

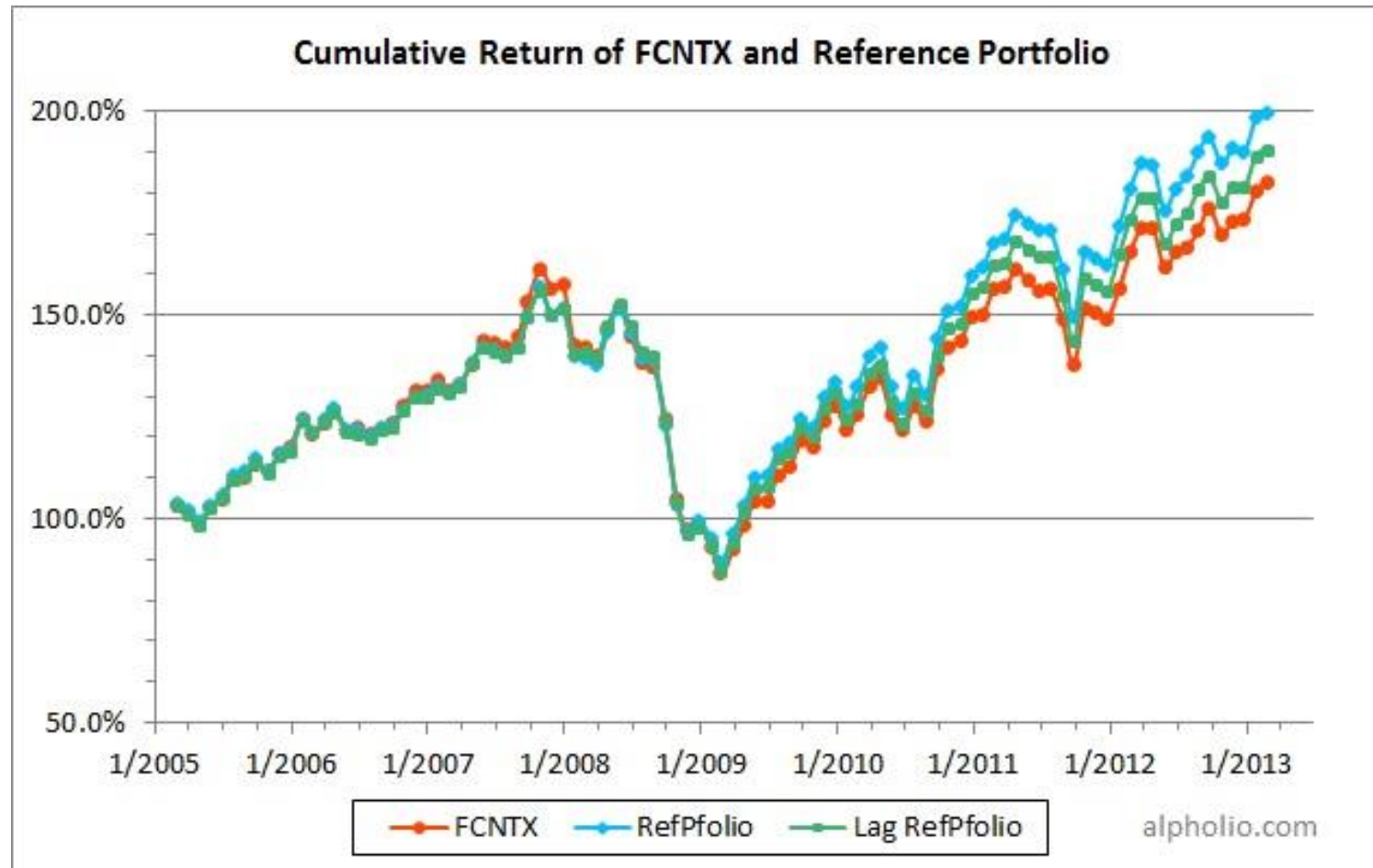
Data-Driven Portfolio Optimization Utilizing Machine Learning

MELINDA HSIEH
RIDER UNIVERSITY

Topics

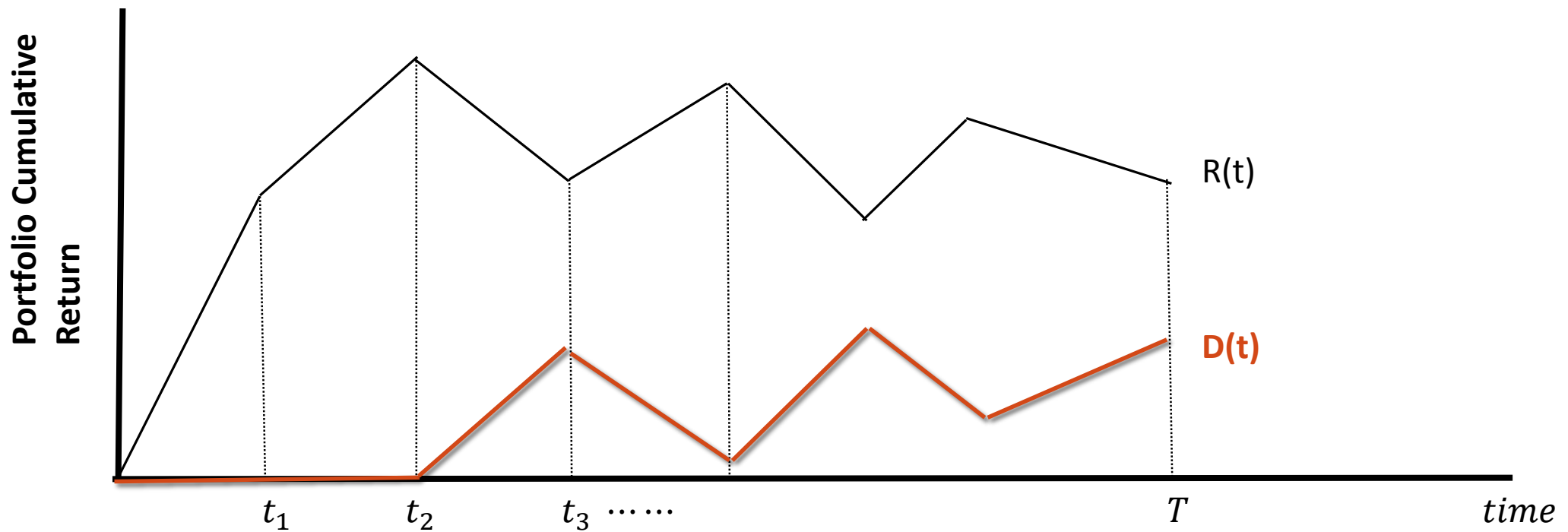
- Data-Driven Portfolio Optimization with Drawdown Constraints
- Prescribe Optimal Portfolio utilizing Machine Learning Methods
- Portfolio Performance Assessment

Portfolio Performance



Drawdown – Peak to Trough decline

$$D(t) = \max_{0 \leq \tau \leq t} R(\tau) - R(t)$$



Portfolio Optimization – with constraint in drawdown

$$\max_{\omega} E[\omega^T \tilde{R}_N]$$

$$\text{s.t.} \quad \max D(\omega, t) \leq C$$

$$C_1 \leq \omega^T \mathbf{1}_N \leq C_2$$

\tilde{R}_N : the annualized cumulative returns of N assets over the time period [0,T]

ω : investment weights of the N assets

Data-Driven Portfolio Optimization – with constraint in drawdown

$$\max_{\omega} \omega^T \hat{R}_N$$

$$\text{s.t.} \quad \max D(\omega, t) \leq C$$

$$C_1 \leq \omega^T \mathbf{1}_N \leq C_2$$

\hat{R}_N : the estimated average annualized cumulative returns of N assets over the time period [0,T]

ω : investment weights of the N assets

Big Data

- An explosion in the availability and accessibility in data
- For example, people frequently browse the internet: shopping, streaming, and searching for key words
- Online activities generate footprints - source of data
- A variety of data sources may be related to stock returns
- Optimal portfolio weights should take into account of these auxiliary variables

Portfolio Optimization with Auxiliary Variables

$$\max_{\omega(\mathbf{x})} \omega^T(\mathbf{x}) E[\tilde{R}_N | \mathbf{x}]$$

$$\text{s.t.} \quad \max D(\omega(\mathbf{x}), t) \leq C$$

$$C_1 \leq \omega(\mathbf{x})^T \mathbf{1}_N \leq C_2$$

$\mathbf{x} \in \mathbf{R}^d$ are auxiliary variables with d-dimensions

Data-Driven Portfolio Optimization with Auxiliary Variables

$$\max_{\omega(x)} \omega^T(x) \hat{R}_N(x)$$

$$\text{s.t.} \quad \max D(\omega(x), t) \leq C$$

$$C_1 \leq \omega(x)^T \mathbf{1}_N \leq C_2$$

$\hat{R}_N(x)$: the estimated conditional mean of annualized cumulative returns of N assets

Searching Optimal Portfolio Weights – Linear Programming Problem

- The portfolio optimization can be represented as a linear programming problem
- Compute average accumulative returns to obtain \hat{R}_N
- The derived optimal portfolio weight is data-driven, depending on the estimated \hat{R}_N

Machine Learning

Apply machine learning methods to estimate the conditional means of returns $\hat{\mathbf{R}}_N(\mathbf{x})$ and derive the optimal portfolio weight $\boldsymbol{\omega}(\mathbf{x})$:

1. Use machine learning methods to classify asset returns into several groups based on the values of auxiliary variables
2. For each group, compute $\hat{\mathbf{R}}_N(\mathbf{x})$
3. Solve the linear programming problem to find the optimal portfolio weight $\boldsymbol{\omega}(\mathbf{x})$

Prescribe Investment Weights

Train data

Machine Learning

$$\begin{pmatrix} \hat{R}_N^{(1)}; x_1^{(1)}, x_2^{(1)}, \dots, x_d^{(1)} \\ \hat{R}_N^{(2)}; x_1^{(2)}, x_2^{(2)}, \dots, x_d^{(2)} \\ \hat{R}_N^{(3)}; x_1^{(3)}, x_2^{(3)}, \dots, x_d^{(3)} \\ \vdots \\ \hat{R}_N^{(T)}; x_1^{(T)}, x_2^{(T)}, \dots, x_d^{(T)} \end{pmatrix}$$

$$\begin{pmatrix} \hat{R}_N^{(1)}; x_1^{(1)}, x_2^{(1)}, \dots, x_d^{(1)} \\ \hat{R}_N^{(3)}; x_1^{(3)}, x_2^{(3)}, \dots, x_d^{(3)} \\ \hat{R}_N^{(4)}; x_1^{(4)}, x_2^{(4)}, \dots, x_d^{(4)} \end{pmatrix}$$

$$\omega(x^{(1)}, x^{(3)}, x^{(4)})$$

$$\begin{pmatrix} \hat{R}_N^{(5)}; x_1^{(5)}, x_2^{(5)}, \dots, x_d^{(5)} \\ \hat{R}_N^{(10)}; x_1^{(10)}, x_2^{(10)}, \dots, x_d^{(10)} \end{pmatrix}$$

$$\omega(x^{(5)}, x^{(10)})$$

$$\begin{pmatrix} \hat{R}_N^{(2)}; x_1^{(2)}, x_2^{(2)}, \dots, x_d^{(2)} \\ \hat{R}_N^{(6)}; x_1^{(6)}, x_2^{(6)}, \dots, x_d^{(6)} \\ \hat{R}_N^{(20)}; x_1^{(20)}, x_2^{(20)}, \dots, x_d^{(20)} \\ \hat{R}_N^{(T)}; x_1^{(T)}, x_2^{(T)}, \dots, x_d^{(T)} \end{pmatrix}$$

$$\omega(x^{(2)}, x^{(6)}, x^{(20)}, x^{(T)})$$

$\hat{R}_N^{(i)}$: cumulative returns observed at time i

KNN - K Nearest Neighborhood

$$\hat{\omega}_{KNN}(\mathbf{x}_0) \in \max_{\mathbf{x}^{(i)} \in \mathfrak{N}_k(\mathbf{x}_0)} \omega^T \hat{R}_N(\mathbf{x}^{(i)})$$

$$\text{s.t.} \quad \max D(\omega, t) \leq C$$

$$C_1 \leq \omega^T \mathbf{1}_N \leq C_2$$

where $\mathfrak{N}_K(\mathbf{x}_0) = \{i = 1, \dots, T : \sum_{j=1}^T I[\|\mathbf{x}_0 - \mathbf{x}^{(i)}\| \geq \|\mathbf{x}_0 - \mathbf{x}^{(j)}\|] \leq k\}$ is the neighborhood of k data points that are closest to the point \mathbf{x}_0

Ctree - Conditional Inference Tree (Hothorn et al (2006))

$$\hat{\omega}_{\text{Ctree}}(\mathbf{x}) \in \max_{i: R(\mathbf{x}^{(i)}) \in R(\mathbf{x})} \omega^T \hat{R}_N(\mathbf{x}^{(i)})$$

$$\text{s.t.} \quad \max D(\omega, t) \leq C$$

$$C_1 \leq \omega^T \mathbf{1}_N \leq C_2$$

where $R(\mathbf{x})$ is the splitting rule implied by a regression tree trained based on the training data

- Ctree applies the inference test to determine if a possible split is significant

Random Forest

$$\hat{\omega}_{RF}^k(\mathbf{x}) \in \max_{i: R^k(\mathbf{x}^{(i)}) \in R^k(\mathbf{x})} \omega^T \hat{R}_N(\mathbf{x}^{(i)})$$

$$\text{s.t.} \quad \max D(\omega, t) \leq C$$

$$C_1 \leq \omega^T \mathbf{1}_N \leq C_2$$

$$\hat{\omega}_{RF}(\mathbf{x}) = \frac{1}{m} \sum_{k=1}^m \hat{\omega}_{RF}^k(\mathbf{x})$$

Simulation

- Let X_1, X_2, X_3 be the market factors that are related to the returns of the underlying assets in the portfolio.
- Assume $X(t) = \{X_1(t), X_2(t), X_3(t)\}$ follows a multivariate ARMA(2,2) process
- Generate 12 time series of returns based on the three market factors

$$R_N(t) = A^T \left(X(t) + \frac{\delta}{4} \right) + (B^T(X(t)))\eta$$

where δ and η are independent Gaussian noises

Portfolio Performance – Validation

For each realization of the simulated returns

- split the series into in-sample (training) and out-of-sample (validation) data.
- Trained the in-sample data based on the machine learning algorithm and prescribe optimal portfolio weights
- Apply the prescribed weights to the out-of-sample (validation) returns and compute the portfolio returns

Performance Metrics

1. Average Cumulative Return

$$\frac{1}{n} \sum_{t=1}^n R_t^v$$

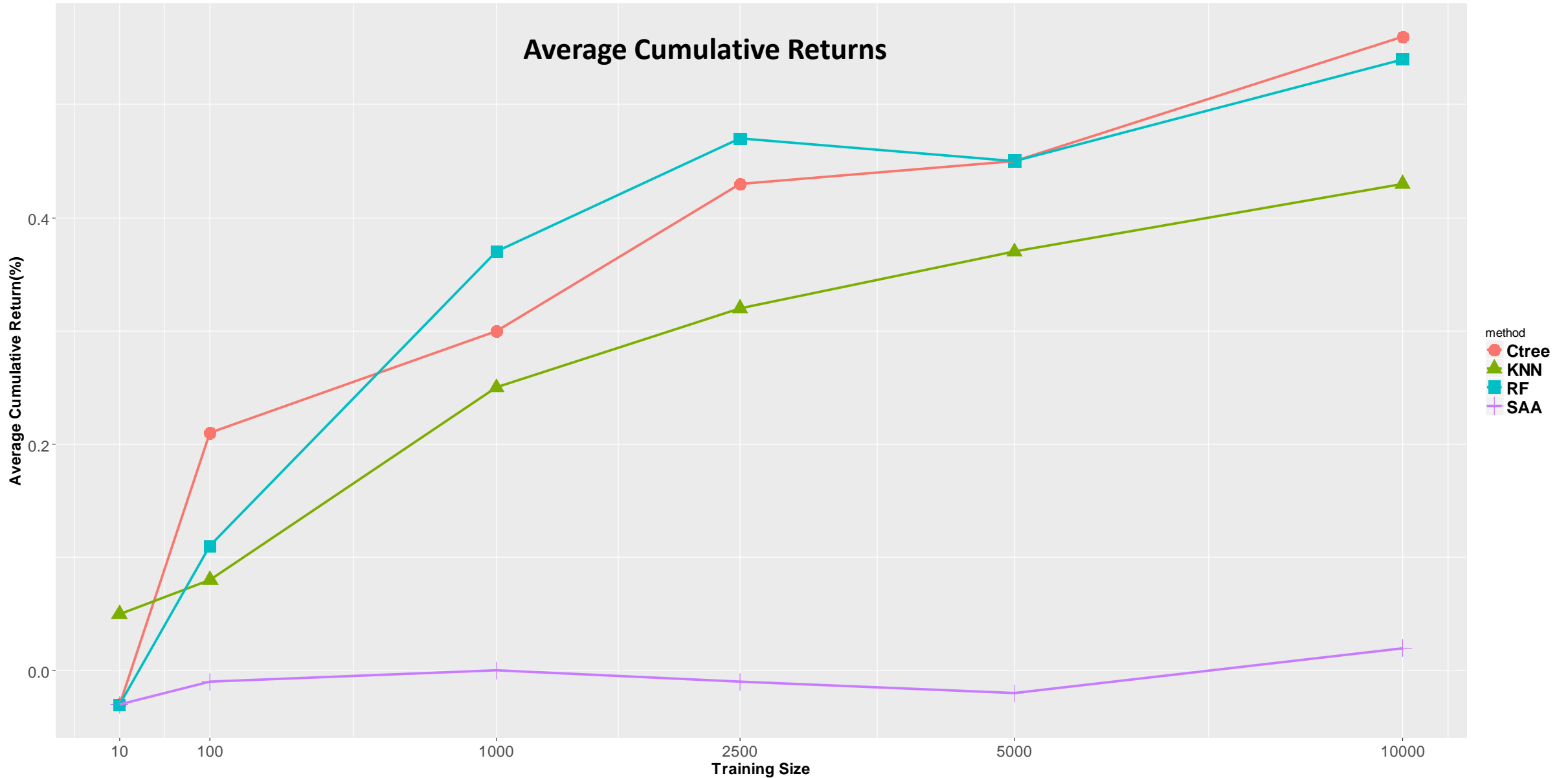
2. Maximum Drawdown

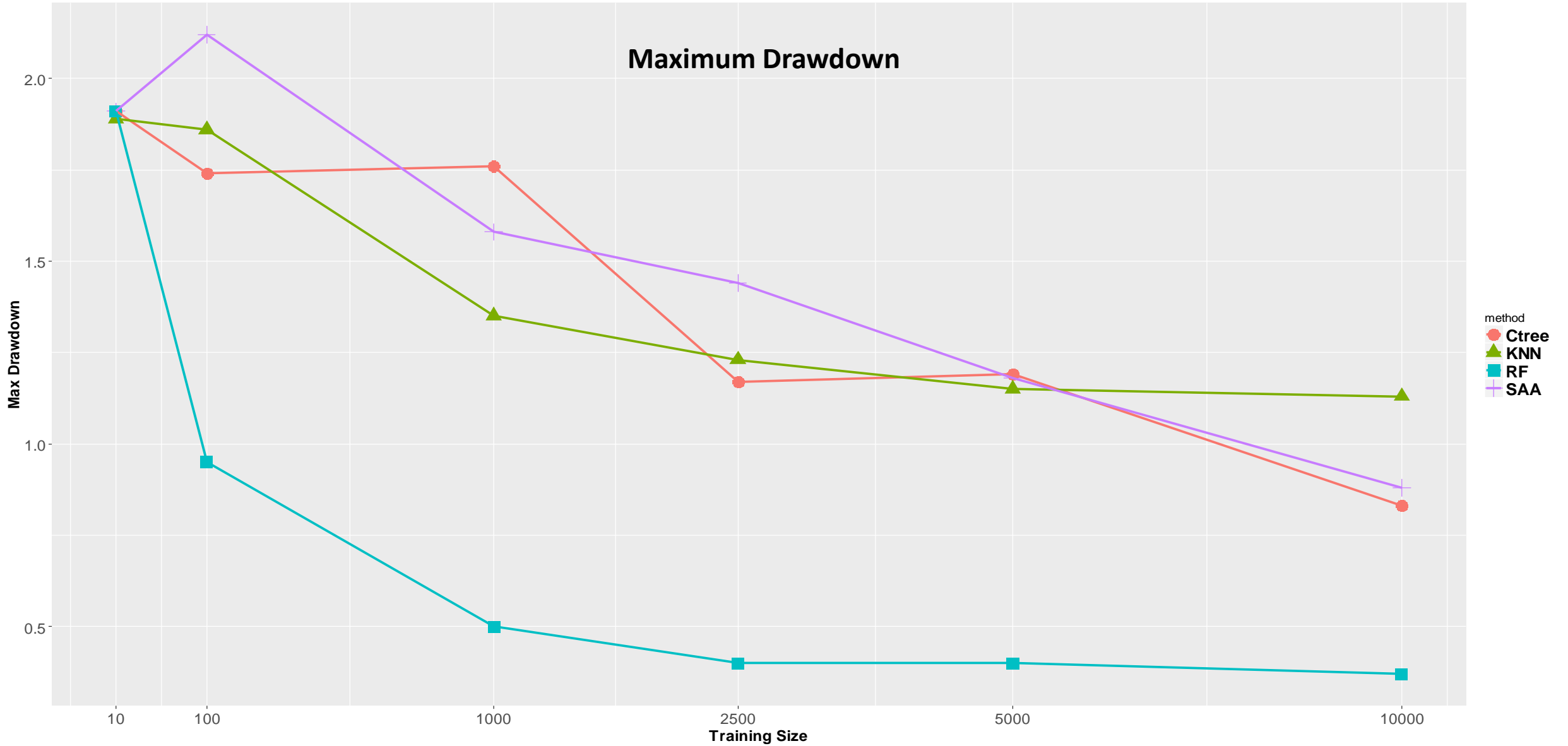
$$\max_{0 \leq t \leq n} \left\{ \max_{0 \leq \tau \leq t} R(\tau) - R(t) \right\}$$

3. Reward Risk

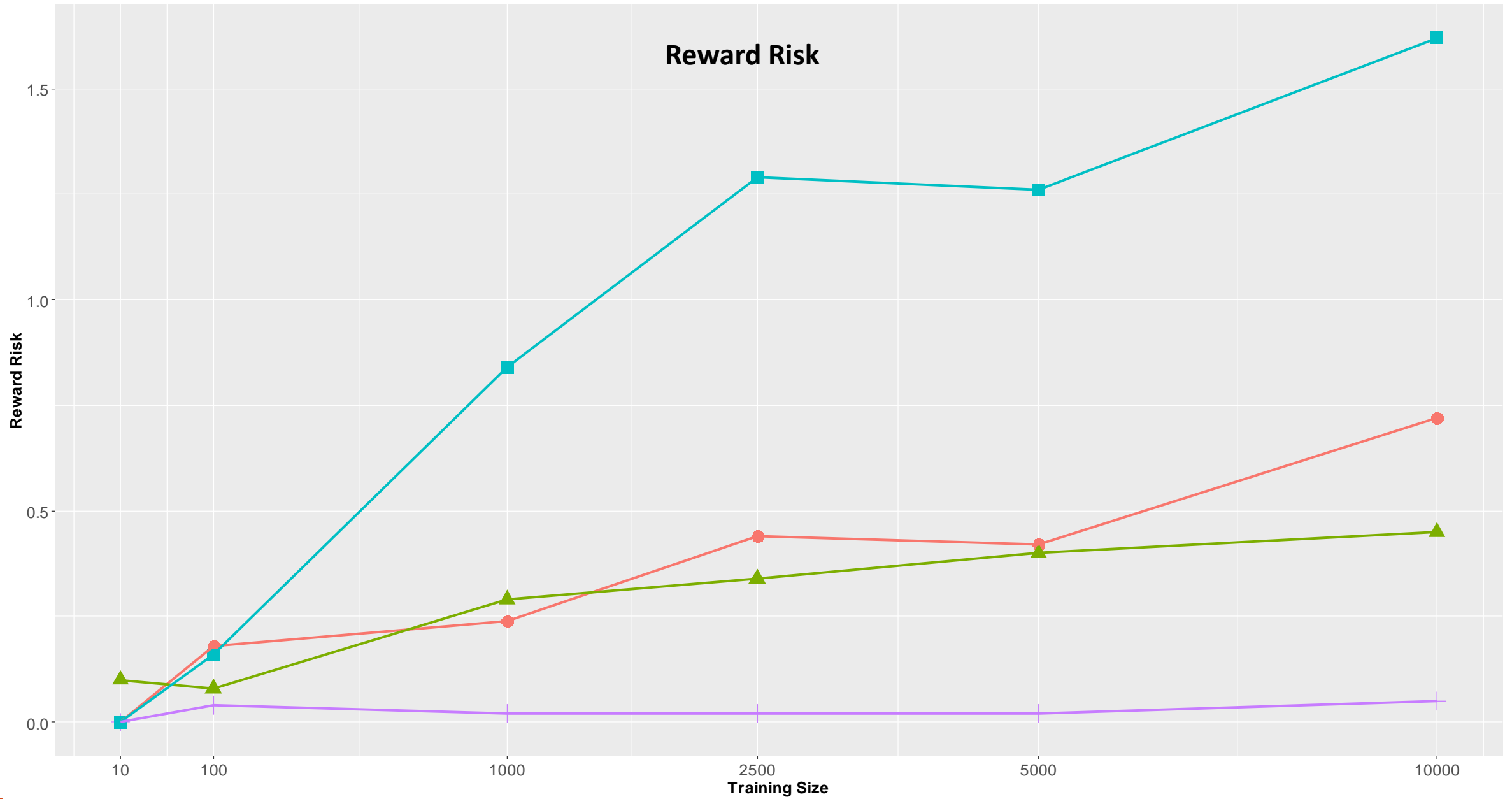
$$\frac{\frac{1}{n} \sum_{t=1}^n R_t^v}{\max_{0 \leq t \leq n} \left\{ \max_{0 \leq \tau \leq t} R(\tau) - R(t) \right\}}$$

Average Cumulative Returns





Reward Risk



Summary

- Machine Learning method improves the out-of sample performance of the data-driven optimal portfolio
- The performance improves as the size of training data increases
- Tree-based approaches such as random forest and Ctree outperform the SAA method which does not incorporate the inputs from auxiliary variables.

Future Work

- how the tuning parameters affect optimal portfolio weights and portfolio performance?
- how the correlations of the auxiliary variables affect optimal portfolio weights and portfolio performance?

Appendix

Searching Optimal Portfolio Weights – Linear Programming Problem

$$\max_{\omega} \omega^T \hat{R}_{T \times N}$$

$$\text{s.t. } u_k - \omega^T \hat{R}_{k \times N} \leq C, 1 \leq k \leq N$$

$$u_k \geq \omega^T \hat{R}_{k \times N}, 1 \leq k \leq N$$

$$u_k \geq u_{k-1}, 1 \leq k \leq N$$

$$u_0 = 0$$

$$C_1 \leq \omega^T \mathbf{1}_N \leq C_2$$

where u_k are auxiliary variables, $1 \leq k \leq N$