

Theoretical Framework for Global Dynamics and Metastability Analysis of Heavy-Tailed Systems

Chang-Han Rhee

Northwestern University

2023 INFORMS Annual Meeting, Phoenix

October 16, 2023

Based on the joint works with

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

How Rare Events Characterize Global Dynamics of Heavy-Tailed Systems

Chang-Han Rhee

Northwestern University

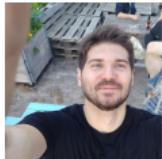
2023 INFORMS Annual Meeting, Phoenix

October 16, 2023

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

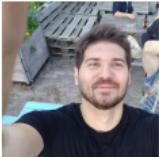


Zhe Su
Northwestern



Bert Zwart
CWI

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

(2023 Nicholson Student Paper Prize, 2nd Place)

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

A Central Mystery of Deep Learning

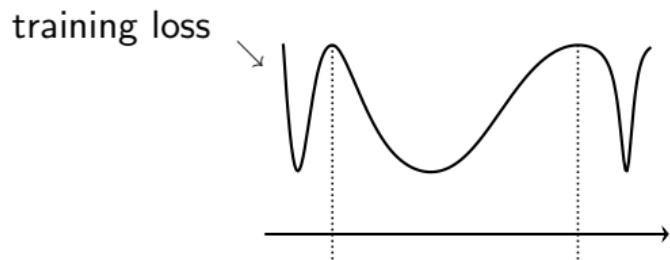
Heavy Tailed Large Deviations and Metastability

Heavy Tailed Large Deviations and Metastability

- Popular explanation: SGD somehow finds flat local minima.

Heavy Tailed Large Deviations and Metastability

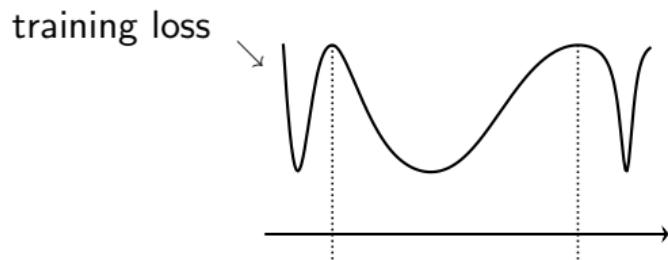
- Popular explanation: SGD somehow finds flat local minima.



Heavy Tailed Large Deviations and Metastability

- Popular explanation: SGD somehow finds flat local minima.

↑ tends to generalize well

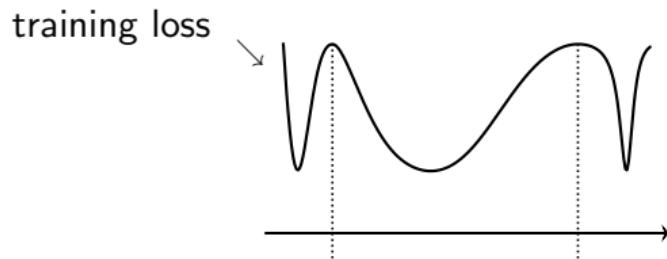


Heavy Tailed Large Deviations and Metastability

- Popular explanation: SGD somehow finds flat local minima.

tends to generalize well

- But how?

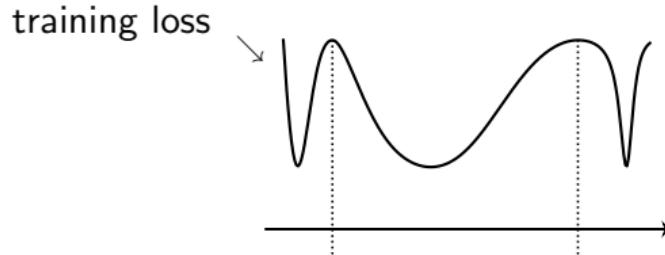


Heavy Tailed Large Deviations and Metastability

- Popular explanation: SGD somehow finds flat local minima.

tends to generalize well

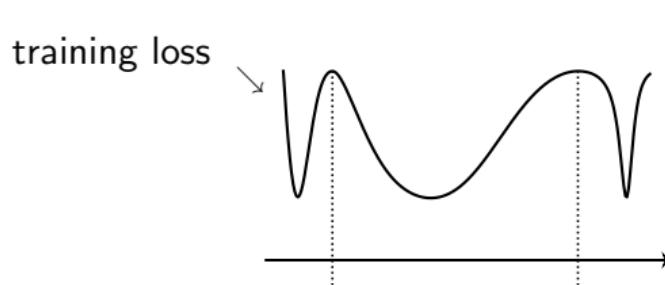
- But how?



It takes more than 10^{400} time steps
for SGD to escape from any of these.

Heavy Tailed Large Deviations and Metastability

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



It takes more than 10^{400} time steps
for SGD to escape from any of these.

Heavy Tailed Large Deviations and Metastability

- Popular explanation: SGD somehow finds flat local minima.
  tends to generalize well
- But how? Previous theoretical attempts failed to explain how.
- Heavy-tails in SGD sheds a new light to the algorithmic regularization.
eg. Simsekli et al. (2019), Hodgkinsons & Mahoney (2020), Wang, Oh, Rhee (2022), etc

Heavy Tailed Large Deviations and Metastability

- Popular explanation: SGD somehow finds flat local minima.
  tends to generalize well
- But how? Previous theoretical attempts failed to explain how.
- Heavy-tails in SGD sheds a new light to the algorithmic regularization.
eg. Simsekli et al. (2019), Hodgkinsons & Mahoney (2020), Wang, Oh, Rhee (2022), etc

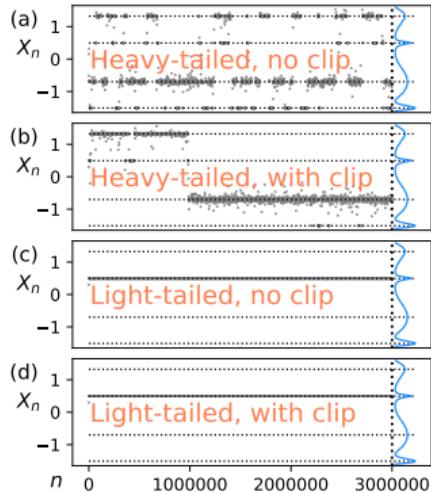
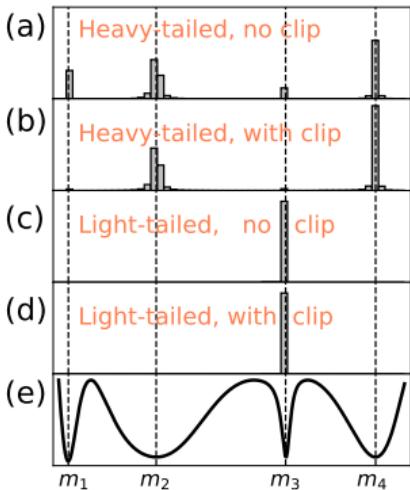
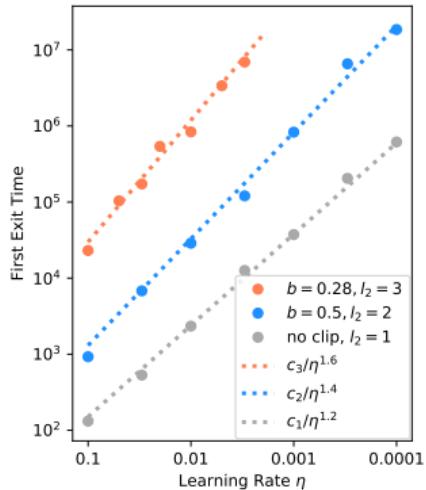
Heavy-tailed SGD escapes local minima and prefers flat local minima.

Heavy Tailed Large Deviations and Metastability

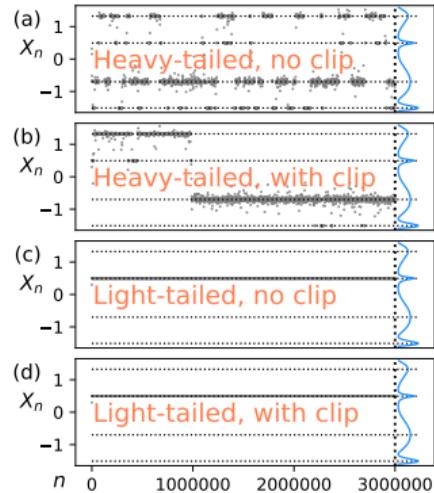
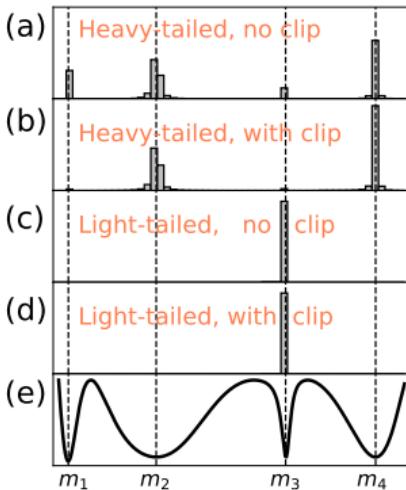
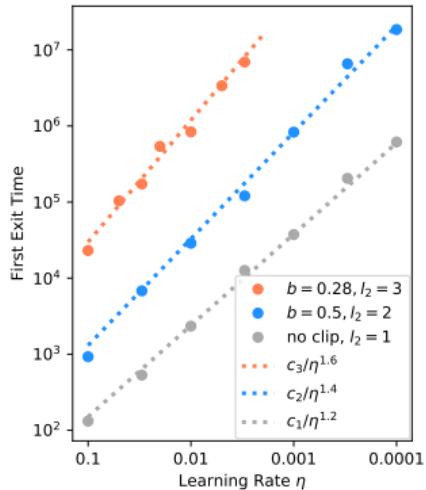
- Popular explanation: SGD somehow finds flat local minima.
  tends to generalize well
- But how? Previous theoretical attempts failed to explain how.
- Heavy-tails in SGD sheds a new light to the algorithmic regularization.
eg. Simsekli et al. (2019), Hodgkinsons & Mahoney (2020), Wang, Oh, Rhee (2022), etc

Large deviations and metastability of SGD has fascinating connection this!

Entirely Different Global Dynamics Depending on Tail Behaviors



Entirely Different Global Dynamics Depending on Tail Behaviors



Explanation?

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



Instagram

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



Instagram

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)



Instagram

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

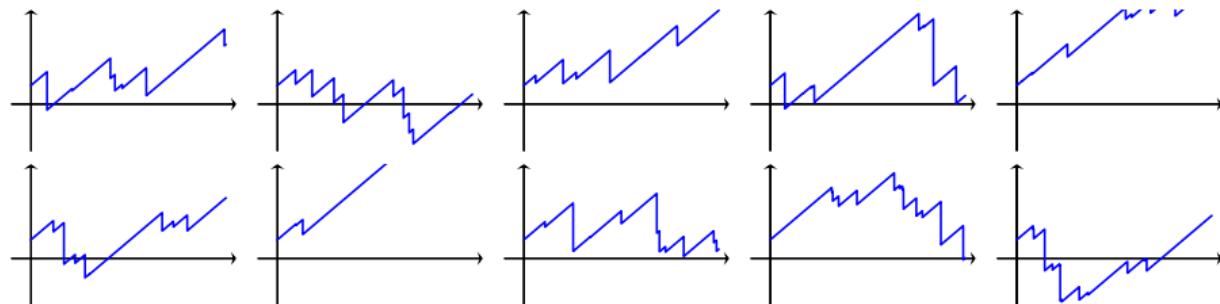
Sample paths of \bar{Y}_n :

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

Typical Scenario

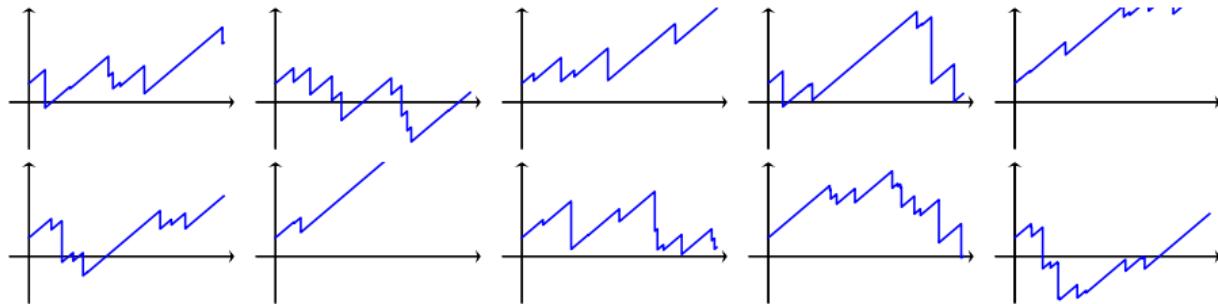
Sample paths of \bar{Y}_n :

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



Typical Scenario

Sample paths of \bar{Y}_n :



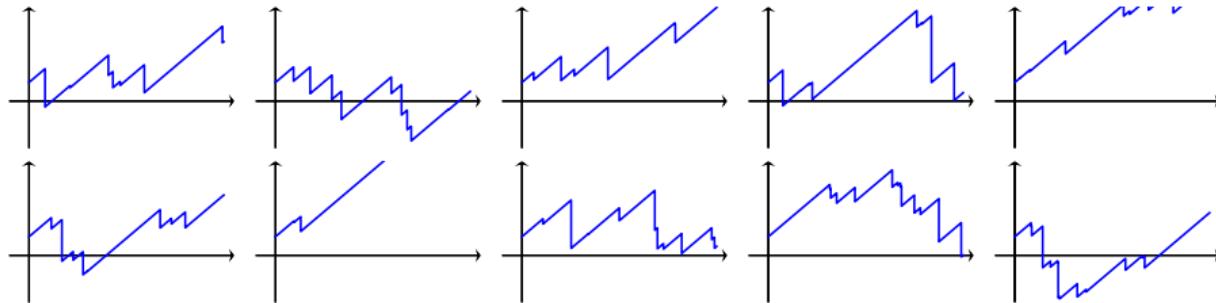
$n=10$ & claim sizes are **light-tailed**

Sample paths of \bar{Y}_n :

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

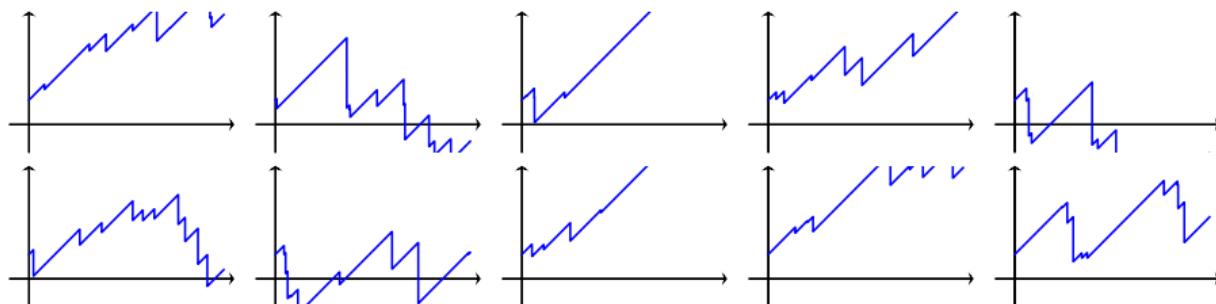
Typical Scenario

Sample paths of \bar{Y}_n :



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

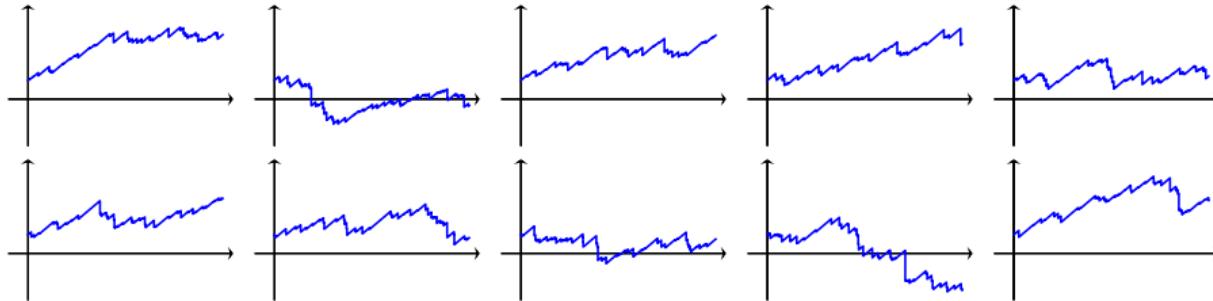
Sample paths of \bar{Y}_n :



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

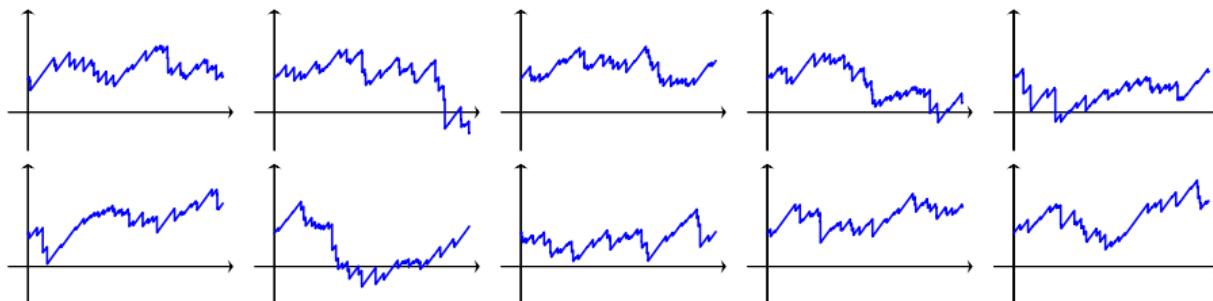
Typical Scenario

Sample paths of \bar{Y}_n :



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

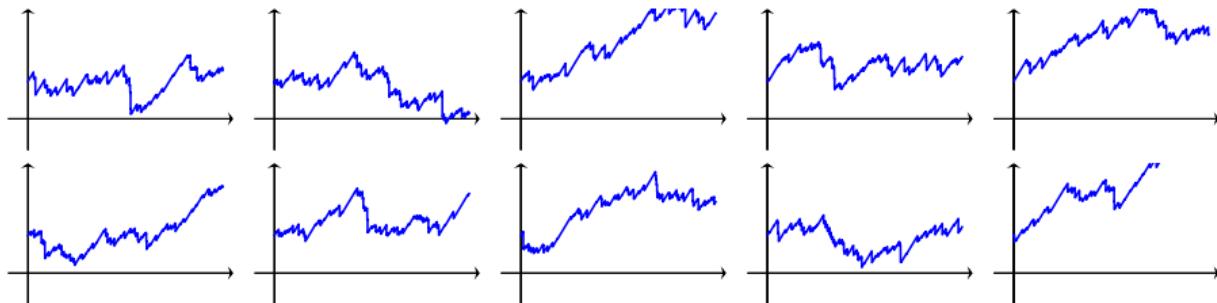
Sample paths of \bar{Y}_n :



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

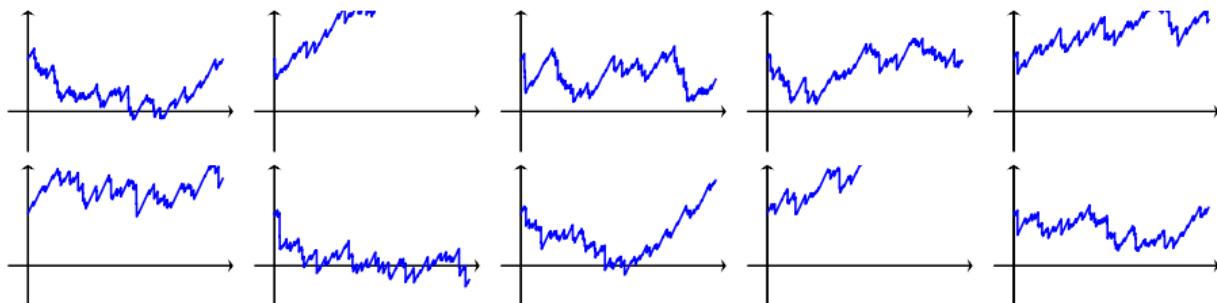
Typical Scenario

Sample paths of \bar{Y}_n :



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

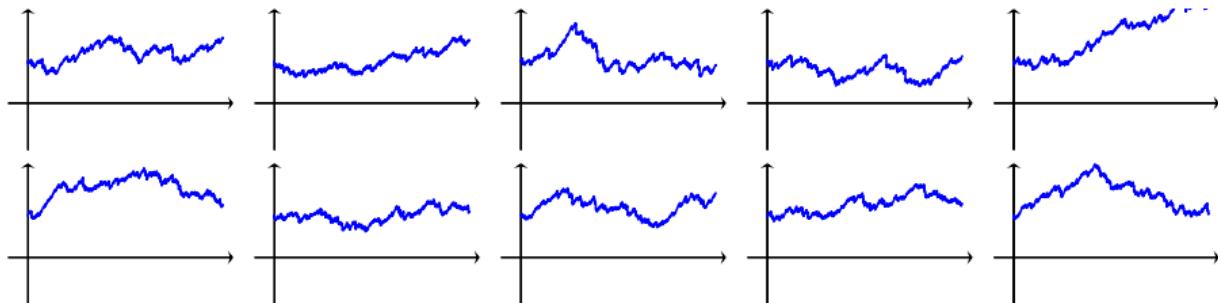
Sample paths of \bar{Y}_n :



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

Typical Scenario

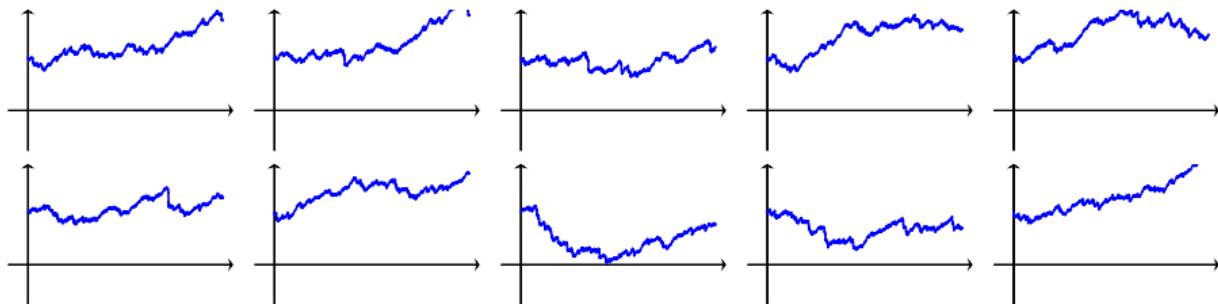
Sample paths of \bar{Y}_n :



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

Sample paths of \bar{Y}_n :

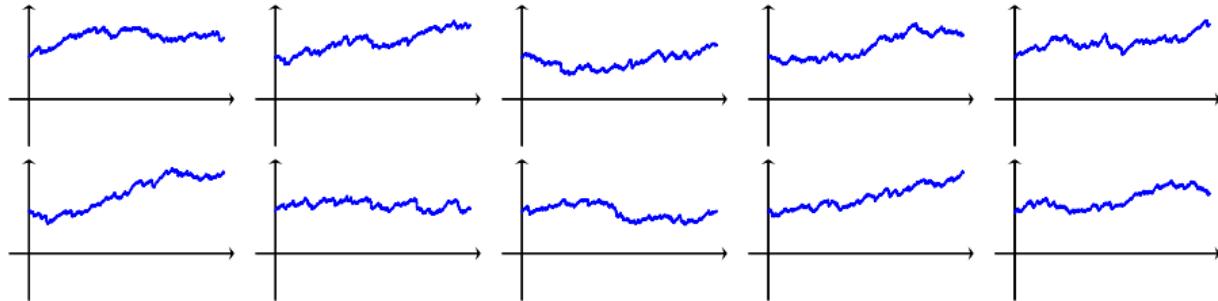
$n=500$ & claim sizes are **heavy-tailed**



Typical Scenario

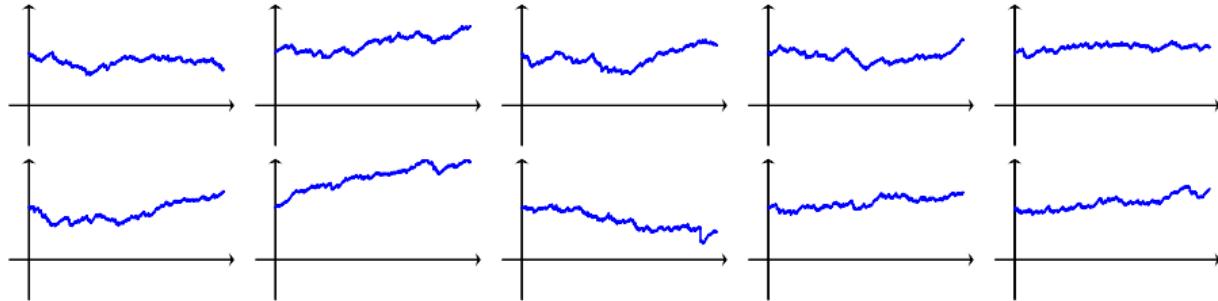
Sample paths of \bar{Y}_n :

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



Sample paths of \bar{Y}_n :

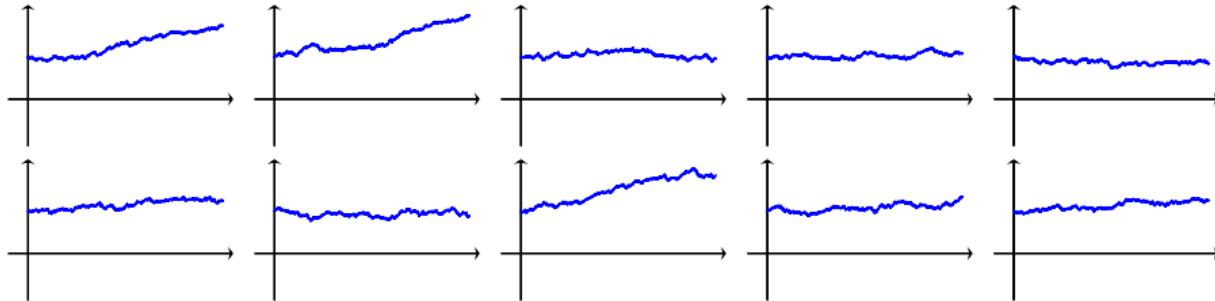
$n=1000$ & claim sizes are **heavy-tailed**



Typical Scenario

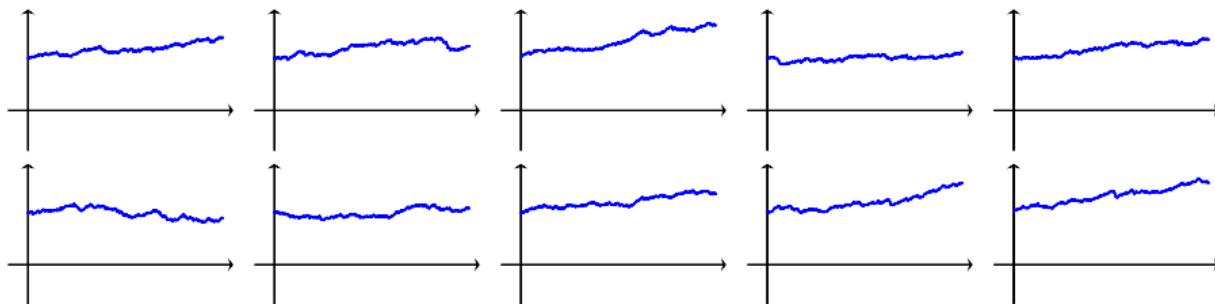
Sample paths of \bar{Y}_n :

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



Sample paths of \bar{Y}_n :

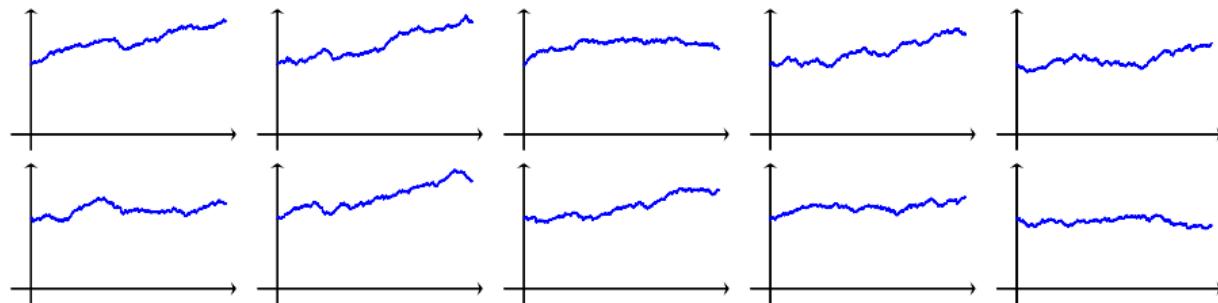
$n=2500$ & claim sizes are **heavy-tailed**



Typical Scenario

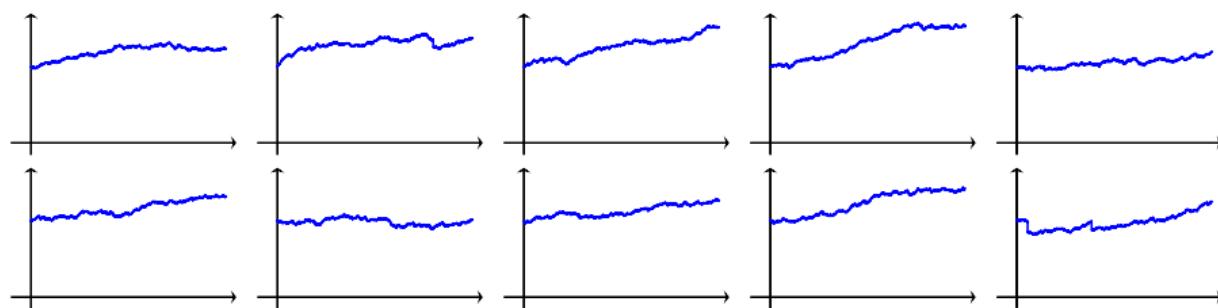
Sample paths of \bar{Y}_n :

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



Sample paths of \bar{Y}_n :

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

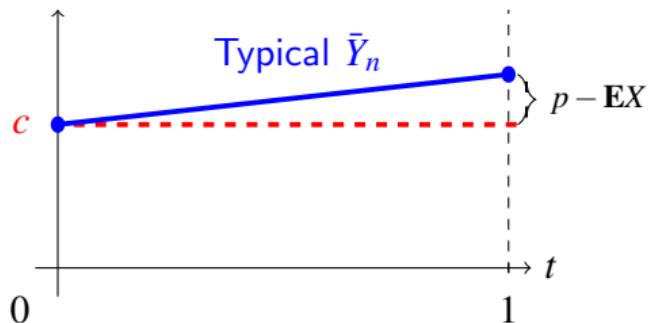


Typical Scenario

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

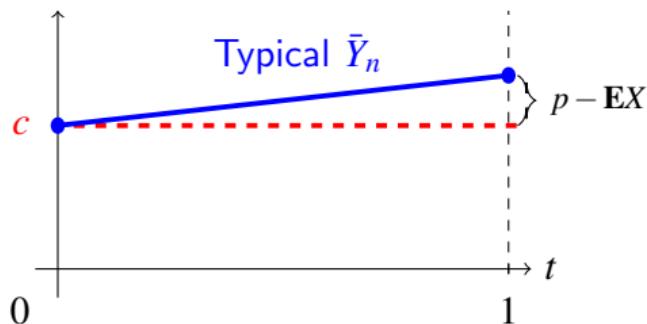
Typical Scenario

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



Typical Scenario

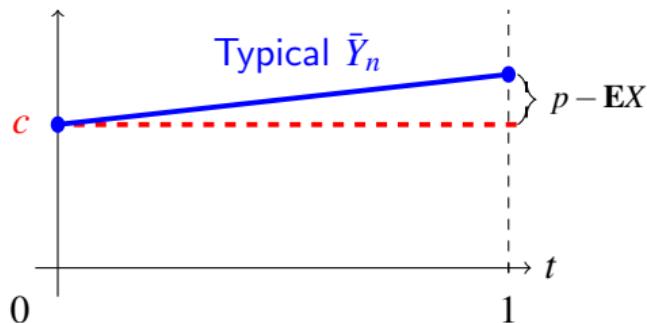
That is, $\bar{Y}_n(t) \approx c + (p - \mathbf{E}X)t$ for large n .



Typically, your insurance business will flourish

Typical Scenario

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

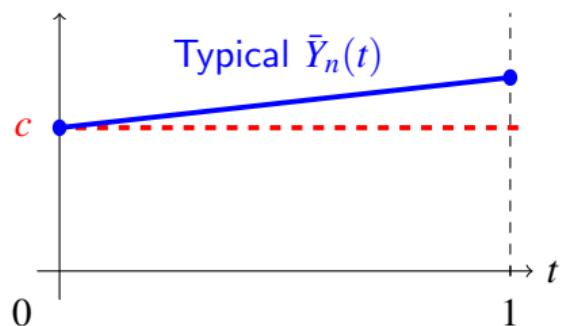


Typically, your insurance business will flourish

regardless of the tail distributions.

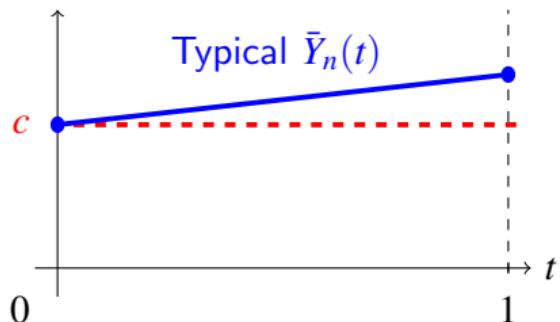
What about atypical cases?

A Rare Event: Bankruptcy



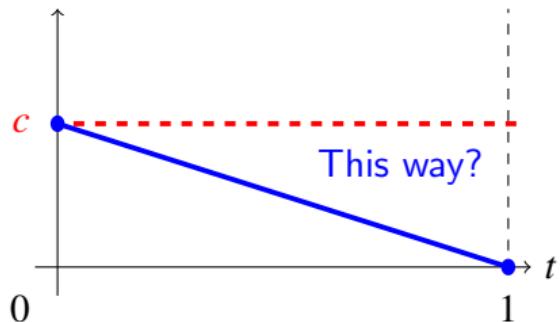
A Rare Event: Bankruptcy

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



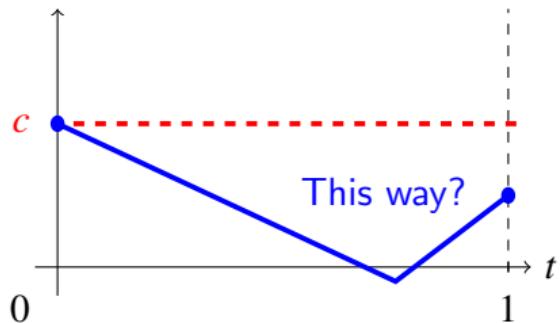
A Rare Event: Bankruptcy

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



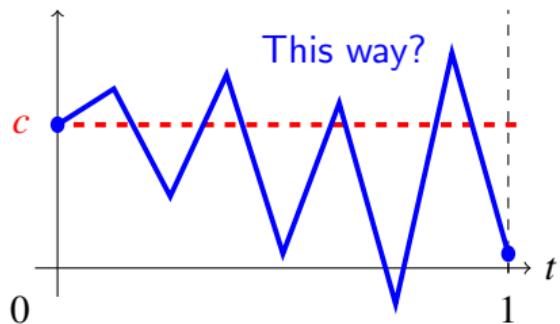
A Rare Event: Bankruptcy

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



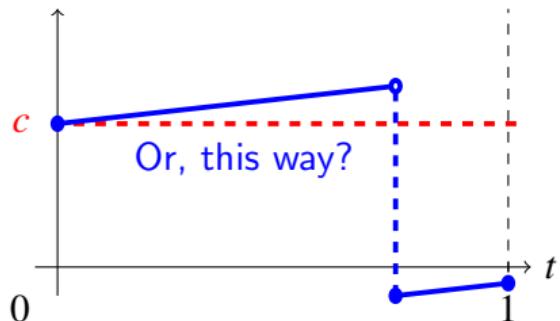
A Rare Event: Bankruptcy

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



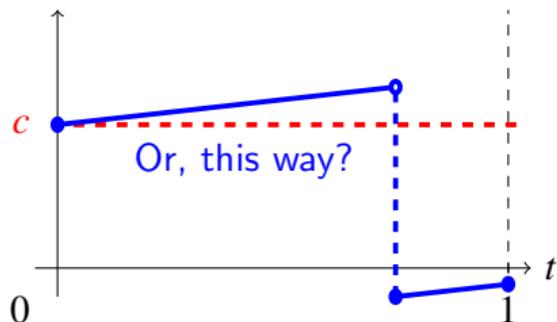
A Rare Event: Bankruptcy

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



A Rare Event: Bankruptcy

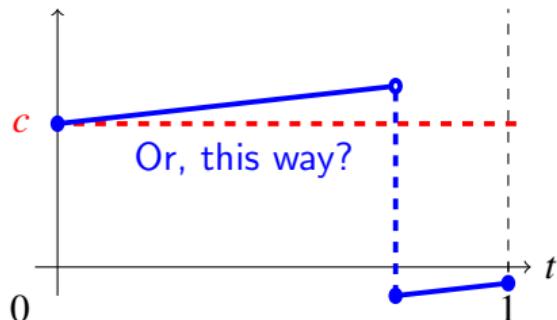
Consider $B \triangleq \{ \bar{Y}_n \text{ falls below } 0 \text{ on } [0, 1] \}$. (i.e., Bankruptcy)



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)

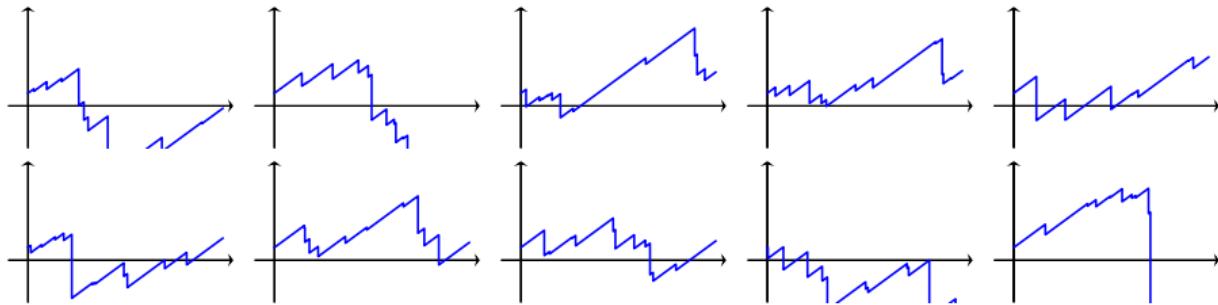


Are we going to see clear patterns?

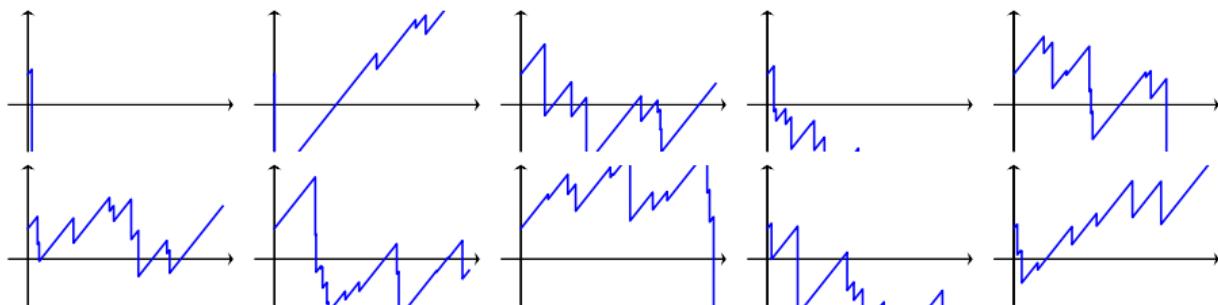
Do they depend on the tail distributions?

A Rare Event: Bankruptcy

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

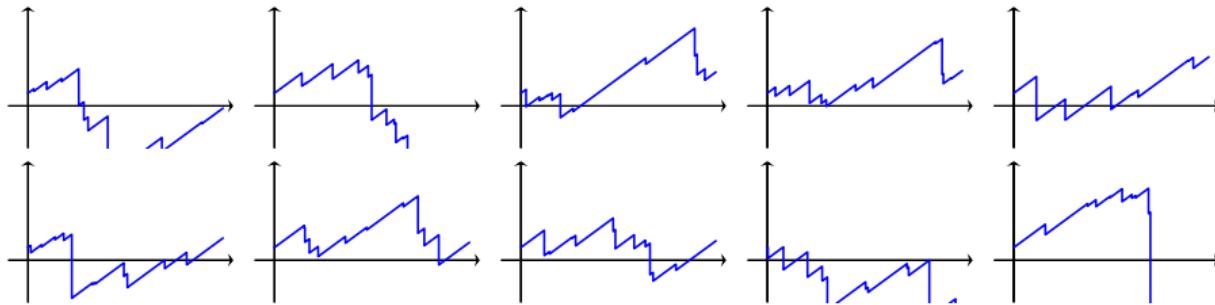


Sample paths of \bar{Y}_{10} conditional on B for **heavy-tailed** claims:

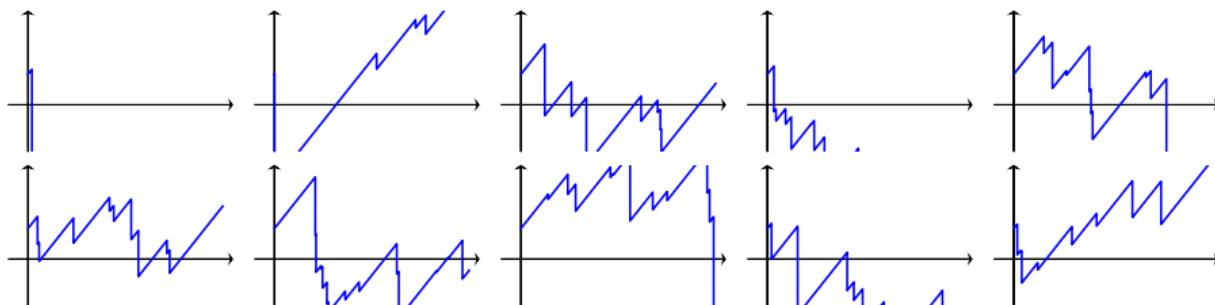


A Rare Event: Bankruptcy

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

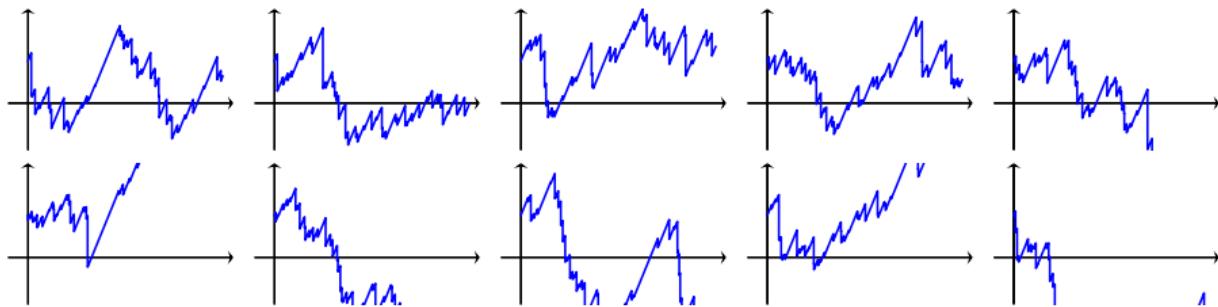


Sample paths of \bar{Y}_{10} conditional on B for **heavy-tailed** claims:

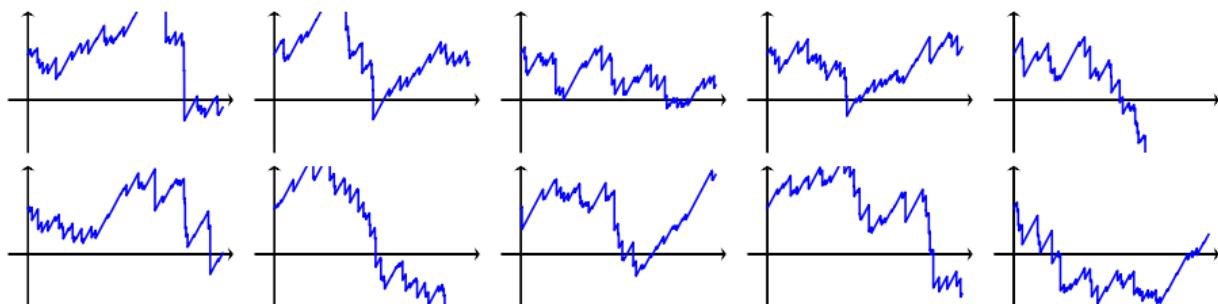


A Rare Event: Bankruptcy

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

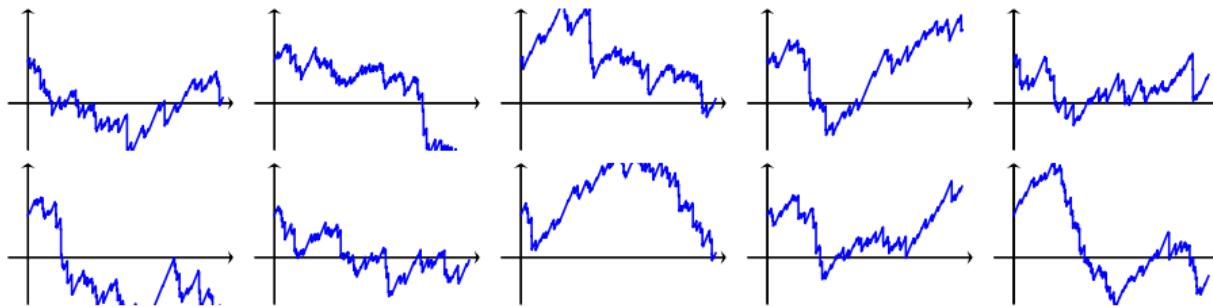


Sample paths of \bar{Y}_{50} conditional on B for **heavy-tailed** claims:

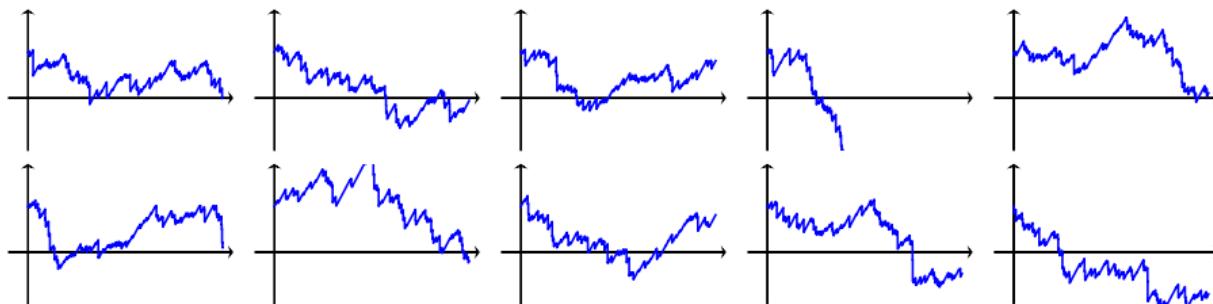


A Rare Event: Bankruptcy

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

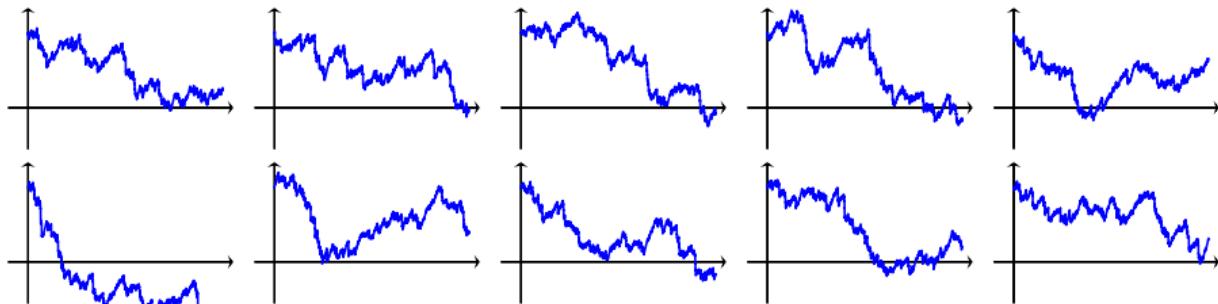


Sample paths of \bar{Y}_{100} conditional on B for **heavy-tailed** claims:



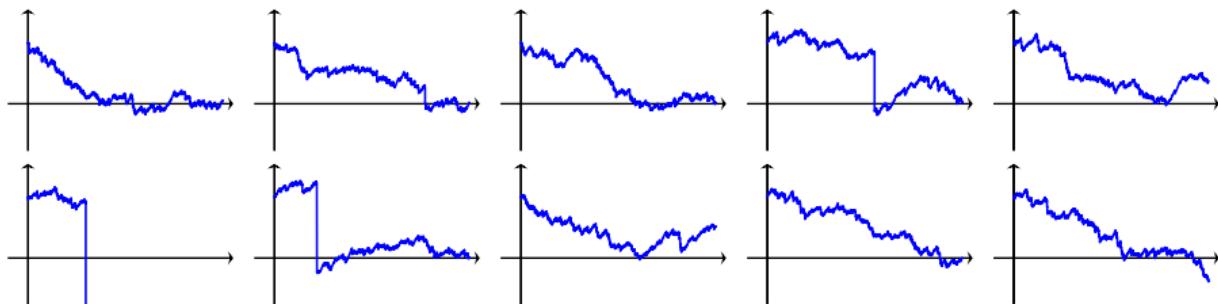
A Rare Event: Bankruptcy

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



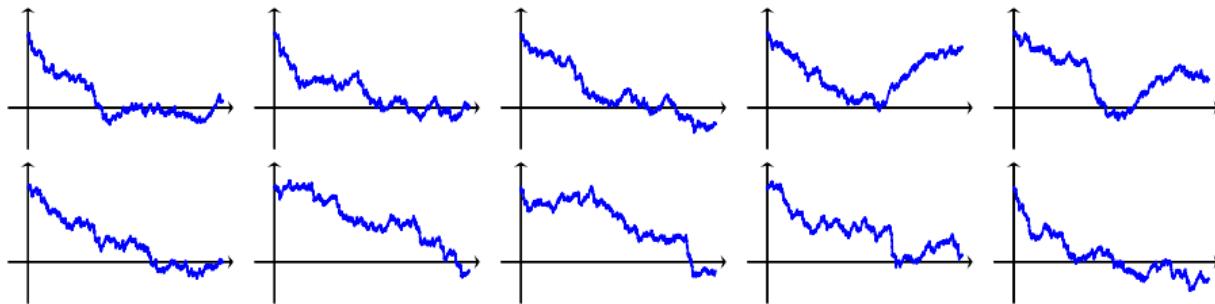
Bankruptcy

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



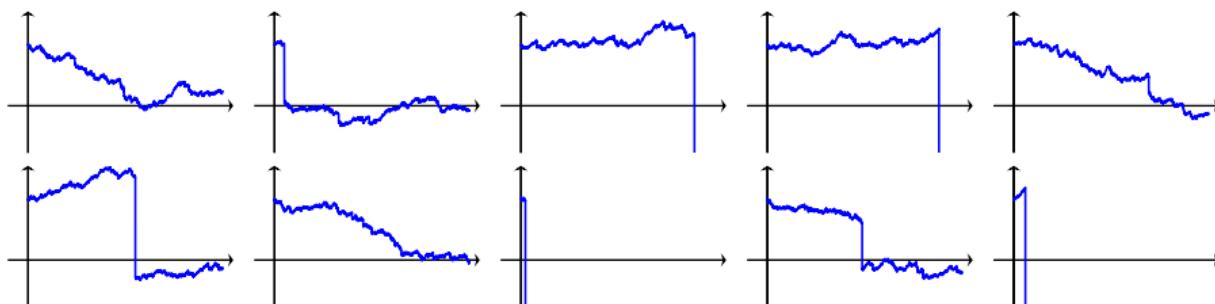
A Rare Event: Bankruptcy

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



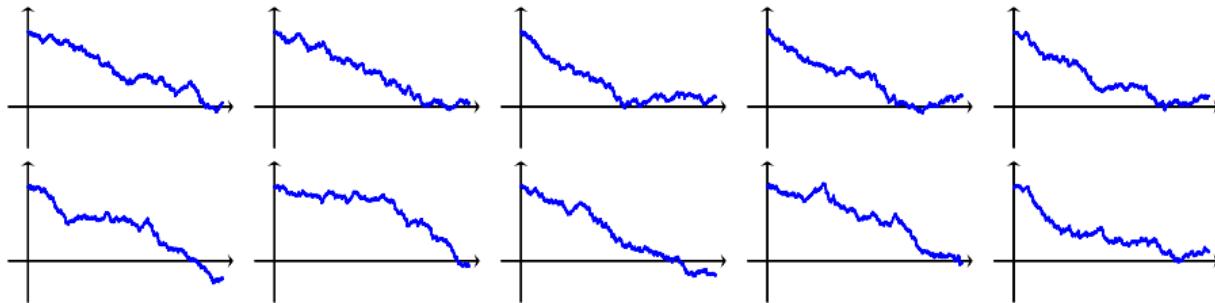
Bankruptcy

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



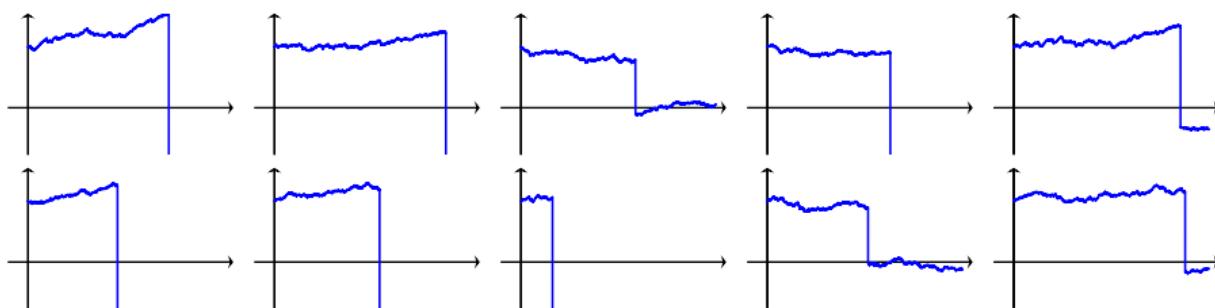
A Rare Event: Bankruptcy

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



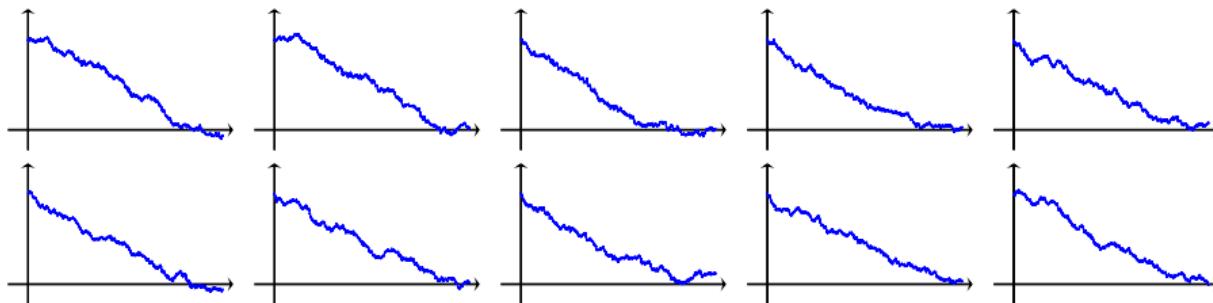
Bankruptcy

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



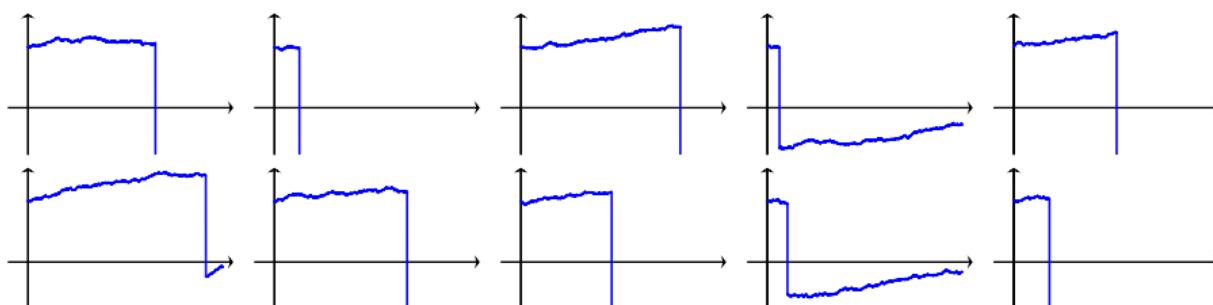
A Rare Event: Bankruptcy

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

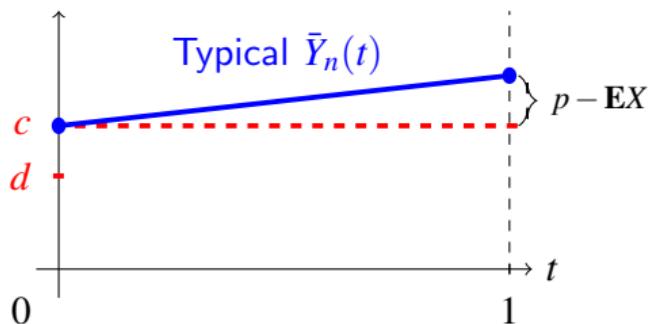


Bankruptcy

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

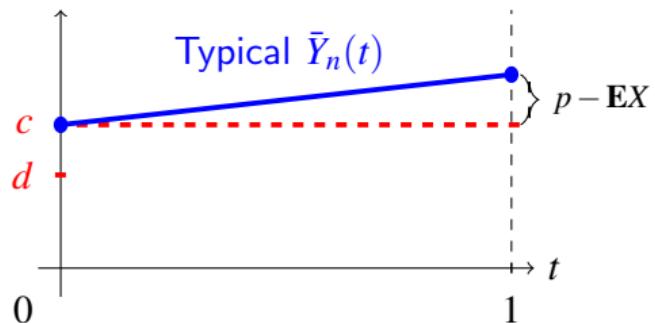


Bankruptcy Despite Reinsurance



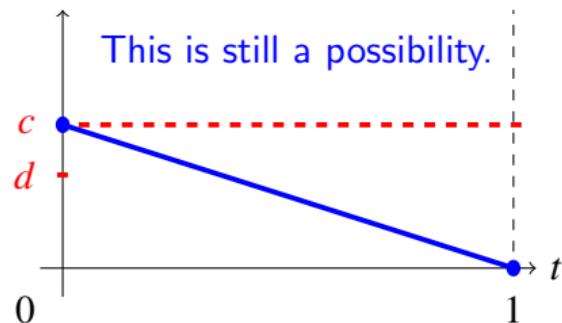
Bankruptcy Despite of Reinsurance

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



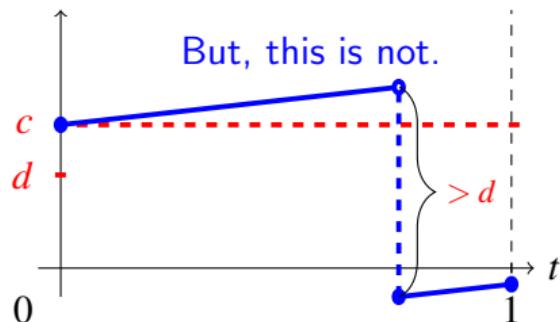
Bankruptcy Despite of Reinsurance

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



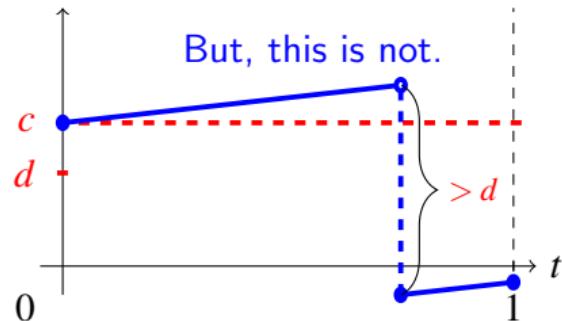
Bankruptcy Despite of Reinsurance

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



Bankruptcy Despite of Reinsurance

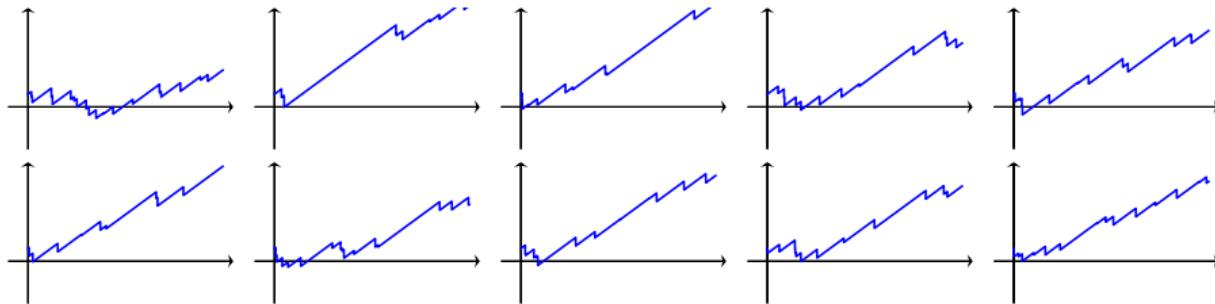
Consider $R \triangleq \{ \bar{Y}_n \text{ falls below } 0 \text{ on } [0,1], \text{ jump sizes } \leq d \}$



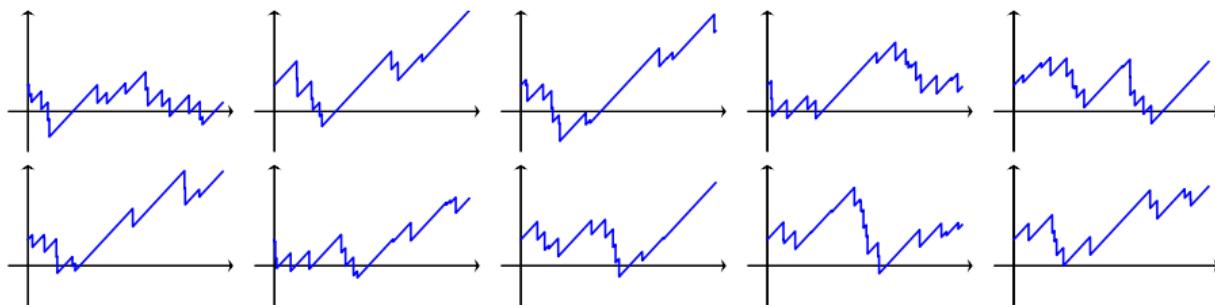
How does the pattern change in this case?

Bankruptcy Despite Reinsurance

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

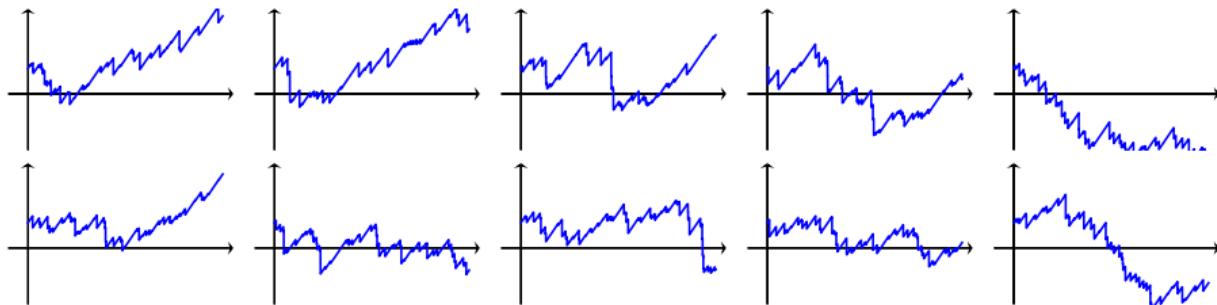


Sample paths of \bar{Y}_{10} conditional on R for **heavy-tailed** claims:

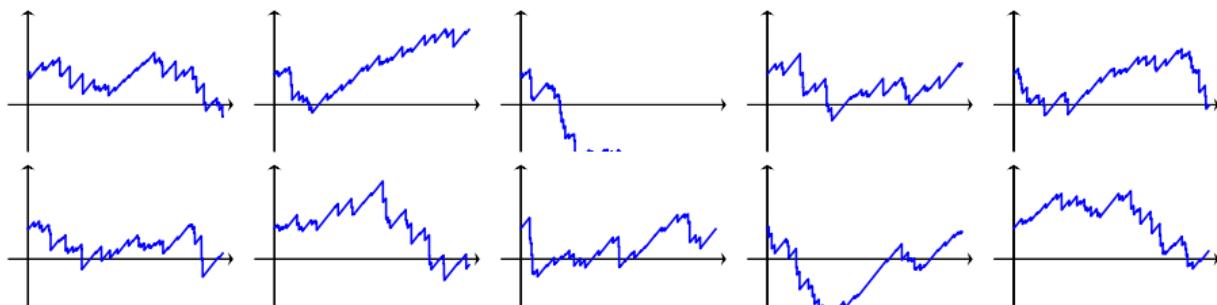


Bankruptcy Despite Reinsurance

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

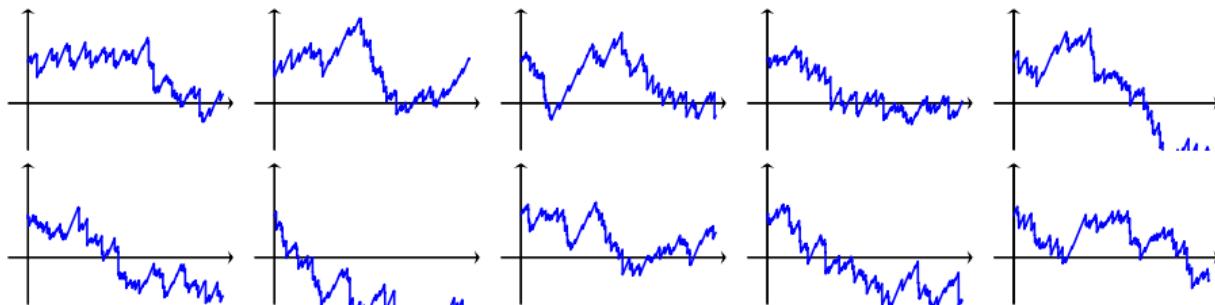


Sample paths of \bar{Y}_{50} conditional on R for **heavy-tailed** claims:

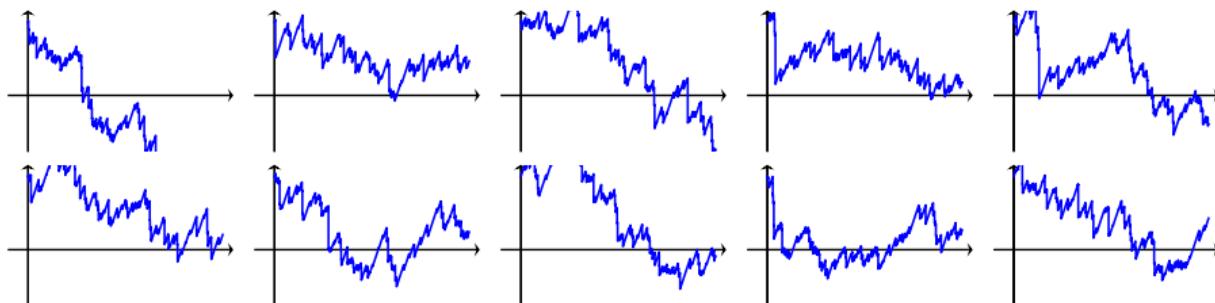


Bankruptcy Despite Reinsurance

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

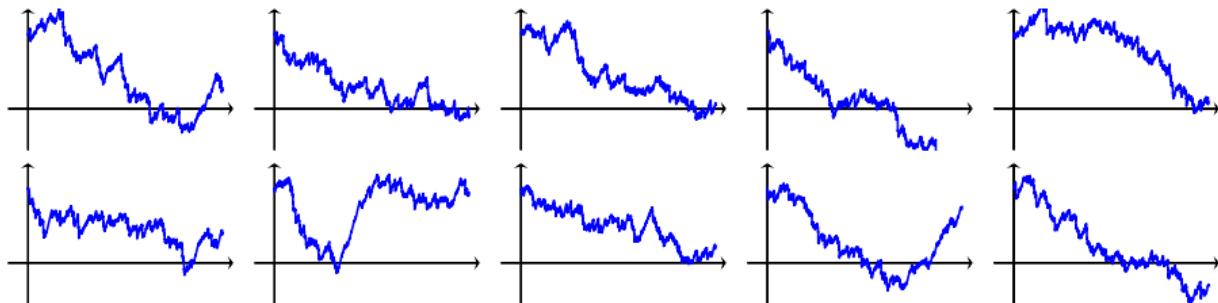


Sample paths of \bar{Y}_{100} conditional on R for **heavy-tailed** claims:

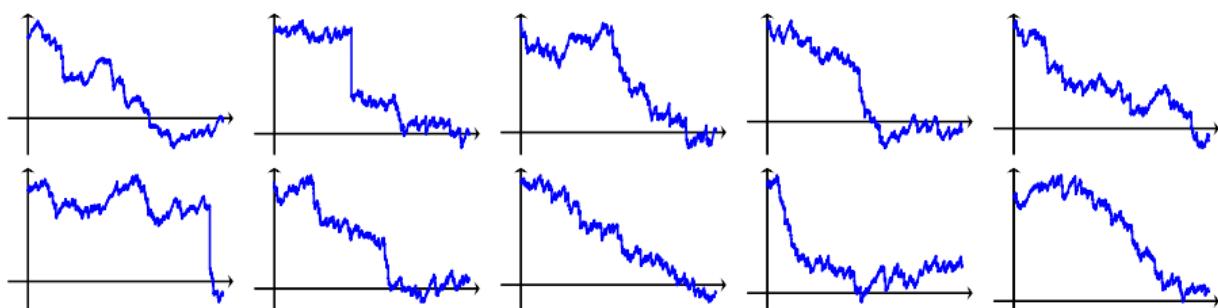


Bankruptcy Despite Reinsurance

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

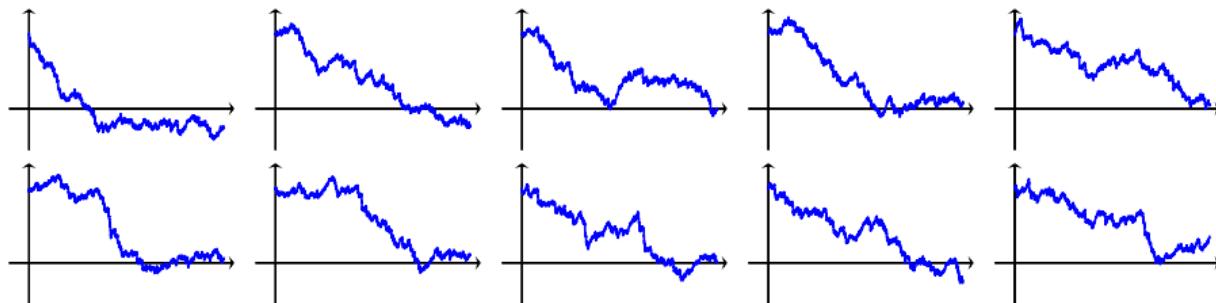


Sample paths of \bar{Y}_{500} conditional on R for **heavy-tailed** claims:

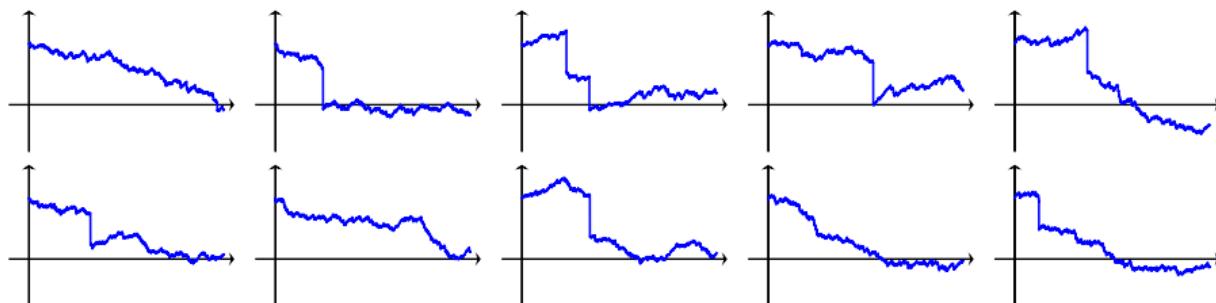


Bankruptcy Despite Reinsurance

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

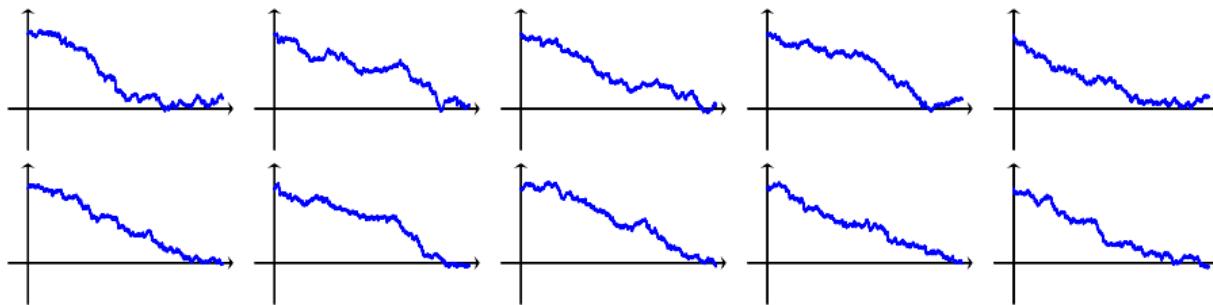


Sample paths of \bar{Y}_{1000} conditional on R for **heavy-tailed** claims:

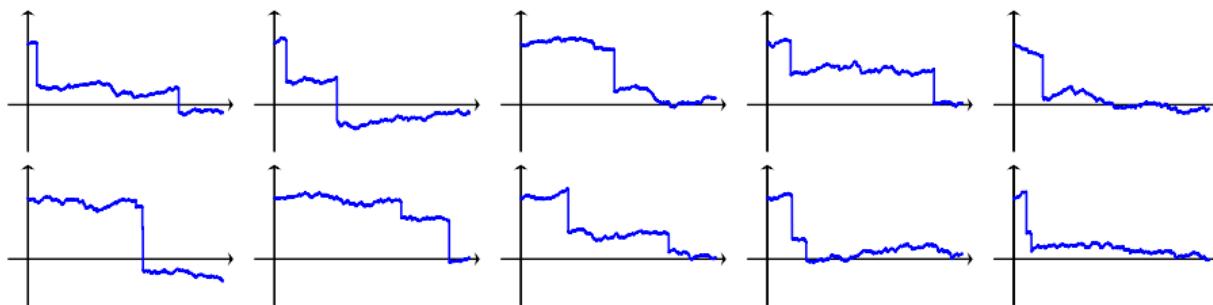


Bankruptcy Despite Reinsurance

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

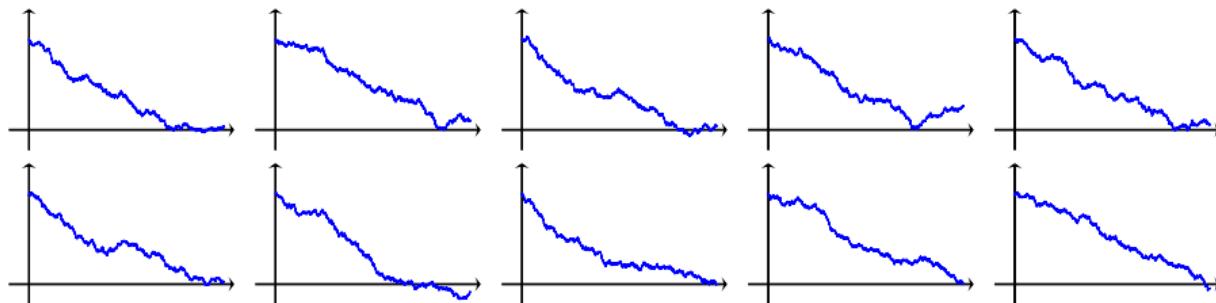


Sample paths of \bar{Y}_{2500} conditional on R for **heavy-tailed** claims:

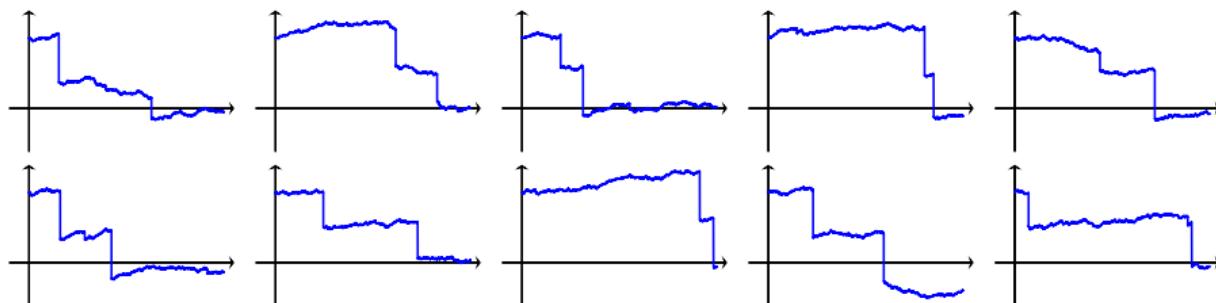


Bankruptcy Despite of Reinsurance

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

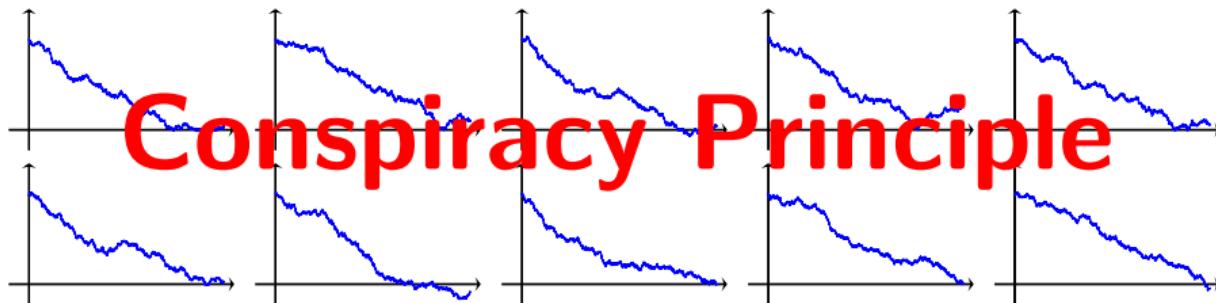


Sample paths of \bar{Y}_{5000} conditional on R for **heavy-tailed** claims:



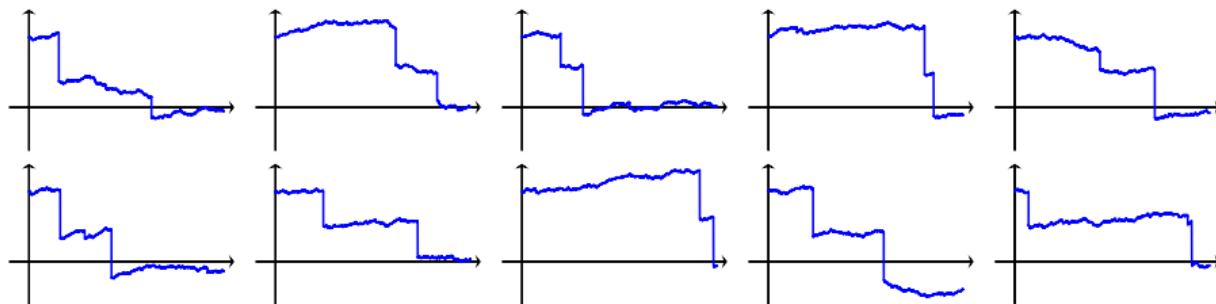
Bankruptcy Despite of Reinsurance

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



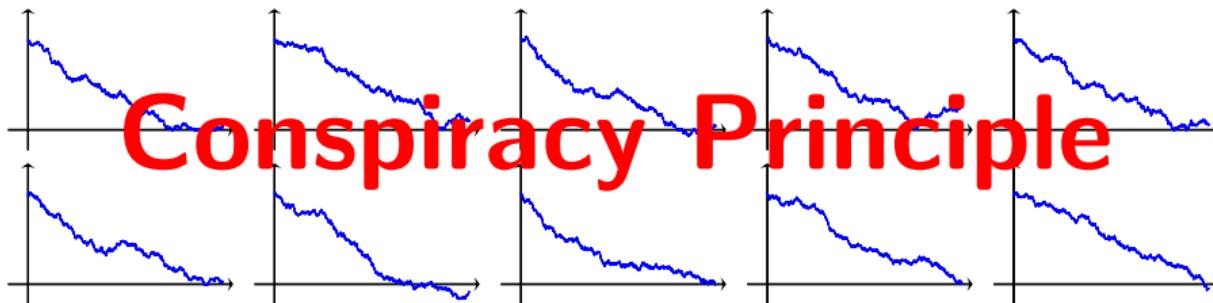
Conspiracy Principle

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

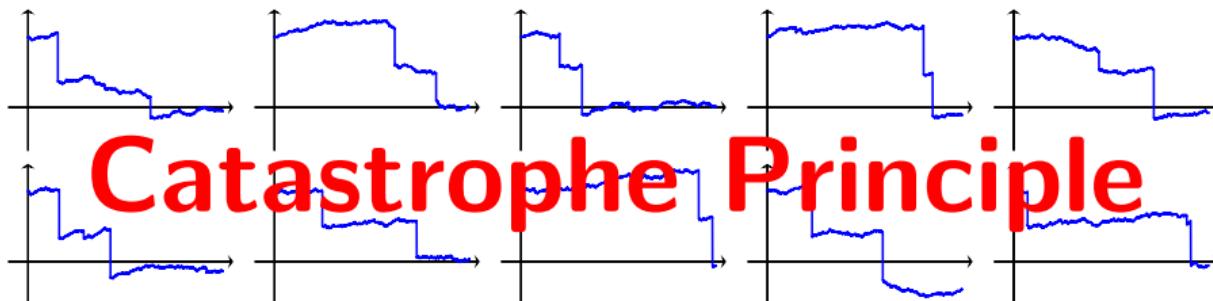


Bankruptcy Despite of Reinsurance

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:



Heavy-Tailed Large Deviations:

Rigorous Characterization of Catastrophe 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

Under certain regularity conditions on A ,

$$\mathcal{L}(\bar{S}_n | \bar{S}_n \in A) \rightarrow \mathcal{L}(\bar{S}_{|A})$$

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

Implication: The Catastrophe Principle

Under certain regularity conditions on A ,

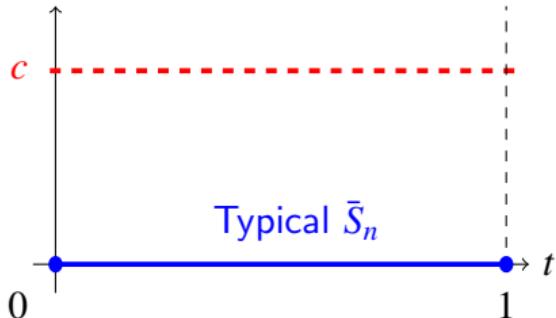
$$\mathcal{L}(\bar{S}_n | \bar{S}_n \in A) \rightarrow \mathcal{L}(\bar{S}_{|A})$$

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

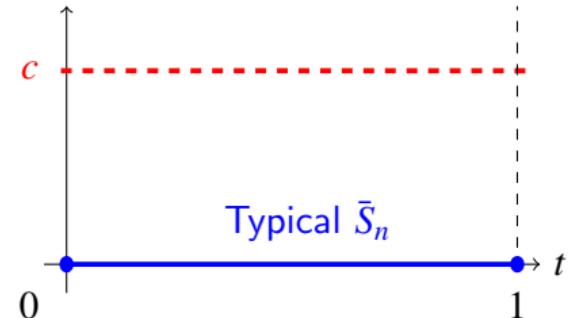
Rigorous Characterization of Catastrophe Principle

Conspiracy vs Catastrophe

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



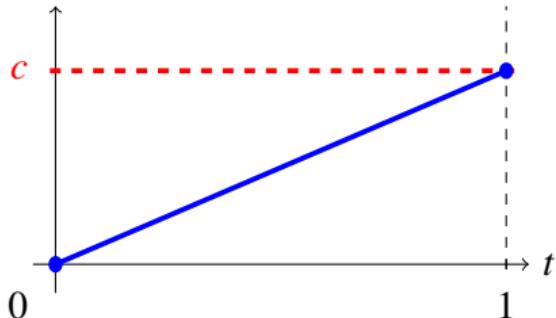
Light-Tailed Claim Size



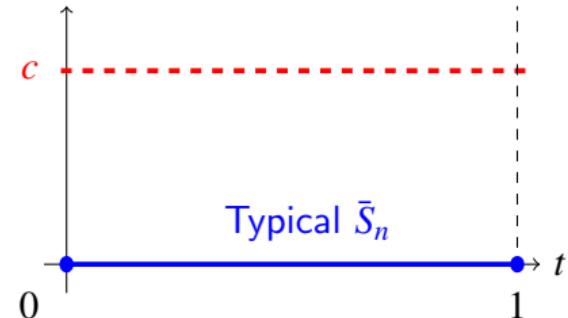
Heavy-Tailed Claim Size

Conspiracy vs Catastrophe

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



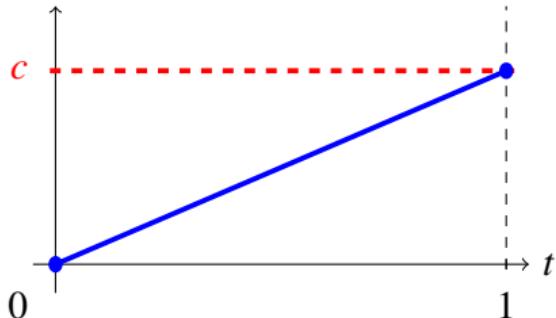
Light-Tailed Claim Size



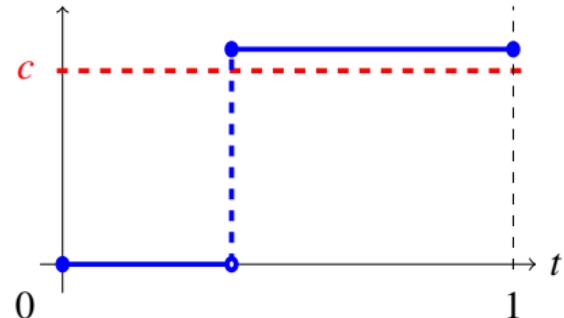
Heavy-Tailed Claim Size

Conspiracy vs Catastrophe

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



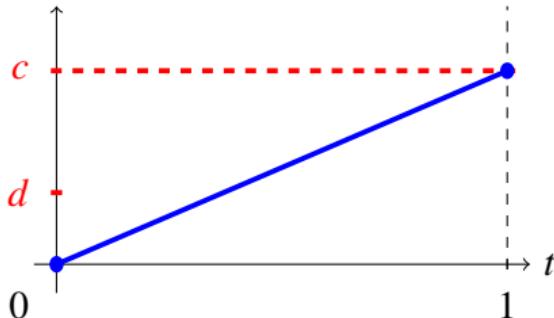
Light-Tailed Claim Size



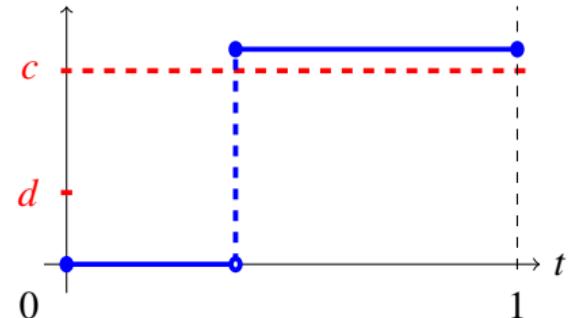
Heavy-Tailed Claim Size

Conspiracy vs Catastrophe

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



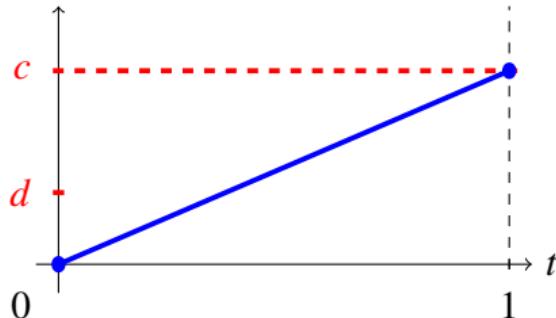
Light-Tailed Claim Size



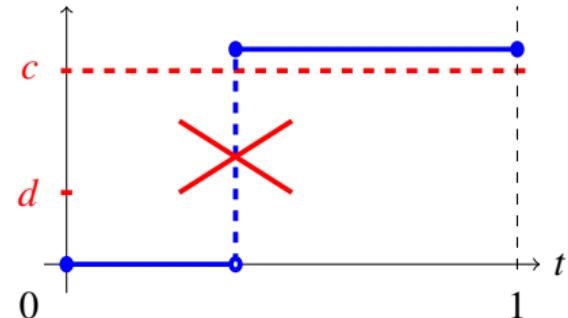
Heavy-Tailed Claim Size

Conspiracy vs Catastrophe

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



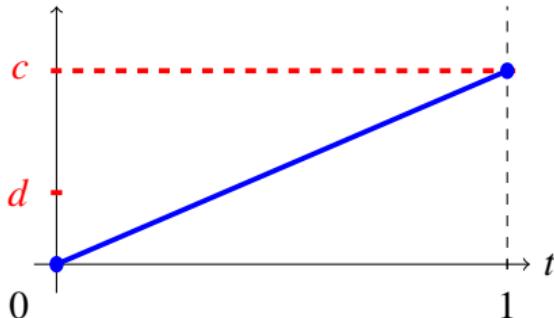
Light-Tailed Claim Size



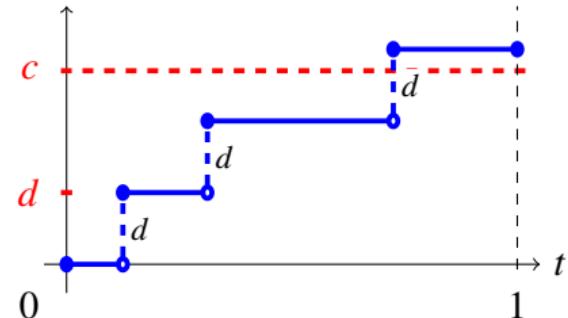
Heavy-Tailed Claim Size

Conspiracy vs Catastrophe

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



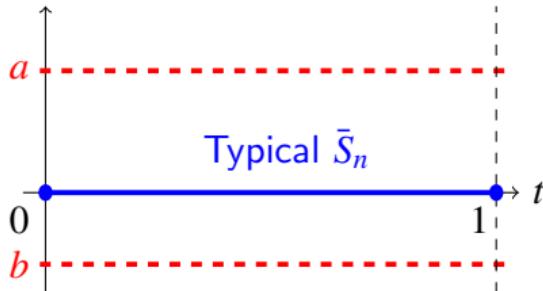
Light-Tailed Claim Size



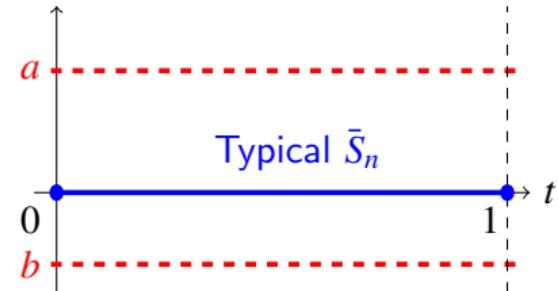
Heavy-Tailed Claim Size

Conspiracy vs Catastrophe

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



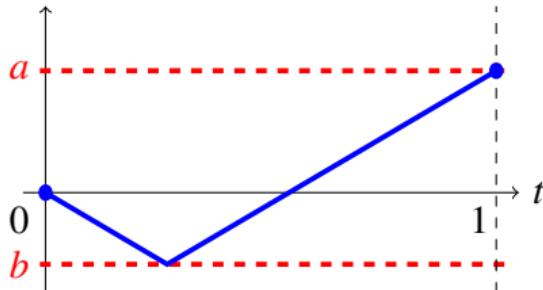
Light-Tailed Increments



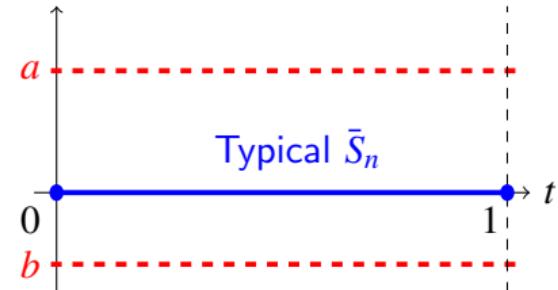
Heavy-Tailed Increments

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



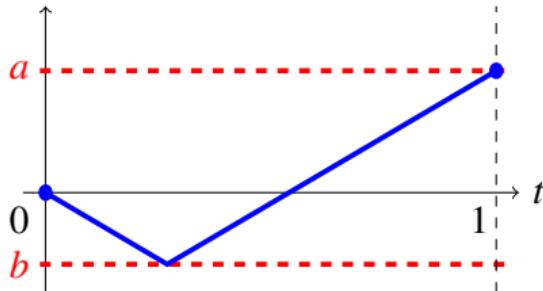
Heavy-Tailed Increments

Conspiracy

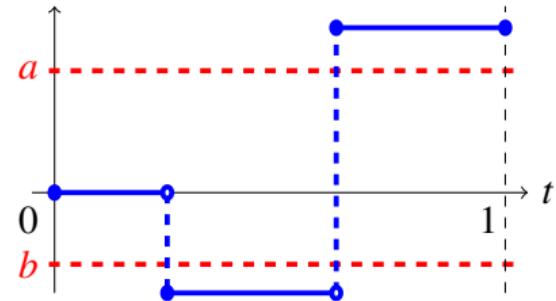
vs

Catastrophe

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



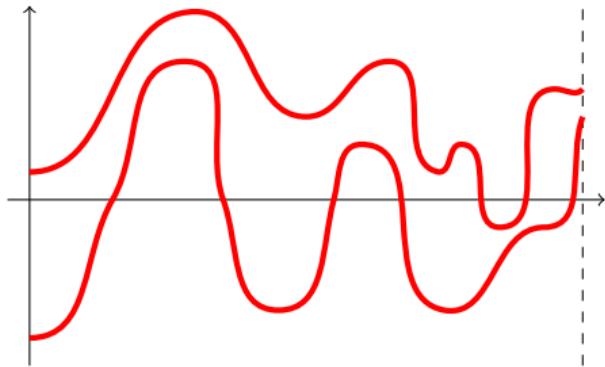
Light-Tailed Increments



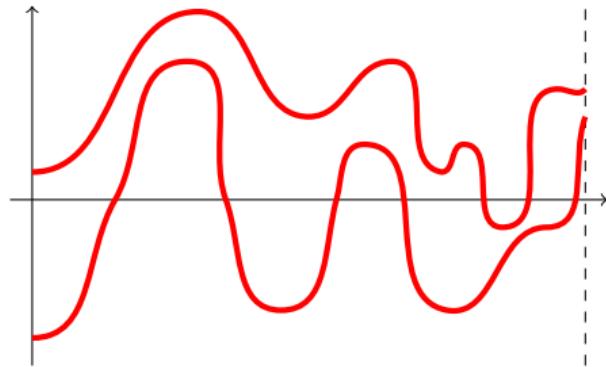
Heavy-Tailed Increments

Conspiracy vs Catastrophe

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



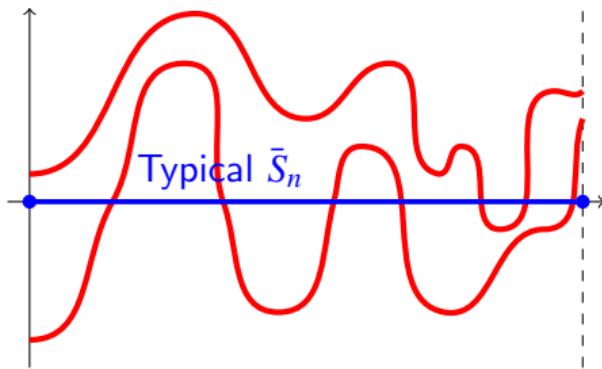
Light-Tailed Increments



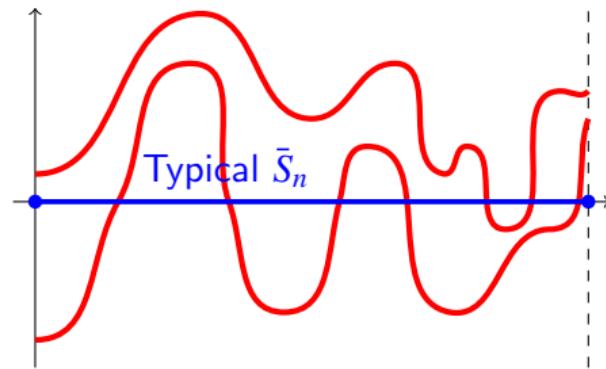
Heavy-Tailed Increments

Conspiracy vs Catastrophe

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



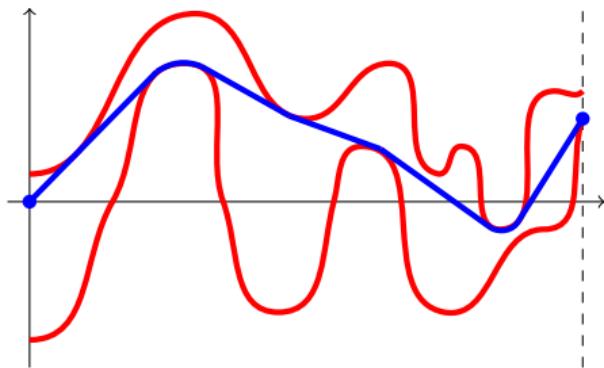
Light-Tailed Increments



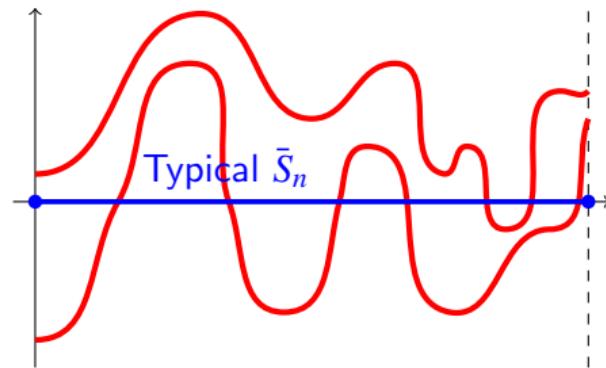
Heavy-Tailed Increments

Conspiracy vs Catastrophe

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



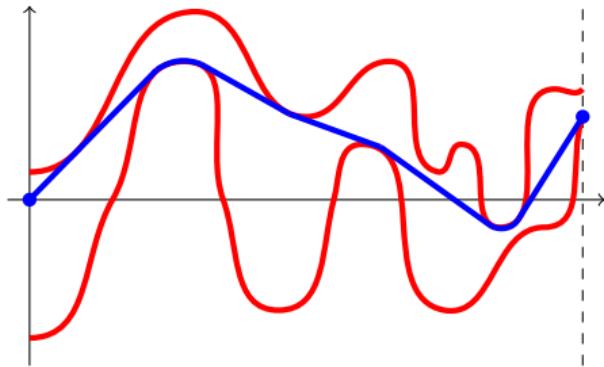
Light-Tailed Increments



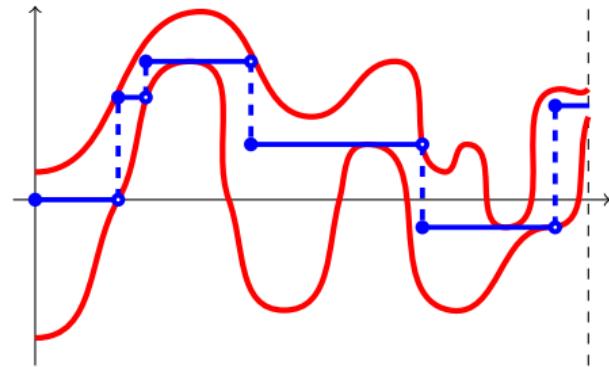
Heavy-Tailed Increments

Conspiracy vs Catastrophe

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



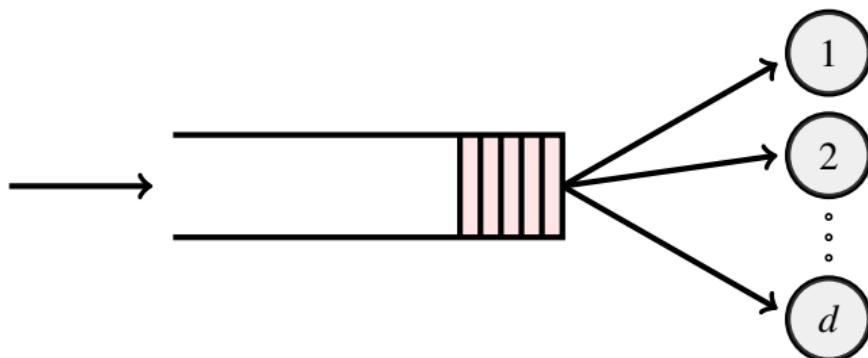
Light-Tailed Increments



Heavy-Tailed Increments

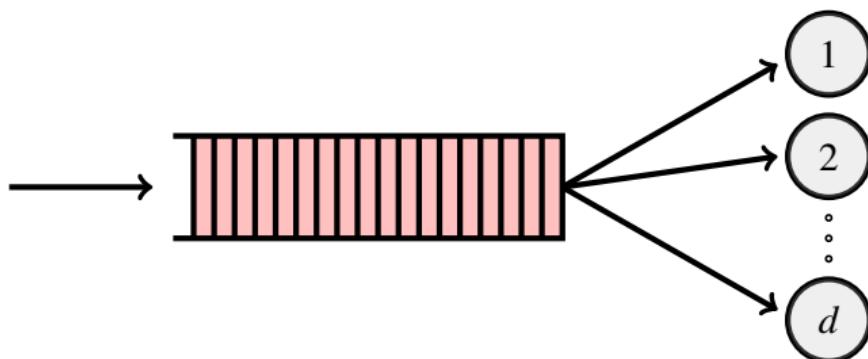
Conspiracy vs Catastrophe

Congestion of Multiple Server Queue:



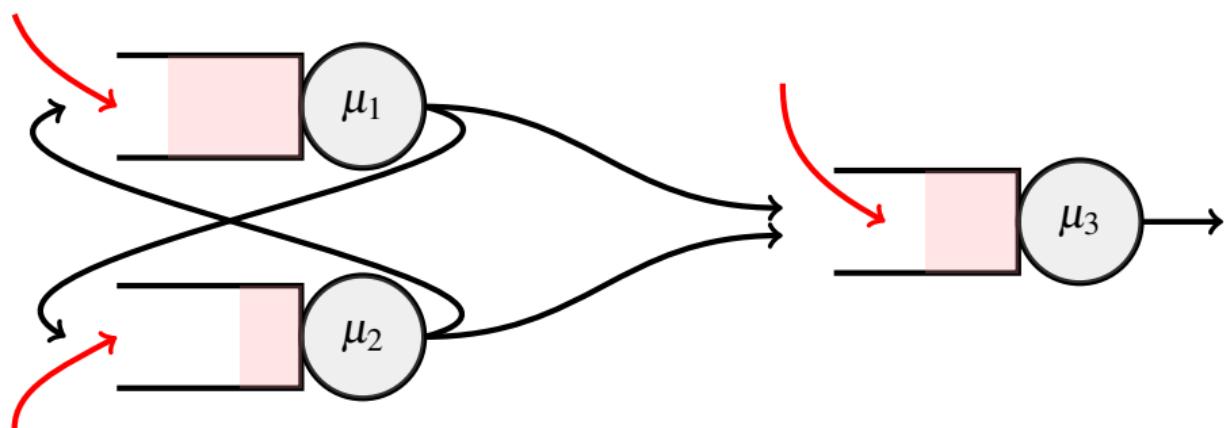
Conspiracy vs Catastrophe

Congestion of Multiple Server Queue:



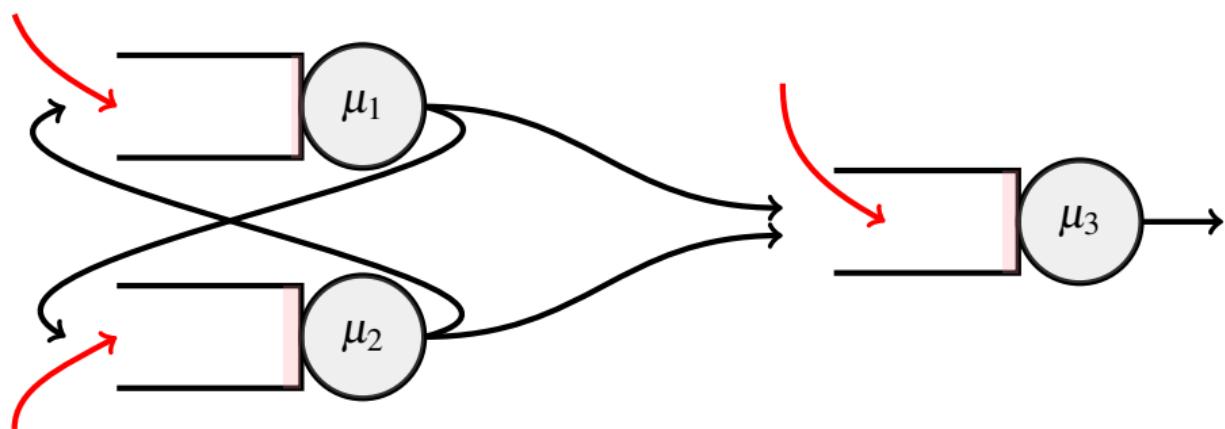
Conspiracy vs Catastrophe

Congestion of Stochastic Fluid Networks:



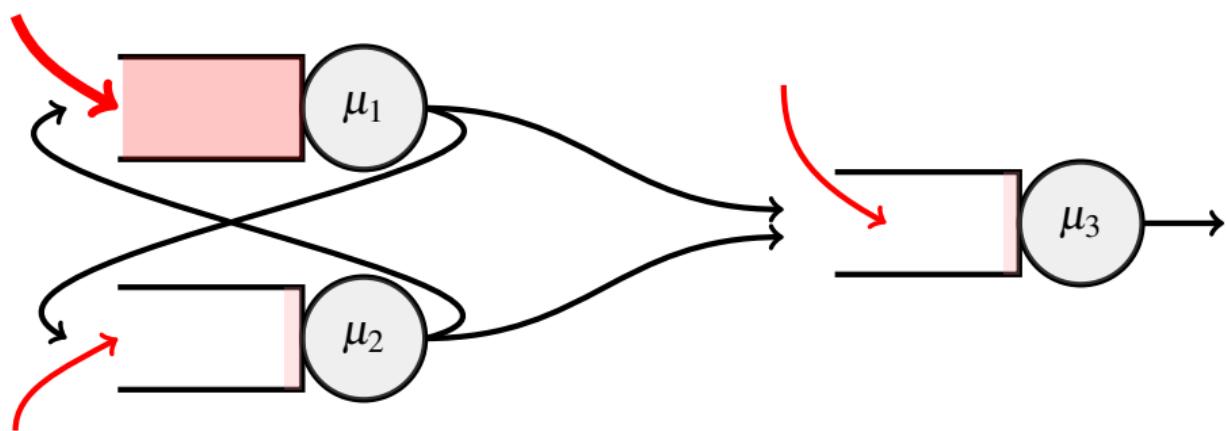
Conspiracy vs Catastrophe

Congestion of Stochastic Fluid Networks:



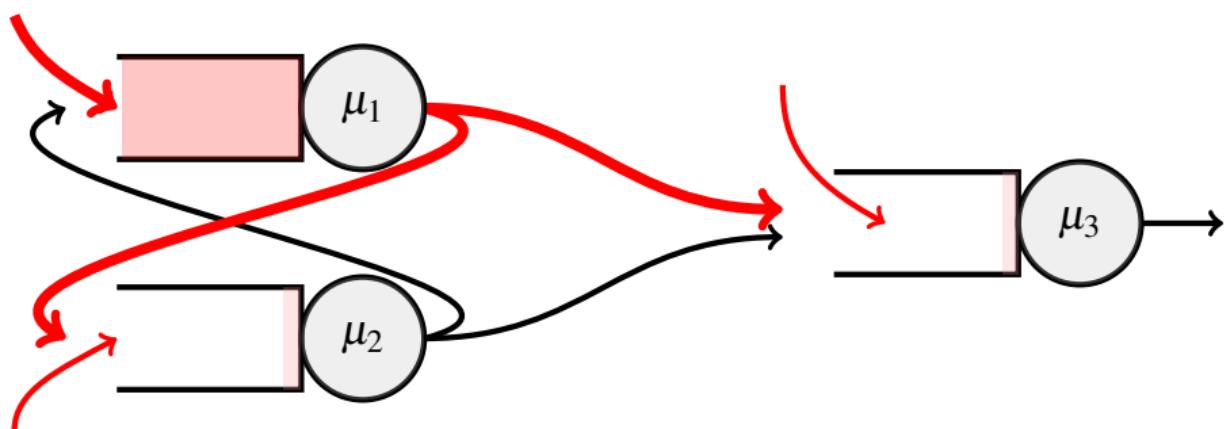
Conspiracy vs Catastrophe

Congestion of Stochastic Fluid Networks:



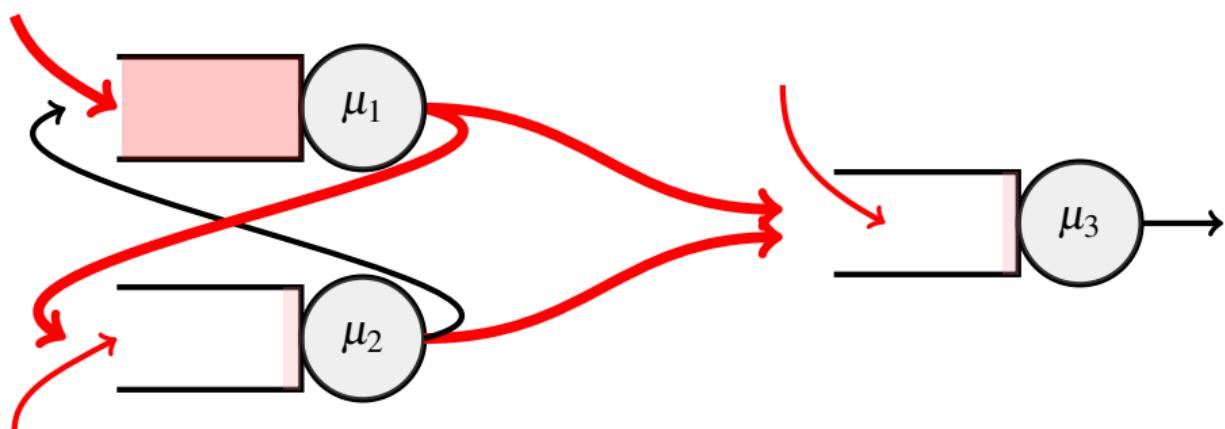
Conspiracy vs Catastrophe

Congestion of Stochastic Fluid Networks:



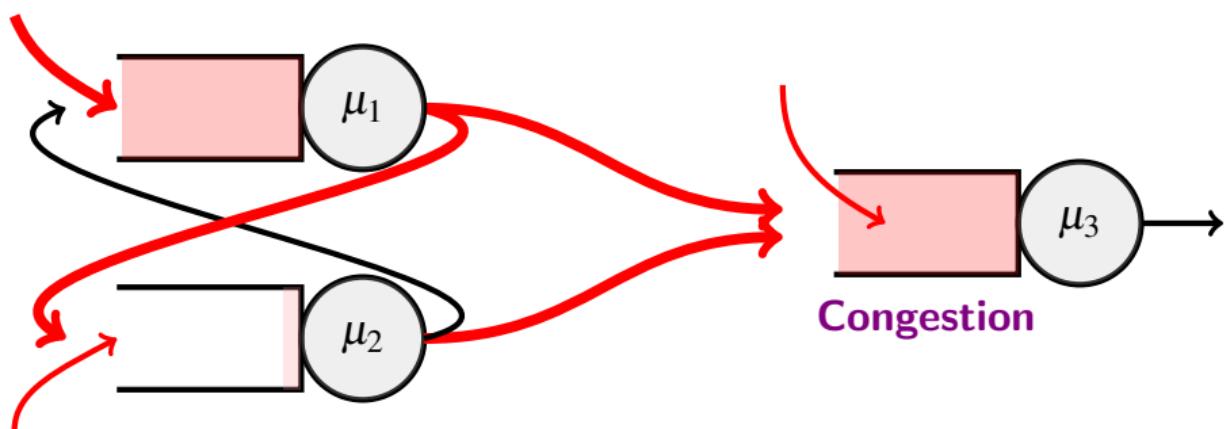
Conspiracy vs Catastrophe

Congestion of Stochastic Fluid Networks:



Conspiracy vs Catastrophe

Congestion of Stochastic Fluid Networks:



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), Su, Wang, R. (2023+)

- Processes with Spatial and Temporal Correlations

Chen, R., Zwart (2023+), Bazhba, Blanchet, R., Zwart (2023+), Su, R. (2023+), Wang, R. (2023+)

Minimal # Jumps added to Typical Paths Characterize the Catastrophe Principle

Open Problems Solved by Heavy-Tailed Large Deviations

- Design of strongly efficient rare-event simulation algorithm
- Multiple server queue length asymptotics
- Ruin probability under the presence of reinsurance contracts

From Large Deviations to Metastability

Stochastic Gradient Descent (SGD)

Minimizing loss function f :

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

Stochastic Gradient Descent (SGD)

Minimizing loss function f :

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

Stochastic Gradient Descent (SGD)

Minimizing loss function f :

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

Stochastic Gradient Descent (SGD)

Minimizing loss function f :

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

Stochastic Gradient Descent (SGD)

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 (SGD)

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 (SGD)

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 (SGD)

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 (SGD)

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

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

Heavy-Tailed Large Deviations for SGD

Theorem (Wang, R., 2023+)

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., 2023+)

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., 2023+)

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., 2023+)

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., 2023+)

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., 2023+)

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., 2023+)

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., 2023+)

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.

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

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

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

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.

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

$\{I(\varepsilon) \subseteq I : \varepsilon > 0\}$: covering of I

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

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

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.

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

$\{I(\varepsilon) \subseteq I : \varepsilon > 0\}$: covering of I

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

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

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.

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

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

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

Also,

$$\limsup_{\eta \downarrow 0} \frac{\sup_{x \in I(\varepsilon)} \mathbf{P}(\tau_{(I(\varepsilon) \setminus A(\varepsilon))^C}^\eta(x) > T/\eta)}{\gamma(\eta)T/\eta} = 0;$$

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.

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

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

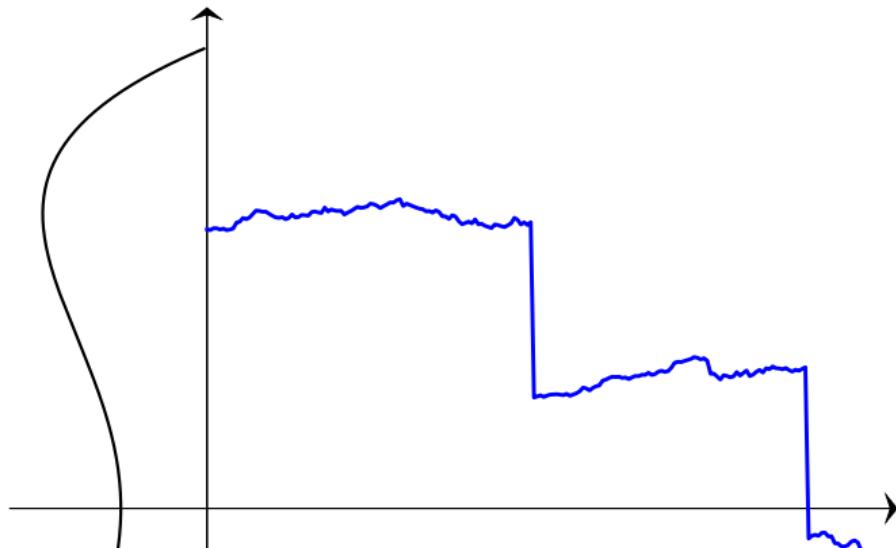
$$\lim_{\varepsilon \rightarrow 0} \lim_{T \rightarrow \infty} \delta_B(\varepsilon, T) = 0.$$

Also,

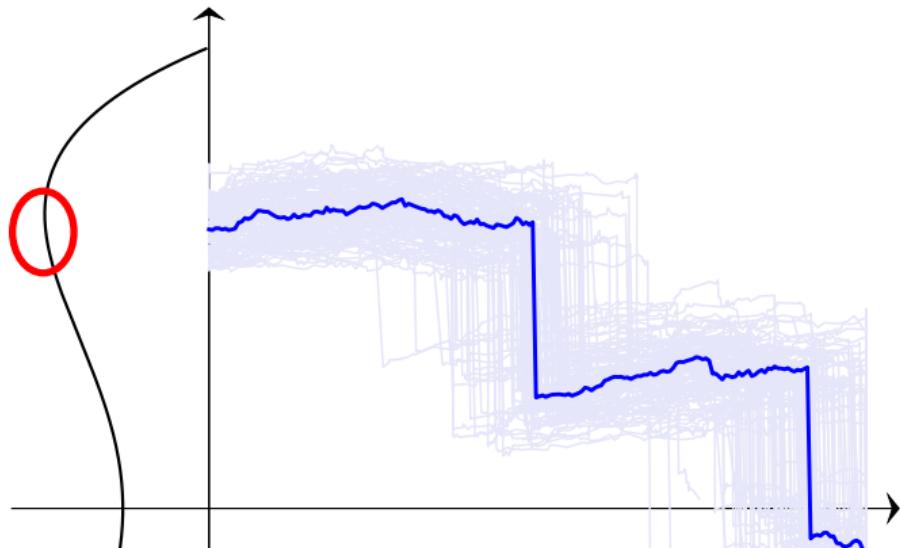
$$\limsup_{\eta \downarrow 0} \frac{\sup_{x \in I(\varepsilon)} \mathbf{P}(\tau_{(I(\varepsilon) \setminus A(\varepsilon))^c}^\eta(x) > T/\eta)}{\gamma(\eta)T/\eta} = 0;$$

$$\liminf_{\eta \downarrow 0} \inf_{x \in I(\varepsilon)} \mathbf{P}(\tau_{A(\varepsilon)}^\eta(x) \leq T/\eta) = 1.$$

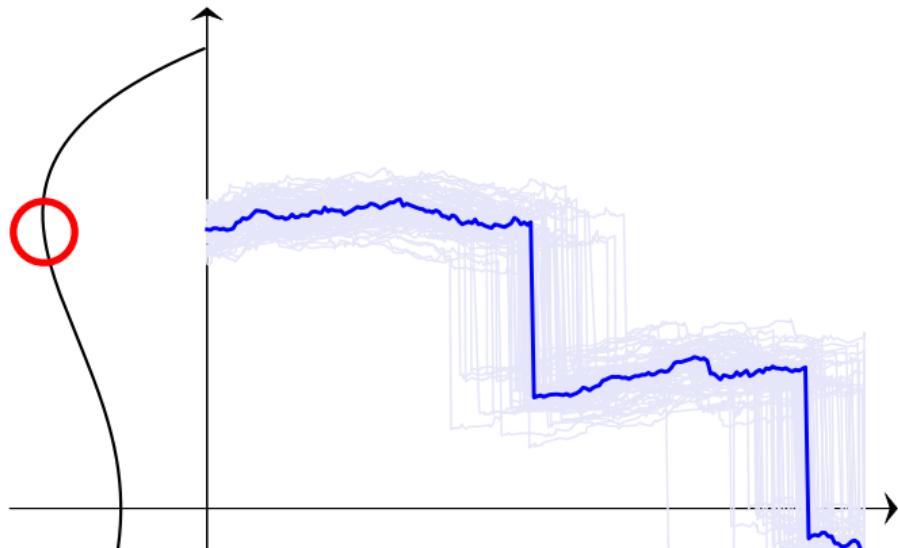
Asymptotic Atoms



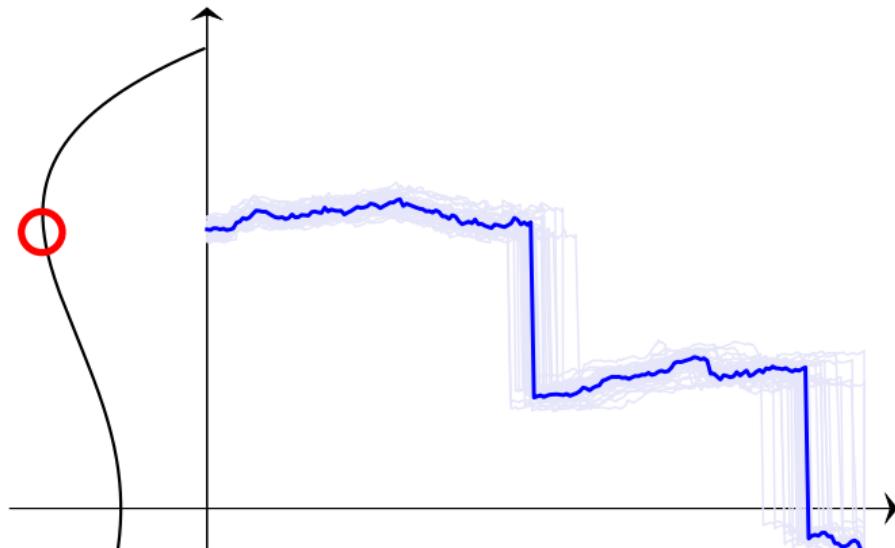
Asymptotic Atoms



Asymptotic Atoms



Asymptotic Atoms



Exit Time and Location under the Presence of Asymptotic Atom

Theorem (Wang, R., 2023+)

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

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

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

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

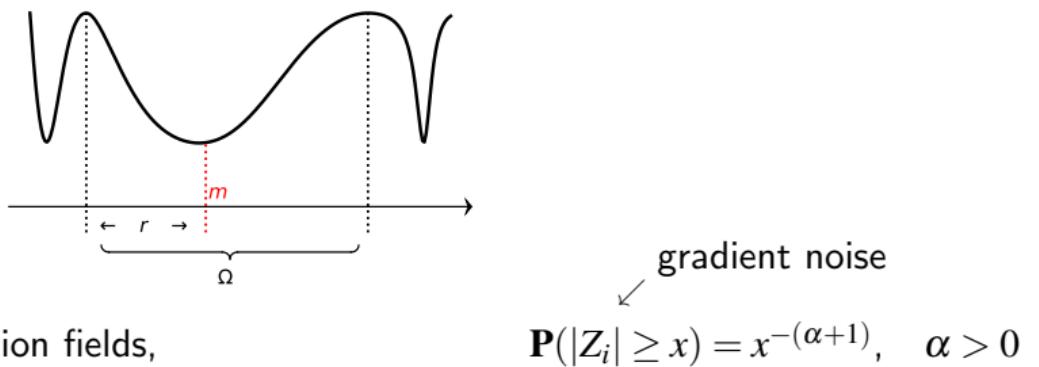
Then, again,

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

where

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

First Exit Time Analysis

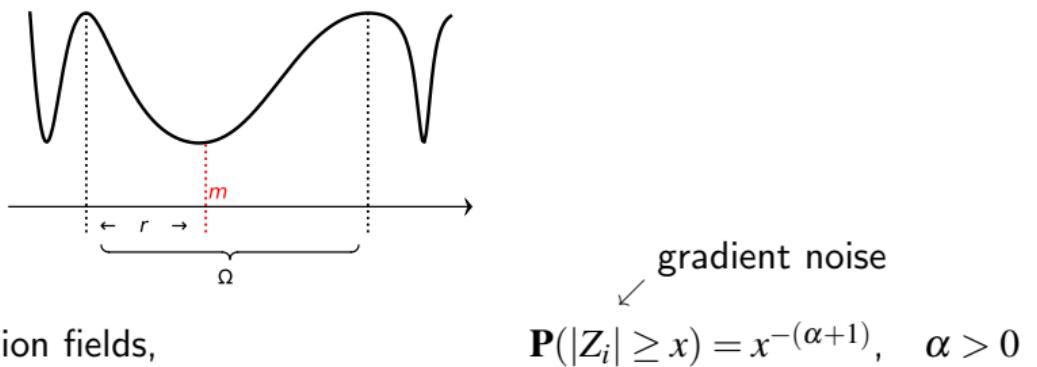


Theorem (Wang, R., 2023+)

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

$$\tau(\eta)\gamma(\eta) \Rightarrow \text{Exp}(1)$$

First Exit Time Analysis



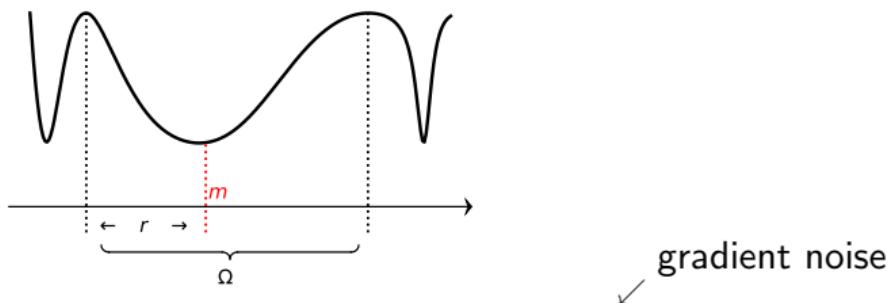
Theorem (Wang, R., 2023+)

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

First Exit Time

$$\tau(\eta)\gamma(\eta) \Rightarrow \text{Exp}(1)$$

First Exit Time Analysis



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

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

Theorem (Wang, R., 2023+)

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

First Exit Time

$$\tau(\eta)\gamma(\eta) \Rightarrow \text{Exp}(1)$$

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

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, R., 2023+)

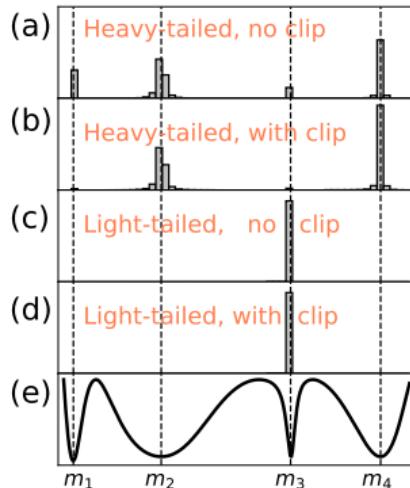
Under certain structural conditions, for any $t > 0$ and $\beta > 1 + \alpha \cdot l^*$,

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

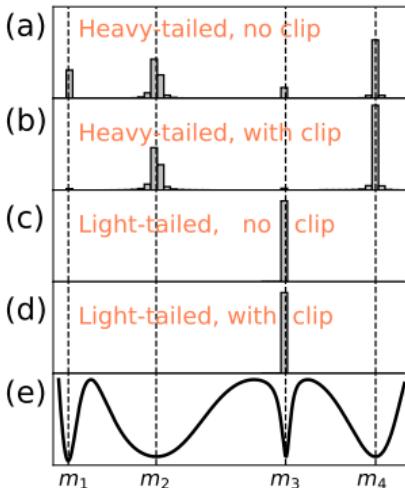
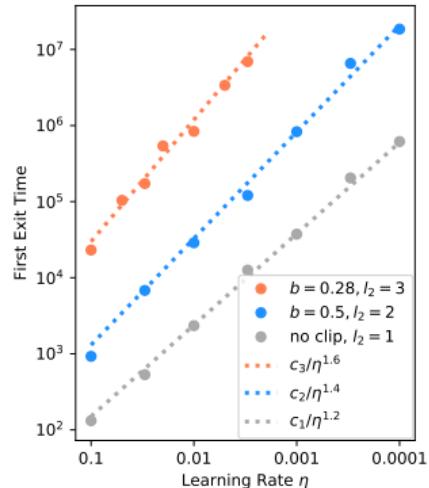
In fact, $W_{\lfloor t/\eta^{1+\alpha \cdot l^*} \rfloor}^\eta$ converges to a Markov jump processes

whose state space consists of wide local minima only.

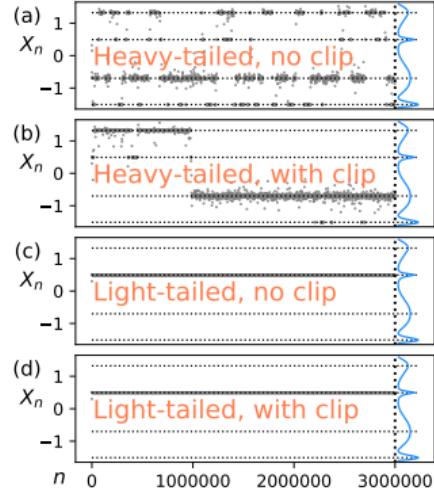
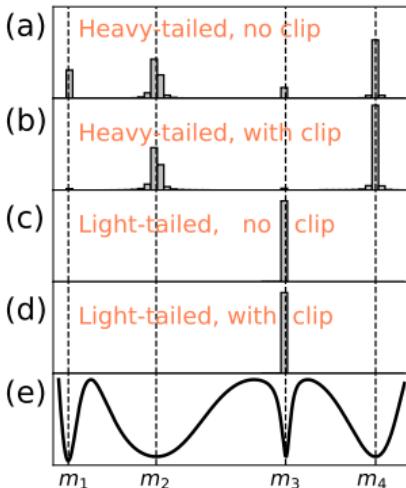
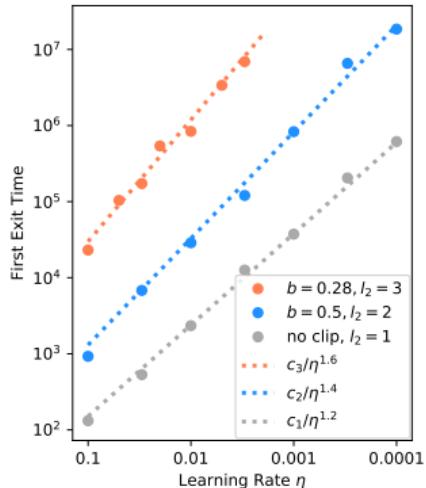
Eliminating Sharp Local Minima with Truncated Heavy-Tails



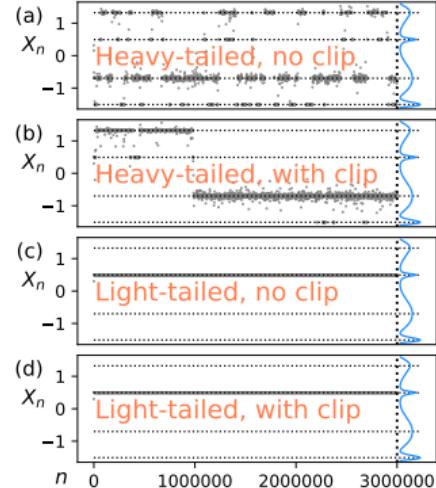
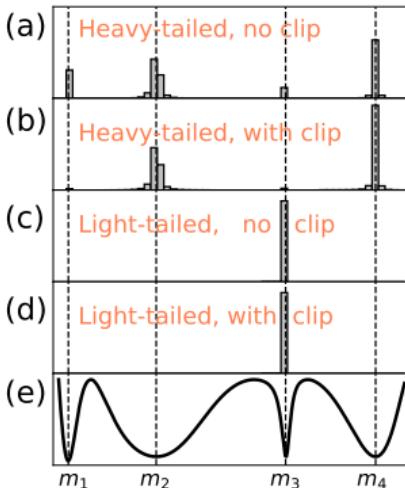
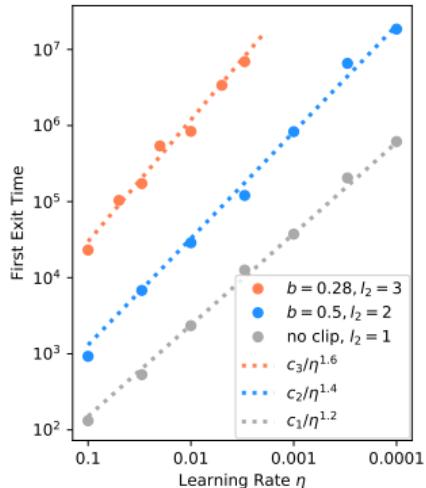
Eliminating Sharp Local Minima with Truncated Heavy-Tails



Eliminating Sharp Local Minima with Truncated Heavy-Tails

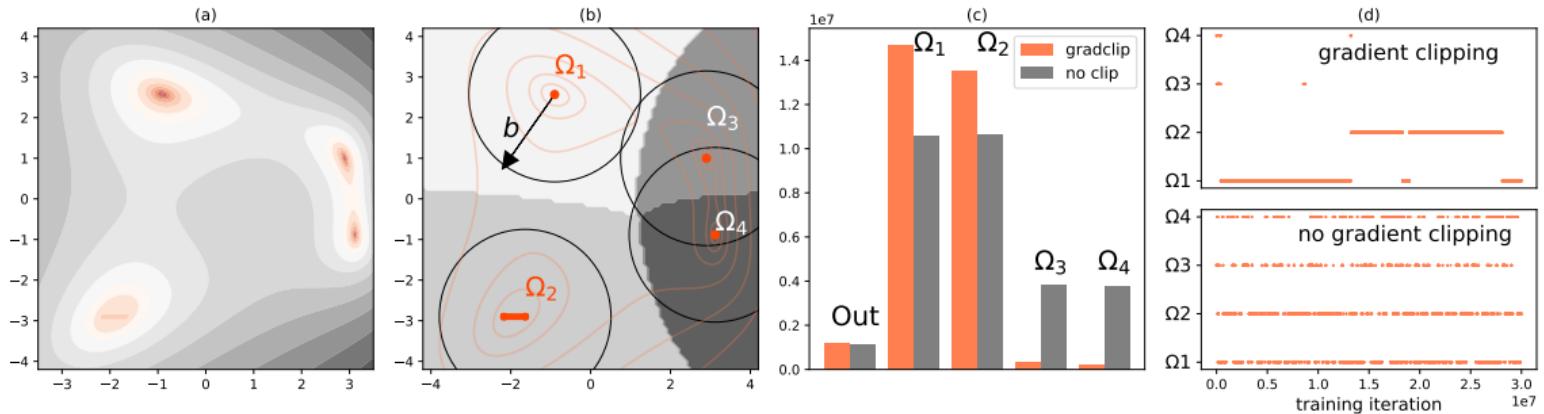


Eliminating Sharp Local Minima with Truncated Heavy-Tails

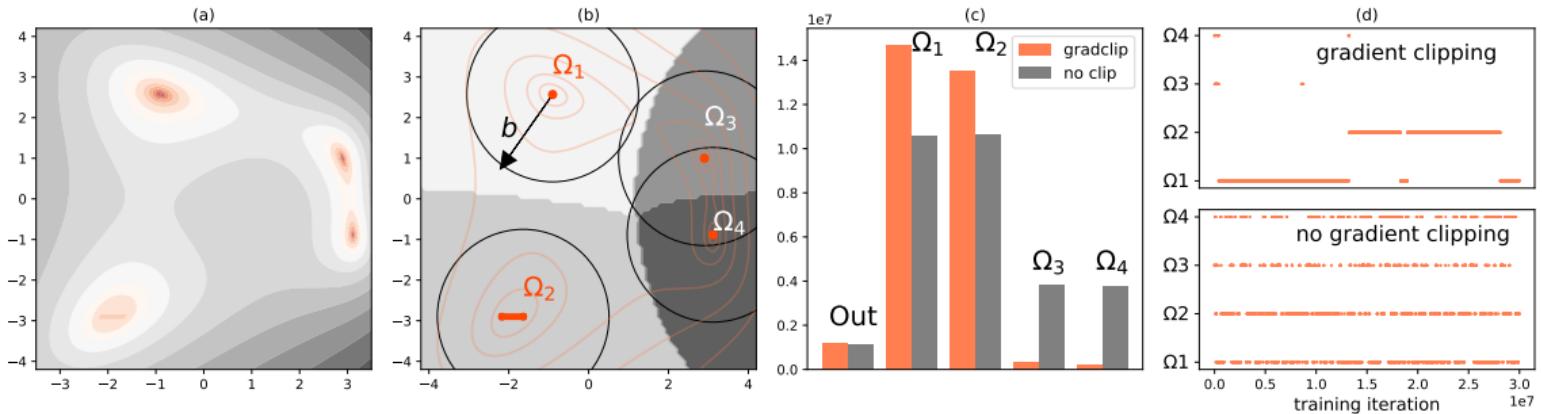


Consistent with what theory predicts!

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!

In Deep Neural Network

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

In Deep Neural Network

Stochastic Gradient

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

In Deep Neural Network

Stochastic Gradient

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

In Deep Neural Network

Stochastic Gradient

Pareto RV

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

In Deep Neural Network

Heavy-Tailed Stochastic Gradient

Pareto RV

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

In Deep Neural Network

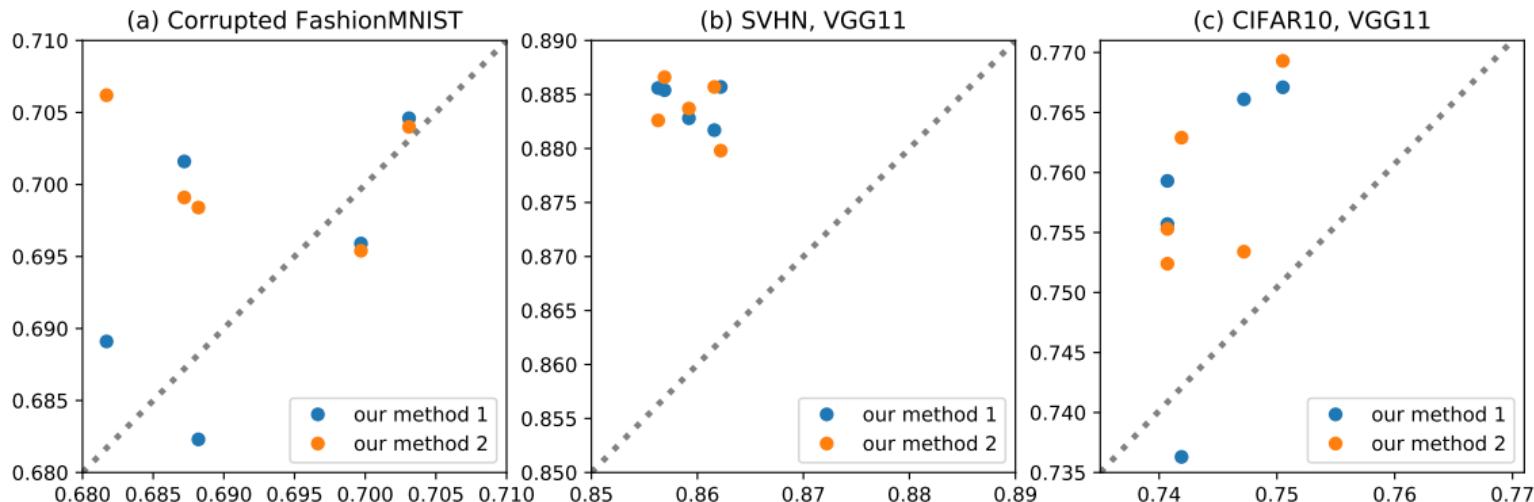
Heavy-Tailed Stochastic Gradient

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



Summary

Summary

- Heavy-tailed large deviations crisply characterize Catastrophe Principle

Summary

- Heavy-tailed large deviations crisply characterize Catastrophe Principle
- New locally uniform large deviations formulation

Summary

- Heavy-tailed large deviations crisply characterize Catastrophe Principle
- New locally uniform large deviations formulation
- Uniform \mathbb{M} -convergence, asymptotic atom, portmanteau theorem

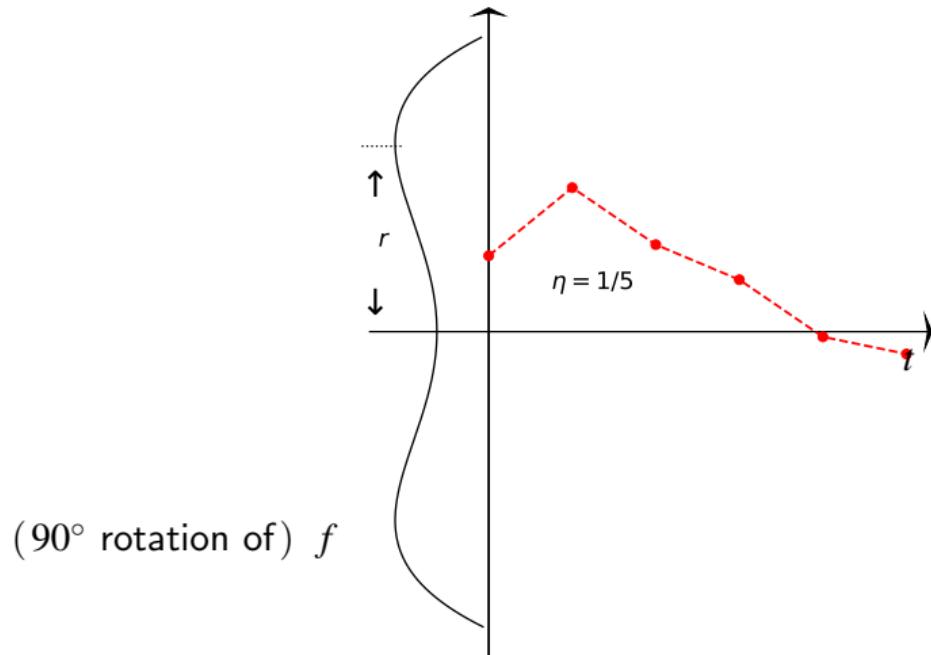
Summary

- Heavy-tailed large deviations crisply characterize Catastrophe Principle
- New locally uniform large deviations formulation
- Uniform \mathbb{M} -convergence, asymptotic atom, portmanteau theorem
- These tools reveal the global dynamics of heavy-tailed dynamical systems

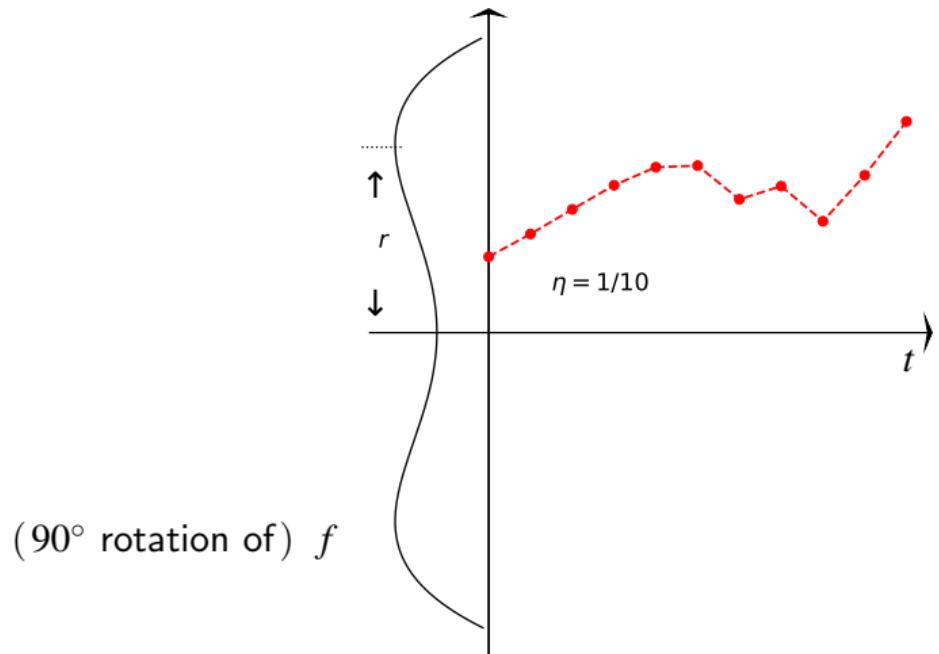
Summary

- Heavy-tailed large deviations crisply characterize Catastrophe Principle
- New locally uniform large deviations formulation
- Uniform \mathbb{M} -convergence, asymptotic atom, portmanteau theorem
- These tools reveal the global dynamics of heavy-tailed dynamical systems
- Tail-Inflation-Truncation strategy improves SGDs' generalization performance in DNN

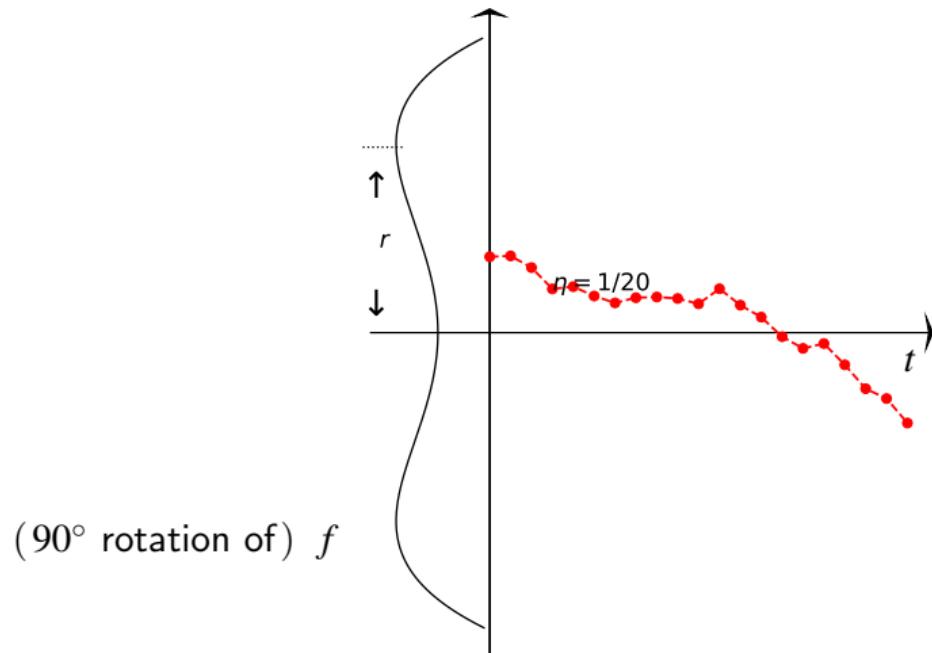
Gradient Flow: Law of Large Numbers



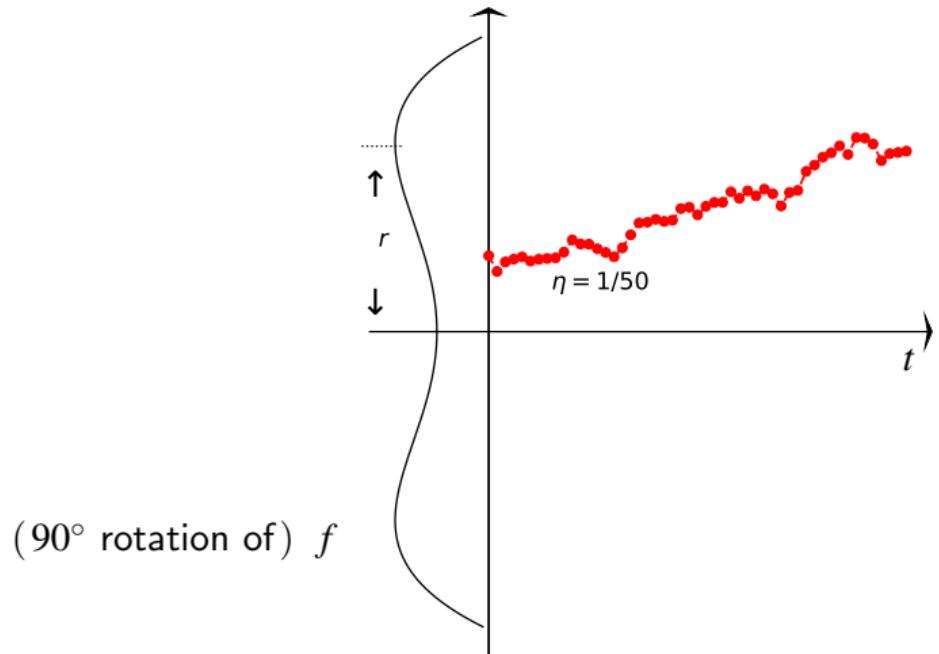
Gradient Flow: Law of Large Numbers



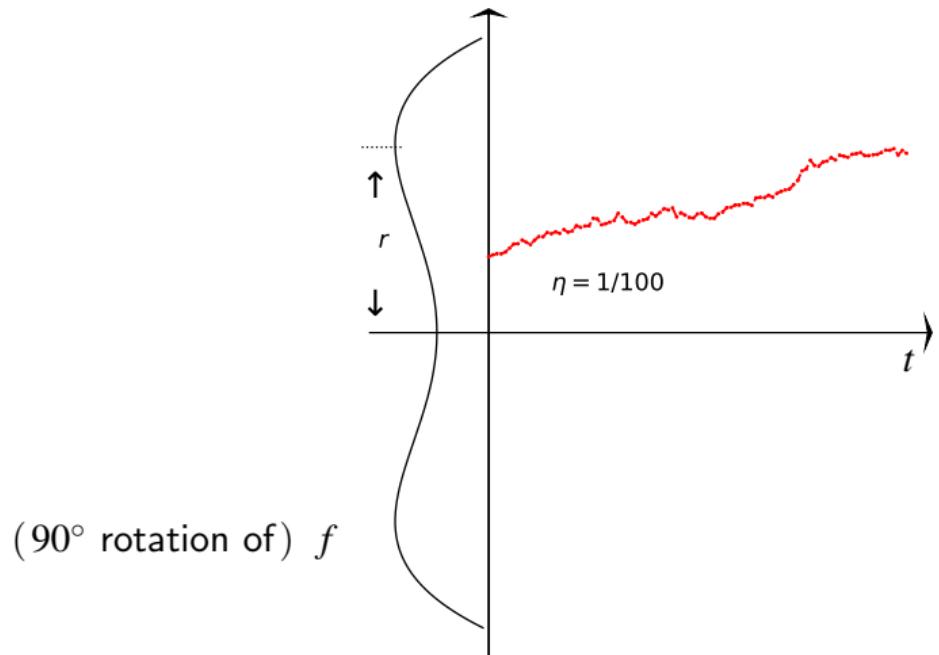
Gradient Flow: Law of Large Numbers



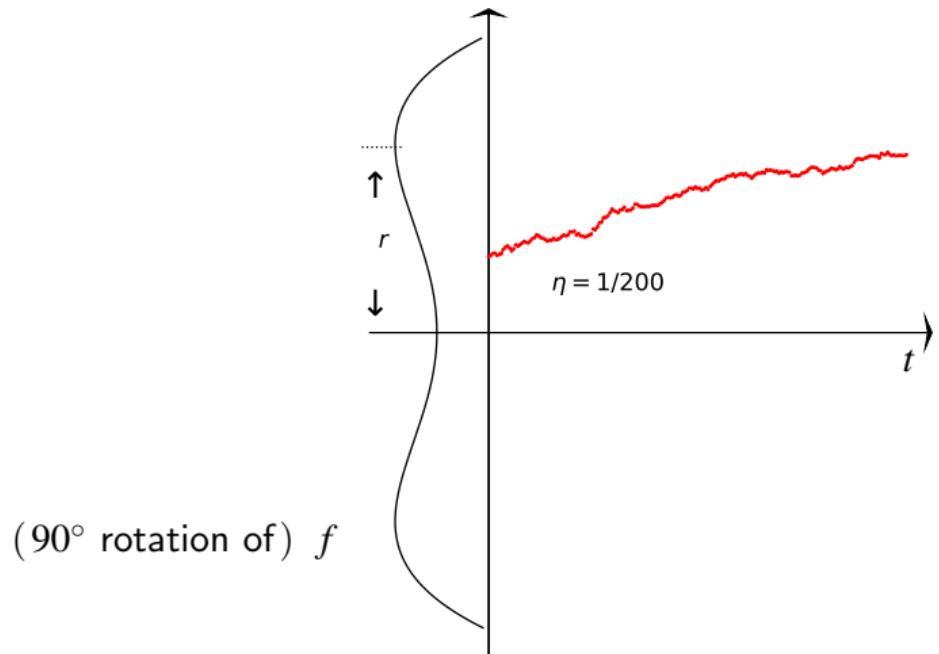
Gradient Flow: Law of Large Numbers



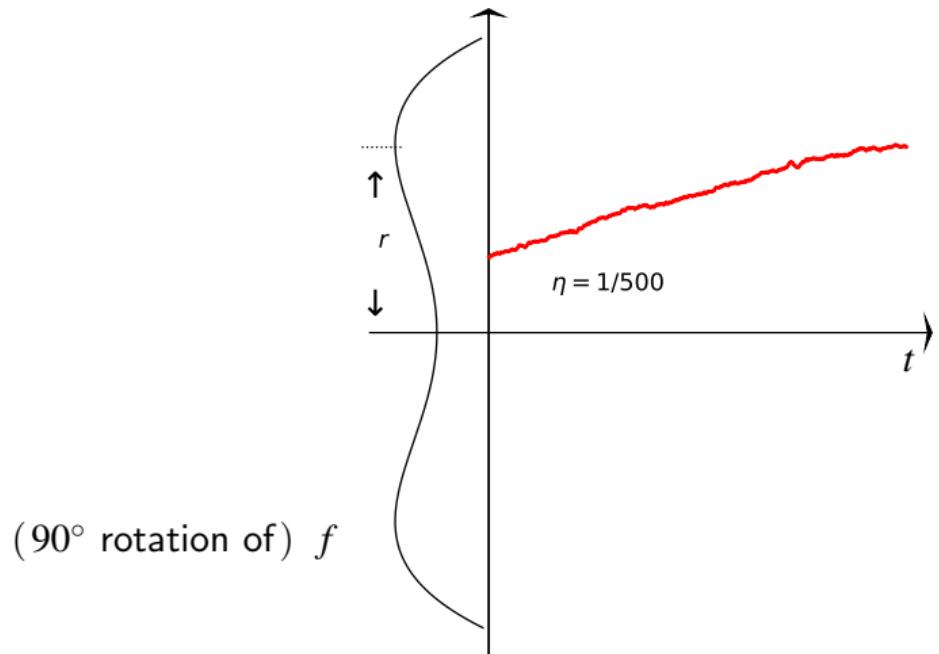
Gradient Flow: Law of Large Numbers



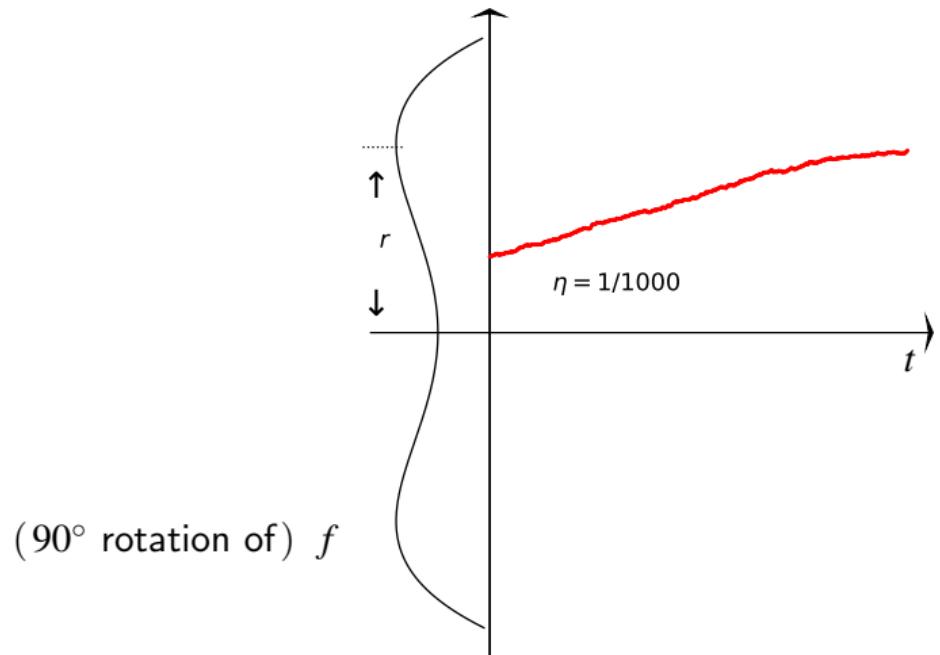
Gradient Flow: Law of Large Numbers



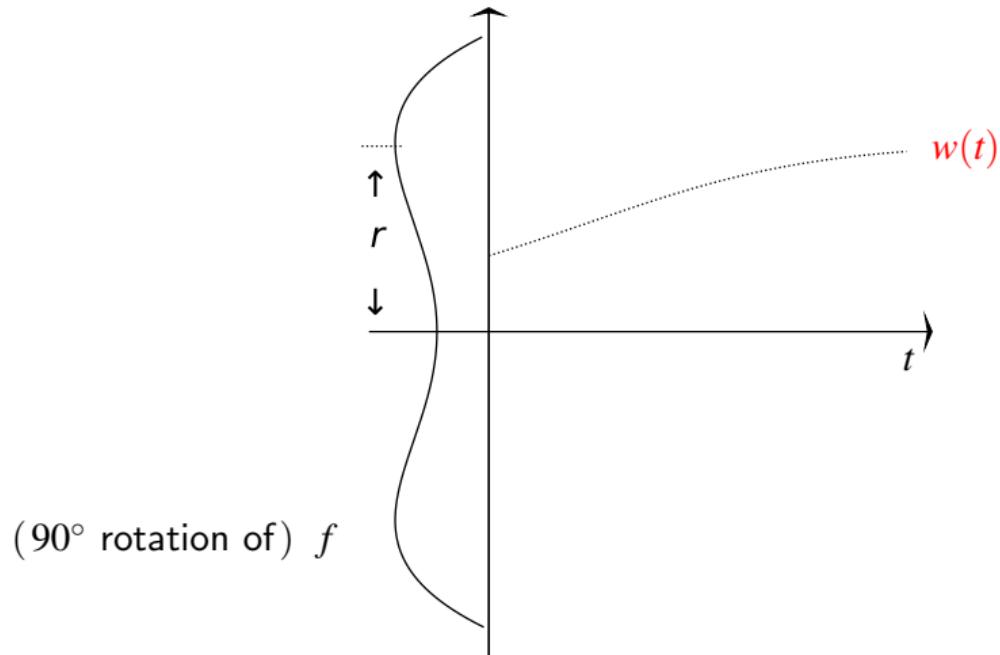
Gradient Flow: Law of Large Numbers



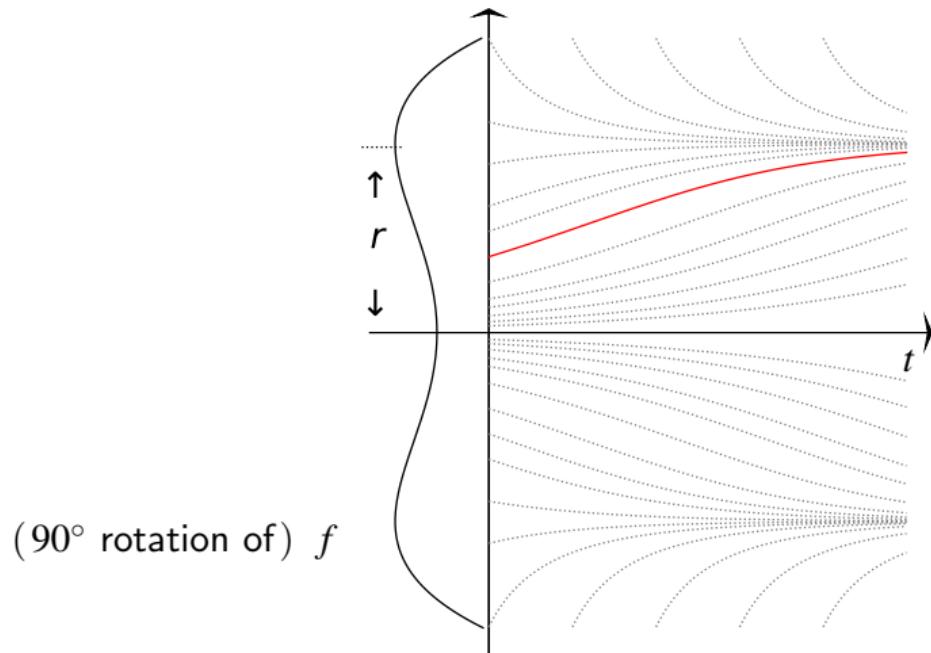
Gradient Flow: Law of Large Numbers



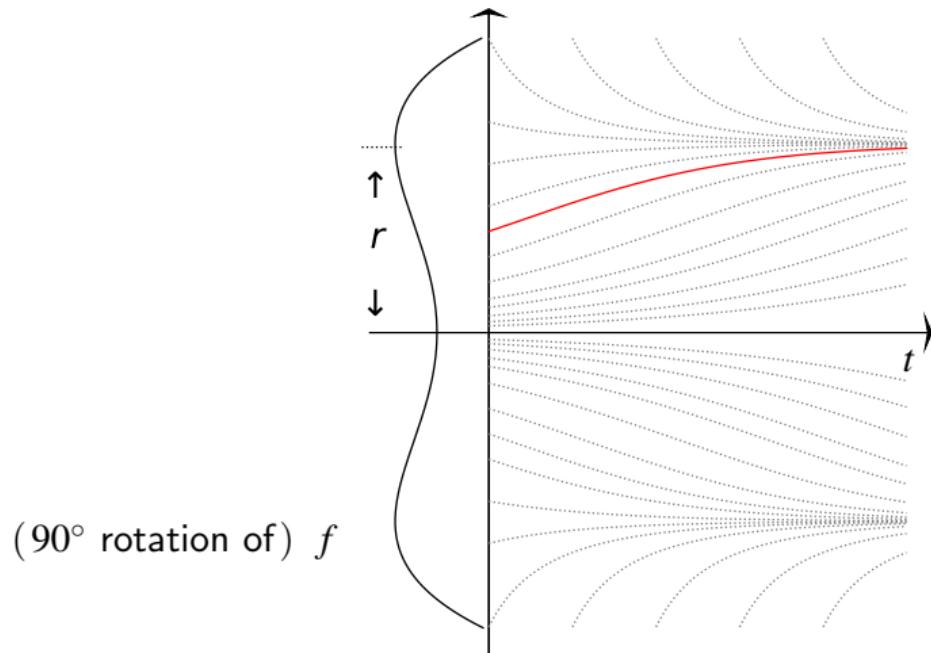
Gradient Flow: Law of Large Numbers



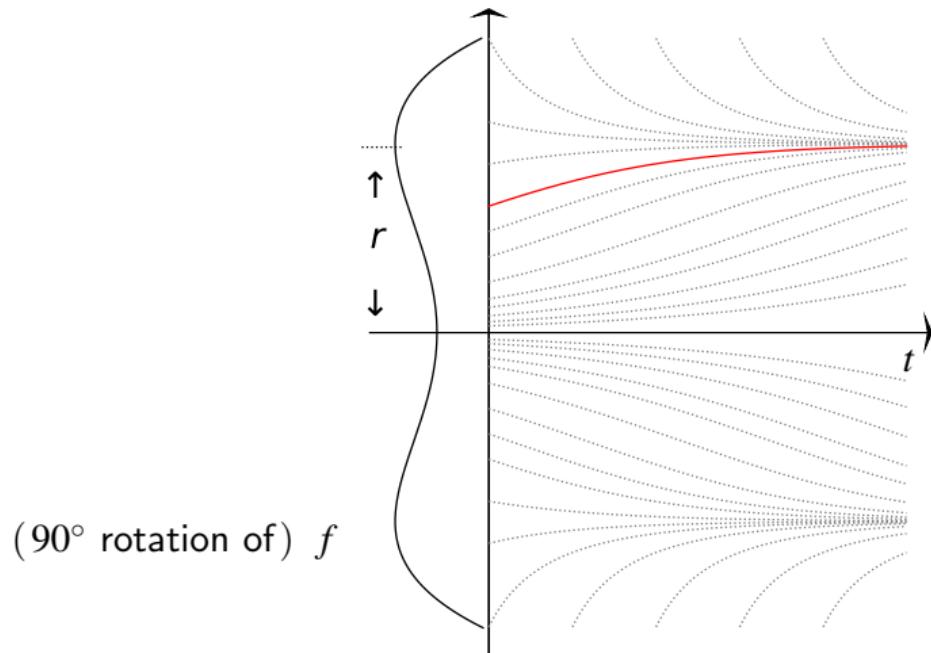
Gradient Flow: Law of Large Numbers



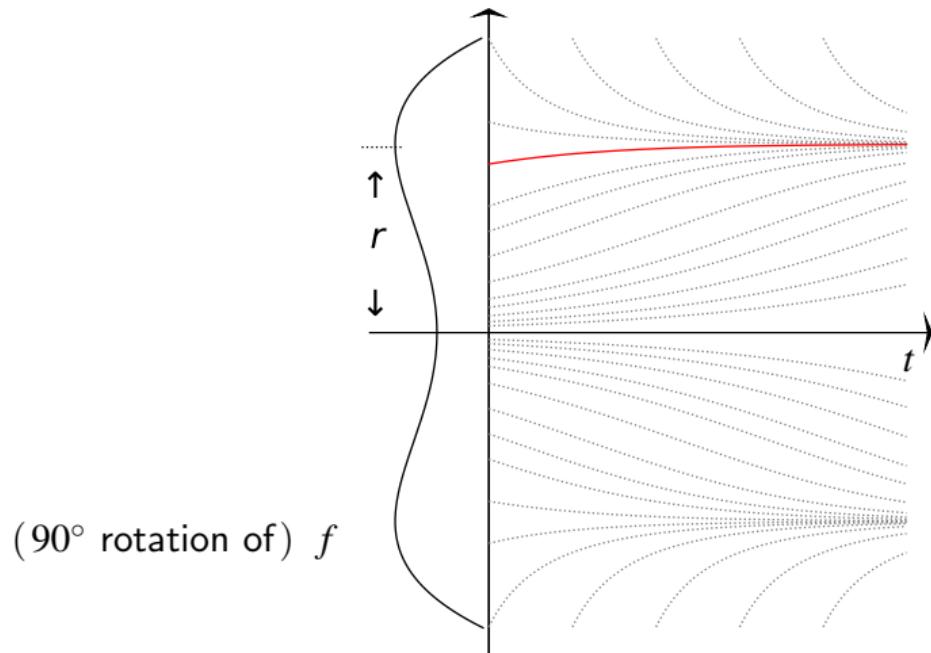
Gradient Flow: Law of Large Numbers



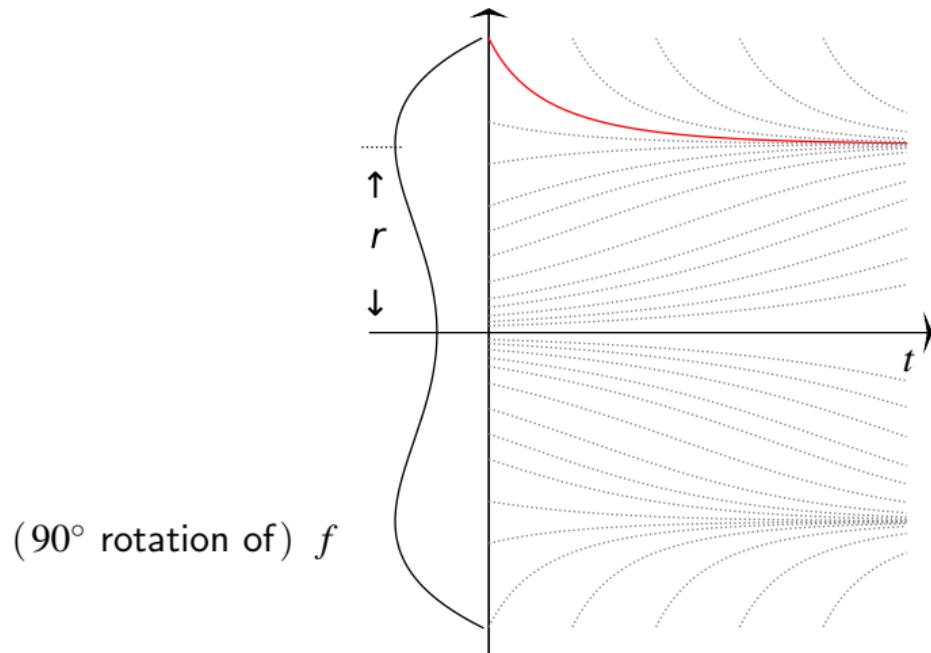
Gradient Flow: Law of Large Numbers



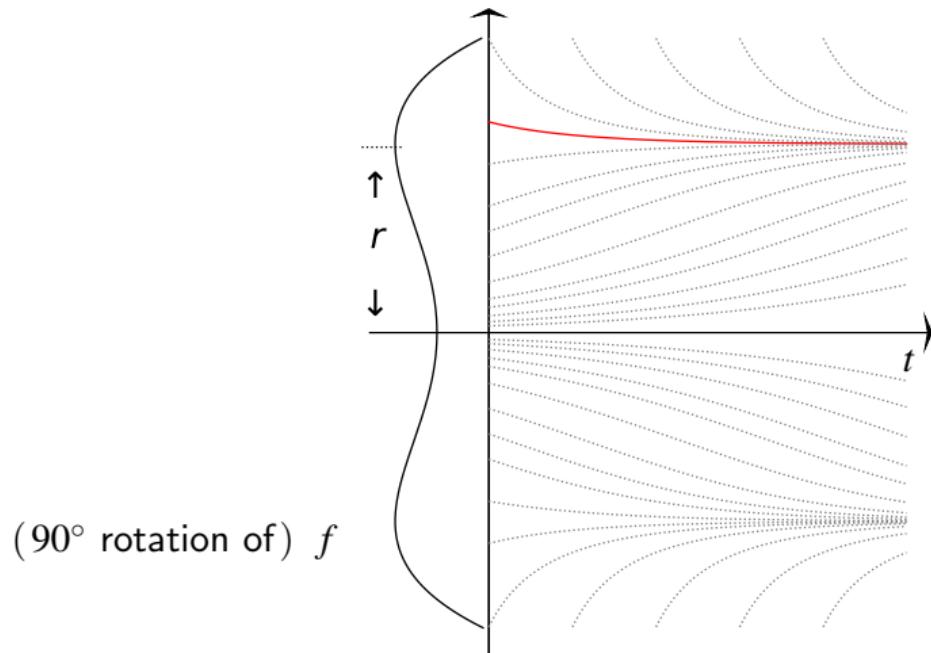
Gradient Flow: Law of Large Numbers



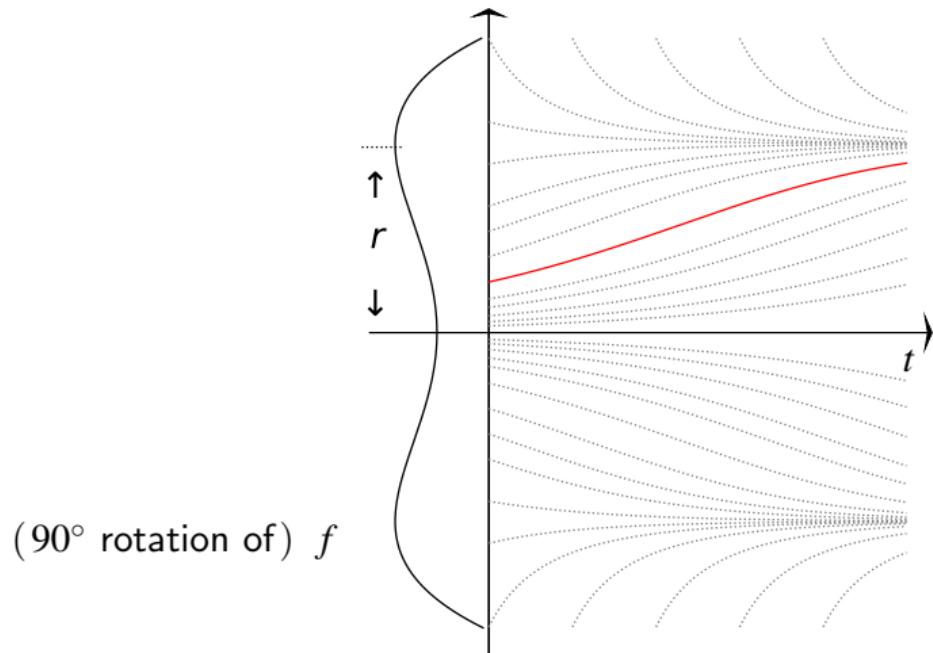
Gradient Flow: Law of Large Numbers



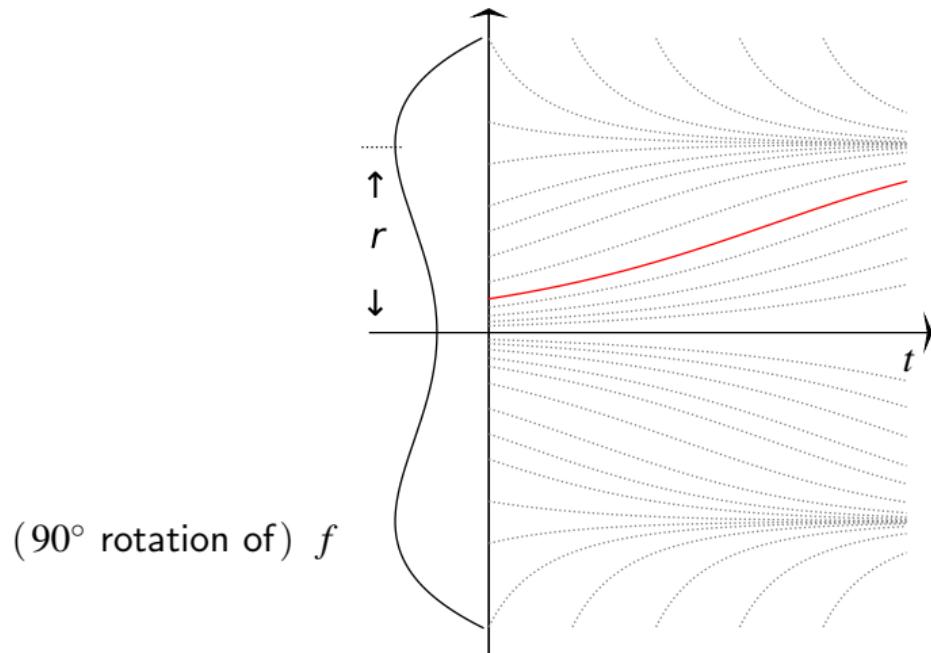
Gradient Flow: Law of Large Numbers



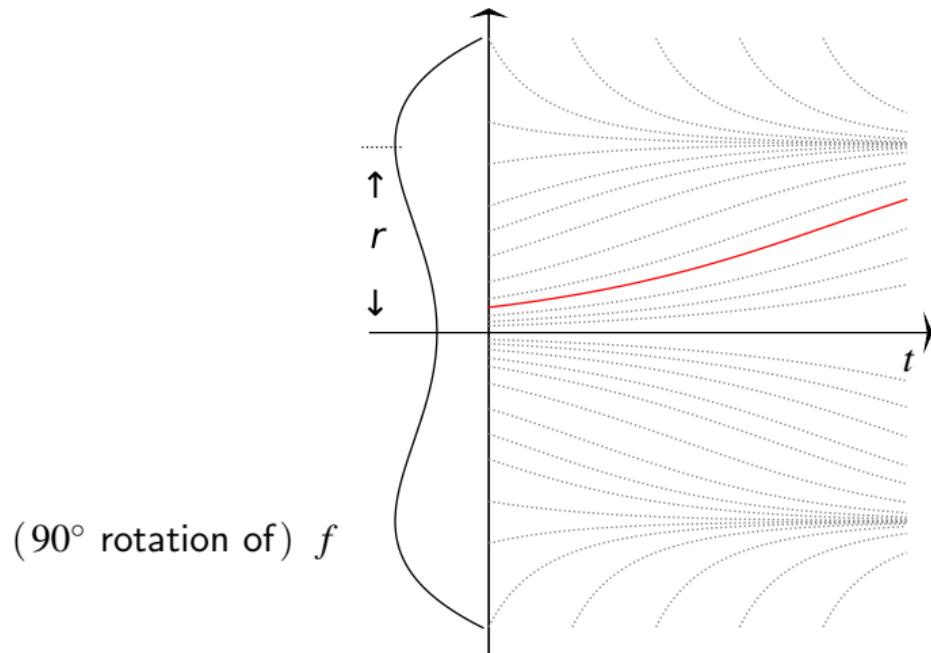
Gradient Flow: Law of Large Numbers



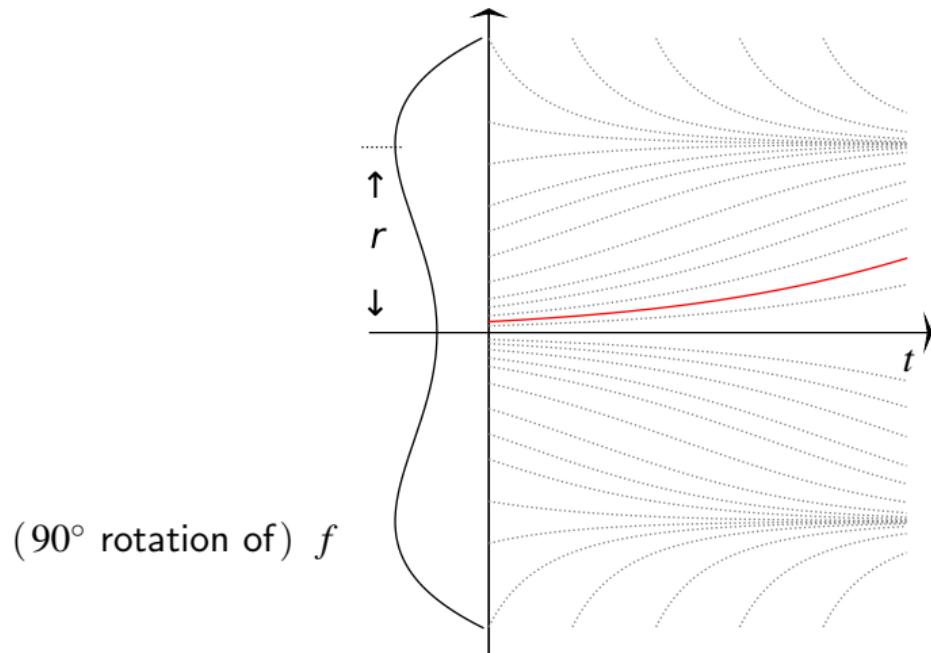
Gradient Flow: Law of Large Numbers



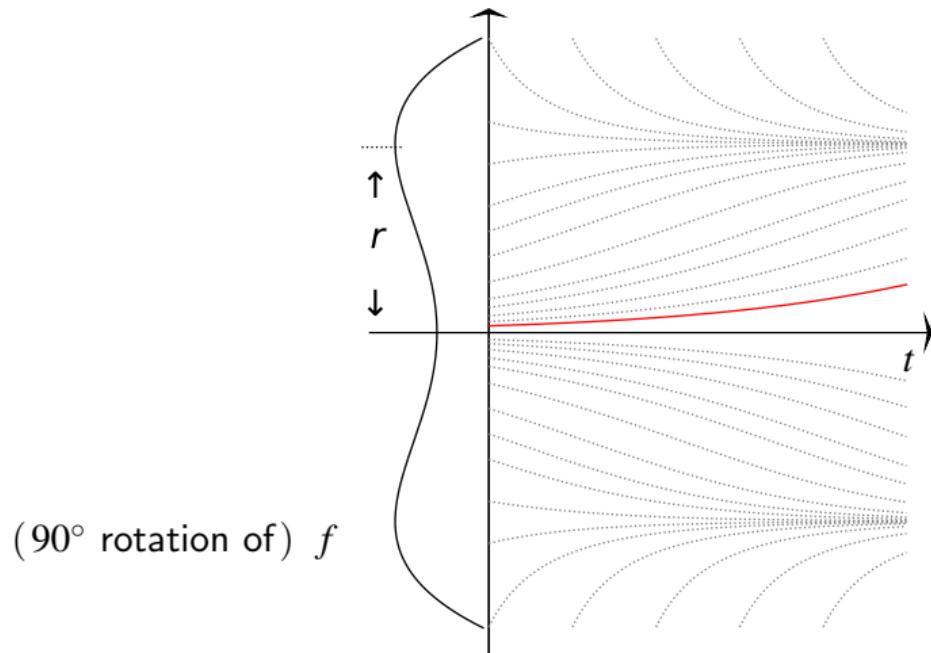
Gradient Flow: Law of Large Numbers



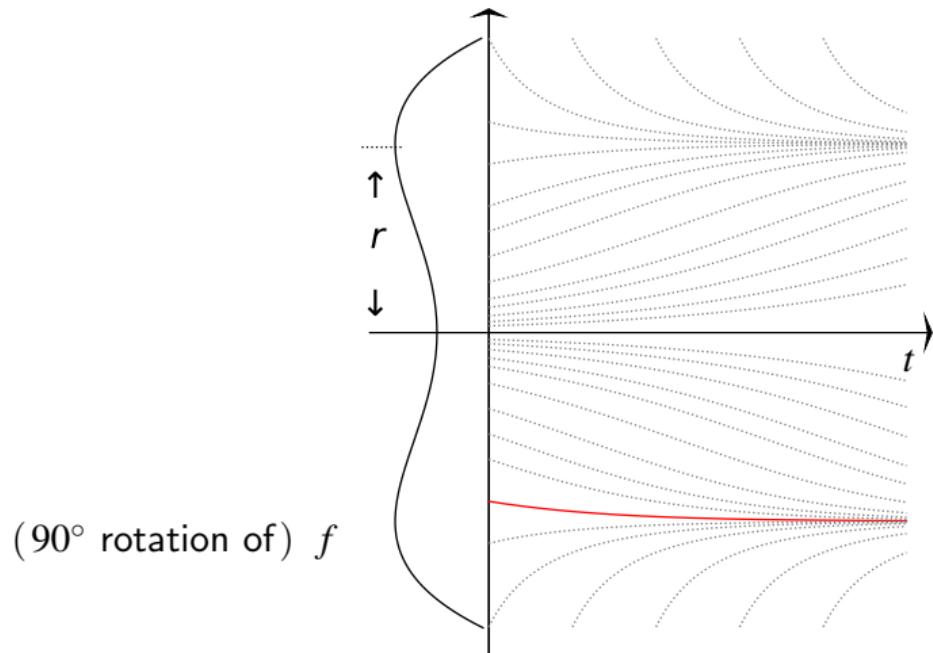
Gradient Flow: Law of Large Numbers



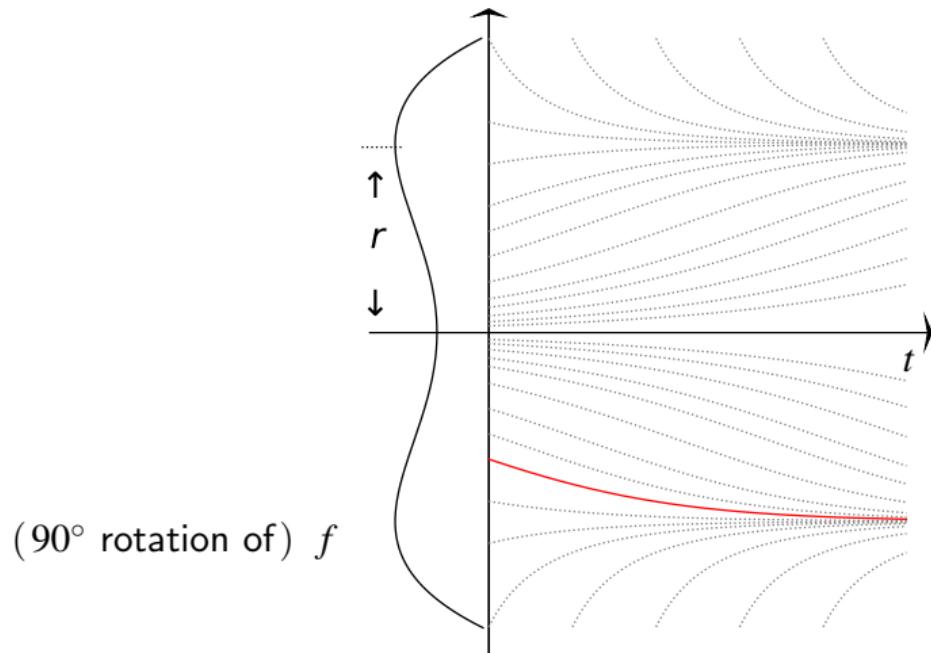
Gradient Flow: Law of Large Numbers



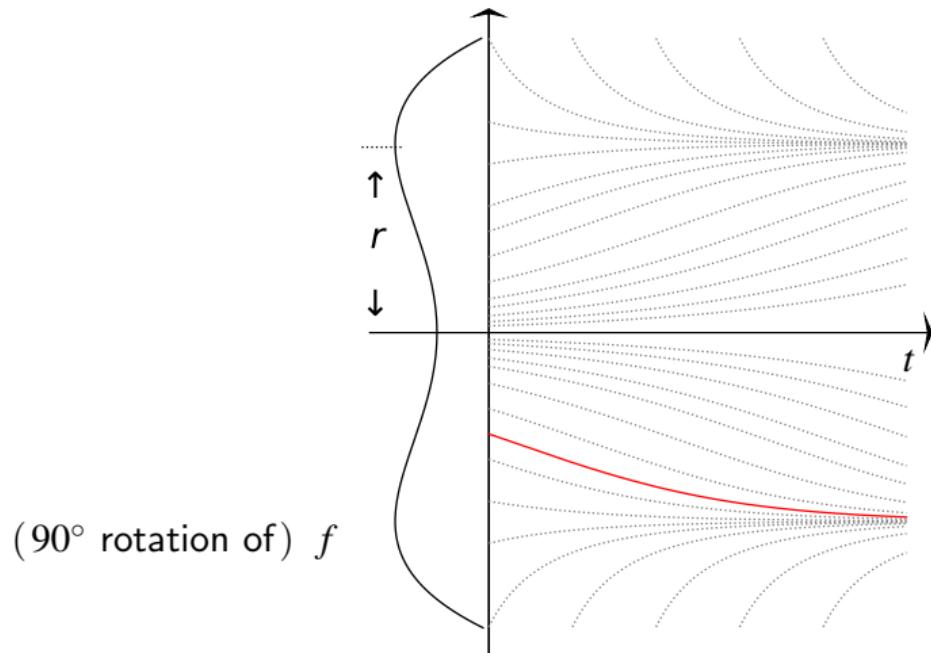
Gradient Flow: Law of Large Numbers



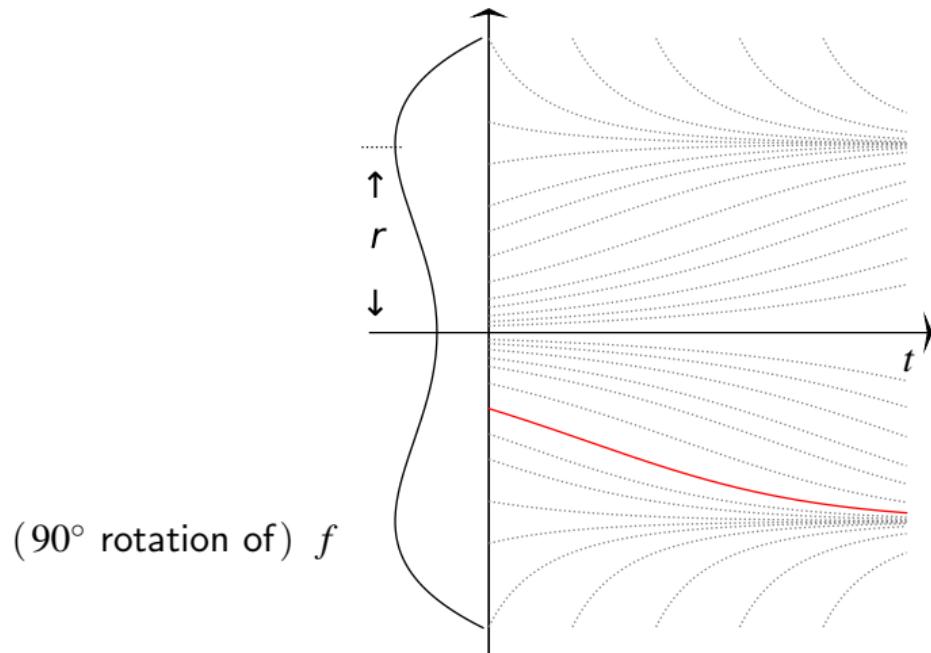
Gradient Flow: Law of Large Numbers



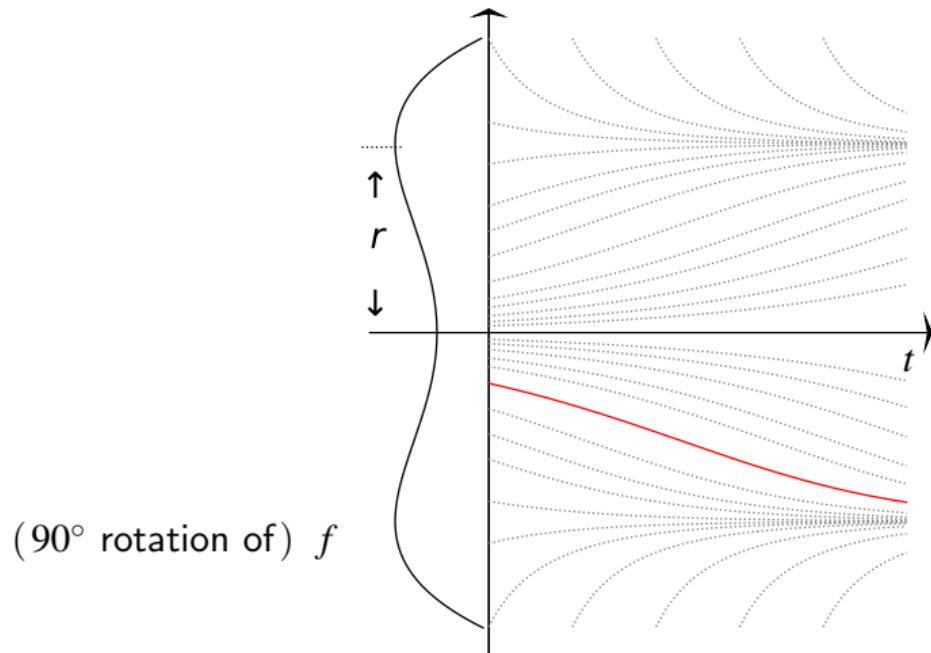
Gradient Flow: Law of Large Numbers



Gradient Flow: Law of Large Numbers



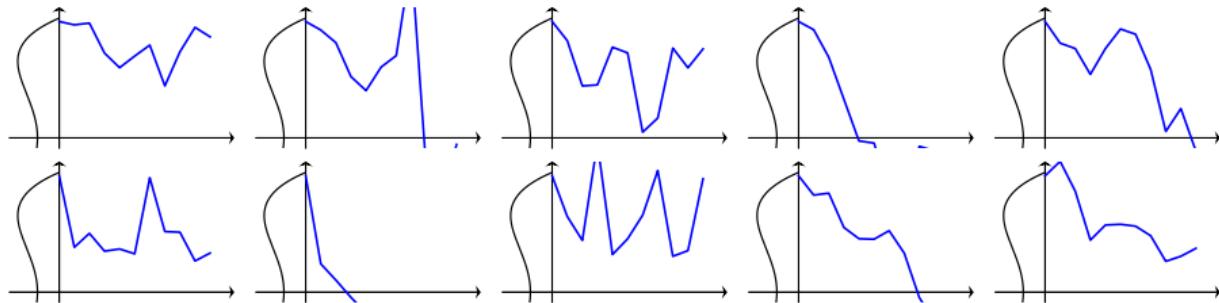
Gradient Flow: Law of Large Numbers



Typical Scenario

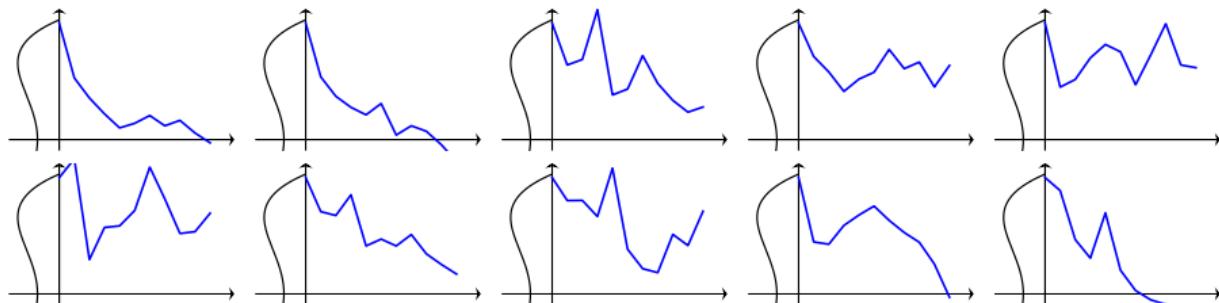
Trajectory of SGD W^η :

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



Trajectory of SGD W^η :

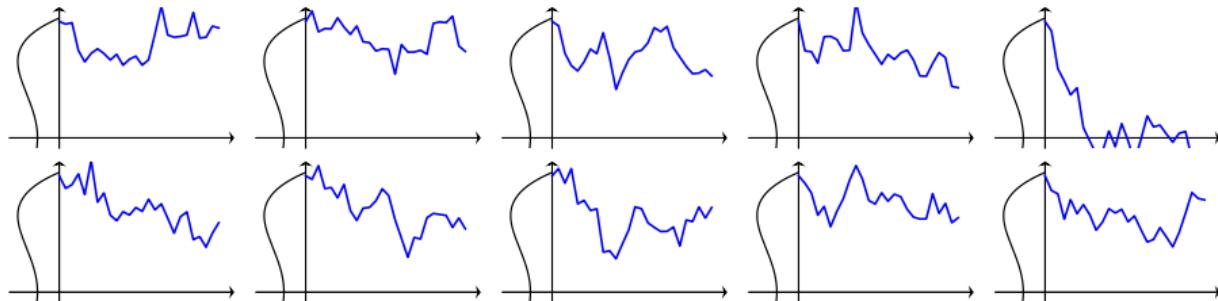
$\eta = 1/10$ & noises are **heavy-tailed**



Typical Scenario

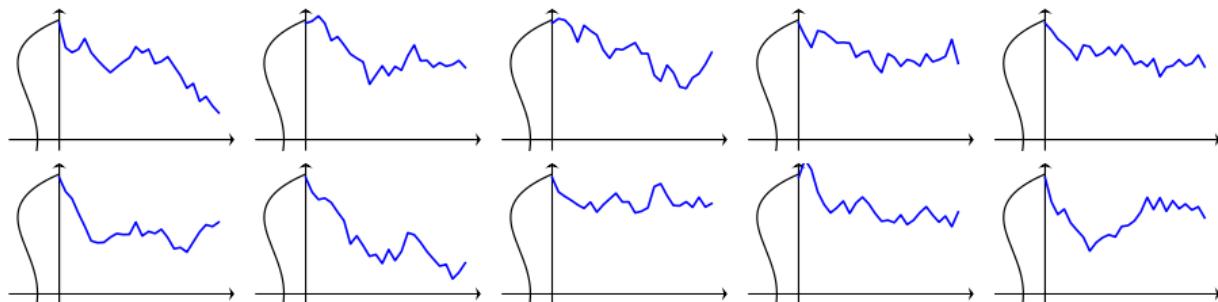
Trajectory of SGD W^η :

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



Trajectory of SGD W^η :

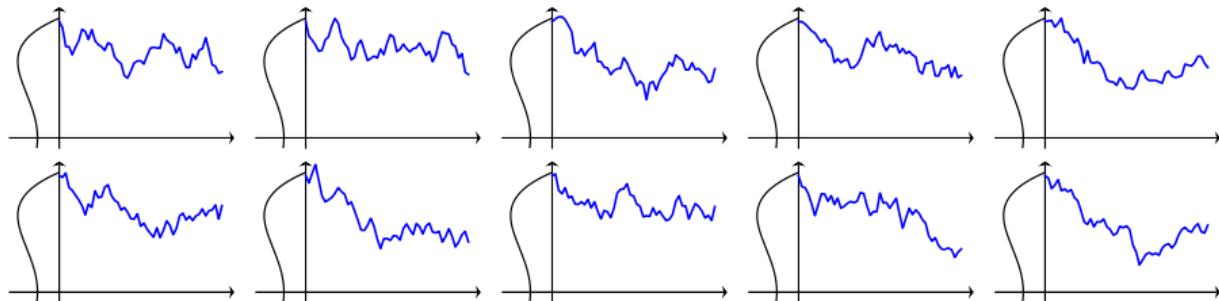
$\eta = 1/25$ & noises are **heavy-tailed**



Typical Scenario

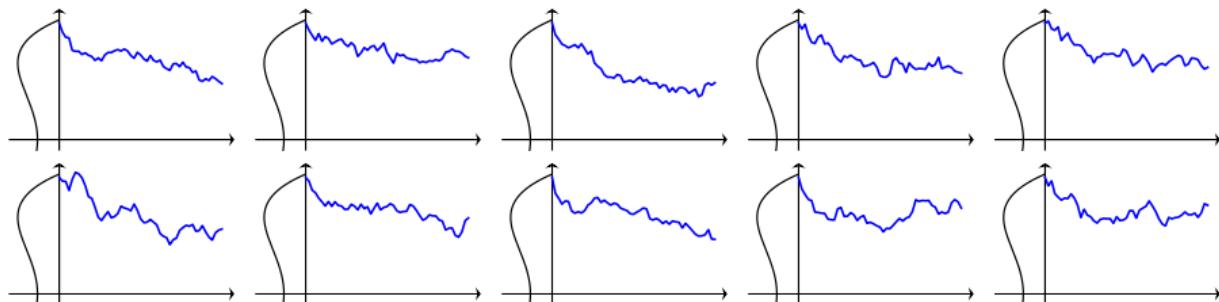
Trajectory of SGD W^η :

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



Trajectory of SGD W^η :

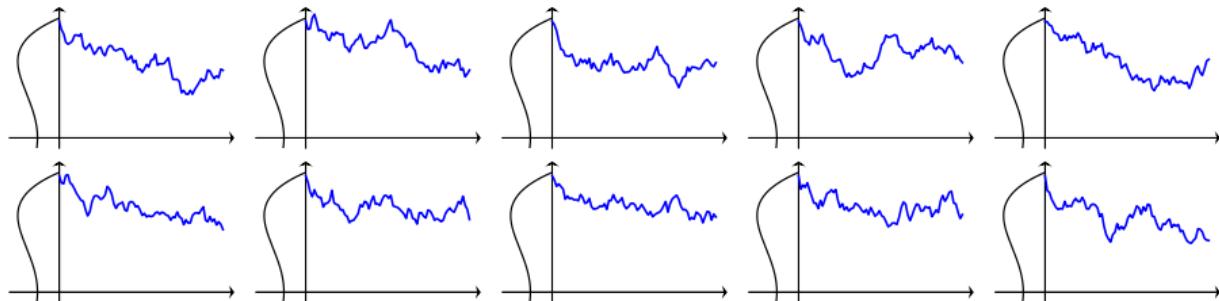
$\eta = 1/50$ & noises are **heavy-tailed**



Typical Scenario

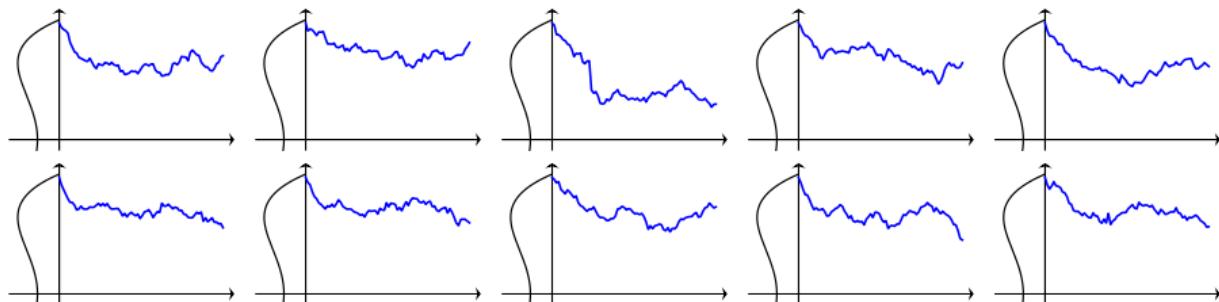
Trajectory of SGD W^η :

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



Trajectory of SGD W^η :

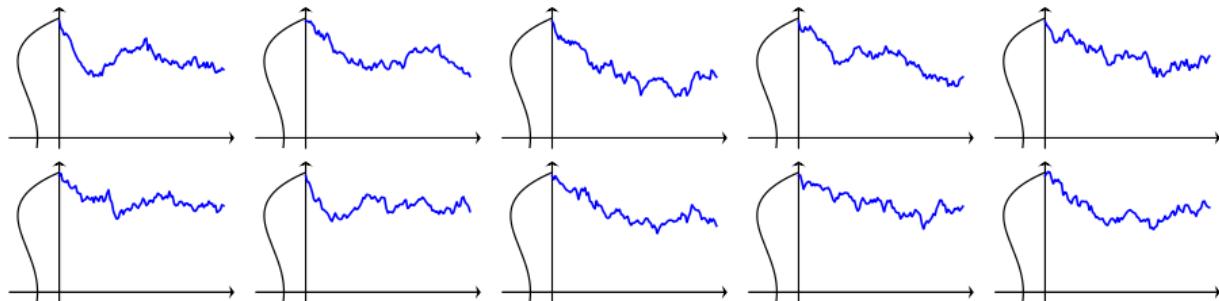
$\eta = 1/75$ & noises are **heavy-tailed**



Typical Scenario

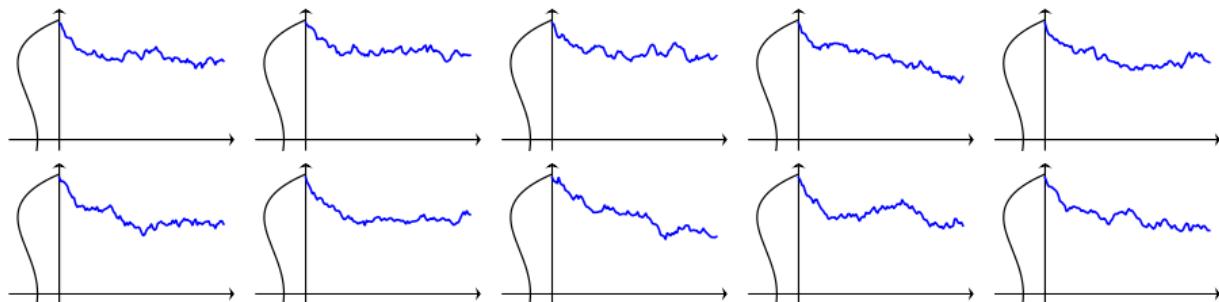
Trajectory of SGD W^η :

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



Trajectory of SGD W^η :

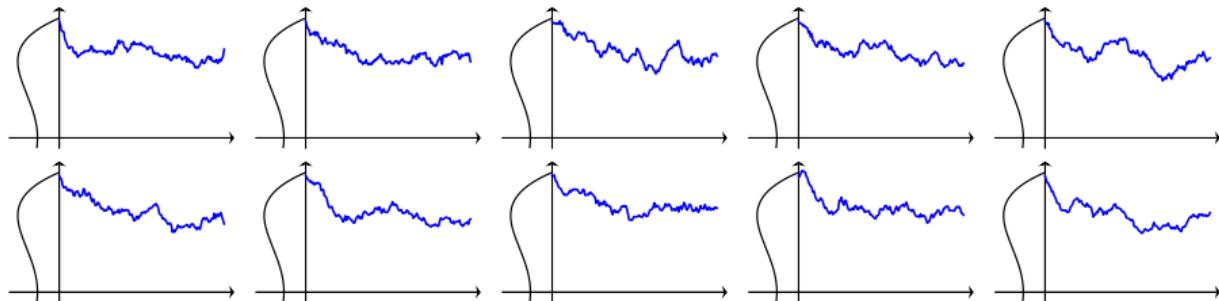
$\eta = 1/100$ & noises are **heavy-tailed**



Typical Scenario

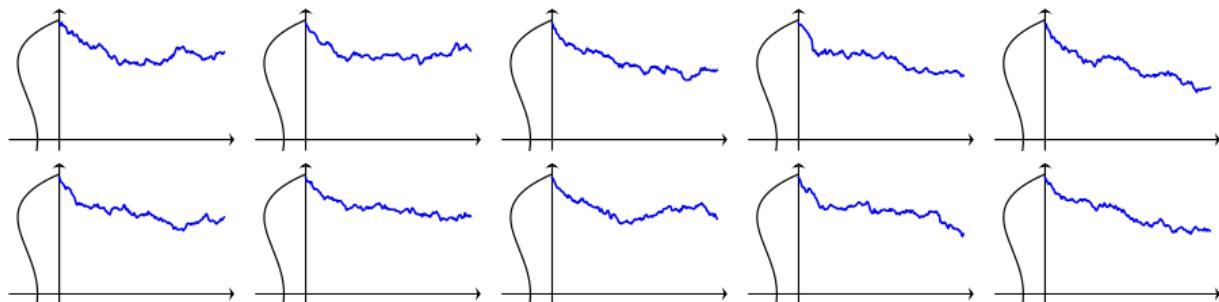
Trajectory of SGD W^η :

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



Trajectory of SGD W^η :

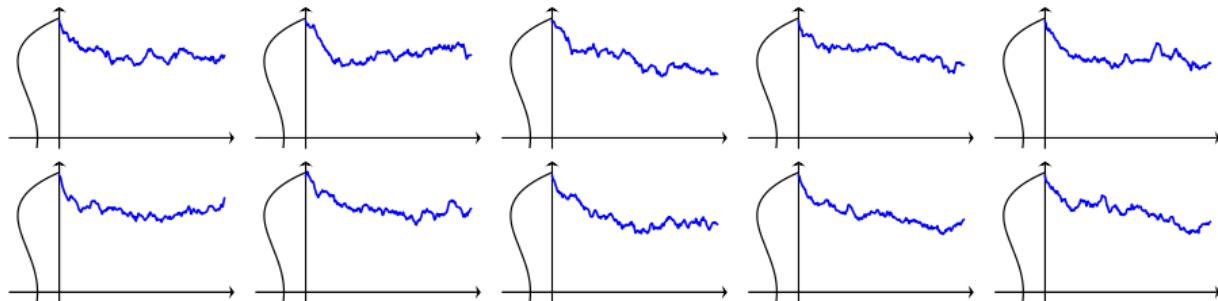
$\eta = 1/150$ & noises are **heavy-tailed**



Typical Scenario

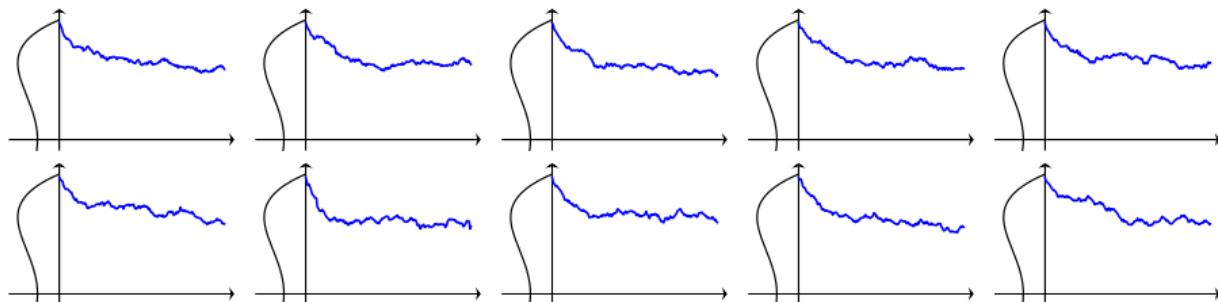
Trajectory of SGD W^η :

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

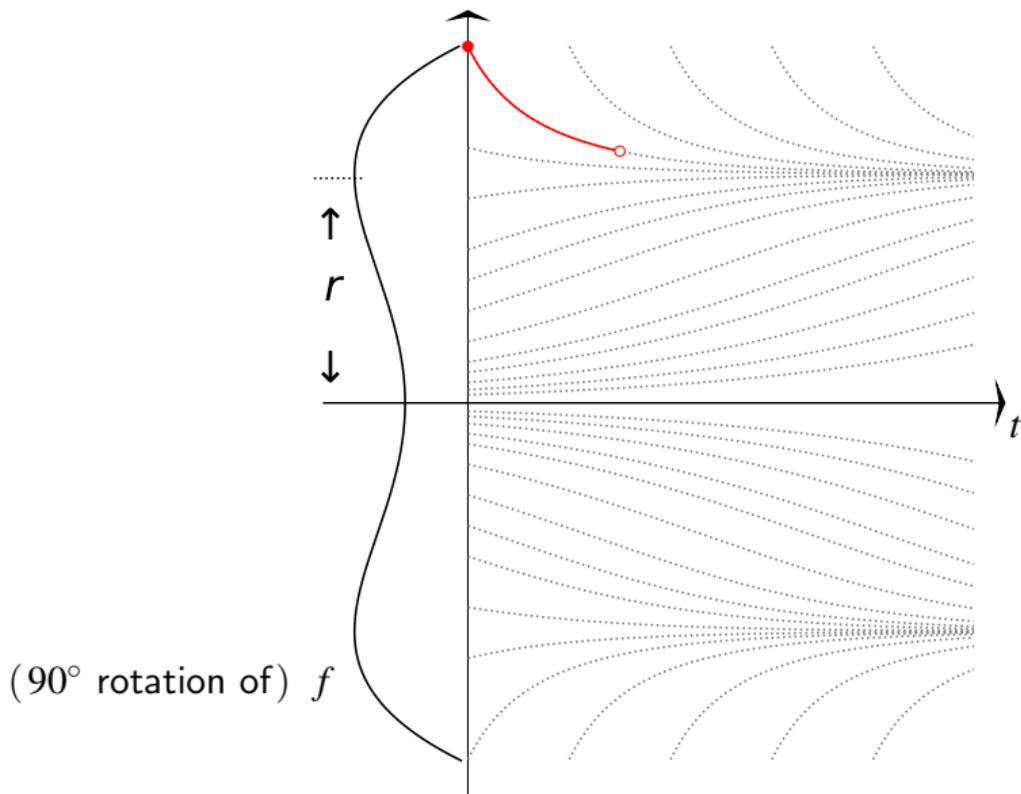


Trajectory of SGD W^η :

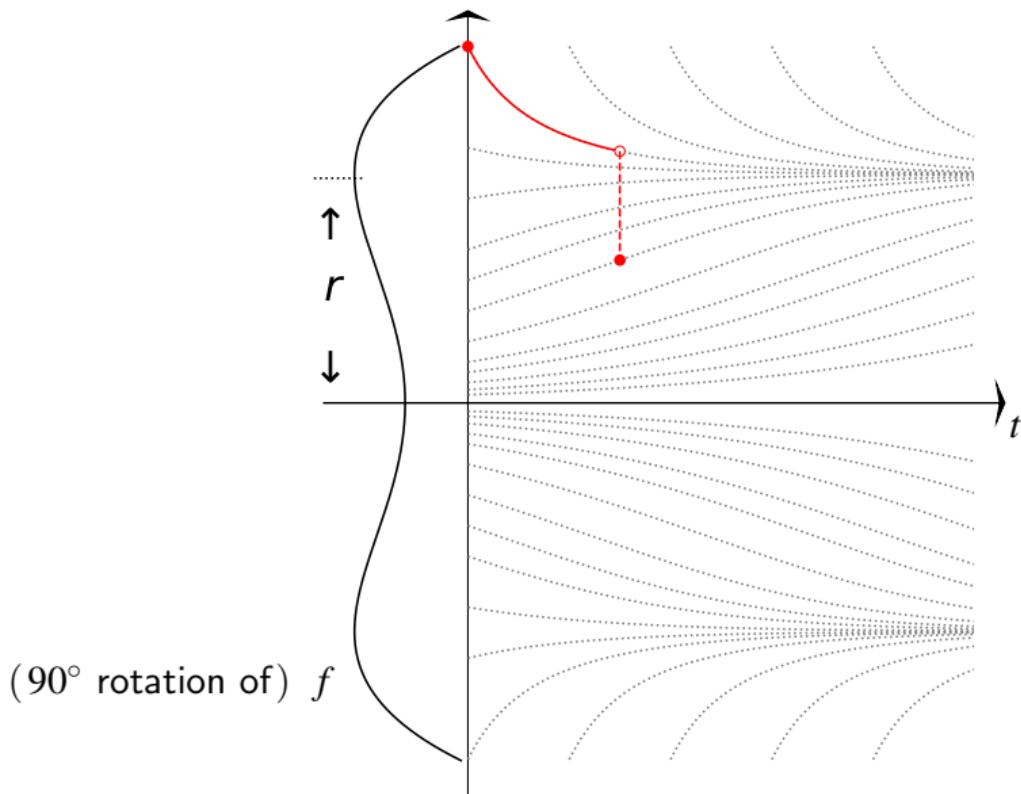
$\eta = 1/200$ & noises are **heavy-tailed**



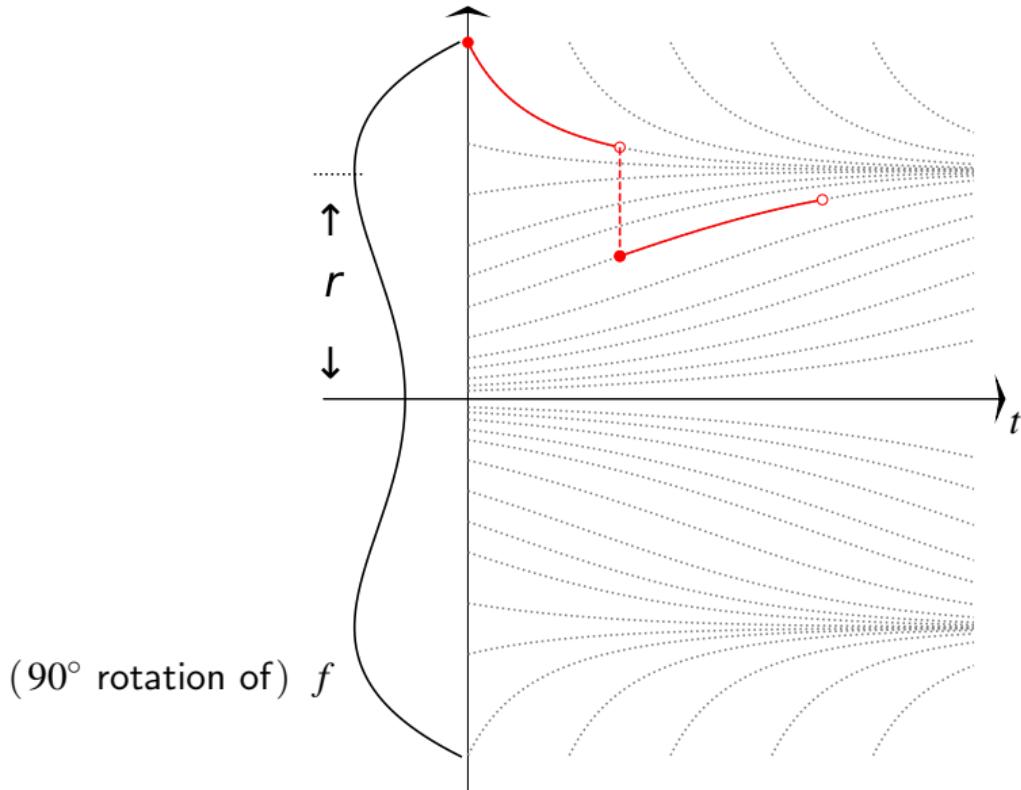
Adding Jumps to $w(\cdot)$: Piecewise Gradient Flow



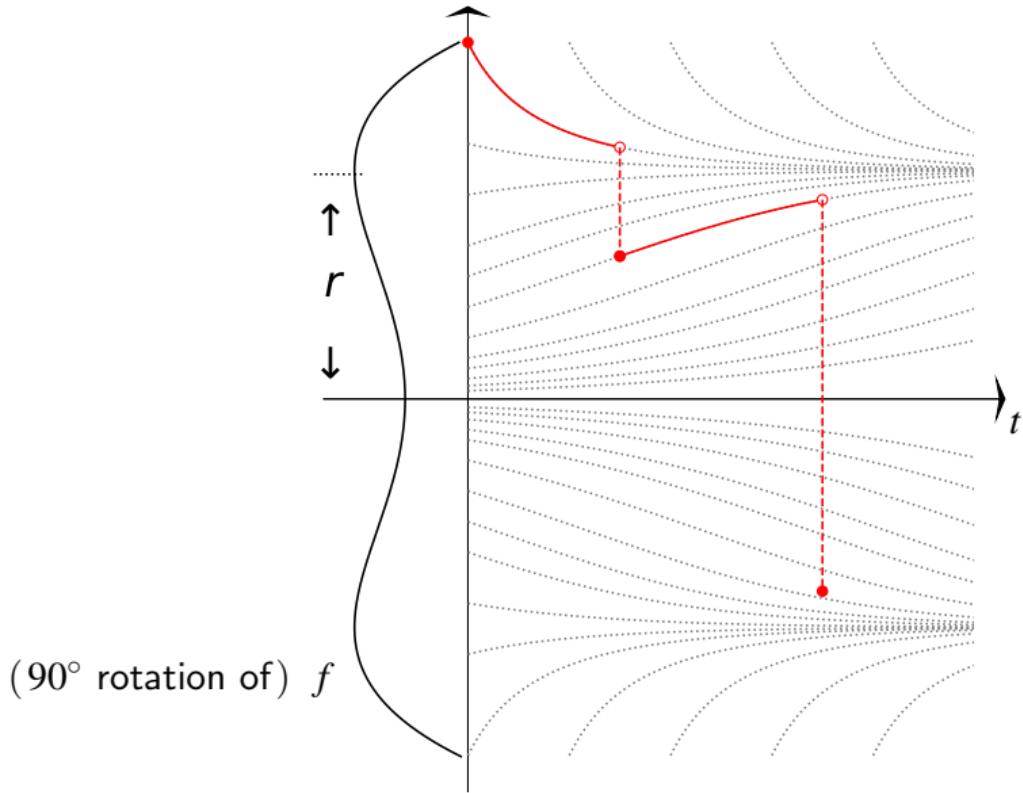
Adding Jumps to $w(\cdot)$: Piecewise Gradient Flow



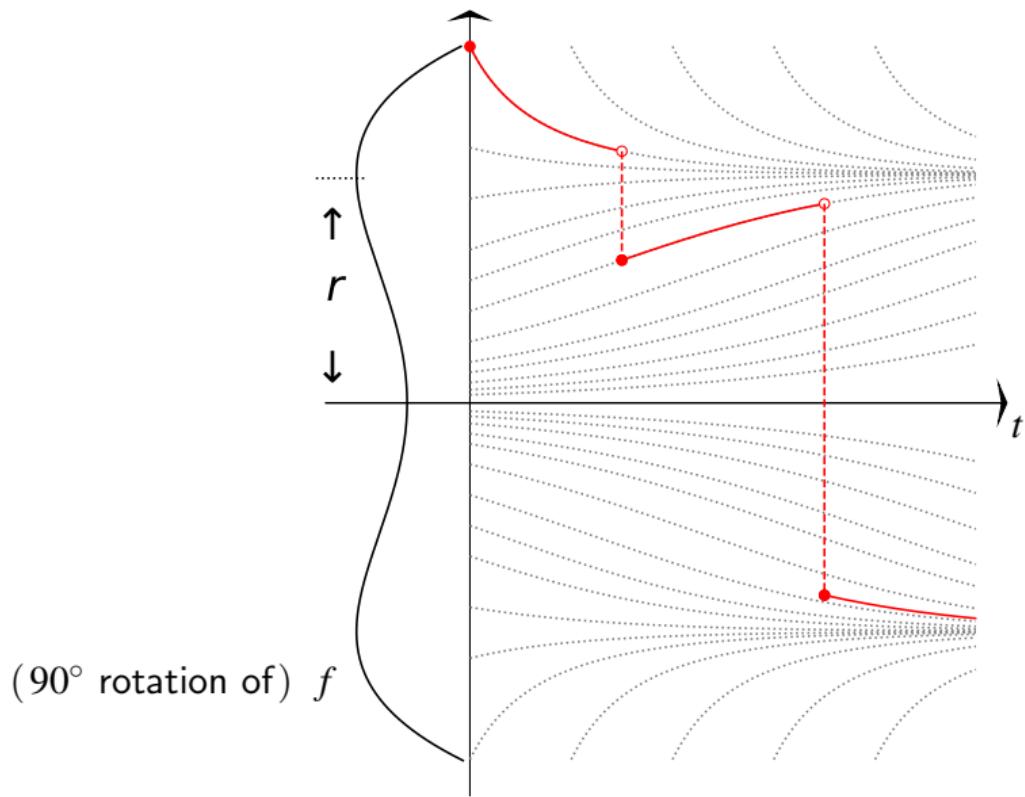
Adding Jumps to $w(\cdot)$: Piecewise Gradient Flow



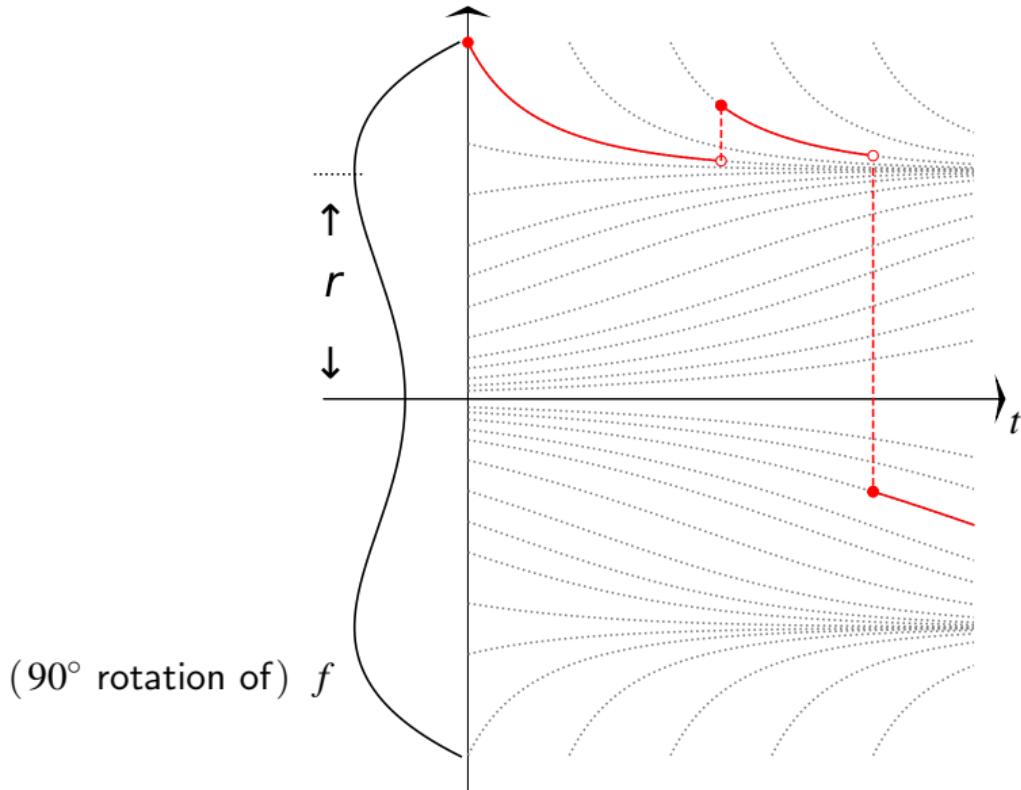
Adding Jumps to $w(\cdot)$: Piecewise Gradient Flow



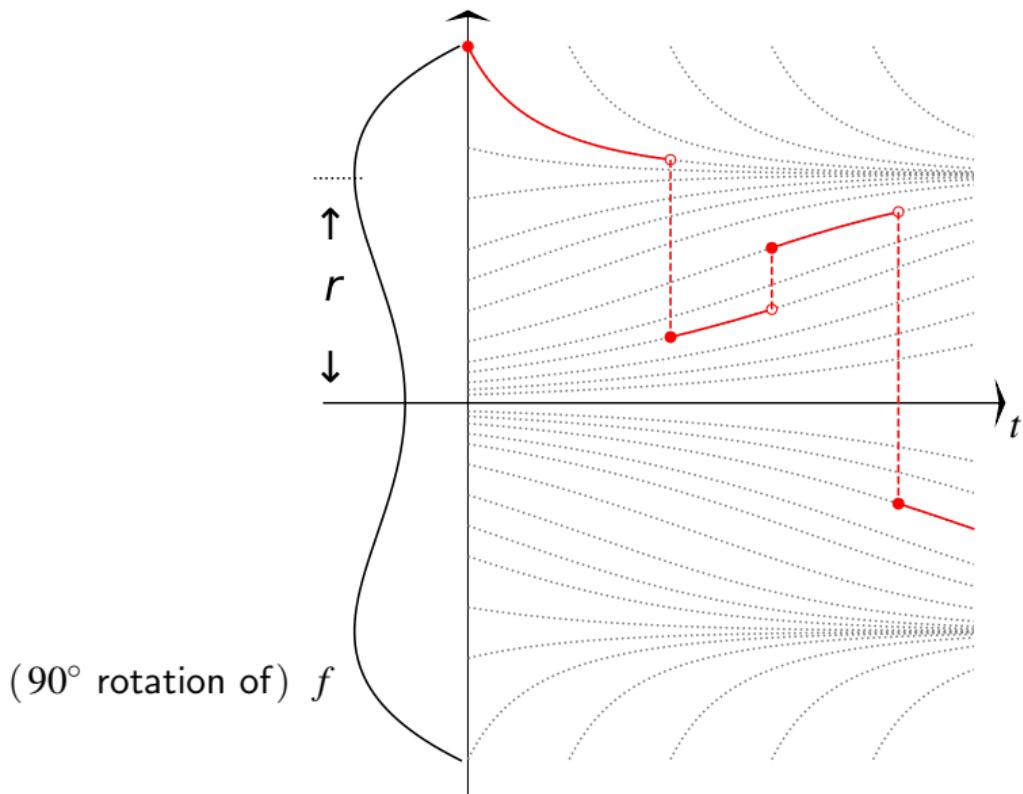
Adding Jumps to $w(\cdot)$: Piecewise Gradient Flow



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

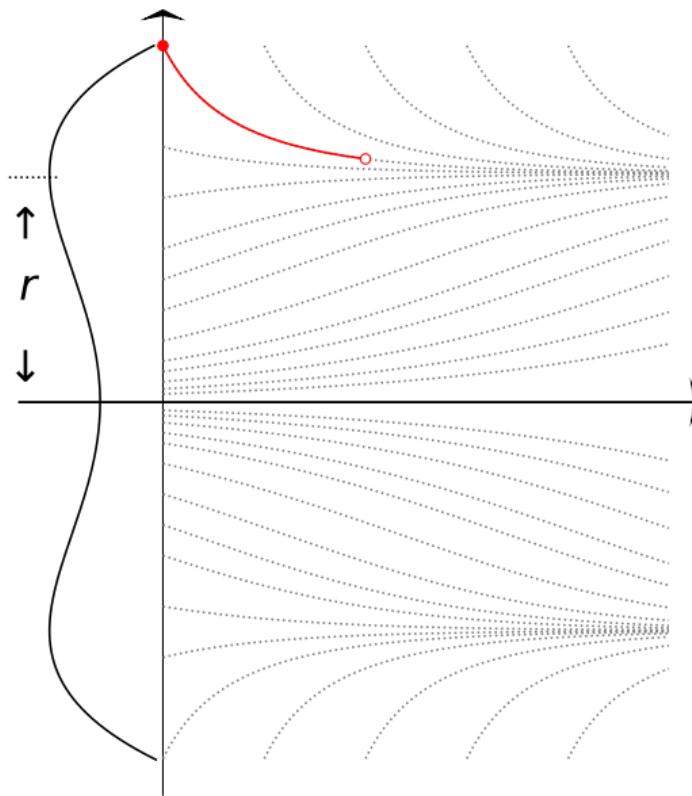


Adding Jumps to $w(\cdot)$: Piecewise Gradient Flow



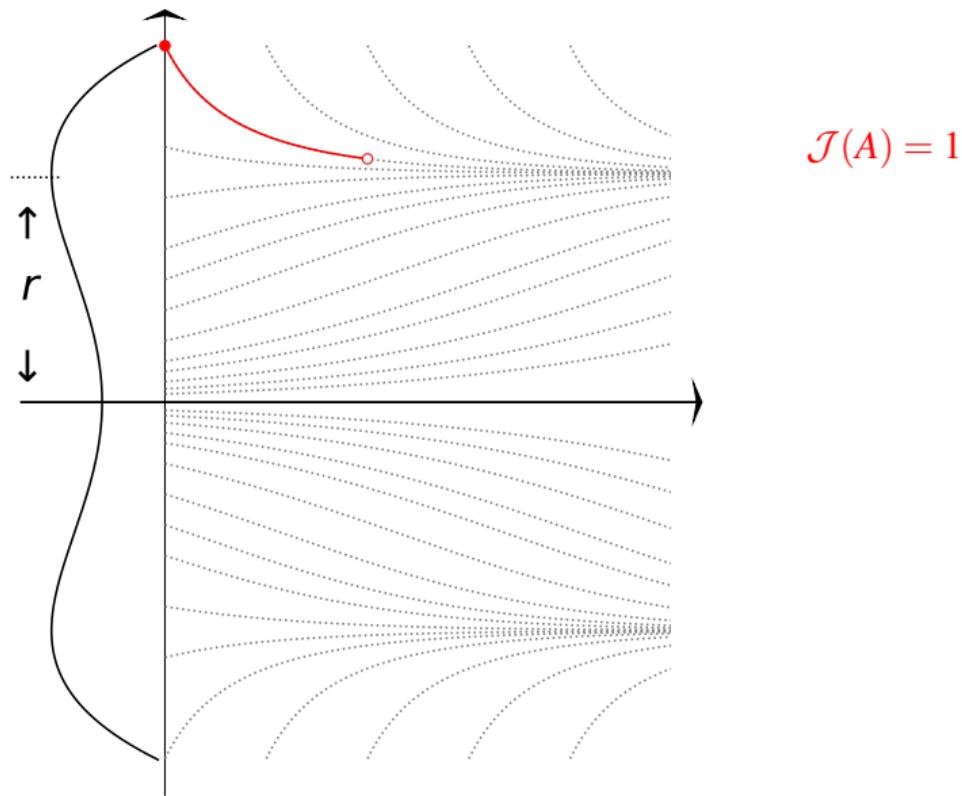
Catastrophe Principle: Most Likely Escape Route

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



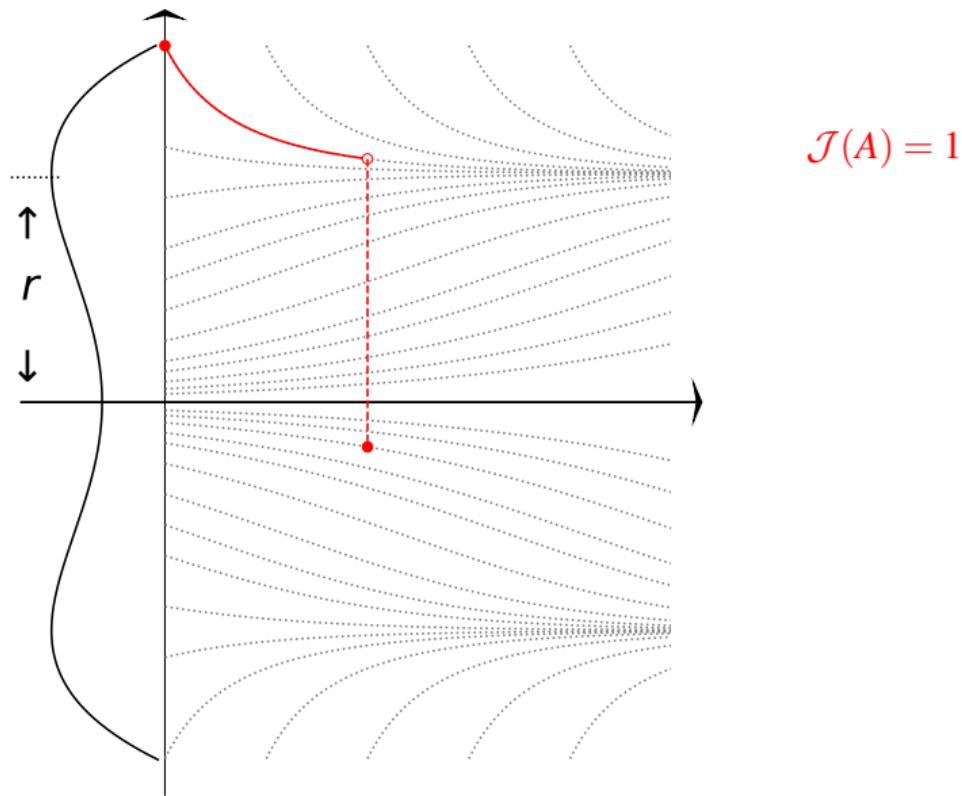
Catastrophe Principle: Most Likely Escape Route

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



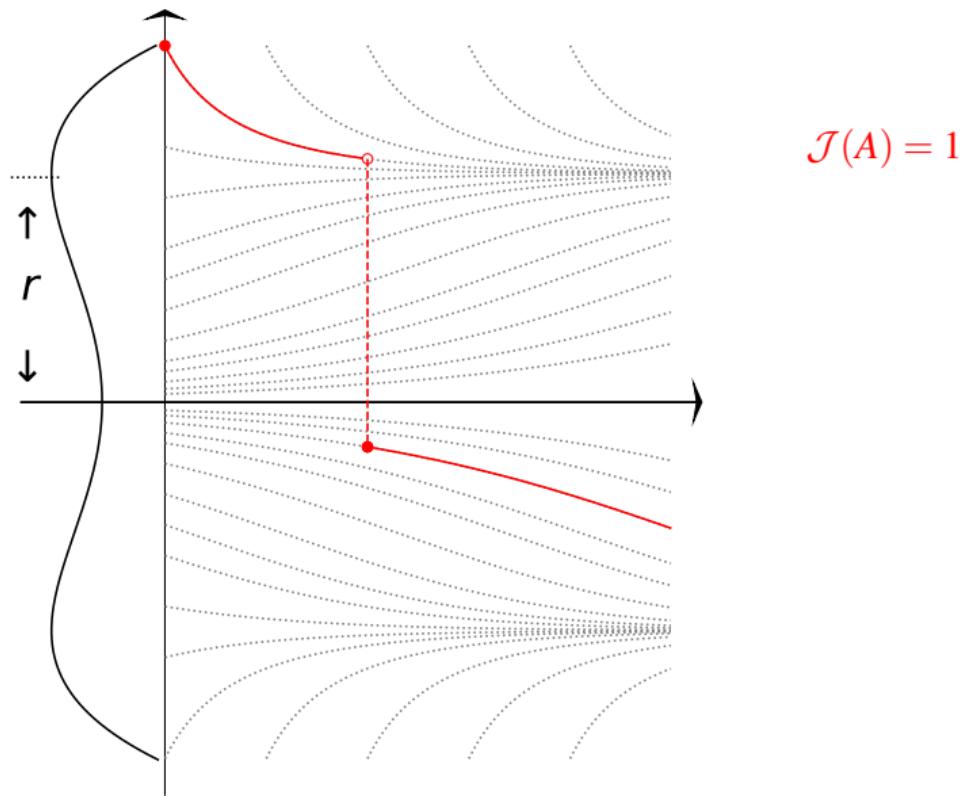
Catastrophe Principle: Most Likely Escape Route

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



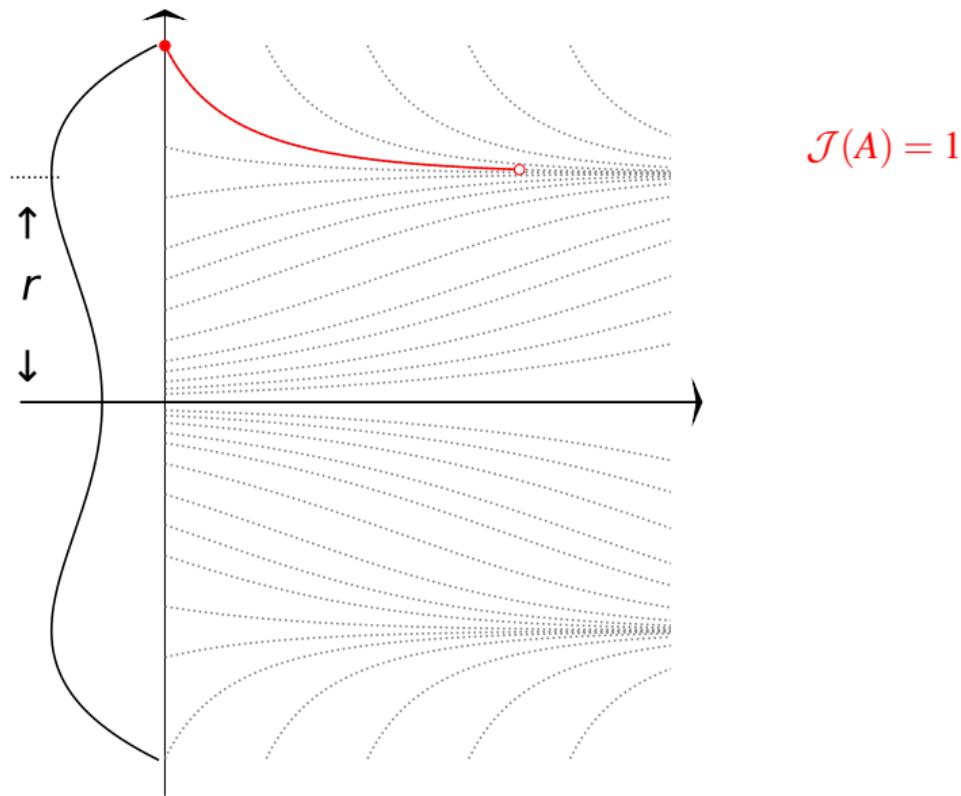
Catastrophe Principle: Most Likely Escape Route

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



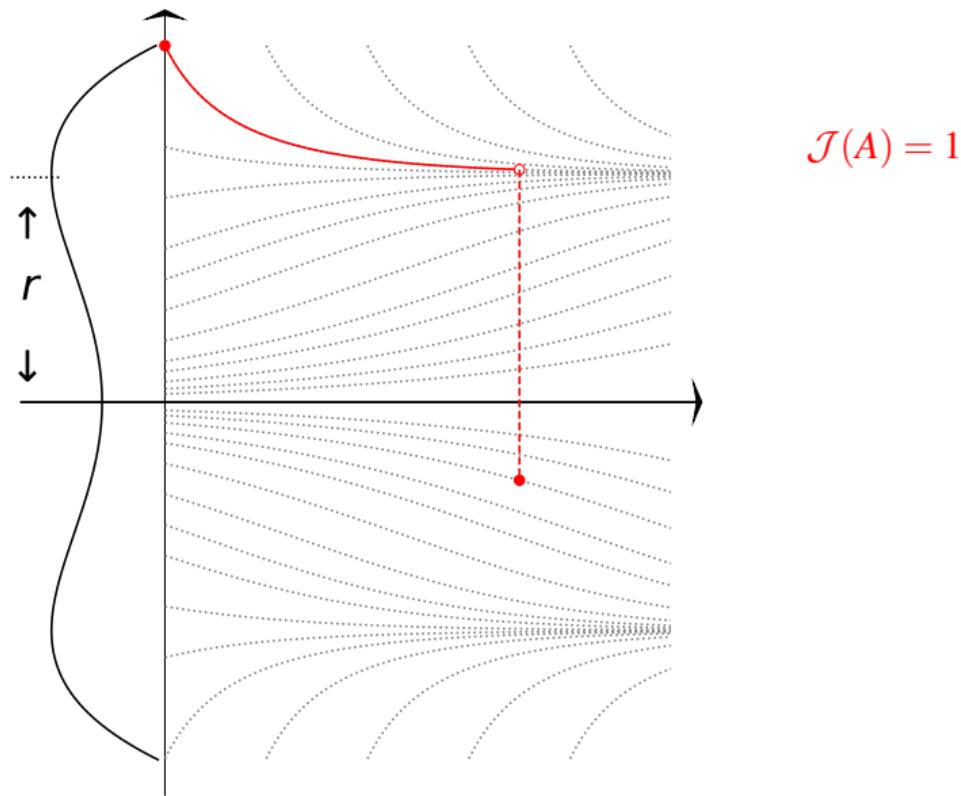
Catastrophe Principle: Most Likely Escape Route

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



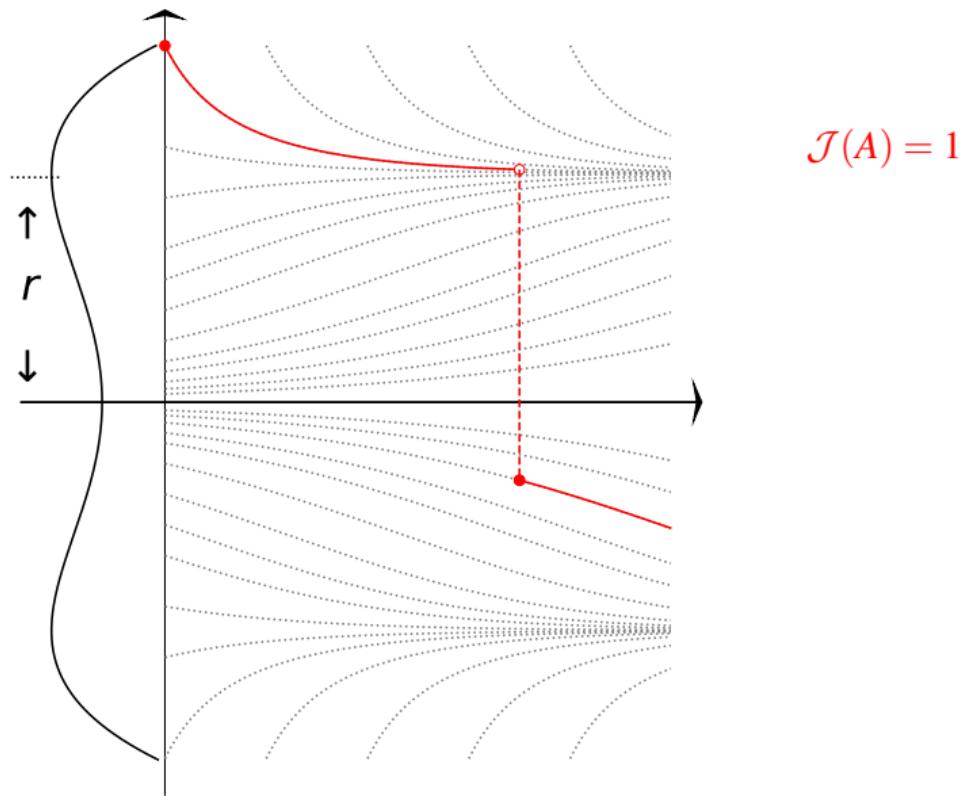
Catastrophe Principle: Most Likely Escape Route

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



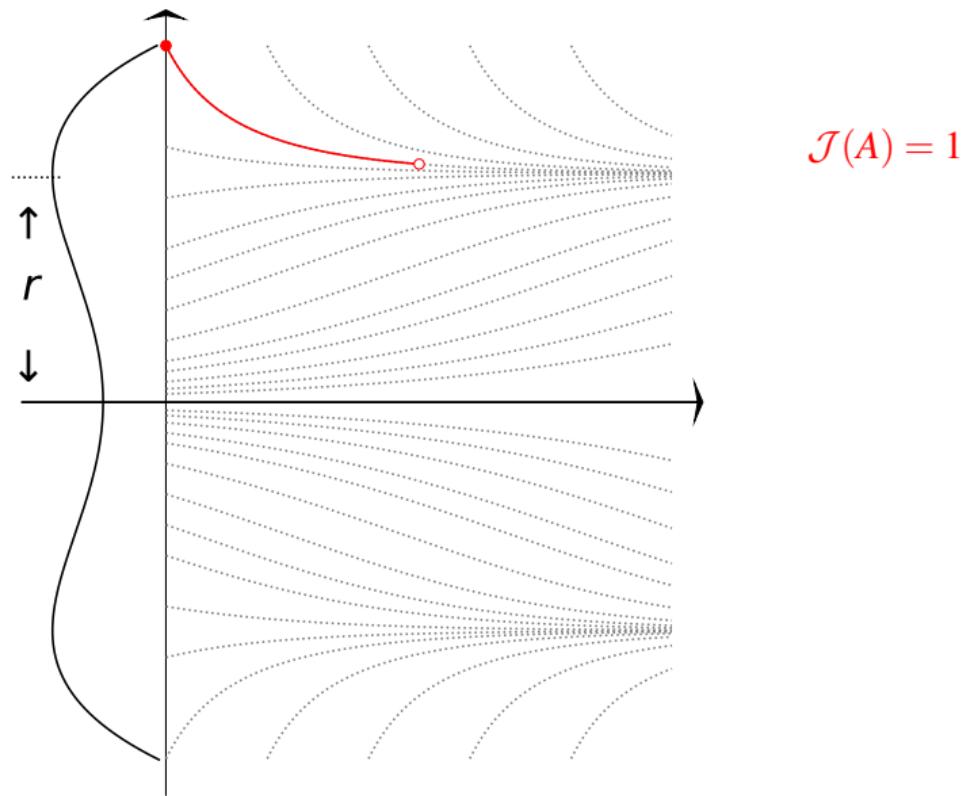
Catastrophe Principle: Most Likely Escape Route

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



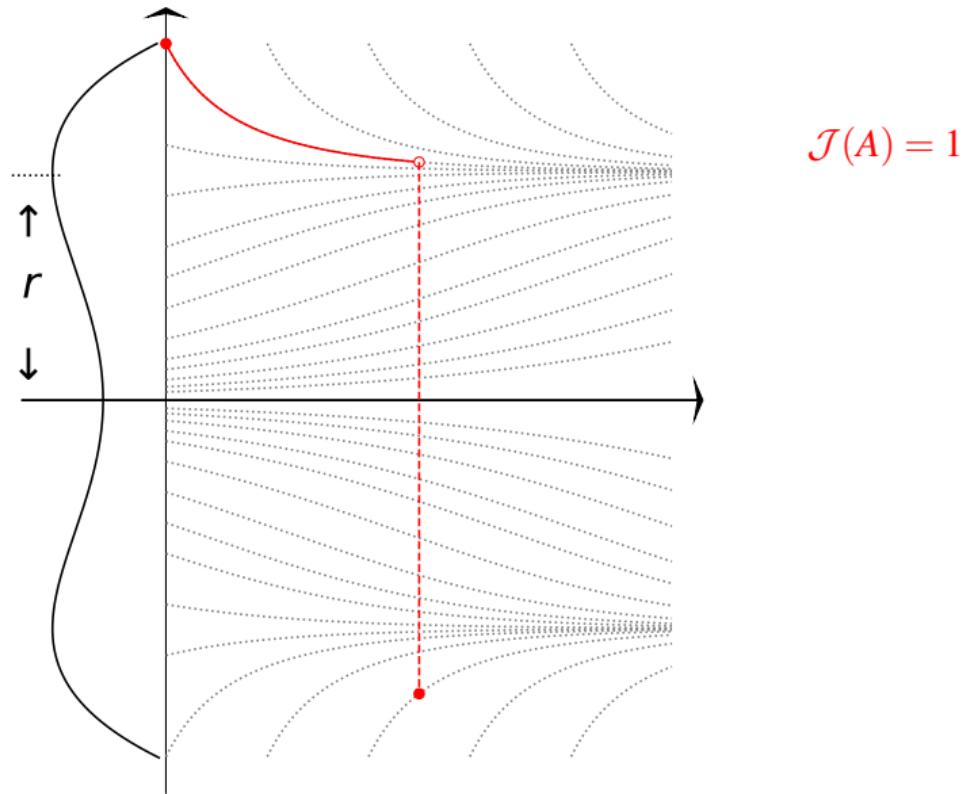
Catastrophe Principle: Most Likely Escape Route

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



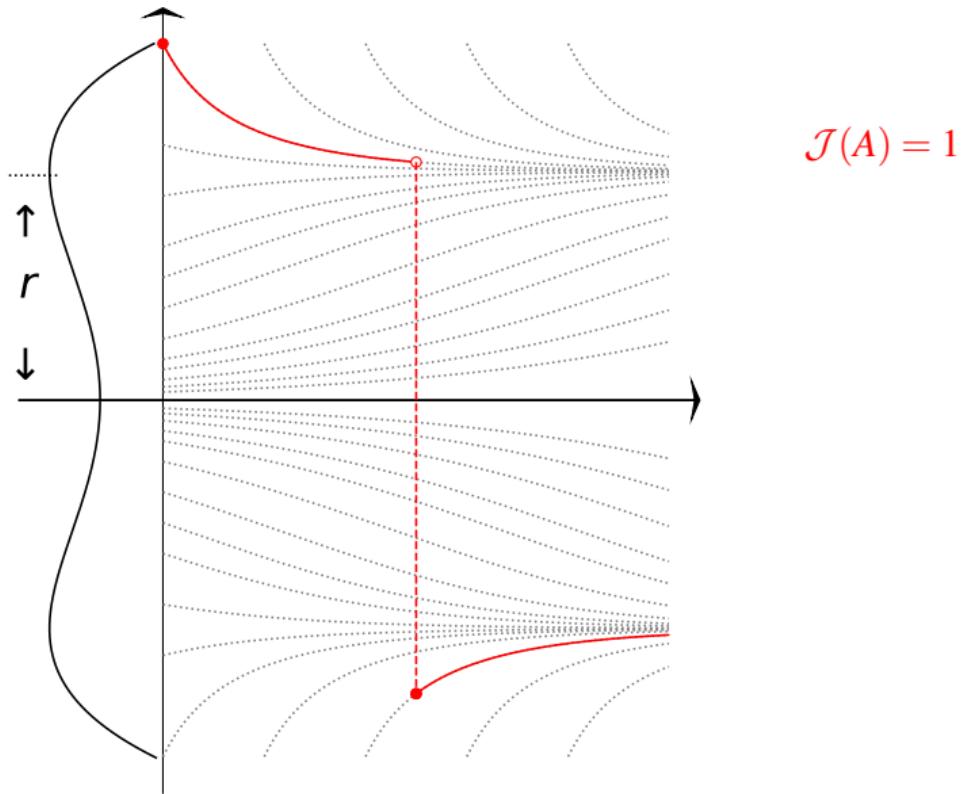
Catastrophe Principle: Most Likely Escape Route

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



Catastrophe Principle: Most Likely Escape Route

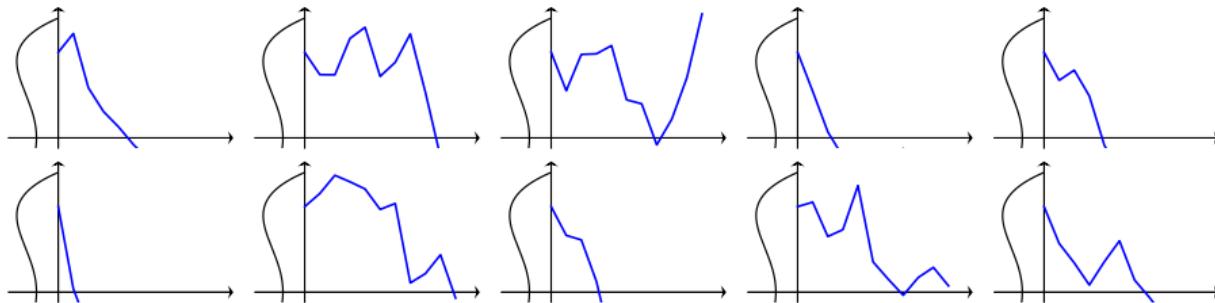
Most Likely Paths for $\{W^\eta \text{ escapes the local minimum}\}$



Catastrophe Principle Dictates SGD's Escape Route

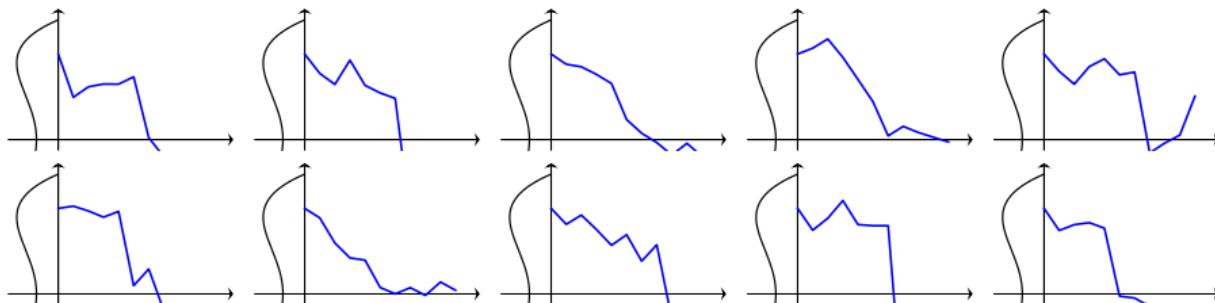
Trajectory of SGD X^η conditional on exit:

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



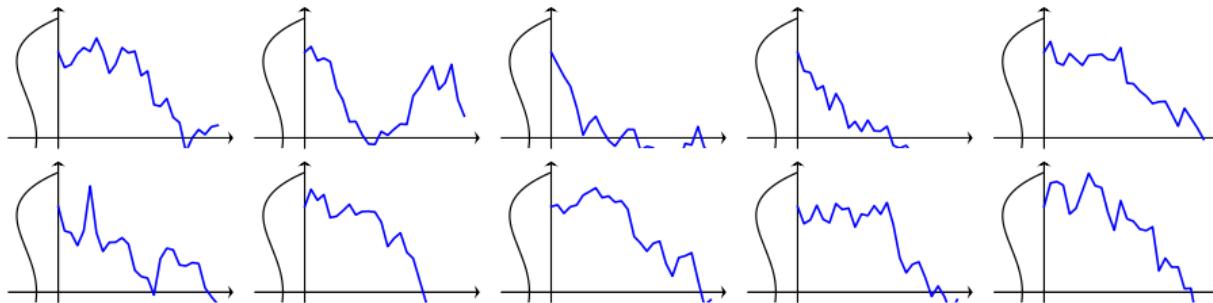
Trajectory of SGD X^η conditional on exit:

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



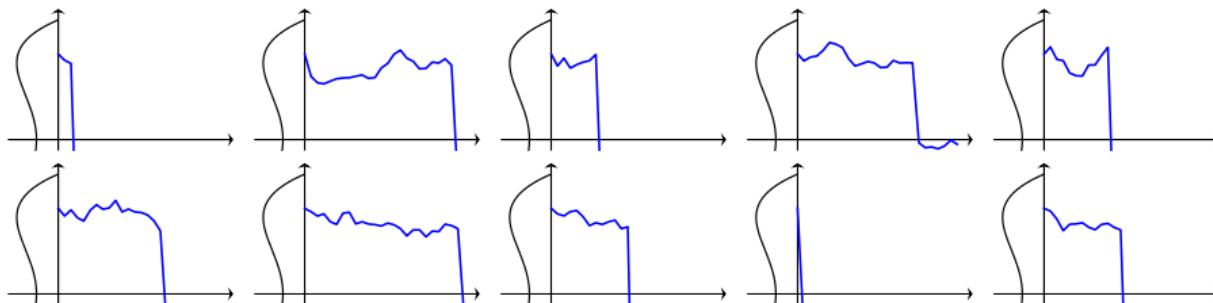
Catastrophe Principle Dictates SGD's Escape Route

Trajectory of SGD X^η conditional on exit:



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

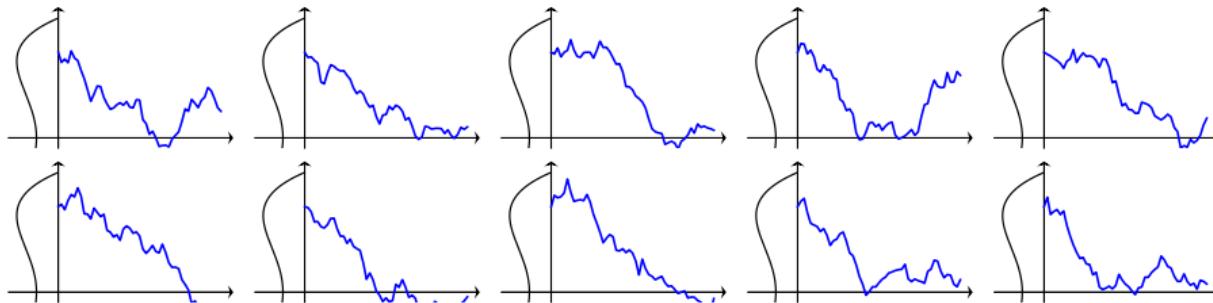
Trajectory of SGD X^η conditional on exit:



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

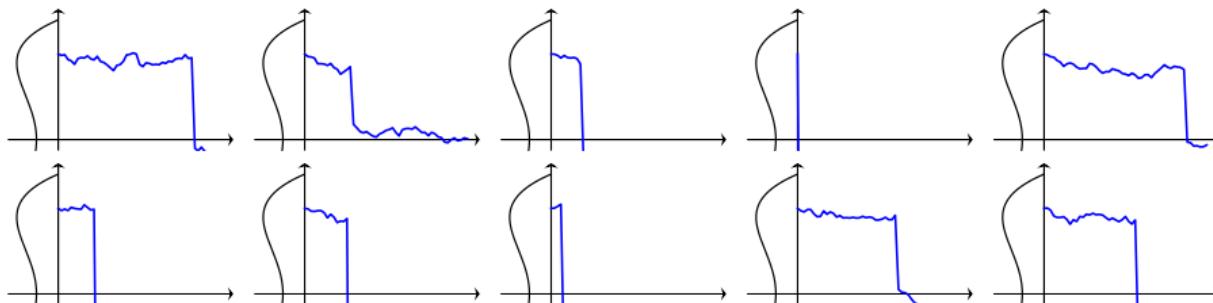
Catastrophe Principle Dictates SGD's Escape Route

Trajectory of SGD X^η conditional on exit:



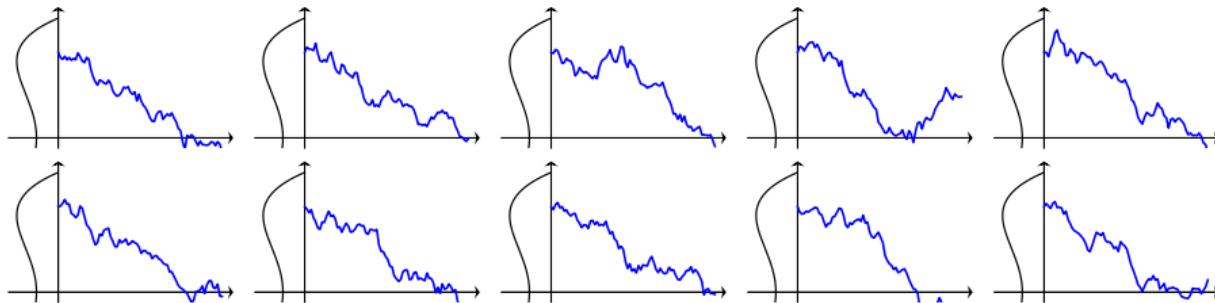
Trajectory of SGD X^η conditional on exit:

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



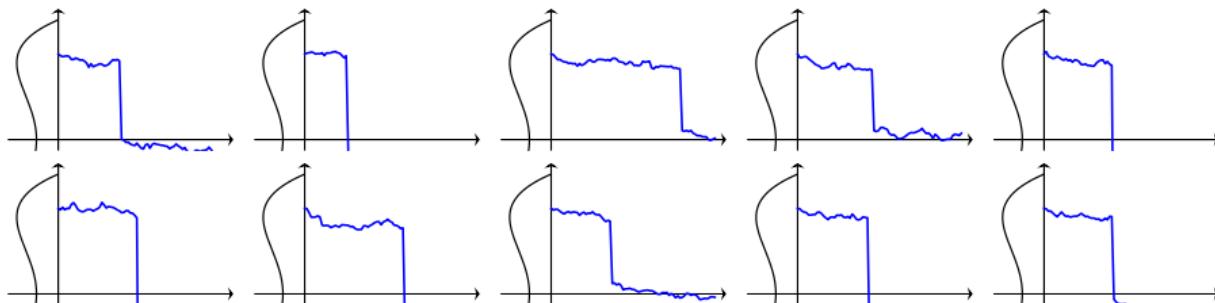
Catastrophe Principle Dictates SGD's Escape Route

Trajectory of SGD X^η conditional on exit:



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

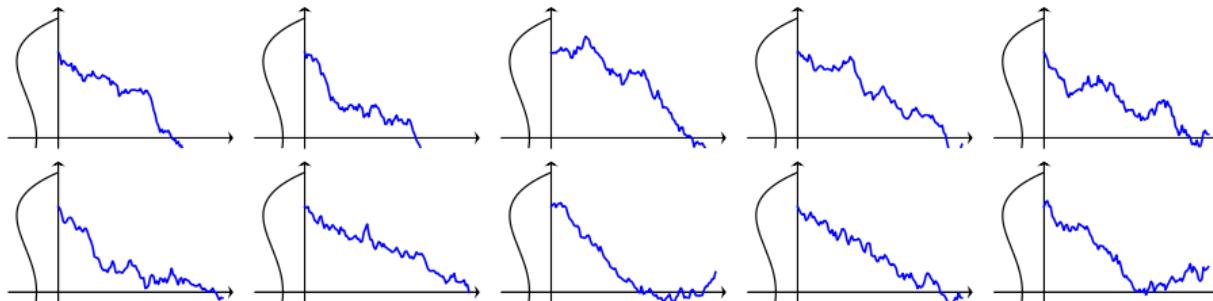
Trajectory of SGD X^η conditional on exit:



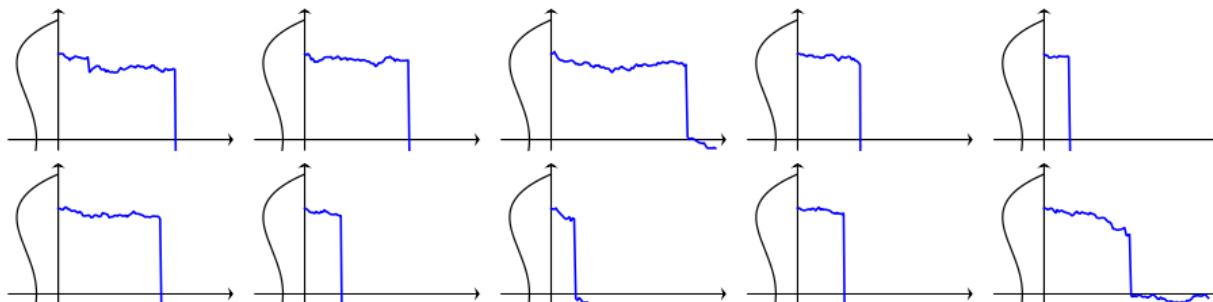
heavy-tailed noises with $\eta = 1/75$

Catastrophe Principle Dictates SGD's Escape Route

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

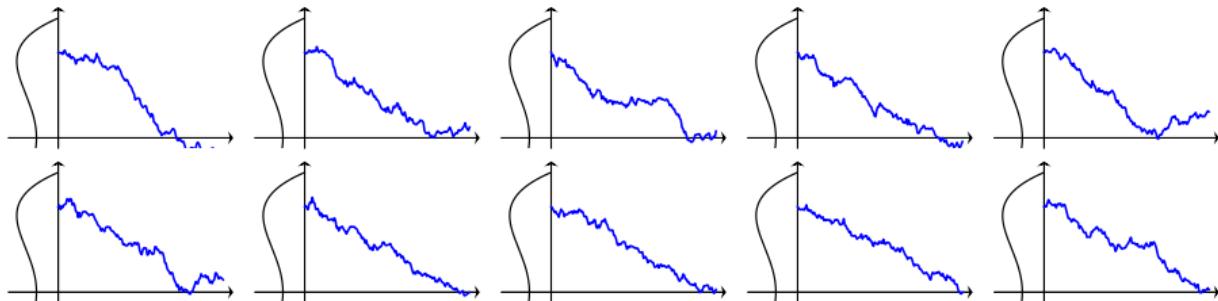


Trajectory of SGD X^η conditional on exit: **heavy-tailed** noises with $\eta = 1/100$

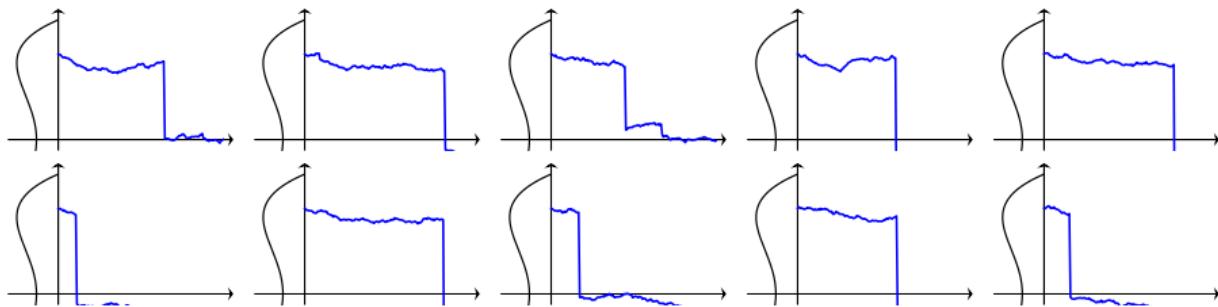


Catastrophe Principle Dictates SGD's Escape Route

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

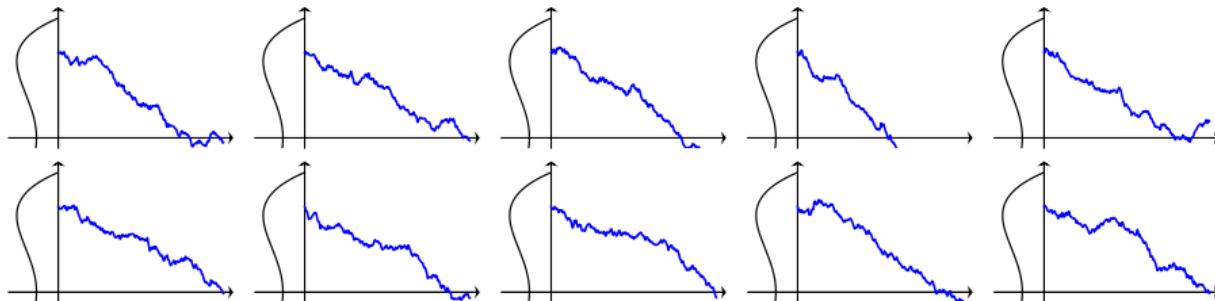


Trajectory of SGD X^η conditional on exit: **heavy-tailed** noises with $\eta = 1/150$

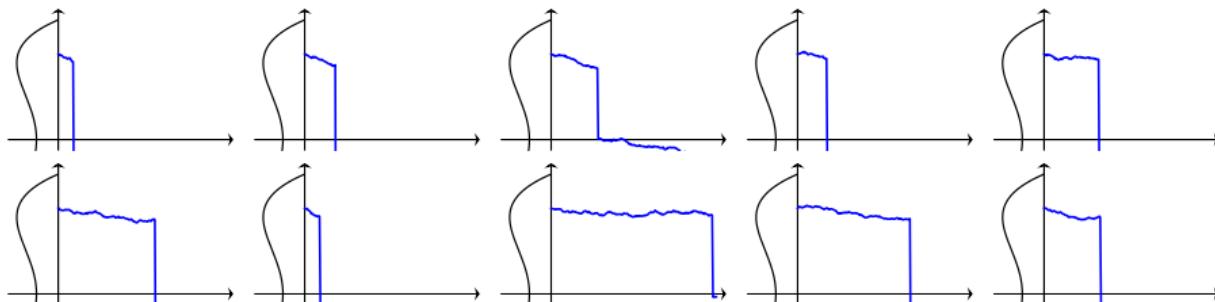


Catastrophe Principle Dictates SGD's Escape Route

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,

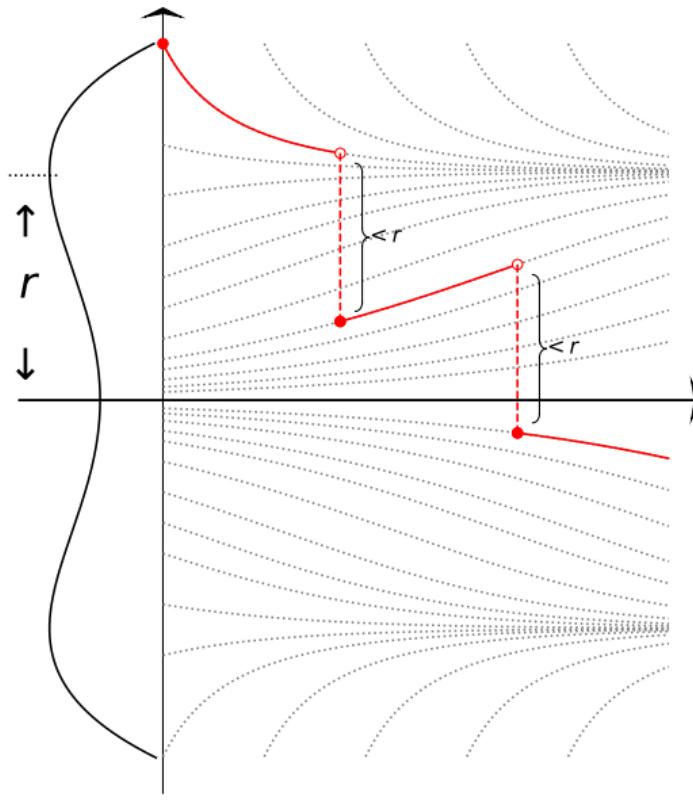
$$W^\eta(\cdot) \rightarrow w(\cdot) \quad \text{as} \quad \eta \rightarrow 0$$

where

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

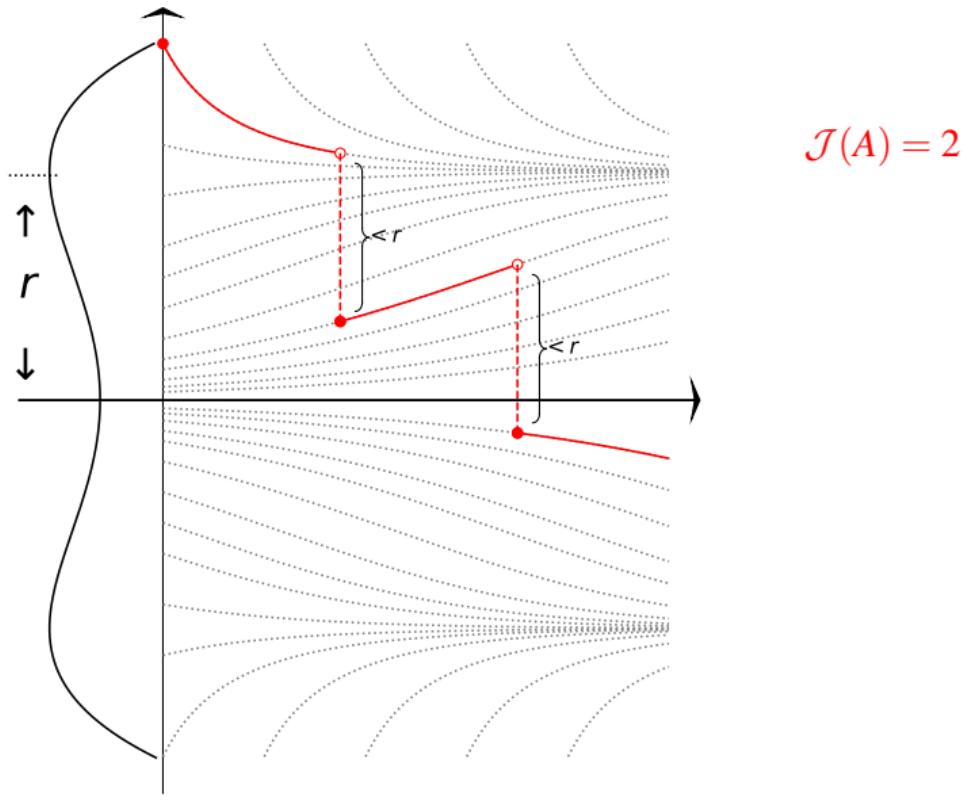
How does $\mathcal{J}(A)$ change?

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



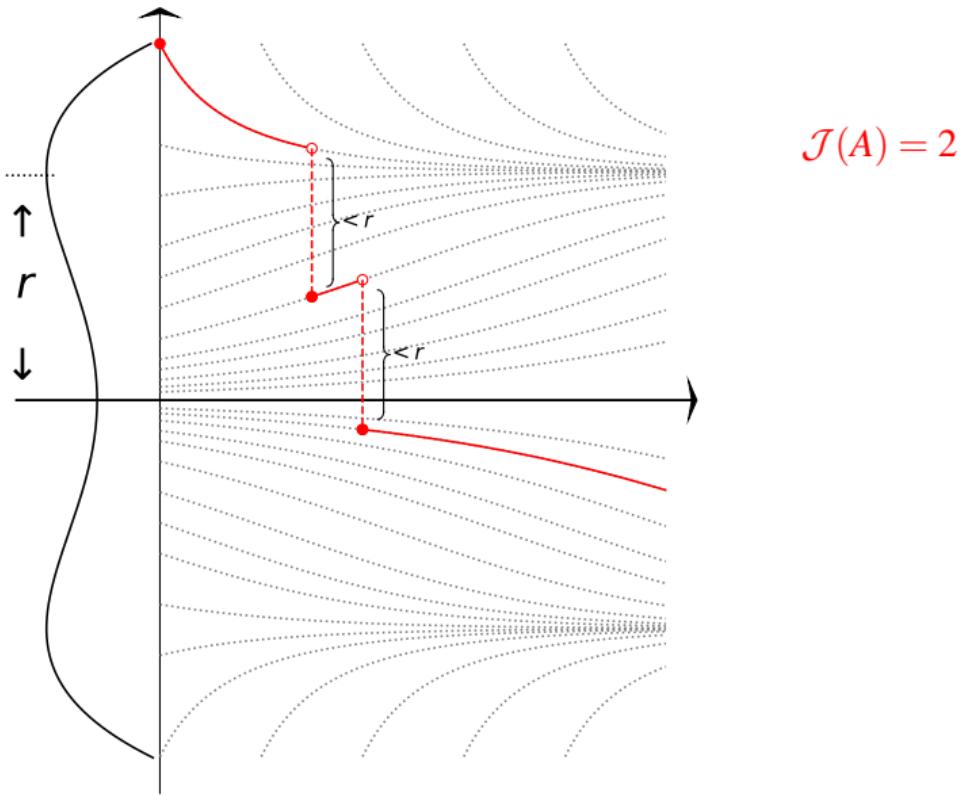
How does $\mathcal{J}(A)$ change?

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



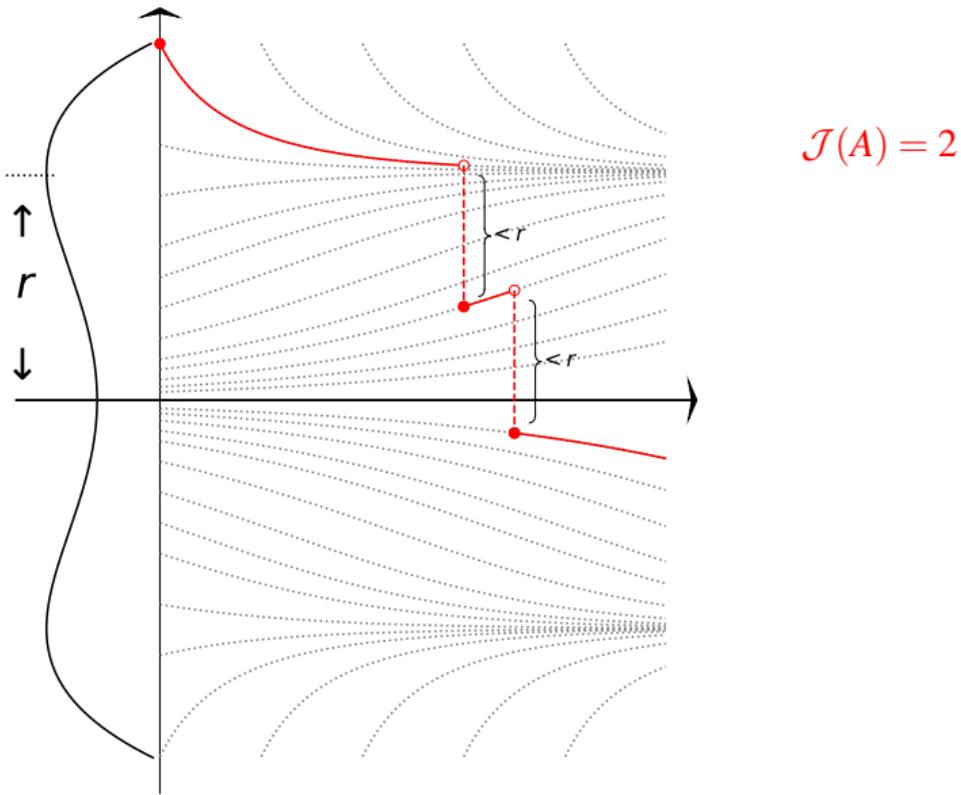
How does $\mathcal{J}(A)$ change?

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



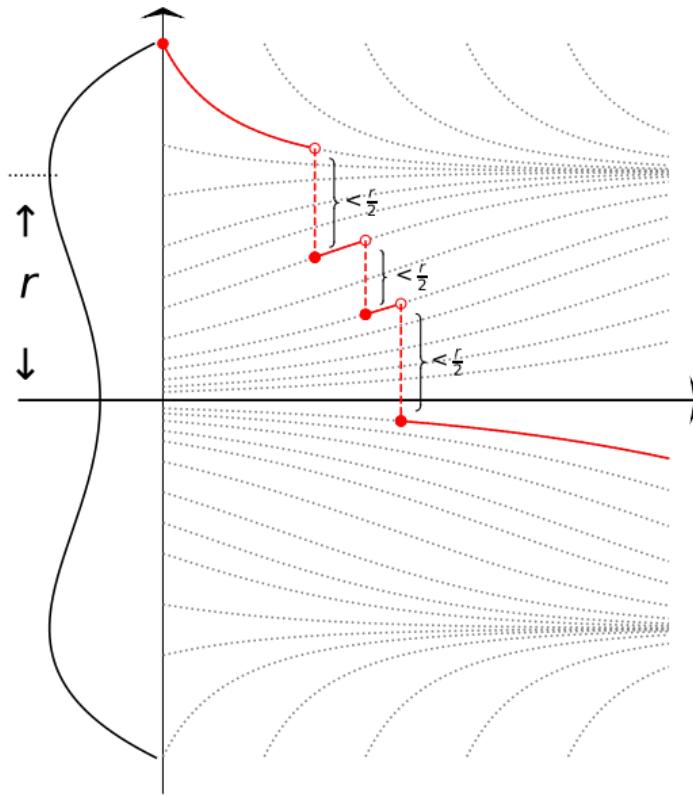
How does $\mathcal{J}(A)$ change?

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



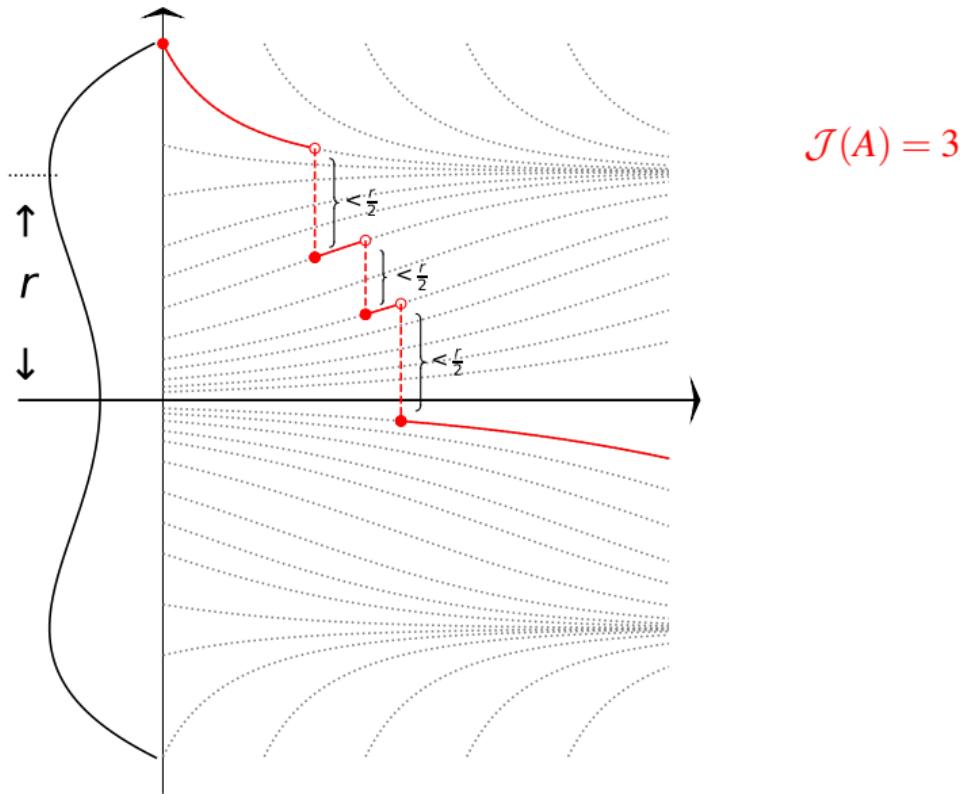
How does $\mathcal{J}(A)$ change?

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



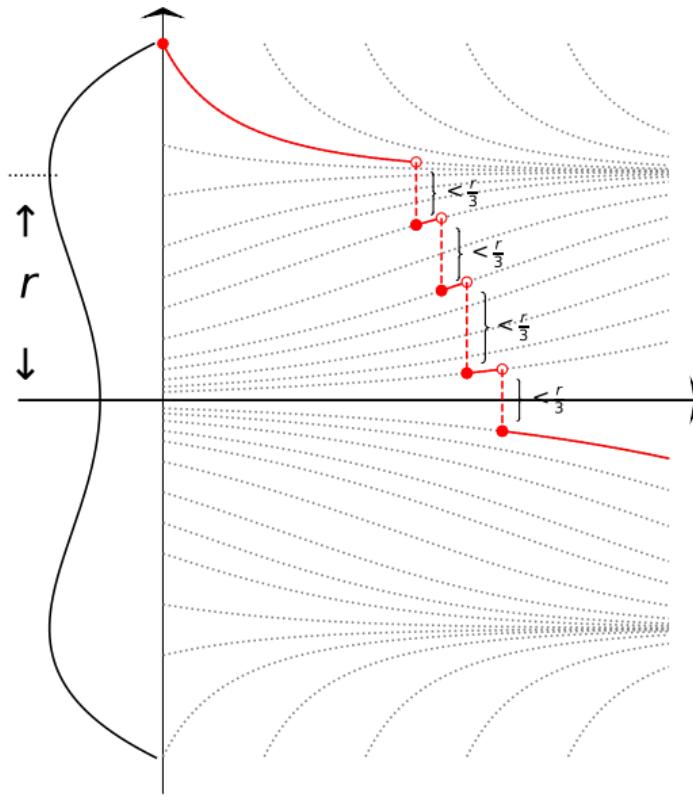
How does $\mathcal{J}(A)$ change?

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



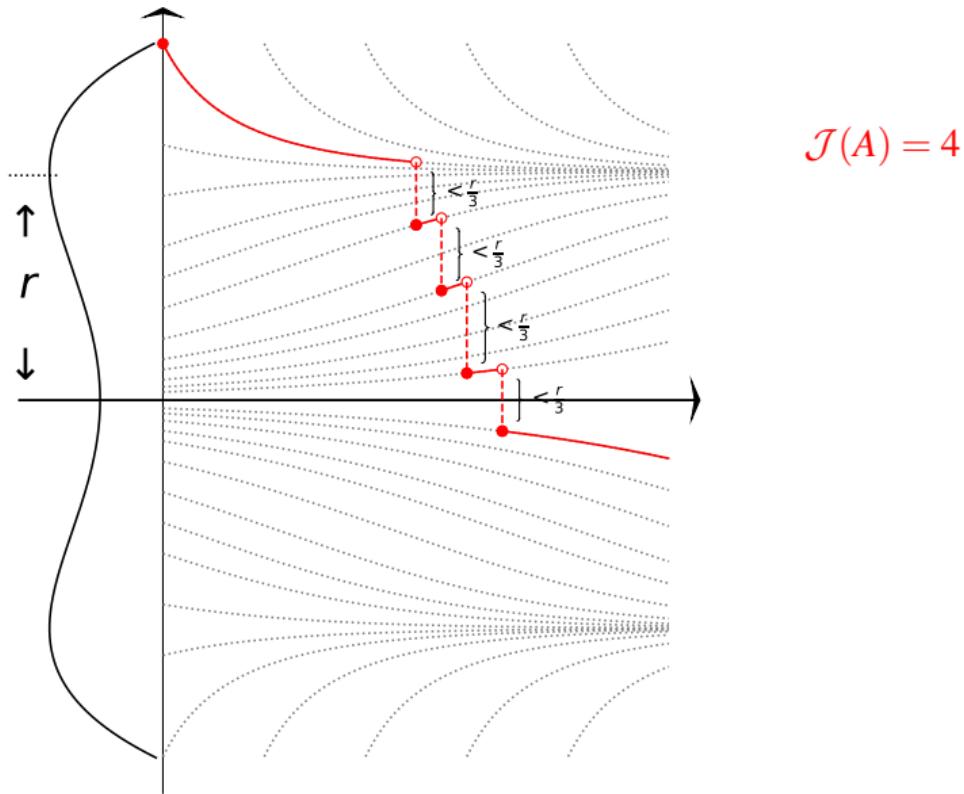
How does $\mathcal{J}(A)$ change?

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



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

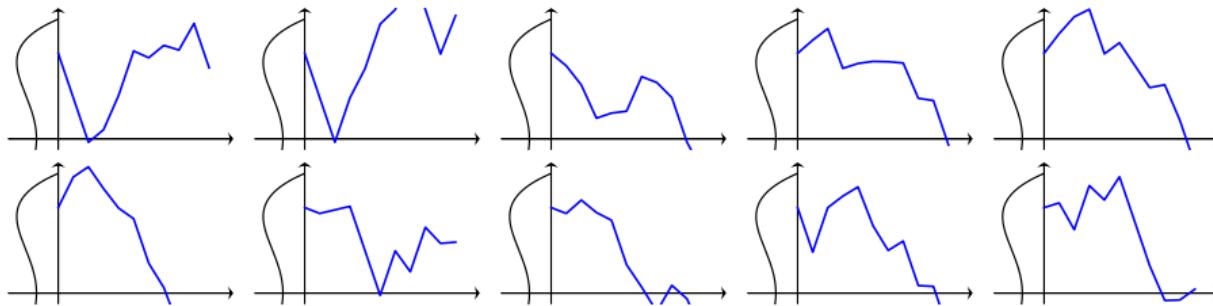
If $r \in (3c, 4c)$



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

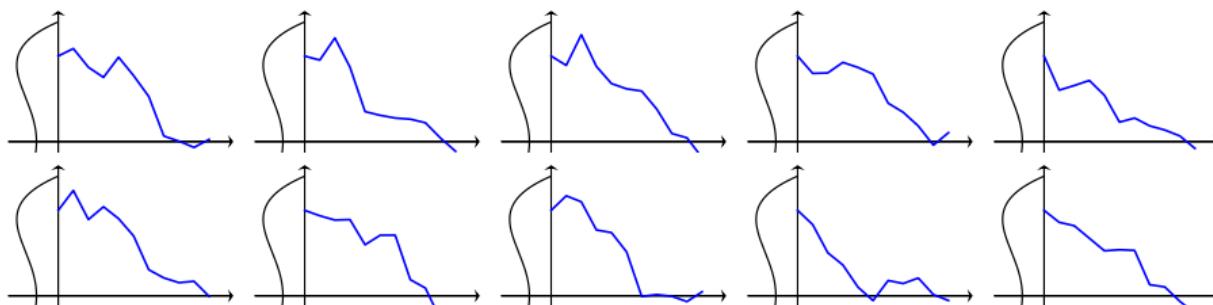
Trajectory of SGD X^η conditional on exit:

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



Trajectory of SGD X^η conditional on exit:

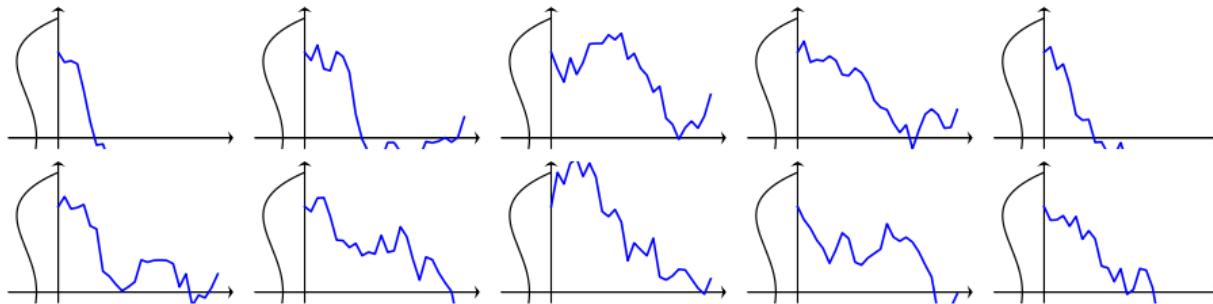
heavy-tailed noises with $\eta = 1/10$



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

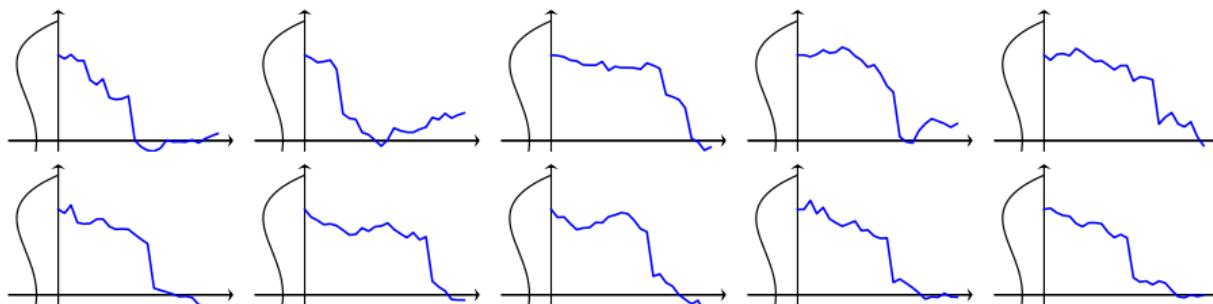
Trajectory of SGD X^η conditional on exit:

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



Trajectory of SGD X^η conditional on exit:

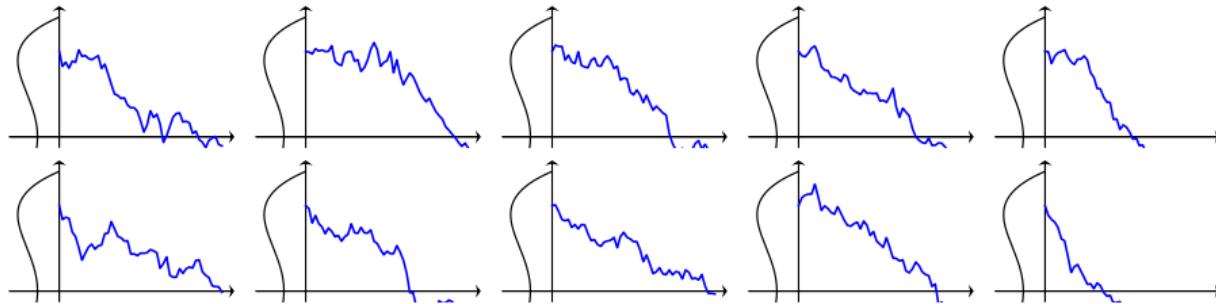
heavy-tailed noises with $\eta = 1/25$



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

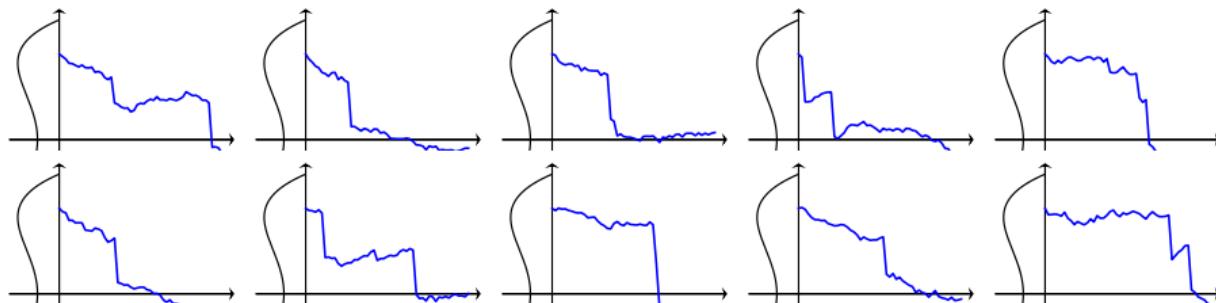
Trajectory of SGD X^η conditional on exit:

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



Trajectory of SGD X^η conditional on exit:

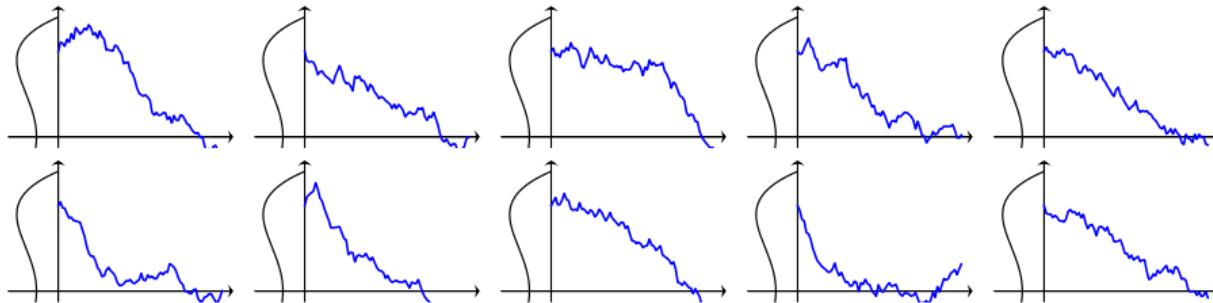
heavy-tailed noises with $\eta = 1/10$



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

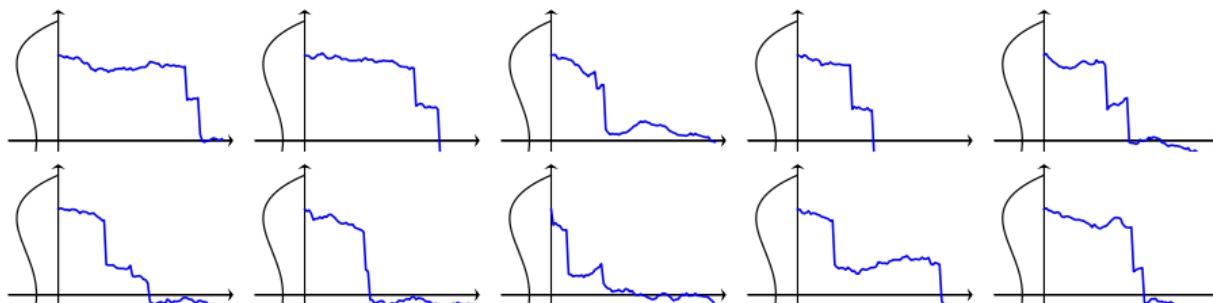
Trajectory of SGD X^η conditional on exit:

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



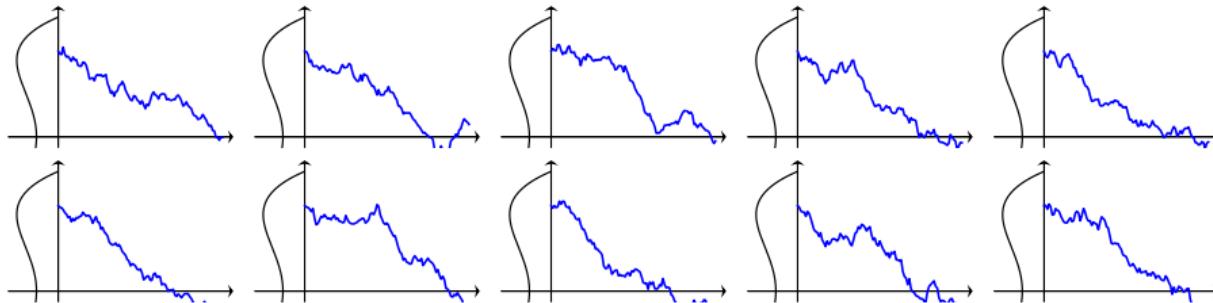
Trajectory of SGD X^η conditional on exit:

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



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

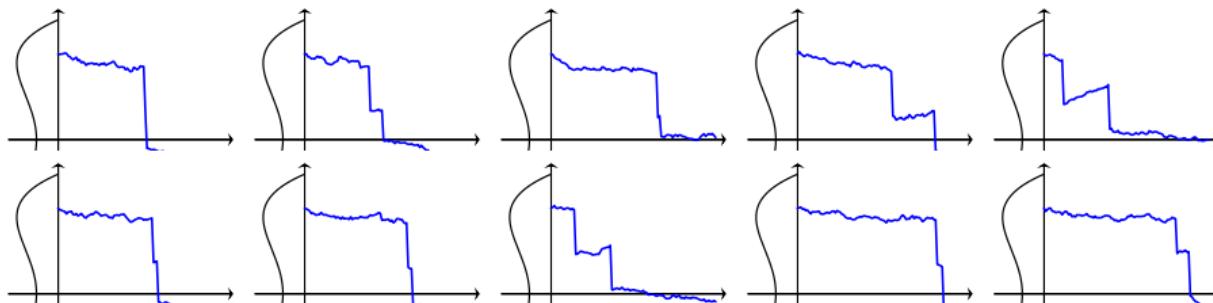
Trajectory of SGD X^η conditional on exit:



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

Trajectory of SGD X^η conditional on exit:

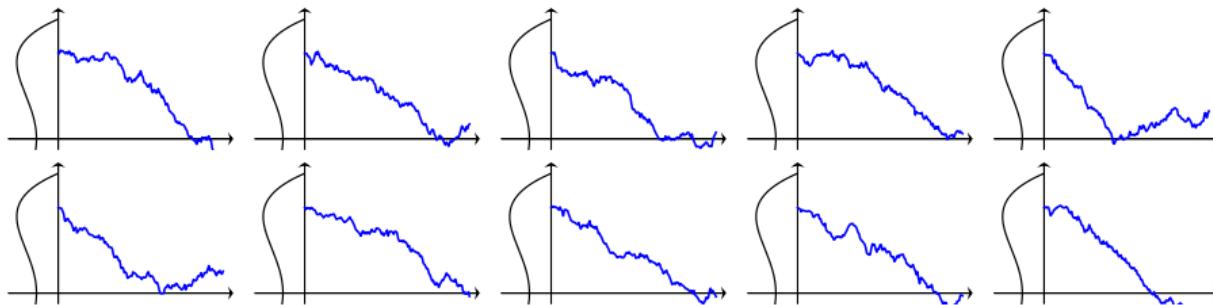
heavy-tailed noises with $\eta = 1/100$



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

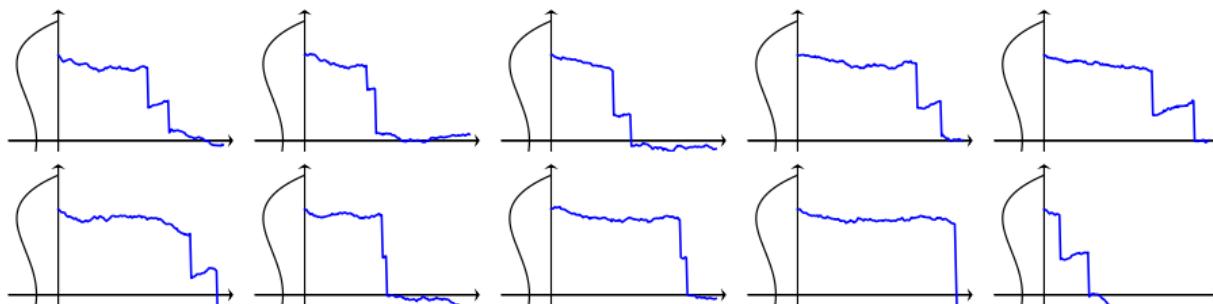
Trajectory of SGD X^η conditional on exit:

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



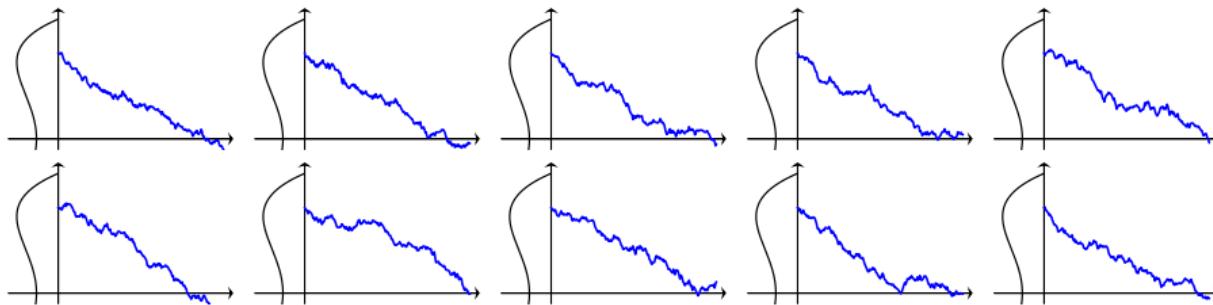
Trajectory of SGD X^η conditional on exit:

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



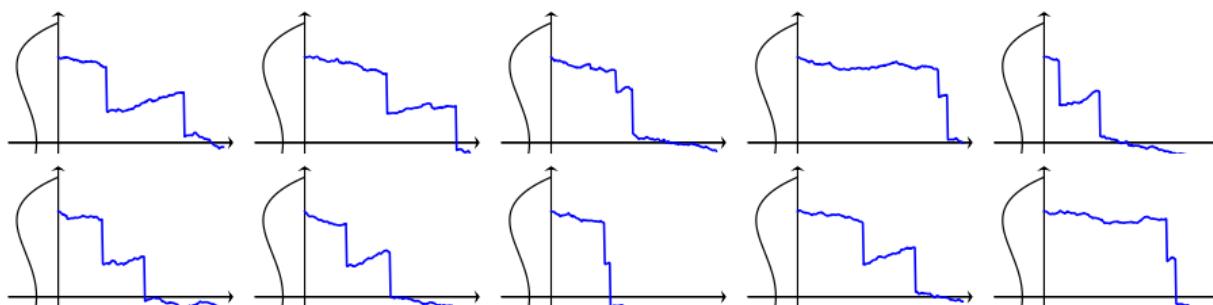
SGD's Escaping Route under Gradient Clipping with $\mathcal{J}(A) = 2$

Trajectory of SGD X^η conditional on exit:



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

Trajectory of SGD X^η conditional on exit:

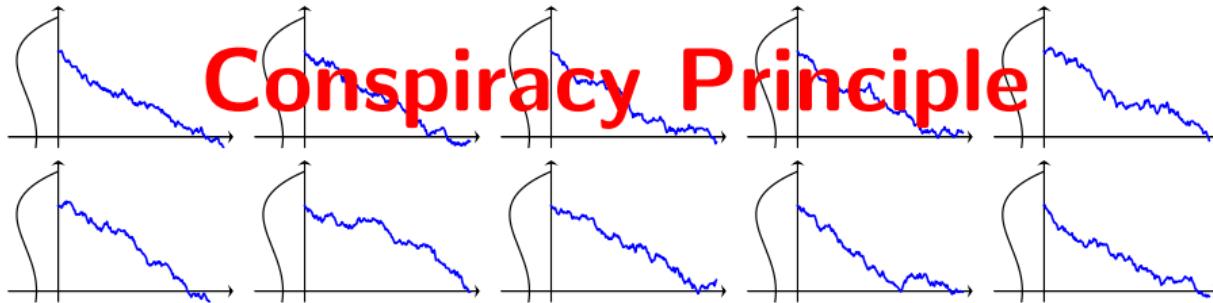


heavy-tailed noises with $\eta = 1/200$

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

Trajectory of SGD X^η conditional on exit:

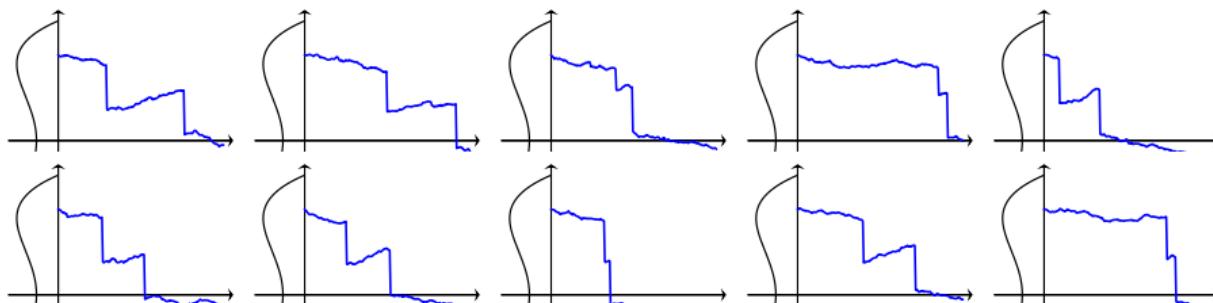
light-tailed noises with $\eta = 1/200$



Conspiracy Principle

Trajectory of SGD X^η conditional on exit:

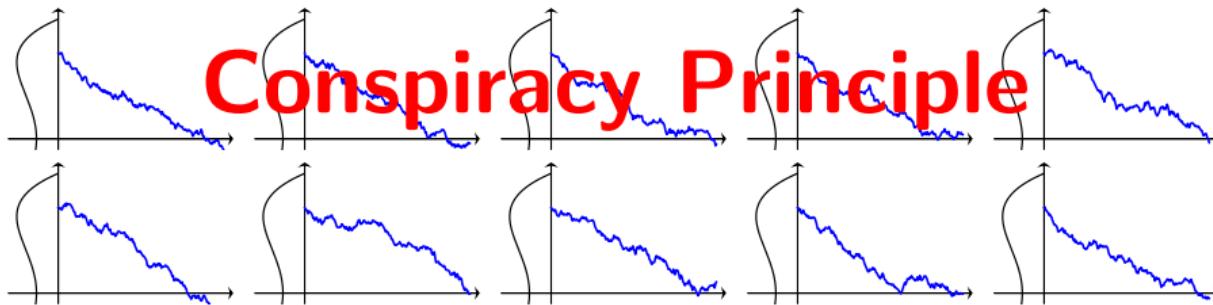
heavy-tailed noises with $\eta = 1/200$



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

Trajectory of SGD X^η conditional on exit:

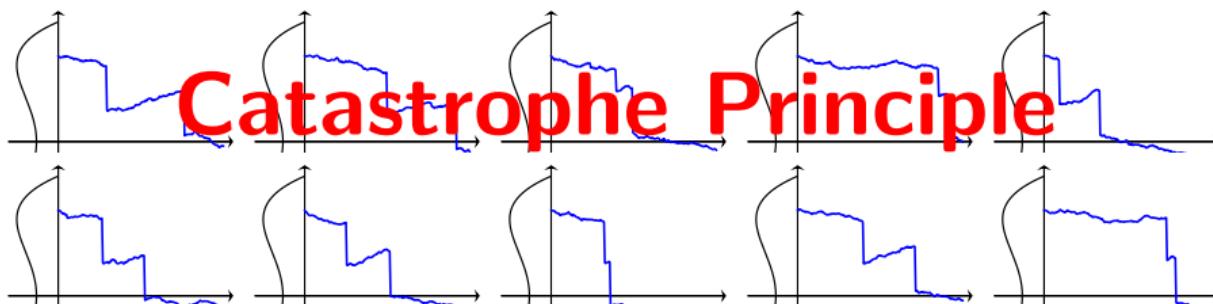
light-tailed noises with $\eta = 1/200$



Conspiracy Principle

Trajectory of SGD X^η conditional on exit:

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



Catastrophe Principle