# First-Order Optimization Methods

## Analysis and design

Bryan Van Scoy

University of Wisconsin–Madison

Nov 1, 2017

**Unconstrained optimization:**

$$
\begin{aligned}
\text{minimize} \quad & f(x) \\
\text{subject to} \quad & x \in \mathbb{R}^d
\end{aligned}
$$

- Need methods which are *fast* and *simple*
- Use *first-order* methods

**Function class**

- quadratic
- smooth strongly convex

**Method**

GM  gradient method

HBM  heavy ball method

FGM  fast gradient method

**Bound**

- $f(x_k) - f(x_\star) \le c_1 \, \rho^k$
- $\|x_k - x_\star\| \le c_2 \, \rho^k$
- $\|\nabla f(x_k)\| \le c_3 \, \rho^k$

function class $+$ method $\implies$ bound

# Method

gradient method
$$x_{k+1} = x_k - \alpha \, \nabla f(x_k)$$

heavy ball method
$$x_{k+1} = (1 + \beta)x_k - \beta x_{k-1} - \alpha \, \nabla f(x_k)$$

fast gradient method
$$x_{k+1} = (1 + \beta)x_k - \beta x_{k-1} - \alpha \, \nabla f\big((1 + \beta)x_k - \beta x_{k-1}\big)$$

# Method

gradient method
$$x_{k+1} = x_k - \alpha \, \nabla f(x_k)$$

heavy ball method
$$x_{k+1} = (1 + \beta)x_k - \beta x_{k-1} - \alpha \, \nabla f(x_k)$$

fast gradient method
$$x_{k+1} = (1 + \beta)x_k - \beta x_{k-1} - \alpha \, \nabla f\big((1 + \beta)x_k - \beta x_{k-1}\big)$$

triple momentum method
$$x_{k+1} = (1 + \beta)x_k - \beta x_{k-1} - \alpha \, \nabla f\big((1 + \gamma)x_k - \gamma x_{k-1}\big)$$

| Method | Parameters |
|---|---|
| GM | $(\alpha, 0, 0)$ |
| HBM (Polyak, 1964) | $(\alpha, \beta, 0)$ |
| FGM (Nesterov, 2004) | $(\alpha, \beta, \beta)$ |
| TMM (Van Scoy, Freeman, Lynch, 2017) | $(\alpha, \beta, \gamma)$ |

# Triple momentum method

$$x_{k+1} = (1 + \beta)x_k - \beta x_{k-1} - \alpha \nabla f\big((1 + \gamma)x_k - \gamma x_{k-1}\big)$$

**Parameters:**

$\rho = 1 - \frac{1}{\sqrt{\kappa}}$

$\alpha = \frac{1+\rho}{L}$

$\beta = \frac{\rho^2}{2-\rho}$                         **Condition ratio** $\kappa := L/m$

$\gamma = \frac{\rho^2}{(1+\rho)(2-\rho)}$

# Triple momentum method

$$x_{k+1} = (1+\beta)x_k - \beta x_{k-1} - \alpha\nabla f\big((1+\gamma)x_k - \gamma x_{k-1}\big)$$

**Parameters:**

$\rho = 1 - \frac{1}{\sqrt{\kappa}}$

$\alpha = \frac{1+\rho}{L}$

$\beta = \frac{\rho^2}{2-\rho}$

$\gamma = \frac{\rho^2}{(1+\rho)(2-\rho)}$

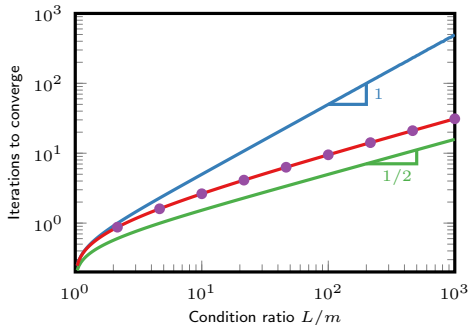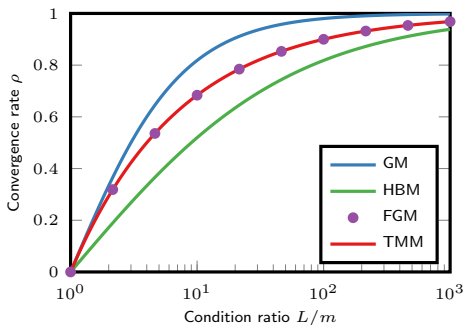**Condition ratio** $\kappa := L/m$

### Theorem (Van Scoy, Freeman, Lynch, 2017)

Suppose $f$ is $L$-smooth and $m$-strongly convex with minimizer $x_\star$. Then for any initial conditions $x_0, x_{-1} \in \mathbb{R}^n$, there exists a constant $c > 0$ such that

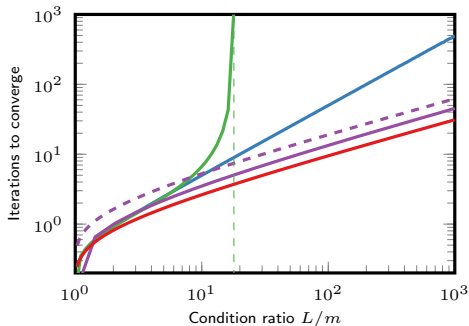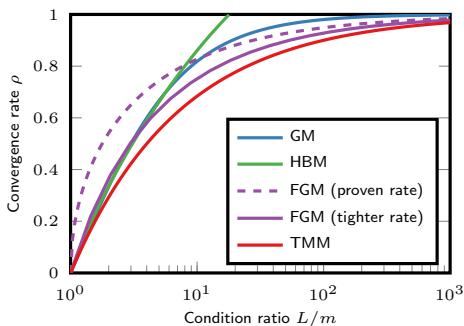$$\|x_k - x_\star\| \le c\,\rho^k \quad \text{for all } k \ge 1.$$

# $f$ quadratic



Convergence rate: $\|x_k - x_\star\| \leq c\,\rho^k$

Iterations to converge $\propto -\dfrac{1}{\log \rho}$

# $f$ smooth strongly convex



- HBM does not converge if $L/m \geq (2 + \sqrt{5})^2 \approx 17.94$
- For FGM, Nesterov proved the rate $\sqrt{1 - \sqrt{m/L}}$ which is loose!
- TMM converges faster than Nesterov's method!

# Simulations

**Objective function:**

$$f(x) = \sum_{i=1}^{p} g(a_i^T x - b_i) + \frac{m}{2} \|x\|^2, \quad x \in \mathbb{R}^n$$

where

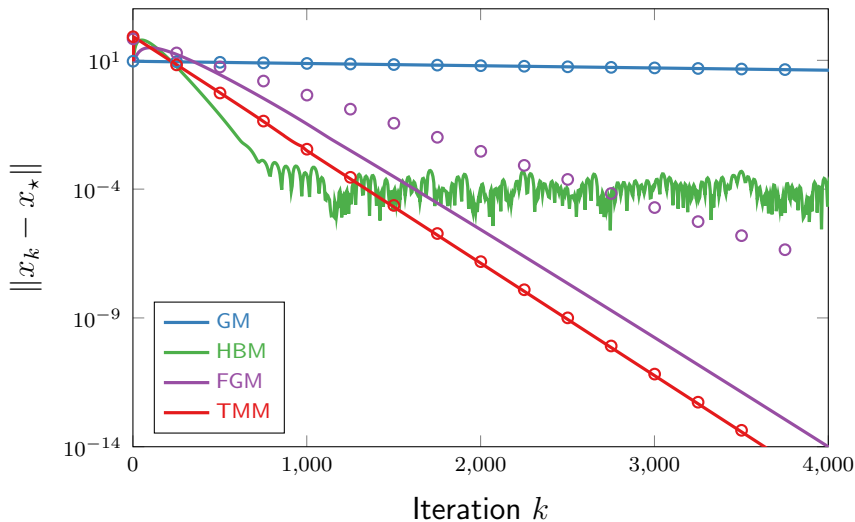$$g(y) = \begin{cases} \frac{1}{2} y^2 e^{-r/y}, & y > 0 \\ 0, & y \leq 0 \end{cases}$$

with $A = [a_1, \ldots, a_p] \in \mathbb{R}^{d \times p}$, $b \in \mathbb{R}^p$, and $\|A\| = \sqrt{L-m}$

$f$ is

- $m$-smooth
- $L$-strongly convex
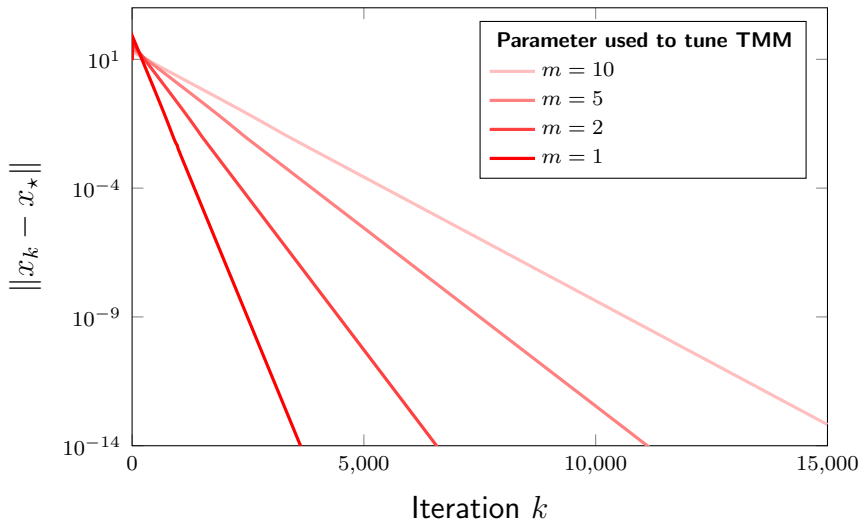- infinitely differentiable (of class $C^\infty$)

# Simulations

**Parameters:** $m = 1$, $L = 10^4$, $d = 100$, $p = 5$, $r = 10^{-6}$

# Robustness to $m$

**Parameters:** $m = 1$, $L = 10^4$, $d = 100$, $p = 5$, $r = 10^{-6}$



Legend within figure:
**Parameter used to tune TMM**
$m = 10$
$m = 5$
$m = 2$
$m = 1$

Vertical axis: $\|x_k - x_\star\|$

Horizontal axis: Iteration $k$

To prove the bound for TMM, using *interpolation*.

To prove the bound for TMM, using *interpolation*.

**Interpolation:** The set $\{y, u, v\}$ is $\mathcal{F}$-interpolable if and only if $u_k = \nabla f(y_k)$ and $v_k = f(y_k)$ for some $f \in \mathcal{F}$ and all $k$.
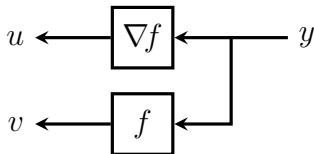
To prove the bound for TMM, using *interpolation*.

> **Interpolation:** The set $\{y, u, v\}$ is $\mathcal{F}$-interpolable if and only if $u_k = \nabla f(y_k)$ and $v_k = f(y_k)$ for some $f \in \mathcal{F}$ and all $k$.



Theorem (Taylor, Hendrickx, Glineur, 2016)

The set $\{y, u, v\}$ is interpolable by an $L$-smooth $m$-strongly convex function if and only if $q_{ij} \geq 0$ for all $i, j$ where

$$q_{ij} := (L - m)(v_i - v_j) - \frac{1}{2}\|u_i - u_j\|^2$$
$$+ (mu_i - Lu_j)^\mathsf{T}(y_i - y_j) - \frac{mL}{2}\|y_i - y_j\|^2.$$

# Sketch of proof for TMM

1. Suppose $f$ is $L$-smooth and $m$-strongly convex. Then the **interpolation conditions** are satisfied; specifically, $q_{ij} \geq 0$ for all $i, j$.

# Sketch of proof for TMM

1. Suppose $f$ is $L$-smooth and $m$-strongly convex. Then the **interpolation conditions** are satisfied; specifically, $q_{ij} \geq 0$ for all $i, j$.

2. Define the **Lyapunov function**

$$V_k := mL \|z_k - x_\star\|^2 + q_{k-1,\star}$$

where $z_k := (1 + \delta)x_k - \delta x_{k-1}$ and $\delta := \frac{\rho^2}{1-\rho^2}$.

# Sketch of proof for TMM

1. Suppose $f$ is $L$-smooth and $m$-strongly convex. Then the **interpolation conditions** are satisfied; specifically, $q_{ij} \geq 0$ for all $i, j$.

2. Define the **Lyapunov function**

$$V_k := mL \|z_k - x_\star\|^2 + q_{k-1,\star}$$

where $z_k := (1 + \delta)x_k - \delta x_{k-1}$ and $\delta := \frac{\rho^2}{1-\rho^2}$.

3. Using the definition of TMM, it is straightforward to verify that

$$V_{k+1} - \rho^2 V_k = -\left[(1 - \rho^2)q_{\star,k} + \rho^2 q_{k-1,k}\right] \leq 0$$

for all $k \geq 1$.

# Sketch of proof for TMM

1. Suppose $f$ is $L$-smooth and $m$-strongly convex. Then the **interpolation conditions** are satisfied; specifically, $q_{ij} \geq 0$ for all $i, j$.

2. Define the **Lyapunov function**

$$V_k := mL \|z_k - x_\star\|^2 + q_{k-1,\star}$$

   where $z_k := (1 + \delta)x_k - \delta x_{k-1}$ and $\delta := \frac{\rho^2}{1-\rho^2}$.

3. Using the definition of TMM, it is straighforward to verify that

$$V_{k+1} - \rho^2 V_k = -\big[(1 - \rho^2)q_{\star,k} + \rho^2 q_{k-1,k}\big] \leq 0$$

   for all $k \geq 1$.

4. Iterating gives the **bound** $V_k \leq \rho^{2(k-1)}V_1$ for $k \geq 1$.

# Numerics

For TMM, we can analyze the convergence rate in closed-form.

What can we say when a closed-form expression for the convergence rate is unknown?

# Numerics

For TMM, we can analyze the convergence rate in closed-form.

What can we say when a closed-form expression for the convergence rate is unknown?

Calculate an upper bound on the convergence rate numerically using:

- Integral Quadratic Constraints
  - Megretzki, Rantzer, 1997
  - Lessard, Recht, Packard, 2016
- Performance Estimation Problem
  - Drori, Teboulle, 2014
  - Taylor, Hendrickx, Glineur, 2016

# Gradient noise

What if the measured gradient is *not* the actual gradient?

$$x_{k+1} = (1 + \beta)x_k - \beta x_{k-1} - \alpha u_k$$
$$y_k = (1 + \gamma)x_k - \gamma x_{k-1}$$

**No noise:** $u = \nabla f(y)$

**Relative gradient noise:** $\|u - \nabla f(y)\|_2 \leq \delta \|\nabla f(y)\|_2$

# Robust momentum method

$$x_{k+1} = (1+\beta)x_k - \beta x_{k-1} - \alpha\nabla f\big((1+\gamma)x_k - \gamma x_{k-1}\big)$$

**Parameters:**

$\rho \in \big[1 - \frac{1}{\sqrt{\kappa}}, 1 - \frac{1}{\kappa}\big]$

$\alpha = \frac{\kappa(1-\rho)^2(1+\rho)}{L}$

$\beta = \frac{\kappa\rho^3}{\kappa-1}$

$\gamma = \frac{\rho^3}{(\kappa-1)(1-\rho)^2(1+\rho)}$

### Theorem (Cyrus, Hu, Van Scoy, Lessard, 2017)

Suppose $f$ is $L$-smooth and $m$-strongly convex with minimizer $x_\star$. Then for any initial conditions $x_0, x_{-1} \in \mathbb{R}^n$, there exists a constant $c > 0$ such that

$$\|x_k - x_\star\| \le c\,\rho^k \quad \text{for all } k \ge 1.$$

# Robust momentum method

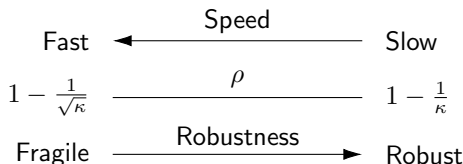$$x_{k+1} = (1+\beta)x_k - \beta x_{k-1} - \alpha \nabla f\big((1+\gamma)x_k - \gamma x_{k-1}\big)$$

**Parameters:**

$\rho \in \left[1 - \frac{1}{\sqrt{\kappa}}, 1 - \frac{1}{\kappa}\right]$

$\alpha = \frac{\kappa(1-\rho)^2(1+\rho)}{L}$

$\beta = \frac{\kappa\rho^3}{\kappa-1}$

$\gamma = \frac{\rho^3}{(\kappa-1)(1-\rho)^2(1+\rho)}$

Fast $\longleftarrow$ Speed $\longrightarrow$ Slow

$1 - \frac{1}{\sqrt{\kappa}}$ —————— $\rho$ —————— $1 - \frac{1}{\kappa}$

Fragile $\xrightarrow{\text{Robustness}}$ Robust

## Theorem (Cyrus, Hu, Van Scoy, Lessard, 2017)

Suppose $f$ is $L$-smooth and $m$-strongly convex with minimizer $x_\star$. Then for any initial conditions $x_0, x_{-1} \in \mathbb{R}^n$, there exists a constant $c > 0$ such that

$$\|x_k - x_\star\| \le c\,\rho^k \quad \text{for all } k \ge 1.$$

# Robust momentum method

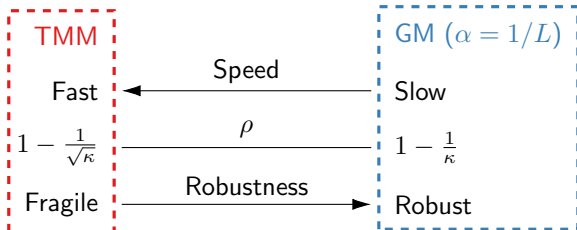$$x_{k+1} = (1+\beta)x_k - \beta x_{k-1} - \alpha \nabla f\big((1+\gamma)x_k - \gamma x_{k-1}\big)$$

**Parameters:**

$\rho \in \left[1 - \frac{1}{\sqrt{\kappa}}, 1 - \frac{1}{\kappa}\right]$

$\alpha = \frac{\kappa(1-\rho)^2(1+\rho)}{L}$

$\beta = \frac{\kappa\rho^3}{\kappa - 1}$
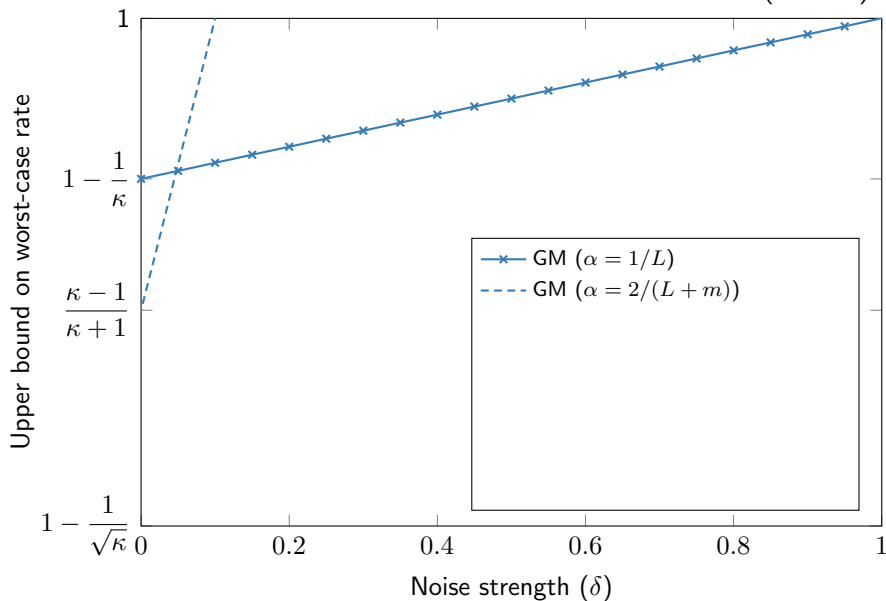
$\gamma = \frac{\rho^3}{(\kappa-1)(1-\rho)^2(1+\rho)}$



| TMM | Speed | GM ($\alpha = 1/L$) |
|---|---|---|
| Fast | ← Speed | Slow |
| $1 - \frac{1}{\sqrt{\kappa}}$ | —— $\rho$ —— | $1 - \frac{1}{\kappa}$ |
| Fragile | Robustness → | Robust |

## Theorem (Cyrus, Hu, Van Scoy, Lessard, 2017)

Suppose $f$ is $L$-smooth and $m$-strongly convex with minimizer $x_\star$. Then for any initial conditions $x_0, x_{-1} \in \mathbb{R}^n$, there exists a constant $c > 0$ such that

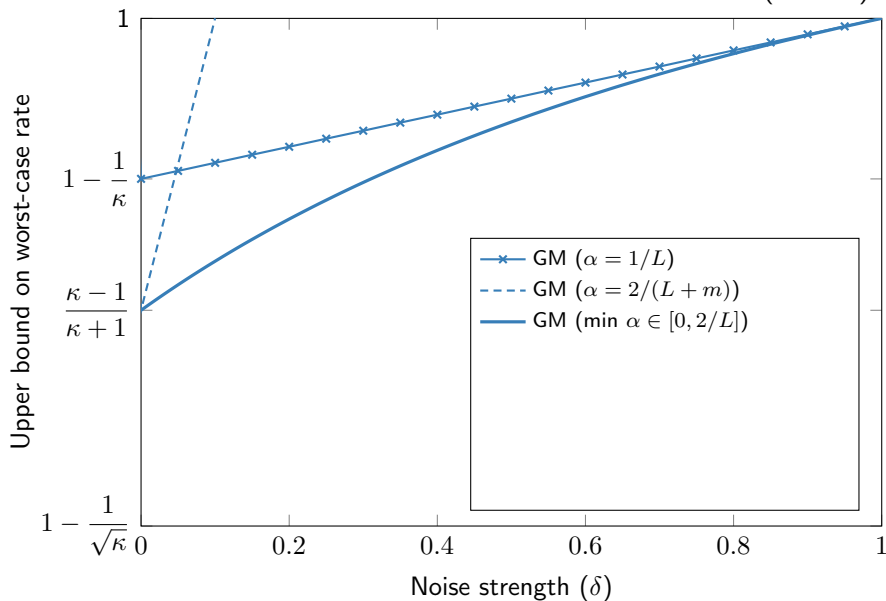$$\|x_k - x_\star\| \le c\,\rho^k \quad \text{for all } k \ge 1.$$
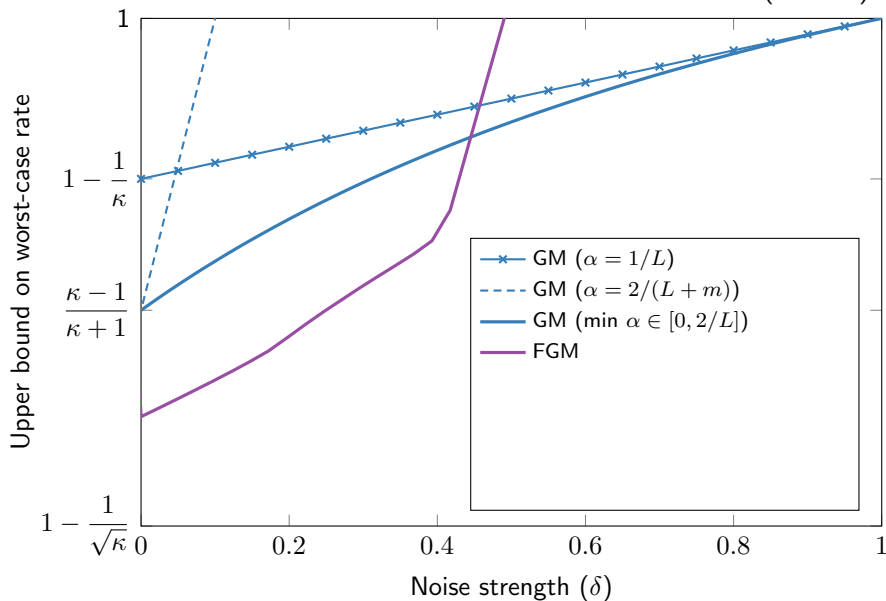
# Trade-off: Speed vs. Robustness

Upper bound on worst-case rate (y-axis), with marked values $1$, $1 - \frac{1}{\kappa}$, $\frac{\kappa - 1}{\kappa + 1}$, $1 - \frac{1}{\sqrt{\kappa}}$.

Noise strength ($\delta$) (x-axis), from $0$ to $1$.

Legend:
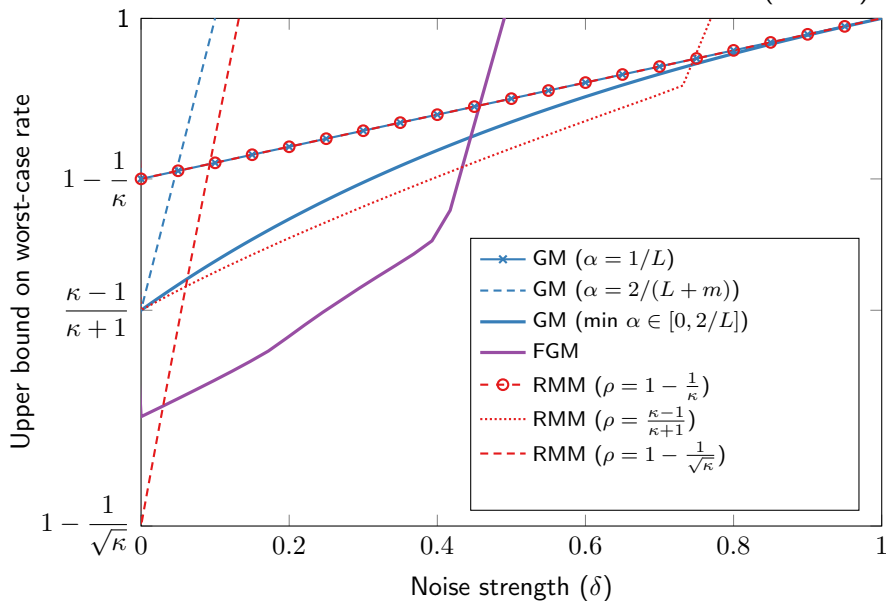- GM ($\alpha = 1/L$)
- GM ($\alpha = 2/(L + m)$)

# Trade-off: Speed vs. Robustness

$(\kappa = 10)$



16

# Trade-off: Speed vs. Robustness

$(\kappa = 10)$
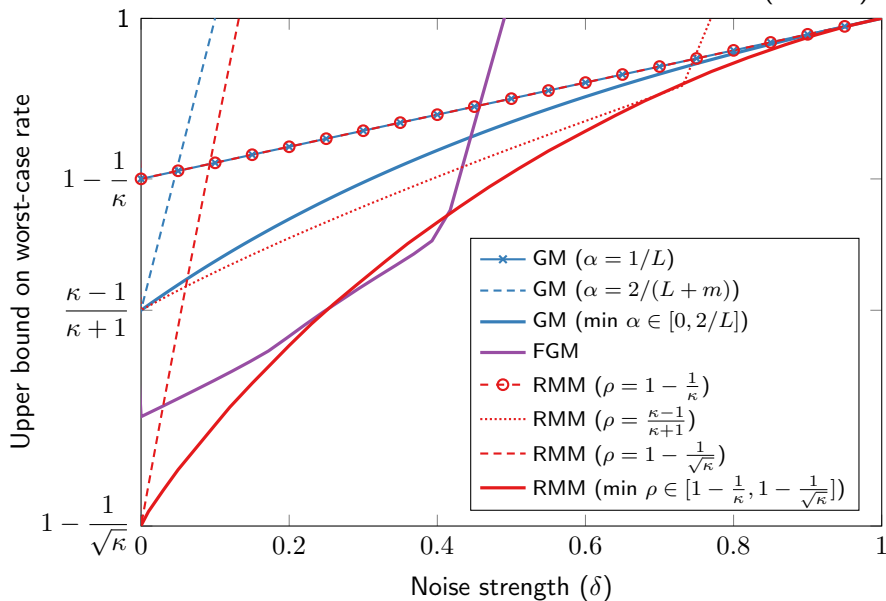


16

# Trade-off: Speed vs. Robustness

$(\kappa = 10)$



16

# Trade-off: Speed vs. Robustness

$(\kappa = 10)$



Legend:
- GM ($\alpha = 1/L$)
- GM ($\alpha = 2/(L+m)$)
- GM (min $\alpha \in [0, 2/L]$)
- FGM
- RMM ($\rho = 1 - \frac{1}{\kappa}$)
- RMM ($\rho = \frac{\kappa-1}{\kappa+1}$)
- RMM ($\rho = 1 - \frac{1}{\sqrt{\kappa}}$)
- RMM (min $\rho \in [1 - \frac{1}{\kappa}, 1 - \frac{1}{\sqrt{\kappa}}]$)

Axes:
- y-axis: Upper bound on worst-case rate
- x-axis: Noise strength ($\delta$)

# Conclusion

## Analysis

- **Numerical:** solve SDP to calculate upper bound on convergence rate
- **Closed-form:** have expressions for convergence rate for some methods and functions classes (such as TMM on smooth strongly convex functions)

## Design

- **Triple momentum method** - Fastest known convergence rate for first-order methods on smooth strongly convex functions
- **Robust momentum method** - Interpolates TMM and GM (with $\alpha = 1/L$) to exploit the trade-off between convergence rate and robustness to gradient noise