

Gap between theory and practice

- For linear model $\mathbf{y} = \mathbf{X}\mathbf{b}^* + \boldsymbol{\varepsilon}$, statistical theories (e.g., error bounds, asymptotic distributions, risk characterization) focus on the optimizer $\hat{\mathbf{b}}$:

$$\hat{\mathbf{b}} \in \arg \min_{\mathbf{b} \in \mathbb{R}^p} \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2 + g(\mathbf{b})$$

- In practice, $\hat{\mathbf{b}}$ cannot be solved exactly; iterative algorithms are used to produce iterates $\hat{\mathbf{b}}^1, \hat{\mathbf{b}}^2, \dots, \hat{\mathbf{b}}^T$ (stop after T iterations).

- $\hat{\mathbf{b}}^T$ can be far from $\hat{\mathbf{b}}$, and the theories about $\hat{\mathbf{b}}$ do not apply to $\hat{\mathbf{b}}^T$.

- There is no guarantee that $\hat{\mathbf{b}}^t$ will get closer to \mathbf{b}^* as t increases.

Q1: How can we quantify the predictive performance of $\hat{\mathbf{b}}^t$ at each iteration?

Q2: How does the performance of $\hat{\mathbf{b}}^t$ depend on the previous iterates?

Q3: How can we perform statistical inference on \mathbf{b}^* using the iterate $\hat{\mathbf{b}}^t$?

Estimation Target

Consider an algorithm of the following form:

$$\hat{\mathbf{b}}^t = \mathbf{g}_t(\hat{\mathbf{b}}^{t-1}, \hat{\mathbf{b}}^{t-2}, \dots, \hat{\mathbf{b}}^2, \hat{\mathbf{b}}^1, \mathbf{v}^{t-1}, \dots, \mathbf{v}^2, \mathbf{v}^1), \quad (1)$$

where $\mathbf{v}^t = \frac{1}{n} \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\hat{\mathbf{b}}^t)$.

Table 1. Examples of several algorithms

GD	$\hat{\mathbf{b}}^t = \hat{\mathbf{b}}^{t-1} + \eta \mathbf{v}^t$
AGD	$\hat{\mathbf{b}}^t = (1 - w_{t-1})(\hat{\mathbf{b}}^{t-1} + \eta \mathbf{v}^{t-1}) + w_{t-1}(\hat{\mathbf{b}}^{t-2} + \eta \mathbf{v}^{t-2})$
ISTA	$\hat{\mathbf{b}}^t = \text{soft}_{\eta\lambda}(\hat{\mathbf{b}}^{t-1} + \eta \mathbf{v}^t)$
FISTA	$\hat{\mathbf{b}}^t = \text{soft}_{\eta\lambda}((1 - w_{t-1})(\hat{\mathbf{b}}^{t-1} + \eta \mathbf{v}^{t-1}) + w_{t-1}(\hat{\mathbf{b}}^{t-2} + \eta \mathbf{v}^{t-2}))$

Estimation target: The generalization error r_t for each $\hat{\mathbf{b}}^t$:

$$r_t \stackrel{\text{def}}{=} \mathbb{E}[(y_{\text{new}} - \mathbf{x}_{\text{new}}^\top \hat{\mathbf{b}}^t)^2 \mid (\mathbf{X}, \mathbf{y})] = \|\Sigma^{1/2}(\hat{\mathbf{b}}^t - \mathbf{b}^*)\|^2 + \sigma^2.$$

Contributions

- Propose a novel estimator of r_t :

$$\underbrace{\|\Sigma^{1/2}(\hat{\mathbf{b}}^t - \mathbf{b}^*)\|^2 + \sigma^2}_{r_t} \approx \frac{1}{n} \underbrace{\left\| \sum_{s=1}^t \hat{w}_{t,s} (\mathbf{y} - \mathbf{X}\hat{\mathbf{b}}^s) \right\|^2}_{\hat{r}_t}$$

- The \hat{r}_t depends on the **weighted residual vector** of all previous iterates.
- The weights $\hat{w}_{t,s}$ can be easily computed using observational quantities.

- Introduce the debiased estimate for the component of each iterate $\hat{\mathbf{b}}^t$:

$$\hat{\mathbf{b}}_j^{t,\text{debias}} \stackrel{\text{def}}{=} \underbrace{\hat{\mathbf{b}}_j^t}_{\text{iterate}} + \underbrace{\sum_{s=1}^t \hat{w}_{t,s} (\mathbf{y} - \mathbf{X}\hat{\mathbf{b}}^s)^\top \mathbf{X}\Sigma^{-1} \mathbf{e}_j / n}_{\text{bias correction}}$$

- Establish asymptotic normality results for \mathbf{b}_j^* using the debiased estimate:

$$\frac{\sqrt{n}(\hat{\mathbf{b}}_j^{t,\text{debias}} - \mathbf{b}_j^*)}{\|\Sigma^{-1/2} \mathbf{e}_j\| \sqrt{\hat{r}_t}} \xrightarrow{d} N(0, 1).$$

Assumptions

A1: The design matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$ has i.i.d. rows from $N(\mathbf{0}, \Sigma)$ with an invertible Σ .

A2: The noise $\boldsymbol{\varepsilon}$ is independent of \mathbf{X} and has i.i.d. entries from $N(0, \sigma^2)$.

A3: The asymptotic regime is $n \rightarrow \infty$ and $p \rightarrow \infty$ with $\frac{p}{n} \leq \gamma \in (0, \infty)$.

A4: The algorithm starts with $\hat{\mathbf{b}}^1 = \mathbf{0}_p$ and \mathbf{g}_t in Equation (1) is ζ -Lipschitz with $\mathbf{g}_t(0) = 0$.

Theorem 1: Estimation of r_t

Assume conditions A1 - A4 hold, then for any $t \in [T]$:

$$\mathbb{E}[|\hat{r}_t - r_t|] \leq \frac{1}{\sqrt{n}} C(\zeta, T, \gamma, \kappa) \text{var}(y_1). \quad (2)$$

Let $\hat{t} \stackrel{\text{def}}{=} \arg \min_{t \in [T]} \hat{r}_t$. For any $c \in (0, 1/2)$, we have:

$$\mathbb{P}\left(\|\Sigma^{1/2}(\hat{\mathbf{b}}^{\hat{t}} - \mathbf{b}^*)\|^2 \leq \min_{s \in [T]} \|\Sigma^{1/2}(\hat{\mathbf{b}}^s - \mathbf{b}^*)\|^2 + \frac{\text{var}(y_1)}{n^{1/2-c}}\right) \geq 1 - \frac{C(\zeta, \gamma, T, \kappa)}{n^c} \rightarrow 1.$$

- The proposed estimator \hat{r}_t is \sqrt{n} -consistent for r_t .

- Minimizing \hat{r}_t can lead to an optimal stopping time with negligible error.

Theorem 2: Inference for \mathbf{b}_j^*

Under Assumptions A1 - A4. There exists a set $J_{n,p} \subset [p]$ with $|J_{n,p}| \geq p - \log p$ such that

$$\frac{\sqrt{n}(\hat{\mathbf{b}}_j^{t,\text{debias}} - \mathbf{b}_j^*)}{\|\Sigma^{-1/2} \mathbf{e}_j\| \sqrt{\hat{r}_t}} \xrightarrow{d} N(0, 1) \quad \text{for any } j \in J_{n,p}. \quad (3)$$

- The asymptotic variance is proportional to \hat{r}_t .

- Suggest picking the t with smallest \hat{r}_t for inference tasks.

Summary

- Proposed a novel **\sqrt{n} -consistent estimator** for the generalization error of iterates along the trajectories of widely used algorithms.

- The form of the estimator depends on a **weighted residual vector** of all previous iterates. The weights are algorithm-specific and can be efficiently calculated.

- The proposed risk estimators can serve as a proxy for the generalization error, aiding in **early stopping** decisions.

- Established a **valid asymptotic normality result** by debiasing each $\hat{\mathbf{b}}^t$, which can be used for statistical inferences.

Reference

Bellec, Pierre C., and Kai Tan. Uncertainty quantification for iterative algorithms in linear models with application to early stopping. arXiv preprint arXiv:2404.17856 (2024).

Numerical experiments

The goals: to confirm (1) $\hat{r}_t \approx r_t$, (2) $\frac{\sqrt{n}(\hat{\mathbf{b}}_j^{t,\text{debias}} - \mathbf{b}_j^*)}{\|\Sigma^{-1/2} \mathbf{e}_j\| \sqrt{\hat{r}_t}} \xrightarrow{d} N(0, 1)$ for all t .

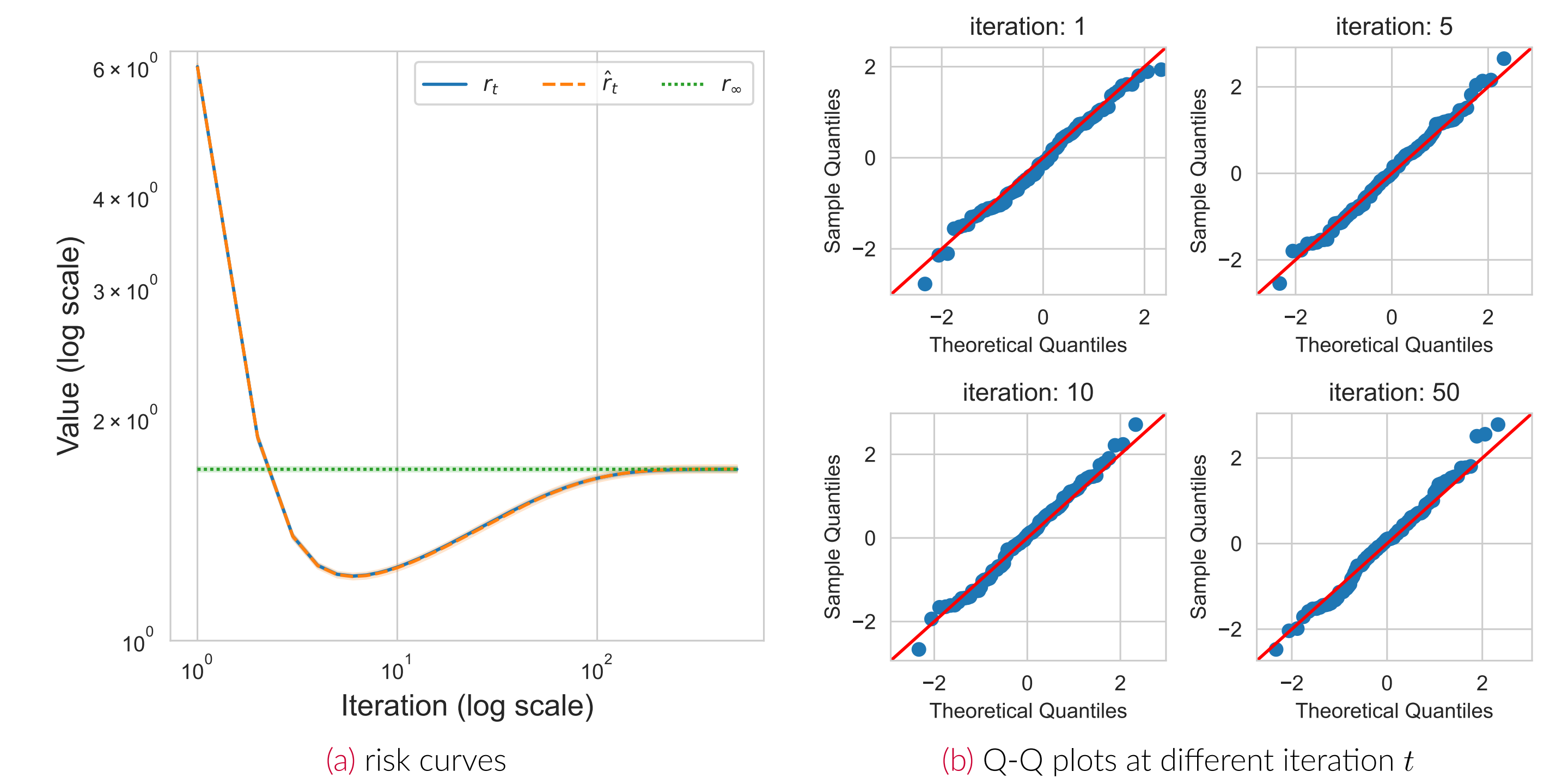


Figure 1. Risk curves and Q-Q plots for GD with $(n, p) = (1200, 500)$

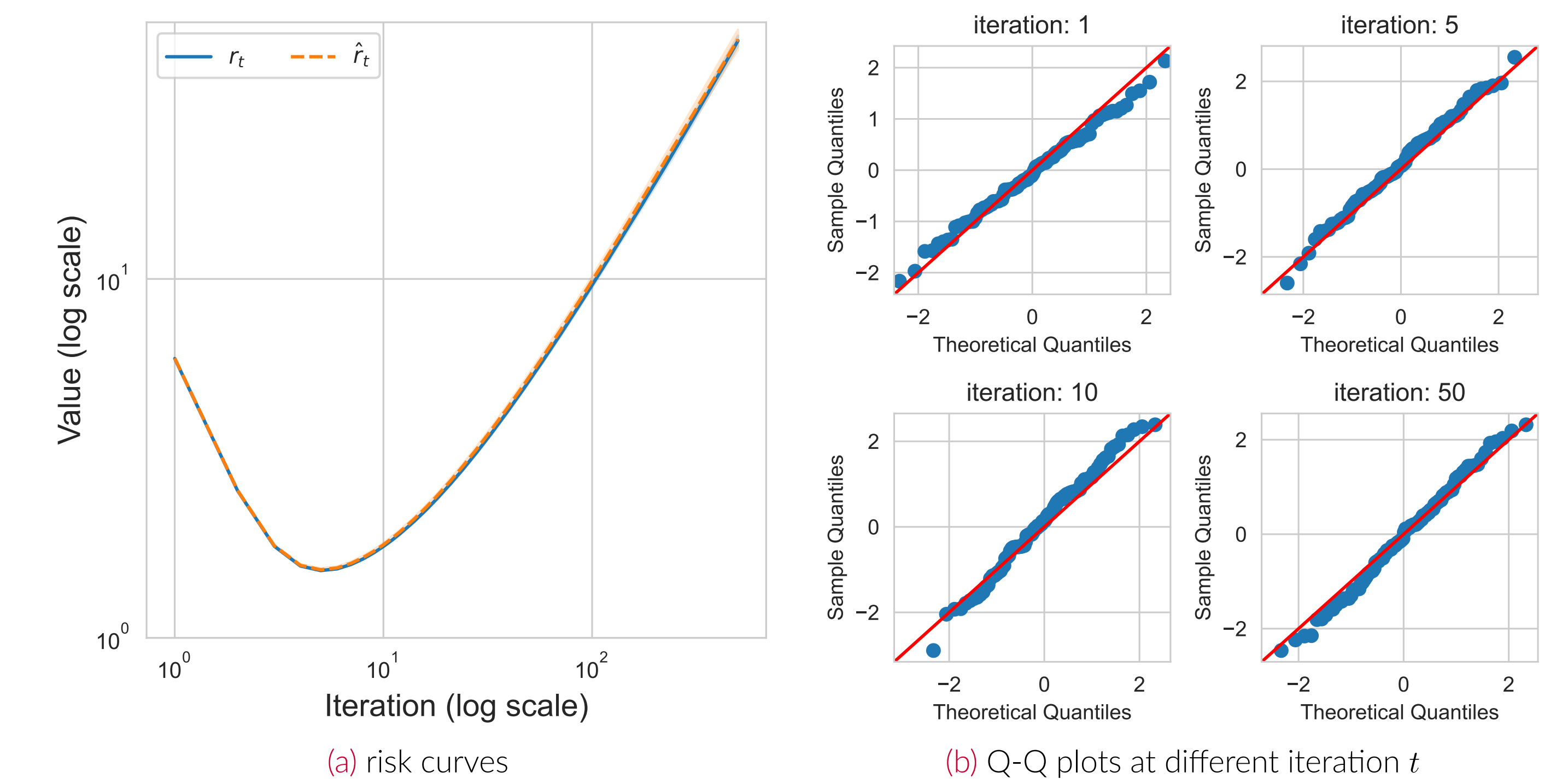


Figure 2. Risk curves and Q-Q plots for AGD with $(n, p) = (1200, 1200)$

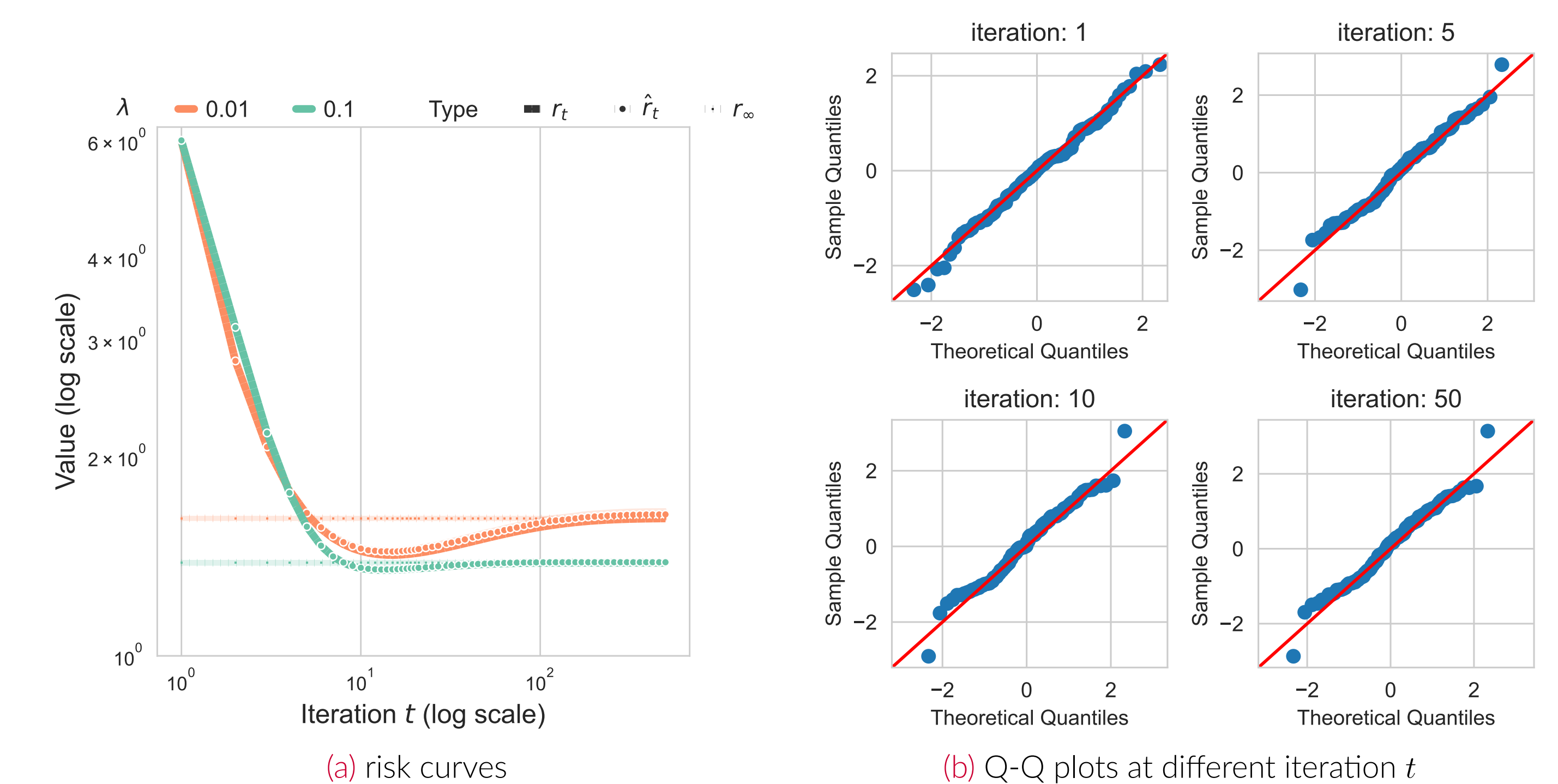


Figure 3. Risk curves and Q-Q plots for ISTA with $(n, p) = (1200, 1500)$