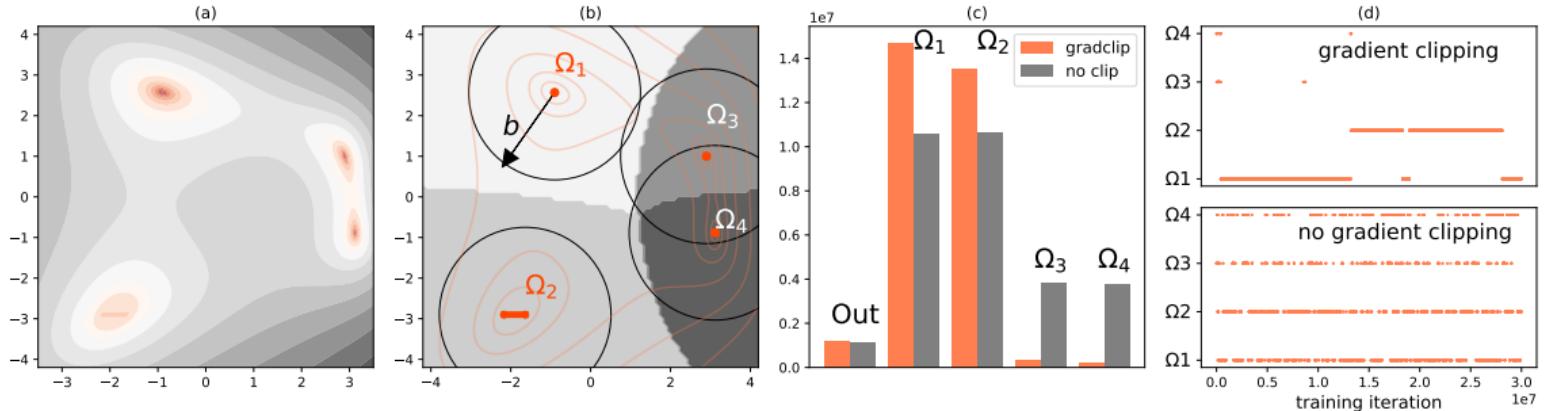


Eliminating Sharp Local Minima with Truncated Heavy-Tails

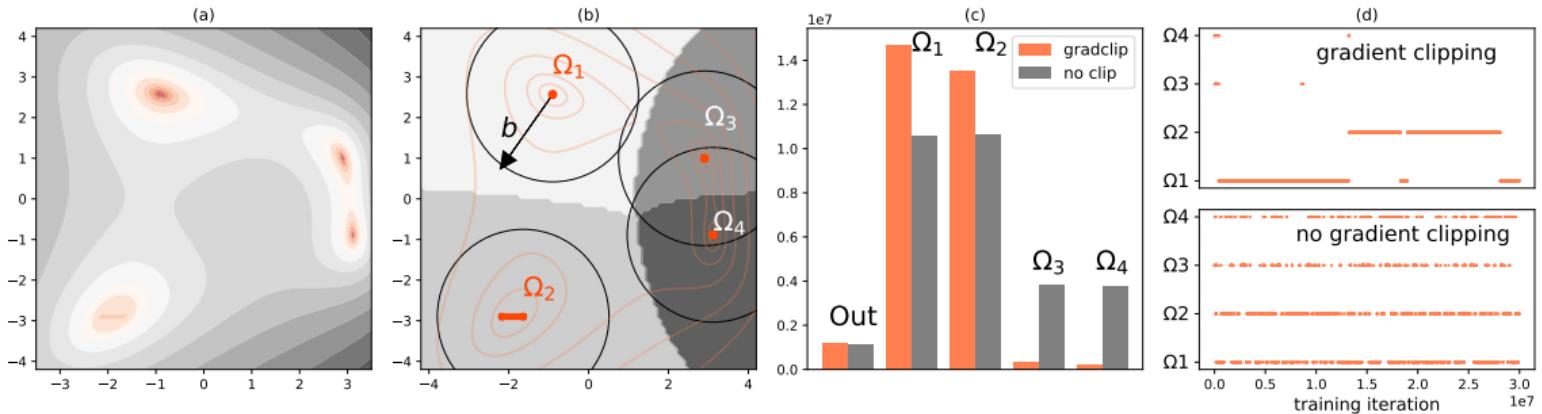


Consistent with what theory predicts!

Same Phenomena in \mathbb{R}^2 with More General Geometry



Same Phenomena in \mathbb{R}^2 with More General Geometry



Again, consistent with what theory predicts!

New Training Strategy: Tail-INflation-Truncation

Tail-Inflation-Truncation Scheme

$$\nabla \tilde{f} = \nabla f_{\text{small batch}}$$

$$W_{k+1} = W_k - \varphi_b(\eta \cdot \nabla \tilde{f}(W_k))$$

Tail-Inflation-Truncation Scheme

Stochastic Gradient

$$\nabla \tilde{f} = \nabla f_{\text{small batch}}$$

$$W_{k+1} = W_k - \varphi_b(\eta \cdot \nabla \tilde{f}(W_k))$$

Tail-Inflation-Truncation Scheme

Stochastic Gradient

$$\nabla \tilde{f}_{\text{our}} = \nabla f_{\text{small batch}} + R \cdot (\nabla f_{\text{large batch}} - \nabla f_{\text{small batch}})$$

$$W_{k+1} = W_k - \varphi_b(\eta \cdot \nabla \tilde{f}(W_k))$$

Tail-Inflation-Truncation Scheme

Stochastic Gradient

Pareto RV

$$\nabla \tilde{f}_{\text{our}} = \nabla f_{\text{small batch}} + R \cdot (\nabla f_{\text{large batch}} - \nabla f_{\text{small batch}})$$

$$W_{k+1} = W_k - \varphi_b(\eta \cdot \nabla \tilde{f}(W_k))$$

Tail-Inflation-Truncation Scheme

Heavy-Tailed Stochastic Gradient

Pareto RV

$$\nabla \tilde{f}_{\text{our}} = \nabla f_{\text{small batch}} + R \cdot (\nabla f_{\text{large batch}} - \nabla f_{\text{small batch}})$$

$$W_{k+1} = W_k - \varphi_b(\eta \cdot \nabla \tilde{f}(W_k))$$

Tail-Inflation-Truncation Scheme

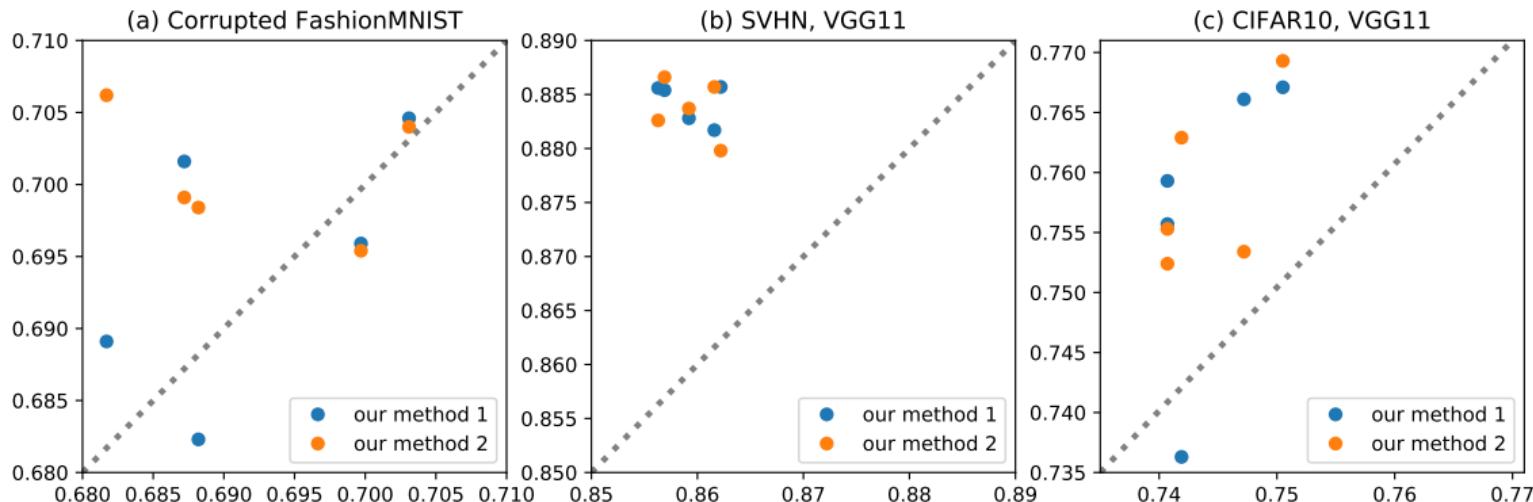
Heavy-Tailed Stochastic Gradient

Pareto RV

$$\nabla \tilde{f}_{\text{our}} = \nabla f_{\text{small batch}} + R \cdot (\nabla f_{\text{large batch}} - \nabla f_{\text{small batch}})$$

Test accuracy	LB	SB	SB+Clip	SB+Noise	Our 1	Our 2
FashionMNIST, LeNet	68.66%	69.20%	68.77%	64.43%	69.47%	70.06%
SVHN, VGG11	82.87%	85.92%	85.95%	38.85%	88.42%	88.37%
CIFAR10, VGG11	69.39%	74.42%	74.38%	40.50%	75.69%	75.87%
Expected Sharpness	LB	SB	SB+Clip	SB+Noise	Our 1	Our 2
FashionMNIST, LeNet	0.032	0.008	0.009	0.047	0.003	0.002
SVHN, VGG11	0.694	0.037	0.041	0.012	0.002	0.005
CIFAR10, VGG11	2.043	0.050	0.039	2.046	0.024	0.037

Improvement is Consistent



Tail-Inflation-Truncation Scheme

CIFAR10-VGG11	SB + Clip	Our 1	Our 2
Test Accuracy	89.54%	90.67%	90.45%
Expected Sharpness	0.167	0.085	0.096
PAC-Bayes Sharpness	1.31×10^4	9×10^3	10^4
Maximal Sharpness	1.66×10^4	1.29×10^4	1.22×10^4
CIFAR100-VGG16	SB + Clip	Our 1	Our 2
Test Accuracy	56.32%	65.44%	62.99%
Expected Sharpness	0.857	0.441	0.479
PAC-Bayes Sharpness	2.49×10^4	1.9×10^4	1.98×10^4
Maximal Sharpness	2.75×10^4	2.12×10^4	2.16×10^4

Does This Actually Work with High-Volume Real-Life Data?

We will know soon:

- Moloco
 - top performing mobile ad tech company
 - 35B+ user impressions per month
- A project to improve Moloco's conversion-rate prediction accuracy

Summary

Summary

- Catastrophe Principle for various heavy-tailed stochastic processes

Summary

- Catastrophe Principle for various heavy-tailed stochastic processes
- Solutions to open problems in queueing theory and rare-event simulation

Summary

- Catastrophe Principle for various heavy-tailed stochastic processes
- Solutions to open problems in queueing theory and rare-event simulation
- Metastability analysis reveals the global dynamics of heavy-tailed dynamical systems

Summary

- Catastrophe Principle for various heavy-tailed stochastic processes
- Solutions to open problems in queueing theory and rare-event simulation
- Metastability analysis reveals the global dynamics of heavy-tailed dynamical systems
- Elimination of sharp local minima from SGD with truncated heavy-tailed gradient noise