

Distributional Reinforcement Learning

Quantile Regression - Implicit Quantile Networks

February 2019

Pantelis R. Vlachas

Seminar on Deep Reinforcement Learning
ETH Zurich

Distribution over returns is **NOT a new idea...**

[1] Sobel, L.M. 1982

[2] Morimura, T. et. al. 2010

But achieved **state-of-the-art results...**

C51 [3] Bellemare, M.G. et. al. 2017

And triggered a lot of **discussion...**

[4] Barth-Maron, G. et. al. 2018

[5] Hessel, M. et. al. 2018

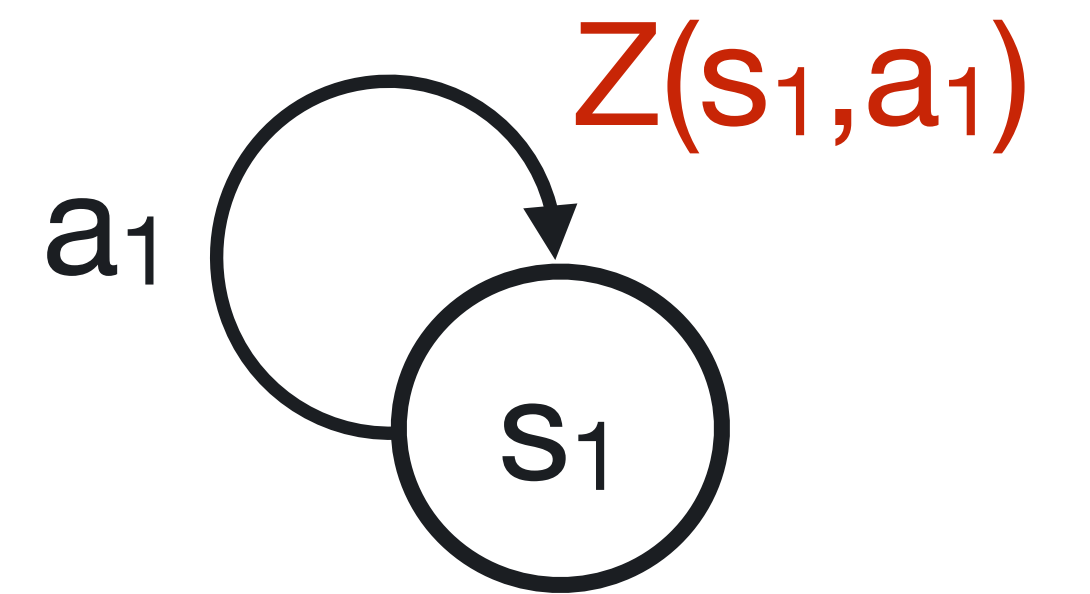
[6] Gruslys, A. et. al. 2018

[7] Rowland, M. et. al. 2018

In this session...

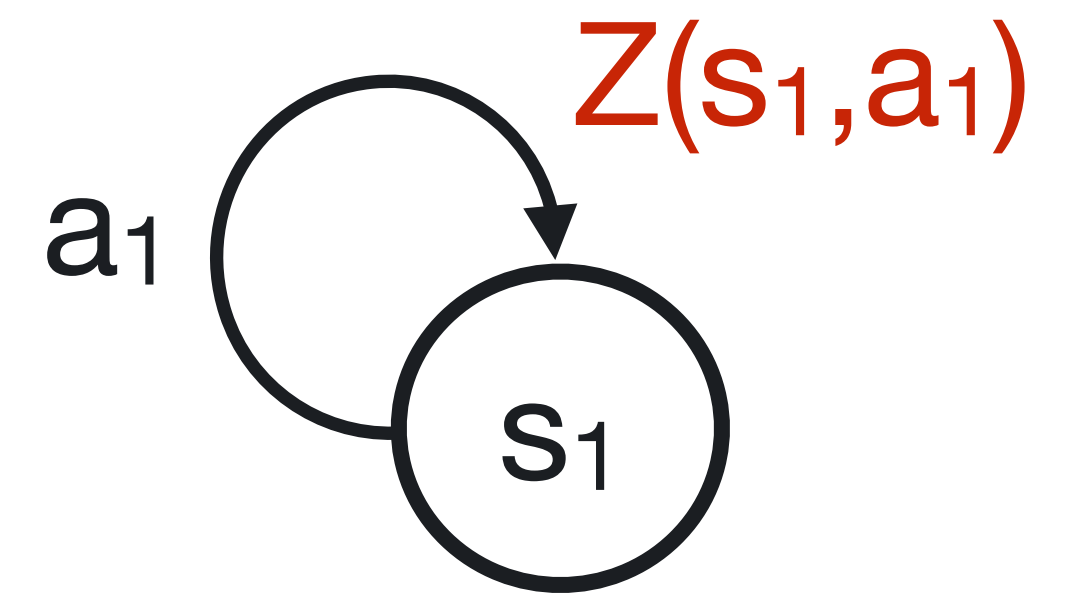
Distributional Reinforcement Learning with (Implicit) Quantile Regression

How to model distributions over returns ?



How to model distributions over returns ?

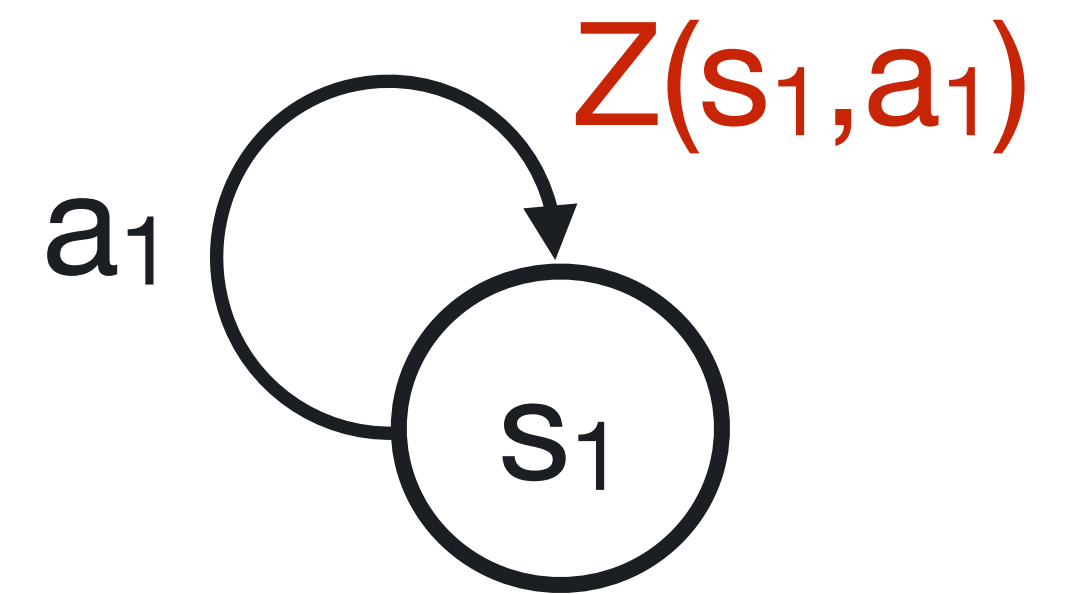
1. **Categorical** of PDF or CDF



How to model distributions over returns ?

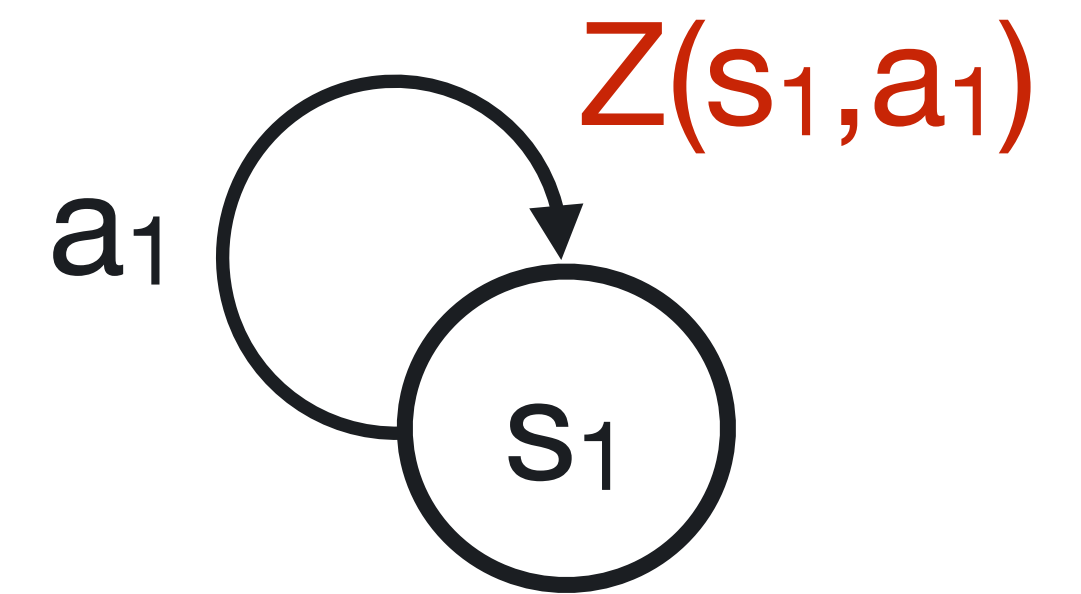
1. **Categorical** of PDF or CDF

2. **Quantiles** of Inverse CDF



How to model distributions over returns ?

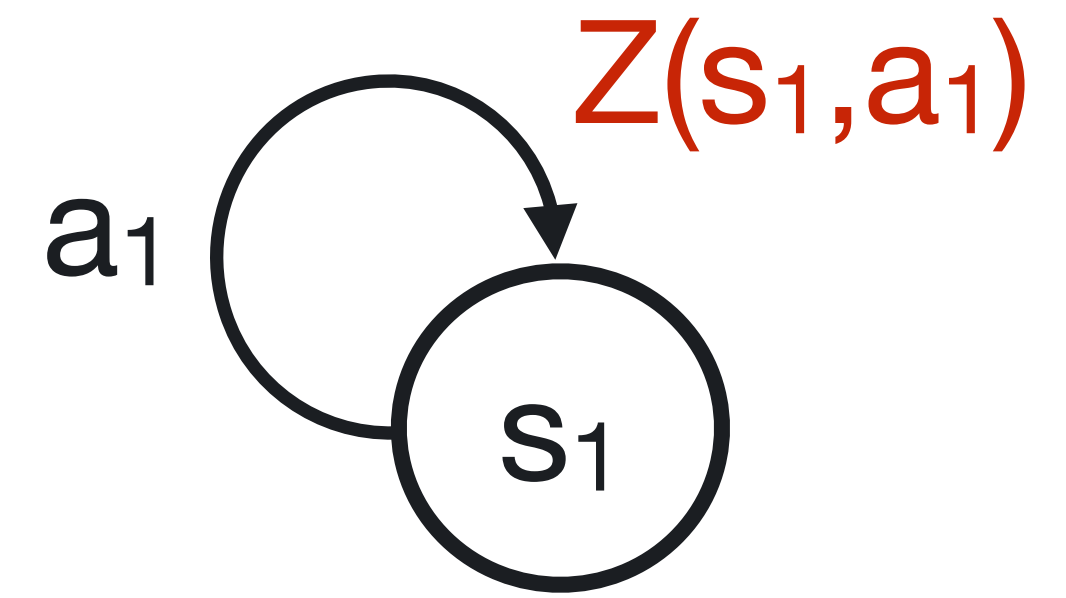
1. **Categorical** of PDF or CDF



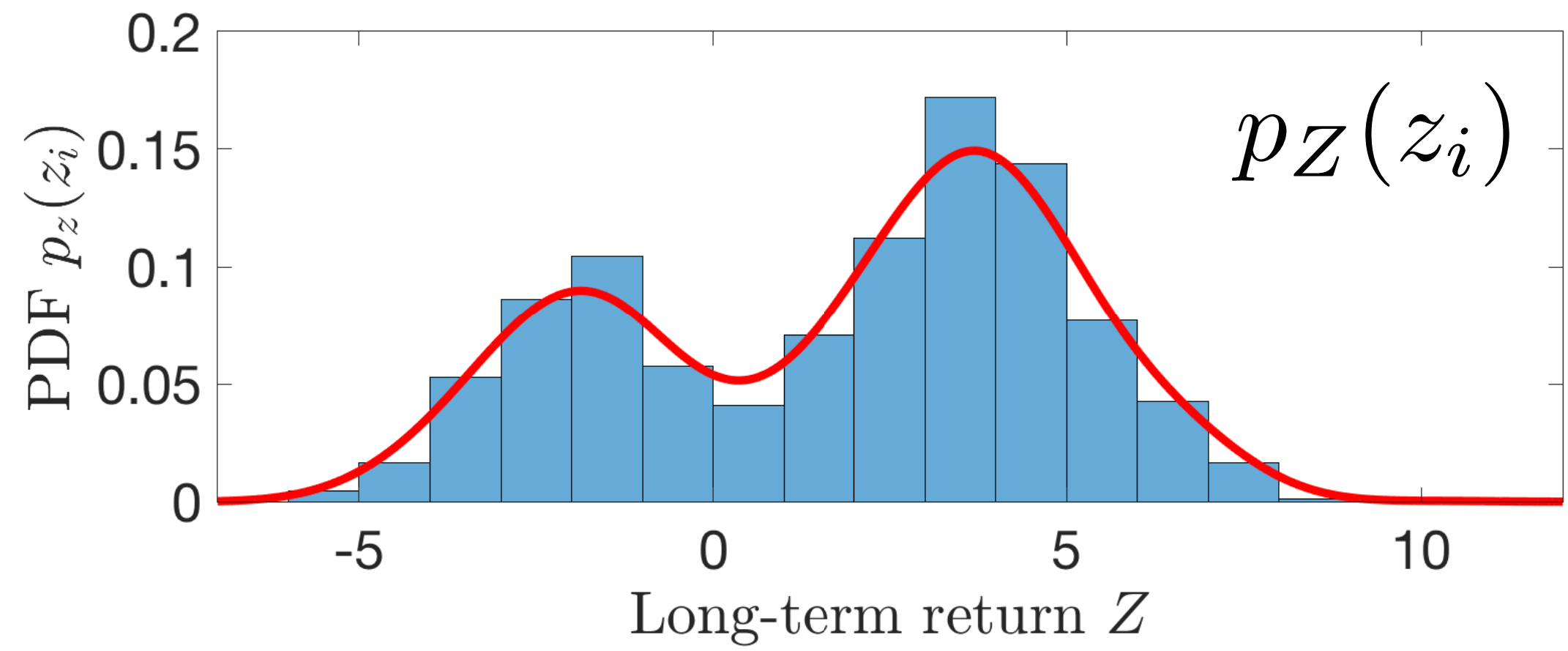
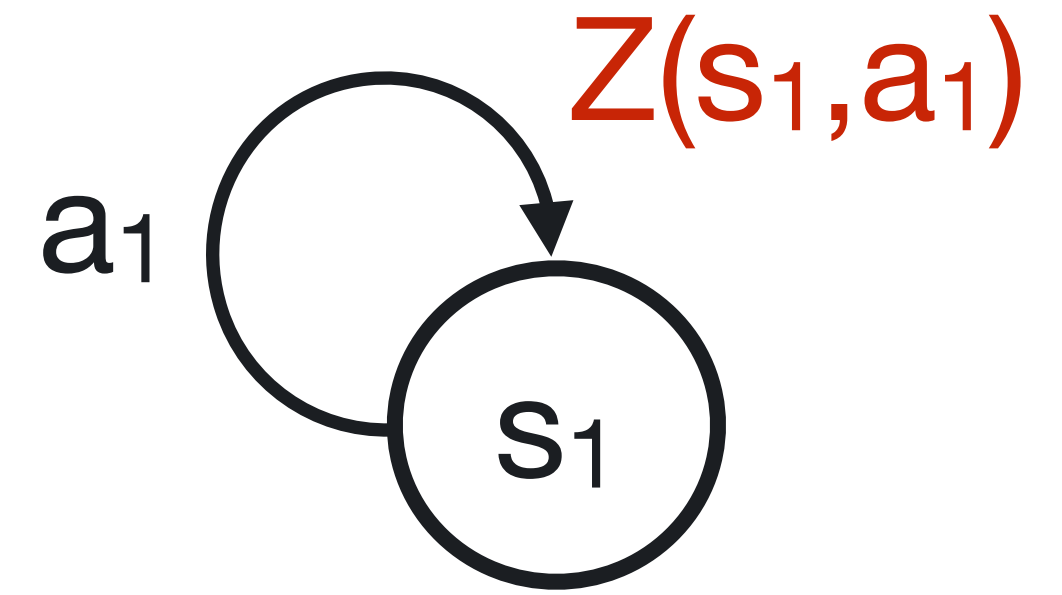
2. **Quantiles** of Inverse CDF

3. **Implicit Quantiles** of Inverse CDF

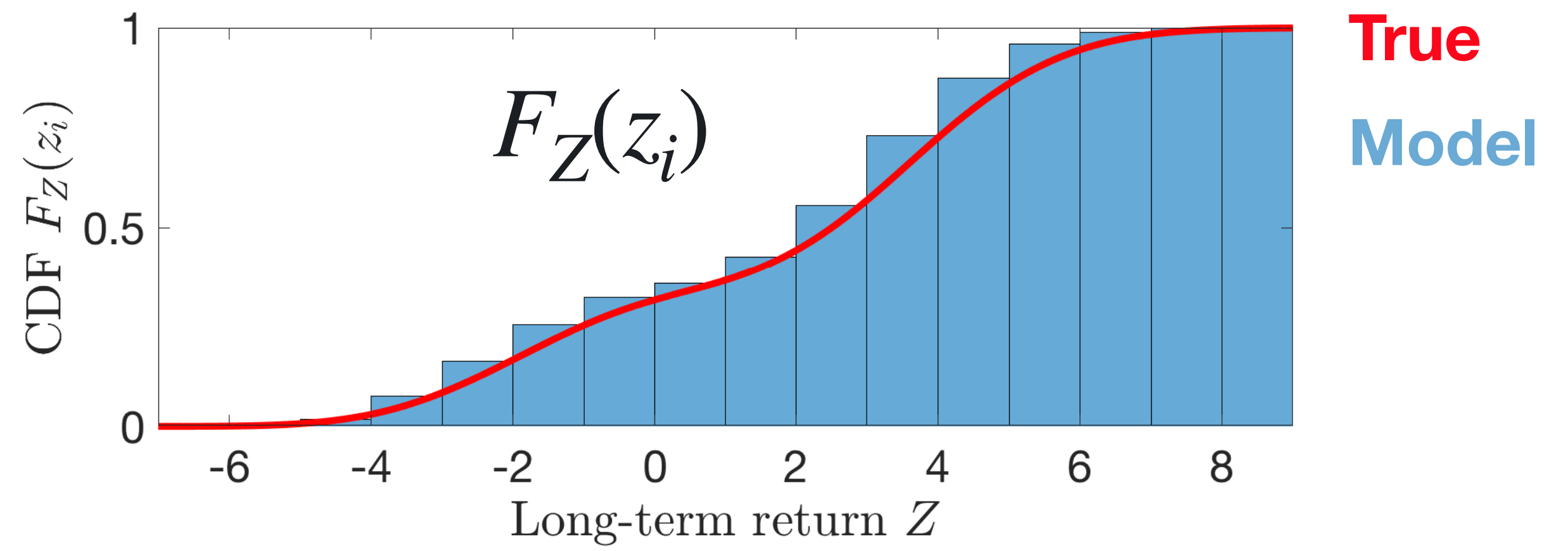
1. Categorical



1. Categorical

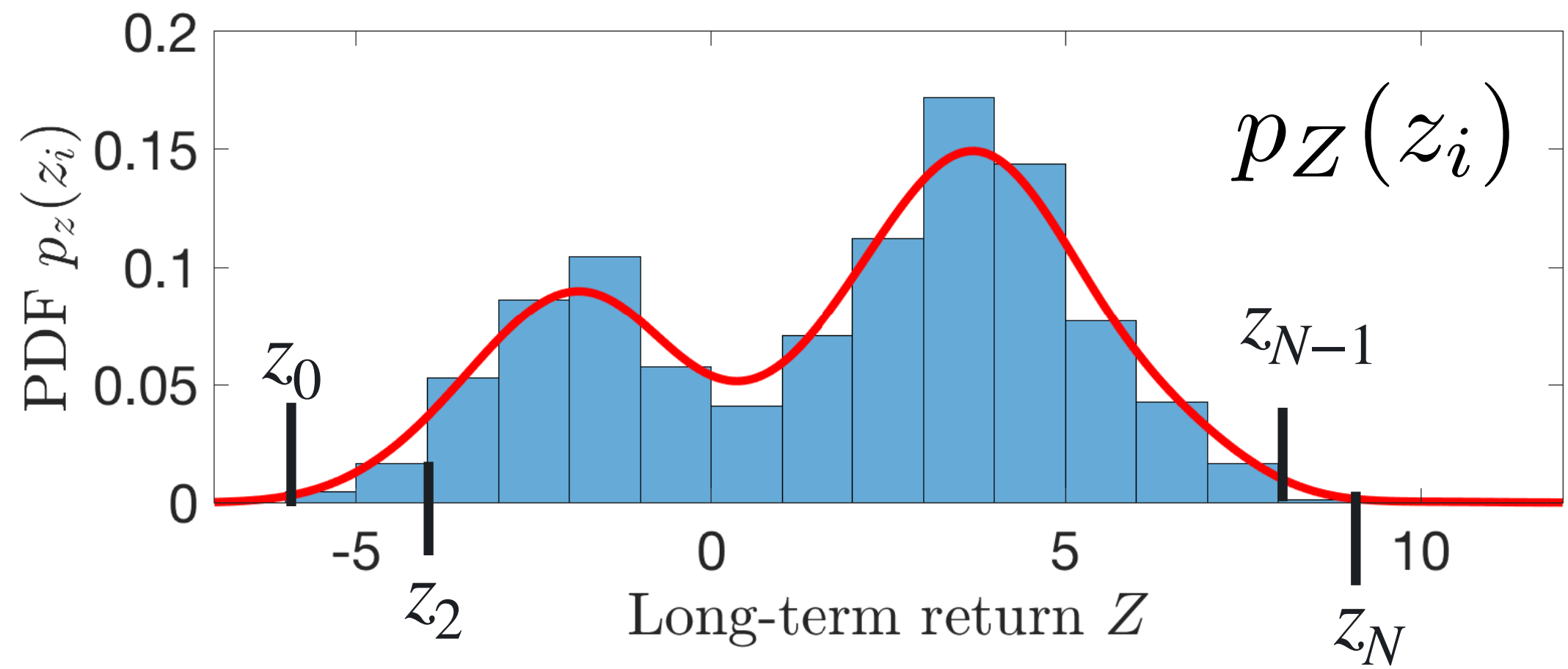
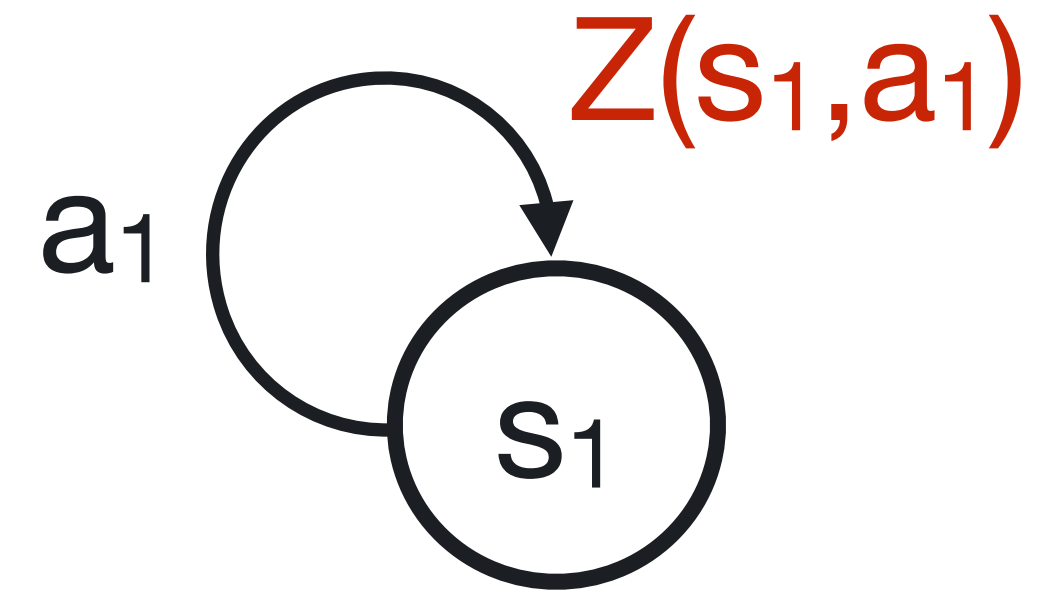


OR

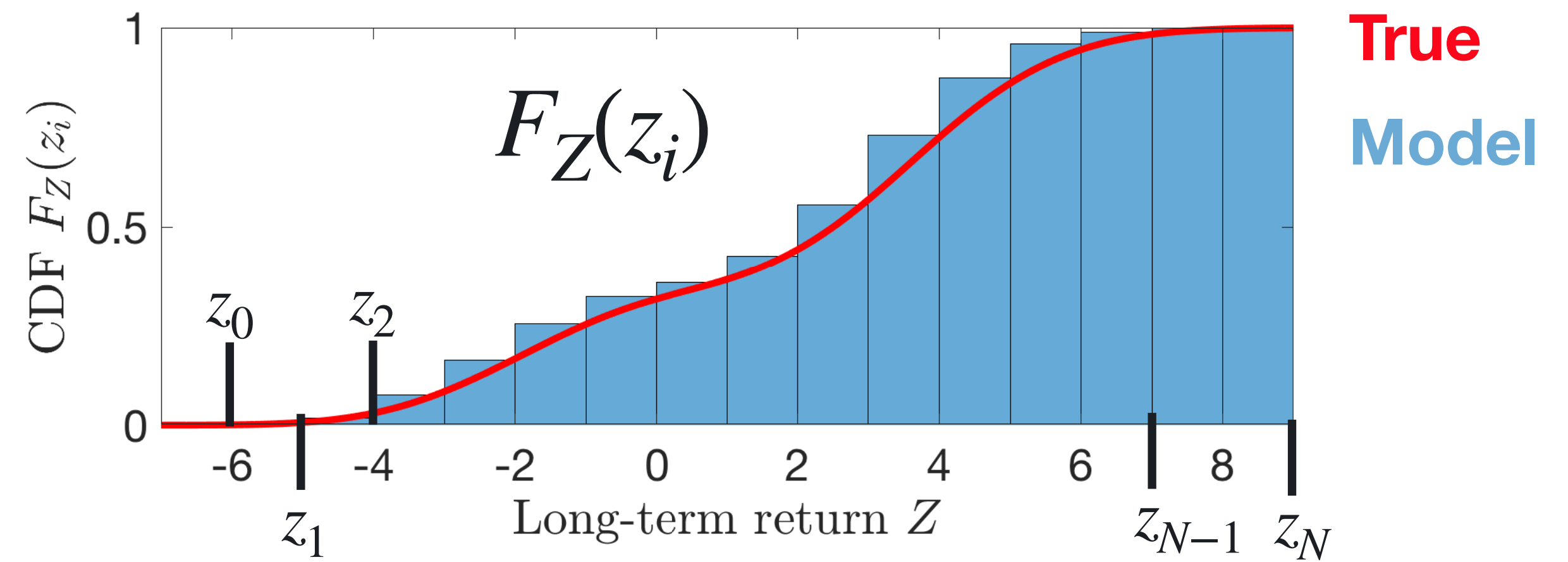


1. Categorical

- Fixed support bins $z_0, z_1, z_2, \dots, z_N$

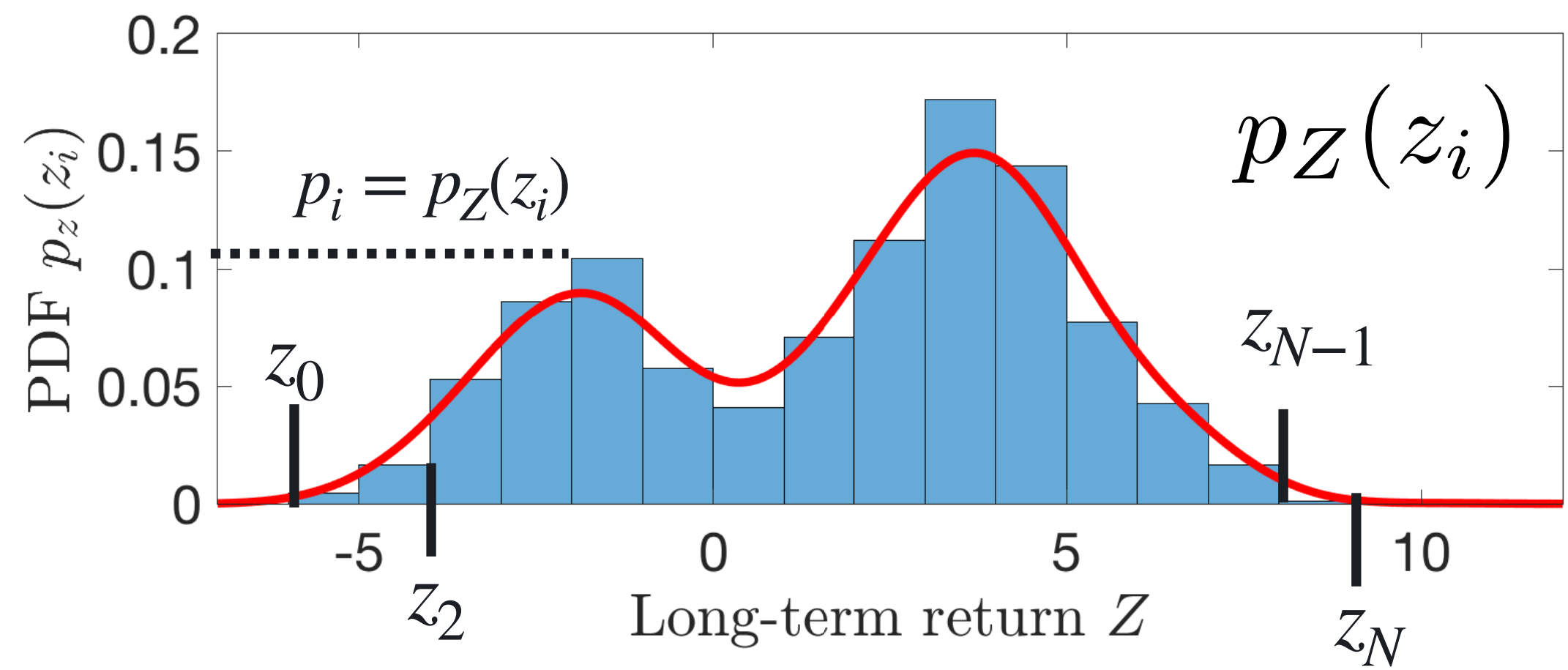
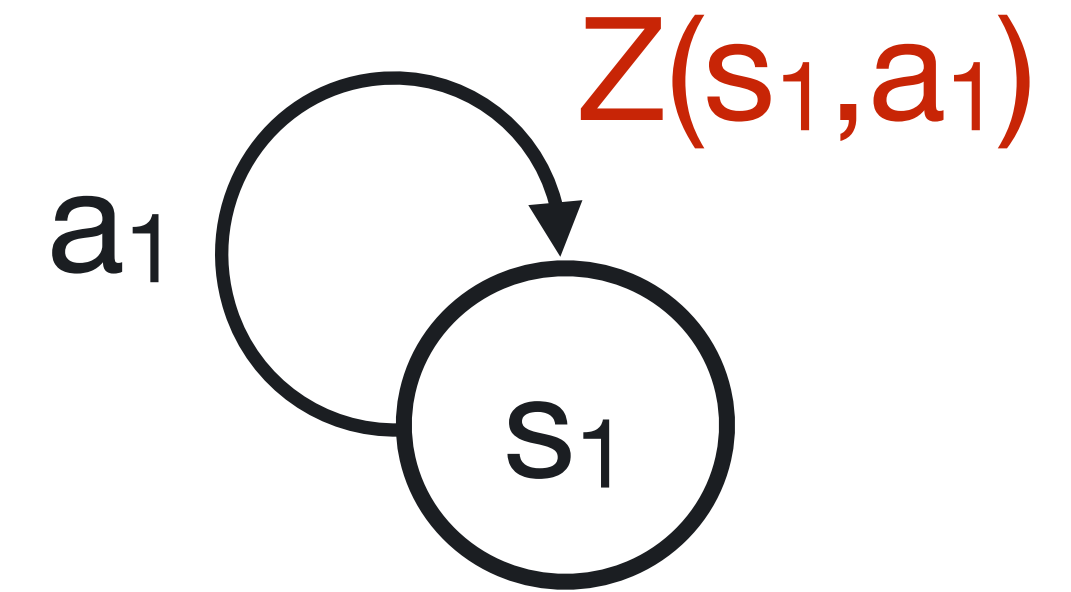


OR

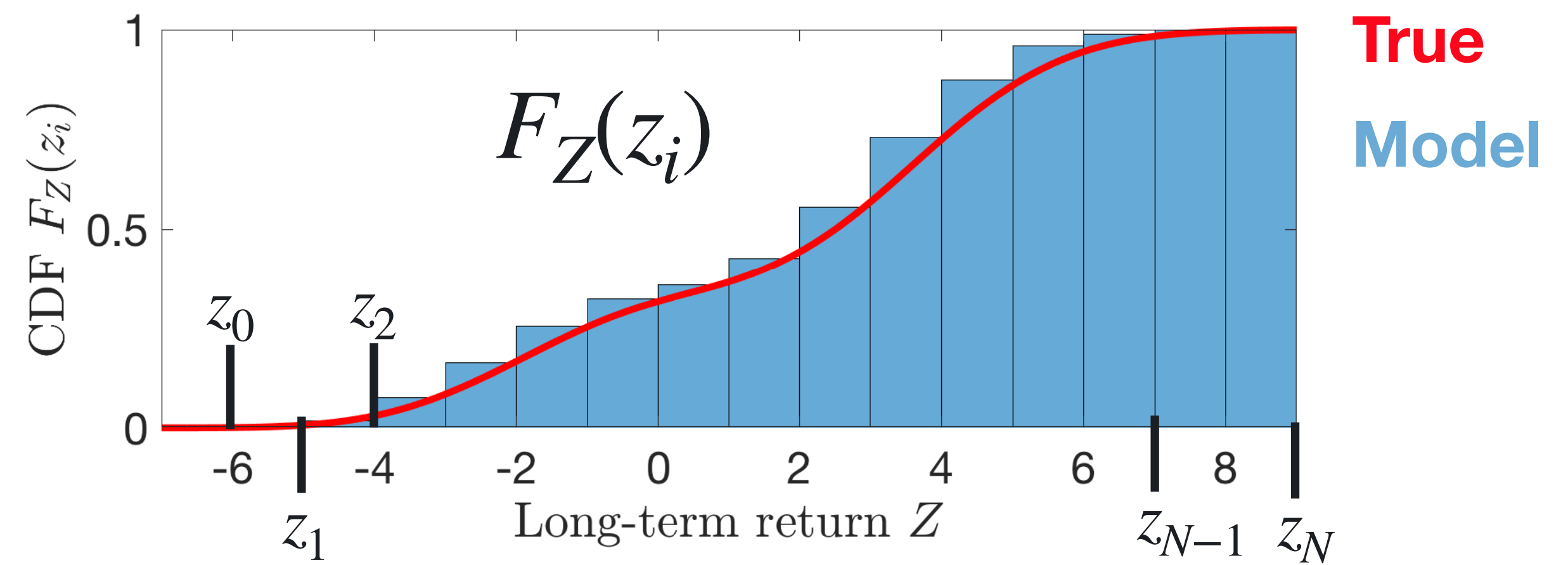


1. Categorical

- Fixed **support** bins $z_0, z_1, z_2, \dots, z_N$
- Learn probabilities/quantiles p_i, τ_i

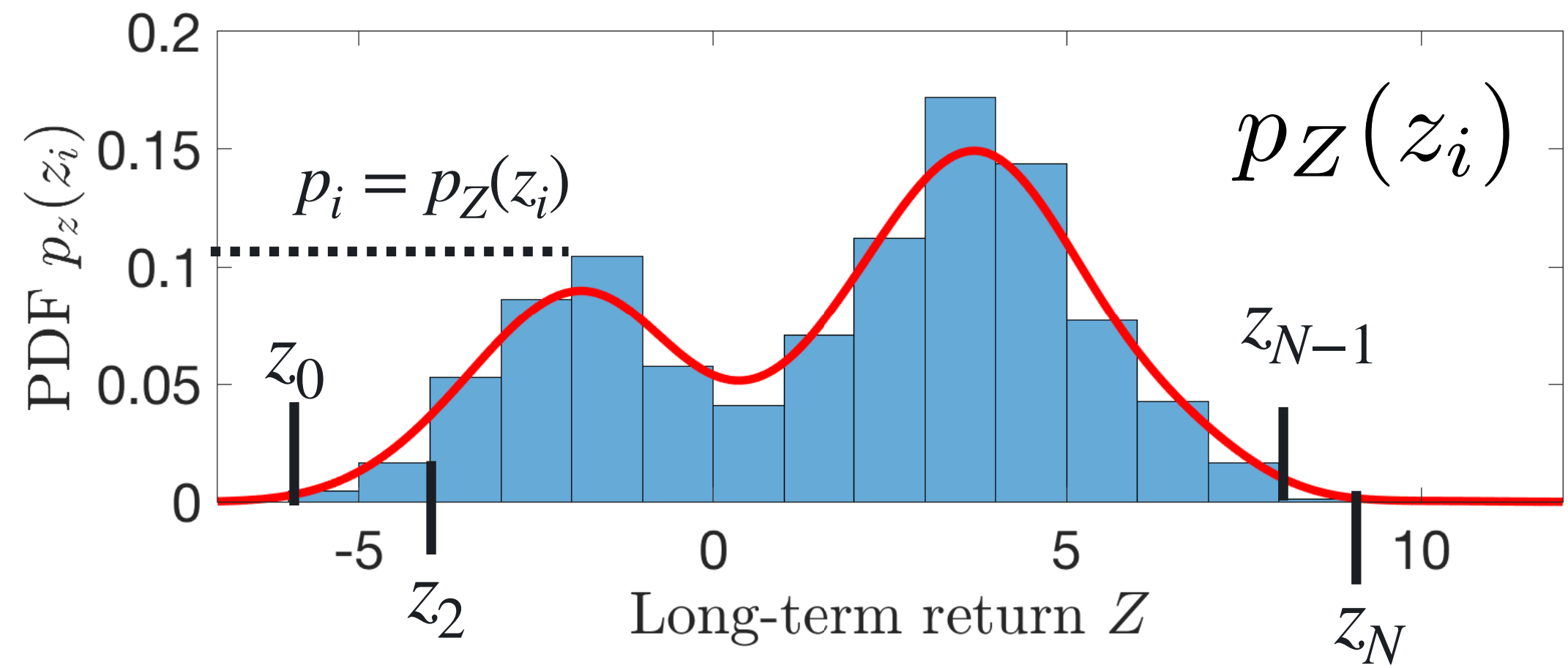
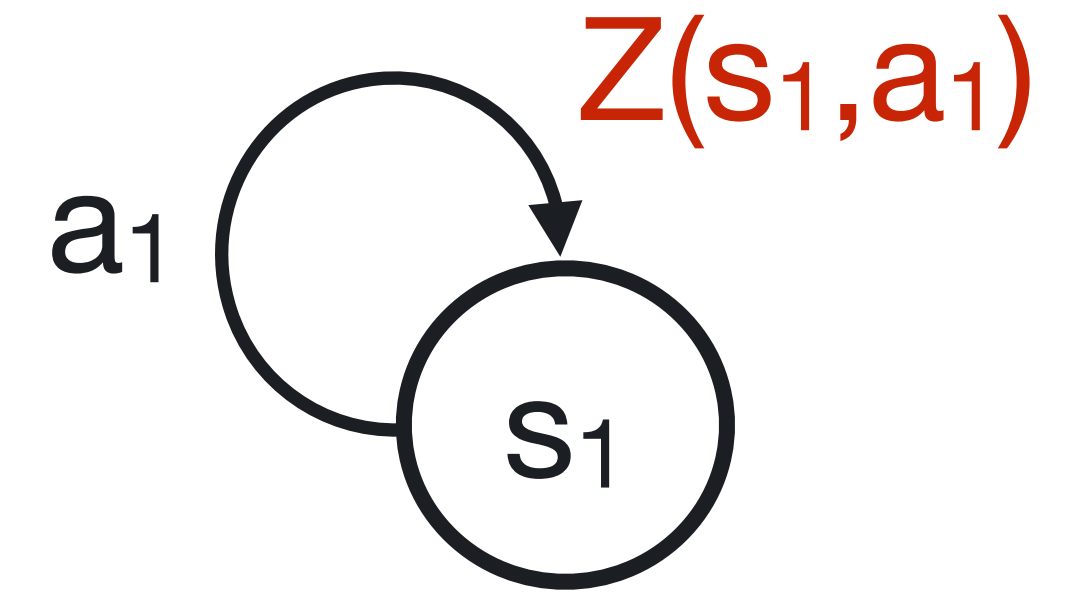


OR

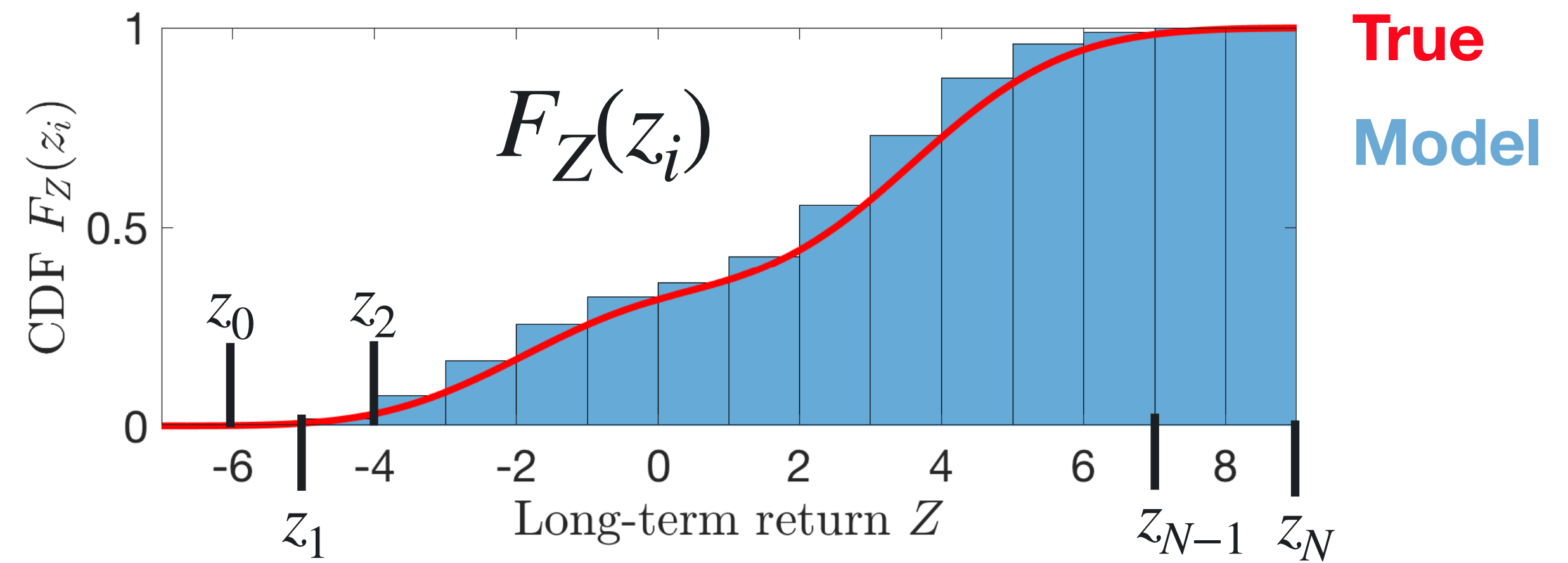


1. Categorical

- Fixed **support** bins $z_0, z_1, z_2, \dots, z_N$
- Learn probabilities/quantiles p_i, τ_i
- Easy to program

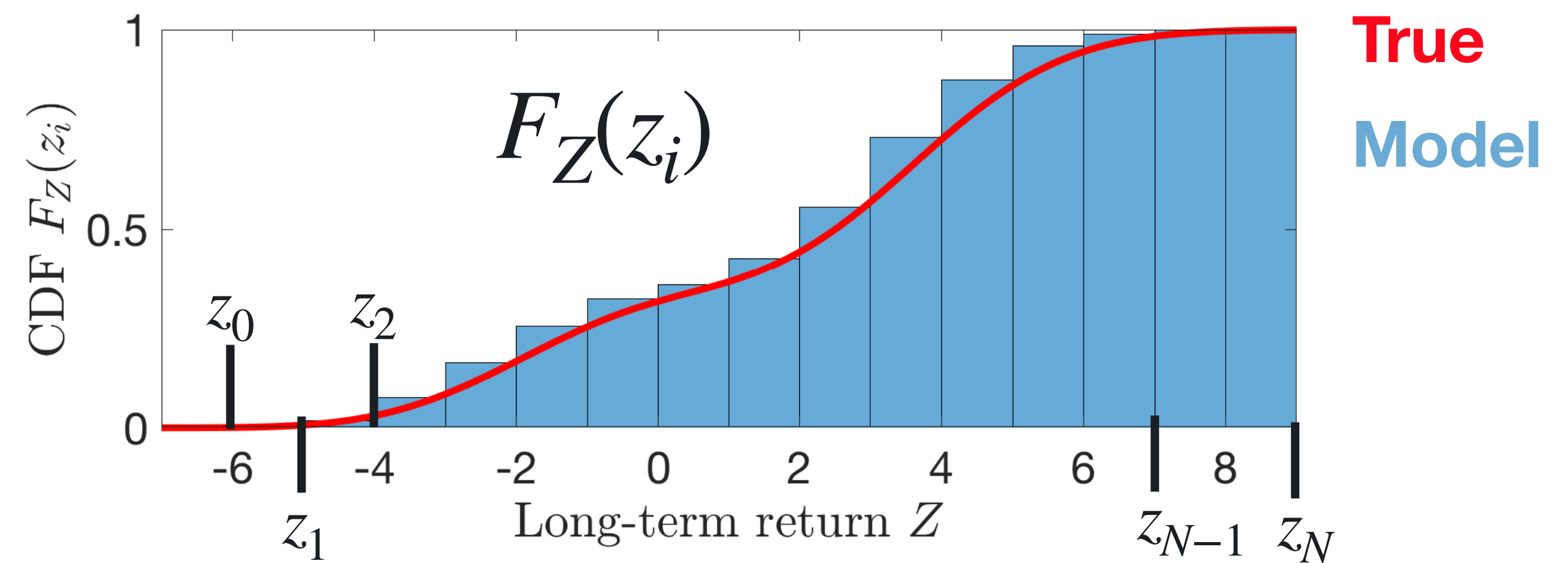
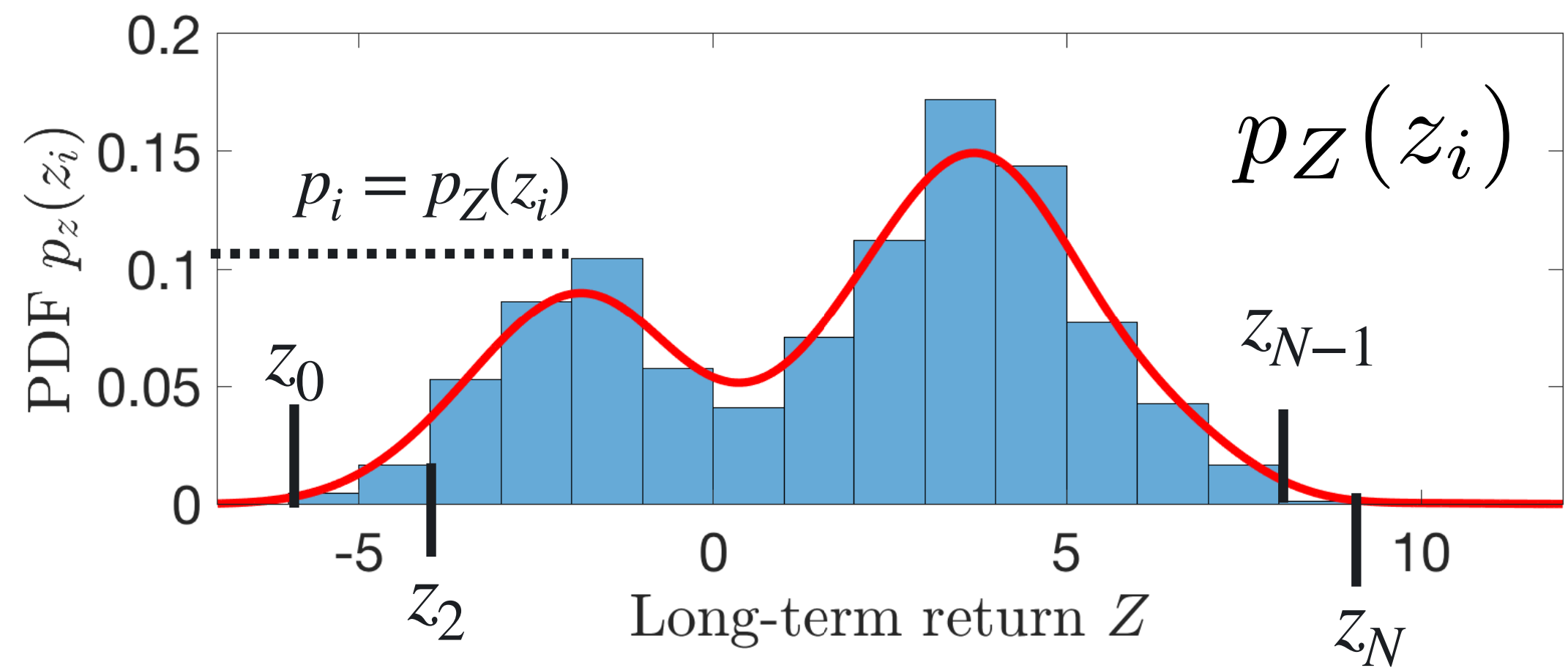
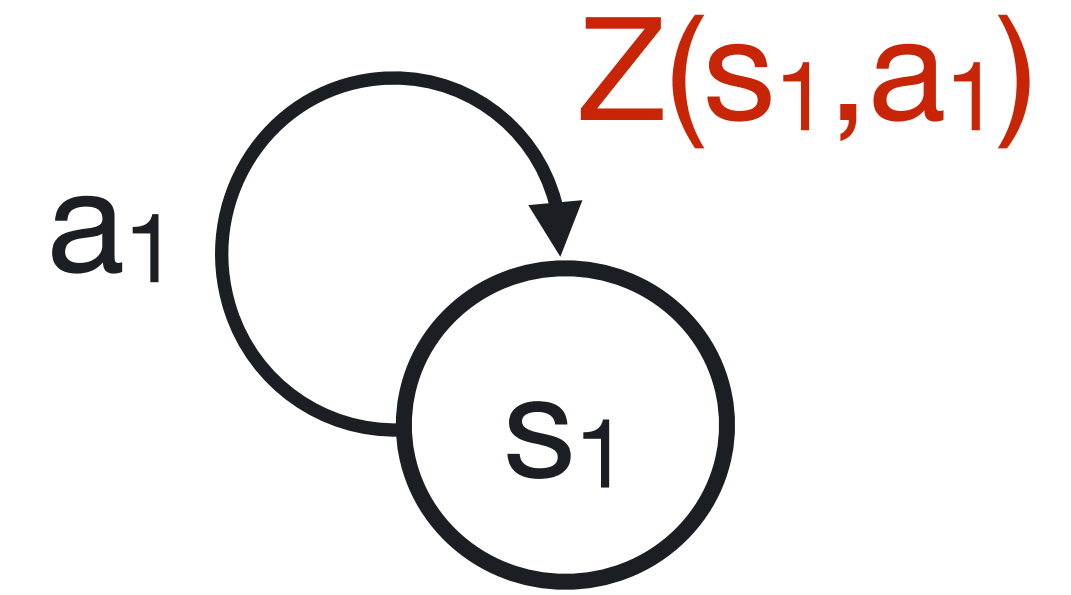


OR



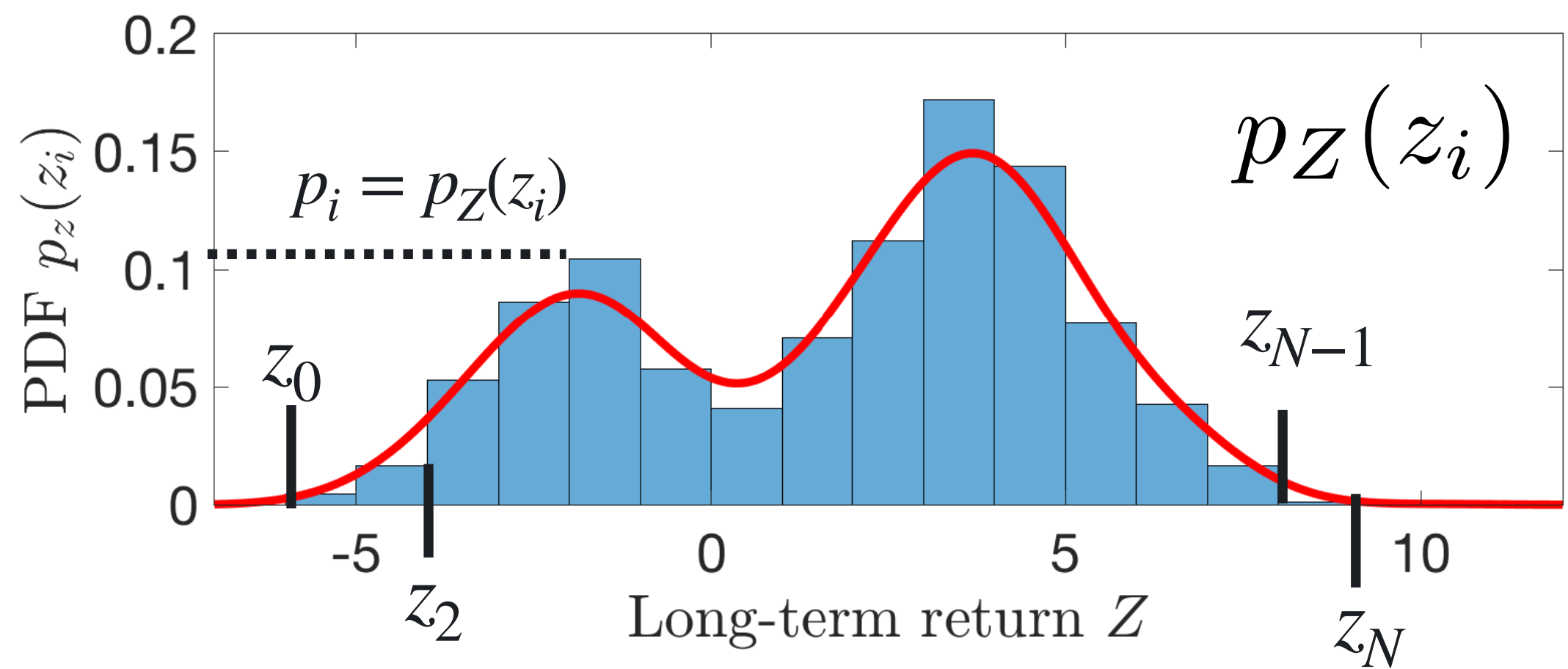
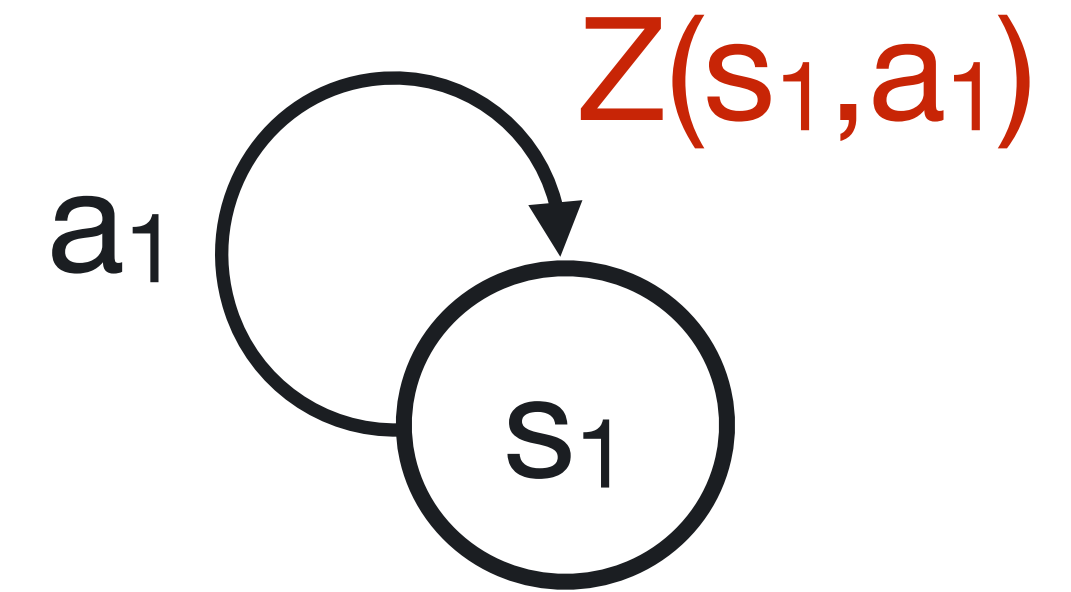
1. Categorical

- Fixed **support** bins $z_0, z_1, z_2, \dots, z_N$
- Learn probabilities/quantiles p_i, τ_i
- Easy to program
- **Value range of Z needs to be known!**

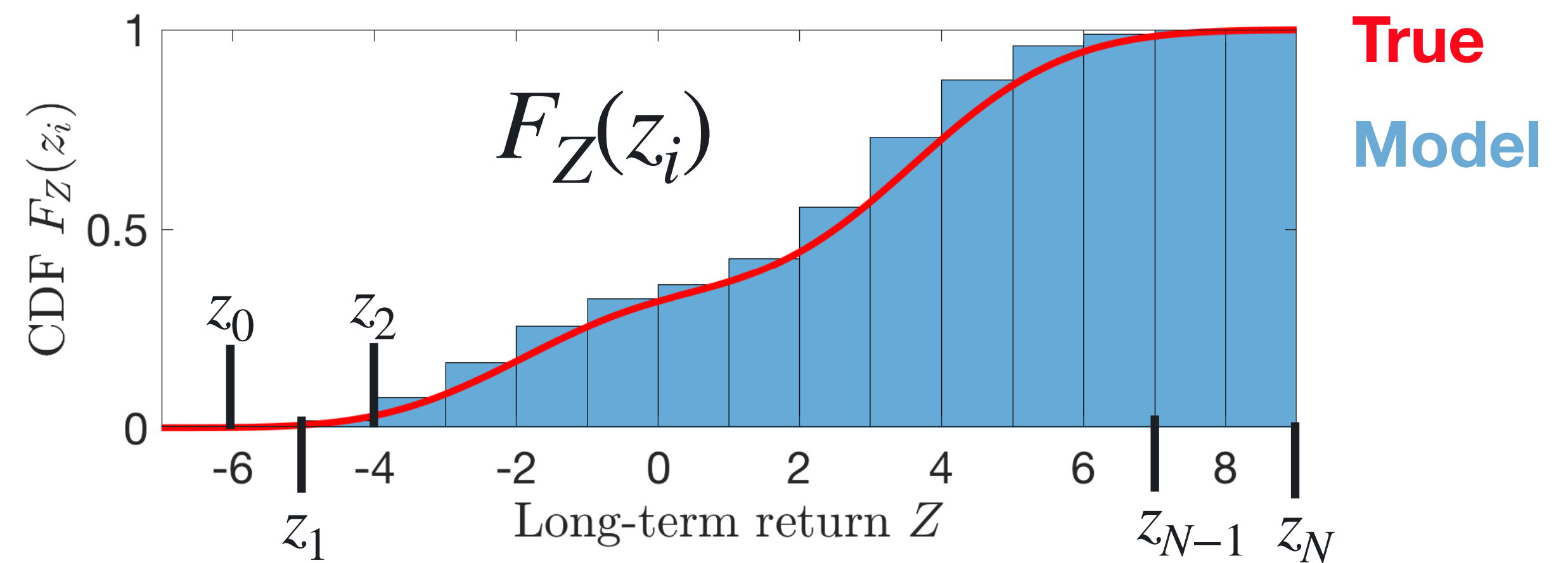


1. Categorical

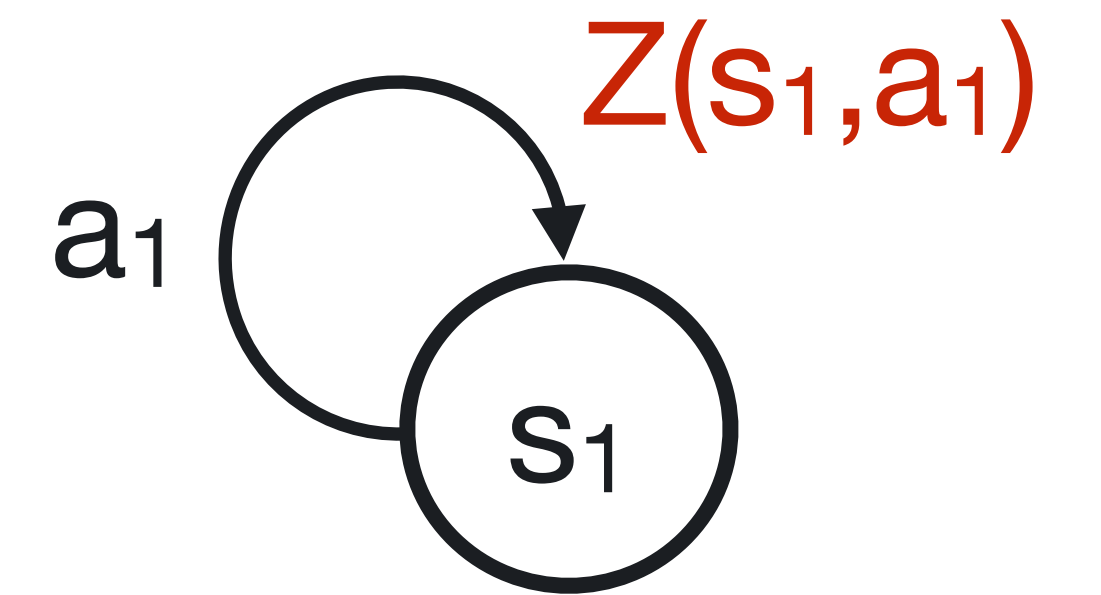
- Fixed **support** bins $z_0, z_1, z_2, \dots, z_N$
- Learn probabilities/quantiles p_i, τ_i
- Easy to program
- Value range of Z needs to be known!
- The bins are fixed (less expressive)



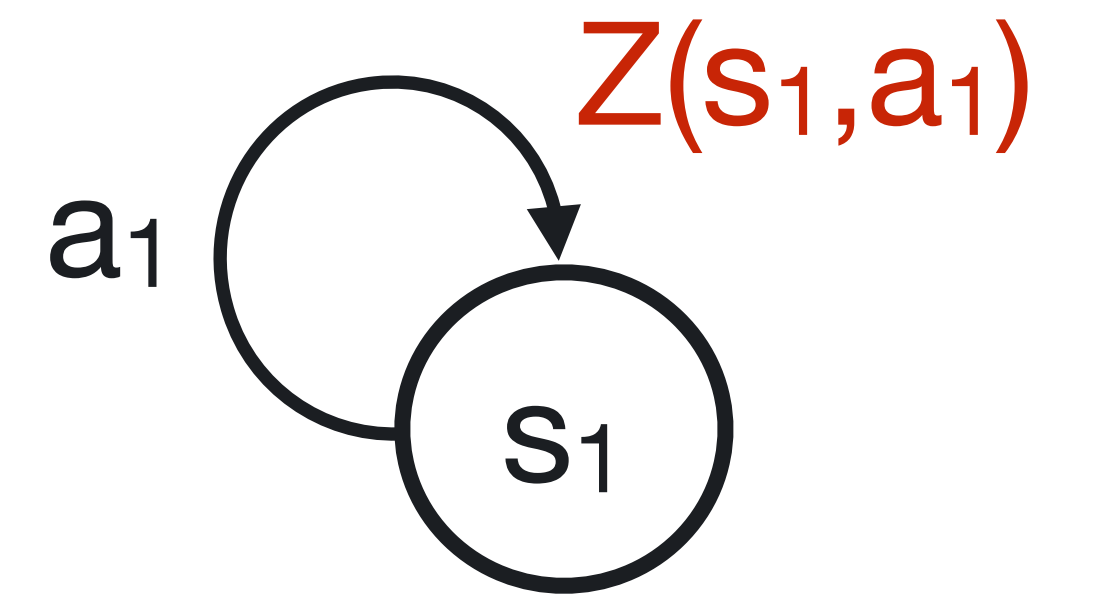
OR



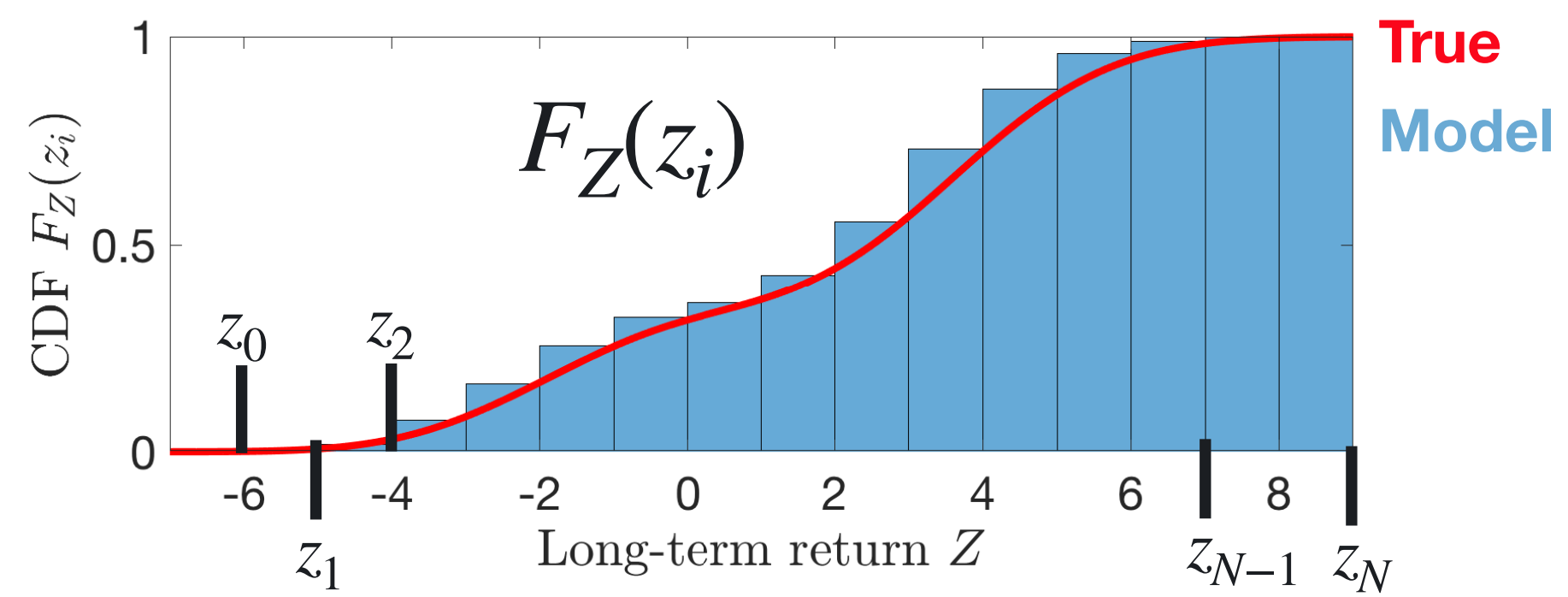
2. Quantiles of Inverse CDF



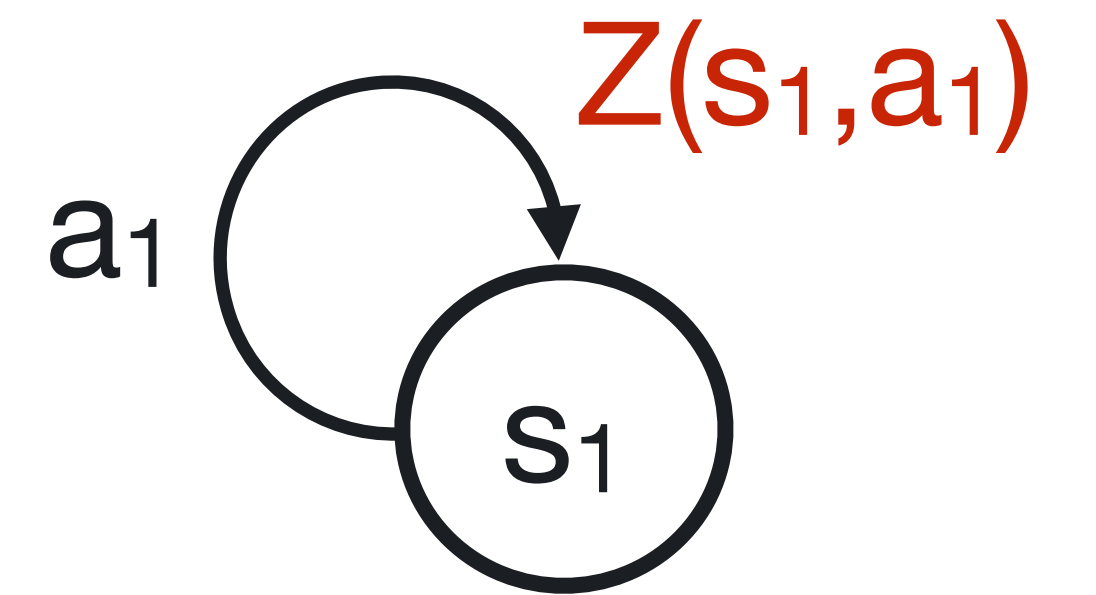
2. Quantiles of Inverse CDF



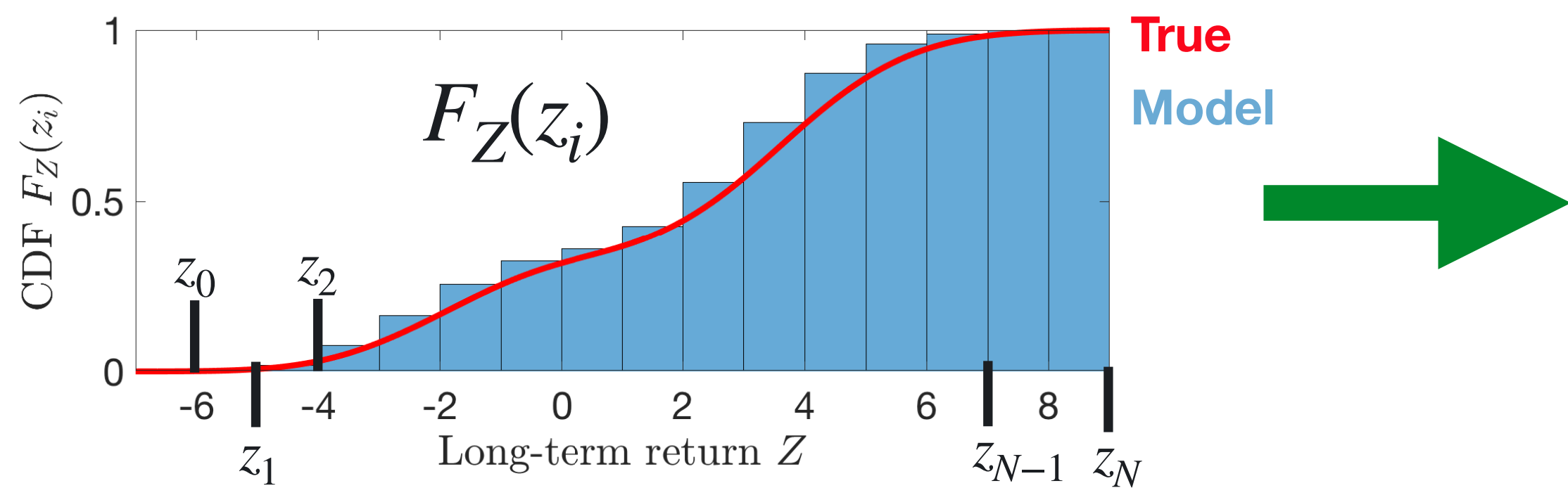
Categorical



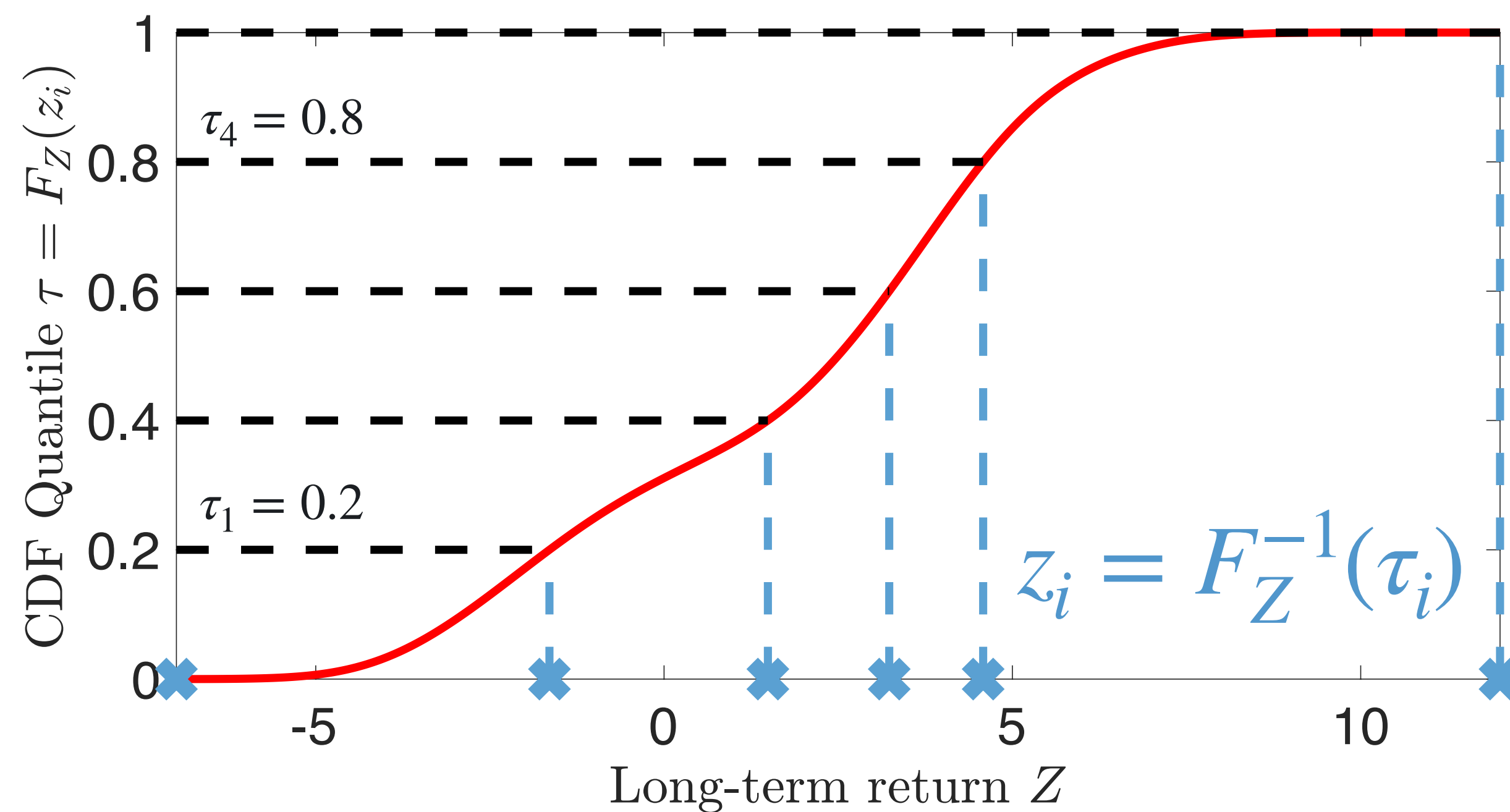
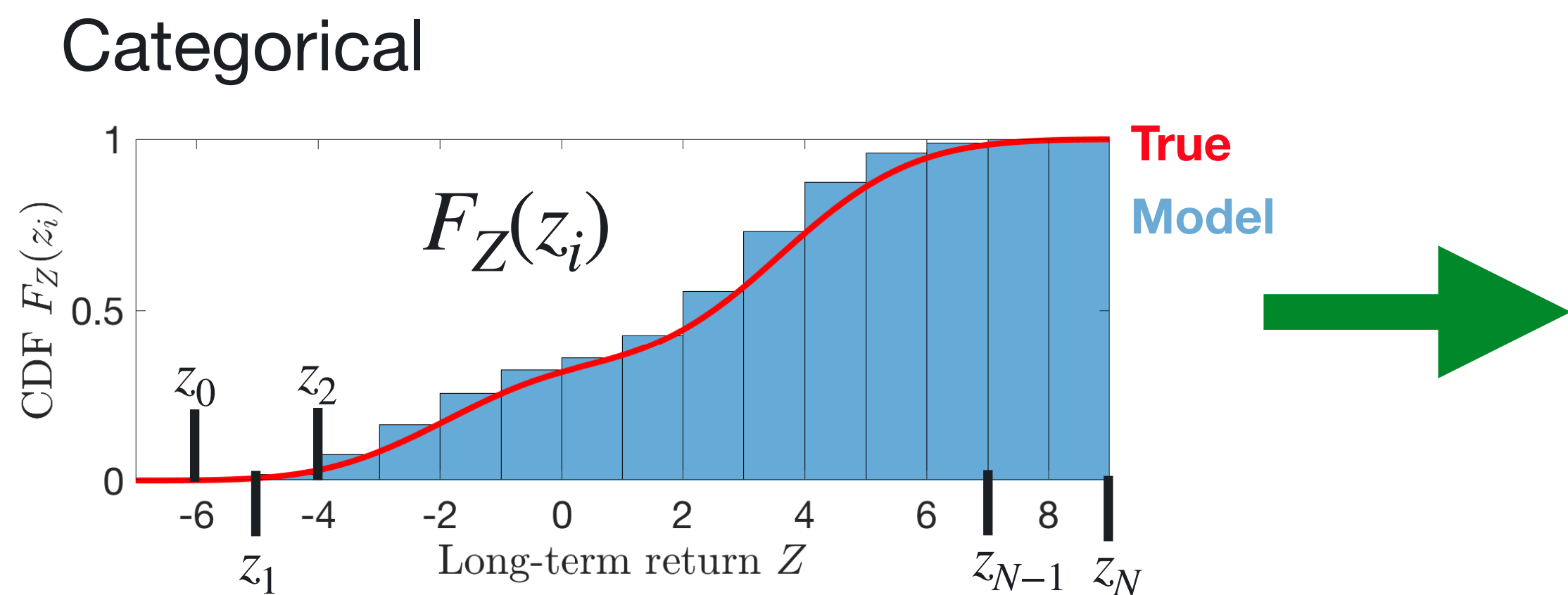
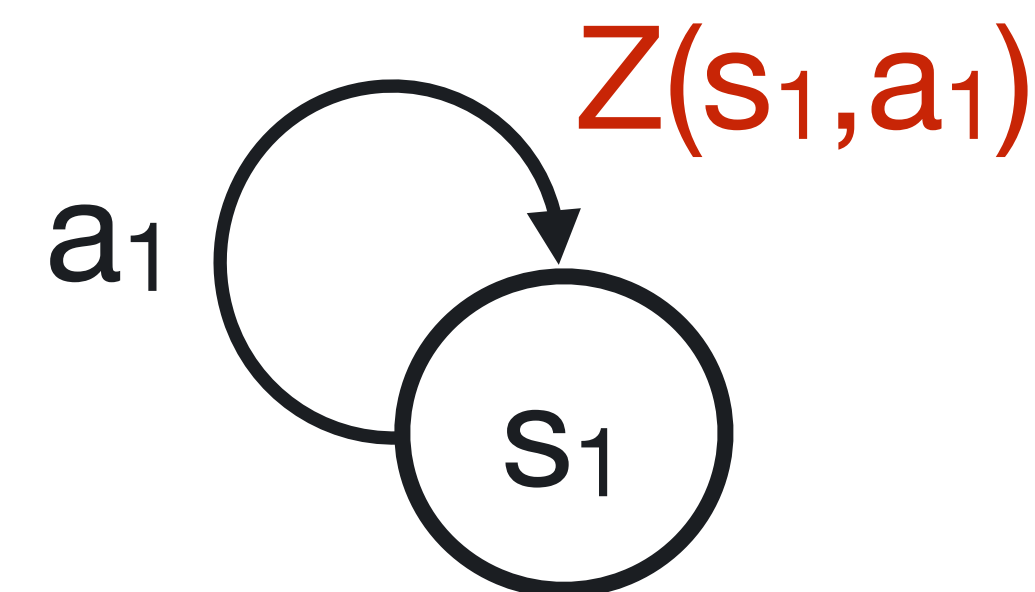
2. Quantiles of Inverse CDF



Categorical

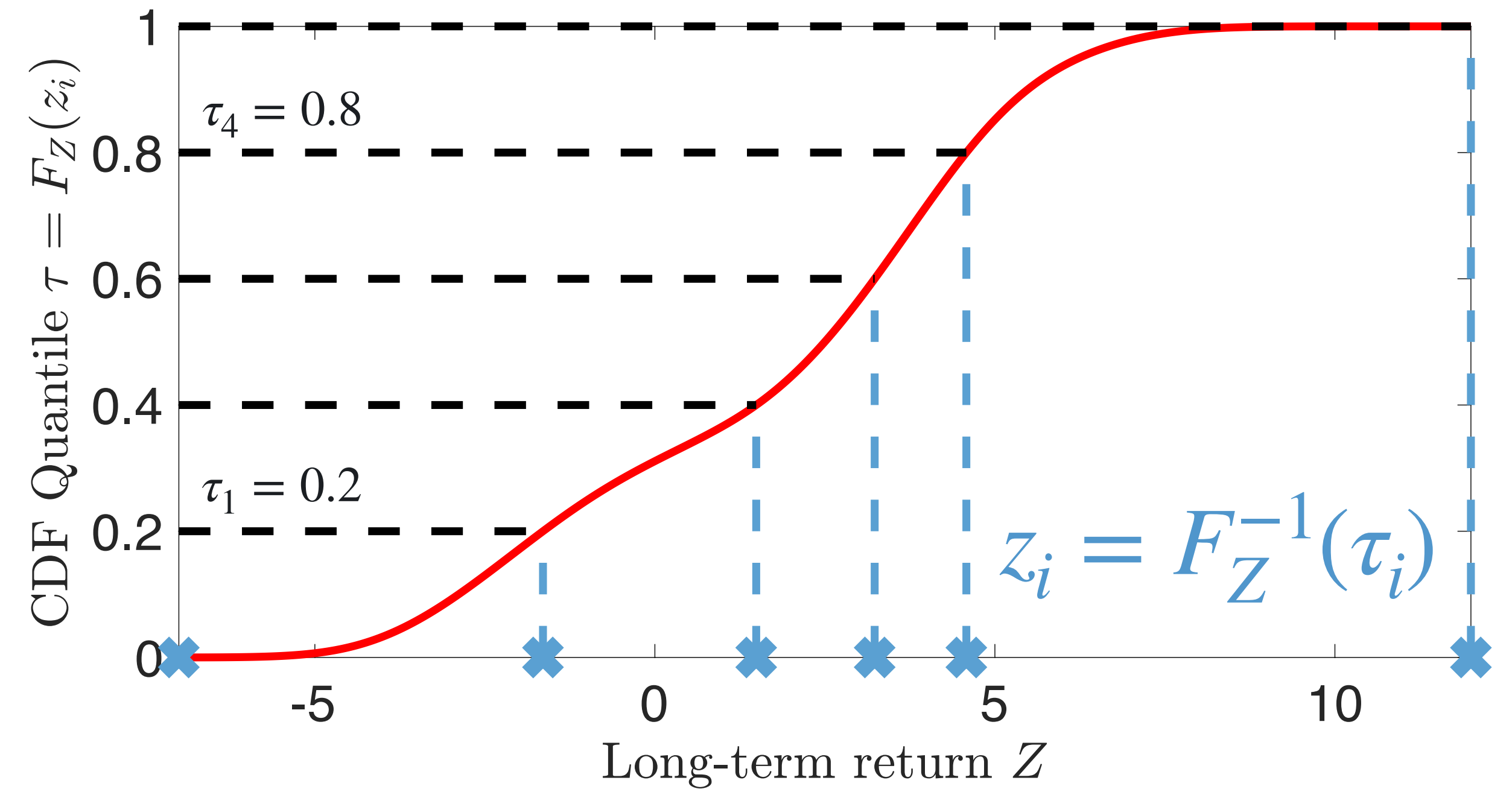
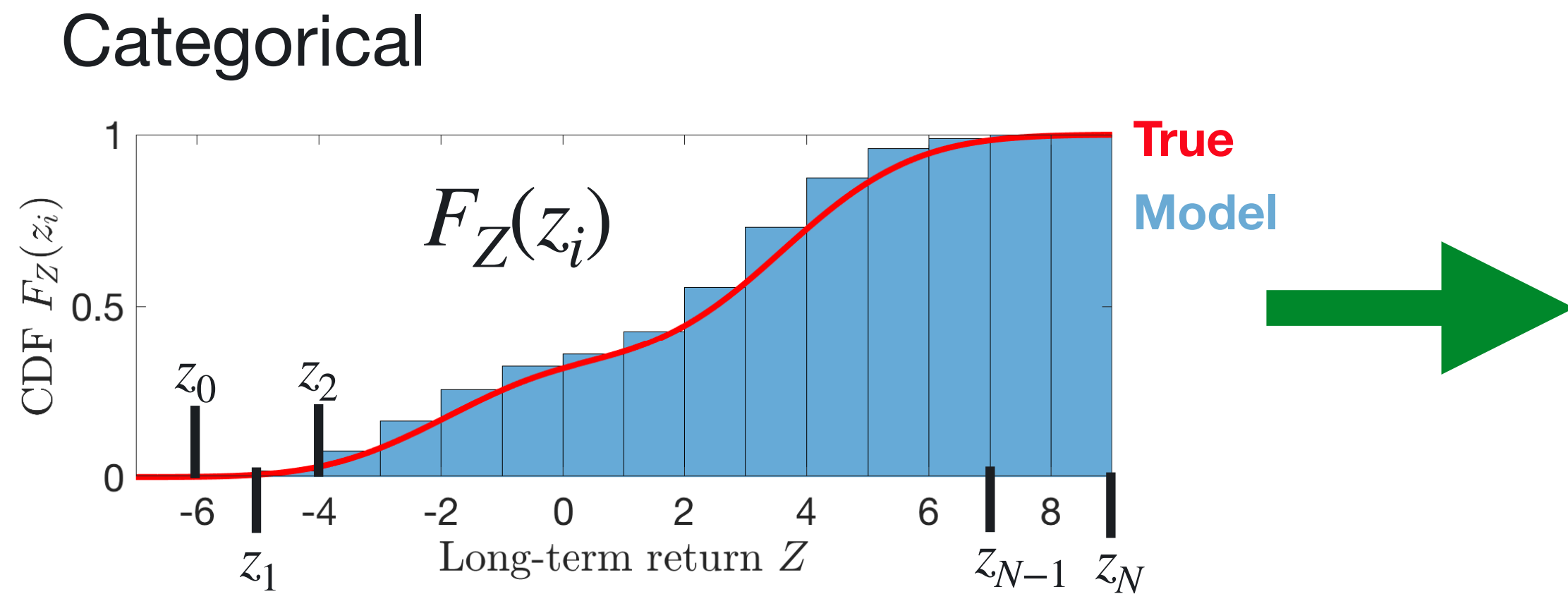
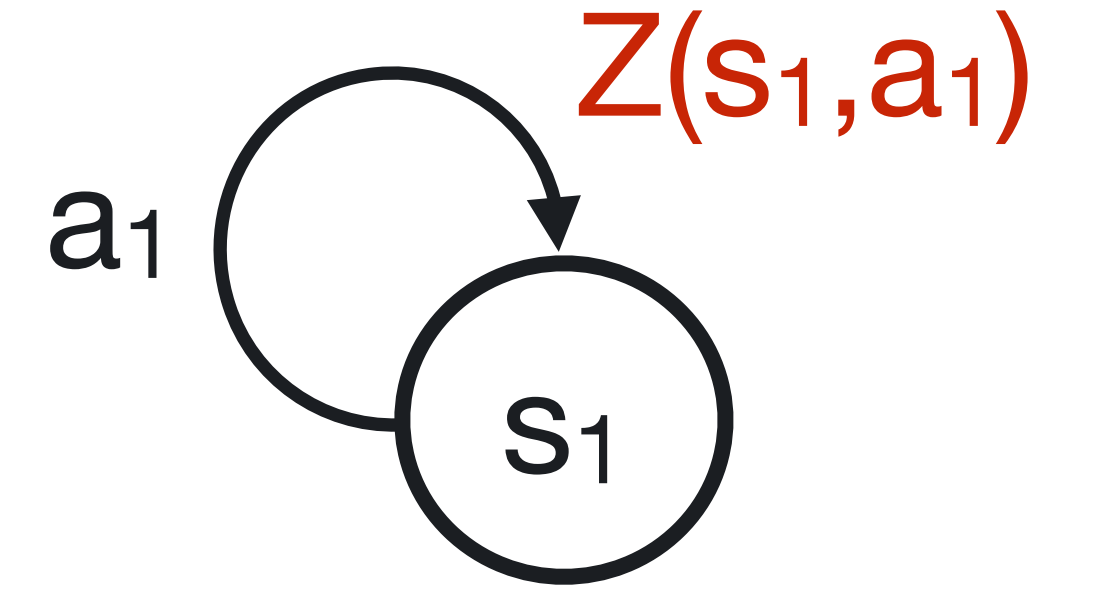


2. Quantiles of Inverse CDF



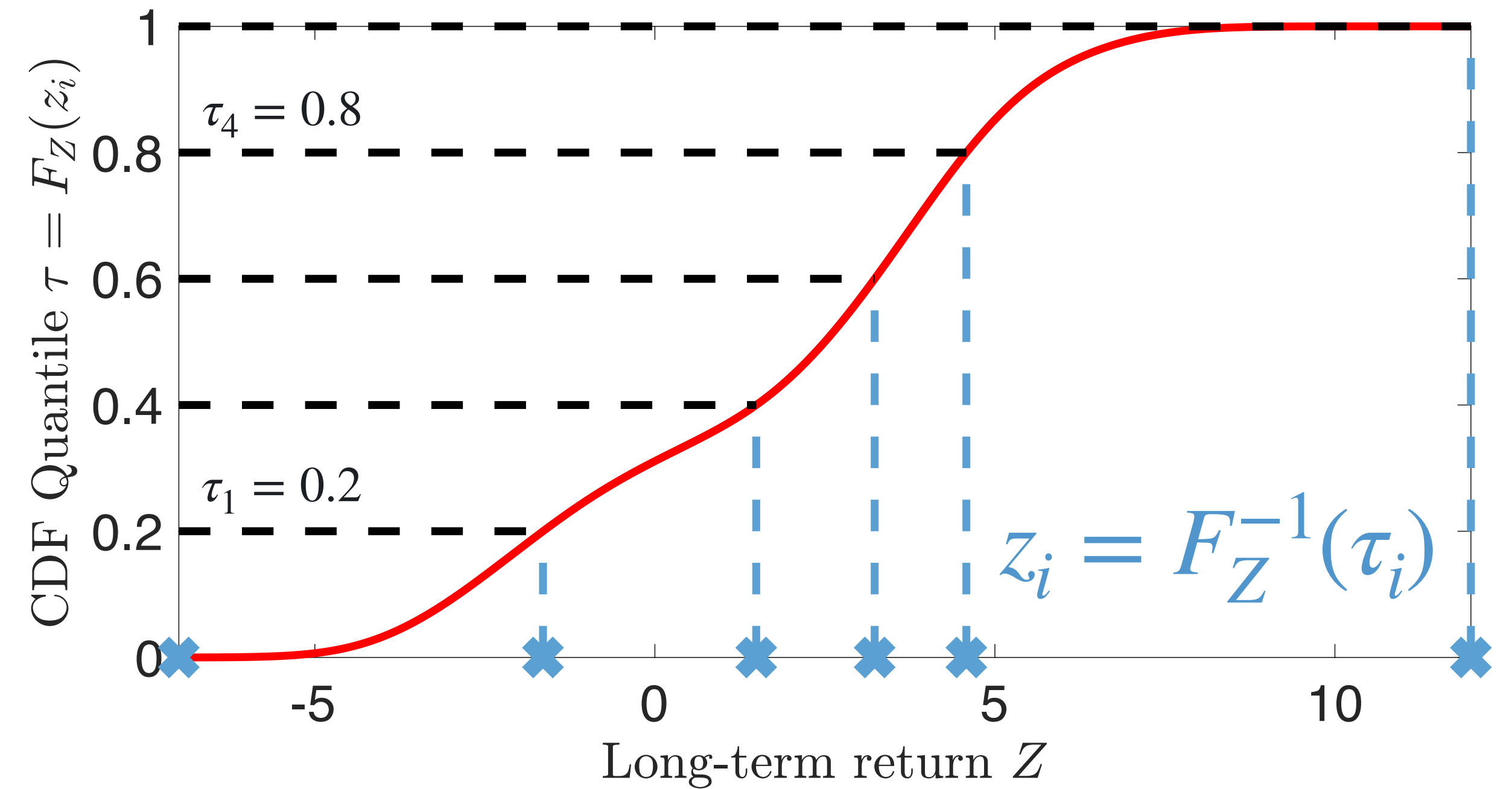
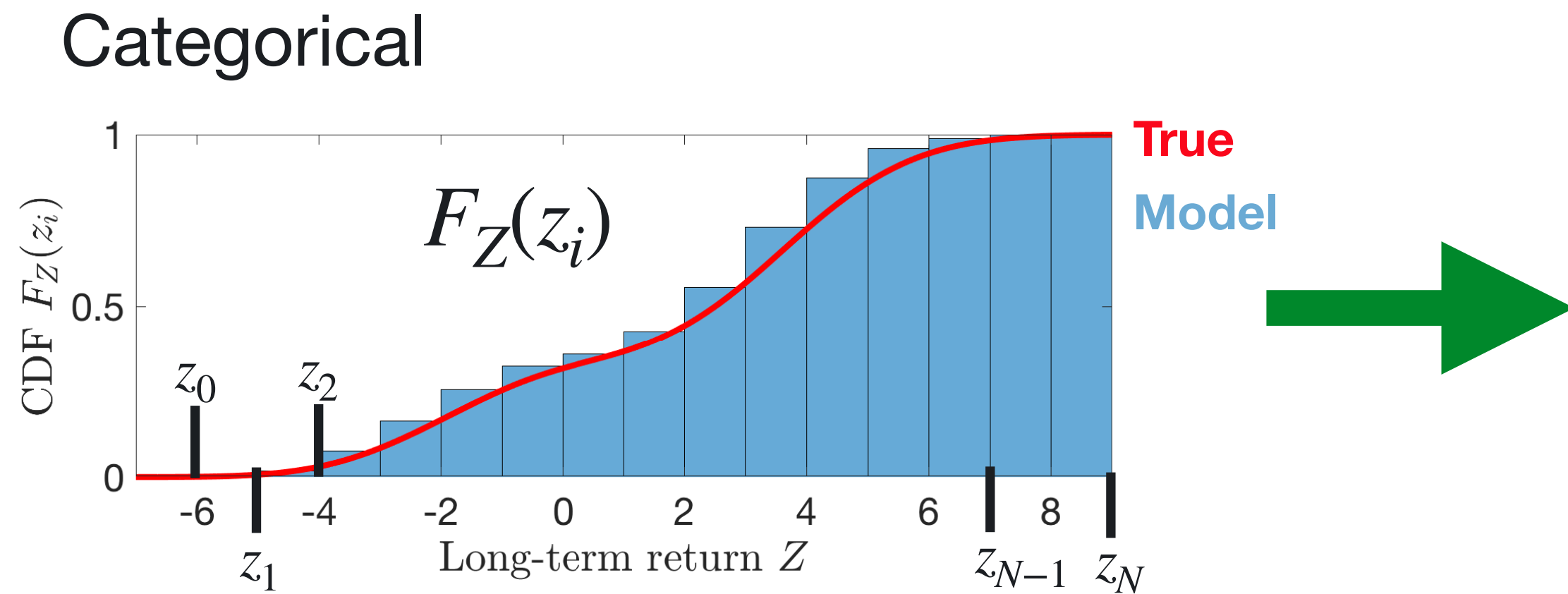
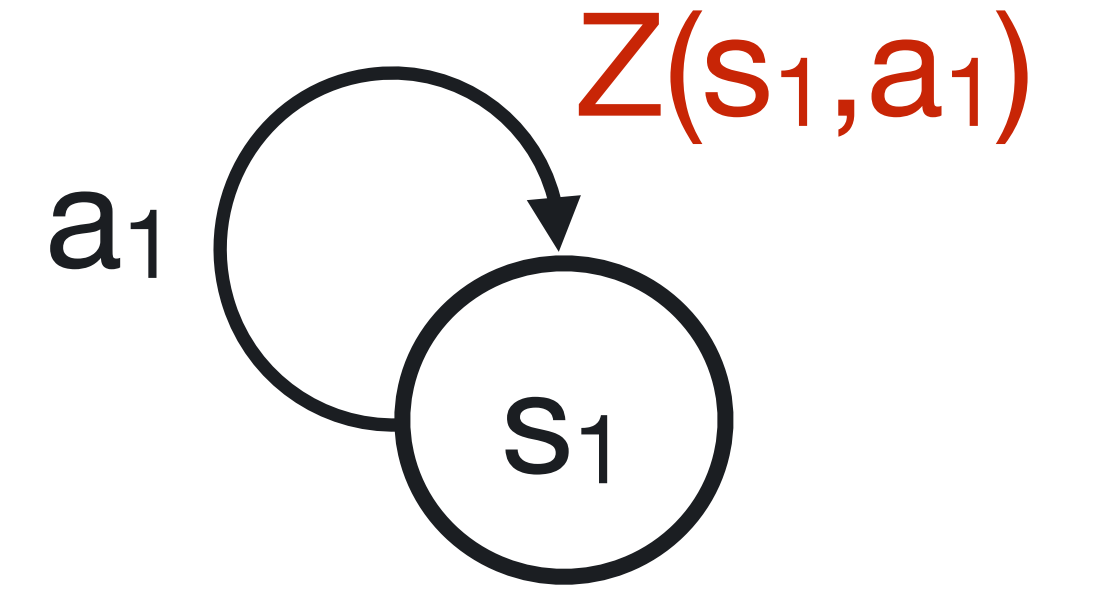
2. Quantiles of Inverse CDF

- Fixed **quantile** bins $\tau_0, \tau_1, \tau_2, \dots, \tau_N$



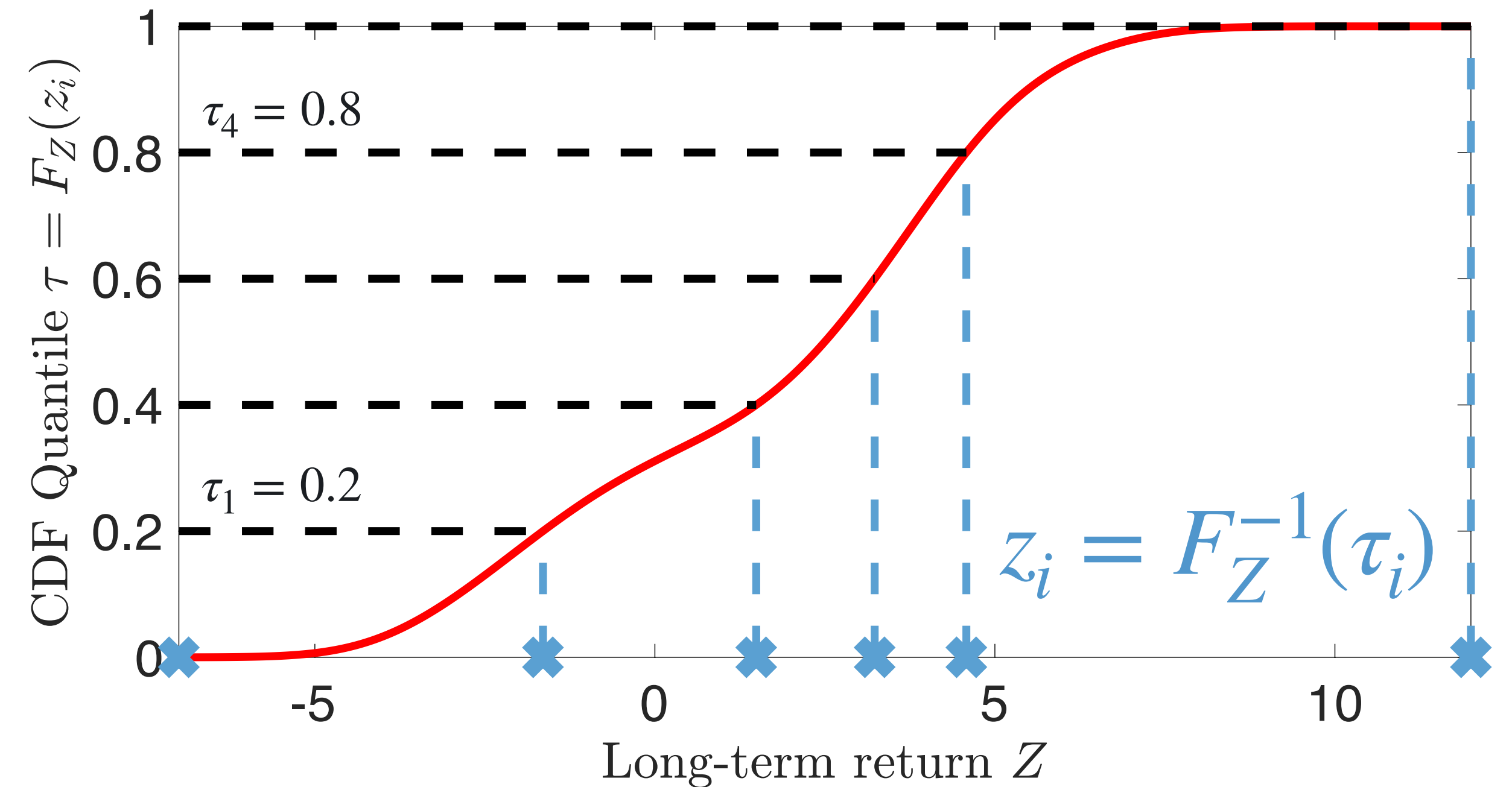
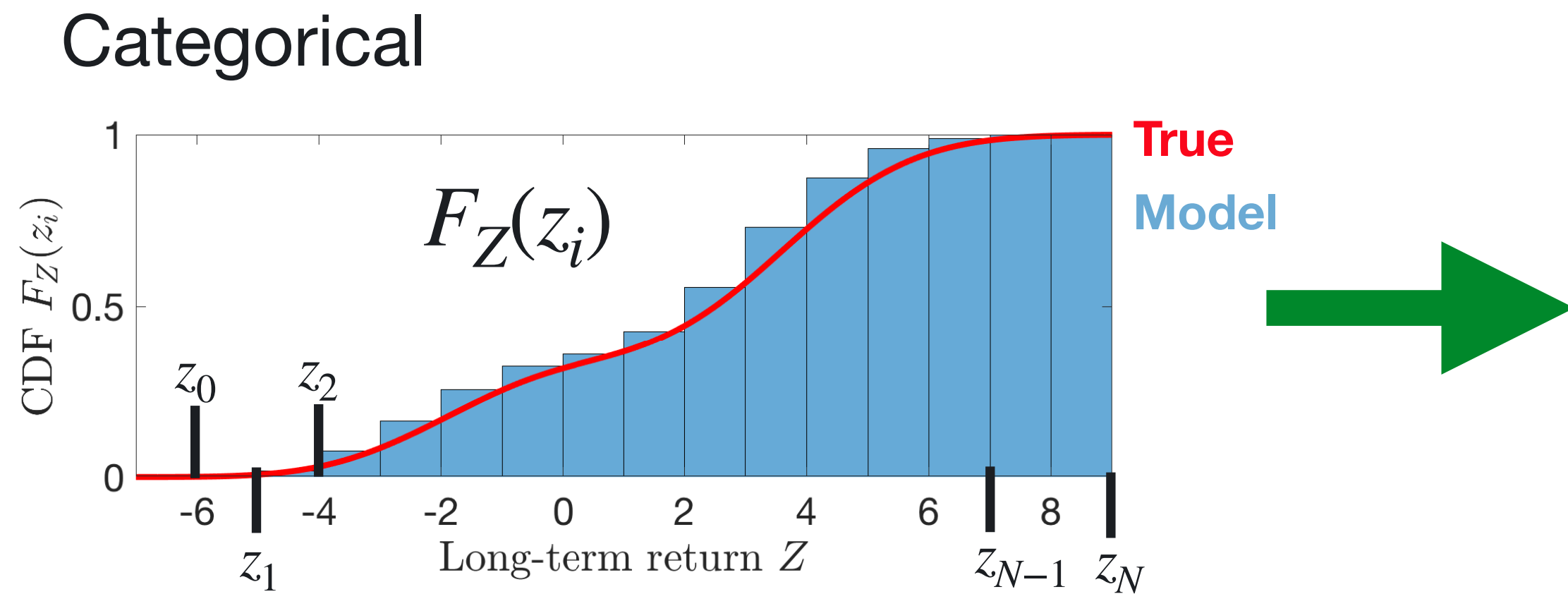
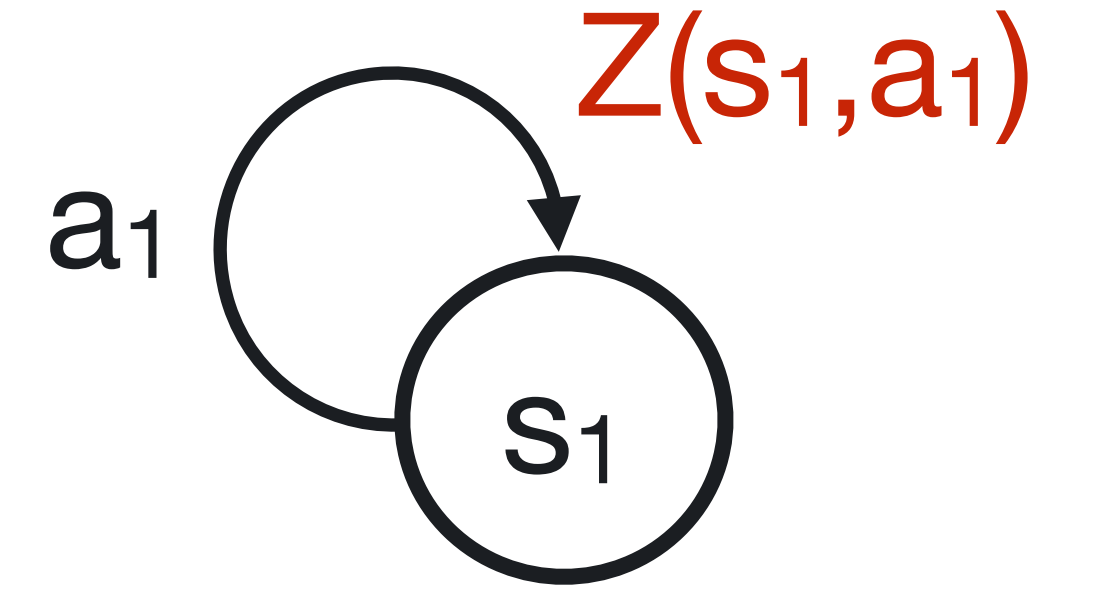
2. Quantiles of Inverse CDF

- Fixed **quantile** bins $\tau_0, \tau_1, \tau_2, \dots, \tau_N$
- Learn support $z_i = F_Z^{-1}(\tau_i)$



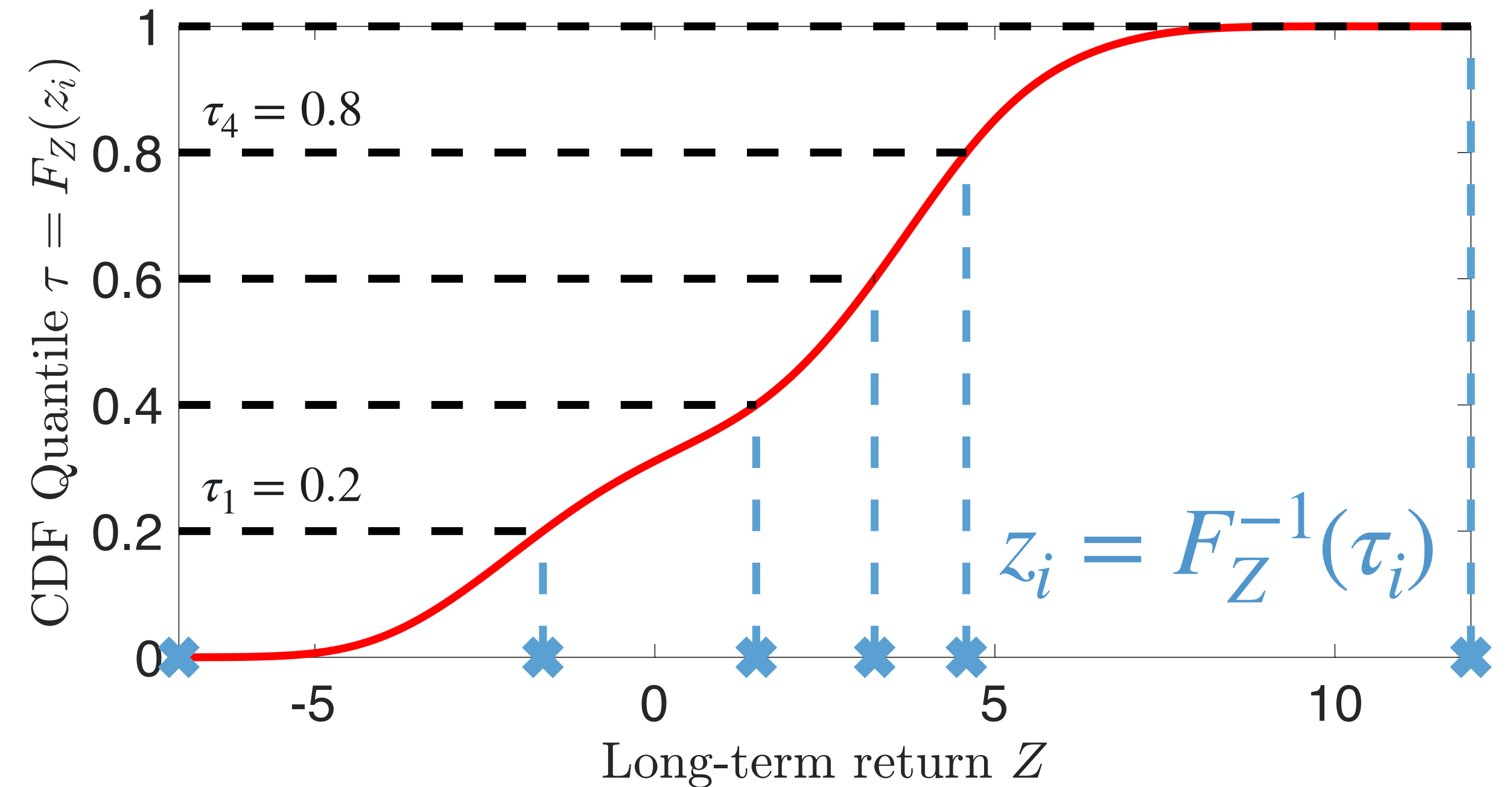
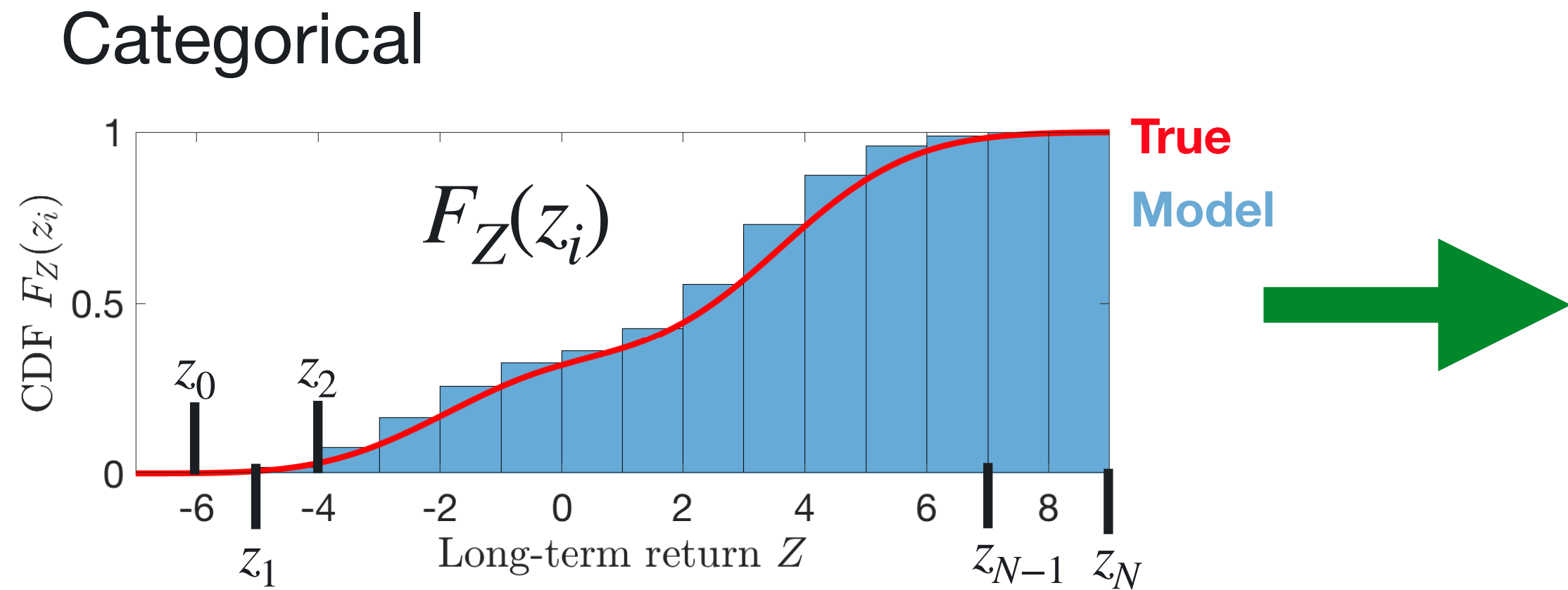
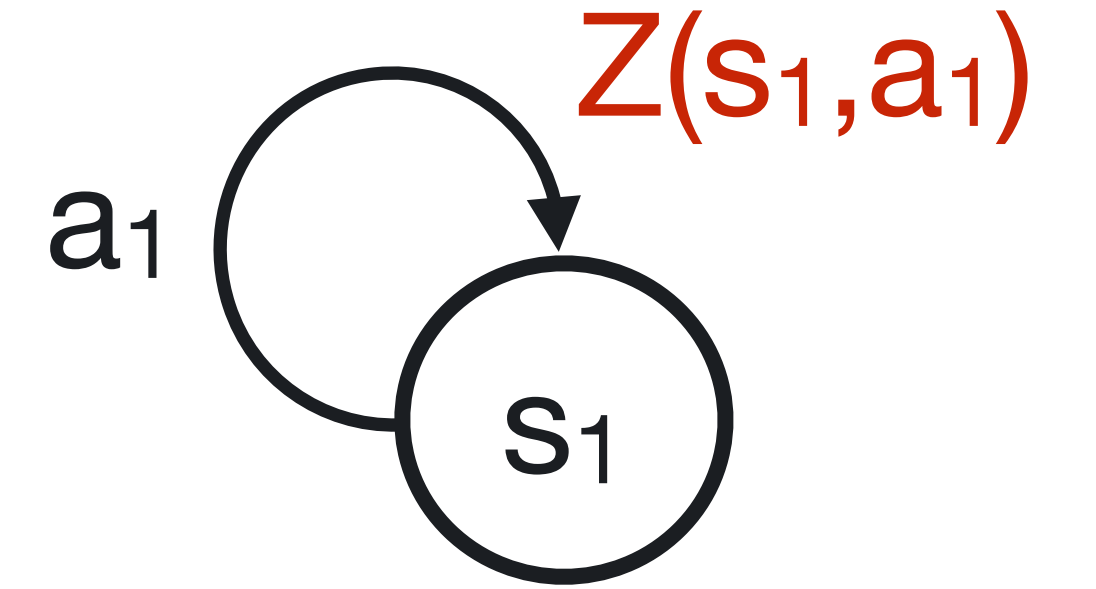
2. Quantiles of Inverse CDF

- Fixed **quantile** bins $\tau_0, \tau_1, \tau_2, \dots, \tau_N$
- Learn support $z_i = F_Z^{-1}(\tau_i)$
- No need for value range of Z



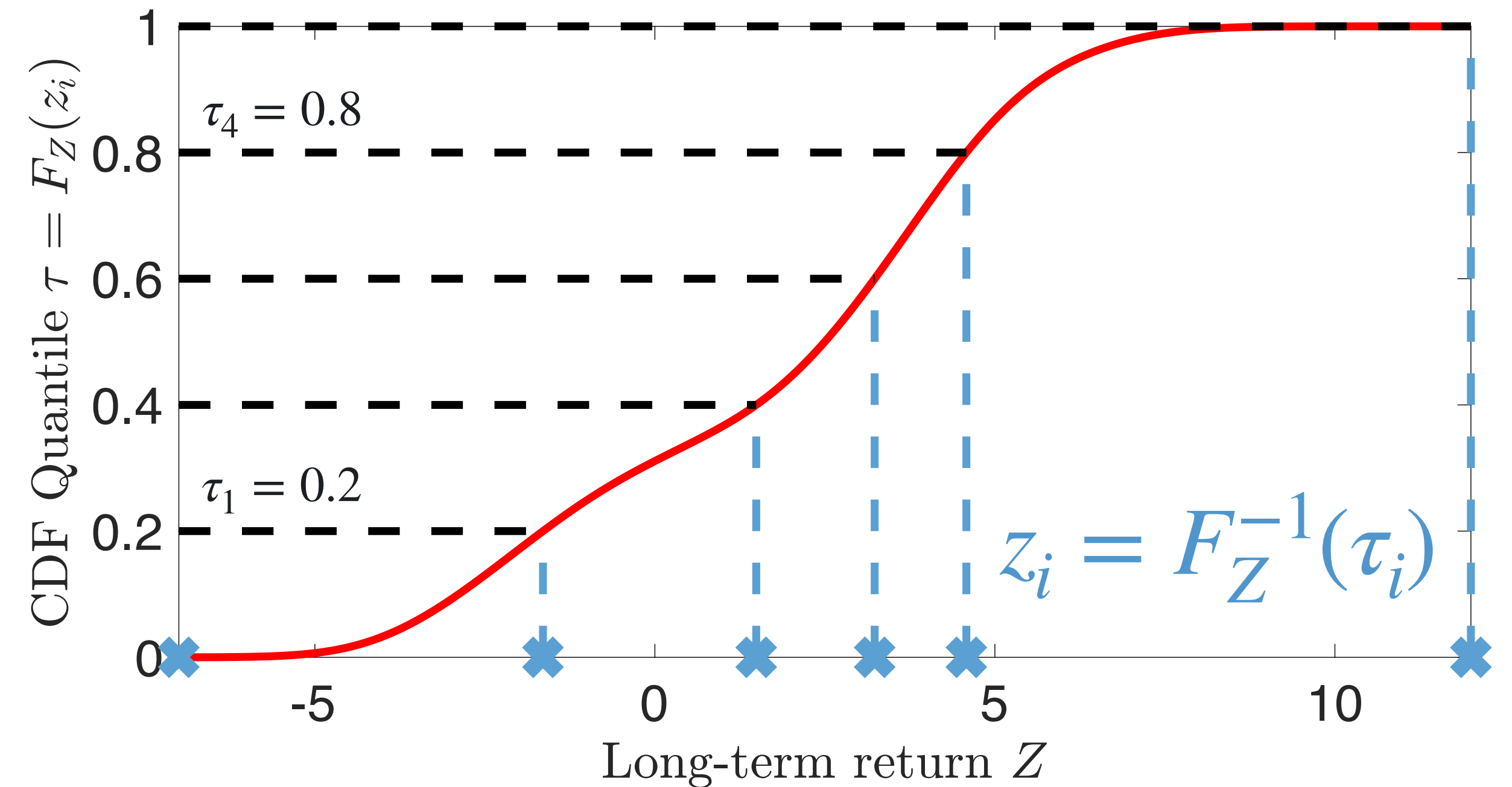
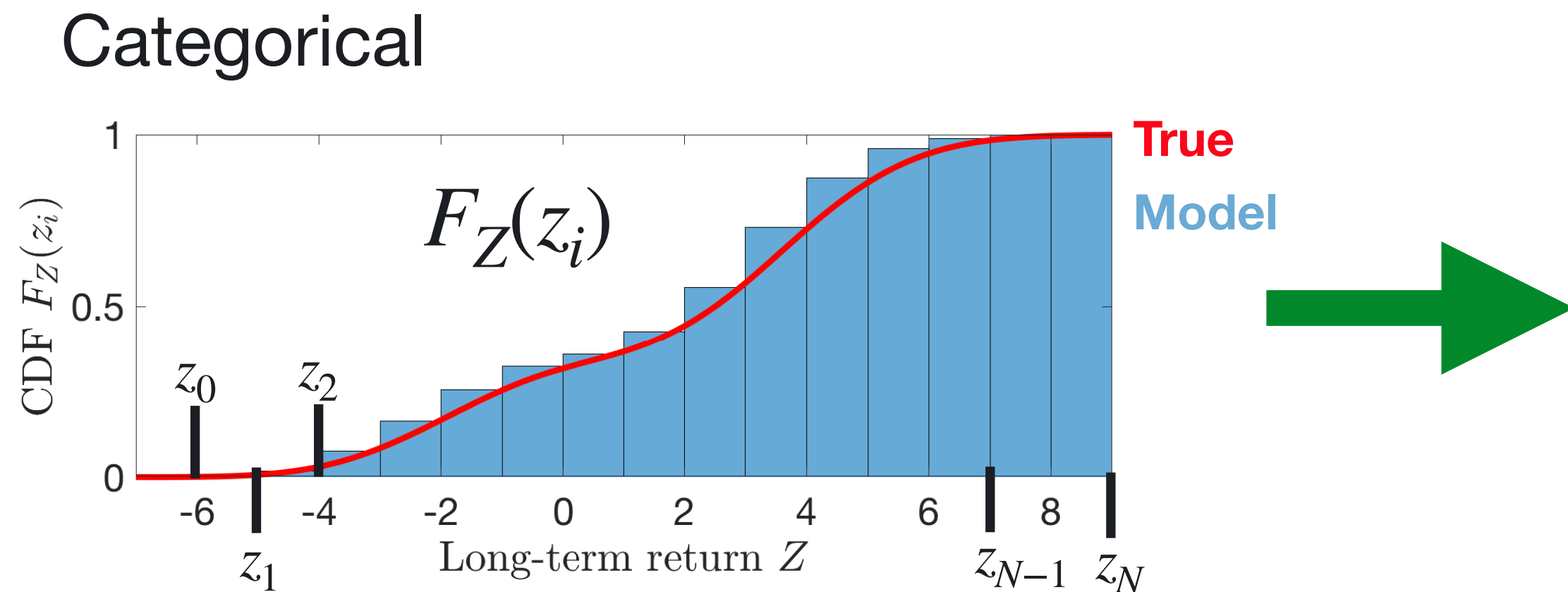
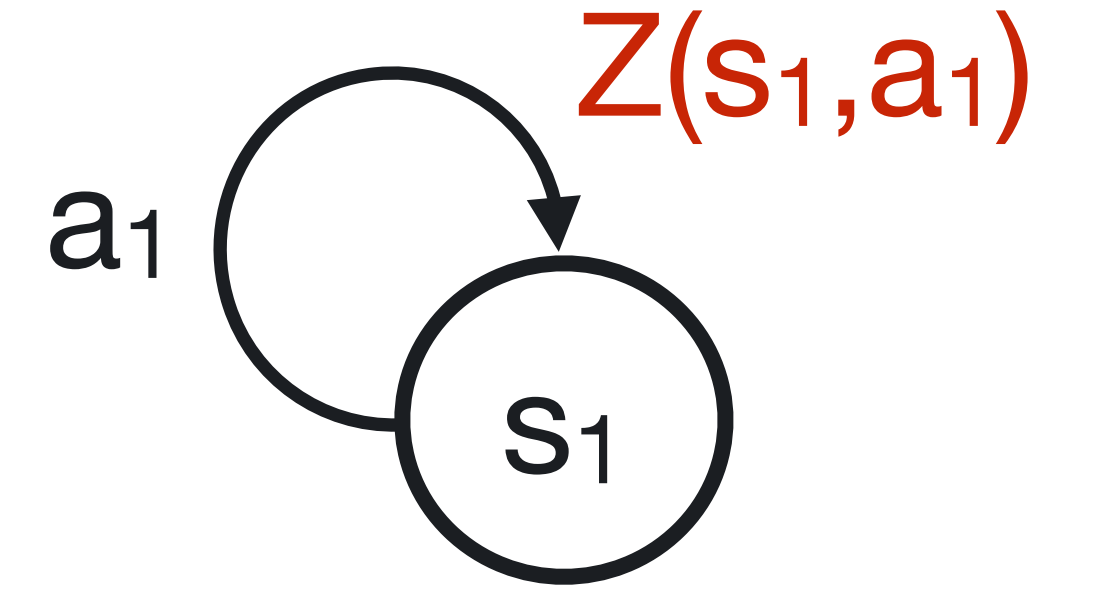
2. Quantiles of Inverse CDF

- Fixed **quantile** bins $\tau_0, \tau_1, \tau_2, \dots, \tau_N$
- Learn support $z_i = F_Z^{-1}(\tau_i)$
- No need for value range of Z
- **Training loss**

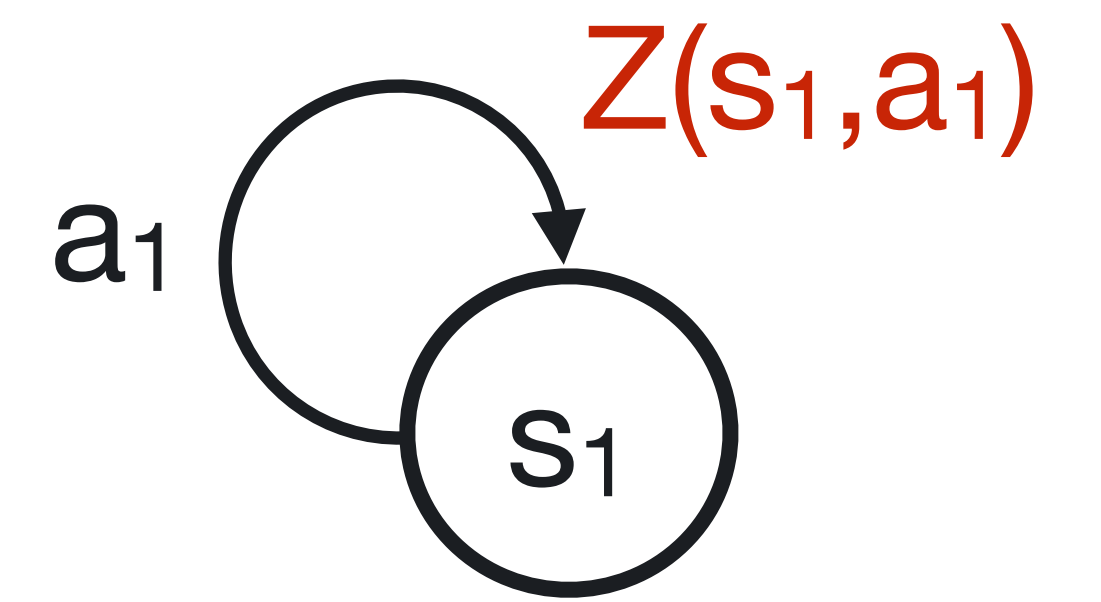


2. Quantiles of Inverse CDF

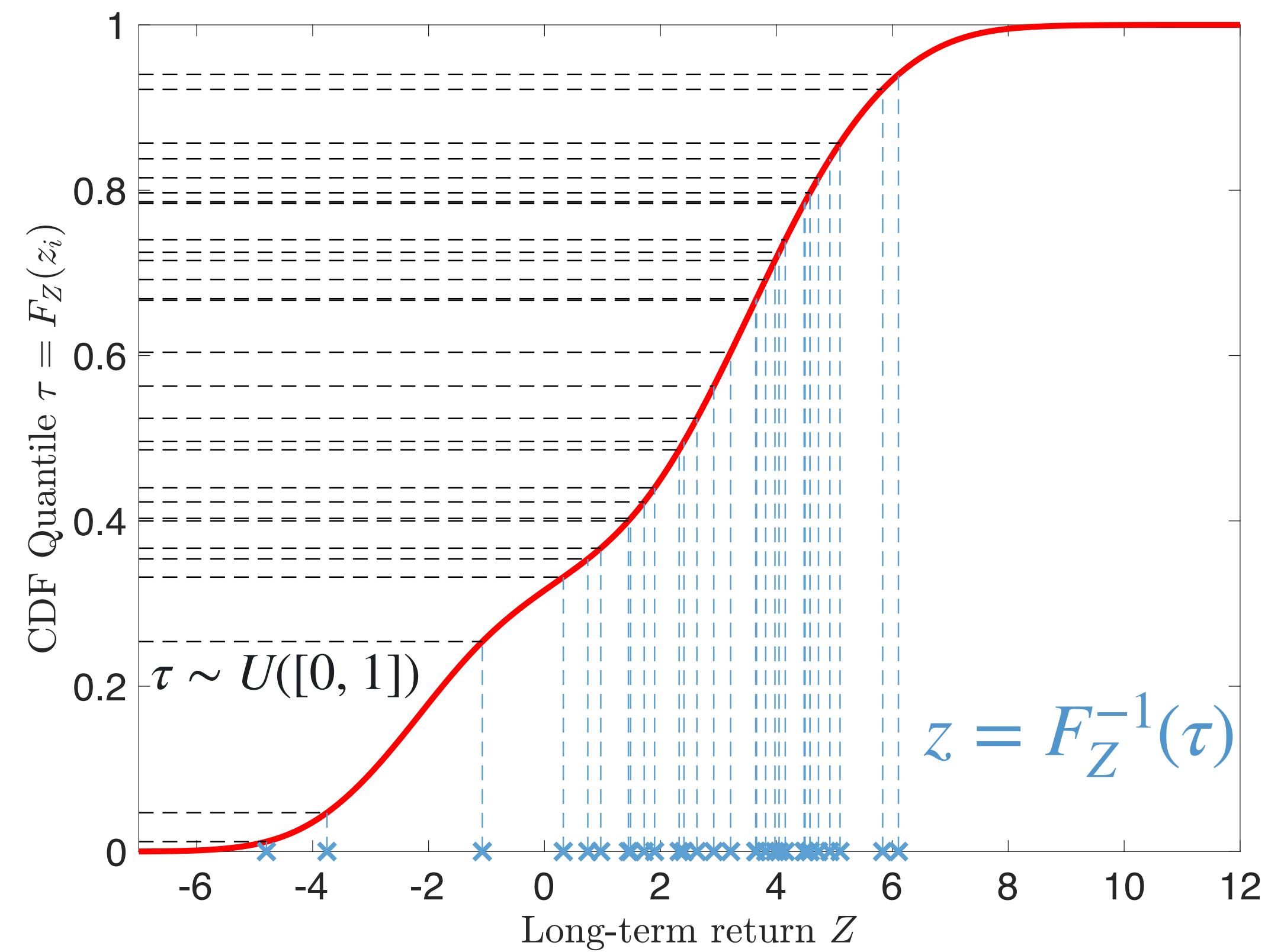
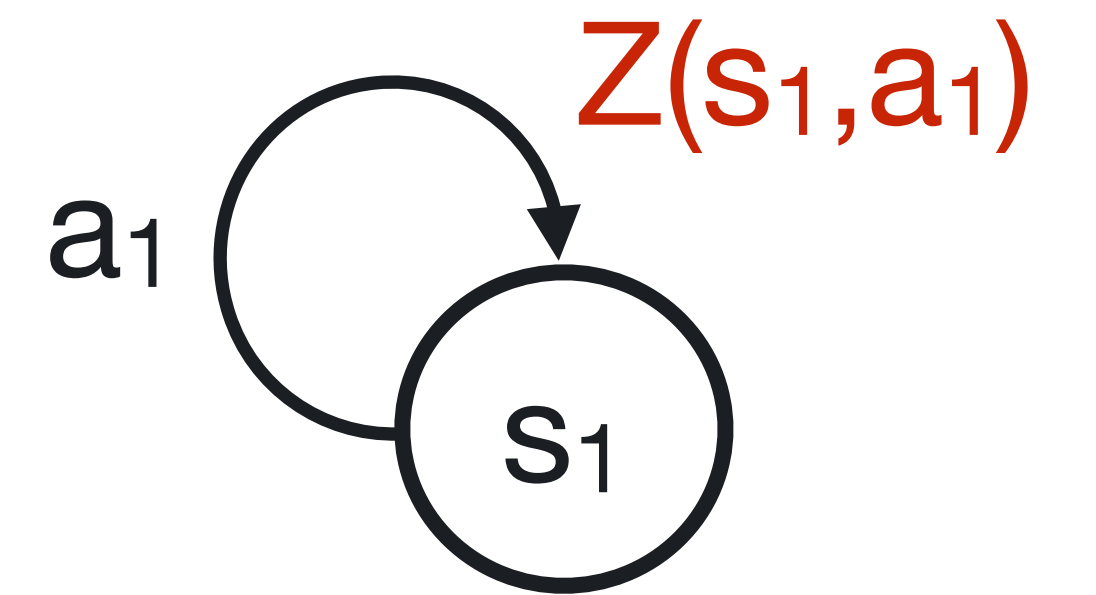
- Fixed **quantile** bins $\tau_0, \tau_1, \tau_2, \dots, \tau_N$
- Learn support $z_i = F_Z^{-1}(\tau_i)$
- No need for value range of Z
- **Training loss**
- **Fixed quantile bins**



3. Implicit Quantile Inverse CDF

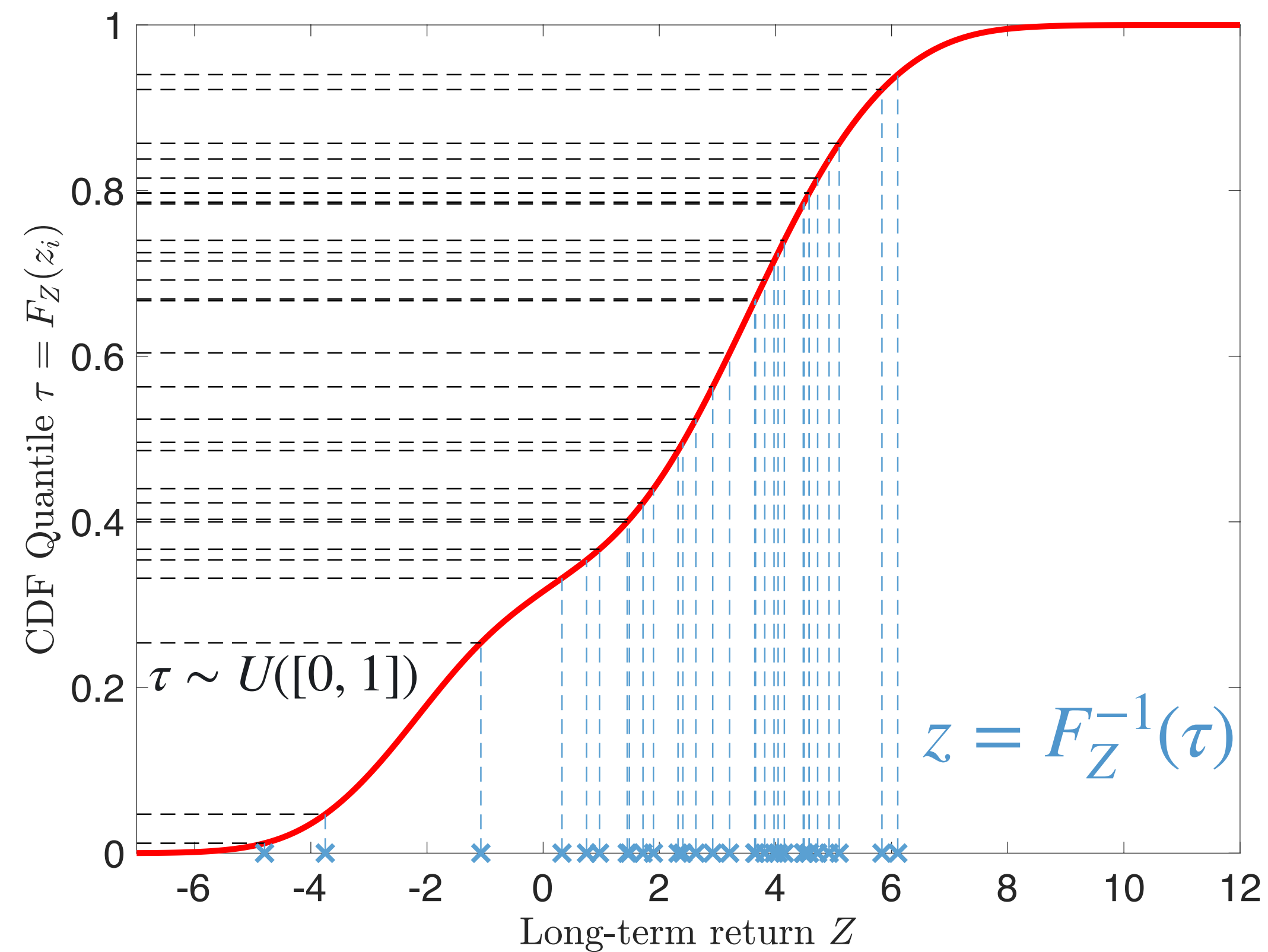
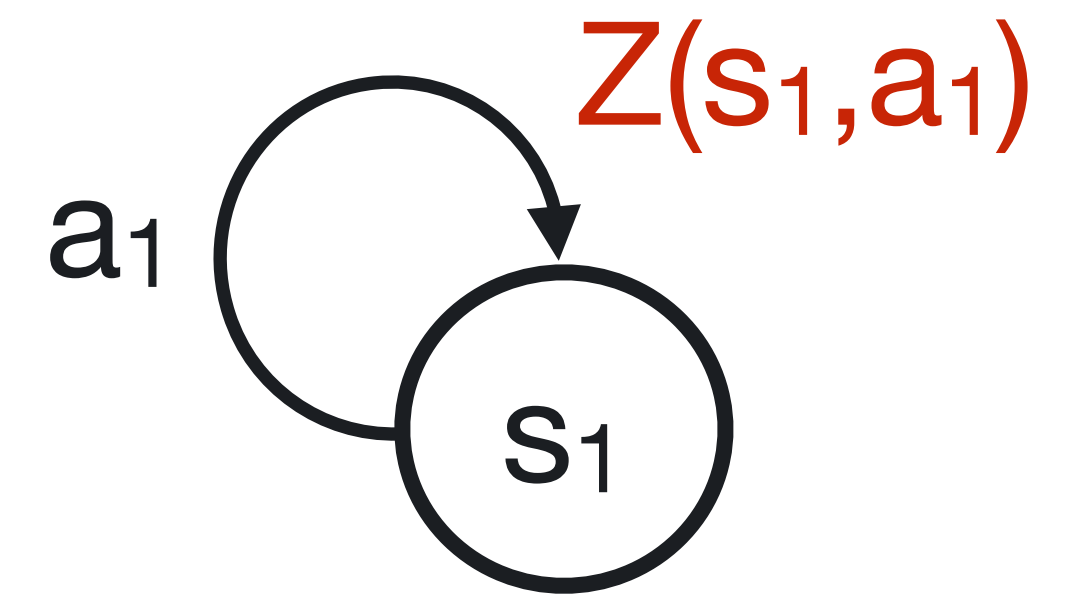


3. Implicit Quantile Inverse CDF



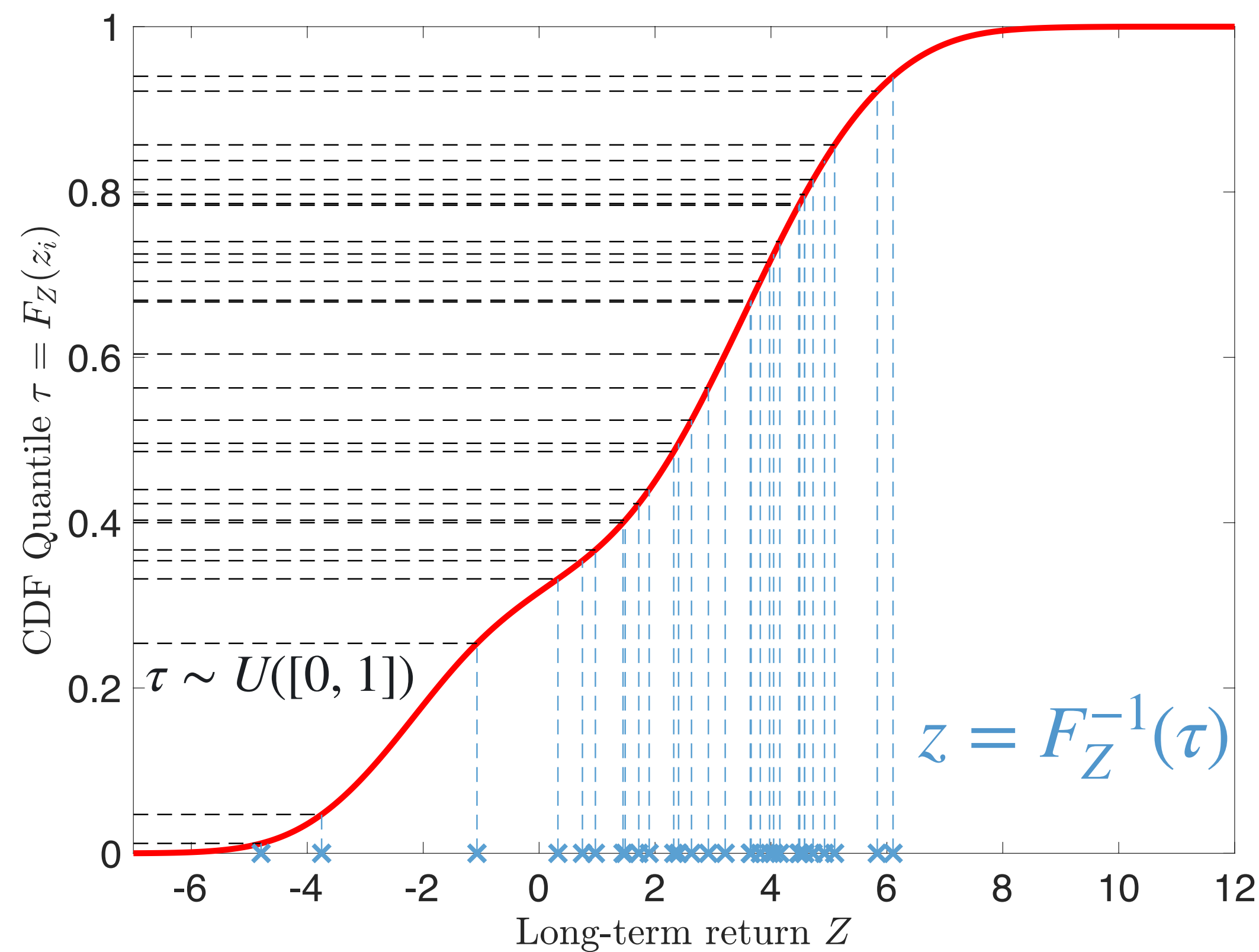
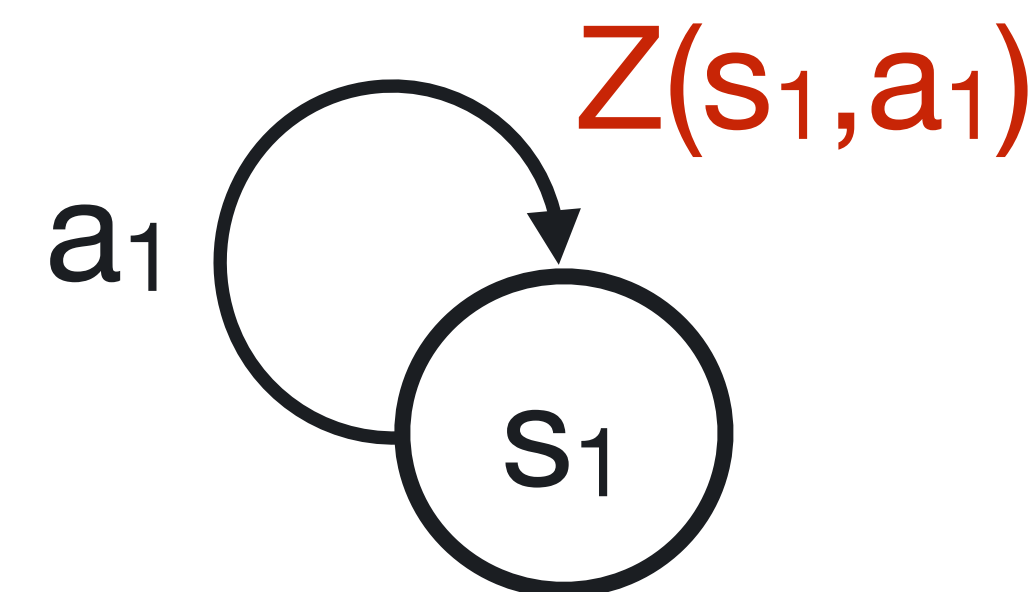
3. Implicit Quantile Inverse CDF

- Quantile sample (input) $\tau \sim U([0, 1])$



3. Implicit Quantile Inverse CDF

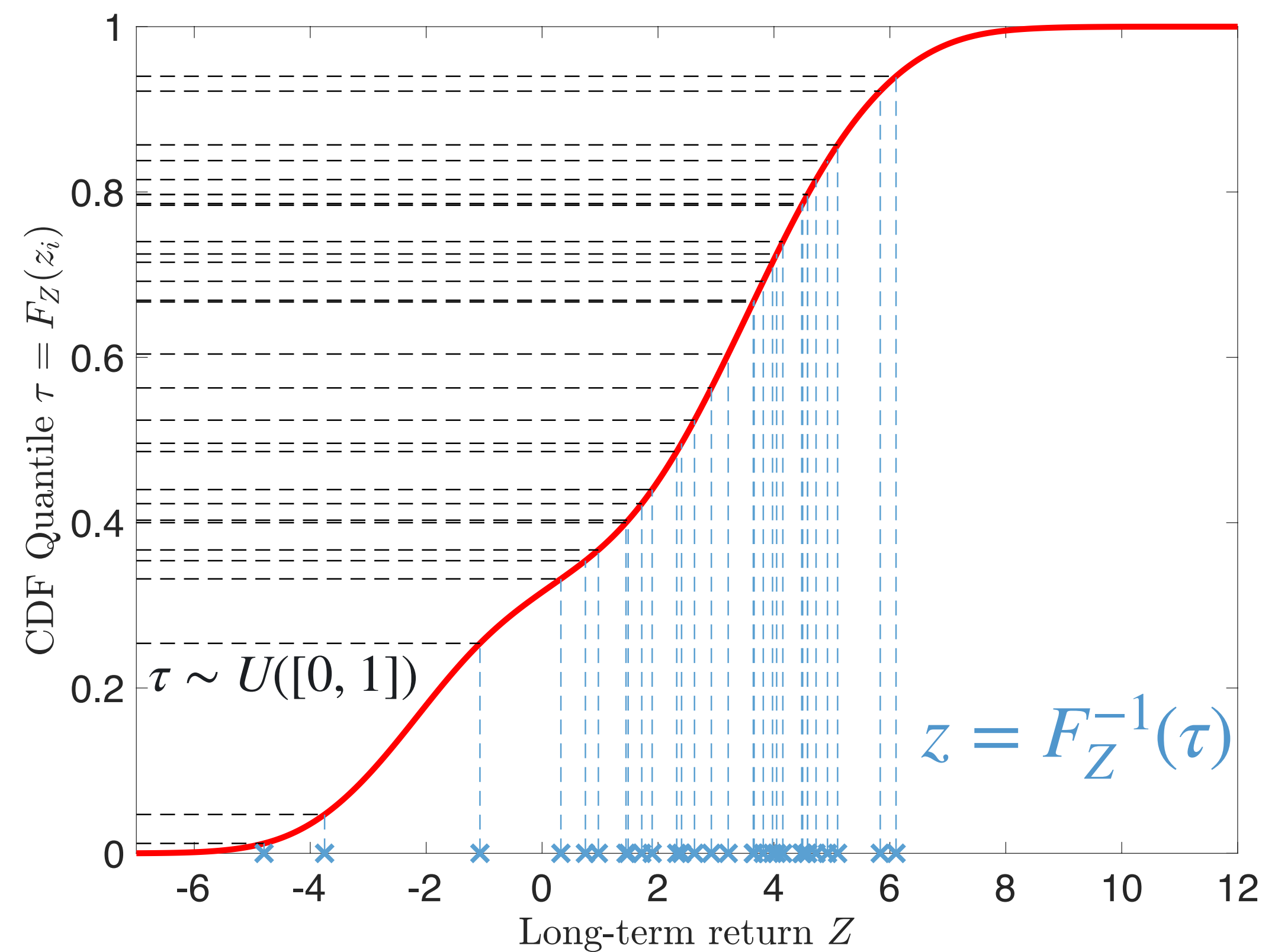
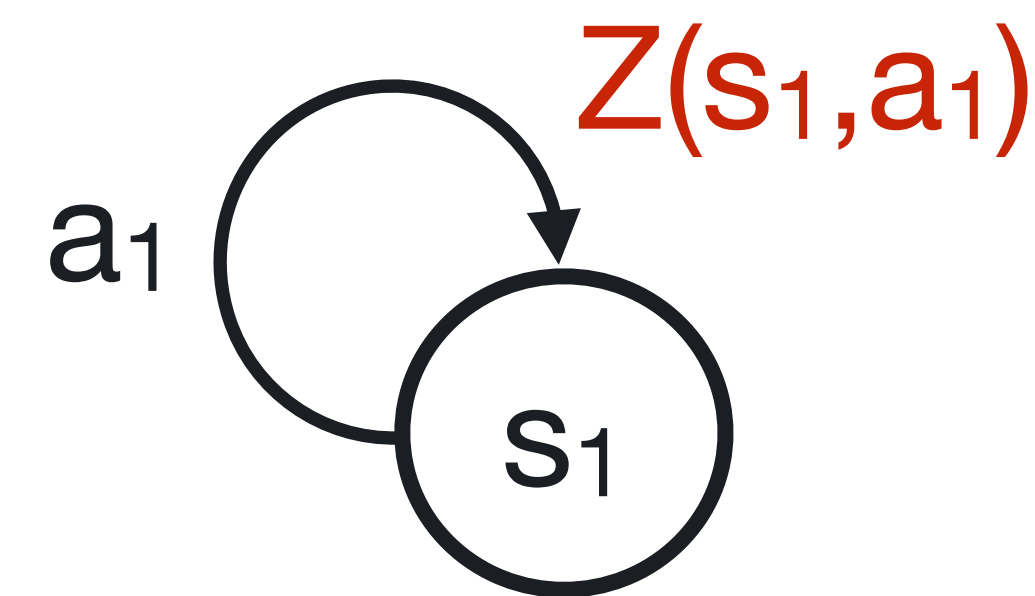
- Quantile sample (input) $\tau \sim U([0, 1])$
- Learn support $z = F_Z^{-1}(\tau)$

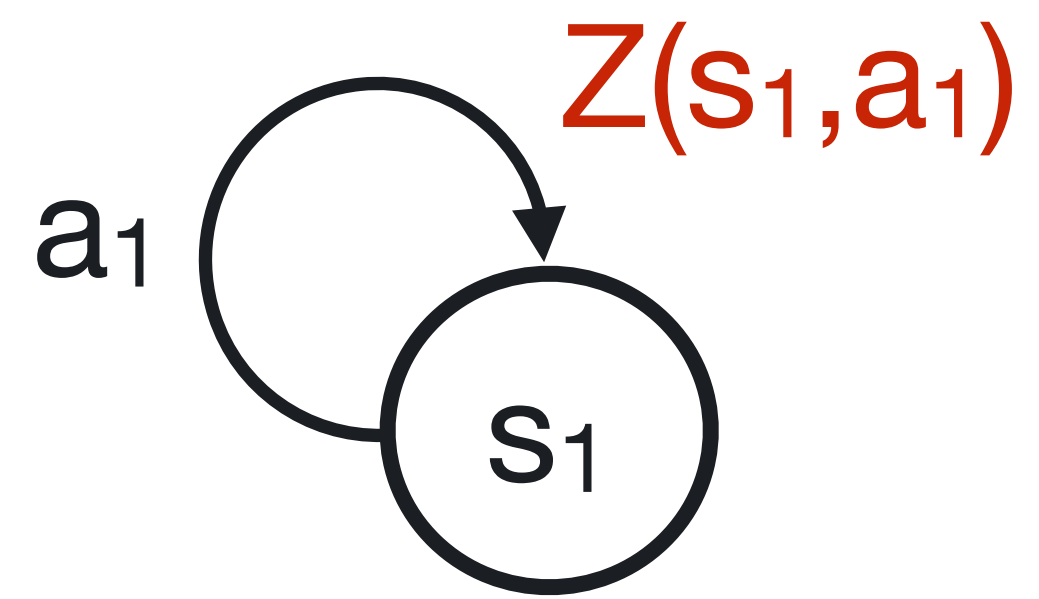


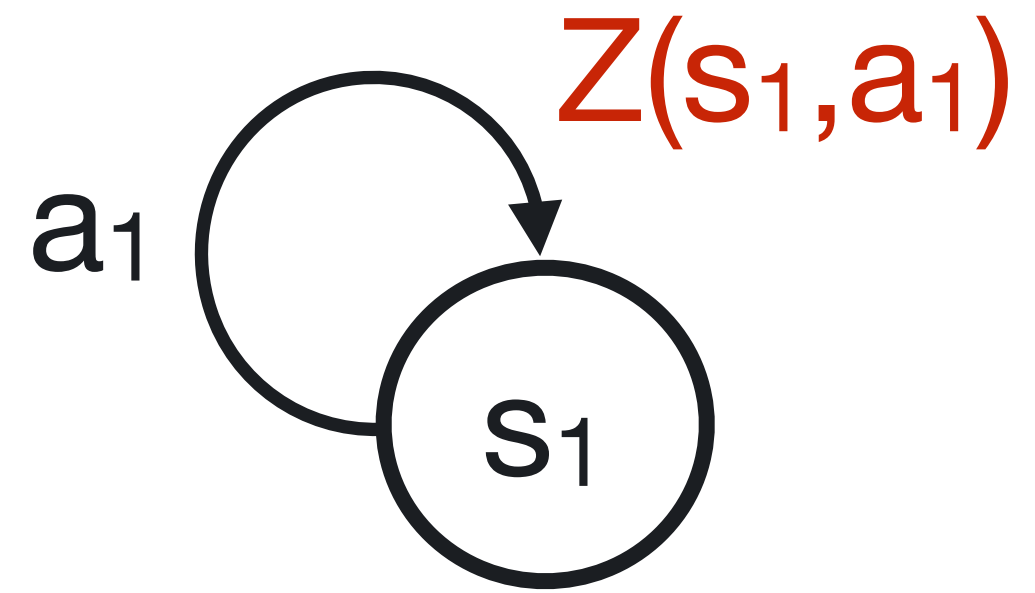
3. Implicit Quantile Inverse CDF

- Quantile sample (input)
- Learn support $z = F_Z^{-1}(\tau)$
- Training loss

$$\tau \sim U([0, 1])$$

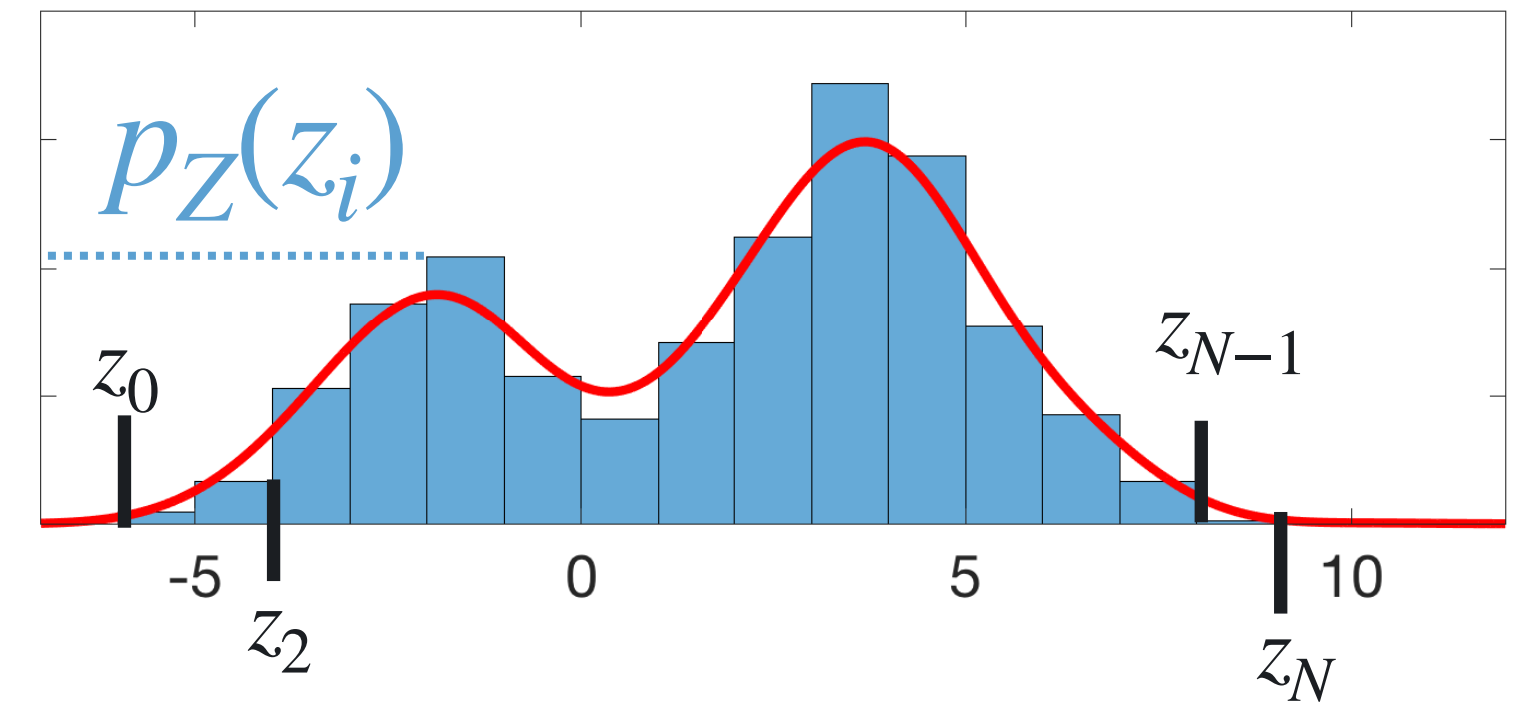


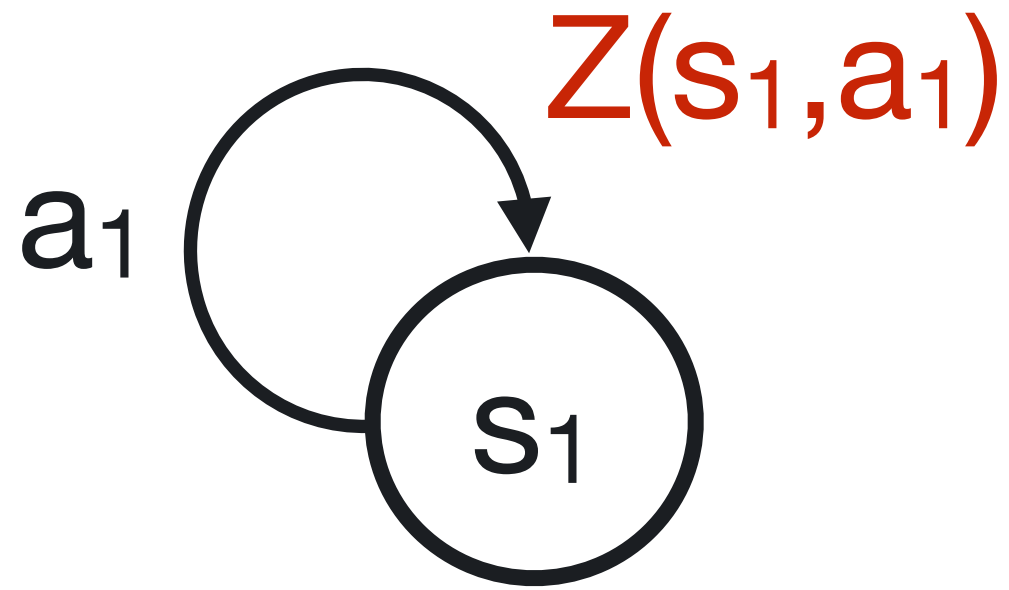




1. Categorical PDF

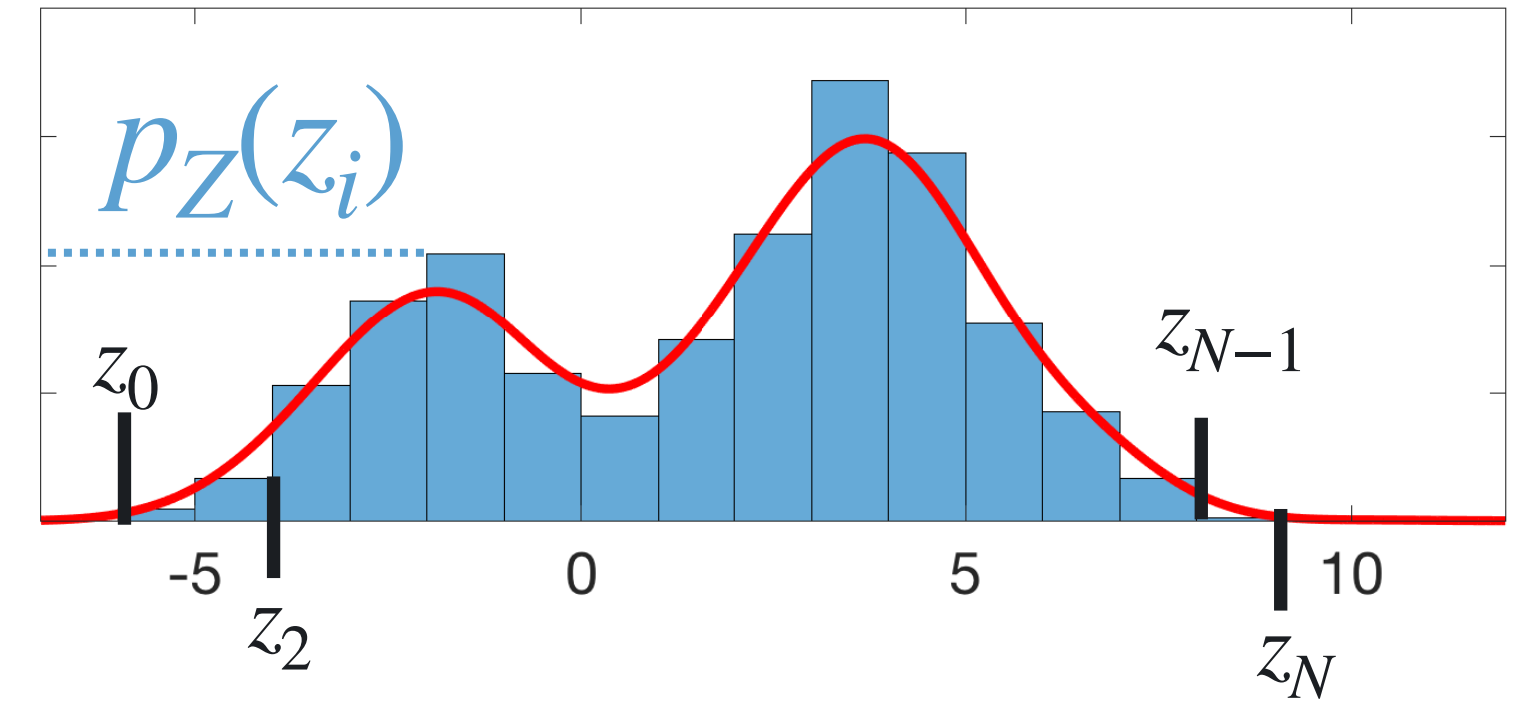
- Fixed support bins z_0, \dots, z_N
- Learn probabilities $p_Z(z_i)$





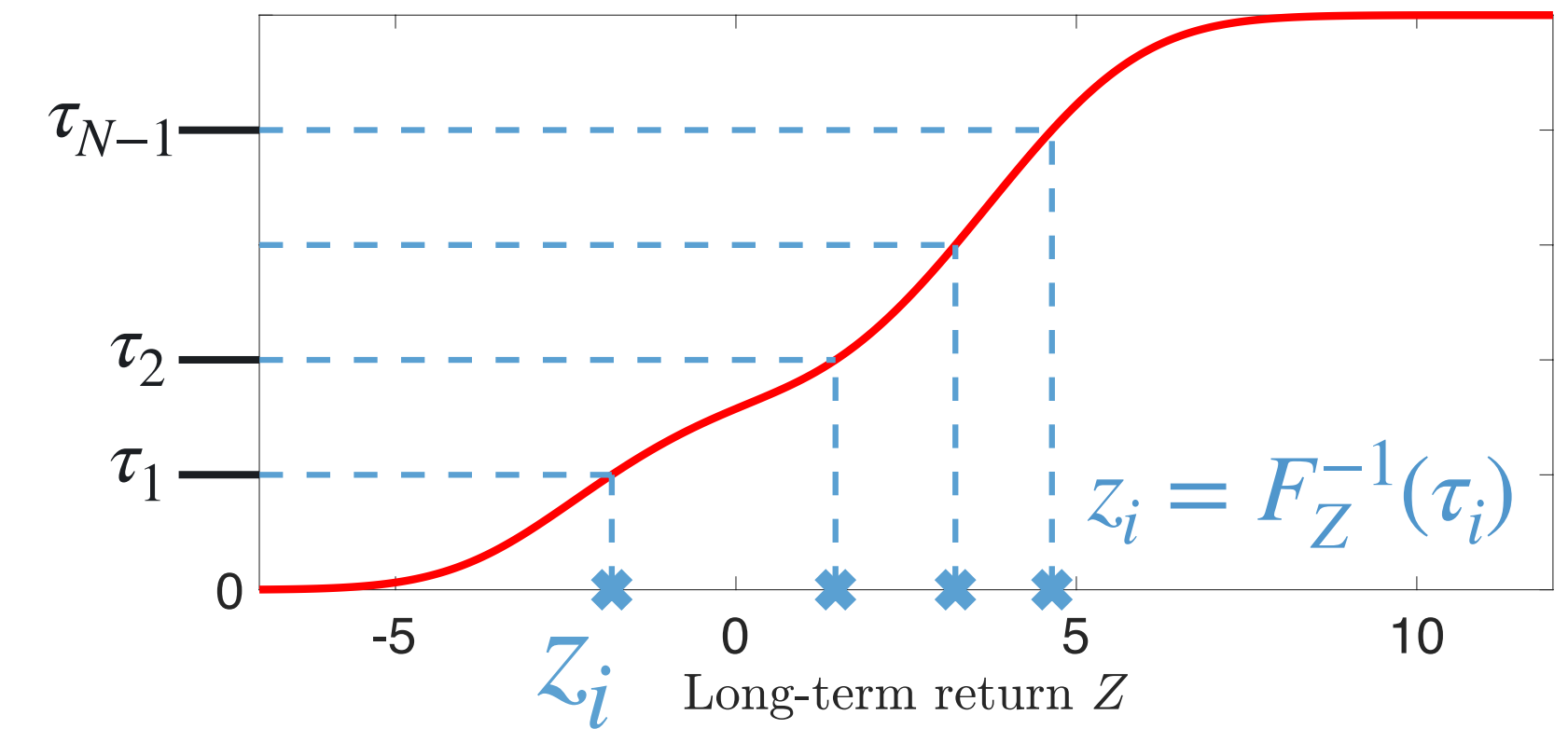
1. Categorical PDF

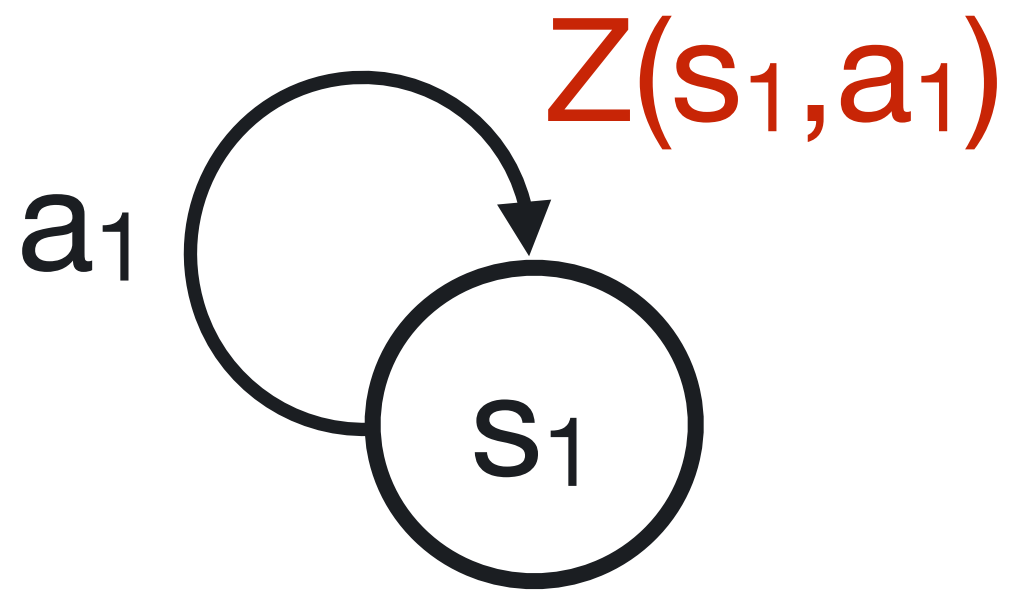
- Fixed support bins z_0, \dots, z_N
- Learn probabilities $p_Z(z_i)$



2. Quantile Inverse CDF

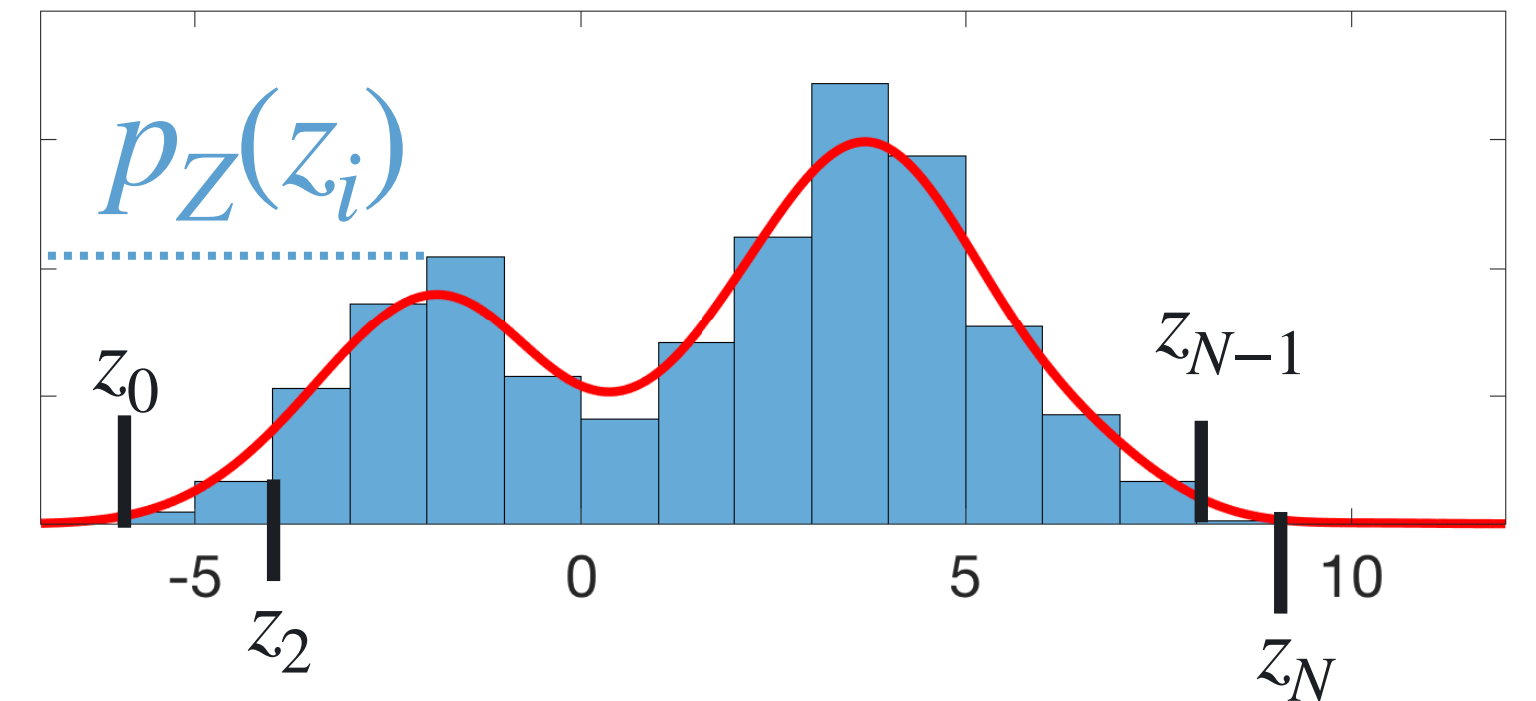
- Fixed quantile bins τ_0, \dots, τ_N
- Learn support values $z_i = F_Z^{-1}(\tau_i)$





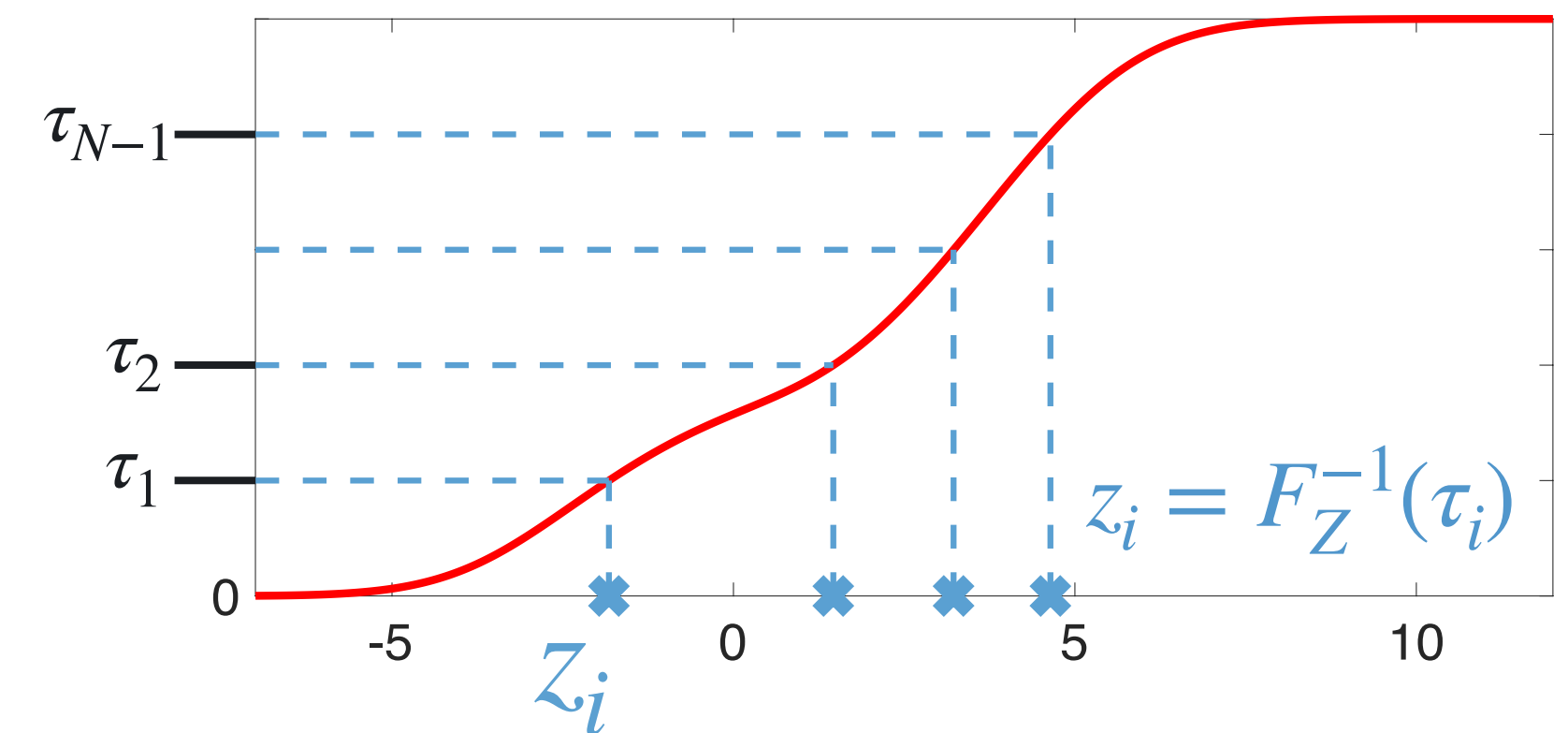
1. Categorical PDF

- Fixed support bins z_0, \dots, z_N
- Learn probabilities $p_Z(z_i)$



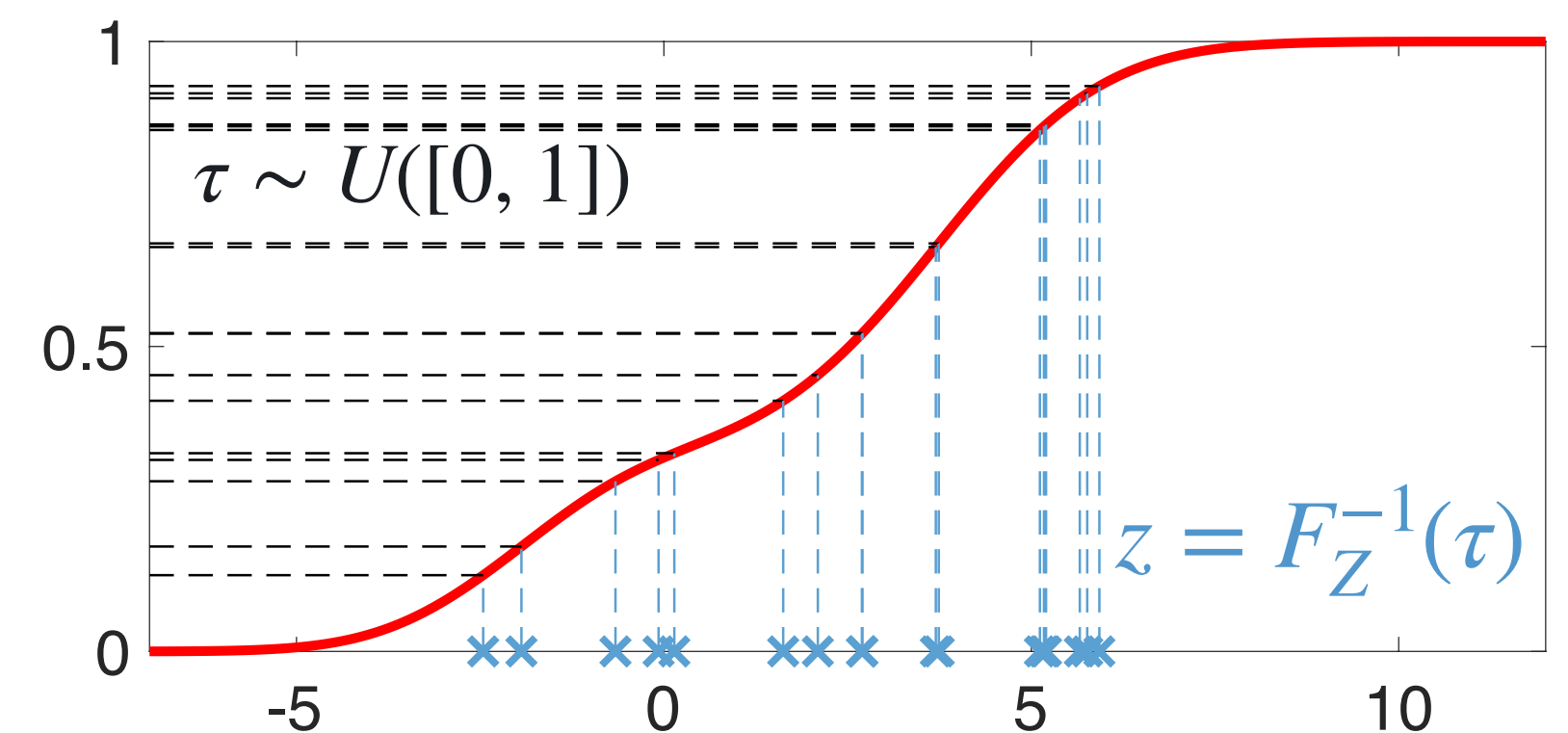
2. Quantile Inverse CDF

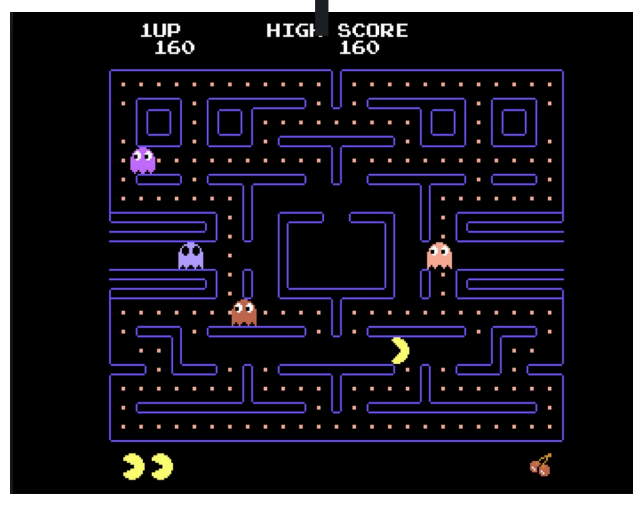
- Fixed quantile bins τ_0, \dots, τ_N
- Learn support values $z_i = F_Z^{-1}(\tau_i)$



3. Implicit Quantile Inverse CDF

- Quantile sampled from uniform $\tau \sim U([0, 1])$
- Learn support values $z = F_Z^{-1}(\tau)$

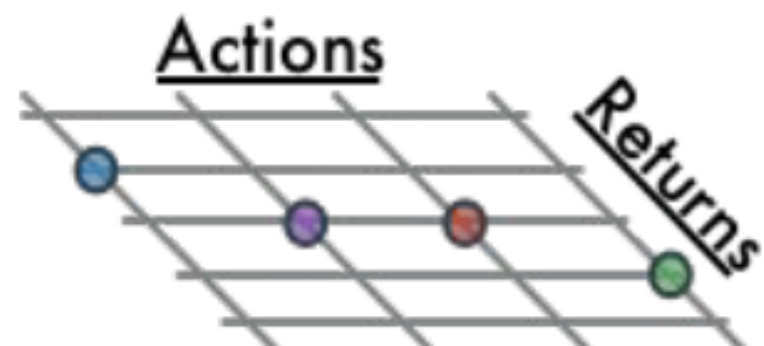




state

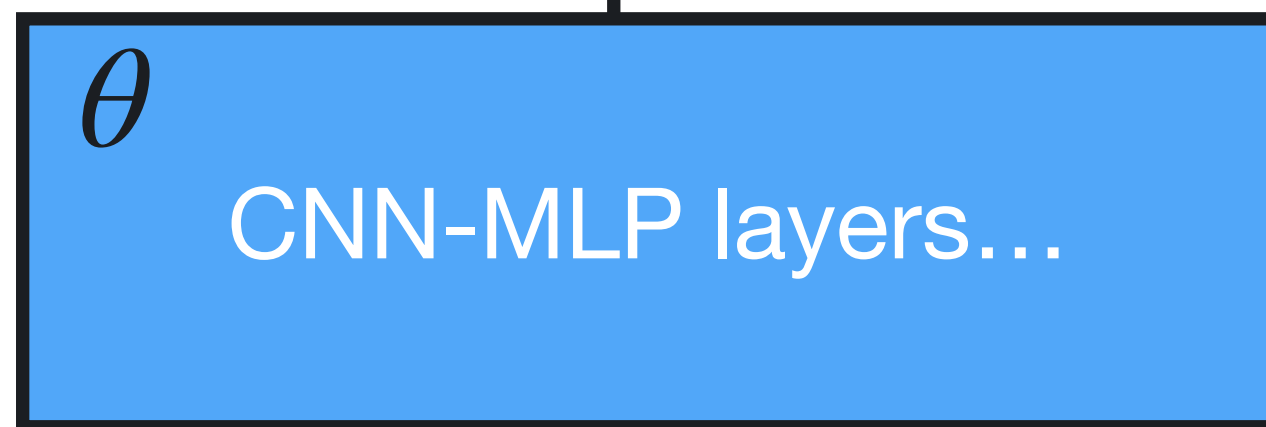
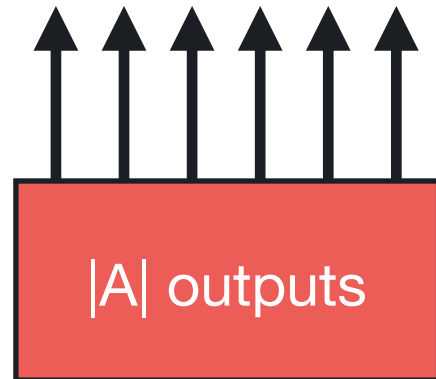
DQN

Mean



return for each action

$$Q(s, a_1) \dots Q(s, a_{|A|})$$

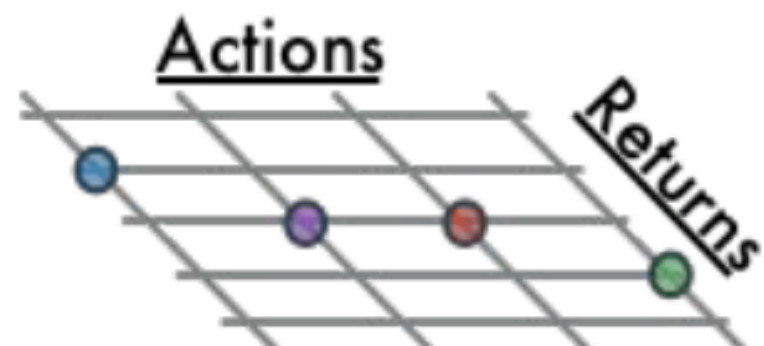


state



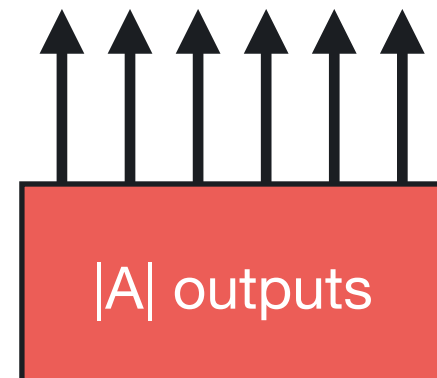
DQN

Mean



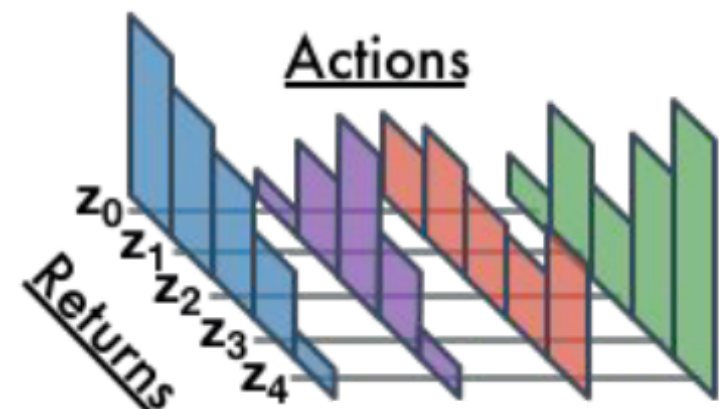
return for each action

$$Q(s, a_1) \dots Q(s, a_{|A|})$$



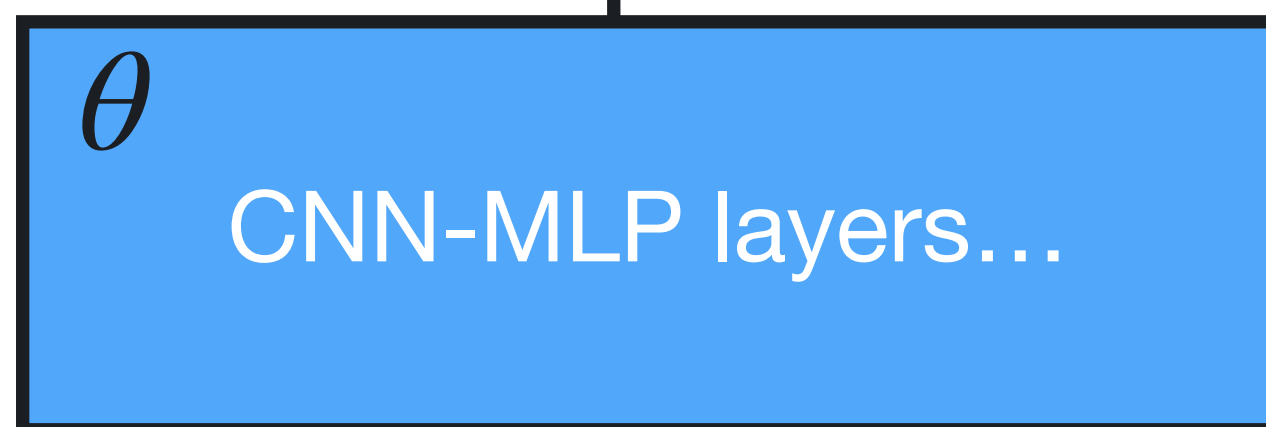
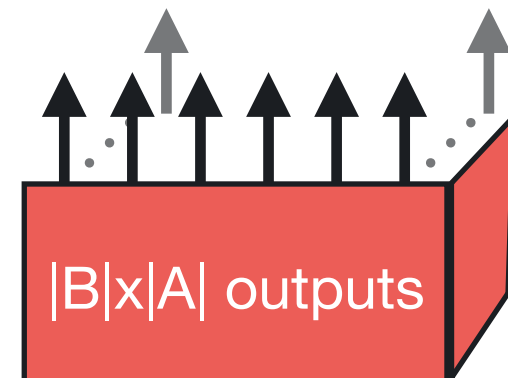
C51

Categorical PDF



bin probabilities for each action

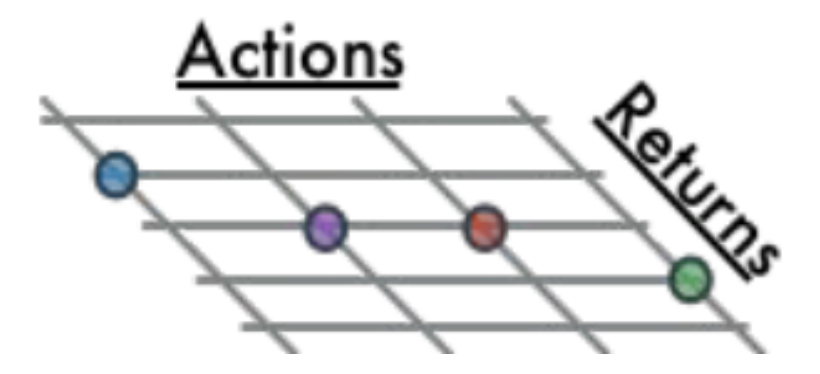
$$\{p_i\}^{a_1} \dots \{p_i\}^{a_{|A|}}$$



state

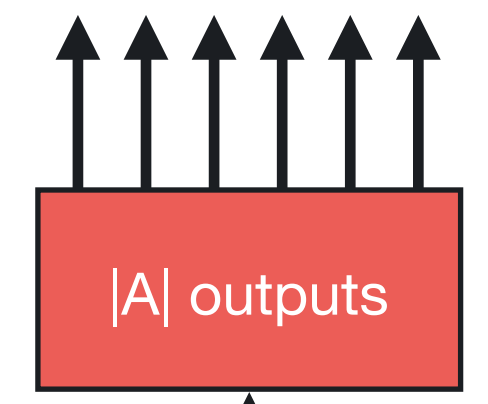
DQN

Mean



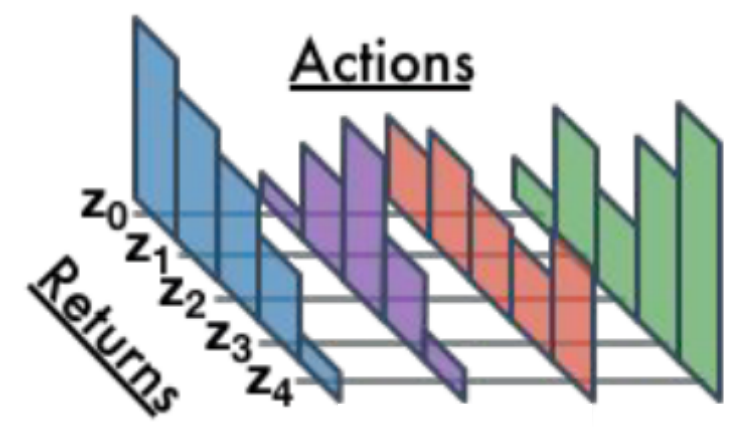
return for each action

$$Q(s, a_1) \dots Q(s, a_{|A|})$$



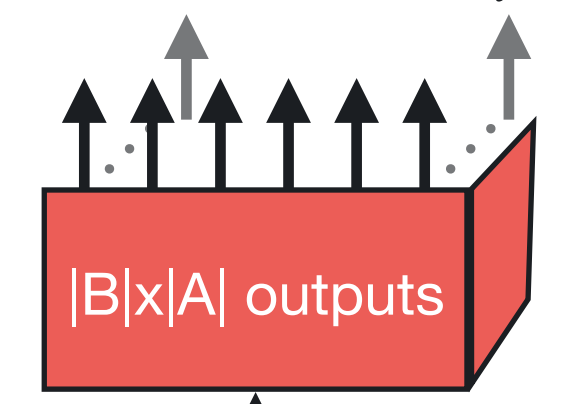
C51

Categorical PDF



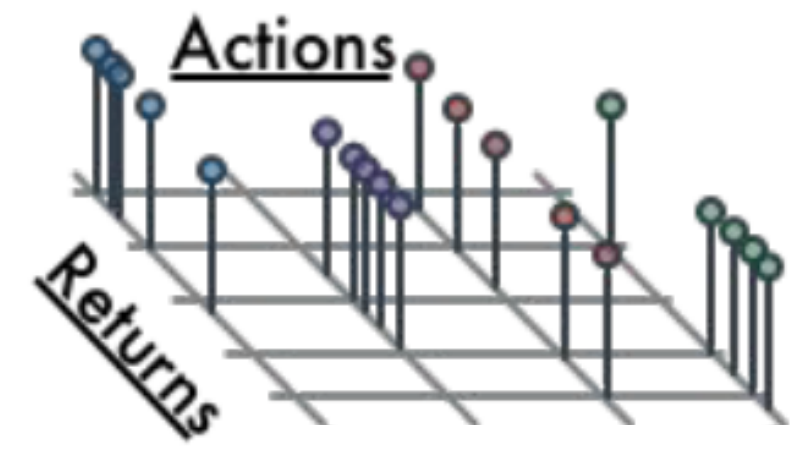
bin probabilities for each action

$$\{p_i\}^{a_1} \dots \{p_i\}^{a_{|A|}}$$



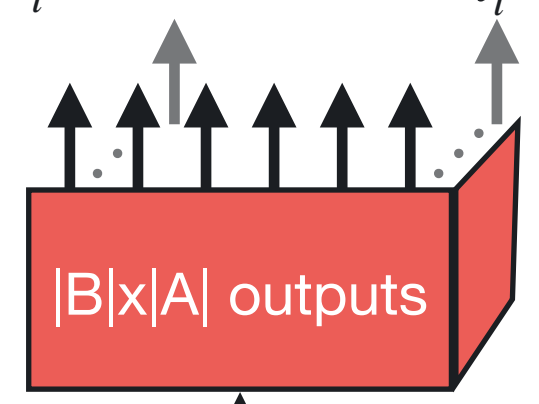
QR-DQN

Quantile Inverse CDF



support value of quantiles for each action

$$\{z_{\tau_i}\}^{a_1} \dots \{z_{\tau_i}\}^{a_{|A|}}$$

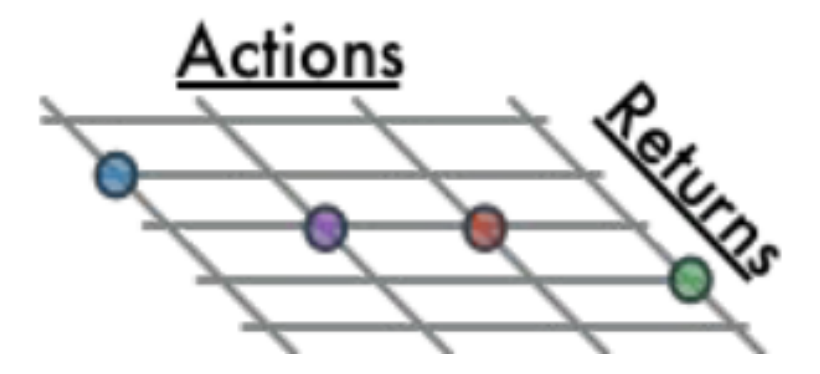


state



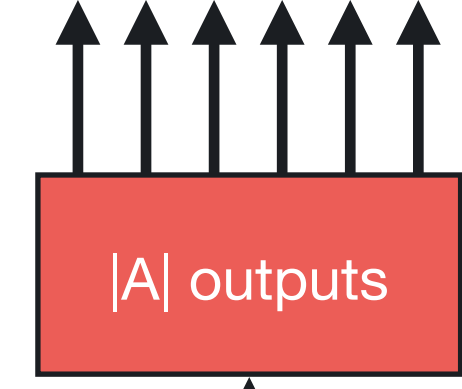
DQN

Mean



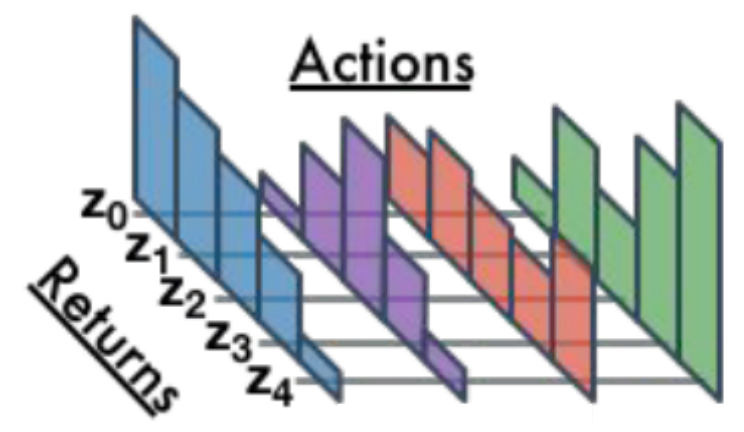
return for each action

$$Q(s, a_1) \dots Q(s, a_{|A|})$$



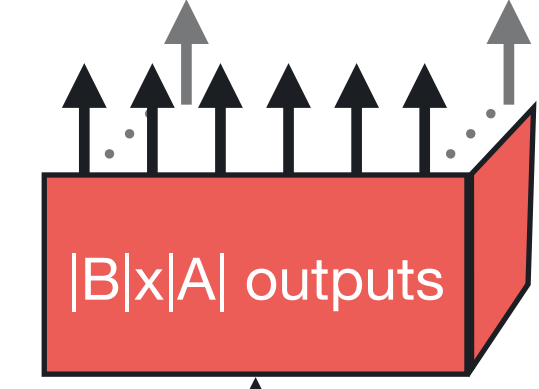
C51

Categorical PDF



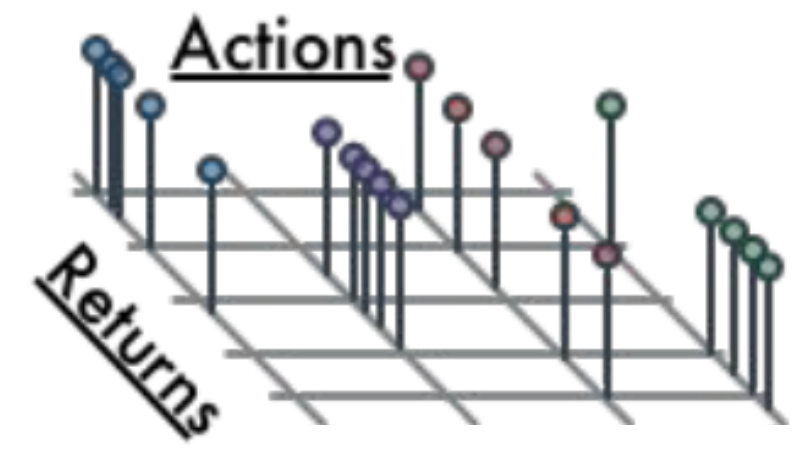
bin probabilities for each action

$$\{p_i\}^{a_1} \dots \{p_i\}^{a_{|A|}}$$



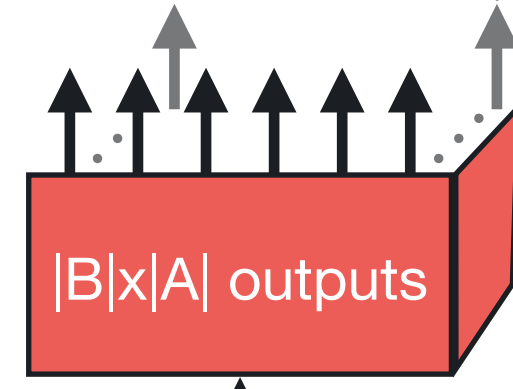
QR-DQN

Quantile Inverse CDF



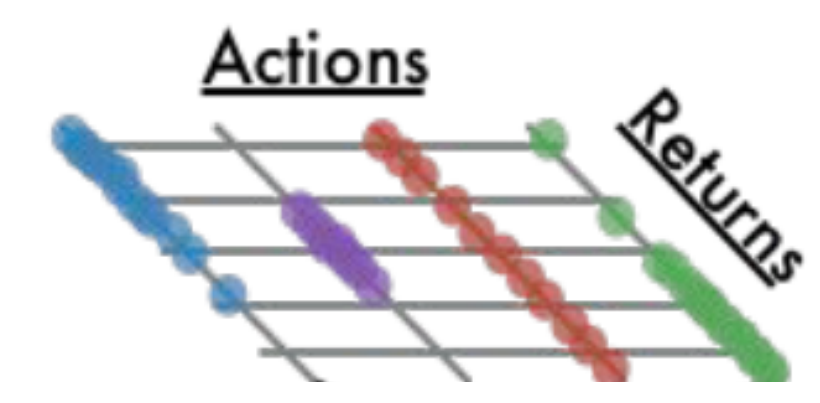
support value of quantiles for each action

$$\{z_{\tau_i}\}^{a_1} \dots \{z_{\tau_i}\}^{a_{|A|}}$$



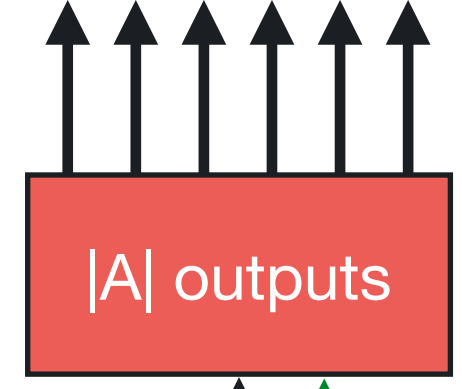
IQN

Implicit Quantile Inverse CDF



support value of sample quantiles for each action

$$z_{\tau}^{a_1} \dots z_{\tau}^{a_{|A|}}$$



$\tau \sim U([0,1])$

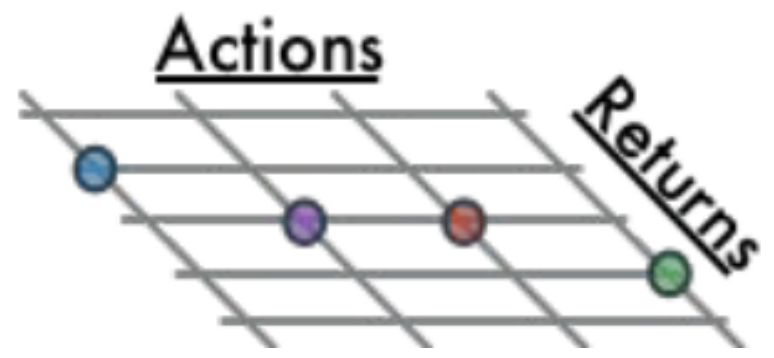


state

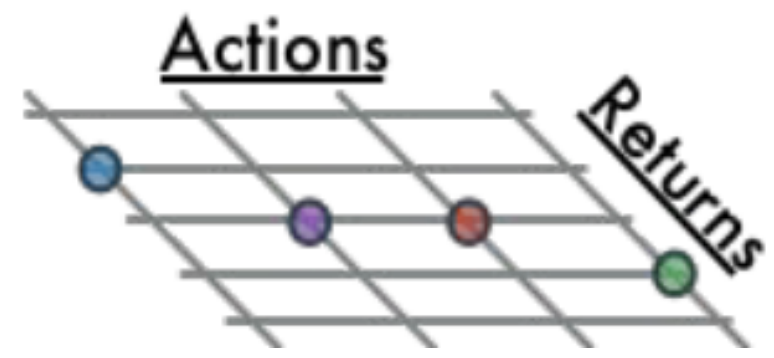


How do we train these networks ?

DQN

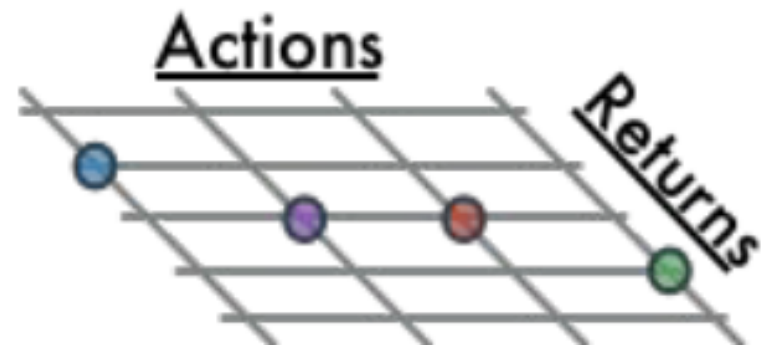


DQN



$$\{s_t, a_t, s_{t+1}, r_t\}$$

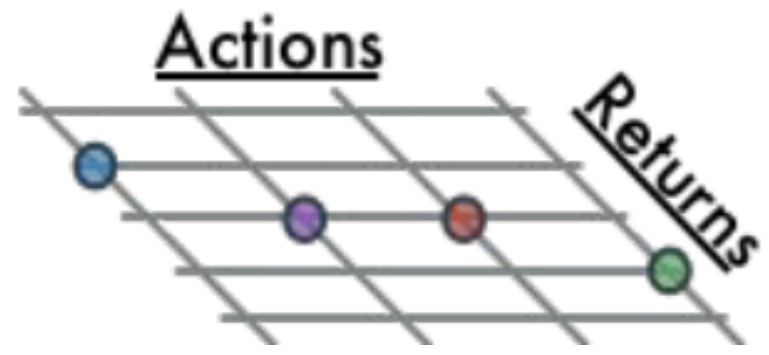
DQN



$$\{s_t, a_t, s_{t+1}, r_t\}$$

$$a^{\star} = \operatorname{argmax}_a Q^{\theta}(s_{t+1}, a)$$

DQN

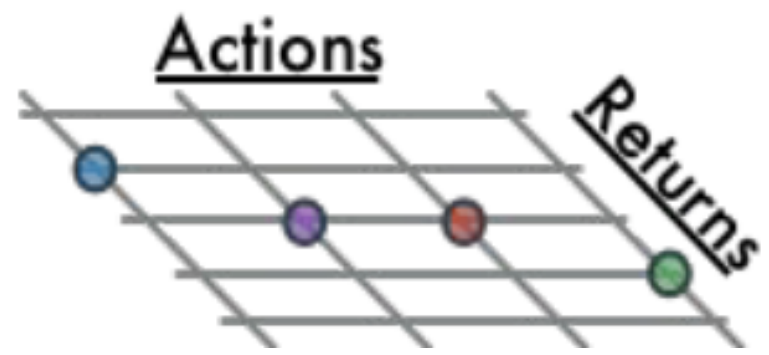


$$\{s_t, a_t, s_{t+1}, r_t\}$$

$$a^{\star} = \operatorname{argmax}_a Q^{\theta}(s_{t+1}, a)$$

$$q' = Q^{\theta}(s_{t+1}, a^{\star})$$

DQN



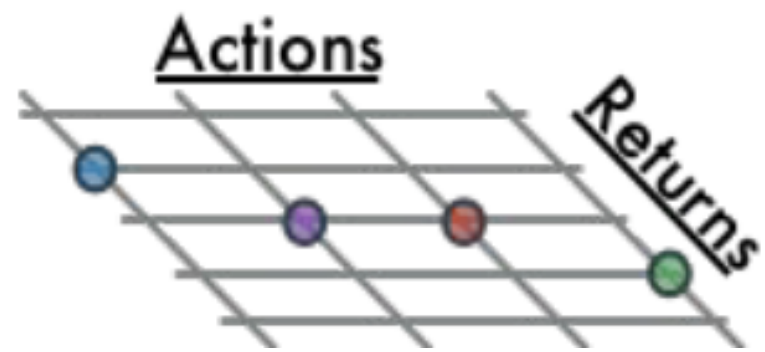
$$\{s_t, a_t, s_{t+1}, r_t\}$$

$$a^{\star} = \operatorname{argmax}_a Q^{\theta}(s_{t+1}, a)$$

$$q' = Q^{\theta}(s_{t+1}, a^{\star})$$

$$q = Q^{\theta}(s_t, a_t)$$

DQN



$$\{s_t, a_t, s_{t+1}, r_t\}$$

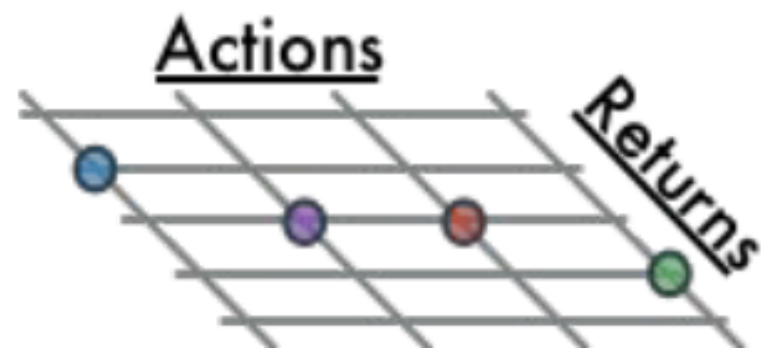
$$a^{\star} = \underset{a}{\operatorname{argmax}} Q^{\theta}(s_{t+1}, a)$$

$$q' = Q^{\theta}(s_{t+1}, a^{\star})$$

$$q = Q^{\theta}(s_t, a_t)$$

$$\delta_t = r_t + \gamma q' - q \quad \text{temp. diff.}$$

DQN



$$\{s_t, a_t, s_{t+1}, r_t\}$$

$$a^{\star} = \operatorname{argmax}_a Q^{\theta}(s_{t+1}, a)$$

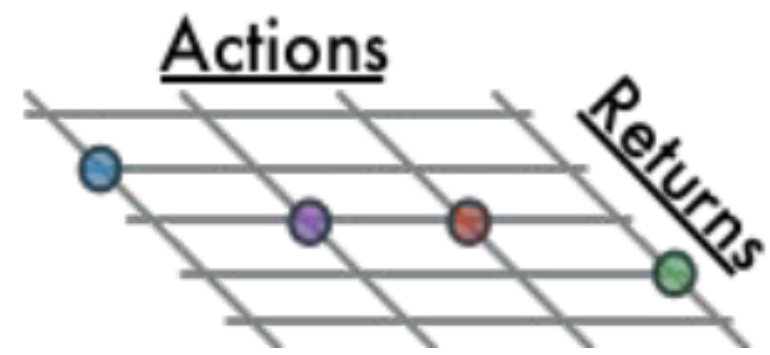
$$q' = Q^{\theta}(s_{t+1}, a^{\star})$$

$$q = Q^{\theta}(s_t, a_t)$$

$$\delta_t = r_t + \gamma q' - q \quad \text{temp. diff.}$$

$$\mathcal{L}_{DQN} = \delta_t^2$$

DQN



$$\{s_t, a_t, s_{t+1}, r_t\}$$

$$a^* = \operatorname{argmax}_a Q^\theta(s_{t+1}, a)$$

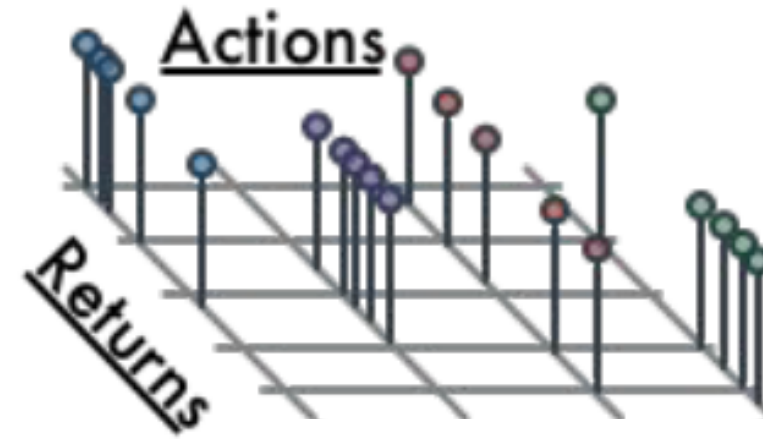
$$q' = Q^\theta(s_{t+1}, a^*)$$

$$q = Q^\theta(s_t, a_t)$$

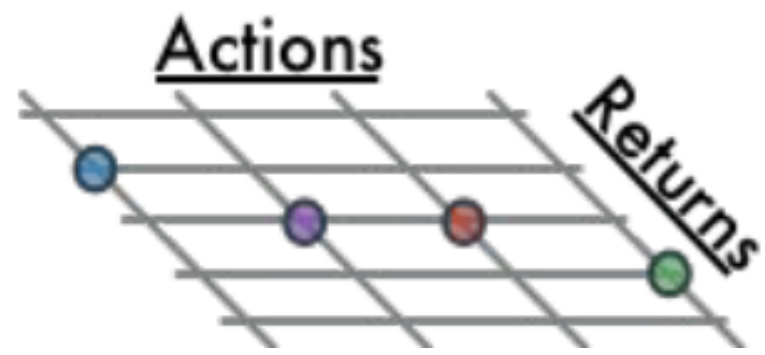
$$\delta_t = r_t + \gamma q' - q \quad \text{temp. diff.}$$

$$\mathcal{L}_{DQN} = \delta_t^2$$

QR-DQN



DQN



$$\{s_t, a_t, s_{t+1}, r_t\}$$

$$a^{\star} = \operatorname{argmax}_a Q^{\theta}(s_{t+1}, a)$$

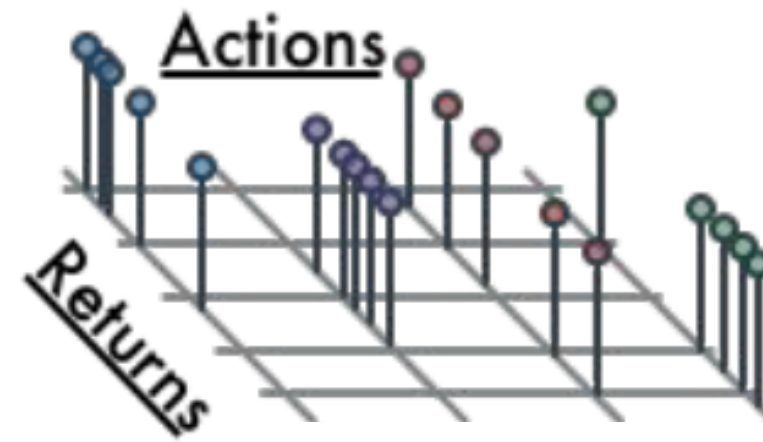
$$q' = Q^{\theta}(s_{t+1}, a^{\star})$$

$$q = Q^{\theta}(s_t, a_t)$$

$$\delta_t = r_t + \gamma q' - q \quad \text{temp. diff.}$$

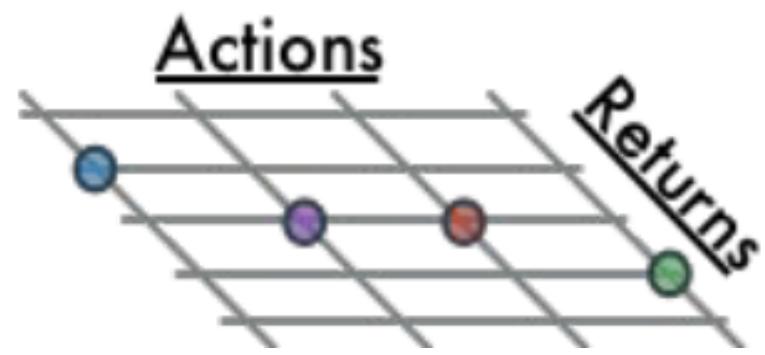
$$\mathcal{L}_{DQN} = \delta_t^2$$

QR-DQN



$$\{s_t, a_t, s_{t+1}, r_t\}$$

DQN



$$\{s_t, a_t, s_{t+1}, r_t\}$$

$$a^* = \operatorname{argmax}_a Q^\theta(s_{t+1}, a)$$

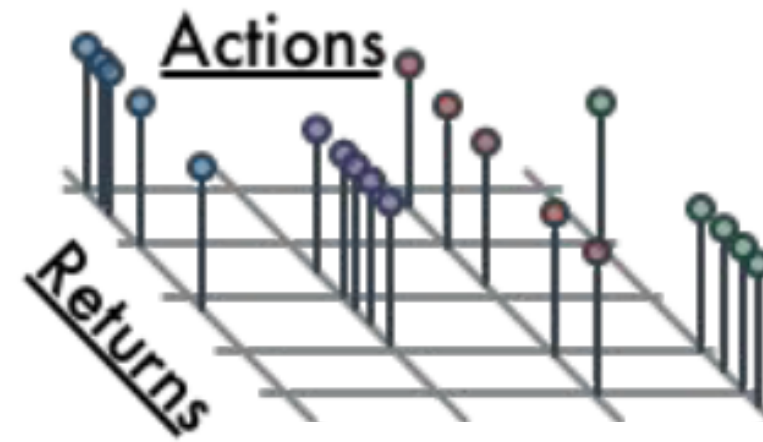
$$q' = Q^\theta(s_{t+1}, a^*)$$

$$q = Q^\theta(s_t, a_t)$$

$$\delta_t = r_t + \gamma q' - q \quad \text{temp. diff.}$$

$$\mathcal{L}_{DQN} = \delta_t^2$$

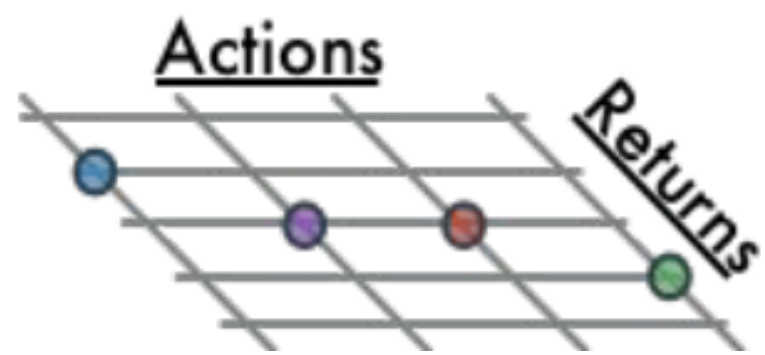
QR-DQN



$$\{s_t, a_t, s_{t+1}, r_t\}$$

$$a^* = \operatorname{argmax}_a \mathbb{E} [Z_\tau^\theta(s_{t+1}, a)]$$

DQN



$$\{s_t, a_t, s_{t+1}, r_t\}$$

$$a^* = \operatorname{argmax}_a Q^\theta(s_{t+1}, a)$$

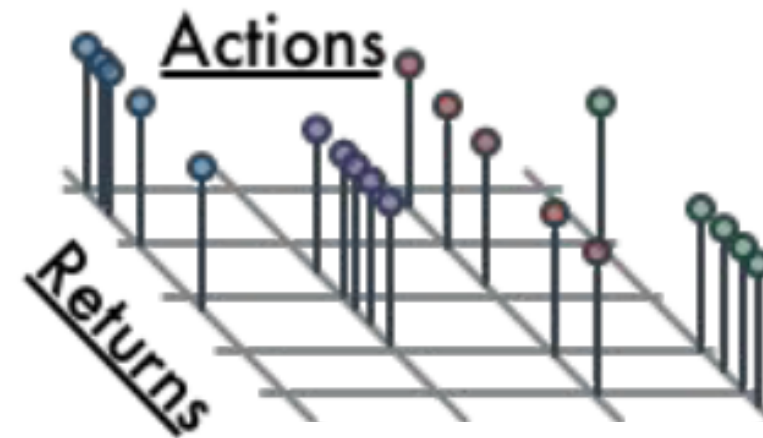
$$q' = Q^\theta(s_{t+1}, a^*)$$

$$q = Q^\theta(s_t, a_t)$$

$$\delta_t = r_t + \gamma q' - q \quad \text{temp. diff.}$$

$$\mathcal{L}_{DQN} = \delta_t^2$$

QR-DQN



$$\{s_t, a_t, s_{t+1}, r_t\}$$

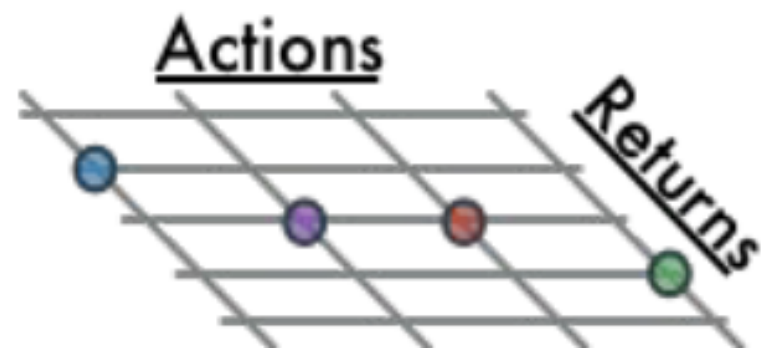
$$a^* = \operatorname{argmax}_a \mathbb{E} [Z_\tau^\theta(s_{t+1}, a)]$$

$$\forall \tau, \tau' \quad \left| \quad z' = Z_{\tau'}^\theta(s_{t+1}, a^*) \right.$$

$$\left. \quad z = Z_\tau^\theta(s_t, a_t) \right.$$

$$\left. \quad \delta_t^{\tau, \tau'} = r_t + \gamma z' - z \right.$$

DQN



$$\{s_t, a_t, s_{t+1}, r_t\}$$

$$a^* = \operatorname{argmax}_a Q^\theta(s_{t+1}, a)$$

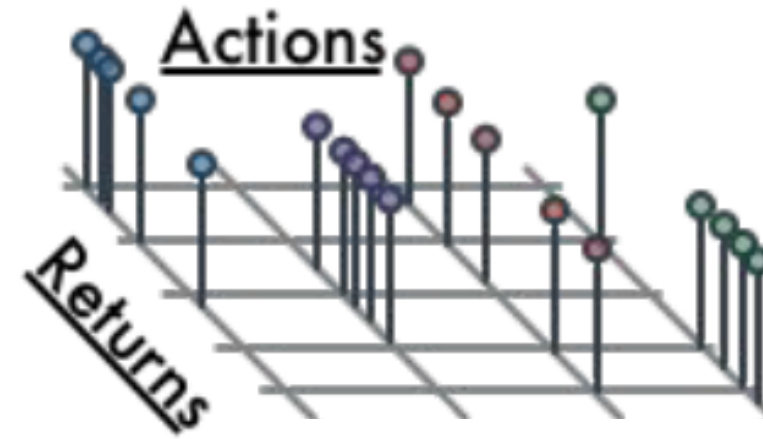
$$q' = Q^\theta(s_{t+1}, a^*)$$

$$q = Q^\theta(s_t, a_t)$$

$$\delta_t = r_t + \gamma q' - q \quad \text{temp. diff.}$$

$$\mathcal{L}_{DQN} = \delta_t^2$$

QR-DQN



$$\{s_t, a_t, s_{t+1}, r_t\}$$

$$a^* = \operatorname{argmax}_a \mathbb{E} [Z_\tau^\theta(s_{t+1}, a)]$$

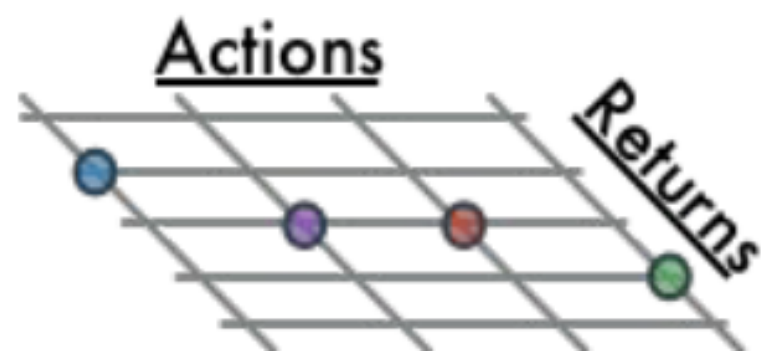
$$\forall \tau, \tau' \quad \left| \begin{array}{l} z' = Z_{\tau'}^\theta(s_{t+1}, a^*) \\ z = Z_\tau^\theta(s_t, a_t) \end{array} \right.$$

$$z = Z_\tau^\theta(s_t, a_t)$$

$$\delta_t^{\tau, \tau'} = r_t + \gamma z' - z$$

$$\mathcal{L}_{QR-DQN} = ?$$

DQN



$$\{s_t, a_t, s_{t+1}, r_t\}$$

$$a^* = \operatorname{argmax}_a Q^\theta(s_{t+1}, a)$$

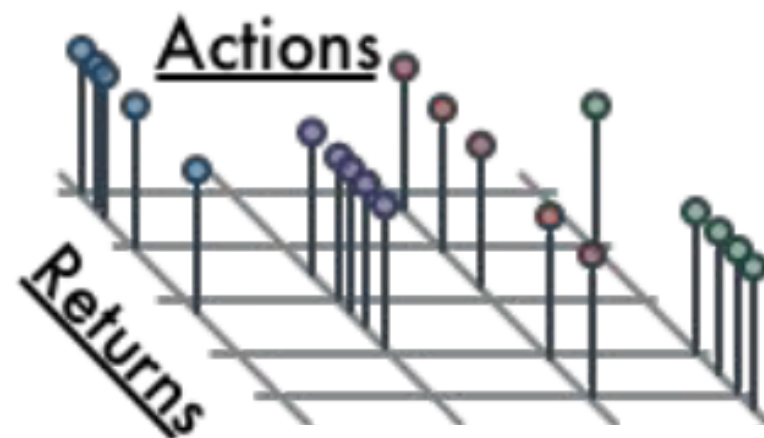
$$q' = Q^\theta(s_{t+1}, a^*)$$

$$q = Q^\theta(s_t, a_t)$$

$$\delta_t = r_t + \gamma q' - q \quad \text{temp. diff.}$$

$$\mathcal{L}_{DQN} = \delta_t^2$$

QR-DQN



$$\{s_t, a_t, s_{t+1}, r_t\}$$

$$a^* = \operatorname{argmax}_a \mathbb{E} [Z_\tau^\theta(s_{t+1}, a)]$$

$$\forall \tau, \tau' \quad \left| \begin{array}{l} z' = Z_{\tau'}^\theta(s_{t+1}, a^*) \\ z = Z_\tau^\theta(s_t, a_t) \\ \delta_t^{\tau, \tau'} = r_t + \gamma z' - z \end{array} \right.$$

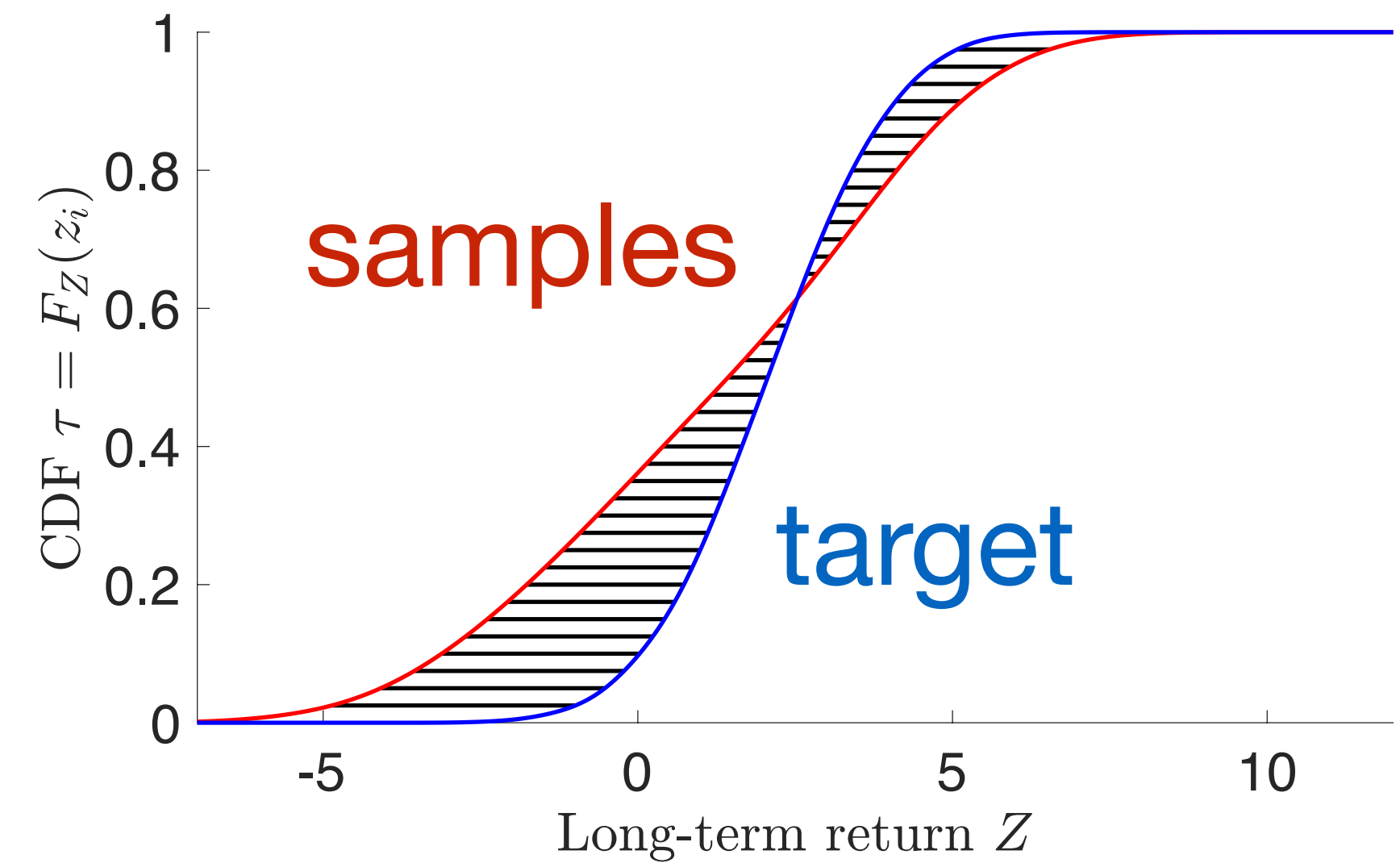
$$\mathcal{L}_{QR-DQN} = ?$$

Projection to Wasserstein metric!

**The distributional Bellman Operator is a contraction
on the Wasserstein metric**

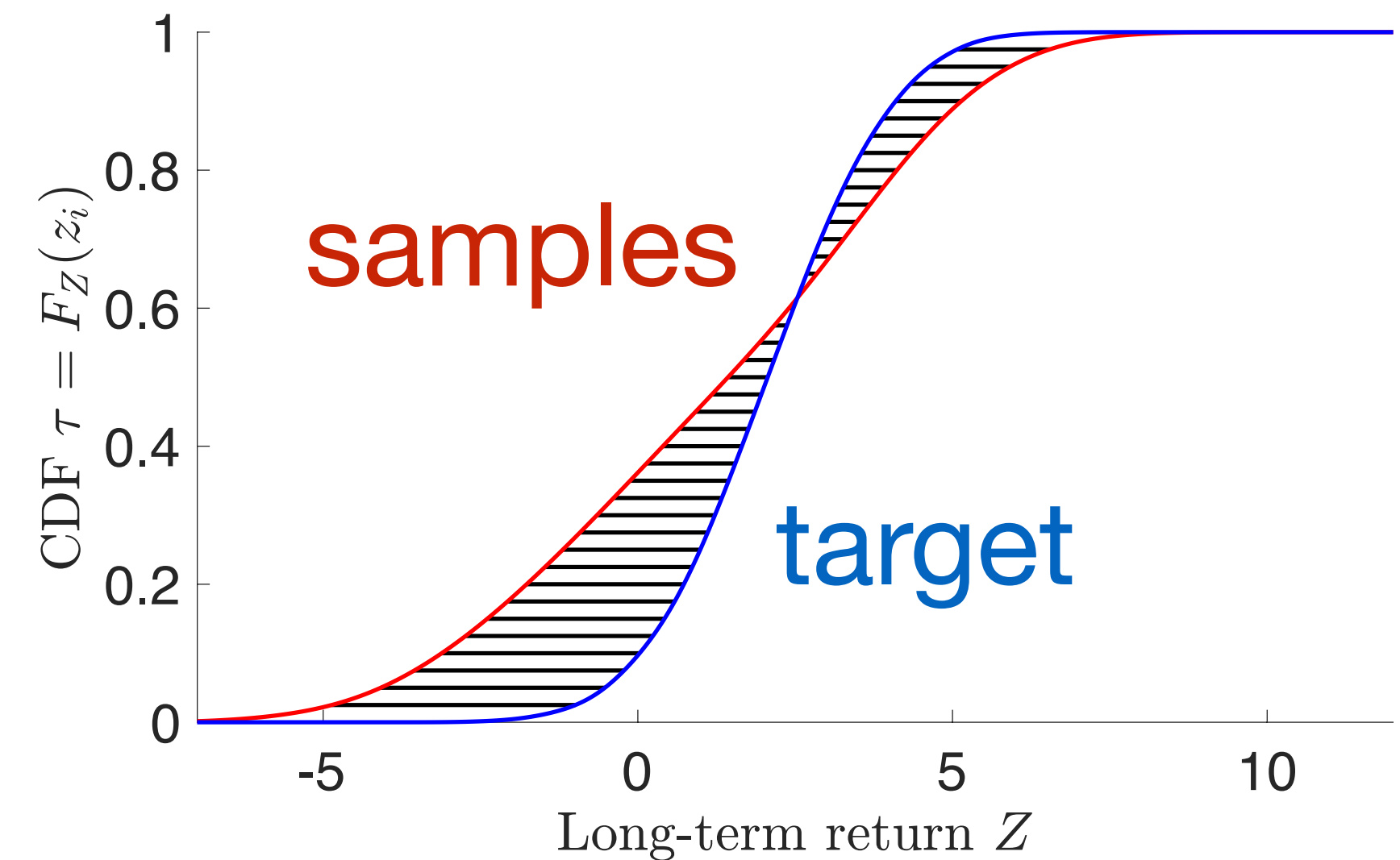
The distributional Bellman Operator is a contraction on the Wasserstein metric

$$w_1(X, Y) = \int_0^1 |F_X^{-1}(\tau) - F_Y^{-1}(\tau)| d\tau$$



The distributional Bellman Operator is a contraction on the Wasserstein metric

$$w_1(X, Y) = \int_0^1 |F_X^{-1}(\tau) - F_Y^{-1}(\tau)| d\tau$$



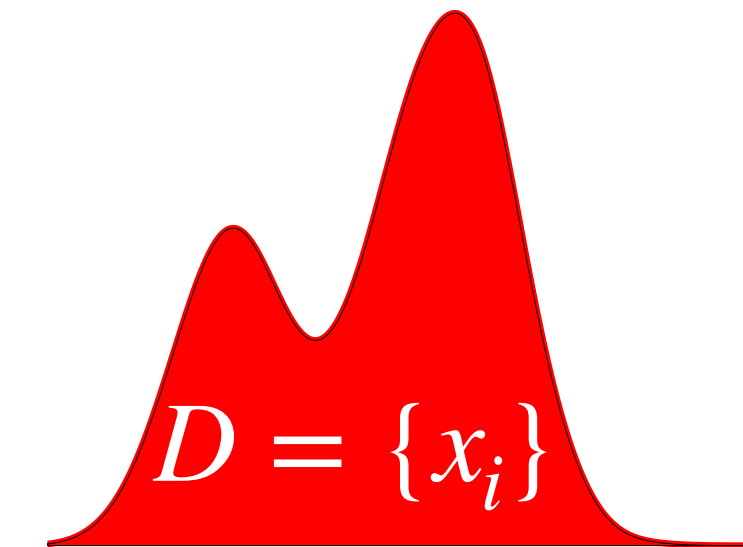
Projection to Wasserstein

=

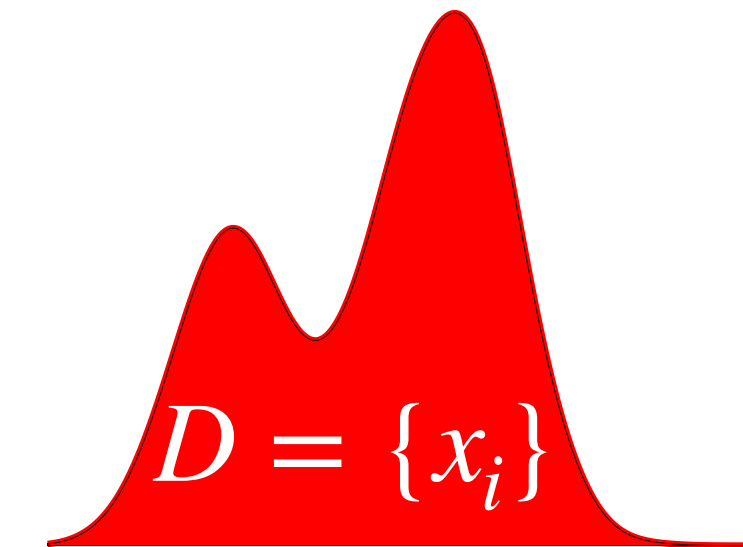
Quantile regression

(uniform quantile grid)

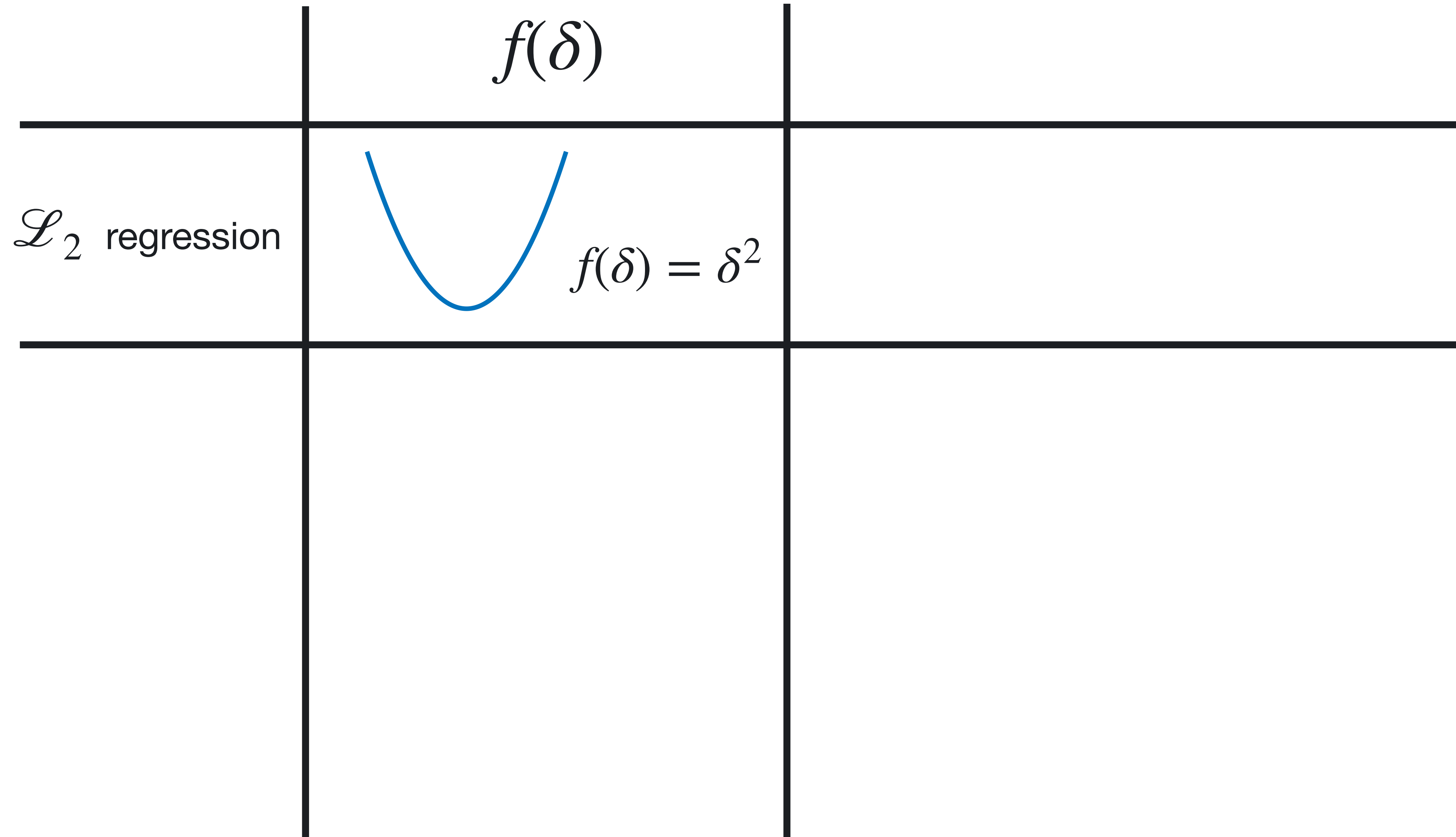
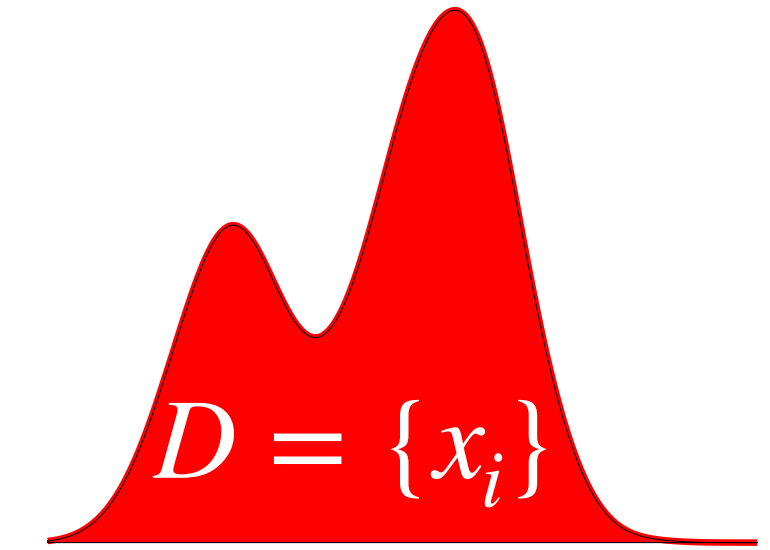
How can we **learn** from data ?



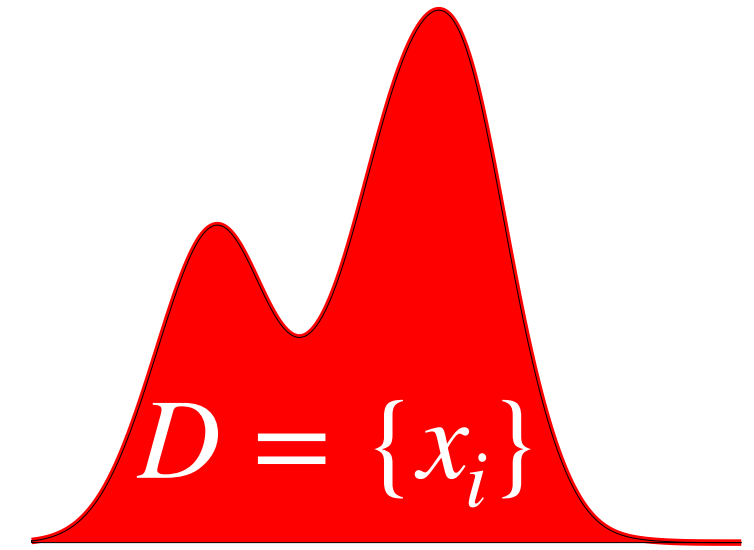
How can we **learn** from data ?

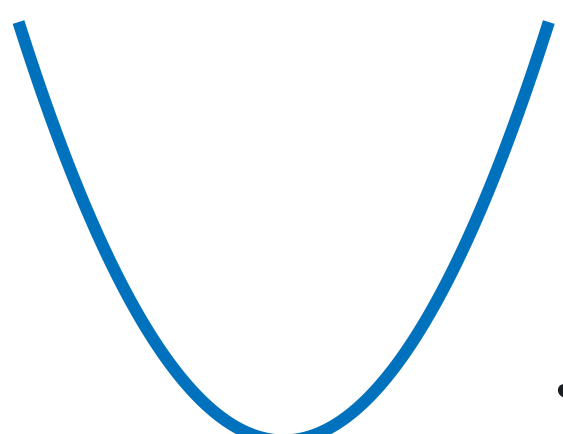


How can we **learn** from data ?



How can we learn from data ?

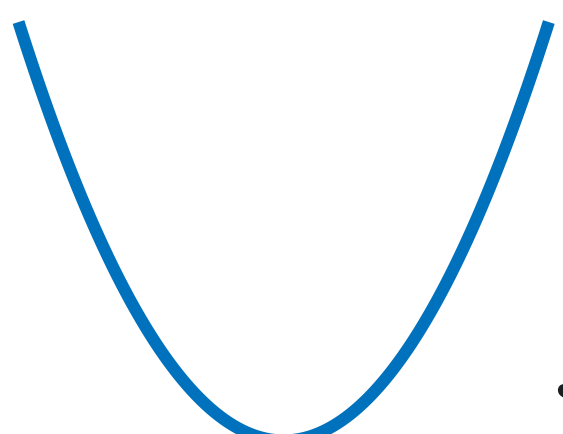


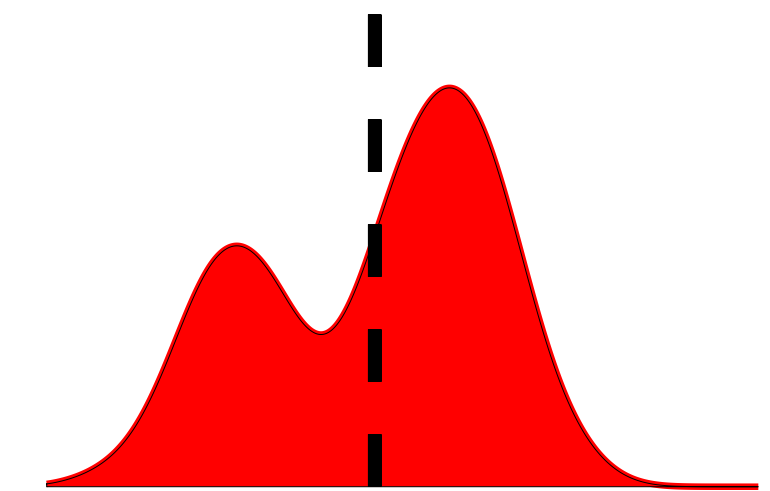
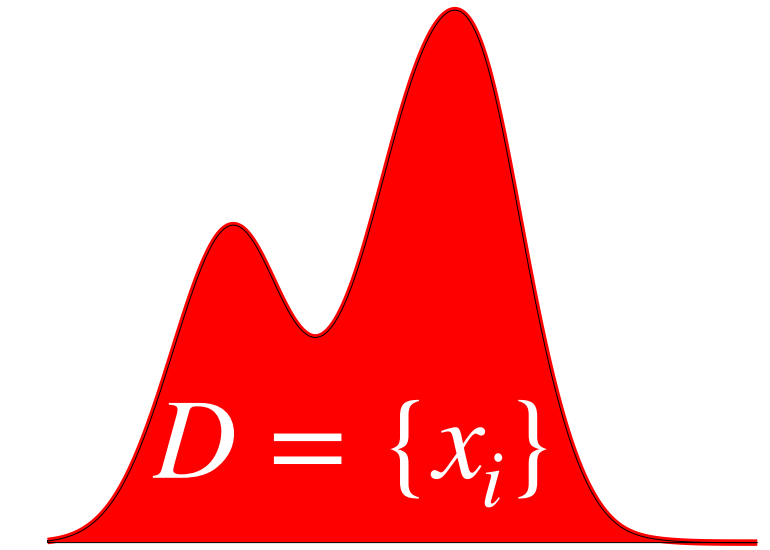
	$f(\delta)$	$x^* = \operatorname{argmin}_x \sum_i f(x_i - x)$
\mathcal{L}_2 regression	 $f(\delta) = \delta^2$	

$$\frac{\partial \sum_i f(x_i - x)}{\partial x} = \frac{\partial \sum_i (x_i - x)^2}{\partial x} = \sum_i 2(x_i - x)$$

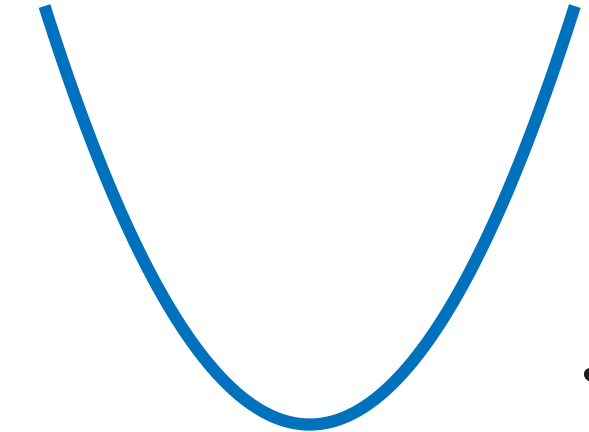
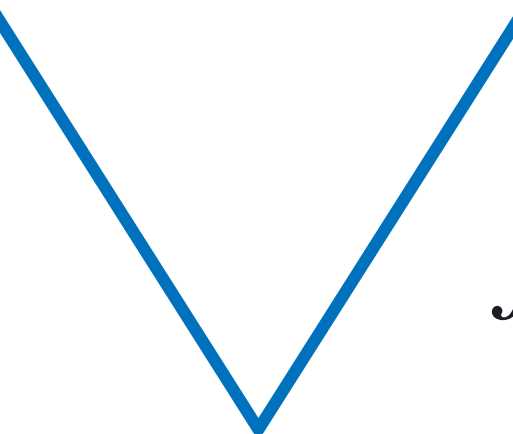
$$\sum_i 2(x_i - x) |_{x^*} \hat{=} 0 \implies x^* = \frac{1}{N} \sum_i x_i \quad \text{Mean}$$

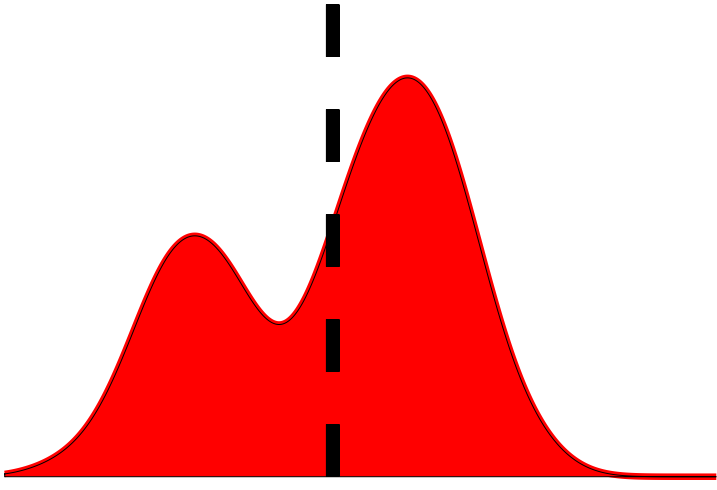
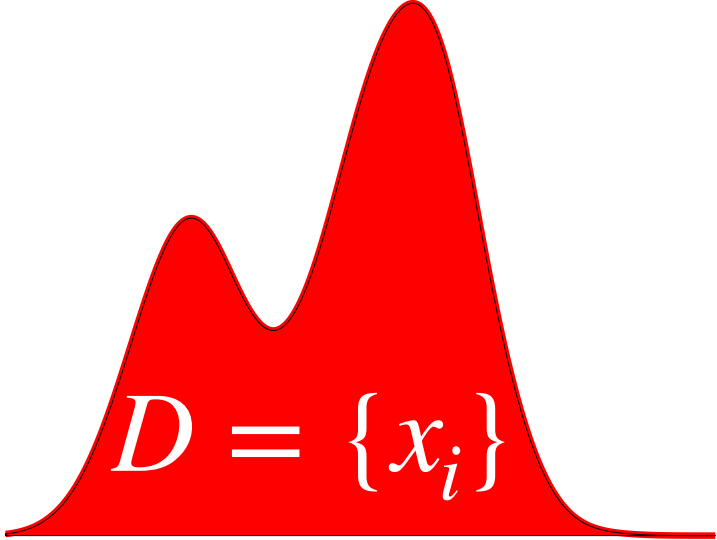
How can we learn from data ?

	$f(\delta)$	$x^\star = \operatorname{argmin}_x \sum_i f(x_i - x)$
\mathcal{L}_2 regression	 $f(\delta) = \delta^2$	$x^\star = \frac{1}{N} \sum_i x_i$ Mean



How can we learn from data ?

	$f(\delta)$	$x^* = \operatorname{argmin}_x \sum_i f(x_i - x)$
\mathcal{L}_2 regression	 $f(\delta) = \delta^2$	$x^* = \frac{1}{N} \sum_i x_i$ Mean
\mathcal{L}_1 regression	 $f(\delta) = \delta $	



$$\frac{\partial \sum_i f(x_i - x)}{\partial x} = \frac{\partial \sum_i |x_i - x|}{\partial x} =$$

$$\frac{\partial \sum_i f(x_i - x)}{\partial x} = \frac{\partial \sum_i |x_i - x|}{\partial x} = \sum_i \left(\mathbf{I}_{x_i \leq x} - \mathbf{I}_{x_i \geq x} \right) = \sum_i \mathbf{I}_{x_i \leq x} - \sum_i \mathbf{I}_{x_i \geq x}$$

$$\frac{\partial \sum_i f(x_i - x)}{\partial x} = \frac{\partial \sum_i |x_i - x|}{\partial x} = \sum_i \left(\mathbf{I}_{x_i \leq x} - \mathbf{I}_{x_i \geq x} \right) = \sum_i \mathbf{I}_{x_i \leq x} - \sum_i \mathbf{I}_{x_i \geq x}$$

$$\sum_i \mathbf{I}_{x_i \leq x} - \sum_i \mathbf{I}_{x_i \geq x} \Big|_{x^\star} \hat{=} 0$$

$$\frac{\partial \sum_i f(x_i - x)}{\partial x} = \frac{\partial \sum_i |x_i - x|}{\partial x} = \sum_i \left(\mathbf{I}_{x_i \leq x} - \mathbf{I}_{x_i \geq x} \right) = \sum_i \mathbf{I}_{x_i \leq x} - \sum_i \mathbf{I}_{x_i \geq x}$$

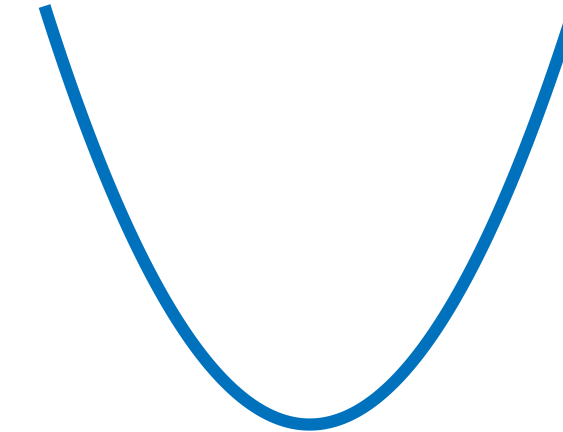

$$\sum_i \mathbf{I}_{x_i \leq x} - \sum_i \mathbf{I}_{x_i \geq x} \Big|_{x^*} \hat{=} 0 \implies \sum_i \mathbf{I}_{x_i \leq x^*} = \sum_i \mathbf{I}_{x_i \geq x^*} \implies$$

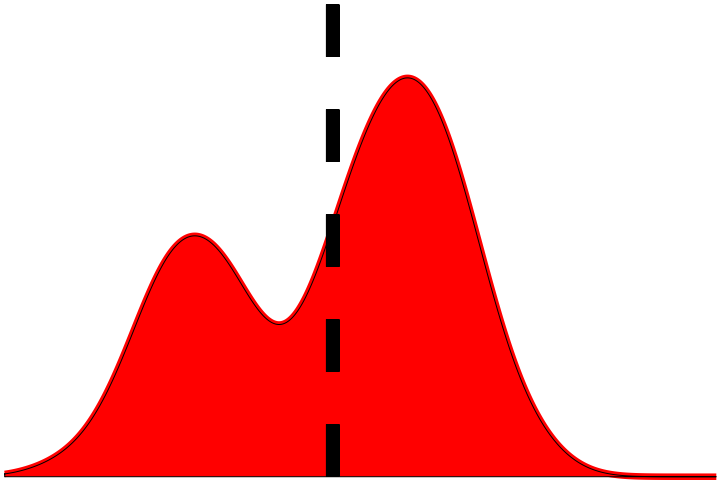
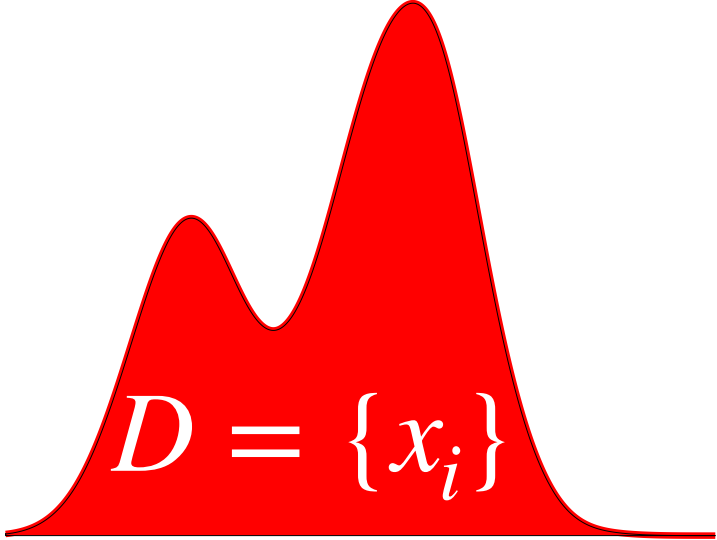
$$\frac{\partial \sum_i f(x_i - x)}{\partial x} = \frac{\partial \sum_i |x_i - x|}{\partial x} = \sum_i \left(\mathbf{I}_{x_i \leq x} - \mathbf{I}_{x_i \geq x} \right) = \sum_i \mathbf{I}_{x_i \leq x} - \sum_i \mathbf{I}_{x_i \geq x}$$

$$\sum_i \mathbf{I}_{x_i \leq x} - \sum_i \mathbf{I}_{x_i \geq x} \Big|_{x^*} \hat{=} 0 \implies \sum_i \mathbf{I}_{x_i \leq x^*} = \sum_i \mathbf{I}_{x_i \geq x^*} \implies$$

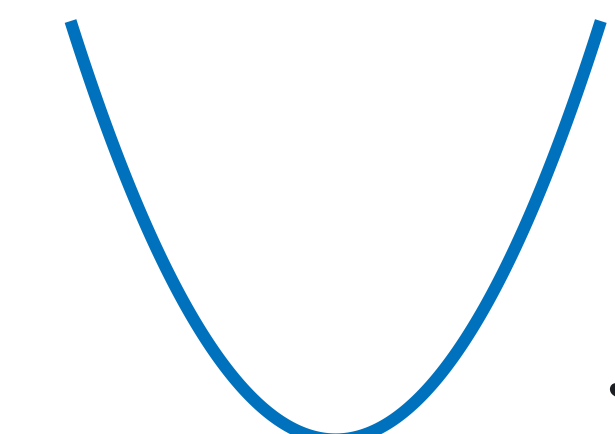
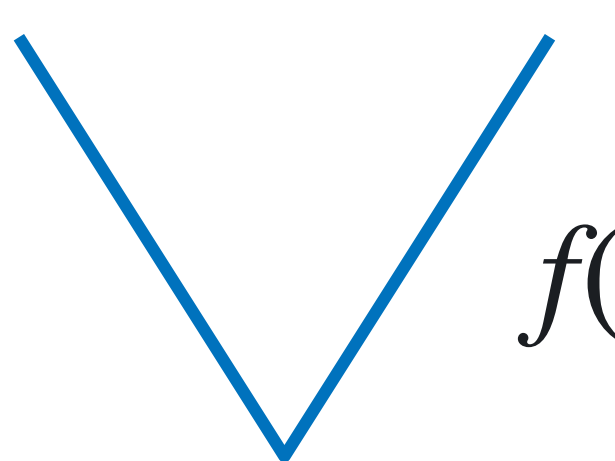
$$\implies \frac{\sum_i \mathbf{I}_{x_i \leq x^*}}{\sum_i \mathbf{I}_{x_i \geq x^*}} = 1 \quad \text{Median}$$

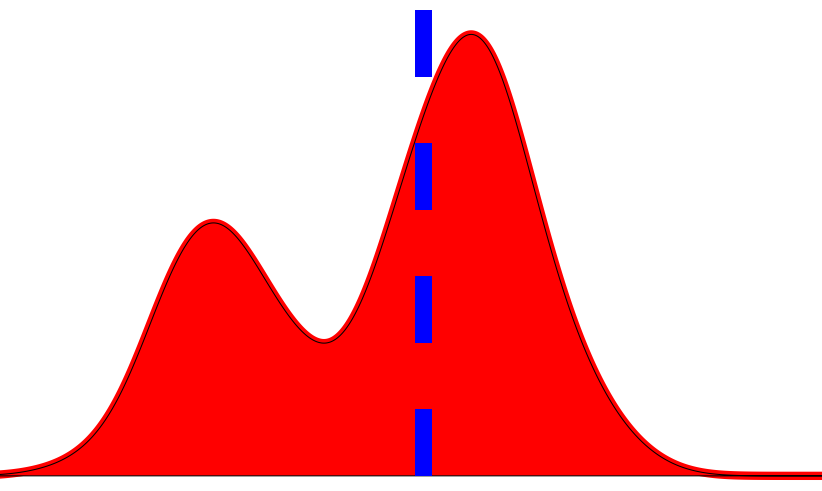
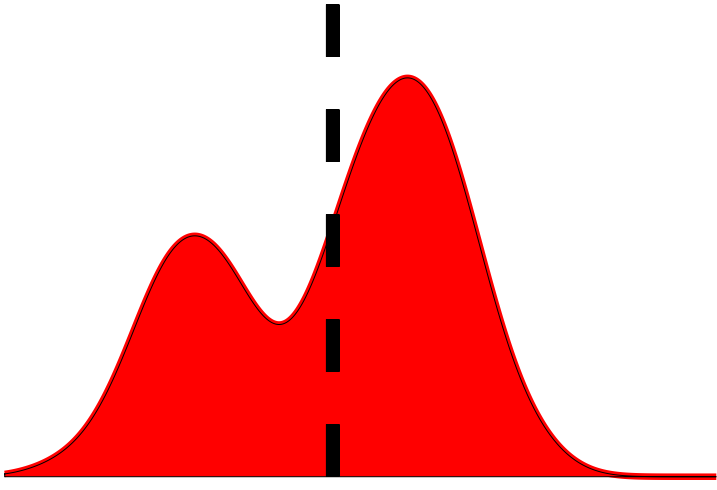
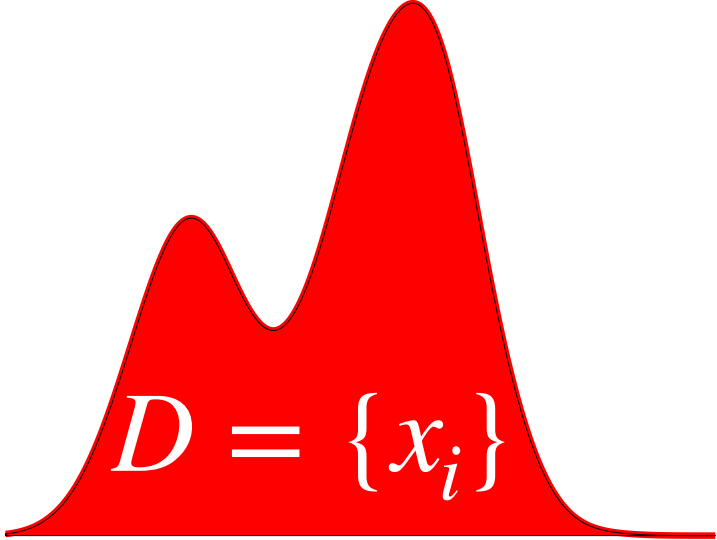
How can we learn from data ?

	$f(\delta)$	$x^* = \operatorname{argmin}_x \sum_i f(x_i - x)$
\mathcal{L}_2 regression	 $f(\delta) = \delta^2$	$x^* = \frac{1}{N} \sum_i x_i$ Mean
\mathcal{L}_1 regression	 $f(\delta) = \delta $	



How can we learn from data ?

	$f(\delta)$	$x^\star = \operatorname{argmin}_x \sum_i f(x_i - x)$	
\mathcal{L}_2 regression	 $f(\delta) = \delta^2$	$x^\star = \frac{1}{N} \sum_i x_i$	Mean
\mathcal{L}_1 regression	 $f(\delta) = \delta $	$x^\star = F_X^{-1}(0.5)$	Median $\tau = 0.5$



$$\frac{\partial \sum_i f(x_i - x)}{\partial x} = \sum_i \left(\mathbf{I}_{x_i \leq x} - \mathbf{I}_{x_i \geq x} \right) = \sum_i \mathbf{I}_{x_i \leq x} - \sum_i \mathbf{I}_{x_i \geq x}$$

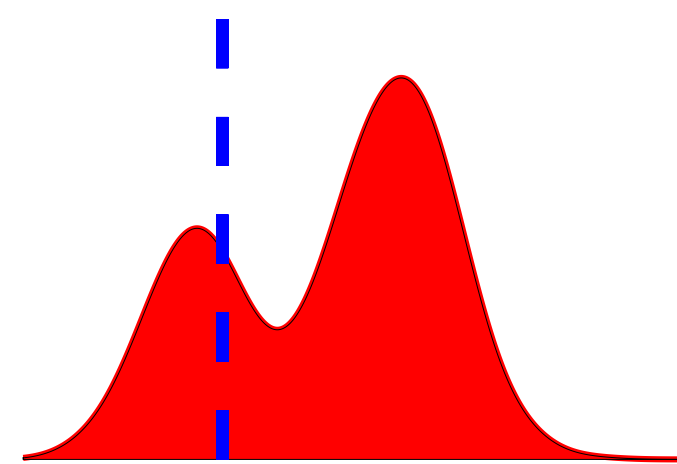
$$\sum_i \mathbf{I}_{x_i \leq x} - \sum_i \mathbf{I}_{x_i \geq x} \Big|_{x^\star} \hat{=} 0 \implies \sum_i \mathbf{I}_{x_i \leq x^\star} = \sum_i \mathbf{I}_{x_i \geq x^\star} \implies$$

$$\implies \frac{\sum_i \mathbf{I}_{x_i \leq x^\star}}{\sum_i \mathbf{I}_{x_i \geq x^\star}} = 1 \quad \text{Median}$$

$$\frac{\partial \sum_i f(x_i - x)}{\partial x} = \sum_i \left(\mathbf{I}_{x_i \leq x} - \mathbf{I}_{x_i \geq x} \right) = \sum_i \mathbf{I}_{x_i \leq x} - \sum_i \mathbf{I}_{x_i \geq x}$$

$$\sum_i \mathbf{I}_{x_i \leq x} - \sum_i \mathbf{I}_{x_i \geq x} \Big|_{x^*} \hat{=} 0 \implies \sum_i \mathbf{I}_{x_i \leq x^*} = \sum_i \mathbf{I}_{x_i \geq x^*} \implies$$

$$\implies \frac{\sum_i \mathbf{I}_{x_i \leq x^*}}{\sum_i \mathbf{I}_{x_i \geq x^*}} = 1$$

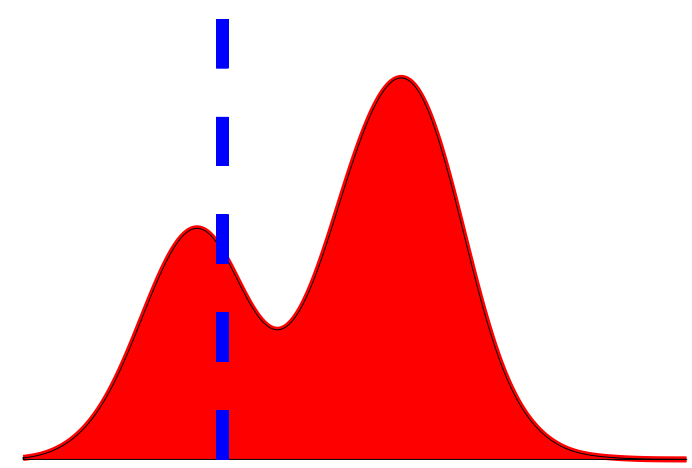


1/4 Quantile?

$$\frac{\partial \sum_i f(x_i - x)}{\partial x} = \sum_i \left(\mathbf{I}_{x_i \leq x} - \mathbf{I}_{x_i \geq x} \right) = \sum_i \mathbf{I}_{x_i \leq x} - \sum_i \mathbf{I}_{x_i \geq x}$$

$$\sum_i \mathbf{I}_{x_i \leq x} - \sum_i \mathbf{I}_{x_i \geq x} \Big|_{x^*} \hat{=} 0 \implies \sum_i \mathbf{I}_{x_i \leq x^*} = \sum_i \mathbf{I}_{x_i \geq x^*} \implies$$

$$\implies \frac{\sum_i \mathbf{I}_{x_i \leq x^*}}{\sum_i \mathbf{I}_{x_i \geq x^*}} = \frac{\frac{1}{4}}{\frac{3}{4}}$$

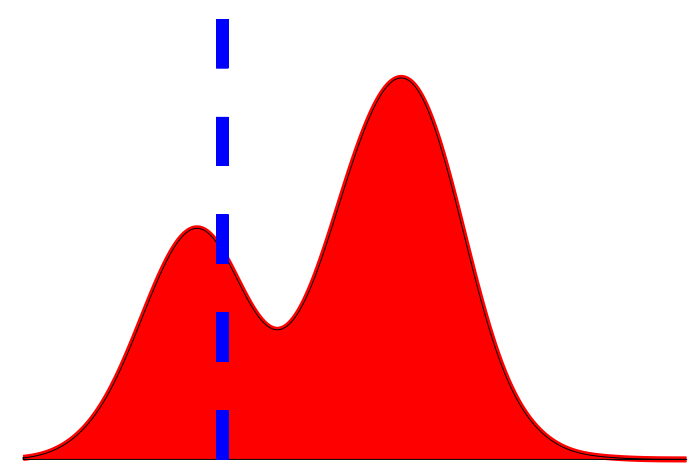


1/4 Quantile?

$$\frac{\partial \sum_i f(x_i - x)}{\partial x} = \sum_i \left(\mathbf{I}_{x_i \leq x} - \mathbf{I}_{x_i \geq x} \right) = \sum_i \mathbf{I}_{x_i \leq x} - \sum_i \mathbf{I}_{x_i \geq x}$$

$$\sum_i \mathbf{I}_{x_i \leq x} - \sum_i \mathbf{I}_{x_i \geq x} \Big|_{x^*} \hat{=} 0 \implies \sum_i \frac{3}{4} \mathbf{I}_{x_i \leq x^*} = \sum_i \frac{1}{4} \mathbf{I}_{x_i \geq x^*} \implies$$

$$\implies \frac{\sum_i \mathbf{I}_{x_i \leq x^*}}{\sum_i \mathbf{I}_{x_i \geq x^*}} = \frac{\frac{1}{4}}{\frac{3}{4}}$$

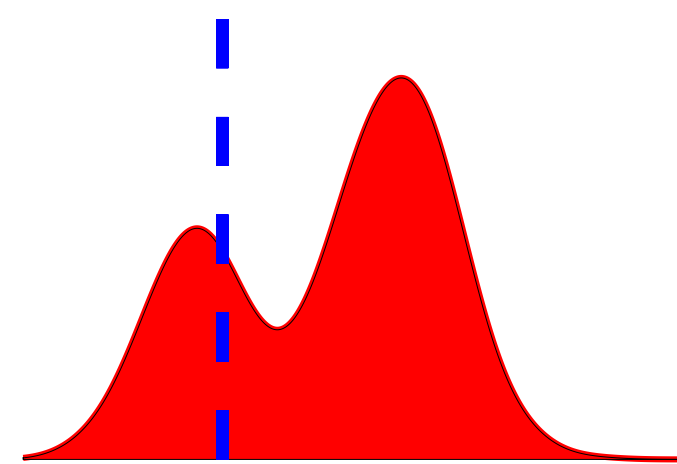


1/4 Quantile?

$$\frac{\partial \sum_i f(x_i - x)}{\partial x} = \sum_i \left(\mathbf{I}_{x_i \leq x} - \mathbf{I}_{x_i \geq x} \right) = \sum_i \mathbf{I}_{x_i \leq x} - \sum_i \mathbf{I}_{x_i \geq x}$$

$$\sum_i \frac{3}{4} \mathbf{I}_{x_i \leq x} - \sum_i \frac{1}{4} \mathbf{I}_{x_i \geq x} \Big|_{x^*} \hat{=} 0 \implies \sum_i \frac{3}{4} \mathbf{I}_{x_i \leq x^*} = \sum_i \frac{1}{4} \mathbf{I}_{x_i \geq x^*} \implies$$

$$\implies \frac{\sum_i \mathbf{I}_{x_i \leq x^*}}{\sum_i \mathbf{I}_{x_i \geq x^*}} = \frac{\frac{1}{4}}{\frac{3}{4}}$$

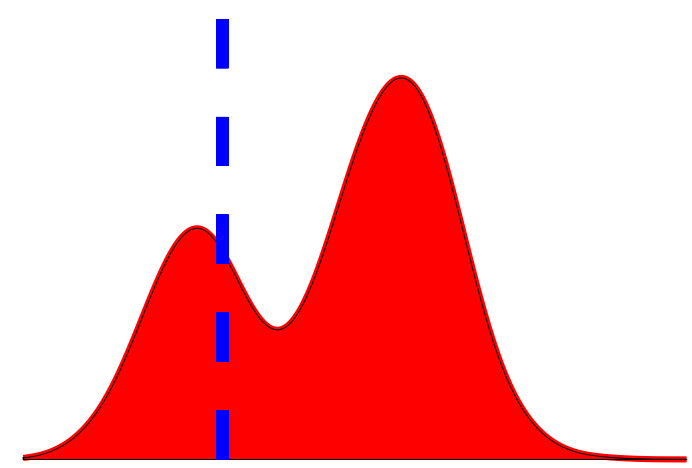


1/4 Quantile?

$$\frac{\partial \sum_i f(x_i - x)}{\partial x} = \sum_i \left(\frac{3}{4} \mathbf{I}_{x_i \leq x} - \frac{1}{4} \mathbf{I}_{x_i \geq x} \right) = \sum_i \frac{3}{4} \mathbf{I}_{x_i \leq x} - \sum_i \frac{1}{4} \mathbf{I}_{x_i \geq x}$$

$$\sum_i \frac{3}{4} \mathbf{I}_{x_i \leq x} - \sum_i \frac{1}{4} \mathbf{I}_{x_i \geq x} \Big|_{x^*} \hat{=} 0 \implies \sum_i \frac{3}{4} \mathbf{I}_{x_i \leq x^*} = \sum_i \frac{1}{4} \mathbf{I}_{x_i \geq x^*} \implies$$

$$\implies \frac{\sum_i \mathbf{I}_{x_i \leq x^*}}{\sum_i \mathbf{I}_{x_i \geq x^*}} = \frac{1/4}{3/4}$$

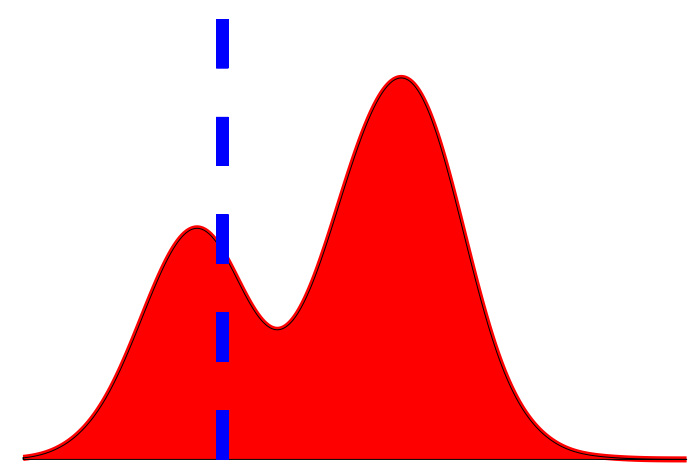


1/4 Quantile?

$$\frac{\partial \sum_i f(x_i - x)}{\partial x} = \sum_i \left(\frac{3}{4} \mathbf{I}_{x_i \leq x} - \frac{1}{4} \mathbf{I}_{x_i \geq x} \right) = \sum_i \frac{3}{4} \mathbf{I}_{x_i \leq x} - \sum_i \frac{1}{4} \mathbf{I}_{x_i \geq x}$$

$$\sum_i \frac{3}{4} \mathbf{I}_{x_i \leq x} - \sum_i \frac{1}{4} \mathbf{I}_{x_i \geq x} \Big|_{x^*} \hat{=} 0 \implies \sum_i \frac{3}{4} \mathbf{I}_{x_i \leq x^*} = \sum_i \frac{1}{4} \mathbf{I}_{x_i \geq x^*} \implies$$

$$\implies \frac{\sum_i \mathbf{I}_{x_i \leq x^*}}{\sum_i \mathbf{I}_{x_i \geq x^*}} = \frac{1/4}{3/4}$$



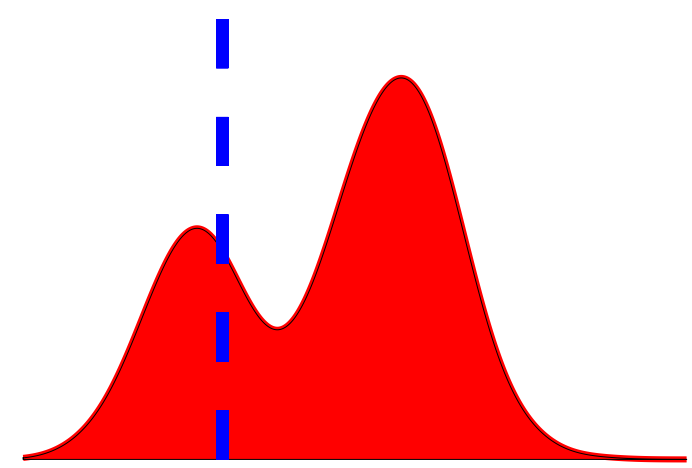
1/4 Quantile?

$$f(x_i - x) = \begin{cases} 1/4 (x_i - x), & x_i \geq x \\ -3/4 (x_i - x), & x_i < x \end{cases}$$

$$\frac{\partial \sum_i f(x_i - x)}{\partial x} = \sum_i \left(\frac{3}{4} \mathbf{I}_{x_i \leq x} - \frac{1}{4} \mathbf{I}_{x_i \geq x} \right) = \sum_i \frac{3}{4} \mathbf{I}_{x_i \leq x} - \sum_i \frac{1}{4} \mathbf{I}_{x_i \geq x}$$

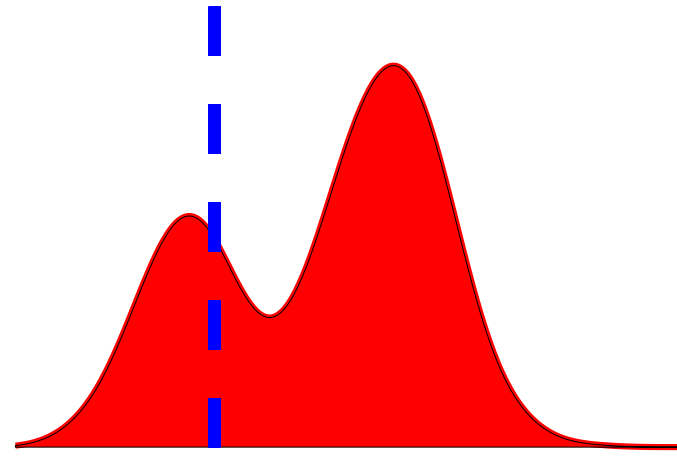
$$\sum_i \frac{3}{4} \mathbf{I}_{x_i \leq x} - \sum_i \frac{1}{4} \mathbf{I}_{x_i \geq x} \Big|_{x^*} \hat{=} 0 \implies \sum_i \frac{3}{4} \mathbf{I}_{x_i \leq x^*} = \sum_i \frac{1}{4} \mathbf{I}_{x_i \geq x^*} \implies$$

$$\implies \frac{\sum_i \mathbf{I}_{x_i \leq x^*}}{\sum_i \mathbf{I}_{x_i \geq x^*}} = \frac{1/4}{3/4}$$



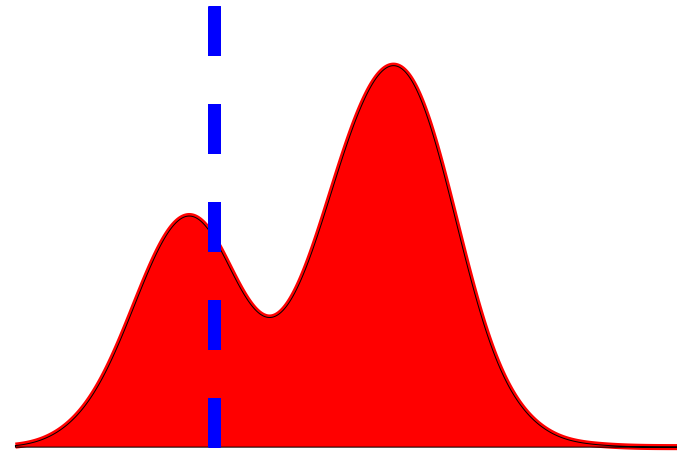
1/4 Quantile?

$$f(\delta) = \begin{cases} 1/4 \delta, & \delta \geq 0 \\ -3/4 \delta, & \delta < 0 \end{cases}$$



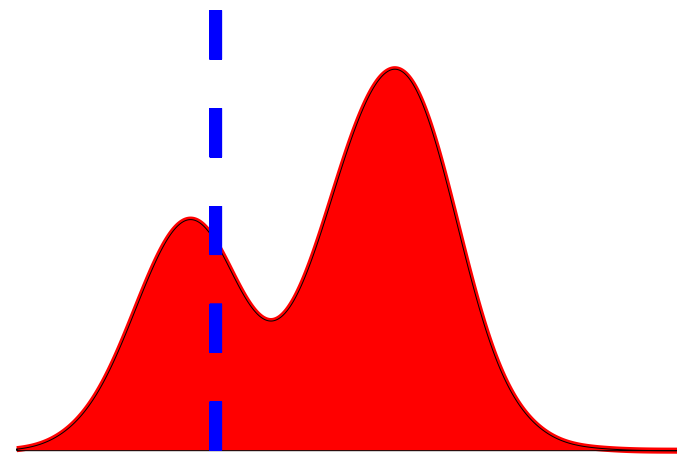
1/4 Quantile?

$$f(\delta) = \begin{cases} 1/4 \delta, & \delta \geq 0 \\ -3/4 \delta, & \delta < 0 \end{cases}$$



1/4 Quantile?

$$f(\delta) = \begin{cases} 1/4 \delta, & \delta \geq 0 \\ -3/4 \delta, & \delta < 0 \end{cases}$$

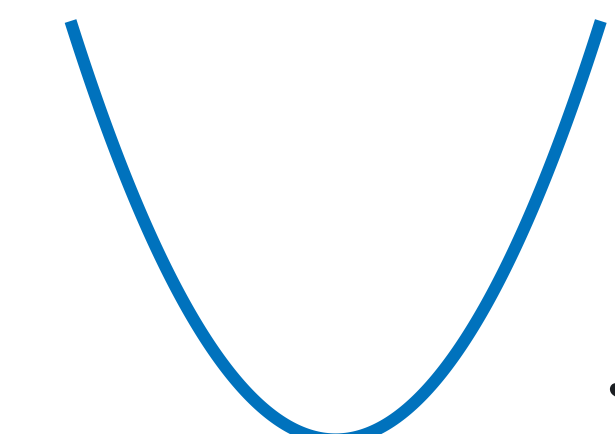
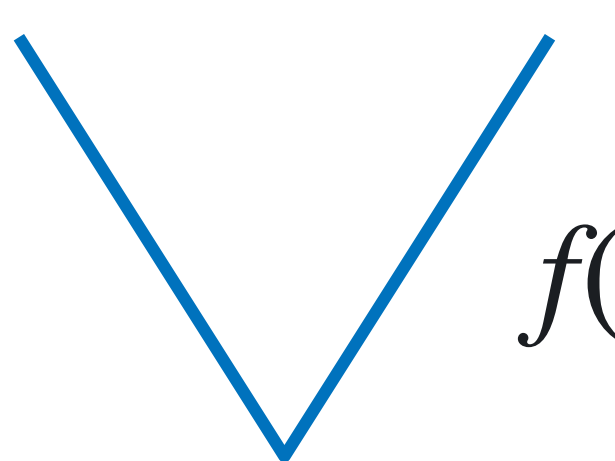


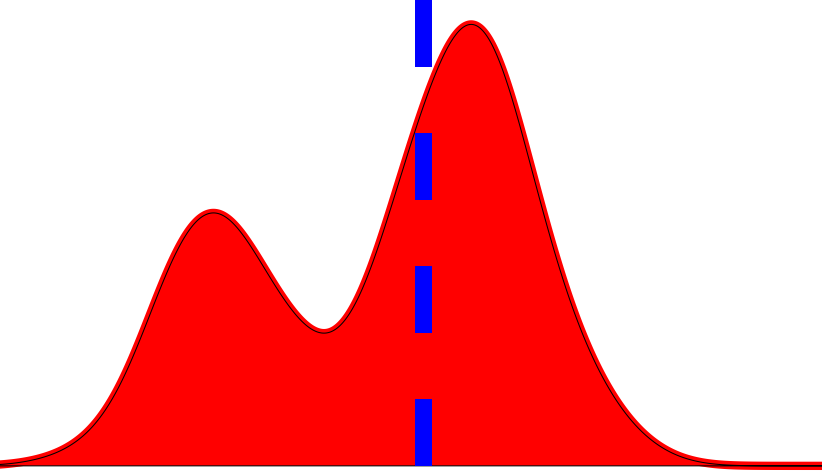
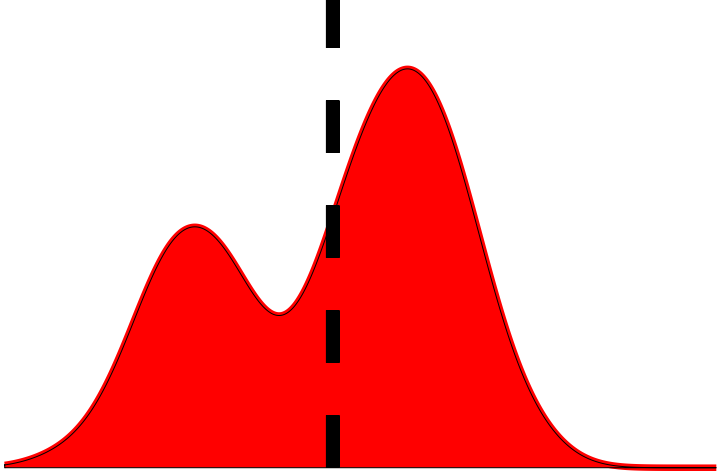
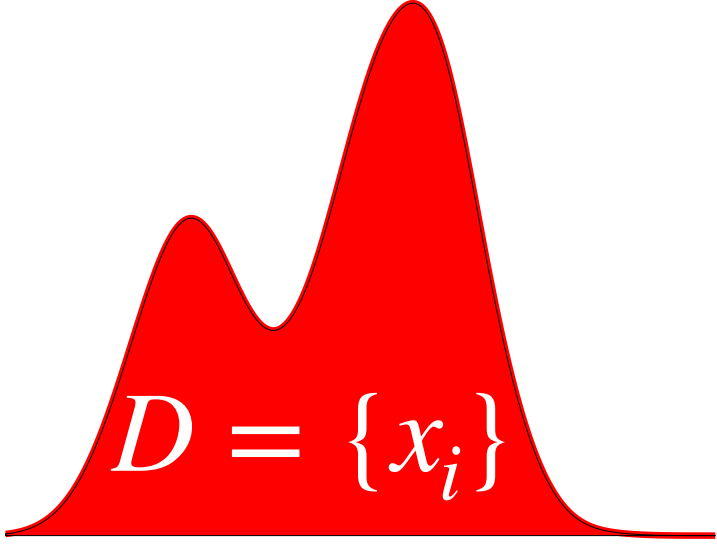
τ -Quantile?

$$f_{\tau}(\delta) = \begin{cases} \tau \delta, & \delta \geq 0 \\ (\tau - 1) \delta, & \delta < 0 \end{cases}$$

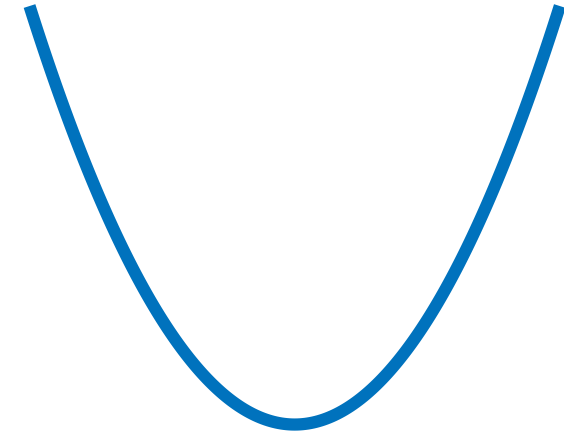

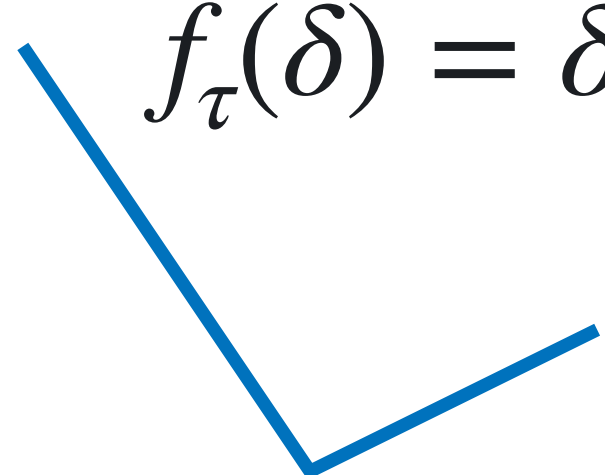
$$f_{\tau}(\delta) = \delta (\tau - \mathbf{I}_{\delta < 0})$$

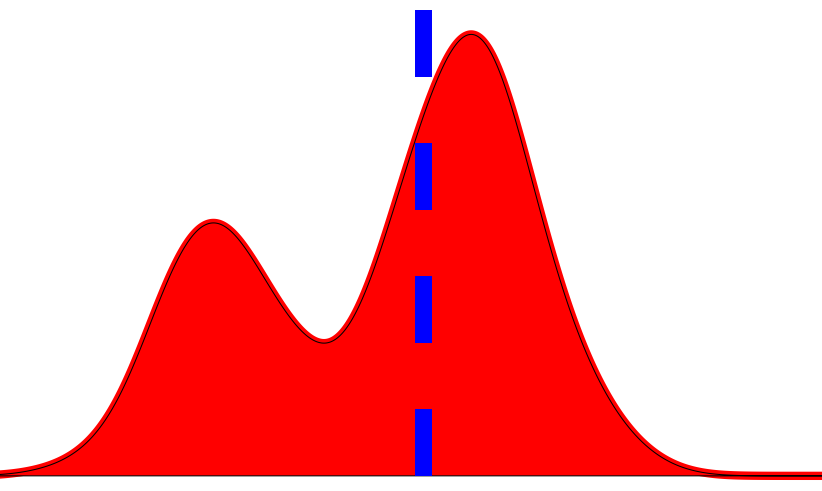
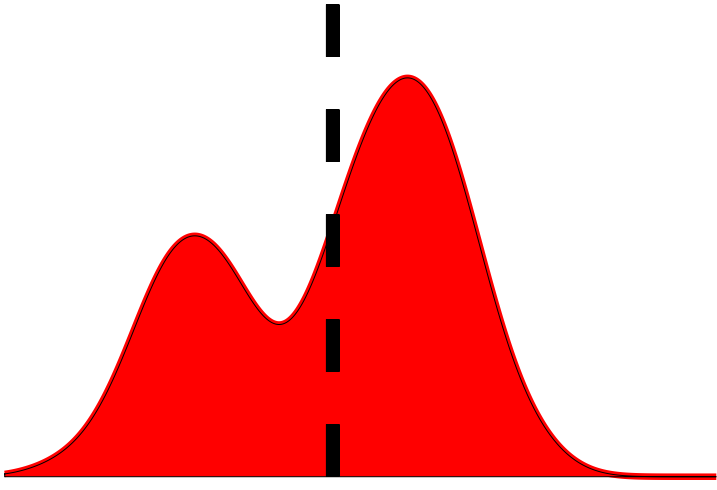
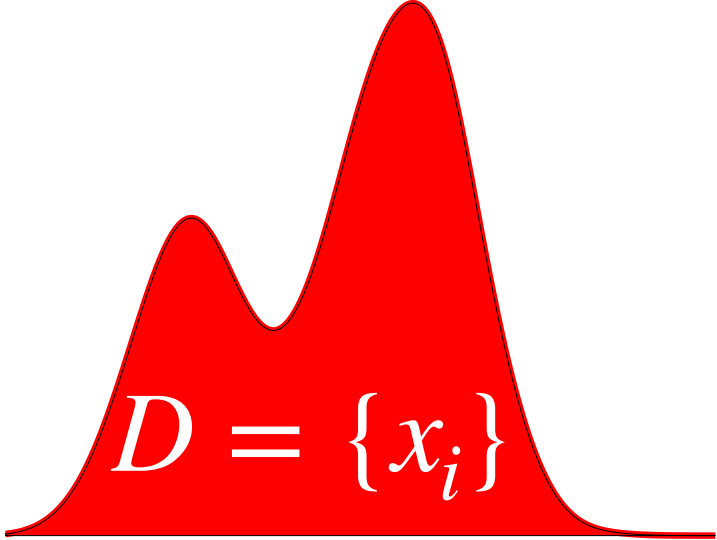
How can we learn from data ?

	$f(\delta)$	$x^* = \operatorname{argmin}_x \sum_i f(x_i - x)$	
\mathcal{L}_2 regression	 $f(\delta) = \delta^2$	$x^* = \frac{1}{N} \sum_i x_i$	Mean
\mathcal{L}_1 regression	 $f(\delta) = \delta $	$x^* = F_X^{-1}(0.5)$	Median $\tau = 0.5$

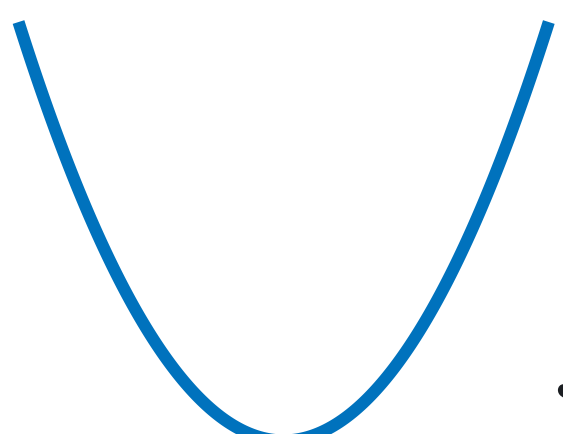
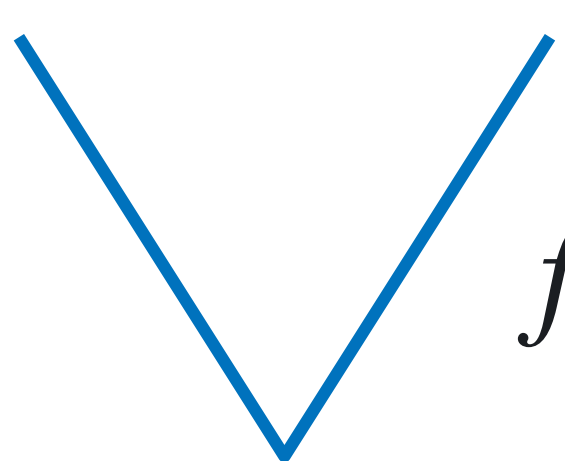
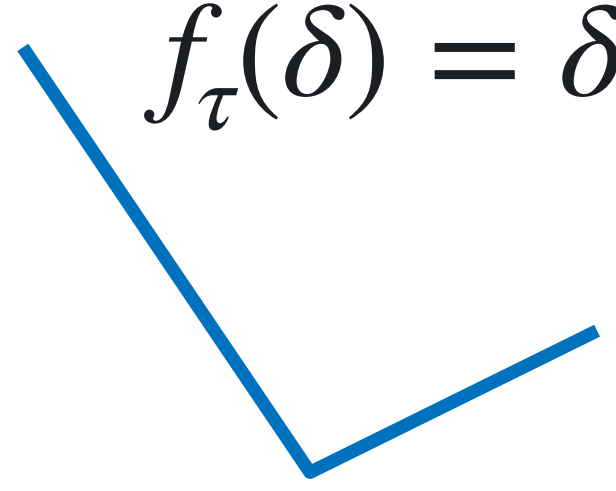


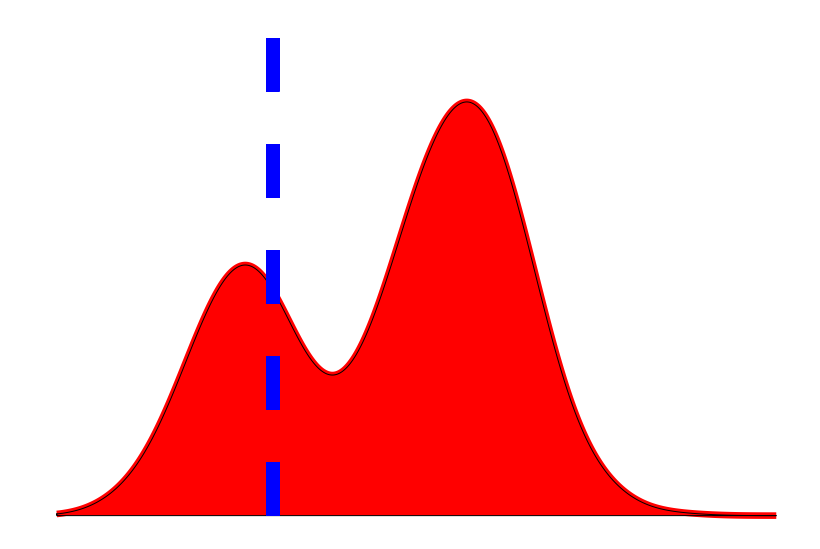
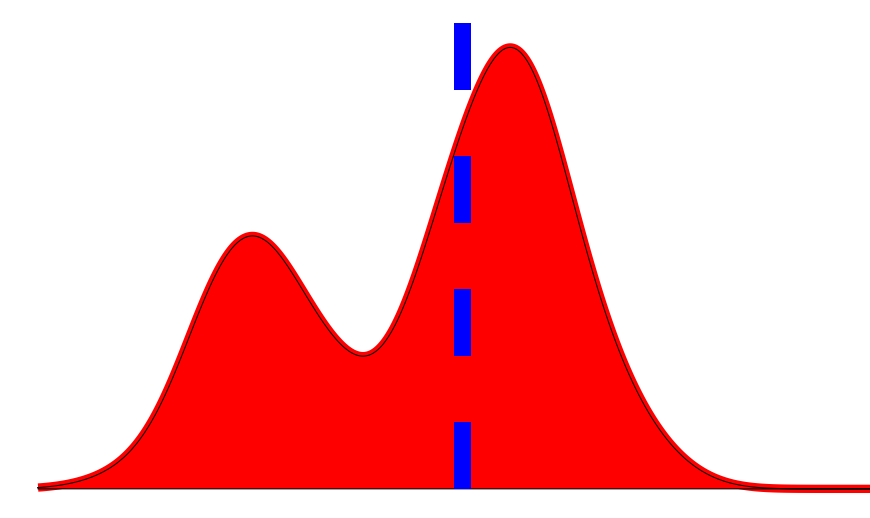
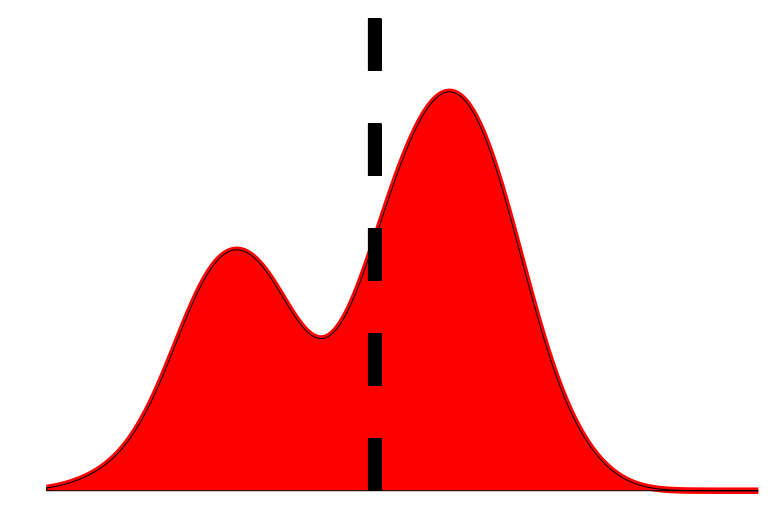
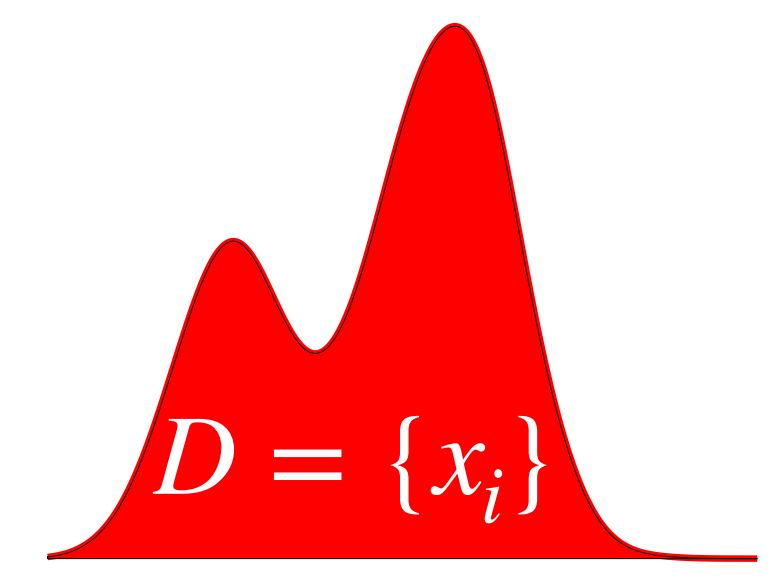
How can we learn from data ?

	$f(\delta)$	$x^* = \operatorname{argmin}_x \sum_i f(x_i - x)$	
\mathcal{L}_2 regression	 $f(\delta) = \delta^2$	$x^* = \frac{1}{N} \sum_i x_i$	Mean
\mathcal{L}_1 regression	 $f(\delta) = \delta $	$x^* = F_X^{-1}(0.5)$	Median $\tau = 0.5$
	 $f_\tau(\delta) = \delta (\tau - \mathbf{I}_{\delta < 0})$		

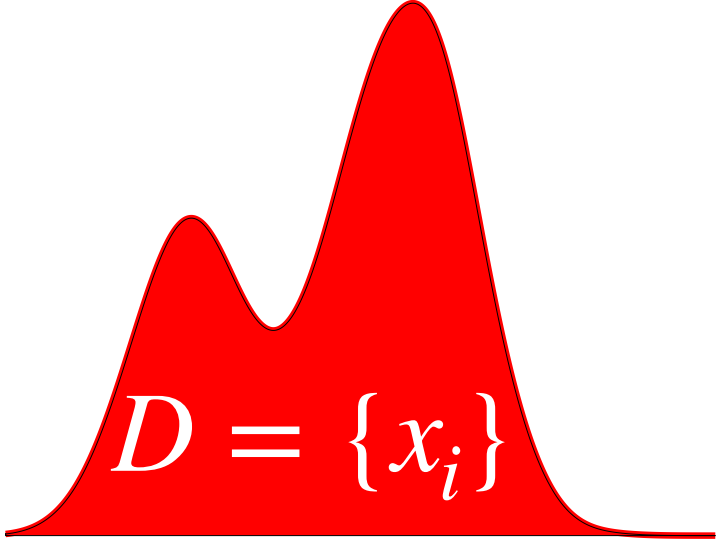


How can we learn from data ?

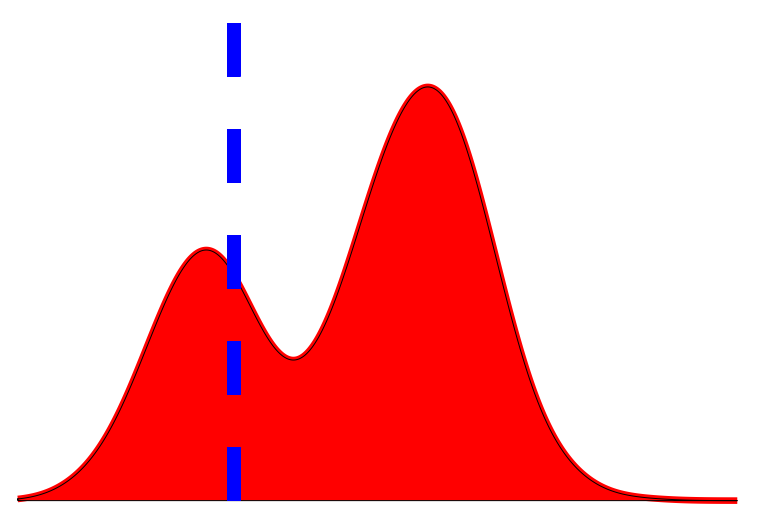
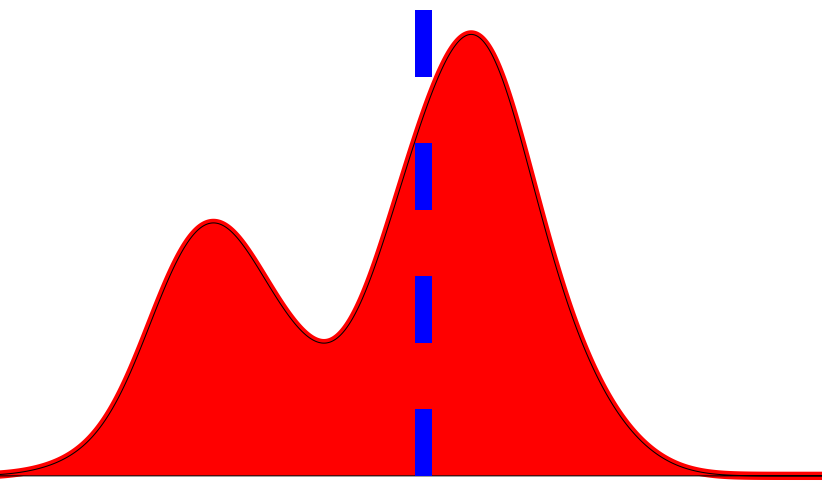
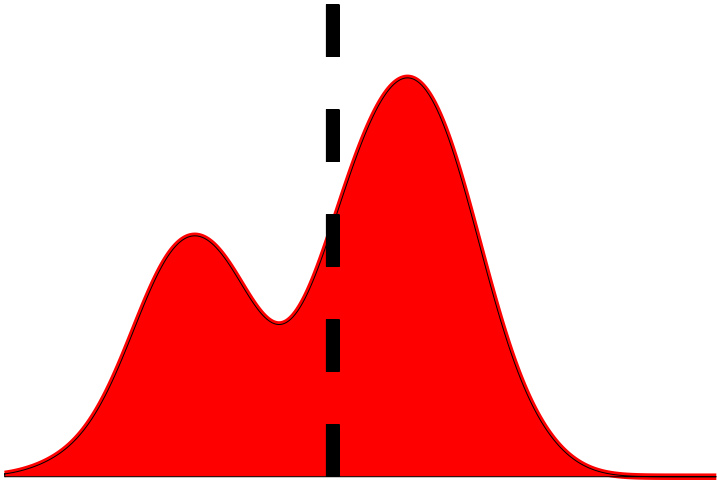
	$f(\delta)$	$x^* = \operatorname{argmin}_x \sum_i f(x_i - x)$	
\mathcal{L}_2 regression	 $f(\delta) = \delta^2$	$x^* = \frac{1}{N} \sum_i x_i$	Mean
\mathcal{L}_1 regression	 $f(\delta) = \delta $	$x^* = F_X^{-1}(0.5)$	Median $\tau = 0.5$
	 $f_\tau(\delta) = \delta (\tau - \mathbf{I}_{\delta < 0})$	$x^* = F_X^{-1}(\tau)$	τ -Quantile



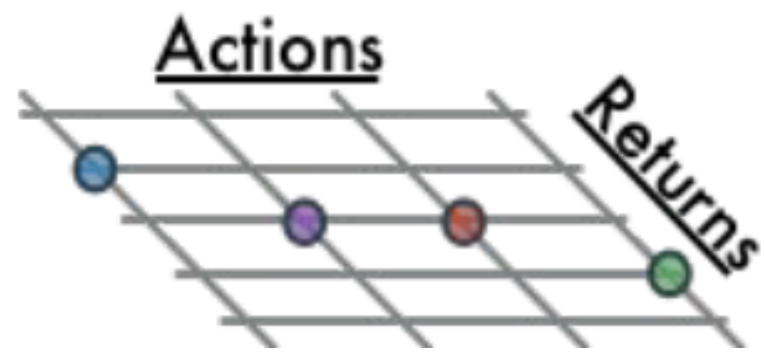
How can we learn from data ?



	$f(\delta)$	$x^* = \operatorname{argmin}_x \sum_i f(x_i - x)$	
\mathcal{L}_2 regression	$f(\delta) = \delta^2$	$x^* = \frac{1}{N} \sum_i x_i$	Mean
\mathcal{L}_1 regression	$f(\delta) = \delta $	$x^* = F_X^{-1}(0.5)$	Median $\tau = 0.5$
Quantile regression	$f_\tau(\delta) = \delta (\tau - \mathbf{I}_{\delta < 0})$	$x^* = F_X^{-1}(\tau)$	τ -Quantile



DQN



$$\{s_t, a_t, s_{t+1}, r_t\}$$

$$a^* = \operatorname{argmax}_a Q^\theta(s_{t+1}, a)$$

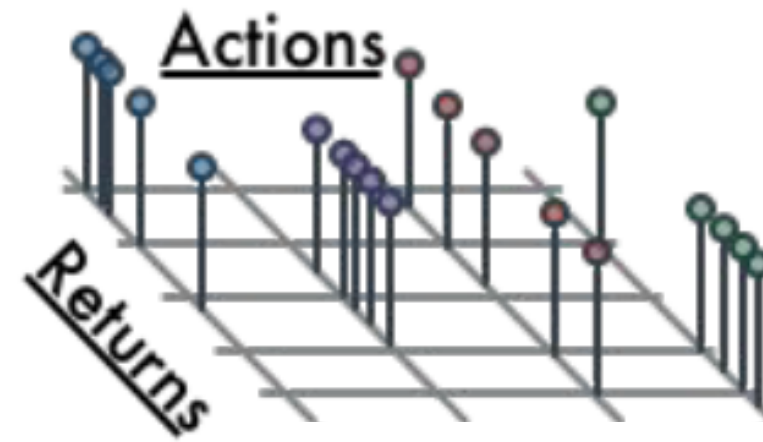
$$q' = Q^\theta(s_{t+1}, a^*)$$

$$q = Q^\theta(s_t, a_t)$$

$$\delta_t = r_t + \gamma q' - q$$

$$\mathcal{L}_{DQN} = \delta_t^2$$

QR-DQN



$$\{s_t, a_t, s_{t+1}, r_t\}$$

$$a^* = \operatorname{argmax}_a \mathbb{E} [Z_\tau^\theta(s_{t+1}, a)]$$

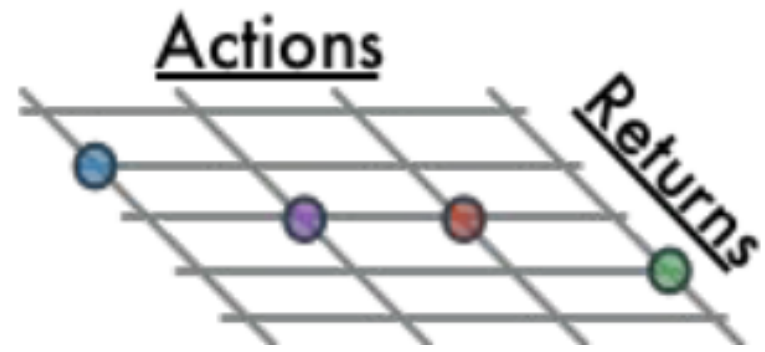
$$\forall \tau, \tau' \left| \begin{array}{l} z' = Z_{\tau'}^\theta(s_{t+1}, a^*) \\ z = Z_\tau^\theta(s_t, a_t) \end{array} \right.$$

$$z = Z_\tau^\theta(s_t, a_t)$$

$$\delta_t^{\tau, \tau'} = r_t + \gamma z' - z$$

$$\mathcal{L}_{QR-DQN} = ?$$

DQN



$$\{s_t, a_t, s_{t+1}, r_t\}$$

$$a^* = \operatorname{argmax}_a Q^\theta(s_{t+1}, a)$$

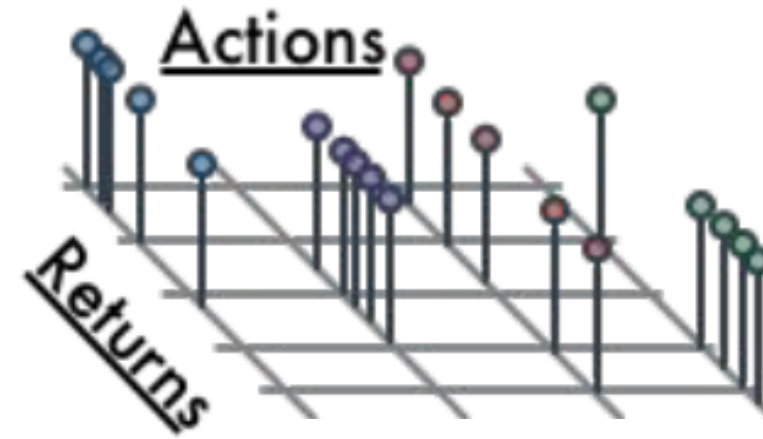
$$q' = Q^\theta(s_{t+1}, a^*)$$

$$q = Q^\theta(s_t, a_t)$$

$$\delta_t = r_t + \gamma q' - q$$

$$\mathcal{L}_{DQN} = \delta_t^2$$

QR-DQN



$$\{s_t, a_t, s_{t+1}, r_t\}$$

$$a^* = \operatorname{argmax}_a \mathbb{E} [Z_\tau^\theta(s_{t+1}, a)]$$

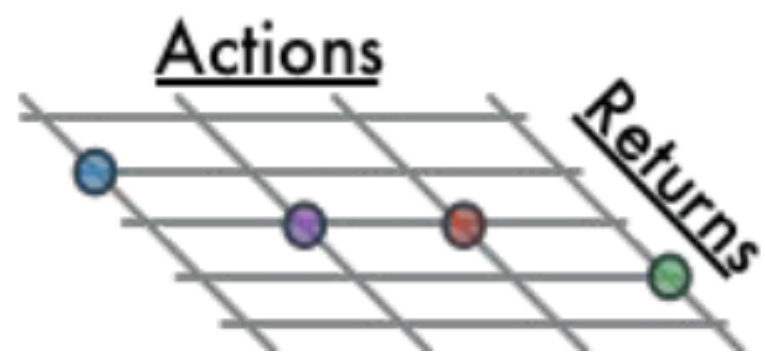
$$\forall \tau, \tau' \left| \begin{array}{l} z' = Z_{\tau'}^\theta(s_{t+1}, a^*) \\ z = Z_\tau^\theta(s_t, a_t) \end{array} \right.$$

$$z = Z_\tau^\theta(s_t, a_t)$$

$$\delta_t^{\tau, \tau'} = r_t + \gamma z' - z$$

$$\mathcal{L}_{QR-DQN} = \sum_{\tau} \sum_{\tau'} \delta_t^{\tau, \tau'} (\tau - \mathbf{I}_{\delta_t^{\tau, \tau'} < 0})$$

DQN



$$\{s_t, a_t, s_{t+1}, r_t\}$$

$$a^* = \operatorname{argmax}_a Q^\theta(s_{t+1}, a)$$

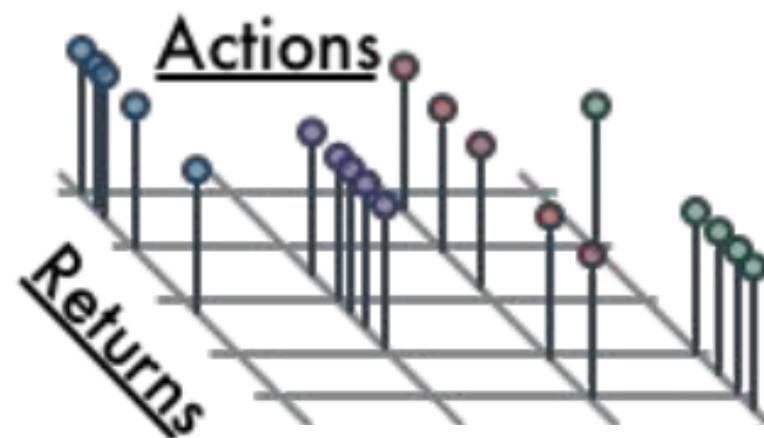
$$q' = Q^\theta(s_{t+1}, a^*)$$

$$q = Q^\theta(s_t, a_t)$$

$$\delta_t = r_t + \gamma q' - q$$

$$\mathcal{L}_{DQN} = \delta_t^2$$

QR-DQN



$$\{s_t, a_t, s_{t+1}, r_t\}$$

$$a^* = \operatorname{argmax}_a \mathbb{E} [Z_\tau^\theta(s_{t+1}, a)]$$

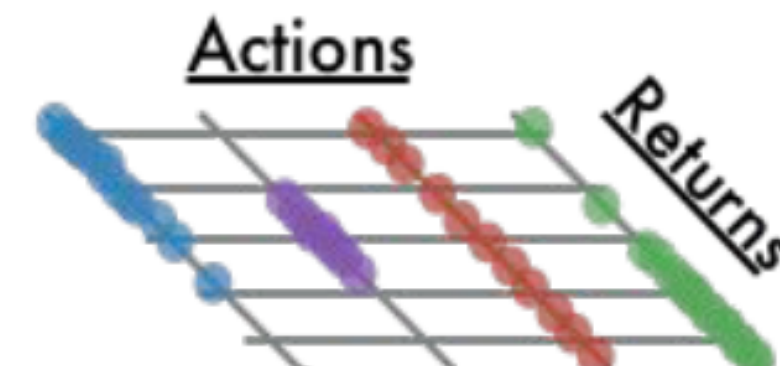
$$\forall \tau, \tau' \left| \begin{array}{l} z' = Z_{\tau'}^\theta(s_{t+1}, a^*) \\ z = Z_\tau^\theta(s_t, a_t) \end{array} \right.$$

$$z = Z_\tau^\theta(s_t, a_t)$$

$$\delta_t^{\tau, \tau'} = r_t + \gamma z' - z$$

$$\mathcal{L}_{QR-DQN} = \sum_{\tau} \sum_{\tau'} \delta_t^{\tau, \tau'} (\tau - \mathbf{I}_{\delta_t^{\tau, \tau'} < 0})$$

IQN



$$\{s_t, a_t, s_{t+1}, r_t\}$$

$$a^* = \operatorname{argmax}_a \mathbb{E} [Z_\tau^\theta(s_{t+1}, a)]$$

$$\forall \tau \sim U([0,1]) \left| \begin{array}{l} z' = Z_{\tau'}^\theta(s_{t+1}, a^*) \\ z = Z_\tau^\theta(s_t, a_t) \end{array} \right.$$

$$z = Z_\tau^\theta(s_t, a_t)$$

$$\delta_t^{\tau, \tau'} = r_t + \gamma z' - z$$

$$\mathcal{L}_{IQN} = \sum_{\tau} \sum_{\tau'} \delta_t^{\tau, \tau'} (\tau - \mathbf{I}_{\delta_t^{\tau, \tau'} < 0})$$

Human normalised score (HNS)

$$\text{score} = \frac{\text{agent} - \text{random}}{\text{human} - \text{random}}$$

Human gap

$$\text{gap} = 1 - \text{clip}(\text{score}, 0, 1) = \begin{cases} 1, & \text{random play} \\ 0, & \text{super-human} \end{cases}$$

RESULTS - HNS on Atari-57

	Mean	Median	(human starts) Median	Human Gap
DQN	228 %	79 %	68 %	0.334
Prio. Duel.	592 %	172 %	128 %	0.178
C51	701 %	178 %	125 %	0.152

RESULTS - HNS on Atari-57

smaller is better!



	Mean	Median	(human starts) Median	Human Gap
DQN	228 %	79 %	68 %	0.334
Prio. Duel.	592 %	172 %	128 %	0.178
C51	701 %	178 %	125 %	0.152

RESULTS - HNS on Atari-57

smaller is better!



	Mean	Median	(human starts) Median	Human Gap
DQN	228 %	79 %	68 %	0.334
Prio. Duel.	592 %	172 %	128 %	0.178
C51	701 %	178 %	125 %	0.152
QR-DQN	864 %	193 %	153 %	0.165

RESULTS - HNS on Atari-57

smaller is better!



	Mean	Median	(human starts) Median	Human Gap
DQN	228 %	79 %	68 %	0.334
Prio. Duel.	592 %	172 %	128 %	0.178
C51	701 %	178 %	125 %	0.152
QR-DQN	864 %	193 %	153 %	0.165
IQN	1019 %	218 %	162 %	0.141

RESULTS - HNS on Atari-57

smaller is better!



	Mean	Median	(human starts) Median	Human Gap
DQN	228 %	79 %	68 %	0.334
Prio. Duel.	592 %	172 %	128 %	0.178
C51	701 %	178 %	125 %	0.152
QR-DQN	864 %	193 %	153 %	0.165
IQN	1019 %	218 %	162 %	0.141
Rainbow	1189 %	230 %	125 %	0.144

QR-DQN
based

RESULTS - HNS on Atari-57

smaller is better!



	Mean	Median	(human starts) Median	Human Gap
DQN	228 %	79 %	68 %	0.334
Prio. Duel.	592 %	172 %	128 %	0.178
C51	701 %	178 %	125 %	0.152
QR-DQN	864 %	193 %	153 %	0.165
IQN	1019 %	218 %	162 %	0.141
Rainbow	1189 %	230 %	125 %	0.144

QR-DQN
based

RESULTS - HNS on Atari-57

smaller is better!



	Mean	Median	(human starts) Median	Human Gap
DQN	228 %	79 %	68 %	0.334
Prio. Duel.	592 %	172 %	128 %	0.178
C51	701 %	178 %	125 %	0.152
QR-DQN	864 %	193 %	153 %	0.165
IQN	1019 %	218 %	162 %	0.141
Rainbow	1189 %	230 %	125 %	0.144

QR-DQN
based

RESULTS - HNS on Atari-57

smaller is better!



	Mean	Median	(human starts) Median	Human Gap
DQN	228 %	79 %	68 %	0.334
Prio. Duel.	592 %	172 %	128 %	0.178
C51	701 %	178 %	125 %	0.152
QR-DQN	864 %	193 %	153 %	0.165
IQN	1019 %	218 %	162 %	0.141
Rainbow	1189 %	230 %	125 %	0.144

QR-DQN
based

RESULTS - HNS on Atari-57

smaller is better!



	Mean	Median	(human starts) Median	Human Gap
DQN	228 %	79 %	68 %	0.334
Prio. Duel.	592 %	172 %	128 %	0.178
C51	701 %	178 %	125 %	0.152
QR-DQN	864 %	193 %	153 %	0.165
IQN	1019 %	218 %	162 %	0.141
QR-DQN based Rainbow	1189 %	230 %	125 %	0.144

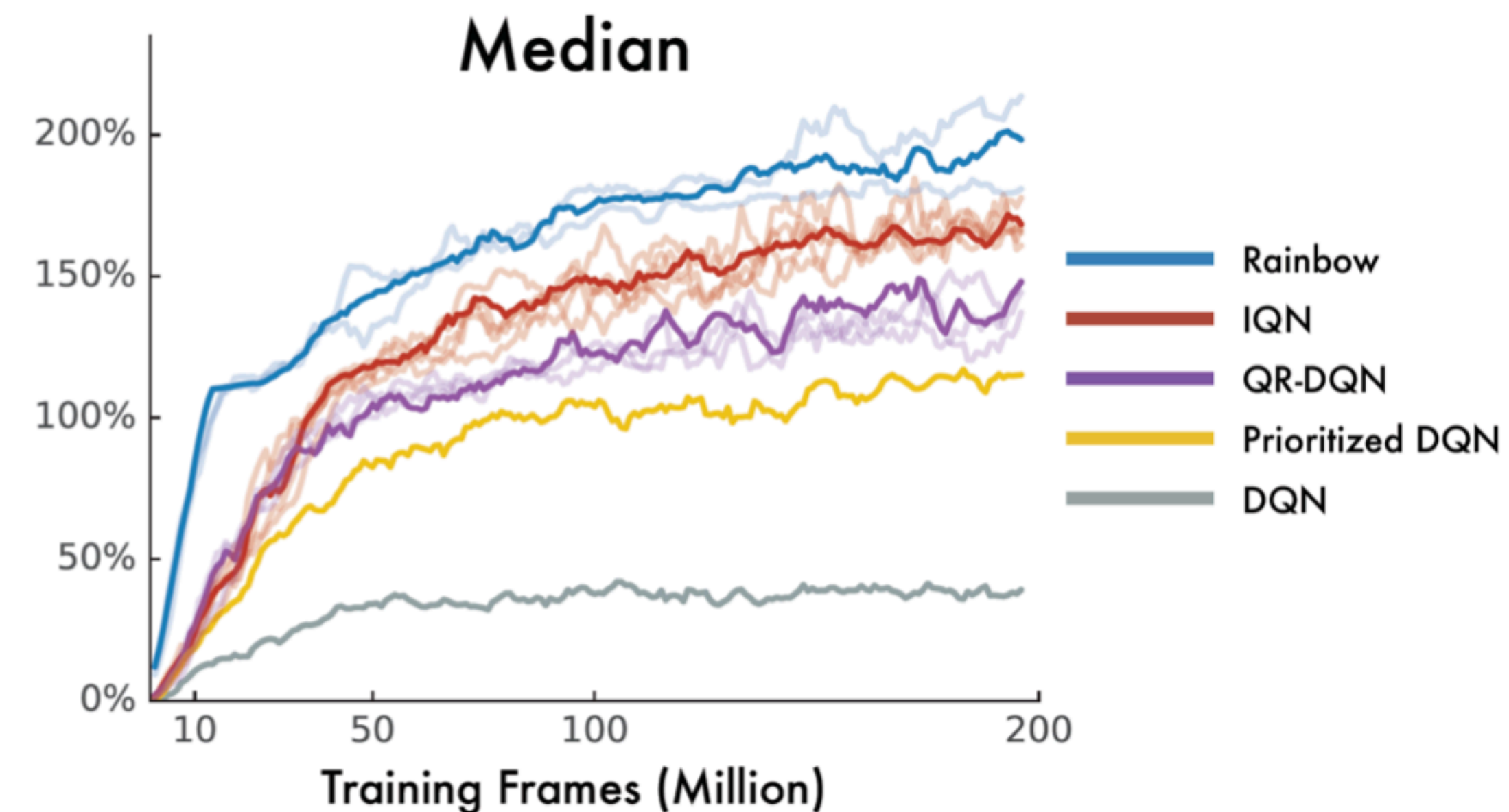
IQN outperforms Rainbow on hardest Atari games!

RESULTS - HNS on Atari-57

smaller is better!



	Mean	Median	(human starts) Median	Human Gap
DQN	228 %	79 %	68 %	0.334
Prio. Duel.	592 %	172 %	128 %	0.178
C51	701 %	178 %	125 %	0.152
QR-DQN	864 %	193 %	153 %	0.165
IQN	1019 %	218 %	162 %	0.141
QR-DQN based Rainbow	1189 %	230 %	125 %	0.144



IQN outperforms Rainbow on hardest Atari games!

RESULTS - HNS on Atari-57

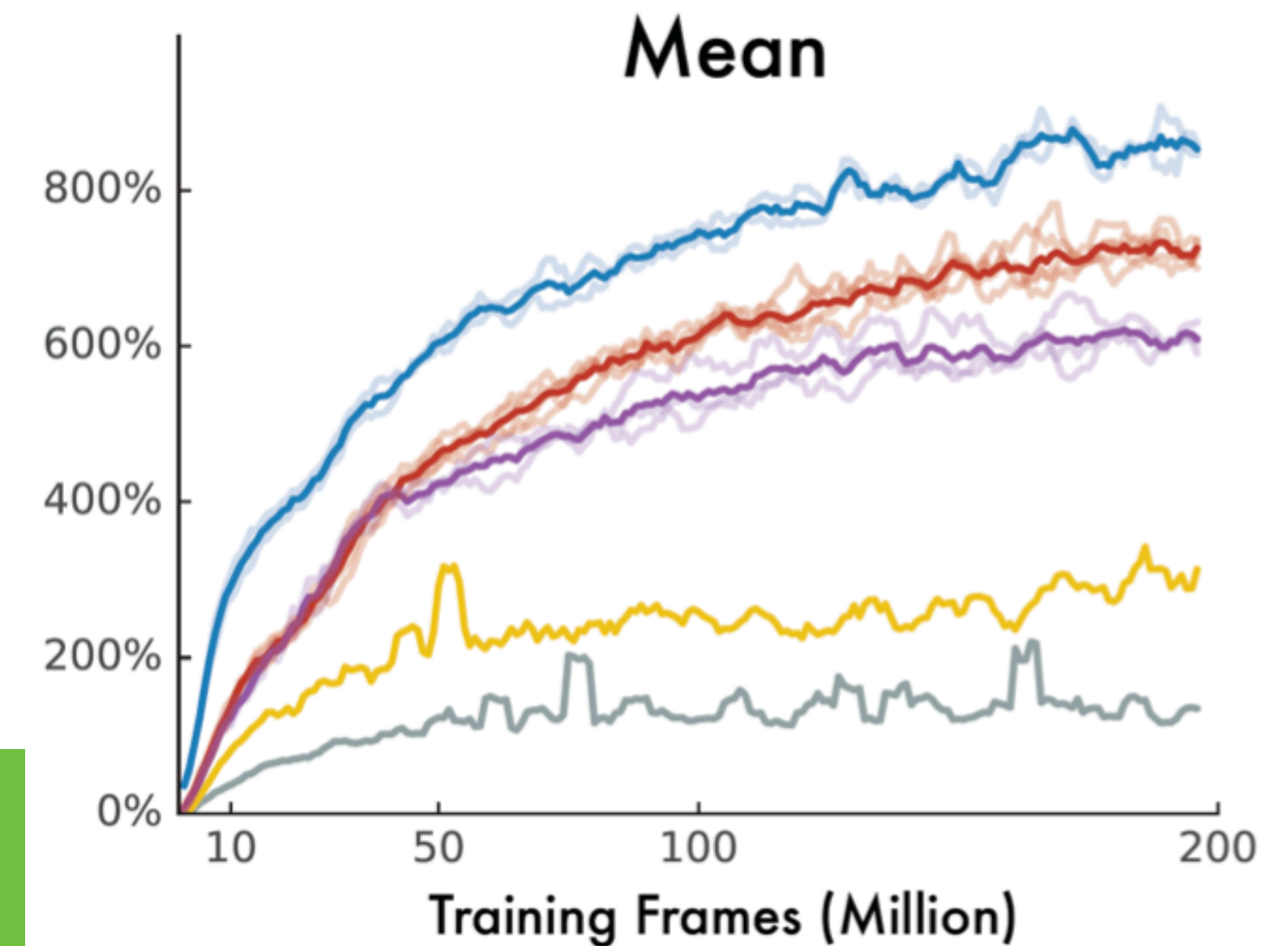
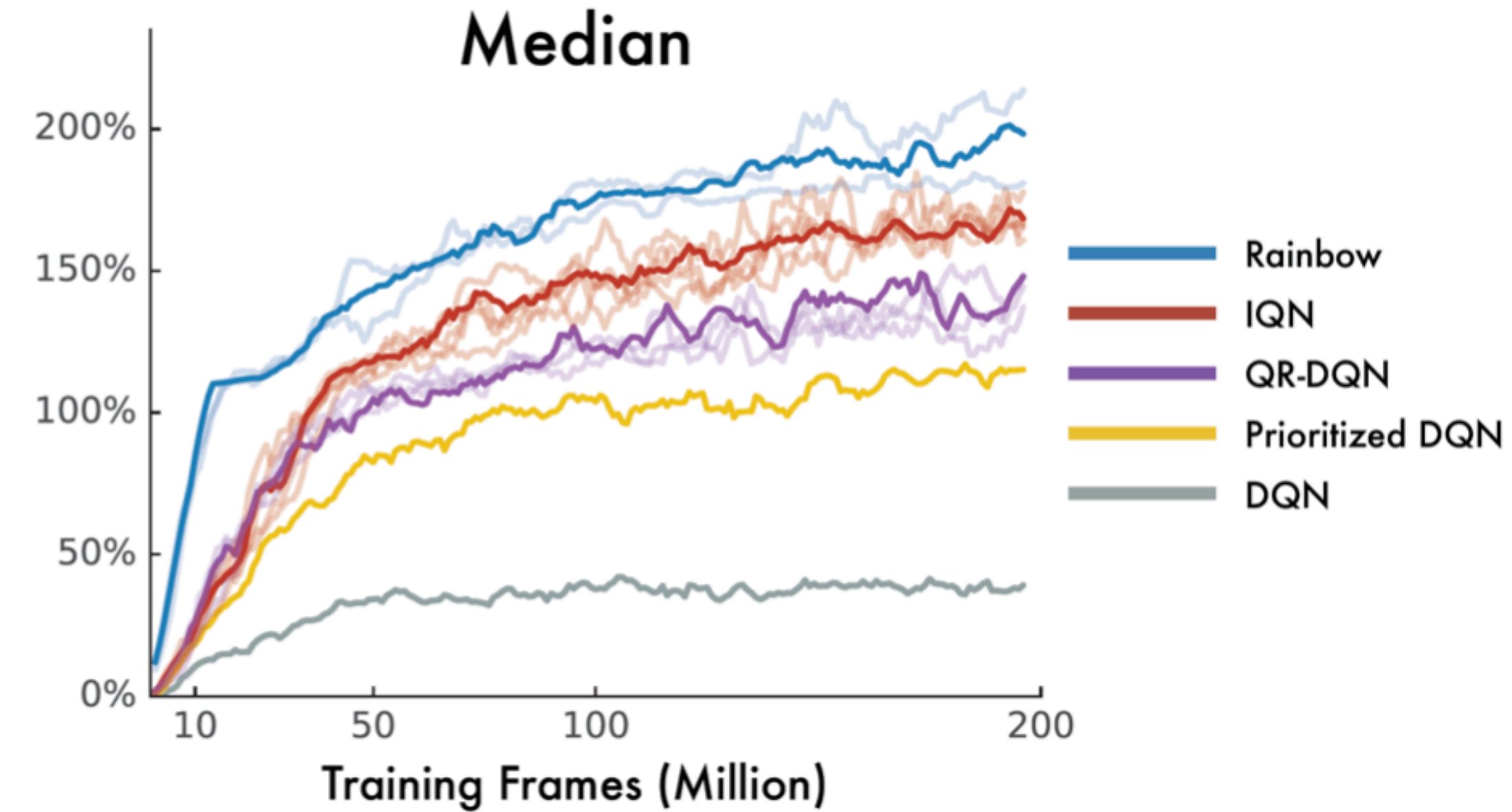
smaller is better!



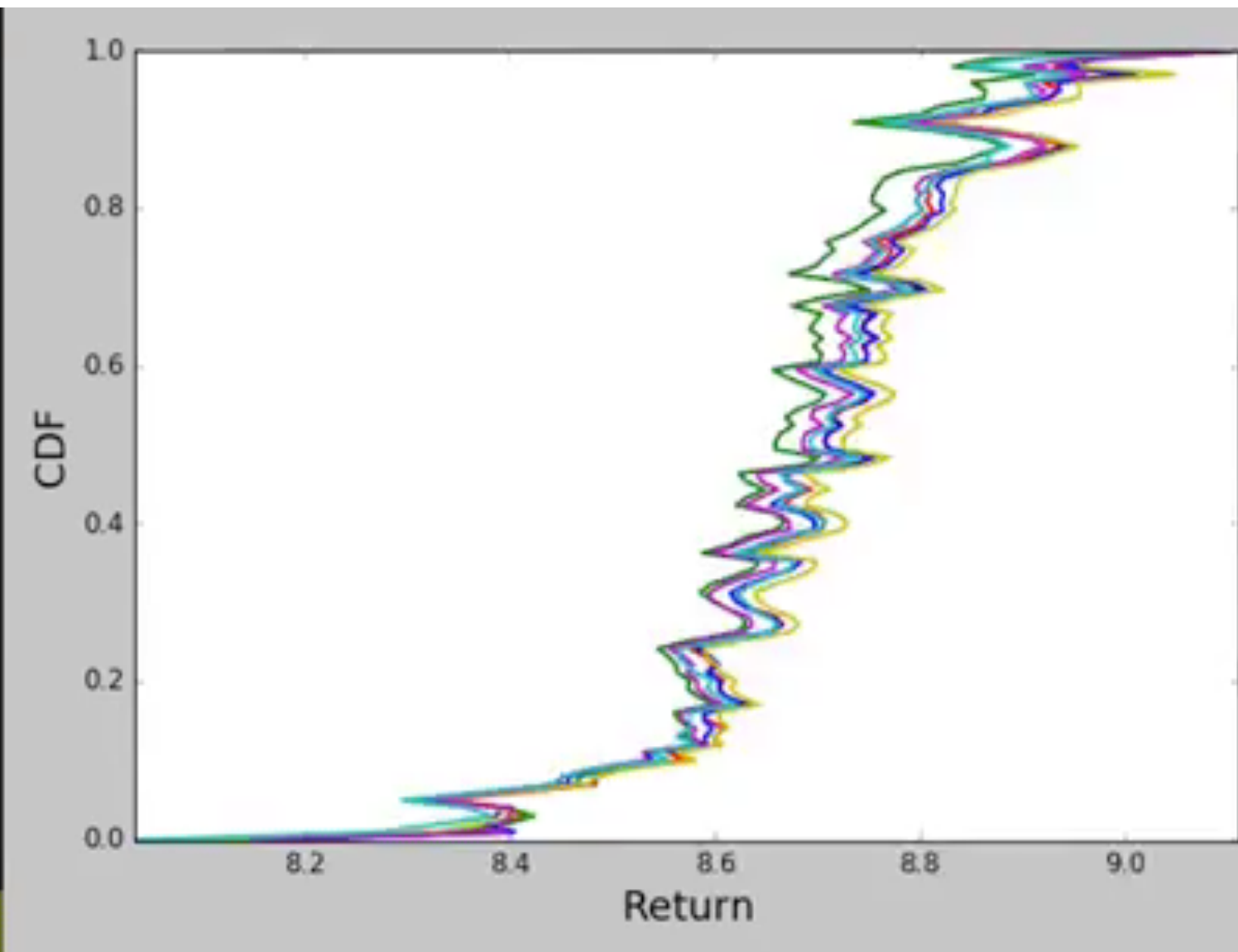
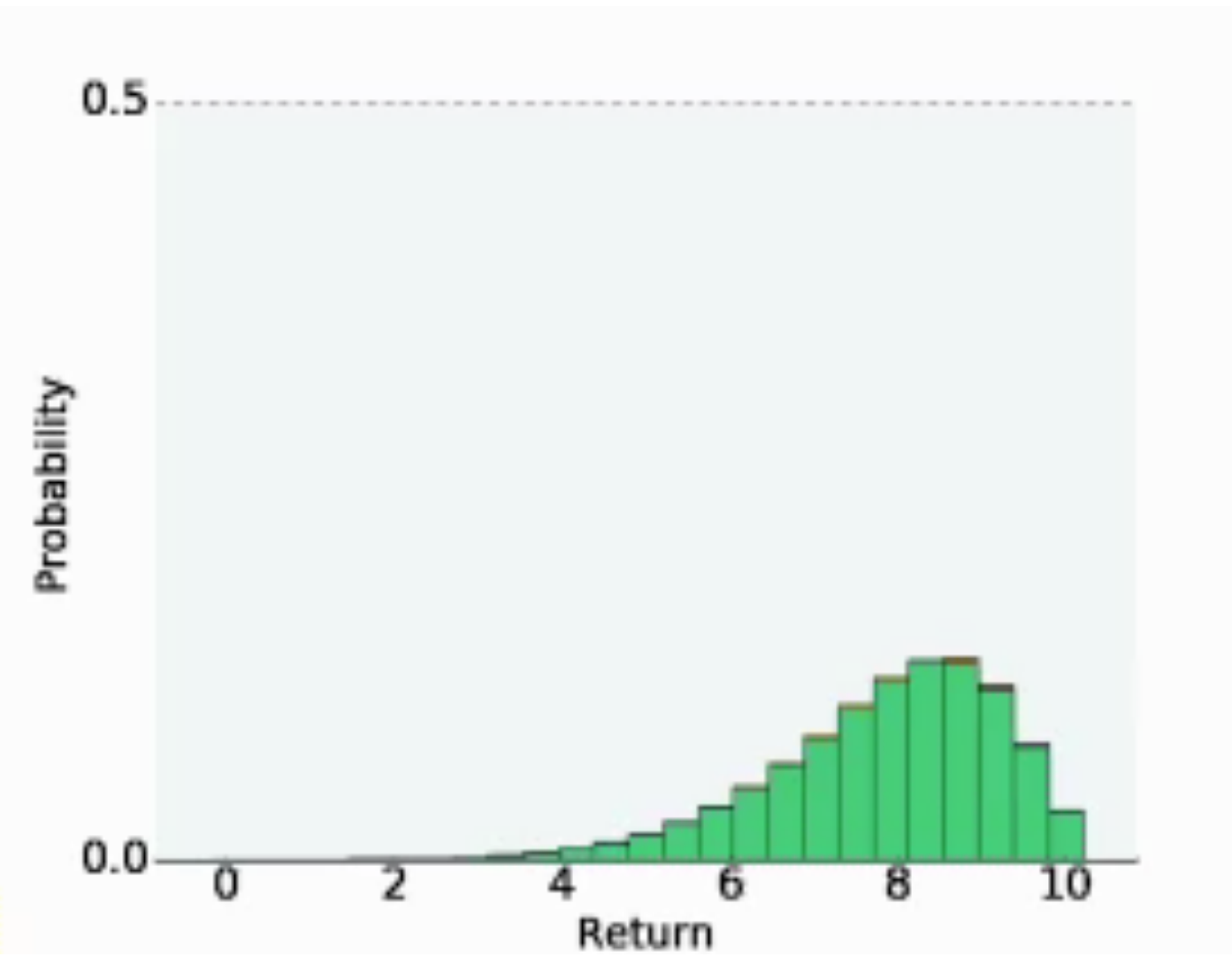
	Mean	Median	Median (human starts)	Human Gap
DQN	228 %	79 %	68 %	0.334
Prio. Duel.	592 %	172 %	128 %	0.178
C51	701 %	178 %	125 %	0.152
QR-DQN	864 %	193 %	153 %	0.165
IQN	1019 %	218 %	162 %	0.141
Rainbow	1189 %	230 %	125 %	0.144

QR-DQN based

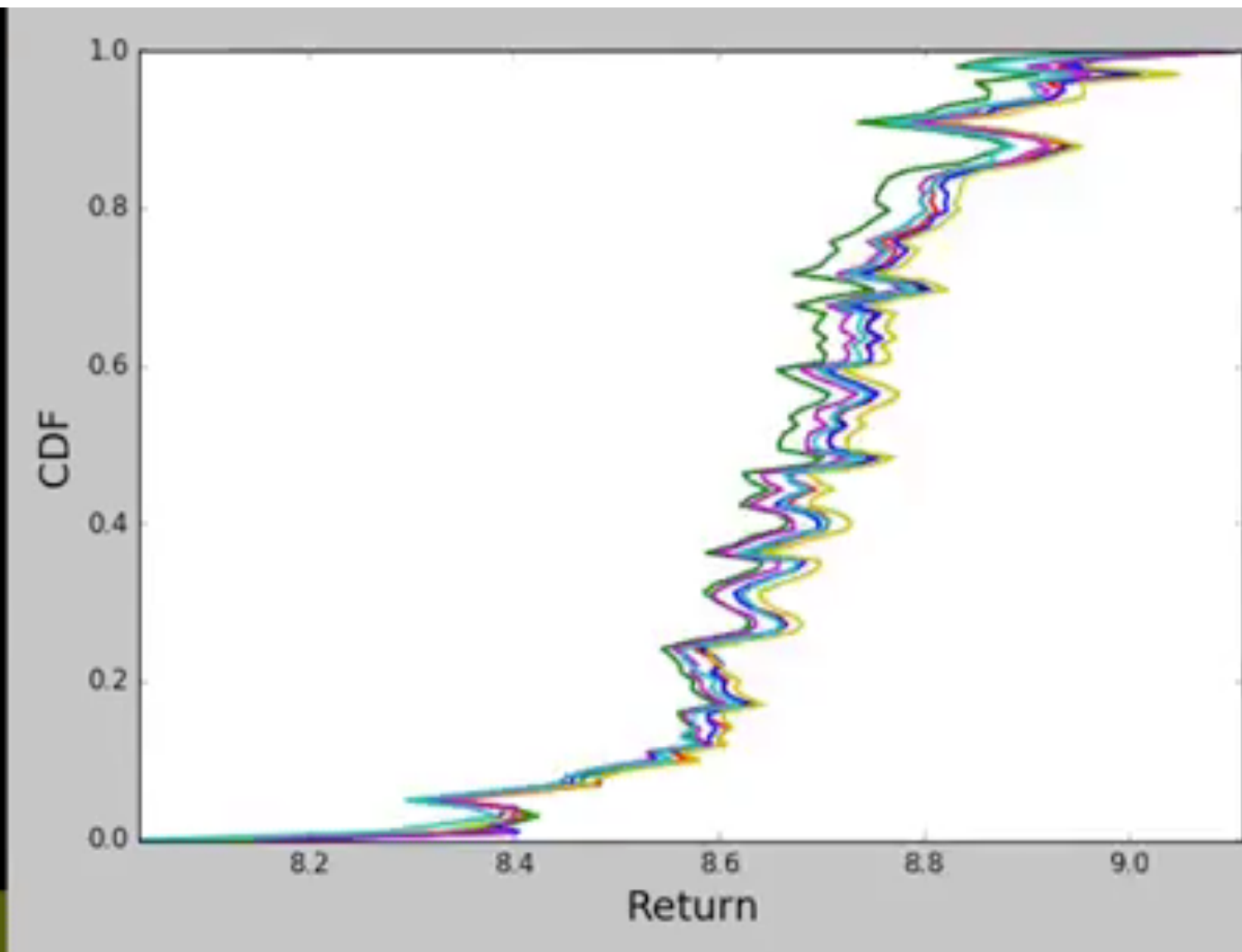
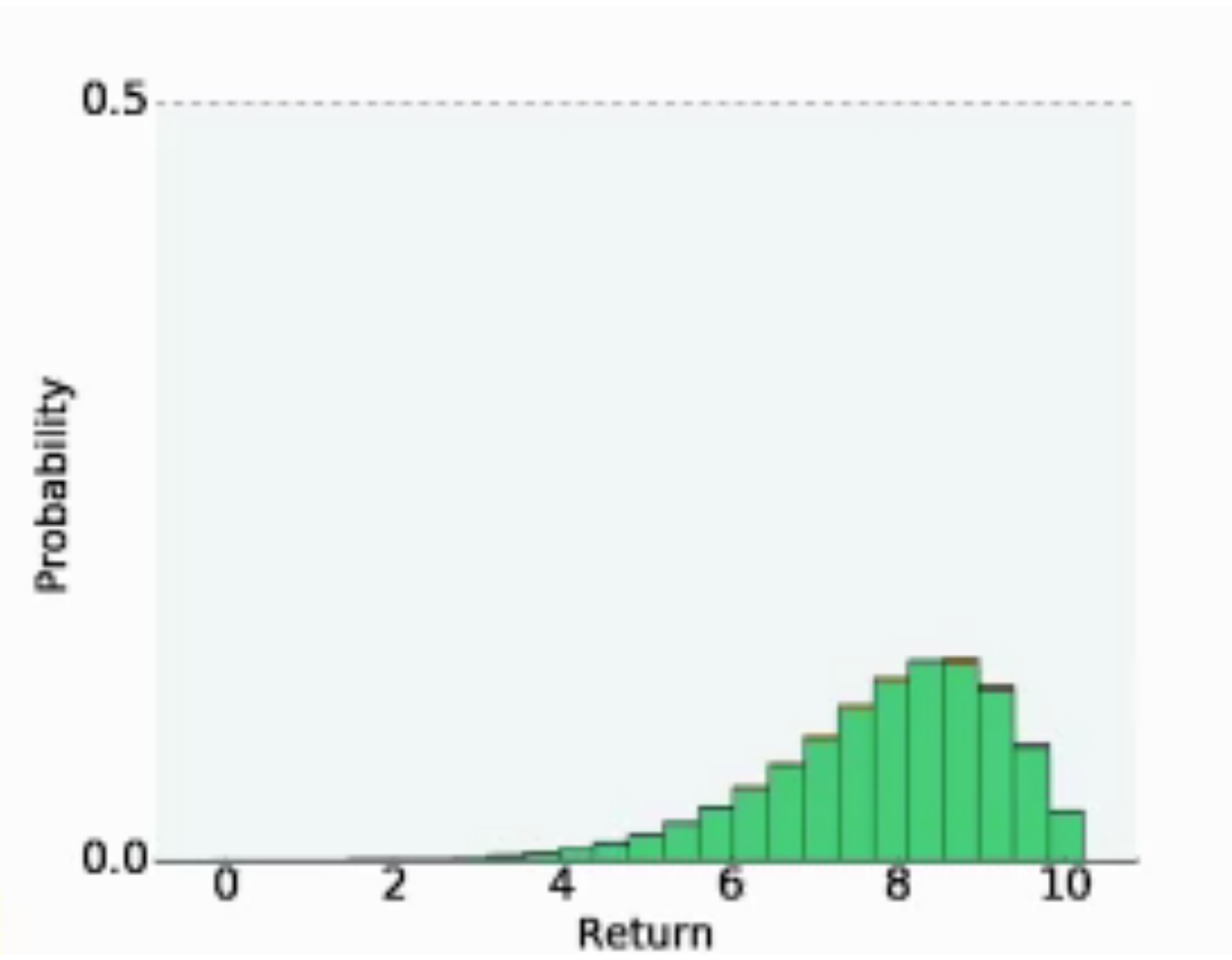
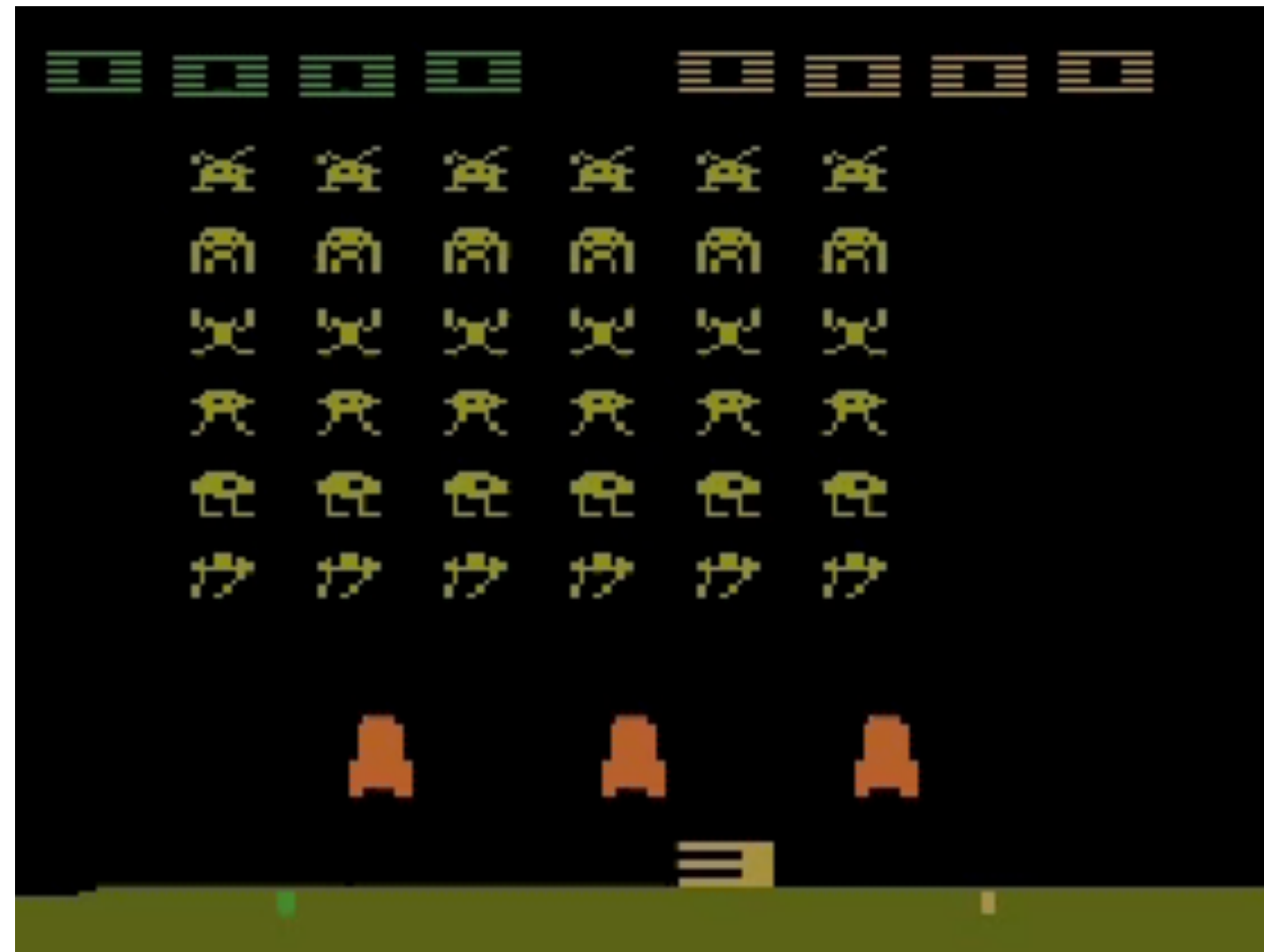
IQN outperforms Rainbow on hardest Atari games!



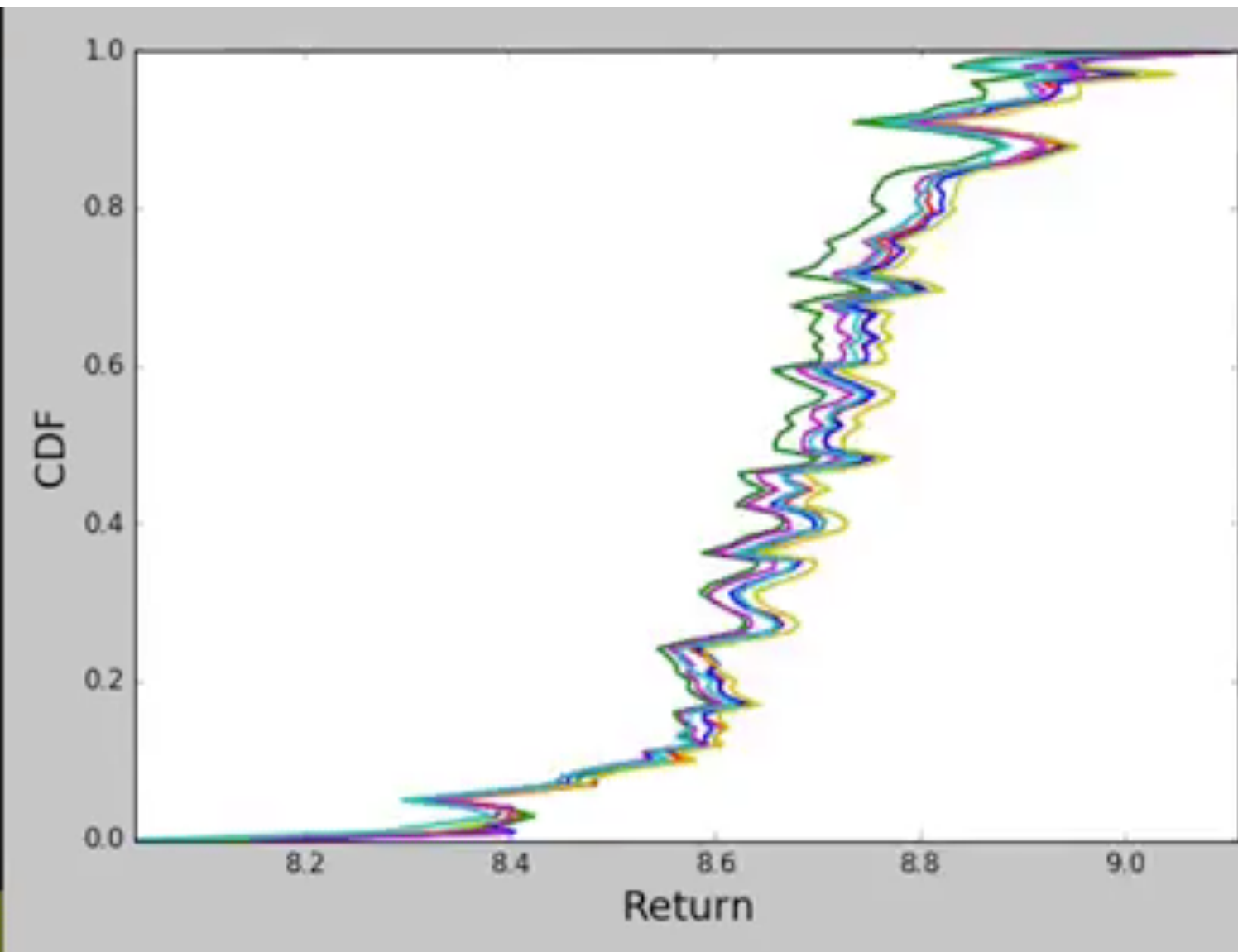
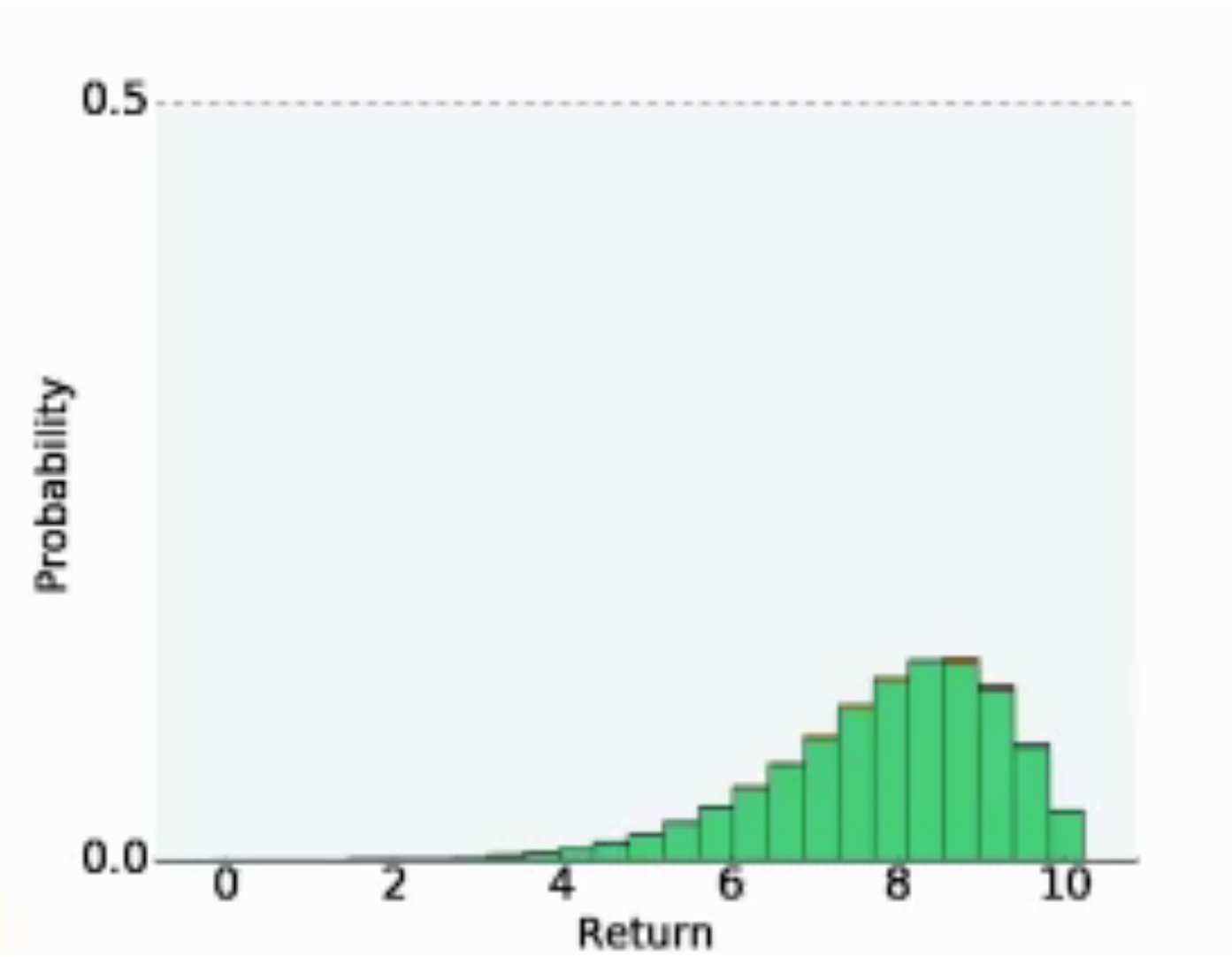
Demo



Demo



Demo



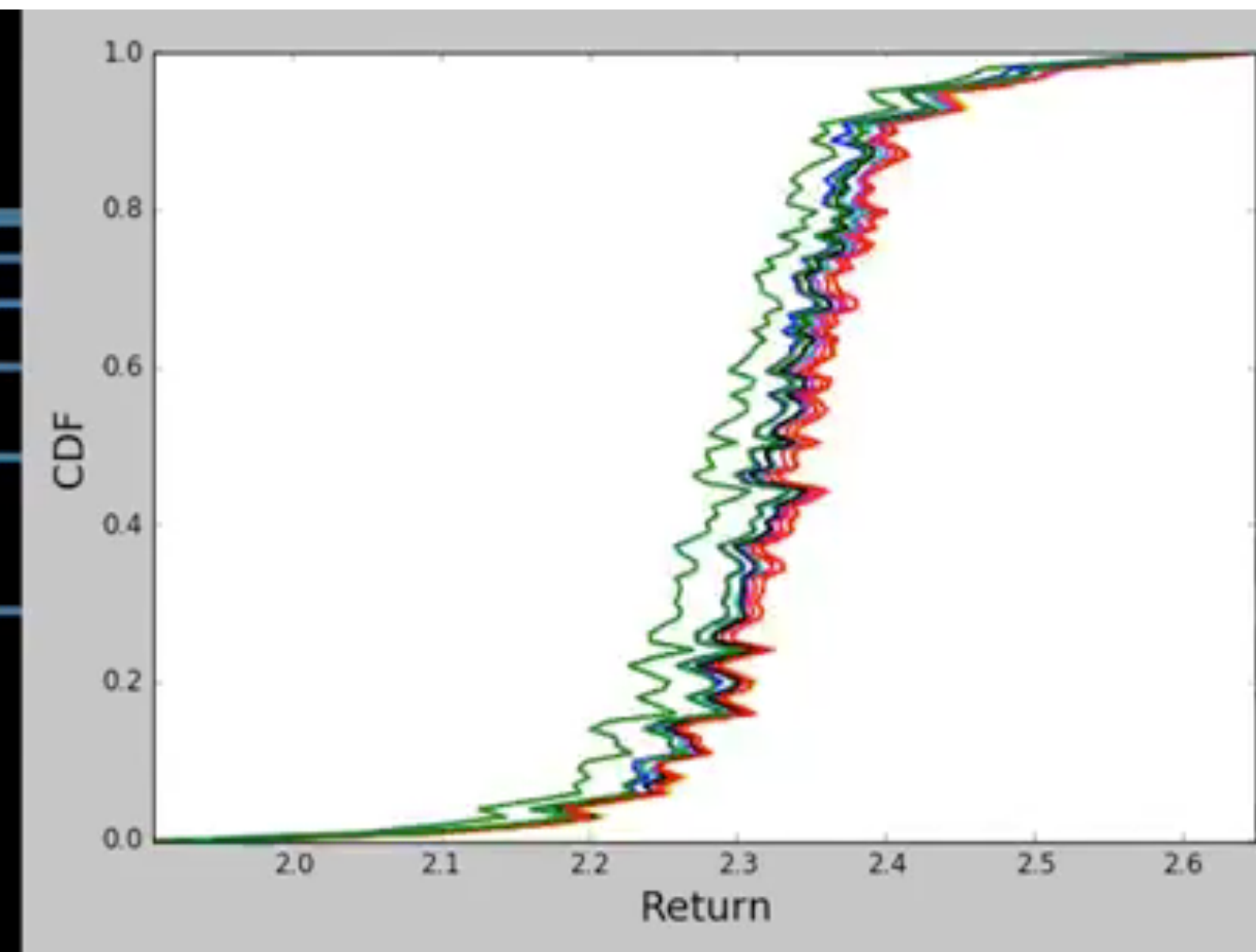
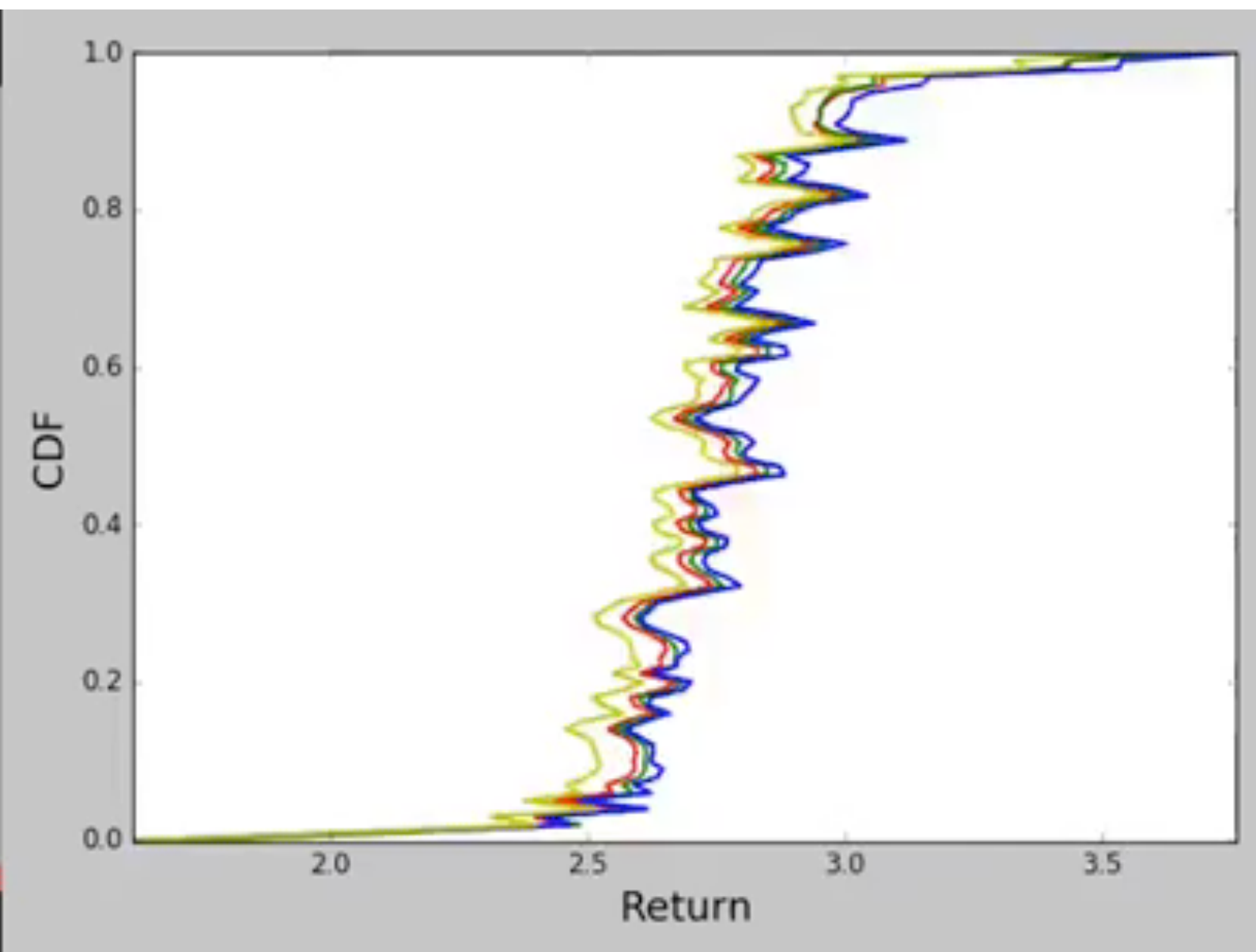
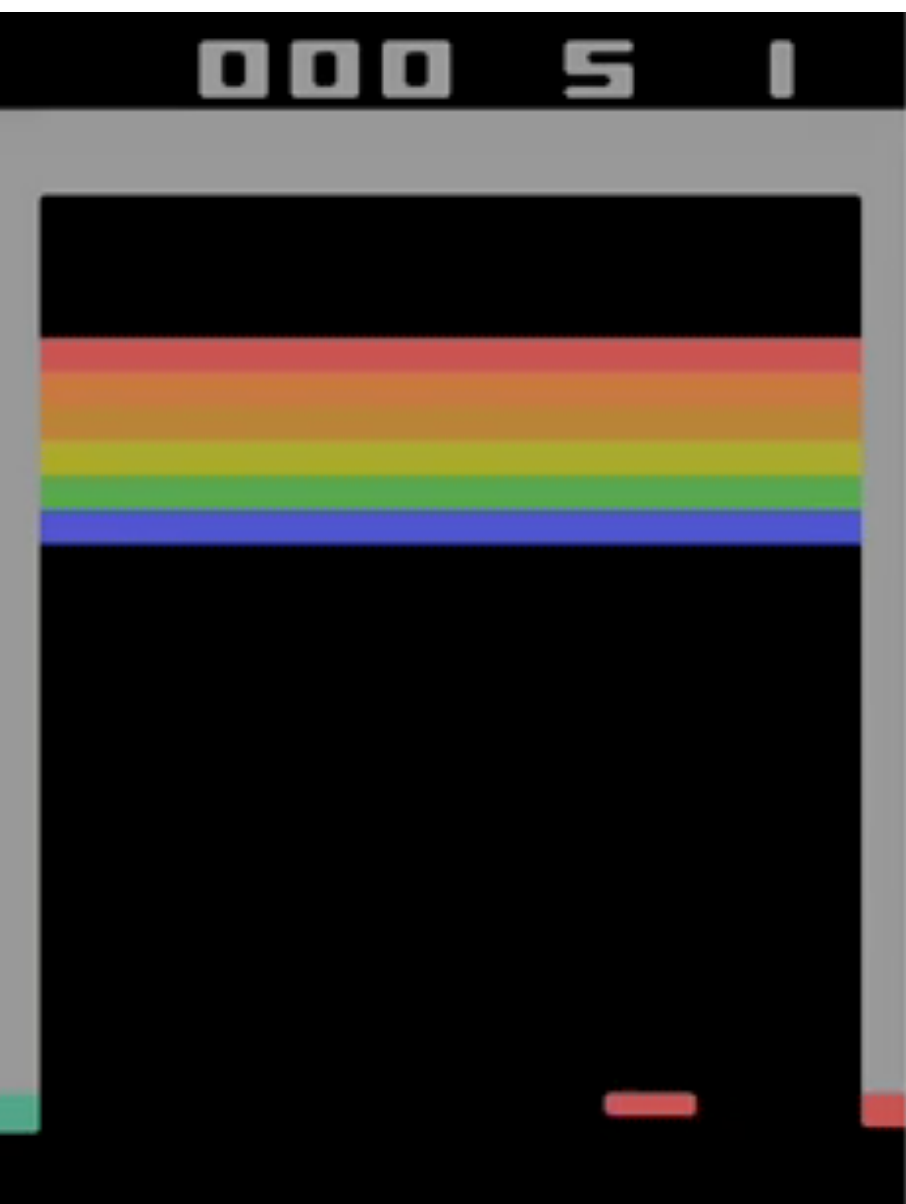
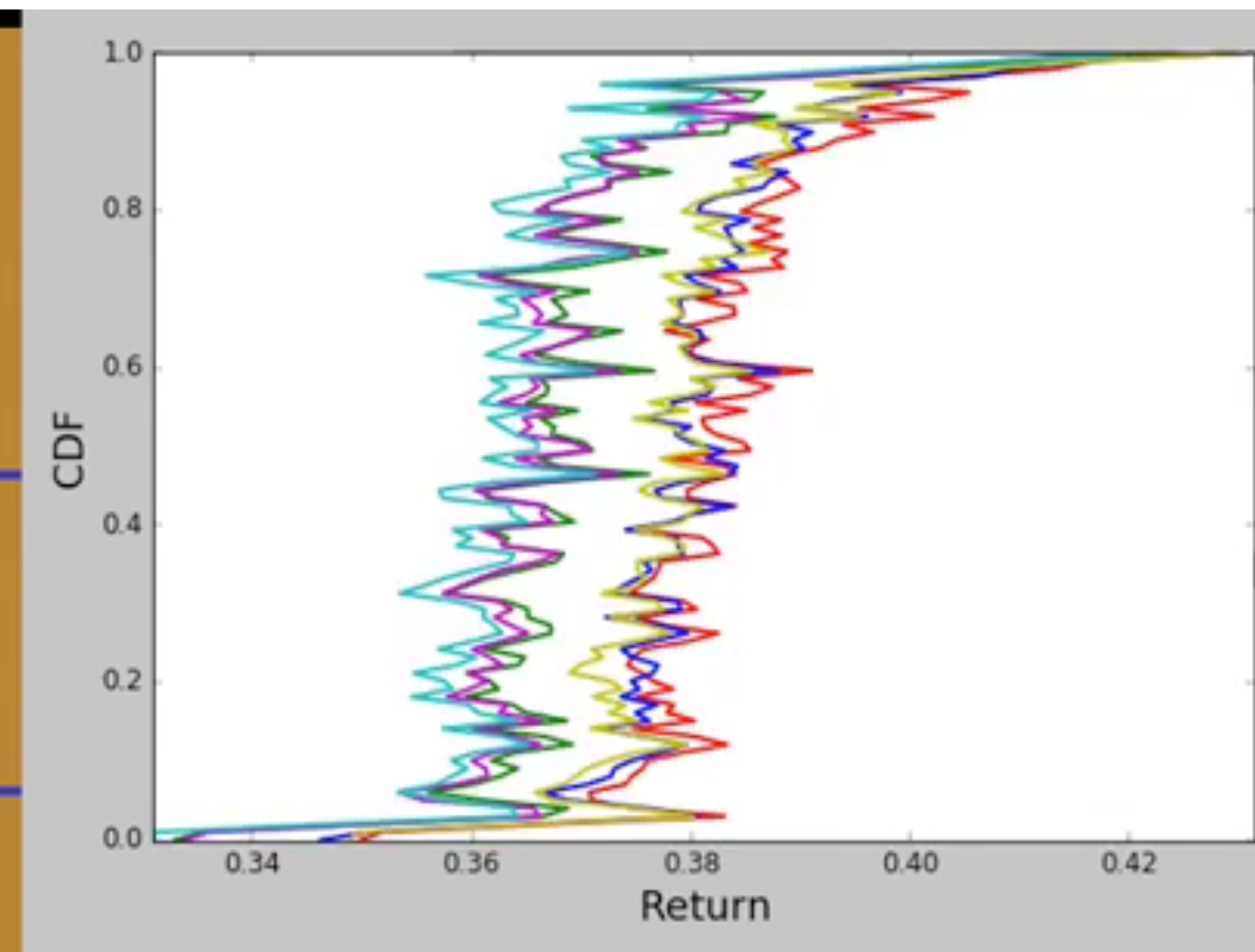
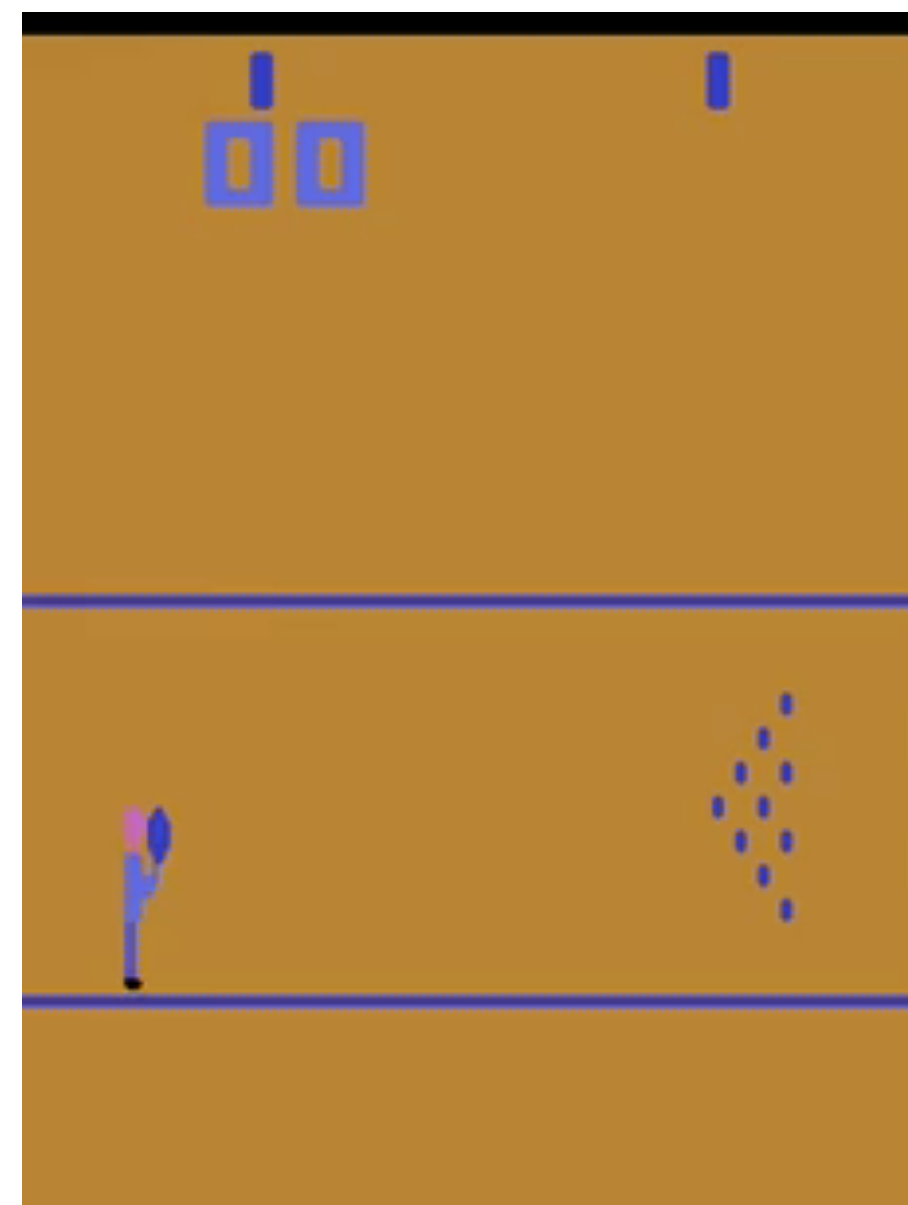
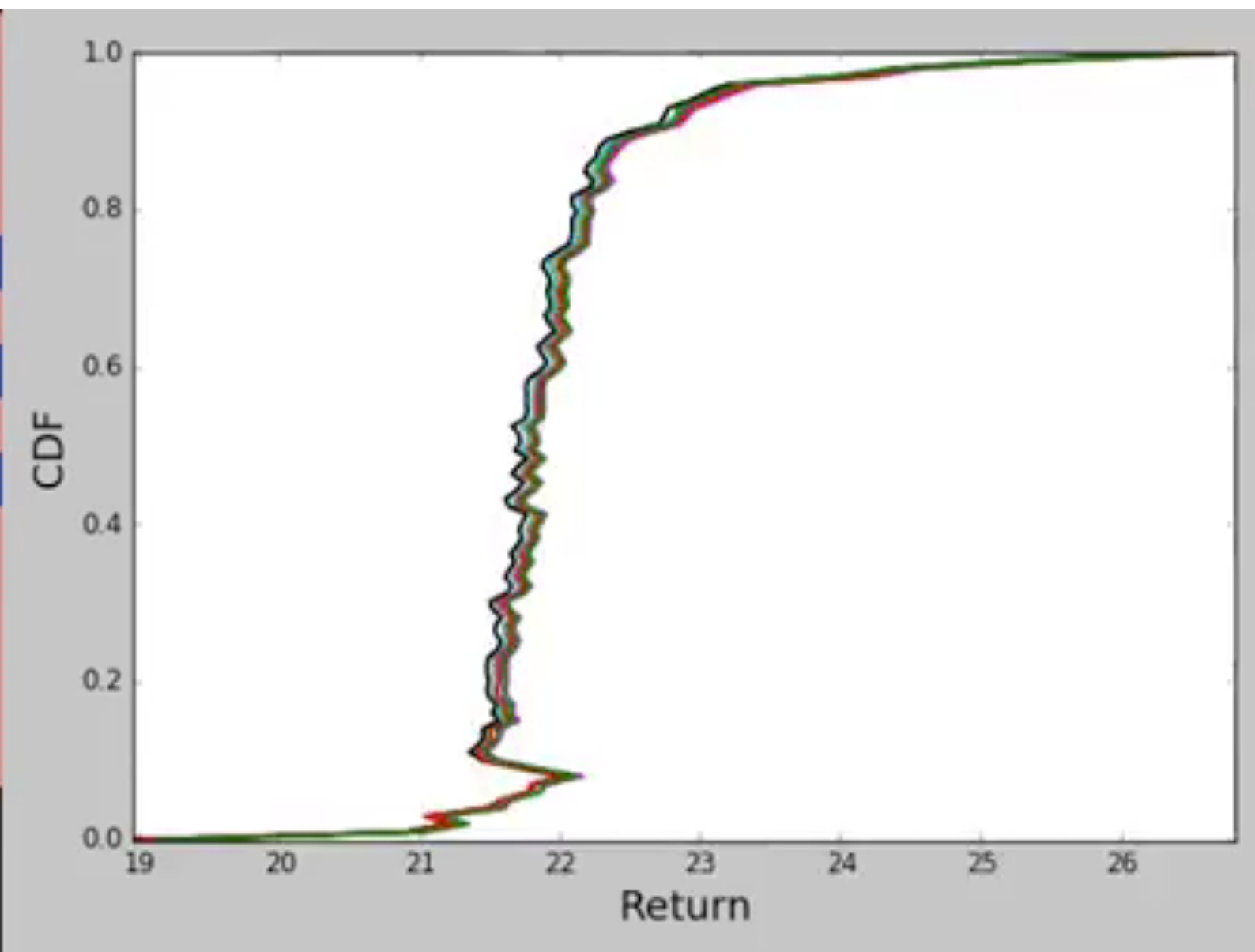
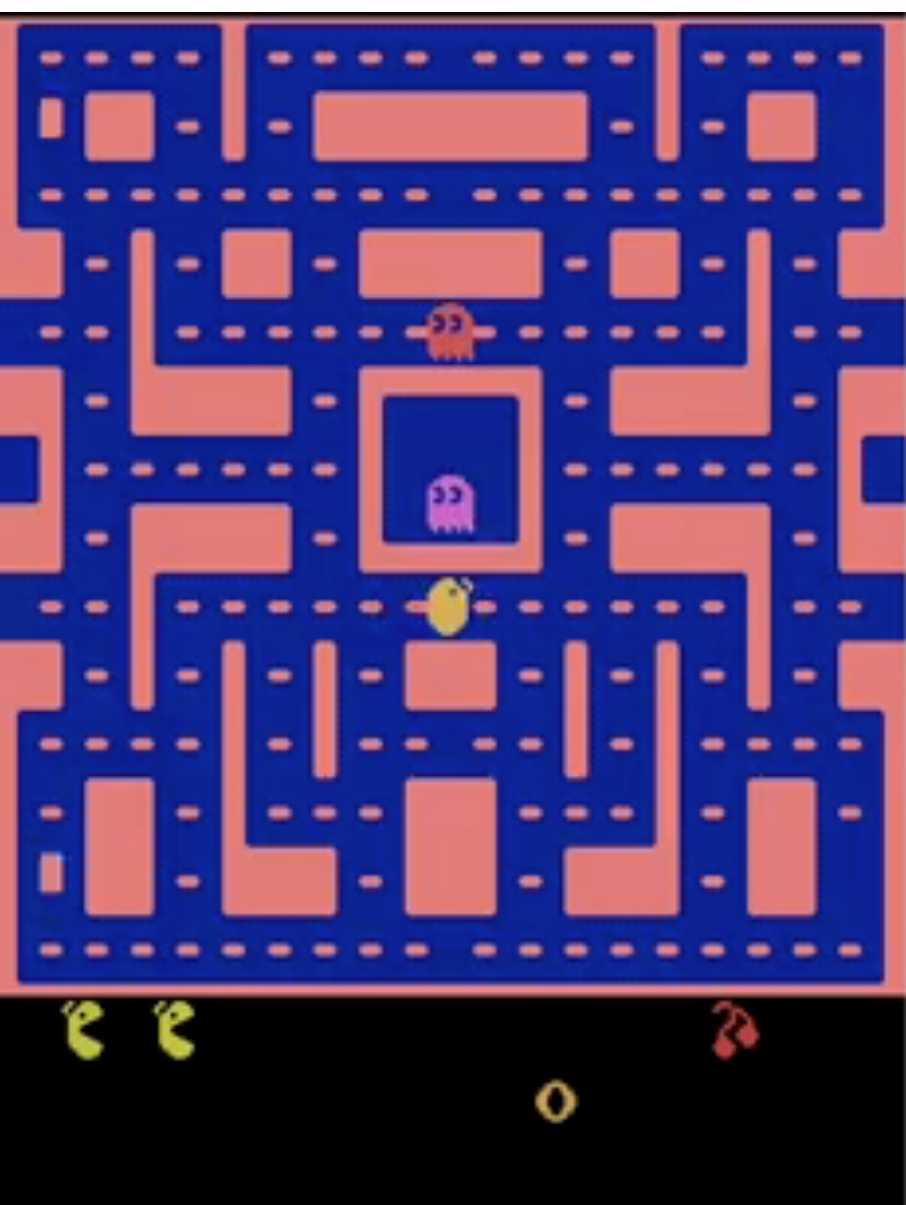
Pitfalls and future outlook

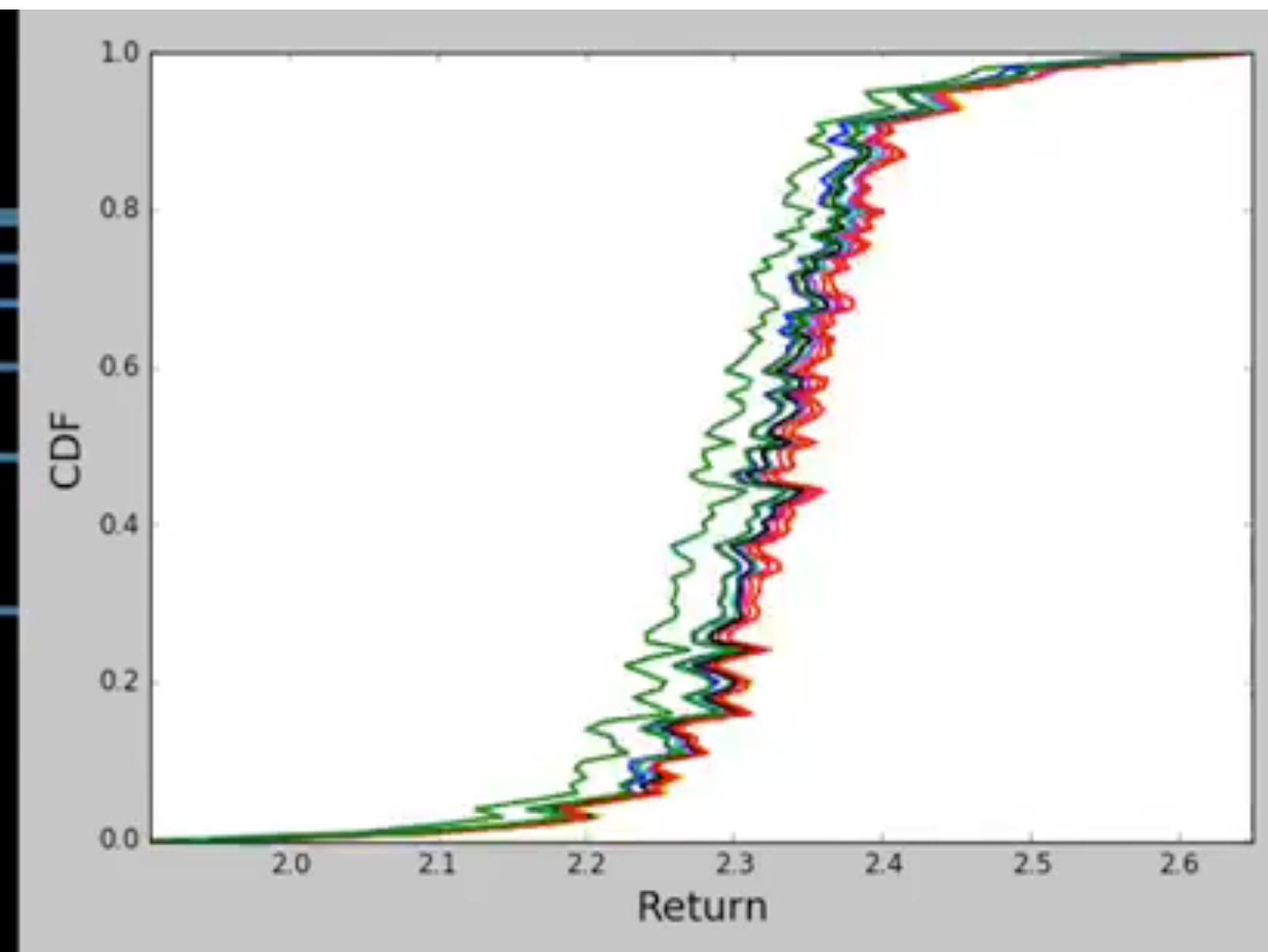
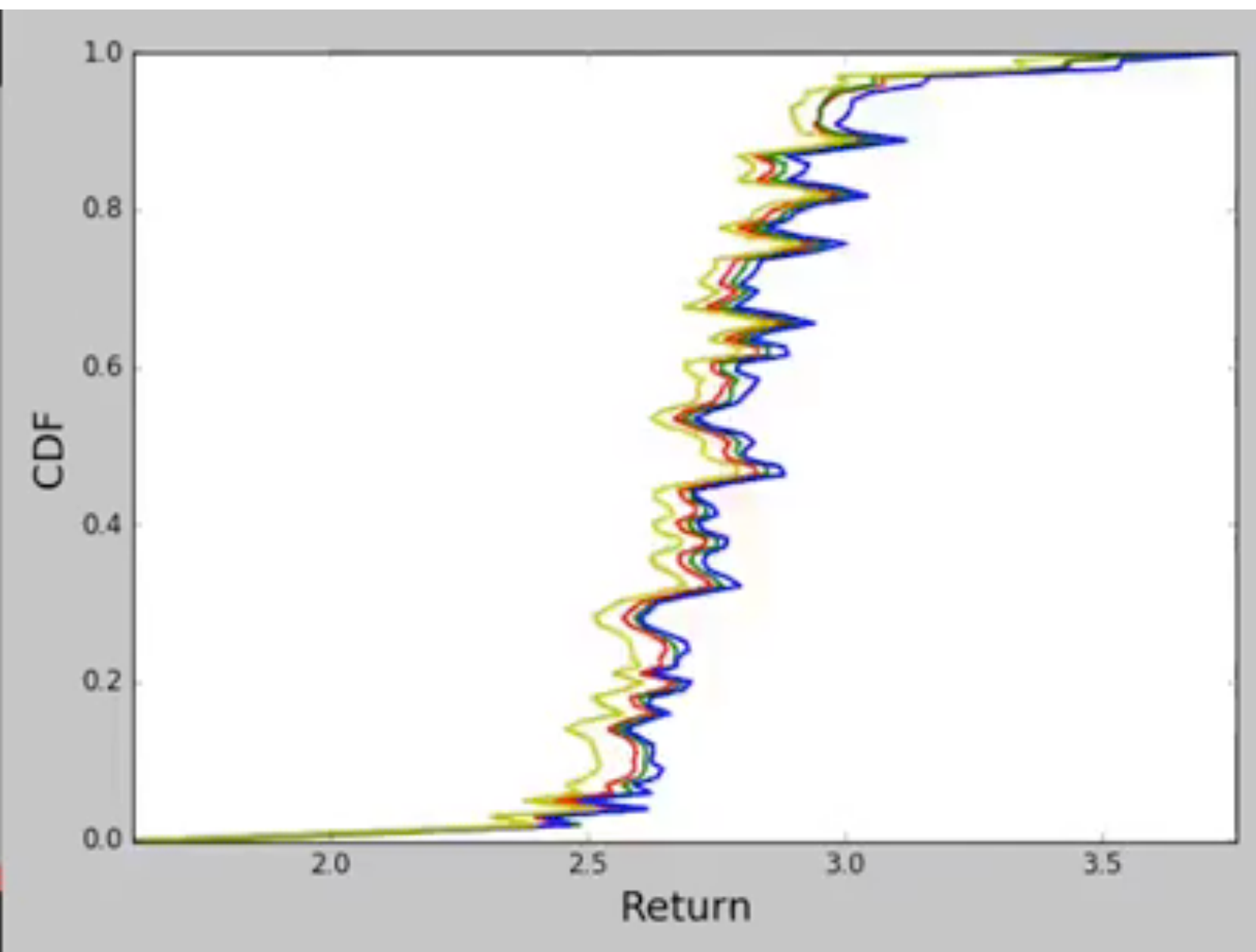
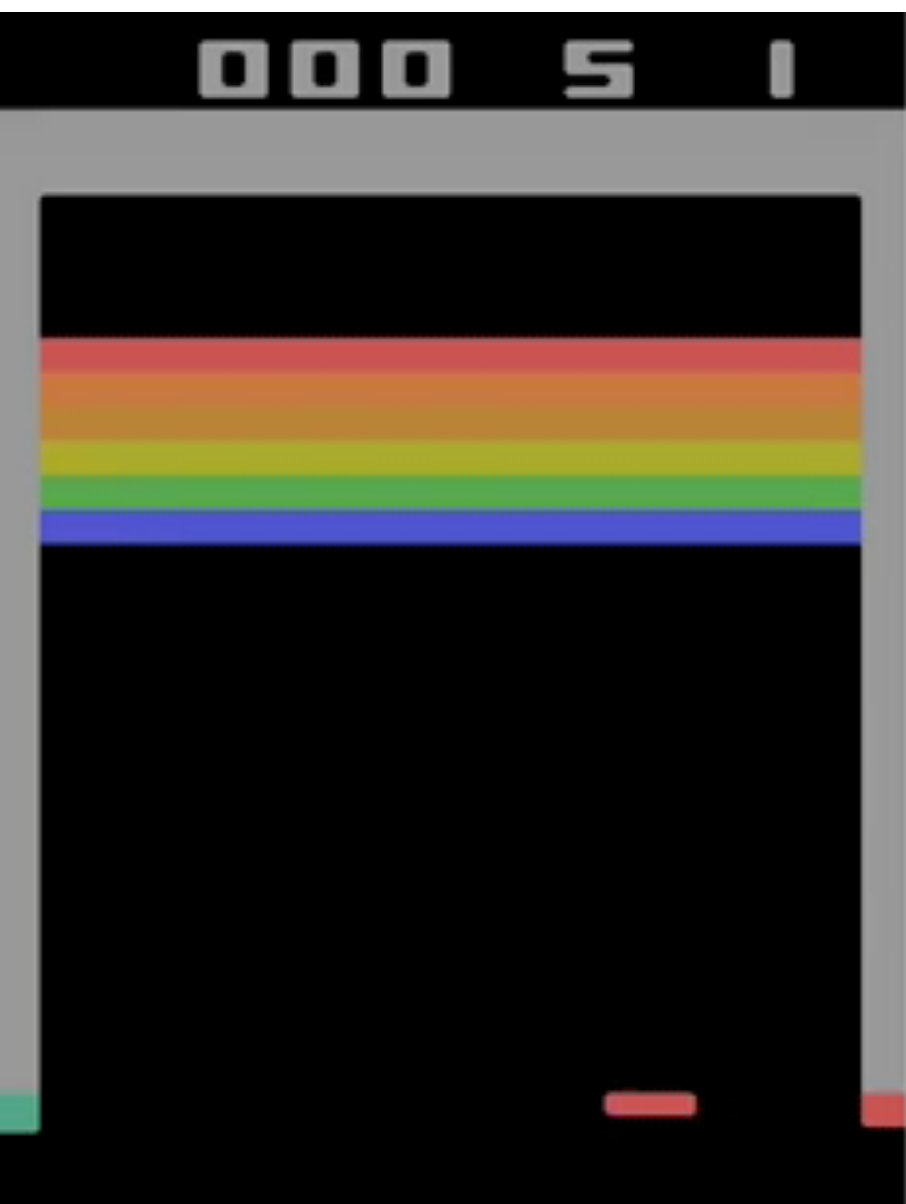
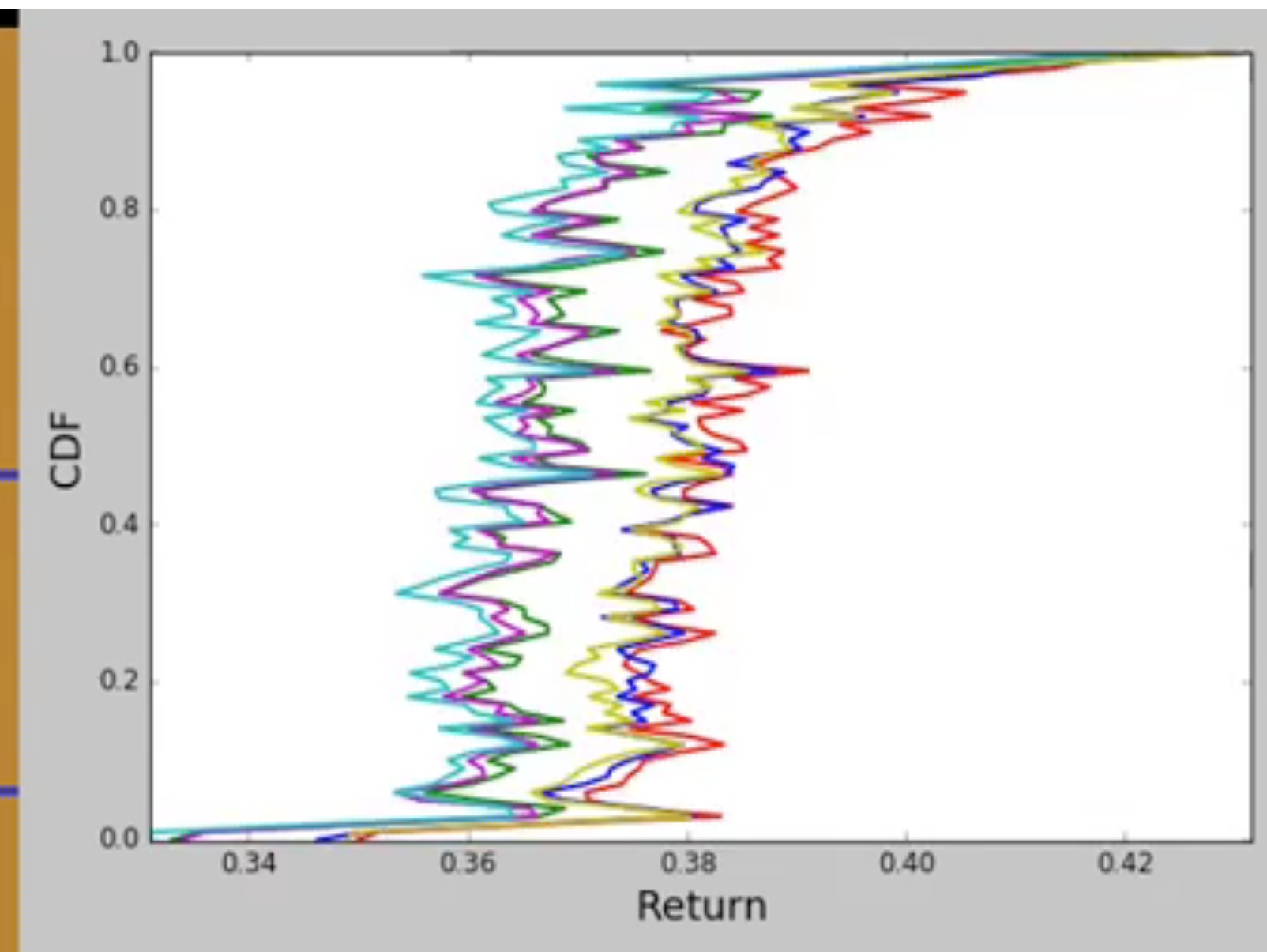
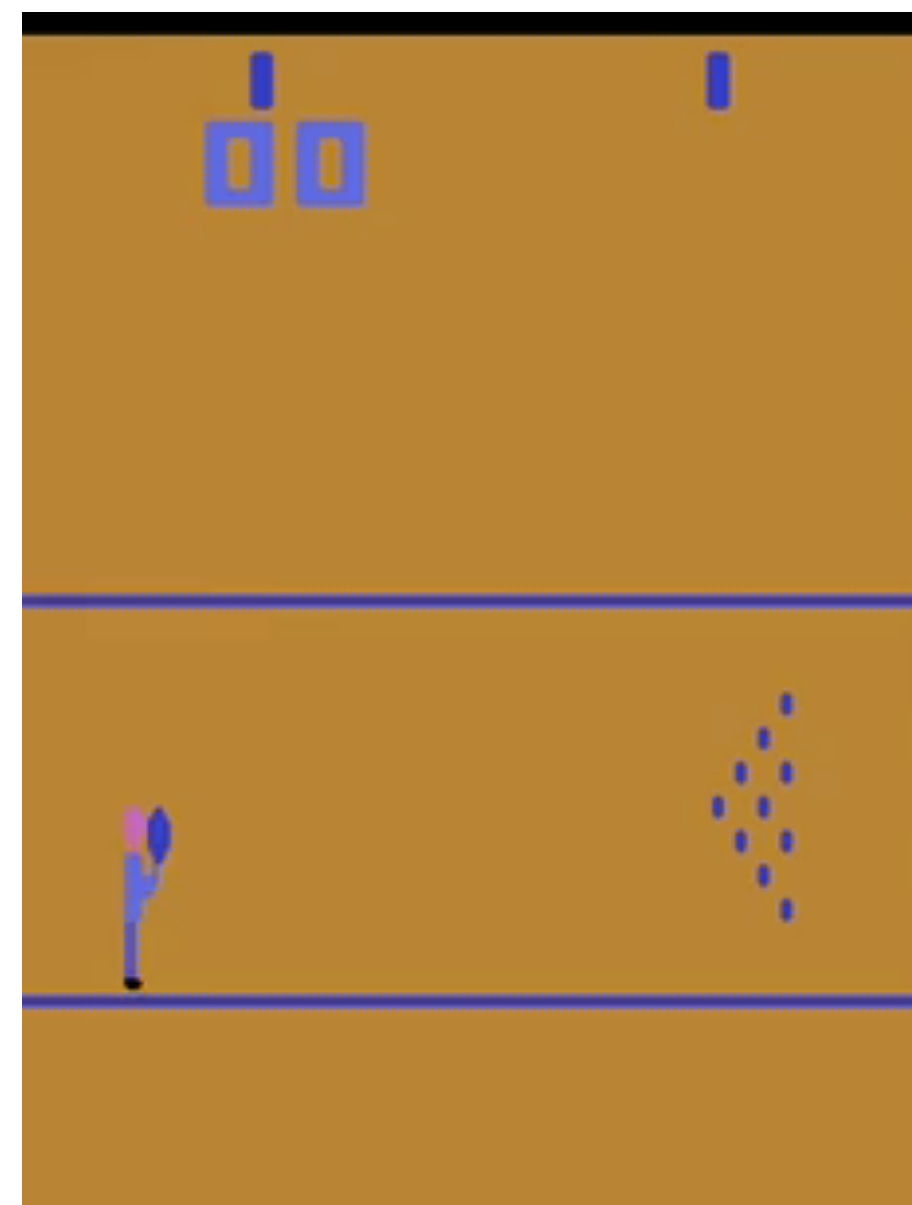
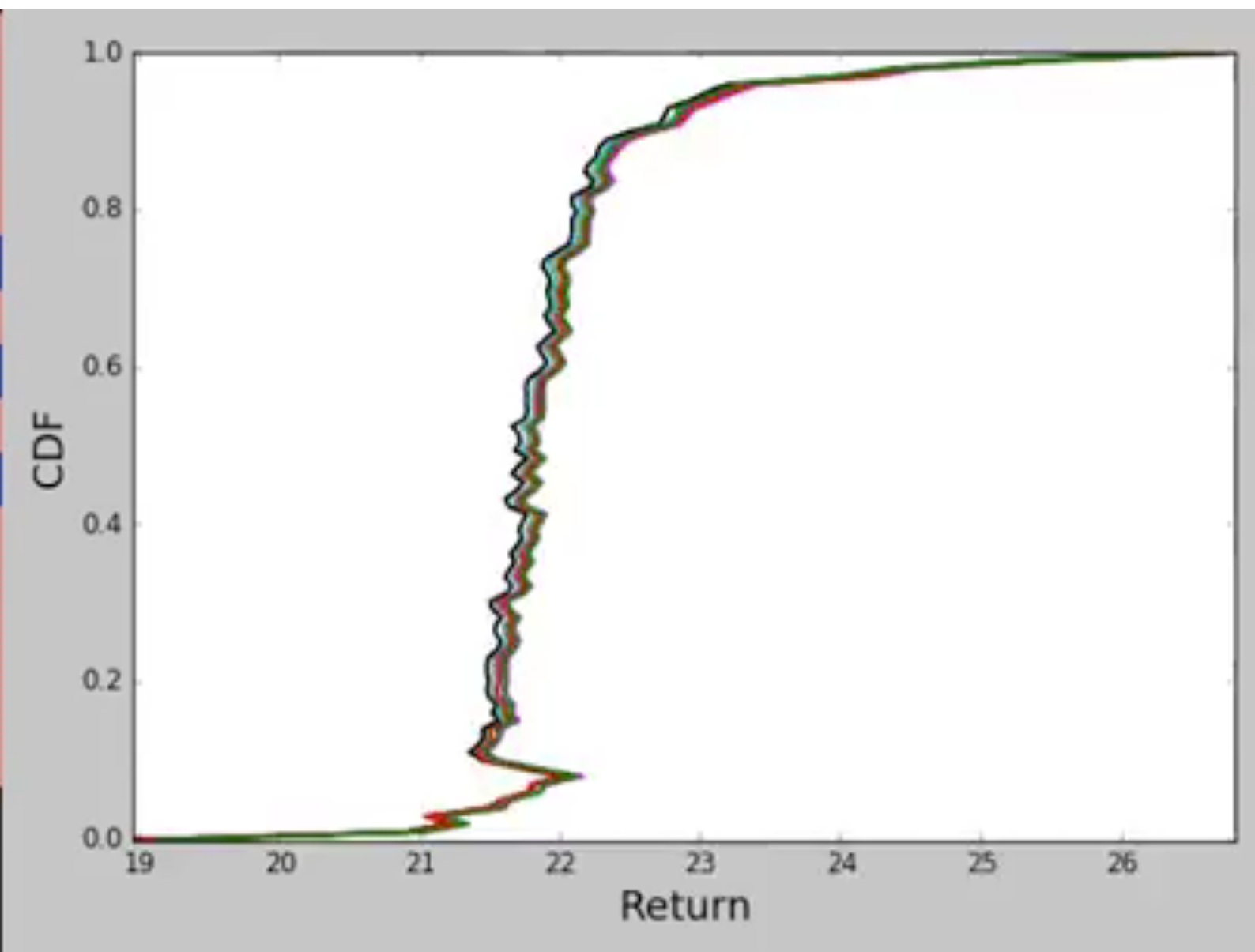
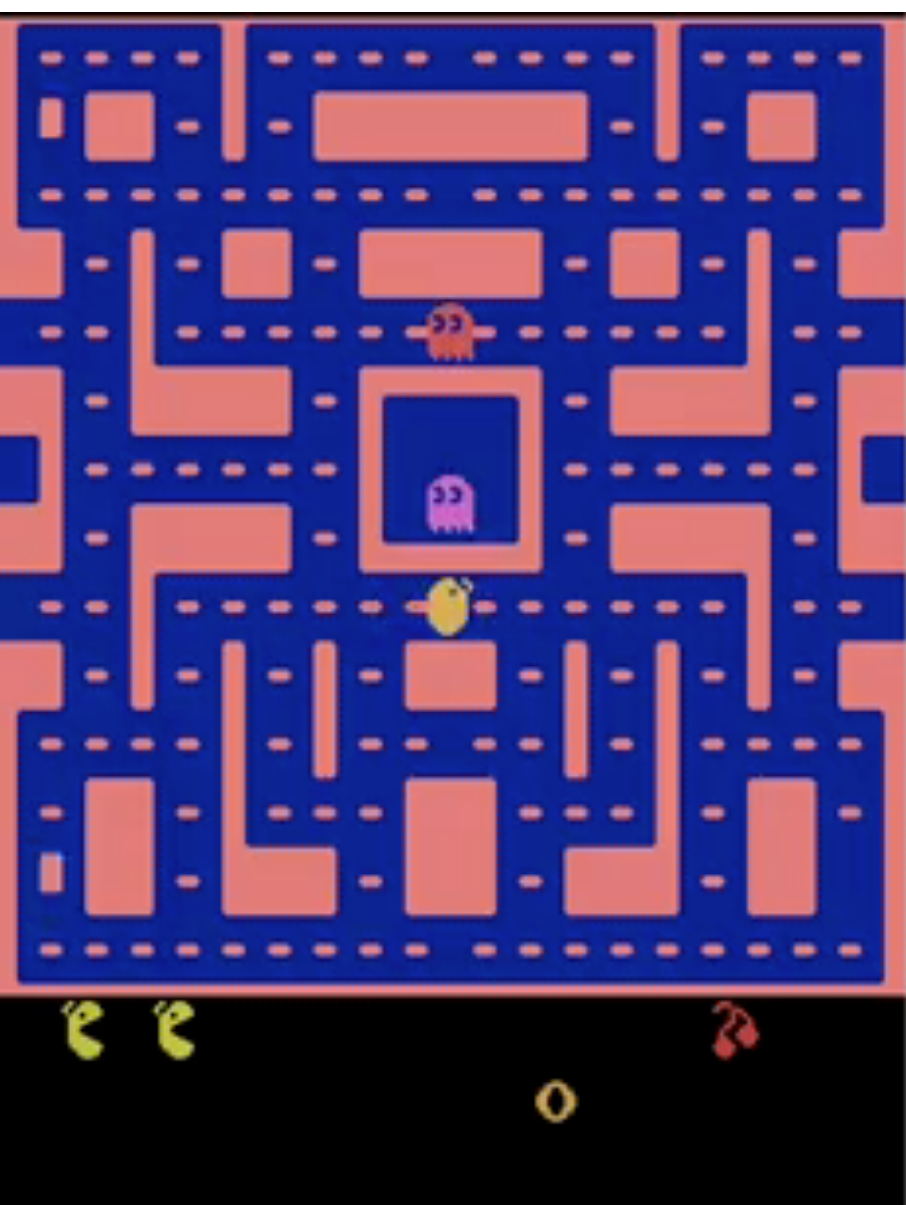
- Pitfalls

- Only the policy **mean** used for evaluation
- Convergence in theory only for **fixed** quantiles
- Evaluated only in **discrete** action environments
- No guarantee that quantiles are ordered

- Future

- Distribution over policies
- Continuous action domains
- Rainbow-IQN
- Solve the quantile ordering problem





References

- [1] Sobel, L.M. 1982. “*The variance of discounted markov decision processes*”
- [2] Morimura T. et. al. 2010. “*Parametric Return Density Estimation for Reinforcement Learning*”
- [3] Bellemare, M.G. et. al. 2017. “*A Distributional Perspective on Reinforcement Learning*”
- [4] Barth-Maron, G. et. al. 2018. “*Distributed Distributional Deterministic Policy Gradients*”
- [5] Hessel, M. et. al. 2018. “*Rainbow: combining improvements in deep reinforcement learning*”
- [6] Gruslys, A. et. al. 2018. “*The Reactor: a fast and sampleefficient actor-critic agent for reinforcement learning*”
- [7] Rowland, M. et. al. 2018. “*An Analysis of Categorical Distributional Reinforcement Learning*”
- [8] Dabney, W. et. al. 2017. “*Distributional Reinforcement Learning with Quantile Regression*”
- [9] Dabney, W. et. al. 2018. “*Implicit Quantile Networks for Distributional Reinforcement Learning*”

IQN - Data efficiency vs. Computation

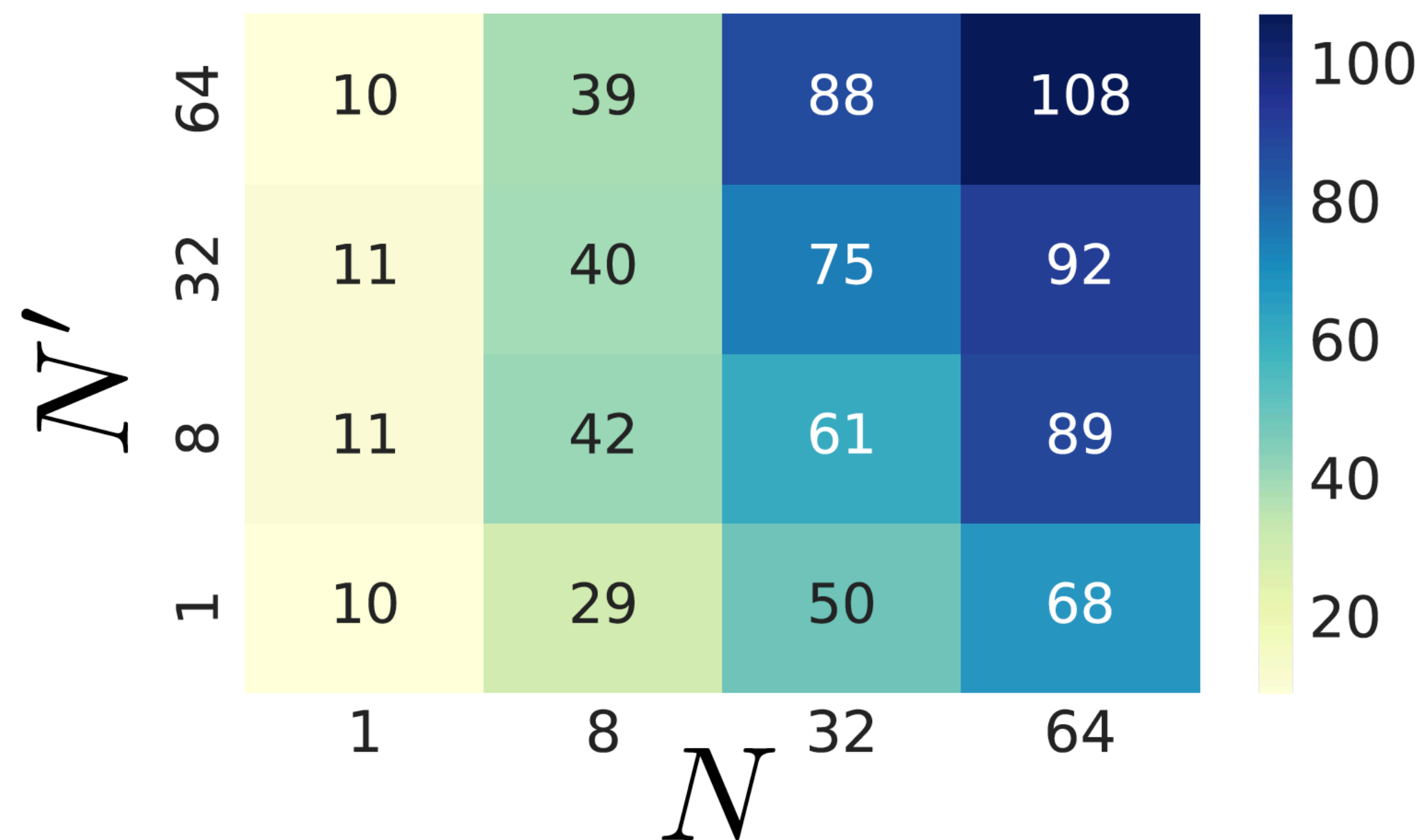
IQN - Data efficiency vs. Computation

$$\mathcal{L}_{IQN} = \sum_{\tau=\tau_1}^{\tau_N} \sum_{\tau'=\tau_1}^{\tau_{N'}} \delta_t^{\tau, \tau'} (\tau - \mathbf{I}_{\delta_t^{\tau, \tau'} < 0})$$

IQN - Data efficiency vs. Computation

$$\mathcal{L}_{IQN} = \sum_{\tau=\tau_1}^{\tau_N} \sum_{\tau'=\tau_1}^{\tau_{N'}} \delta_t^{\tau, \tau'} (\tau - \mathbf{I}_{\delta_t^{\tau, \tau'} < 0})$$

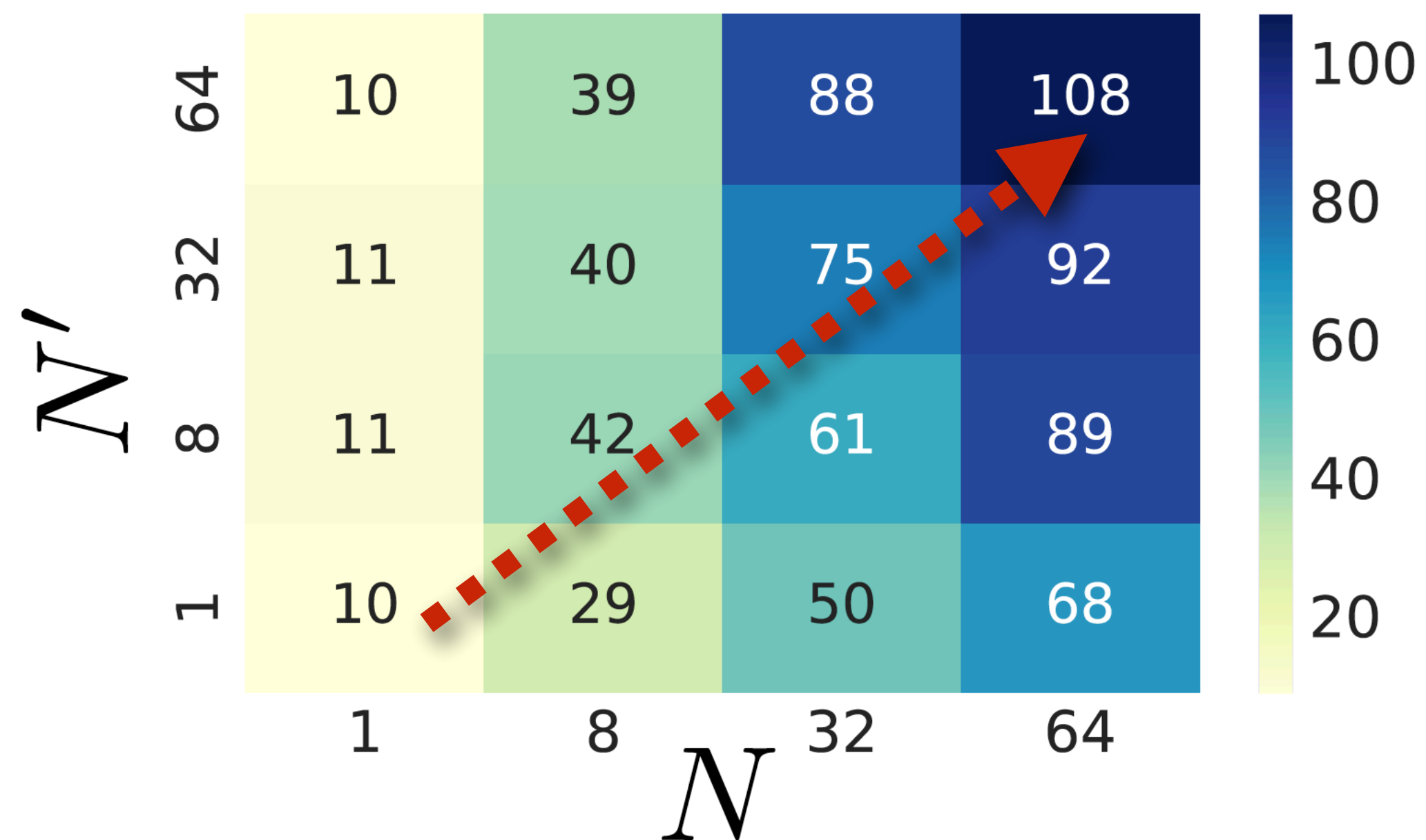
HNS on **first** 10 million frames



IQN - Data efficiency vs. Computation

$$\mathcal{L}_{IQN} = \sum_{\tau=\tau_1}^{\tau_N} \sum_{\tau'=\tau_1}^{\tau_{N'}} \delta_t^{\tau, \tau'} (\tau - \mathbf{I}_{\delta_t^{\tau, \tau'} < 0})$$

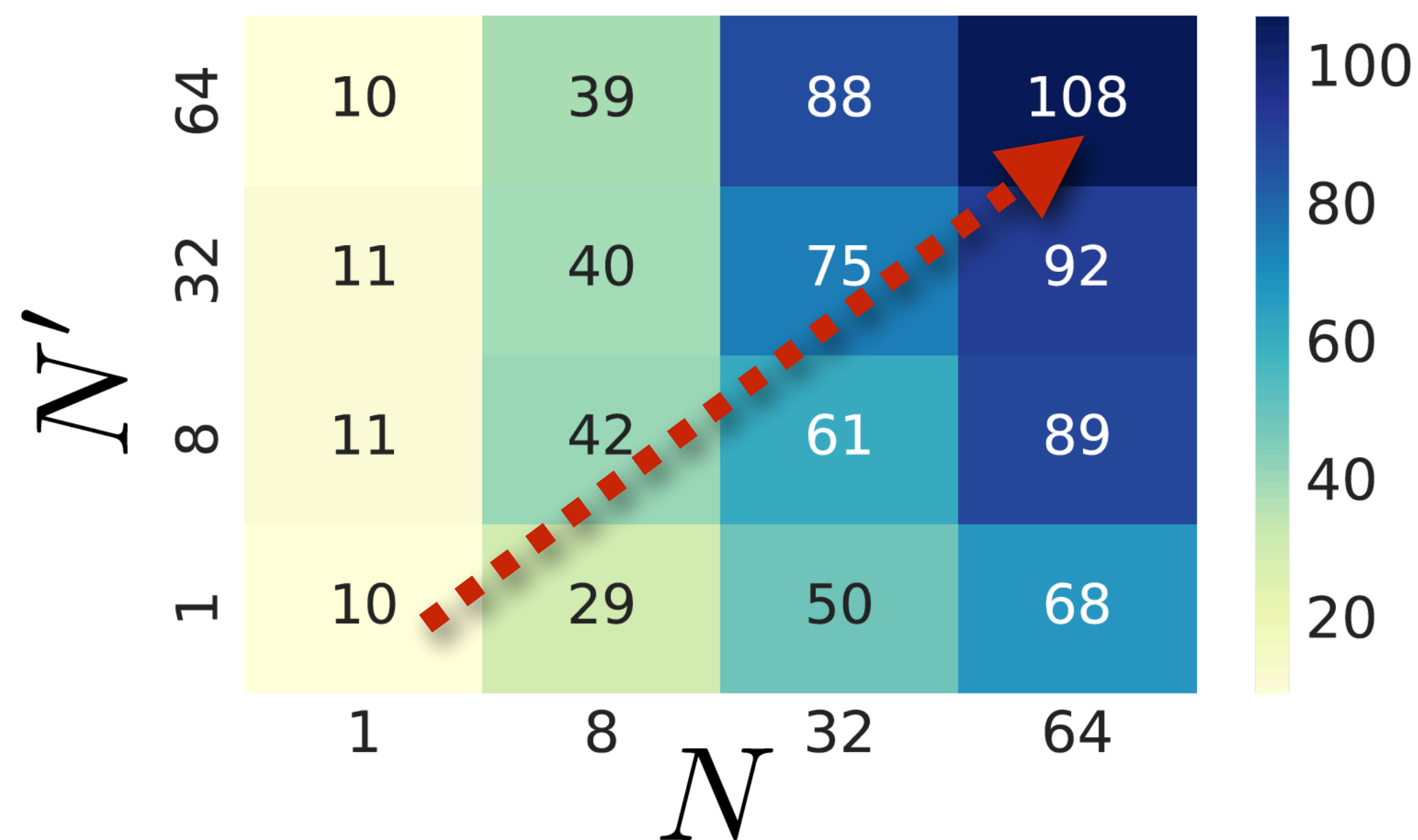
HNS on **first** 10 million frames



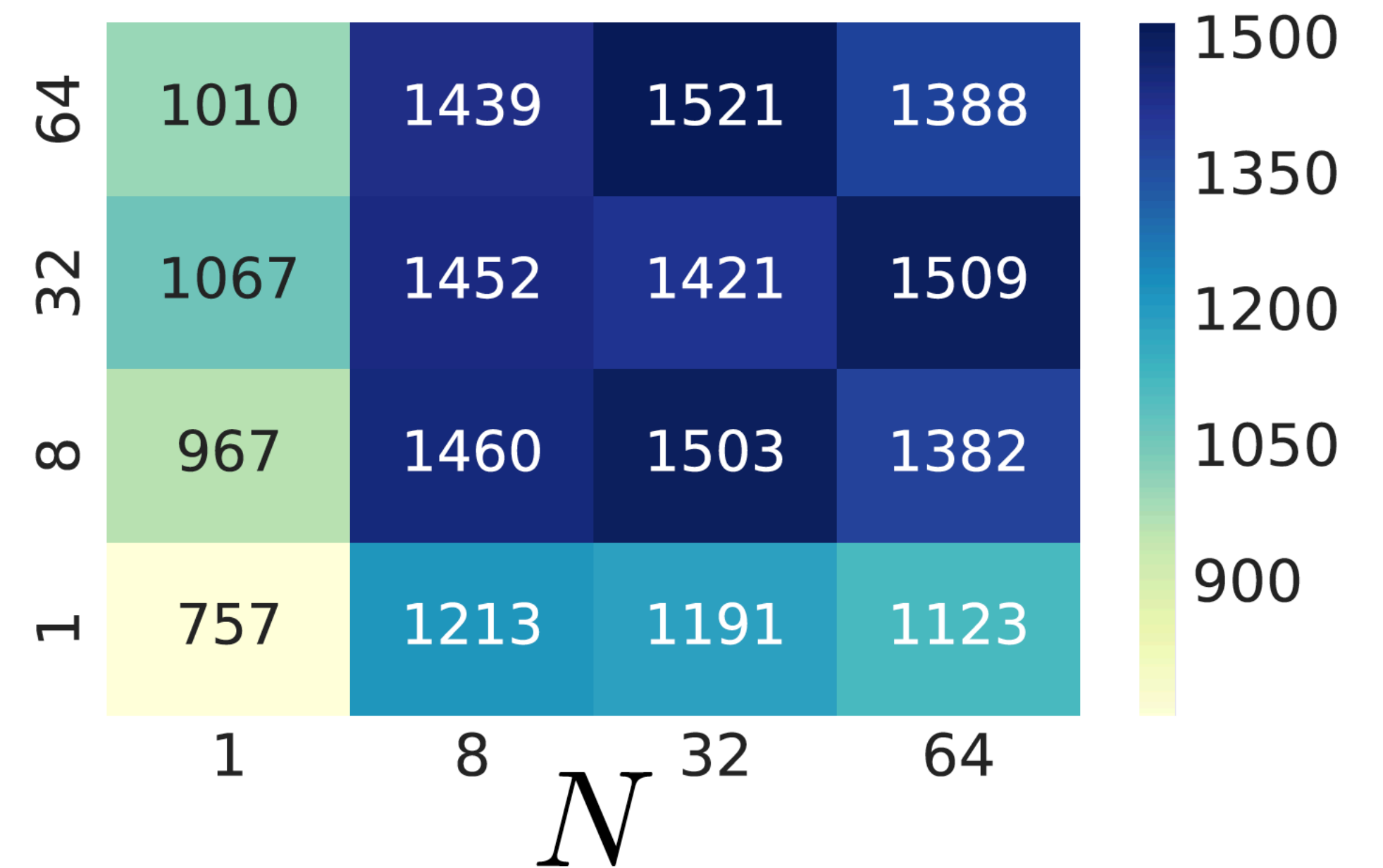
IQN - Data efficiency vs. Computation

$$\mathcal{L}_{IQN} = \sum_{\tau=\tau_1}^{\tau_N} \sum_{\tau'=\tau_1}^{\tau_{N'}} \delta_t^{\tau, \tau'} (\tau - \mathbf{I}_{\delta_t^{\tau, \tau'} < 0})$$

HNS on **first** 10 million frames



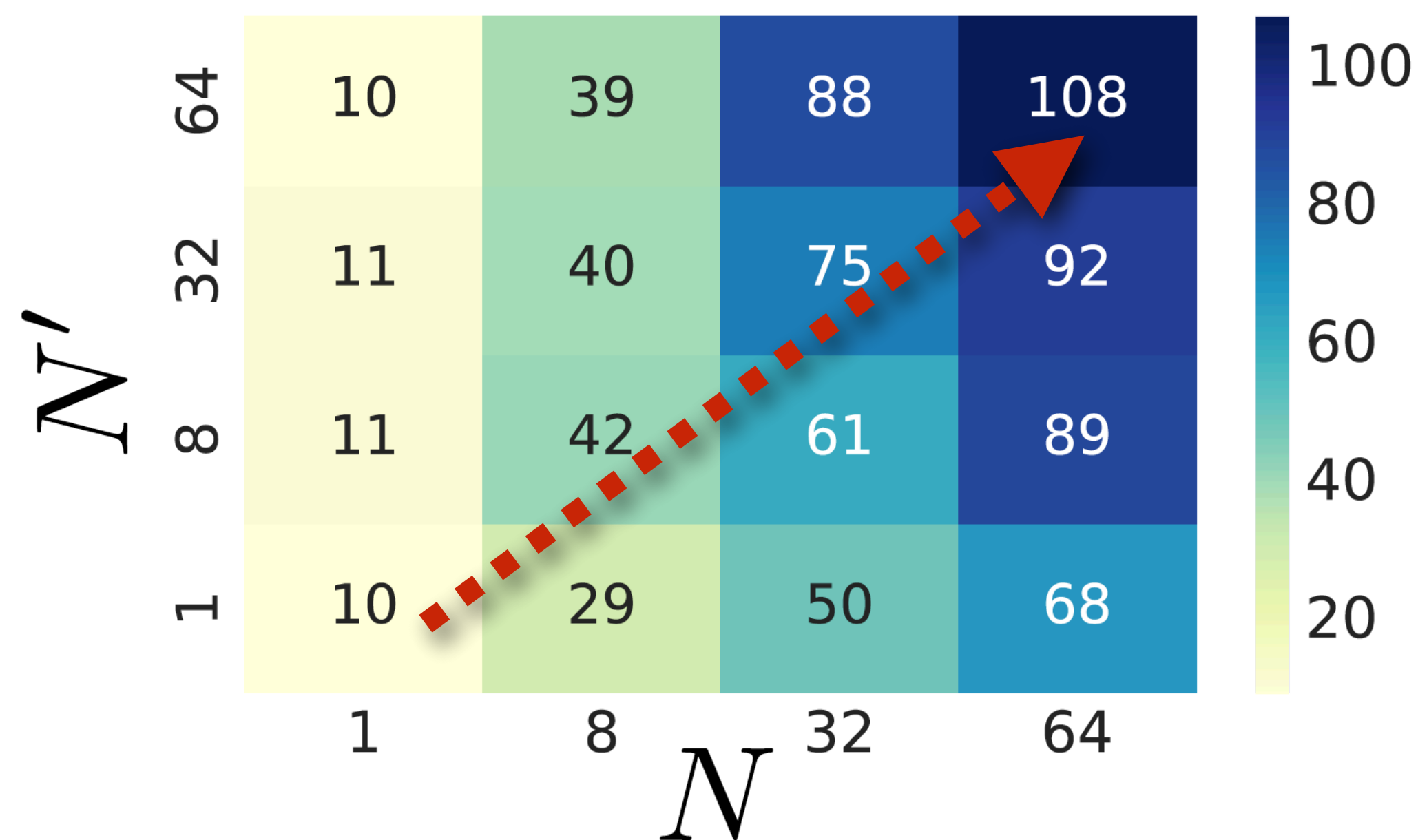
HNS on **last** 10 million frames



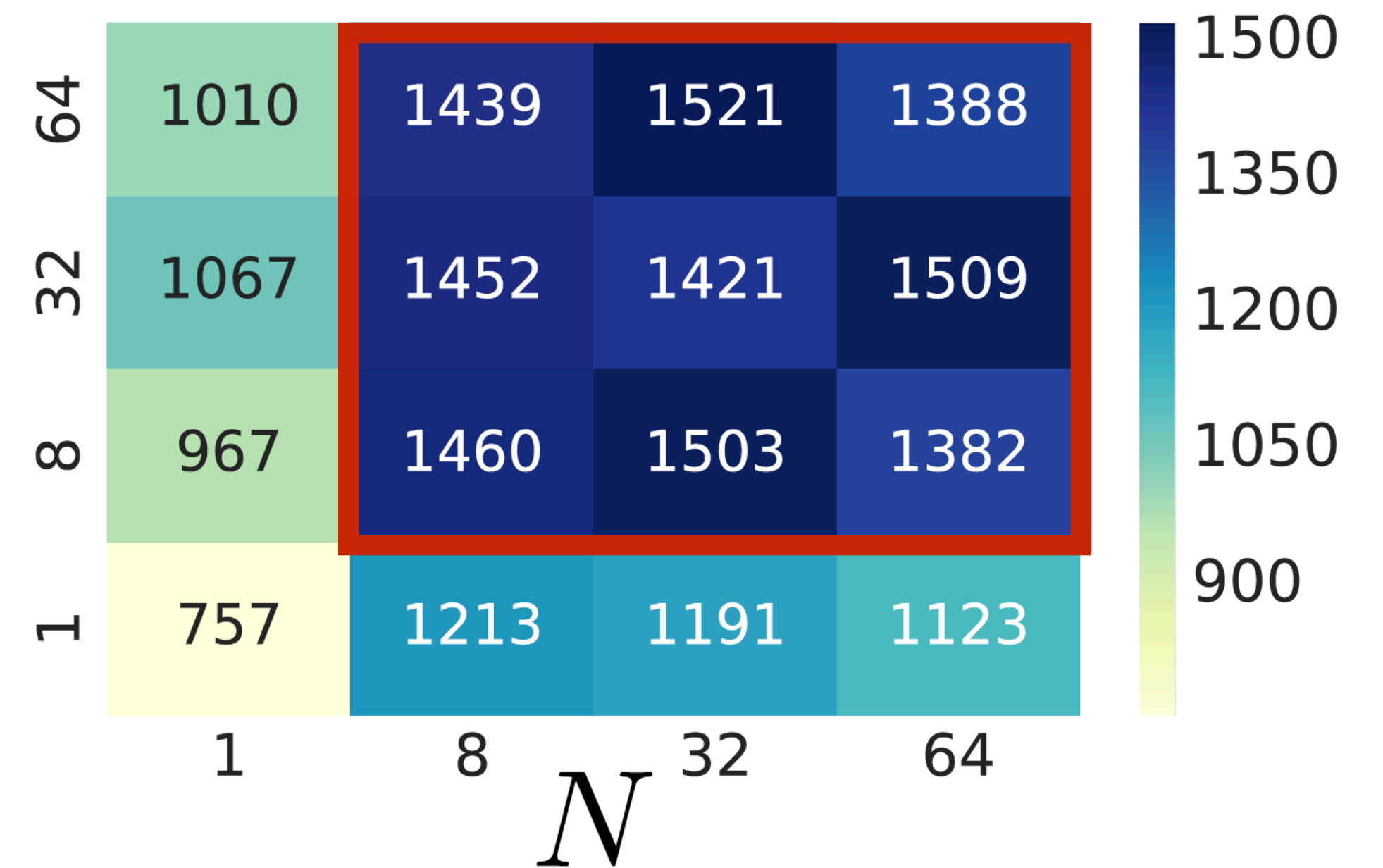
IQN - Data efficiency vs. Computation

$$\mathcal{L}_{IQN} = \sum_{\tau=\tau_1}^{\tau_N} \sum_{\tau'=\tau_1}^{\tau_{N'}} \delta_t^{\tau, \tau'} (\tau - \mathbf{I}_{\delta_t^{\tau, \tau'} < 0})$$

HNS on **first** 10 million frames



HNS on **last** 10 million frames

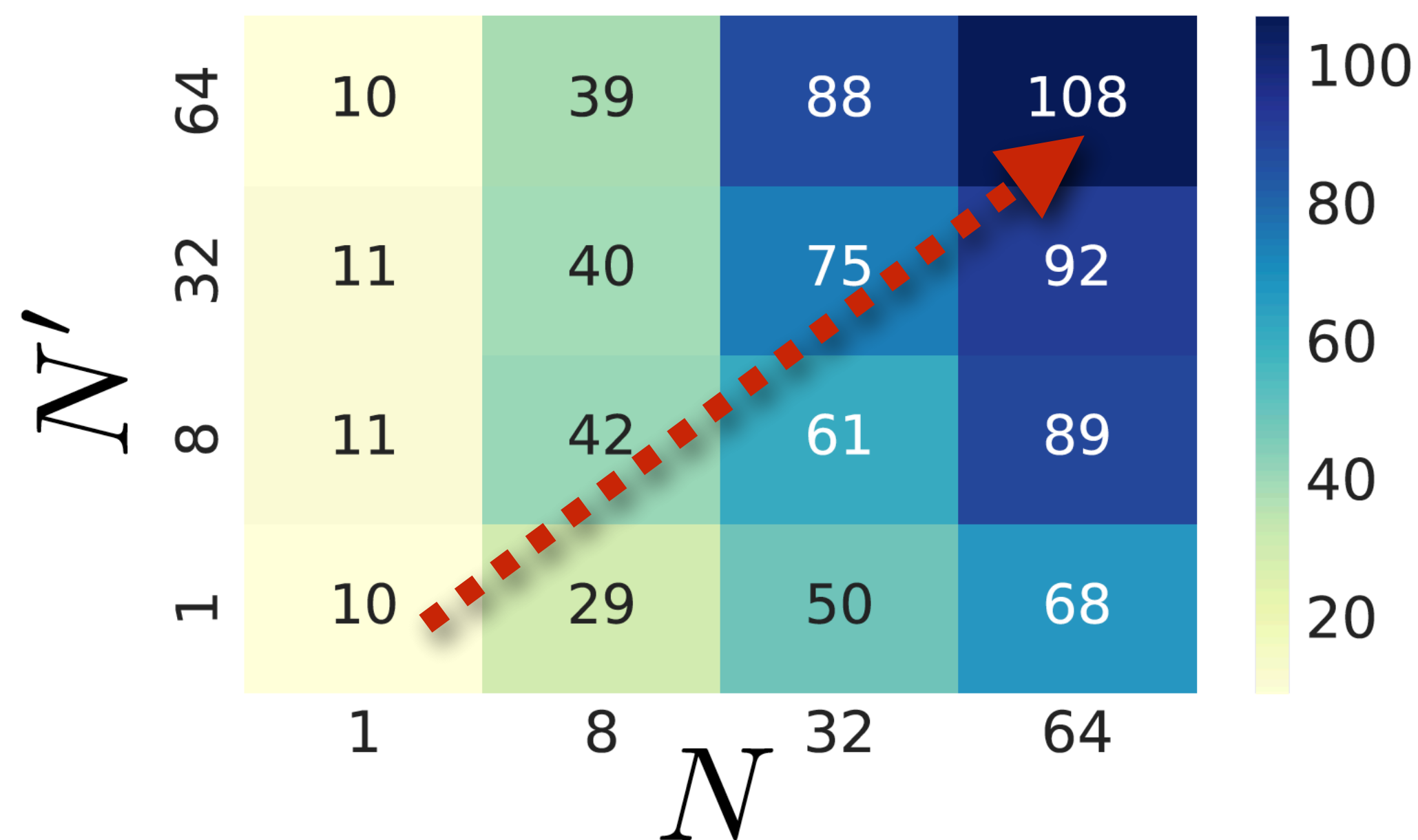


IQN - Data efficiency vs. Computation

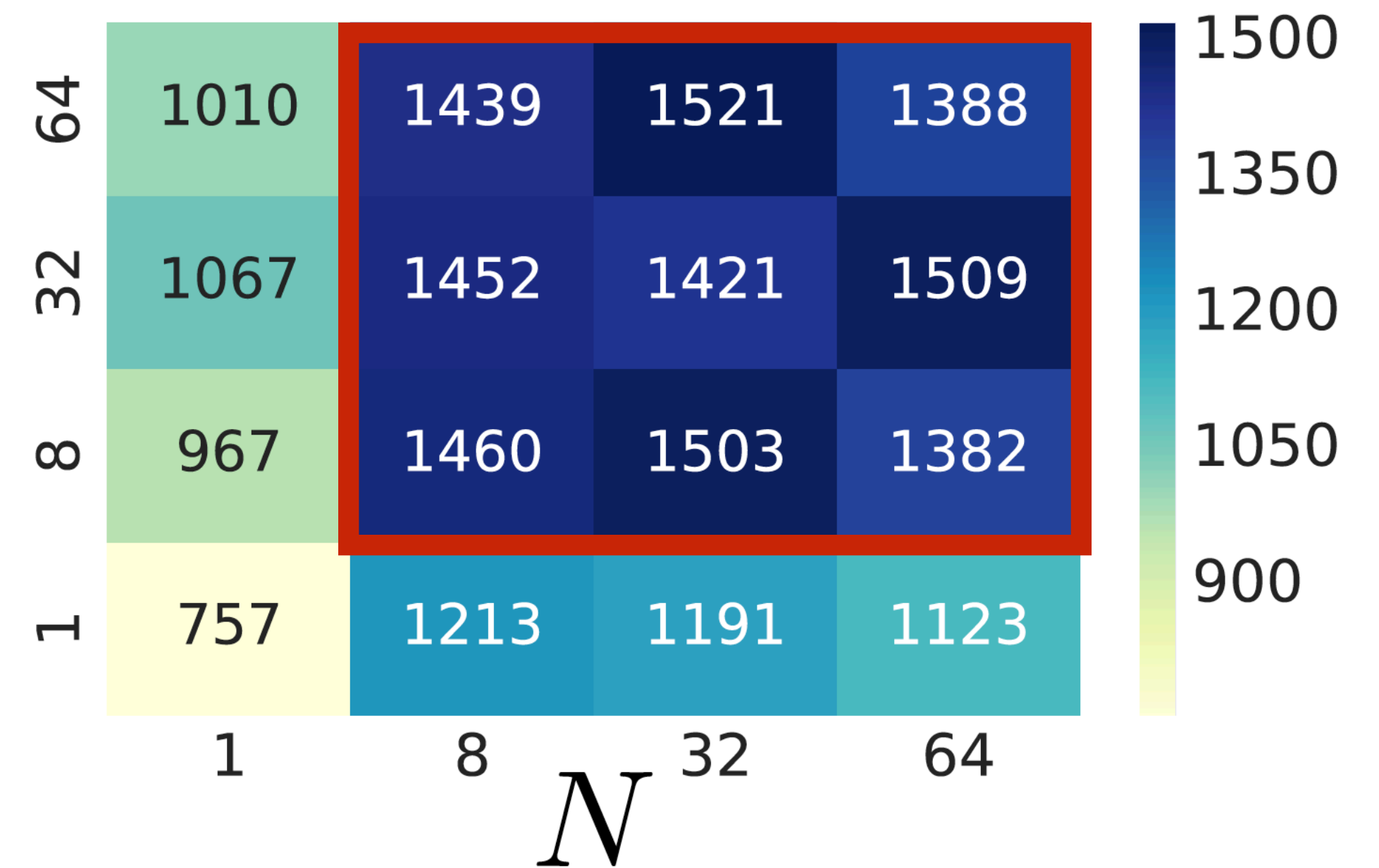
$$\mathcal{L}_{IQN} = \sum_{\tau=\tau_1}^{\tau_N} \sum_{\tau'=\tau_1}^{\tau_{N'}} \delta_t^{\tau, \tau'} (\tau - \mathbf{I}_{\delta_t^{\tau, \tau'} < 0})$$

- Total frames 200 million
- Averaged over 6 Atari games
- DQN (**32, 253**)
- QR-DQN (144, 1243)

HNS on **first** 10 million frames



HNS on **last** 10 million frames

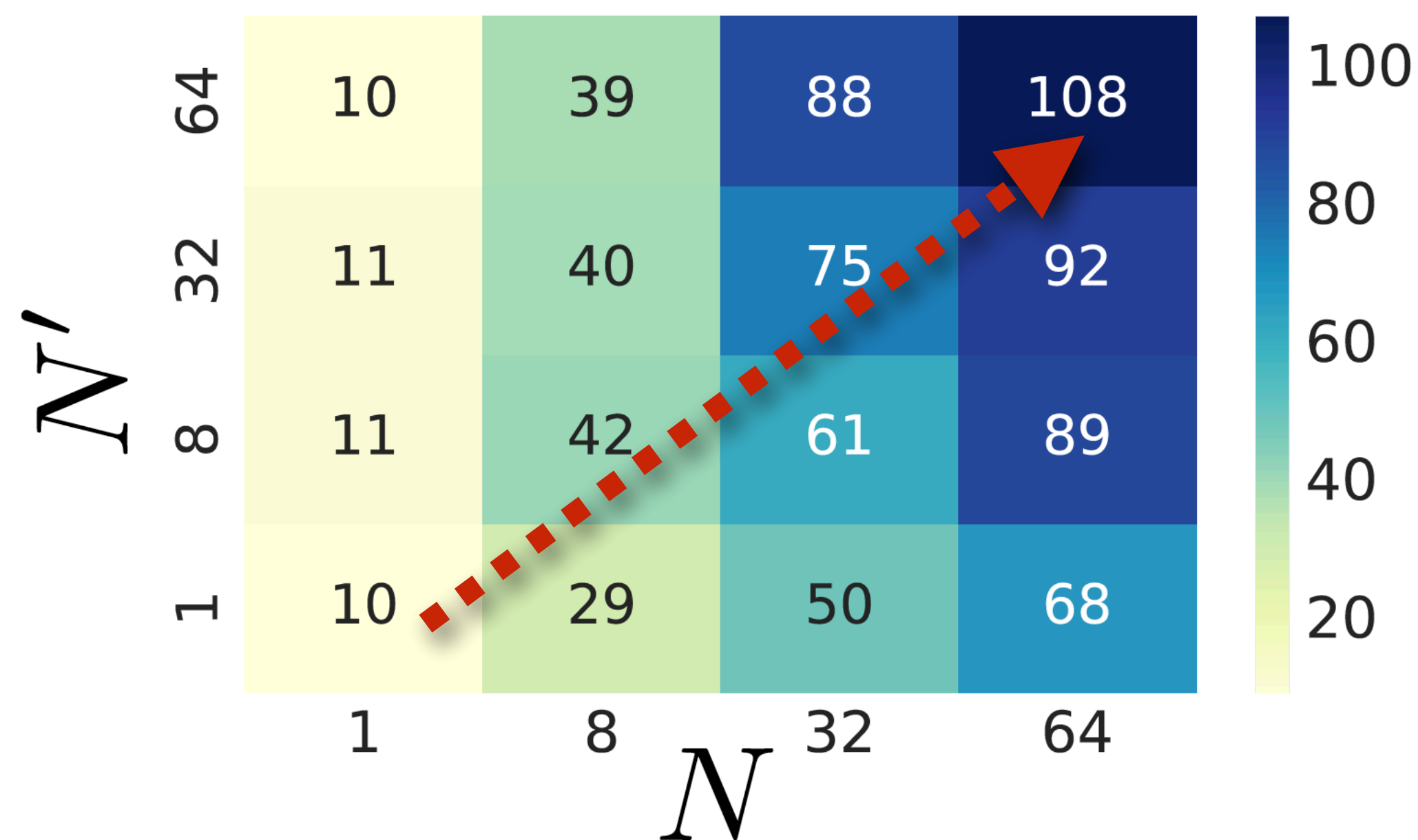


IQN - Data efficiency vs. Computation

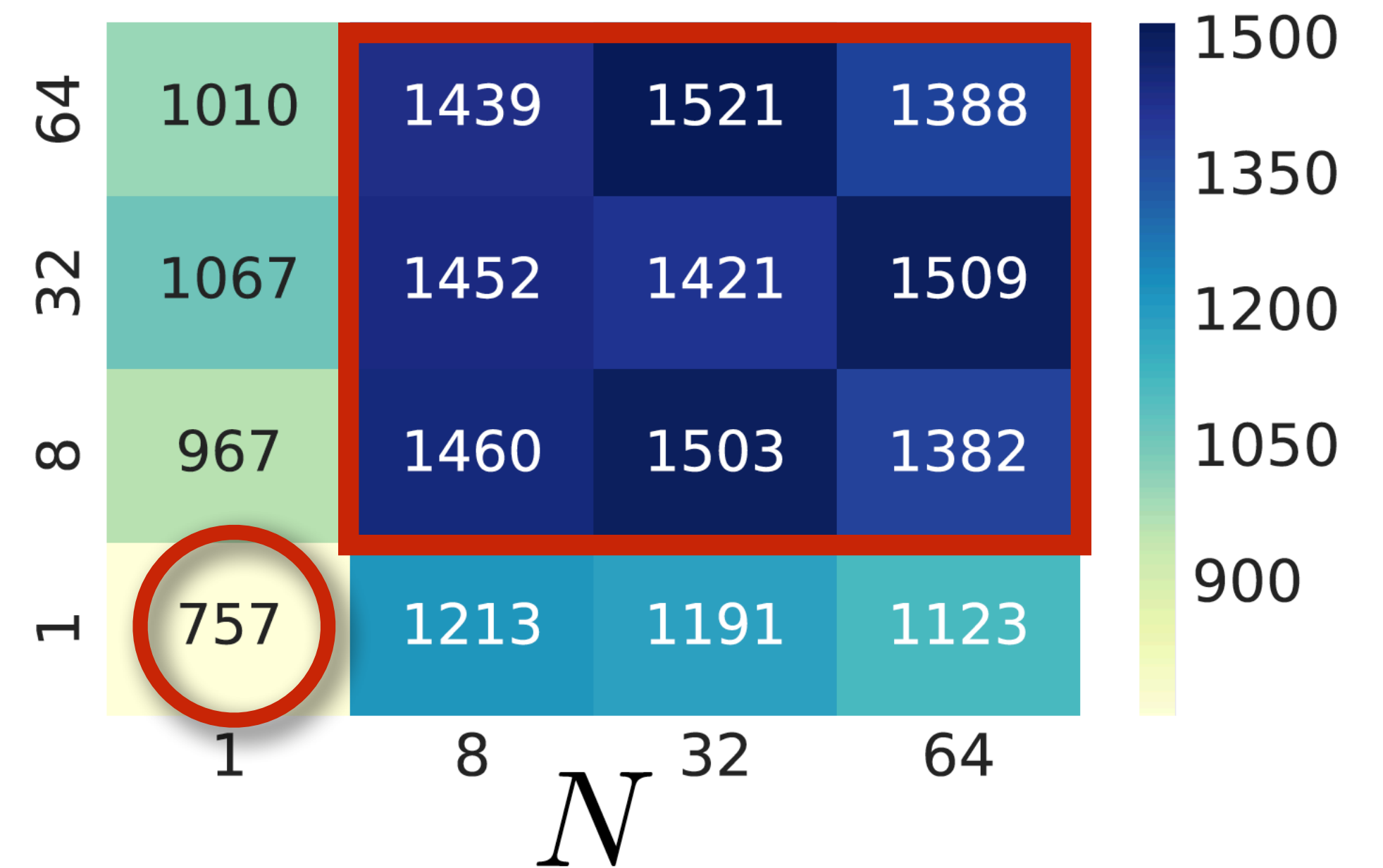
$$\mathcal{L}_{IQN} = \sum_{\tau=\tau_1}^{\tau_N} \sum_{\tau'=\tau_1}^{\tau_{N'}} \delta_t^{\tau, \tau'} (\tau - \mathbf{I}_{\delta_t^{\tau, \tau'} < 0})$$

- Total frames 200 million
- Averaged over 6 Atari games
- DQN (**32, 253**)
- QR-DQN (144, 1243)

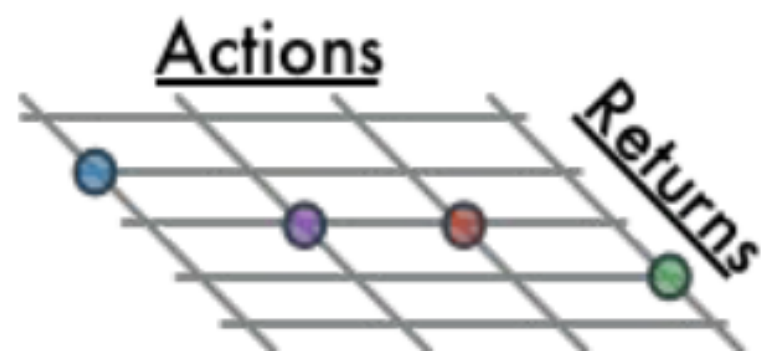
HNS on **first** 10 million frames



HNS on **last** 10 million frames



DQN



$$\{s_t, a_t, s_{t+1}, r_t\}$$

$$a^* = \operatorname{argmax}_a Q^\theta(s_{t+1}, a)$$

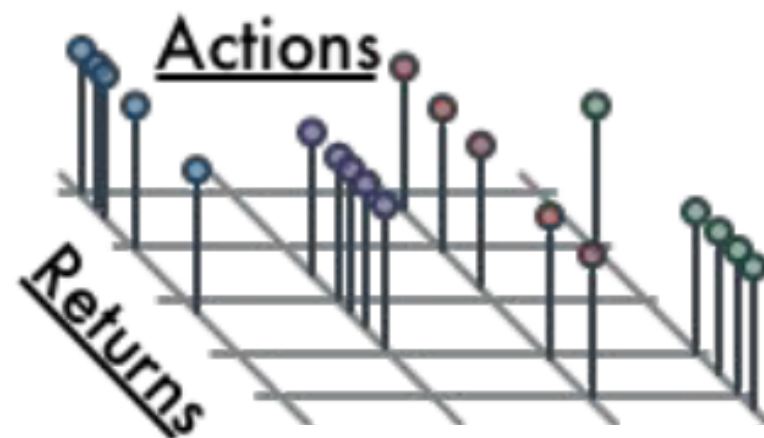
$$q' = Q^\theta(s_{t+1}, a^*)$$

$$q = Q^\theta(s_t, a_t)$$

$$\delta_t = r_t + \gamma q' - q$$

$$\mathcal{L}_{DQN} = \delta_t^2$$

QR-DQN



$$\{s_t, a_t, s_{t+1}, r_t\}$$

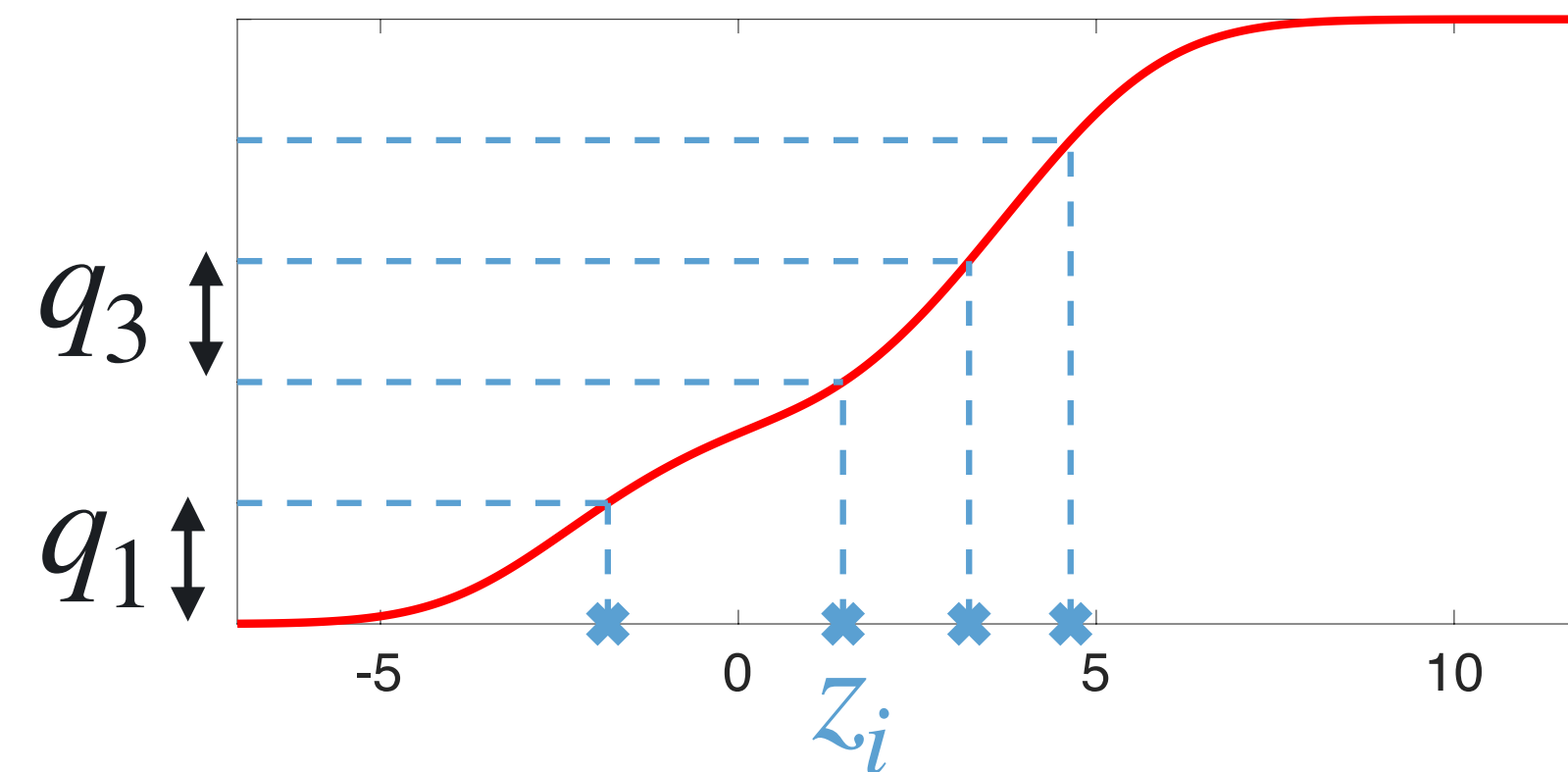
$$a^* = \operatorname{argmax}_a \mathbb{E}[Z_\tau^\theta(s_{t+1}, a)]$$

$$\forall \tau, \tau' \left| \begin{array}{l} z' = Z_{\tau'}^\theta(s_{t+1}, a^*) \\ z = Z_\tau^\theta(s_t, a_t) \end{array} \right.$$

$$z = Z_\tau^\theta(s_t, a_t)$$

$$\delta_t^{\tau, \tau'} = r_t + \gamma z' - z$$

$$\mathcal{L}_{QR-DQN} = ?$$



$$Q(x', a') = \sum_j q_j z_j^\theta(x', a'), \quad \forall a'$$

$$a^* = \operatorname{argmax}_{a'} Q(x', a')$$