#### Noise Robustness of Optimization Algorithms via Statistical Queries

#### Cristóbal Guzmán

guzman@cwi.nl





CWI Scientific Meeting 27-11-2015

• In algorithms: Faster is better, right?

• In algorithms: Faster is better, right?

• Sometimes, other objectives are as important

• In algorithms: Faster is better, right?

• Sometimes, other objectives are as important

• This time, we care about robustness

#### Example: Gradient Methods

Goal: Find  $f^* = \min_{x \in X} f(x)$ 

where  $f: X \to \mathbb{R}$  smooth convex function,  $X \subseteq \mathbb{R}^d$  compact convex set





• Gradient Descent [Euler:1770]

 $x^{t+1} = x^t - \gamma_t \nabla f(x^t)$ 

• Gradient Descent [Euler:1770]

 $x^{t+1} = x^t - \gamma_t \nabla f(x^t) \qquad \Rightarrow \quad f(x^T) - f^* = O(1/T)$ 

• Gradient Descent [Euler:1770]

 $x^{t+1} = x^t - \gamma_t \nabla f(x^t) \qquad \Rightarrow \quad f(x^T) - f^* = O(1/T)$ 

• Accelerated Gradient Method [Nesterov: 1983]

$$x^{t} = y^{t} - \gamma_{t} \nabla f(y^{t})$$
$$y^{t+1} = x^{t} + \frac{\alpha_{t} - 1}{\alpha_{t+1}} (x^{t} - x^{t-1})$$

• Gradient Descent [Euler:1770]

 $x^{t+1} = x^t - \gamma_t \nabla f(x^t) \qquad \Rightarrow \quad f(x^T) - f^* = O(1/T)$ 

• Accelerated Gradient Method [Nesterov: 1983]

$$\begin{array}{rcccc}
x^{t} &=& y^{t} - \gamma_{t} \nabla f(y^{t}) \\
y^{t+1} &=& x^{t} + \frac{\alpha_{t} - 1}{\alpha_{t+1}} (x^{t} - x^{t-1}) & f(x^{T}) - f^{*} = O(1/T^{2})
\end{array}$$

#### Example: Gradient Methods



#### Example: Gradient Methods



• There are other important reasons why acceleration could be harmful

- There are other important reasons why acceleration could be harmful
  - Overfitting, privacy leakage, etc.

- There are other important reasons why acceleration could be harmful
  - Overfitting, privacy leakage, etc.
- Goal: A theory of optimization algorithms that incorporates noise sensitivity



- There are other important reasons why acceleration could be harmful
  - Overfitting, privacy leakage, etc.



 $\bullet$  Unknown distribution  ${\cal D}$  supported on  ${\cal W}$ 



- $\bullet$  Unknown distribution  ${\cal D}$  supported on  ${\cal W}$
- $\bullet$  Algorithm has oracle access to  ${\cal D}$





- $\bullet$  Unknown distribution  ${\cal D}$  supported on  ${\cal W}$
- $\bullet$  Algorithm has oracle access to  ${\cal D}$



An oracle query can be emulated by  $1/\tau^2$  samples

- $\bullet$  Unknown distribution  ${\cal D}$  supported on  ${\cal W}$
- $\bullet$  Algorithm has oracle access to  ${\cal D}$





• Efficiency ~ small number of queries

- $\bullet$  Unknown distribution  ${\cal D}$  supported on  ${\cal W}$
- $\bullet$  Algorithm has oracle access to  ${\cal D}$



#### • Goals

- Efficiency ~ small number of queries
- Noise tolerance ~ au as large as possible

# Why Statistical Queries?

SQ algorithms have additional properties, or provide the basis for designing algorithms with additional features:

## Why Statistical Queries?

SQ algorithms have additional properties, or provide the basis for designing algorithms with additional features:

- Noise tolerance [Kearns: 1994]
- Differential privacy [Blum, Dwork, McSherry, Nissim:2005], [Kasiviswanathan, Lee, Nissim, Raskhodnikova, Smith, 2011]
- Distributed computation [Balcan, Blum, Fine, Mansour:2012]
- Generalization in adaptive data analysis [Dwork, Feldman, Hardt, Pitassi, Reingold, Roth:2014]

Stochastic convex optimization

 $\min_{x \in X} \mathop{\mathbb{E}}_{\mathbf{w} \sim \mathcal{D}} [f(x, \mathbf{w})]$ 

Stochastic convex optimization

$$\min_{x \in X} \mathop{\mathbb{E}}_{\mathbf{w} \sim \mathcal{D}} [f(x, \mathbf{w})]$$

- Statistical queries
  - Function value:  $\phi(\cdot) = f(x, \cdot)$

Stochastic convex optimization

$$\min_{x \in X} \mathop{\mathbb{E}}_{\mathbf{w} \sim \mathcal{D}} [f(x, \mathbf{w})]$$

- Statistical queries
  - Function value:  $\phi(\cdot) = f(x, \cdot)$

• Gradient: 
$$\phi(\cdot) = \frac{\partial f(x, \cdot)}{\partial x_i}$$
  $i = 1, \dots, d$ 

Stochastic convex optimization

$$\min_{x \in X} \mathop{\mathbb{E}}_{\mathbf{w} \sim \mathcal{D}} [f(x, \mathbf{w})]$$

- Statistical queries
  - Function value:  $\phi(\cdot) = f(x, \cdot)$

• Gradient: 
$$\phi(\cdot) = \frac{\partial f(x, \cdot)}{\partial x_i}$$
  $i = 1, \dots, d$ 

• Difficulty: estimation in 2-norm accumulates errors

$$g_i = \mathop{\mathbb{E}}_{\mathbf{w} \sim \mathcal{D}} \left[ \frac{\partial f(x, \mathbf{w})}{\partial x_i} \right] \pm \tau \quad \Rightarrow \quad \|\nabla \mathbb{E}_{\mathbf{w}}[f(x, \mathbf{w})] - g\|_2 \approx \sqrt{d\tau}$$

[Feldman, G., Vempala:2015]

Q: Is it possible to avoid noise accumulation in gradient estimation?



[Feldman, G., Vempala:2015]

Q: Is it possible to avoid noise accumulation in gradient estimation?



[Feldman, G., Vempala:2015]

Q: Is it possible to avoid noise accumulation in gradient estimation?



[Feldman, G., Vempala:2015]

Q: Is it possible to avoid noise accumulation in gradient estimation?



With high probability  $\|\nabla \mathbb{E}_{\mathbf{w}}[f(x, \mathbf{w})] - g\|_2 \approx O(\sqrt{\log d})\tau$ 

# Further Consequences

[Feldman, G., Vempala:2015]

New statistical query algorithms for optimization and machine learning

# Further Consequences

[Feldman, G., Vempala:2015]

New statistical query algorithms for optimization and machine learning

- Gradient methods: mirror-descent, accelerated method, strongly convex minimization
- Polynomial-time algorithms: center of gravity, interior-point method
- Improved algorithms for high-dimensional classification (Perceptron) and regression
- Improved differentially-private convex optimization algorithms

# Further Consequences

[Feldman, G., Vempala:2015]

New statistical query algorithms for optimization and machine learning

- Gradient methods: mirror-descent, accelerated method, strongly convex minimization
- Polynomial-time algorithms: center of gravity, interior-point method
- Improved algorithms for high-dimensional classification (Perceptron) and regression
- Improved differentially-private convex optimization algorithms

We provide new structural lower bounds for convex optimization

Thank you