The Gradient Method: Past and Present

Amir Beck

School of Mathematical Sciences, Tel Aviv University

Workshop in honor of Uri Rothblum, Tel Aviv University, January 22, 2018

Amir Beck - Tel Aviv University The G

The Gradient Method: Past and Present

The Gradient Method

The problem.

```
\min\{f(\mathbf{x}):\mathbf{x}\in\mathbb{R}^n\}
```

f differentiable.

The Gradient Method

The problem.

```
\min\{f(\mathbf{x}):\mathbf{x}\in\mathbb{R}^n\}
```

f differentiable.

The Gradient Method

$$\mathbf{x}^{k+1} = \mathbf{x}^k - t_k \nabla f(\mathbf{x}^k)$$

 $t_k > 0$ - chosen stepsize.

The Gradient Method

The problem.

```
\min\{f(\mathbf{x}):\mathbf{x}\in\mathbb{R}^n\}
```

f differentiable.

The Gradient Method

$$\mathbf{x}^{k+1} = \mathbf{x}^k - t_k \nabla f(\mathbf{x}^k)$$

 $t_k > 0$ - chosen stepsize.

- What is the starting point?
- What stepsize should be taken?
- What is the stopping criteria?

Stepsize Selection Rules

• constant stepsize - $t_k = \overline{t}$ for any k.

¹also referred to as Armijo

Stepsize Selection Rules

- constant stepsize $t_k = \overline{t}$ for any k.
- exact stepsize t_k is a minimizer of f along the ray $\mathbf{x}_k t \nabla f(\mathbf{x}^k)$:

$$t_k \in \operatorname*{argmin}_{t \geq 0} f(\mathbf{x}^k - t \nabla f(\mathbf{x}^k)).$$

¹also referred to as Armijo

Stepsize Selection Rules

- constant stepsize $t_k = \overline{t}$ for any k.
- exact stepsize t_k is a minimizer of f along the ray $\mathbf{x}_k t \nabla f(\mathbf{x}^k)$:

$$t_k \in \operatorname*{argmin}_{t \geq 0} f(\mathbf{x}^k - t
abla f(\mathbf{x}^k)).$$

▶ **backtracking**¹ - The method requires three parameters: $s > 0, \alpha \in (0, 1), \beta \in (0, 1)$. Here we start with an initial stepsize $t_k = s$. While

$$f(\mathbf{x}^k) - f(\mathbf{x}^k - t_k \nabla f(\mathbf{x}^k)) < \alpha t_k \| \nabla f(\mathbf{x}^k) \|^2.$$

set $t_k := \beta t_k$

Sufficient Decrease Property:

$$f(\mathbf{x}^k) - f(\mathbf{x}^k - t_k \nabla f(\mathbf{x}^k)) \ge \alpha t_k \|\nabla f(\mathbf{x}^k)\|^2.$$

¹also referred to as Armijo

Gradient Method as Steepest Descent

• $-\nabla f(\mathbf{x}^k)$ is a descent direction:

 $f'(\mathbf{x}^k; -\nabla f(\mathbf{x}^k)) = -\nabla f(\mathbf{x}^k)^T \nabla f(\mathbf{x}^k) = -\|\nabla f(\mathbf{x}^k)\|^2 < 0.$

In addition for being a descent direction, minus the gradient is also the steepest descent direction method.

Gradient Method as Steepest Descent

• $-\nabla f(\mathbf{x}^k)$ is a descent direction:

 $f'(\mathbf{x}^k; -\nabla f(\mathbf{x}^k)) = -\nabla f(\mathbf{x}^k)^T \nabla f(\mathbf{x}^k) = -\|\nabla f(\mathbf{x}^k)\|^2 < 0.$

In addition for being a descent direction, minus the gradient is also the steepest descent direction method.

Lemma. Let f be a differentiable function and let $\mathbf{x} \in \mathbb{R}^n$ satisfy $\nabla f(\mathbf{x}) \neq \mathbf{0}$. Then an optimal solution of

 $\min_{\mathbf{d}} \{ f'(\mathbf{x}; \mathbf{d}) : \|\mathbf{d}\| = 1 \}$

is $\mathbf{d} = -\nabla f(\mathbf{x}) / \|\nabla f(\mathbf{x})\|.$

Convergence(?) of the Gradient Method

 $\min\{f(\mathbf{x}):\mathbf{x}\in\mathbb{R}^n\}$

Standard conditions:

- *f* is bounded below and differentiable.
- f is L-smooth meaning that

 $\|
abla f(\mathbf{x}) -
abla f(\mathbf{y})\| \le L \|\mathbf{x} - \mathbf{y}\| \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n.$

Convergence(?) of the Gradient Method

 $\min\{f(\mathbf{x}):\mathbf{x}\in\mathbb{R}^n\}$

Standard conditions:

- *f* is bounded below and differentiable.
- f is L-smooth meaning that

 $\|
abla f(\mathbf{x}) -
abla f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\| \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n.$

Can be proved:

- Descent method: $f(\mathbf{x}^{k+1}) < f(\mathbf{x}^k)$
- ▶ Accumulation pts. of the sequence generated by GM are stationary points $(\nabla f(\mathbf{x}^*) = 0)$ (constant stepsize $t_k \equiv \overline{t} \in (0, \frac{2}{L})$, backtracking or exact minimization)
- ▶ If f is convex, convergence to a global optimal solution.

Gradient Method - the Oldest Continuous Optimization Method?

Méthode generales pour la résolution des systèmes d'equations simultanées, 1847



Augustin Louis Cauchy

1789-1857

Suggested the method for solving sets of nonlinear equations

$$f_i(\mathbf{x}) = 0, i = 1, 2, \dots, m \Rightarrow \min_{\mathbf{x}} \sum_{i=1}^m f_i(\mathbf{x})^2$$

- Not a particularly rigorous paper...
- Modern optimization starts only 100 years afterwards (simplex for LP)



Gradient-Based Algorithms

Widely used in applications....

- **Clustering Analysis:** The k-means algorithm
- Neuro-computing: The backpropagation algorithm
- Statistical Estimation: The EM (Expectation-Maximization) algorithm.
- Machine Learning: SVM, Regularized regression, etc...
- Signal and Image Processing: Sparse Recovery, Denoising and Deblurring Schemes, Total Variation minimization...
- Matrix minimization Problems....and much more...

The Zig-Zag Property

Zig-Zagging: directions produced by the gradient method with exact minimization are perpendicular.

$$\langle \nabla f(\mathbf{x}^k), \nabla f(\mathbf{x}^{k+1}) \rangle = 0$$



The Zig-Zag Property

Zig-Zagging: directions produced by the gradient method with exact minimization are perpendicular.

$$\langle \nabla f(\mathbf{x}^k), \nabla f(\mathbf{x}^{k+1}) \rangle = 0$$



Main disadvantage: gradient method is rather slow.

Amir Beck - Tel Aviv University

The Gradient Method: Past and Present

The Zig-Zag Property

Zig-Zagging: directions produced by the gradient method with exact minimization are perpendicular.

$$\langle \nabla f(\mathbf{x}^k), \nabla f(\mathbf{x}^{k+1}) \rangle = 0$$



Main disadvantage: gradient method is rather slow. **Advantages:** requires minimal information $(f, \nabla f)$, "cheap" iterative scheme, suitable for large-scale problems.

Amir Beck - Tel Aviv University

The Gradient Method: Past and Present

The Condition Number

Rate of convergence of the gradient method depends on the condition number of the matrix \(\nabla^2 f(\mathbf{x}^*)\):

$$\kappa(\nabla^2 f(\mathbf{x}^*)) = \frac{\sigma_{\max}(\nabla^2 f(\mathbf{x}^*))}{\sigma_{\min}(\nabla^2 f(\mathbf{x}^*))}$$

- Ill-conditioned problems high condition number
- Well-conditioned problems small condition number

A Severely III-Condition Function - Rosenbrock

min {
$$f(x_1, x_2) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$
 }.

condition number: 2508

A Severely III-Condition Function - Rosenbrock

$$\min \left\{ f(x_1, x_2) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2 \right\}.$$

condition number: 2508



6890(!!!) iterations.

Amir Beck - Tel Aviv University

The Gradient Method: Past and Present

Scaled Gradient Method

$$\mathbf{x}^{k+1} = \mathbf{x}^k - t_k \mathbf{D}_k \nabla f(\mathbf{x}^k)$$

 $t_k > 0$ - chosen stepsize. $\mathbf{D}_k \succ \mathbf{0}$

Scaled Gradient Method

$$\mathbf{x}^{k+1} = \mathbf{x}^k - t_k \mathbf{D}_k \nabla f(\mathbf{x}^k)$$

 $t_k > 0$ - chosen stepsize. $\mathbf{D}_k \succ \mathbf{0}$

Since $\mathbf{D}_k \succ \mathbf{0}$ - still a descent directions method.

Scaled Gradient Method

$$\mathbf{x}^{k+1} = \mathbf{x}^k - t_k \mathbf{D}_k \nabla f(\mathbf{x}^k)$$

 $t_k > 0$ - chosen stepsize. $\mathbf{D}_k \succ \mathbf{0}$

- Since $\mathbf{D}_k \succ \mathbf{0}$ still a descent directions method.
- Same as the gradient method employed after the change of variables
 x = **D**_k^{1/2}**y**

Scaled Gradient Method

$$\mathbf{x}^{k+1} = \mathbf{x}^k - t_k \mathbf{D}_k \nabla f(\mathbf{x}^k)$$

 $t_k > 0$ - chosen stepsize. $\mathbf{D}_k \succ \mathbf{0}$

- Since $\mathbf{D}_k \succ \mathbf{0}$ still a descent directions method.
- Same as the gradient method employed after the change of variables
 x = D_k^{1/2}y
- Convergence is related to the condition number of $\mathbf{D}_k^{-1/2} \nabla^2 f(\mathbf{x}^k) \mathbf{D}_k^{-1/2}$

Scaled Gradient Method

$$\mathbf{x}^{k+1} = \mathbf{x}^k - t_k \mathbf{D}_k \nabla f(\mathbf{x}^k)$$

 $t_k > 0$ - chosen stepsize. $\mathbf{D}_k \succ \mathbf{0}$

- Since $\mathbf{D}_k \succ \mathbf{0}$ still a descent directions method.
- Same as the gradient method employed after the change of variables
 x = **D**_k^{1/2}**y**
- ► Convergence is related to the condition number of D_k^{-1/2}∇²f(x^k)D_k^{-1/2}
- "best" choice $\mathbf{D}_k = \nabla^2 f(\mathbf{x}^k)^{-1}$. pure Newton's method:

$$\mathbf{x}^{k+1} = \mathbf{x}^k -
abla^2 f(\mathbf{x}^k)^{-1}
abla f(\mathbf{x}^k)$$

Popular and "cheap" choice: D_k diagonal (diagonal scaling)

Gradient Method





Nonlinear least squares:

$$\min_{\mathbf{x}\in\mathbb{R}^n}\sum_{i=1}^m(f_i(\mathbf{x})-c_i)^2$$

 f_1, f_2, \ldots, f_m - differentiable.

Nonlinear least squares:

 $\min_{\mathbf{x}\in\mathbb{R}^n}\sum_{i=1}^m(f_i(\mathbf{x})-c_i)^2$

 f_1, f_2, \ldots, f_m - differentiable.

Given the kth iterate \mathbf{x}^k , the next iterate is chosen to minimize the sum of squares of the linearized terms, that is,

$$\mathbf{x}^{k+1} = \operatorname*{argmin}_{\mathbf{x} \in \mathbb{R}^n} \left\{ \sum_{i=1}^m \left[f_i(\mathbf{x}^k) + \nabla f_i(\mathbf{x}^k)^T (\mathbf{x} - \mathbf{x}^k) - c_i \right]^2 \right\}.$$

Nonlinear least squares:

 $\min_{\mathbf{x}\in\mathbb{R}^n}\sum_{i=1}^m(f_i(\mathbf{x})-c_i)^2$

 f_1, f_2, \ldots, f_m - differentiable.

Given the *k*th iterate \mathbf{x}^k , the next iterate is chosen to minimize the sum of squares of the linearized terms, that is,

$$\mathbf{x}^{k+1} = \operatorname*{argmin}_{\mathbf{x} \in \mathbb{R}^n} \left\{ \sum_{i=1}^m \left[f_i(\mathbf{x}^k) + \nabla f_i(\mathbf{x}^k)^T (\mathbf{x} - \mathbf{x}^k) - c_i \right]^2 \right\}.$$

The general step requires to solve a linear least squares problem at each iteration.

Nonlinear least squares:

 $\min_{\mathbf{x}\in\mathbb{R}^n}\sum_{i=1}^m(f_i(\mathbf{x})-c_i)^2$

 f_1, f_2, \ldots, f_m - differentiable.

Given the *k*th iterate \mathbf{x}^k , the next iterate is chosen to minimize the sum of squares of the linearized terms, that is,

$$\mathbf{x}^{k+1} = \operatorname*{argmin}_{\mathbf{x} \in \mathbb{R}^n} \left\{ \sum_{i=1}^m \left[f_i(\mathbf{x}^k) +
abla f_i(\mathbf{x}^k)^T (\mathbf{x} - \mathbf{x}^k) - c_i
ight]^2
ight\}.$$

- The general step requires to solve a linear least squares problem at each iteration.
- ► Actually a scaled gradient method with D_k = (J(x^k)^TJ(x^k))⁻¹ (J(·) Jacobian)

Amir Beck - Tel Aviv University

The Gradient Method: Past and Present





Problems with Newton's Method

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \nabla^2 f(\mathbf{x}^k)^{-1} \nabla f(\mathbf{x}^k)$$

• $\nabla^2 f(\mathbf{x}^k)$ difficult to compute and/or problematic to solve the system $\nabla^2 f(\mathbf{x}^k) \mathbf{z} = \nabla f(\mathbf{x}^k)$

Problems with Newton's Method

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \nabla^2 f(\mathbf{x}^k)^{-1} \nabla f(\mathbf{x}^k)$$

- ∇²f(x^k) difficult to compute and/or problematic to solve the system ∇²f(x^k)z = ∇f(x^k)
- $\nabla^2 f(\mathbf{x}^k)$ might be singular
- ▷ ∇² f(x^k) might not be positive definite (Newton's direction not a descent direction...)
Problems with Newton's Method

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \nabla^2 f(\mathbf{x}^k)^{-1} \nabla f(\mathbf{x}^k)$$

- $\nabla^2 f(\mathbf{x}^k)$ difficult to compute and/or problematic to solve the system $\nabla^2 f(\mathbf{x}^k) \mathbf{z} = \nabla f(\mathbf{x}^k)$
- $\nabla^2 f(\mathbf{x}^k)$ might be singular
- ▷ ∇² f(x^k) might not be positive definite (Newton's direction not a descent direction...)
- Convergence extremely problematic: requires a lot of assumptions that are usually not satisfied.

Problems with Newton's Method

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \nabla^2 f(\mathbf{x}^k)^{-1} \nabla f(\mathbf{x}^k)$$

- $\nabla^2 f(\mathbf{x}^k)$ difficult to compute and/or problematic to solve the system $\nabla^2 f(\mathbf{x}^k) \mathbf{z} = \nabla f(\mathbf{x}^k)$
- $\nabla^2 f(\mathbf{x}^k)$ might be singular
- ▷ ∇² f(x^k) might not be positive definite (Newton's direction not a descent direction...)
- Convergence extremely problematic: requires a lot of assumptions that are usually not satisfied.
- main advantage: quadratic rate of convergence (under very restrictive conditions...)

Problems with Newton's Method

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \nabla^2 f(\mathbf{x}^k)^{-1} \nabla f(\mathbf{x}^k)$$

- $\nabla^2 f(\mathbf{x}^k)$ difficult to compute and/or problematic to solve the system $\nabla^2 f(\mathbf{x}^k) \mathbf{z} = \nabla f(\mathbf{x}^k)$
- $\nabla^2 f(\mathbf{x}^k)$ might be singular
- ▷ ∇² f(x^k) might not be positive definite (Newton's direction not a descent direction...)
- Convergence extremely problematic: requires a lot of assumptions that are usually not satisfied.
- main advantage: quadratic rate of convergence (under very restrictive conditions...)

Btw, pure Newton's is a utopian method. Better to incorporate a stepsize (damped Newton).

Classics from the 70's - Trying to Mend Newton

• Trust-Region Methods

$$\mathbf{x}^{k+1} \in \operatorname{argmin} \left\{ m(\mathbf{x}; \mathbf{x}^k) : \|\mathbf{x} - \mathbf{x}^k\| \leq \Delta_k
ight\}$$

where $m(\mathbf{x}; \mathbf{x}^k)$ is a model of f around \mathbf{x}^k , e.g., $m(\mathbf{x}; \mathbf{x}^k) \equiv f(\mathbf{x}^k) + \nabla f(\mathbf{x}^k)^T (\mathbf{x} - \mathbf{x}^k) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^k)^T \nabla^2 f(\mathbf{x}^k) (\mathbf{x} - \mathbf{x}^k)$ Classics from the 70's - Trying to Mend Newton

• Trust-Region Methods

$$\mathbf{x}^{k+1} \in \operatorname{argmin}\left\{m(\mathbf{x};\mathbf{x}^k): \|\mathbf{x}-\mathbf{x}^k\| \leq \Delta_k
ight\}$$

where $m(\mathbf{x}; \mathbf{x}^k)$ is a model of f around \mathbf{x}^k , e.g., $m(\mathbf{x}; \mathbf{x}^k) \equiv f(\mathbf{x}^k) + \nabla f(\mathbf{x}^k)^T (\mathbf{x} - \mathbf{x}^k) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^k)^T \nabla^2 f(\mathbf{x}^k) (\mathbf{x} - \mathbf{x}^k)$

• Quasi-Newton Try to mimic the Hessian without actually forming it. e.g., BFGS

$$\mathbf{x}^{k+1} = \mathbf{x}^k - t_k \mathbf{D}_k^{-1}
abla f(\mathbf{x}^k)$$

 \mathbf{D}_k is chosen to satisfy the QN condition

$$\mathbf{D}_k(\mathbf{x}^k - \mathbf{x}^{k-1}) =
abla f(\mathbf{x}^k) -
abla f(\mathbf{x}^{k-1})$$

Amir Beck - Tel Aviv University

The Gradient Method: Past and Present

Classics from the 70's - Trying to Mend Newton

• Trust-Region Methods

$$\mathbf{x}^{k+1} \in \operatorname{argmin}\left\{m(\mathbf{x};\mathbf{x}^k): \|\mathbf{x}-\mathbf{x}^k\| \leq \Delta_k
ight\}$$

where $m(\mathbf{x}; \mathbf{x}^k)$ is a model of f around \mathbf{x}^k , e.g., $m(\mathbf{x}; \mathbf{x}^k) \equiv f(\mathbf{x}^k) + \nabla f(\mathbf{x}^k)^T (\mathbf{x} - \mathbf{x}^k) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^k)^T \nabla^2 f(\mathbf{x}^k) (\mathbf{x} - \mathbf{x}^k)$

• Quasi-Newton Try to mimic the Hessian without actually forming it. e.g., BFGS

$$\mathbf{x}^{k+1} = \mathbf{x}^k - t_k \mathbf{D}_k^{-1}
abla f(\mathbf{x}^k)$$

 D_k is chosen to satisfy the QN condition

$$\mathbf{D}_k(\mathbf{x}^k - \mathbf{x}^{k-1}) = \nabla f(\mathbf{x}^k) - \nabla f(\mathbf{x}^{k-1})$$

D_{k+1} "simply" constructed from D_k
 Computation of D_k⁻¹ requires only O(n²) flops (linear algebra tricks)

Amir Beck - Tel Aviv University









So far...

Classical algorithms for solving



So far...

Classical algorithms for solving

The problem.

 $\min\{f(\mathbf{x}):\mathbf{x}\in\mathbb{R}^n\}$

f differentiable.

What happens if f is nonsmooth? e.g.,

$$f(\mathbf{x}) = \sum_{i=1}^{m} |\mathbf{a}_i^T \mathbf{x} - b_i|, \max_{i=1,\dots,m} |\mathbf{a}_i^T \mathbf{x} - b_i|....$$

Amir Beck - Tel Aviv University

The Gradient Method: Past and Present

Wolfe's Example

False hope: What happens if the method never encounters non-differentiability points?

Wolfe's Example

False hope: What happens if the method never encounters non-differentiability points?

• Let $\gamma > 1$ and consider

$$f(x_1, x_2) = \begin{cases} \sqrt{x_1^2 + \gamma x_2^2}, & |x_2| \le x_1, \\ \frac{x_1 + \gamma |x_2|}{\sqrt{1 + \gamma}}, & \text{else.} \end{cases}$$

▶ *f* is differentiable at all points except for the ray $\{(x_1, 0) : x_1 \leq 0\}$.

Amir Beck - Tel Aviv University

The Gradient Method: Past and Present

Wolfe's Example

The gradient method with exact line search converges to a non-optimal point.

 $\mbox{Conclusion:}$ cannot ignore non-differentiability \rightarrow extend the notion of the gradient



Amir Beck - Tel Aviv University

The Subgradient Method. Shor (63) Polyak (65) $\mathbf{x}^{k+1} = \mathbf{x}^k - t_k f'(\mathbf{x}^k)$

Replace the gradient $\nabla f(\mathbf{x})$ by a subgradient $f'(\mathbf{x}) \in \partial f(\mathbf{x})$ (vectors that correspond to underestimators of the function)



The Gradient Method: Past and Present

Projected Subgradient Method

Model: f - convex. C - closed convex

$\min\{f(\mathbf{x}):\mathbf{x}\in C\}$

Projected Subgradient Method

Model: f - convex. C - closed convex

$\min\{f(\mathbf{x}):\mathbf{x}\in C\}$

Projected Subgradient Method (PSM): Shor (63), Polyak (65)

$$\mathbf{x}^{k} = P_{\mathcal{C}}(\mathbf{x}^{k-1} - t_{k}f'(\mathbf{x}_{k-1})), \quad f'(\mathbf{x}_{k-1}) \in \partial f(\mathbf{x}^{k-1})$$

 $t_k > 0$ - stepsize, $P_C(\cdot)$ - orthogonal projection operator.

Orthogonal Projection Operator:

$$P_C(\mathbf{x}) = \text{ closest point in } C \text{ to } \mathbf{x} = \underset{\mathbf{y} \in C}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{x}\|.$$

Amir Beck - Tel Aviv University

Projected Subgradient Method

Model: f - convex. C - closed convex

$\min\{f(\mathbf{x}):\mathbf{x}\in C\}$

Projected Subgradient Method (PSM): Shor (63), Polyak (65)

$$\mathbf{x}^{k} = P_{\mathcal{C}}(\mathbf{x}^{k-1} - t_{k}f'(\mathbf{x}_{k-1})), \quad f'(\mathbf{x}_{k-1}) \in \partial f(\mathbf{x}^{k-1})$$

 $t_k > 0$ - stepsize, $P_C(\cdot)$ - orthogonal projection operator.

Orthogonal Projection Operator: $P_C(\mathbf{x}) = \text{ closest point in } C \text{ to } \mathbf{x} = \arg\min_{\mathbf{x}} ||\mathbf{y} - \mathbf{x}||.$

SPM is not a descent method.

►
$$t_k \propto \frac{1}{\sqrt{k}} \Rightarrow f_{\text{best}}^k := \min_{1 \le s \le k} f(\mathbf{x}_s) \to f_{\text{opt}}$$

Amir Beck - Tel Aviv University

The Gradient Method: Past and Present

 $\mathbf{v} \in C$

 $z^* = P_C(x)$

Rate of Convergence of SPM

A typical result: assume C convex compact. Take

$$t_{k} = \frac{\text{Diam}(C)}{\sqrt{k}}; \text{ Diam}(C) := \max_{\mathbf{x}, \mathbf{y} \in C} \|\mathbf{x} - \mathbf{y}\| < \infty,$$

Then,
$$\min_{1 \le s \le k} f(\mathbf{x}_{s}) - f_{*} \le O(1)M \frac{\text{Diam}(C)}{\sqrt{k}}$$

• Thus, to find an approximate ε solution: $O(1/\varepsilon^2)$

Amir Beck - Tel Aviv University

The Gradient Method: Past and Present

Rate of Convergence of SPM

A typical result: assume C convex compact. Take

$$t_{k} = \frac{\text{Diam}(C)}{\sqrt{k}}; \text{ Diam}(C) := \max_{\mathbf{x}, \mathbf{y} \in C} \|\mathbf{x} - \mathbf{y}\| < \infty,$$

Then,
$$\min_{1 \le s \le k} f(\mathbf{x}_{s}) - f_{*} \le O(1)M \frac{\text{Diam}(C)}{\sqrt{k}}$$

- Thus, to find an approximate ε solution: $O(1/\varepsilon^2)$
- Key Advantages: rate nearly *independent* of problem's dimension. Simple, when projections are easy to compute...
- Main Drawback of SPM: too slow...needs $k \ge \varepsilon^{-2}$ iterations.
- Can we improve the situation of SPM? by exploiting the structure/geometry of the constraint set C.







Mirror Descent: Non-Euclidean Version of SD

- Originated from functional analytic arguments in infinite dimensional setting between primal-dual spaces. Nemirovsky and Yudin (83)
- In (B.-Teboulle-2003) it was shown that the (MDA) can be simply viewed as a Non-Euclidean projected subgradient method.

The Idea

Another representation of the projected subgradient method:

$$\mathbf{x}^{k+1} = \operatorname*{argmin}_{\mathbf{x} \in C} \left\{ f(\mathbf{x}^k) + \langle f'(\mathbf{x}^k), \mathbf{x} - \mathbf{x}^k \rangle + \frac{1}{2t_k} \|\mathbf{x} - \mathbf{x}^k\|_2^2 \right\}$$

Next iterate is a minimizer of the linear approximation regularized by a prox term.

The Idea

Another representation of the projected subgradient method:

$$\mathbf{x}^{k+1} = \operatorname*{argmin}_{\mathbf{x} \in C} \left\{ f(\mathbf{x}^k) + \langle f'(\mathbf{x}^k), \mathbf{x} - \mathbf{x}^k \rangle + \frac{1}{2t_k} \|\mathbf{x} - \mathbf{x}^k\|_2^2 \right\}$$

Next iterate is a minimizer of the linear approximation regularized by a prox term.

The Idea: Replace the Euclidean distance by a non-Euclidean function:

$$\mathbf{x}^{k+1} = \operatorname*{argmin}_{\mathbf{x} \in C} \left\{ f(\mathbf{x}) + \langle f'(\mathbf{x}^k), \mathbf{x} - \mathbf{x}^k \rangle + \frac{1}{t_k} D(\mathbf{x}, \mathbf{x}^k) \right\}$$

The Idea

Another representation of the projected subgradient method:

$$\mathbf{x}^{k+1} = \operatorname*{argmin}_{\mathbf{x} \in C} \left\{ f(\mathbf{x}^k) + \langle f'(\mathbf{x}^k), \mathbf{x} - \mathbf{x}^k \rangle + \frac{1}{2t_k} \|\mathbf{x} - \mathbf{x}^k\|_2^2 \right\}$$

Next iterate is a minimizer of the linear approximation regularized by a prox term.

The Idea: Replace the Euclidean distance by a non-Euclidean function:

$$\mathbf{x}^{k+1} = \operatorname*{argmin}_{\mathbf{x} \in C} \left\{ f(\mathbf{x}) + \langle f'(\mathbf{x}^k), \mathbf{x} - \mathbf{x}^k \rangle + \frac{1}{t_k} D(\mathbf{x}, \mathbf{x}^k) \right\}$$

What should we expect from $D(\cdot, \cdot)$?

- Take into account the structure of the constraints and "easy to compute".
- "distance-like": $D(\mathbf{u}, \mathbf{v}) \ge 0$ and equal zero iff $\mathbf{u} = \mathbf{v}$.
- Popular choice: Bregman distance

 $D(\mathbf{u}, \mathbf{v}) = B_{\omega}(\mathbf{u}, \mathbf{v}) = \omega(\mathbf{u}) - \omega(\mathbf{v}) - \nabla \omega(\mathbf{v})^{T}(\mathbf{u} - \mathbf{v})$ strongly convex w.r.t. to an arbitrary norm.

Demo - Trust Topology Design

Design a truss of a given total weight capable to withstand a collection of forces acting on the nodes. Simplex constraints.

$$\min_{\mathbf{t}} \{ \mathbf{f}^{\mathsf{T}} A^{-1}(\mathbf{t}) \mathbf{f} : \mathbf{t} \in \Delta_n \}$$

Comparing PSM with mirror descant $(\omega(\mathbf{x}) = \frac{1}{2} ||\mathbf{x}||_2^2, \sum_{i=1}^n x_i \log x_i)$

$$x_i^{k+1} = \frac{x_i^k e^{-t_k f_i'(\mathbf{x}^k)}}{\sum_{j=1}^n x_j^k e^{-t_k f_j'(\mathbf{x}^k)}}, \quad i = 1, 2, \dots, n.$$

Theoretically the efficiency estimate is still of the order $O(1/\sqrt{k})$ but the constants can be improved by using non-Euclidean distances.







Dual Projected Subgradient Method Model: $f_{opt} = \min_{x} f(x)$

$$egin{aligned} & \mathsf{f}(\mathbf{x}) \ & \mathsf{s.t.} & \mathbf{g}(\mathbf{x}) \leq \mathbf{0}, \ & \mathbf{x} \in X. \end{aligned}$$

Assumptions:

- (A) $X \subseteq \mathbb{R}^n$ is convex.
- (B) $f : \mathbb{R}^n \to \mathbb{R}$ is convex.
- (C) $\mathbf{g}(\cdot) = (g_1(\cdot), g_2(\cdot), \dots, g_m(\cdot))^T$, where $g_1, g_2, \dots, g_m : \mathbb{R}^n \to \mathbb{R}$ are convex.
- (D) For any $\lambda \in \mathbb{R}^m_+$, the problem $\min_{\mathbf{x} \in X} \{f(\mathbf{x}) + \lambda^T \mathbf{g}(\mathbf{x})\}$ attains an optimal solution.

The Lagrangian of the problem is given by

$$L(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{g}(\mathbf{x}).$$

Amir Beck - Tel Aviv University

The Dual Problem

(D)
$$q_{\mathrm{opt}} \equiv \max\{q(\boldsymbol{\lambda}) : \boldsymbol{\lambda} \in \mathbb{R}^m_+\},$$

where

$$q(\boldsymbol{\lambda}) = \min_{\mathbf{x} \in X} f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{g}(\mathbf{x}).$$

Under the assumptions, strong duality holds, meaning that f_{opt} = q_{opt} and the optimal solution of the dual problem is attained.

The Method

Main Observation: To compute a subgradient of -q at λ :

- Find $\mathbf{x}_{\lambda} \in \underset{\mathbf{x} \in X}{\operatorname{argmin}} L(\mathbf{x}, \boldsymbol{\lambda}).$
- ► $-\mathbf{g}(\mathbf{x}_{\lambda}) \in \partial(-q)(\boldsymbol{\lambda}).$

The Dual Projected Subgradient Method Initialization: pick $\lambda^0 \in \mathbb{R}^m_+$ arbitrarily. General step: for any k = 0, 1, 2, ...,

(a) pick a positive number γ_k .

(b) compute $\mathbf{x}^k \in \underset{\mathbf{x} \in X}{\operatorname{argmin}} \left\{ f(\mathbf{x}) + (\boldsymbol{\lambda}^k)^T \mathbf{g}(\mathbf{x}) \right\}.$

(c) if $\mathbf{g}(\mathbf{x}^k) = \mathbf{0}$, then **terminate** with an output \mathbf{x}^k ; otherwise,

$$oldsymbol{\lambda}^{k+1} = \left[oldsymbol{\lambda}^k + \gamma_k rac{\mathbf{g}(\mathbf{x}^k)}{\|\mathbf{g}(\mathbf{x}^k)\|_2}
ight]_+$$

 $O(1/\sqrt{k})$ rate of convergence can be shown Amir Beck - Tel Aviv University The Gradient Method: Past and Present

$O(1/\varepsilon^2)$ Rate of Convergence in Nonsmooth Convex Optimization

- SPM,MD and dual projected subgradient are all O(1/ε²), O(1/√k) methods. Can we do better?
- According to lower complexity bounds, the answer is No!
- However, by exploiting the structure of the functions, we can do better. For example, if assuming some smoothness properties...
Polynomial versus Gradient-Based Methods (80's and 90's)

- Rise of Polynomial methods for convex programming: ellipsoid, interior point methods.
- Convex problems are polynomially solvable within ε accuracy:

Running Time $\leq Poly$ (Problem's size, # of accuracy digits).

- Theoretically: large scale problems can be solved to high accuracy with polynomial methods, such as IPM.
- Practically: Running time is dimension-dependent and grows nonlinearly with problem's dimension. For IPM which are Newton's type methods: ~ O(n³).
- Thus, a "single iteration" of IPM can last forever!
- 2000-... Gradient-based method have become popular again due to increasing size of applications.

One way to deal with the large or even huge-scale size of the new arising applications is use decomposition. For example...

- One way to deal with the large or even huge-scale size of the new arising applications is use decomposition. For example...
- Consider the problem

$$\min_{\mathbf{x}} \left\{ f(\mathbf{x}) \equiv \sum_{i=1}^{m} f_i(\mathbf{x}) \right\}$$

where f_1, f_2, \ldots, f_m are all convex functions. Suppose that *m* is huge.

- One way to deal with the large or even huge-scale size of the new arising applications is use decomposition. For example...
- Consider the problem

$$\min_{\mathbf{x}} \left\{ f(\mathbf{x}) \equiv \sum_{i=1}^{m} f_i(\mathbf{x}) \right\}$$

where f_1, f_2, \ldots, f_m are all convex functions. Suppose that *m* is huge. • The subgradient method is very expansive to execute:

$$\mathbf{x}^{k+1} = \mathbf{x}^k - t_k \left(\sum_{i=1}^m f'_i(\mathbf{x}^k) \right).$$

Amir Beck - Tel Aviv University

The Gradient Method: Past and Present

- One way to deal with the large or even huge-scale size of the new arising applications is use decomposition. For example...
- Consider the problem

$$\min_{\mathbf{x}} \left\{ f(\mathbf{x}) \equiv \sum_{i=1}^{m} f_i(\mathbf{x}) \right\}$$

where f_1, f_2, \ldots, f_m are all convex functions. Suppose that *m* is huge. • The subgradient method is very expansive to execute:

$$\mathbf{x}^{k+1} = \mathbf{x}^k - t_k \left(\sum_{i=1}^m f'_i(\mathbf{x}^k) \right).$$

 Instead, we can use the stochastic projected subgradient method that exploits only one randomlly chosen subgradient at each iteration (decomposition)

$$\mathbf{x}^{k+1} = \mathbf{x}^k - t_k f'_{i_k}(\mathbf{x}^k)$$

 i_k - randomly chosen Amir Beck - Tel Aviv University The Gradient Method: Past and Present





The General Composite Model

We will be interested in the following model:

$$P) \qquad \min\{F(\mathbf{x}) \equiv f(\mathbf{x}) + g(\mathbf{x}) : \mathbf{x} \in \mathbb{E}\}.$$

• $f : \mathbb{R}^n \to \mathbb{R}$ is an L_f -smooth convex functin:

 $\|
abla f(\mathbf{x}) -
abla f(\mathbf{y})\| \leq L_f \|\mathbf{x} - \mathbf{y}\| \quad ext{for every } \mathbf{x}, \mathbf{y} \in \mathbb{R}^n,$

- g: ℝⁿ → ℝ ∪ {∞} extended valued convex function which is nonsmooth.
- Problem (P) is solvable, i.e., X_{*} := argmin f ≠ Ø, and for x^{*} ∈ X_{*} we set F_{opt} := F(x^{*}).

Amir Beck - Tel Aviv University The Gradient Method: Past and Present

Special Cases of the General Model

• g = 0 - smooth unconstrained convex minimization.

 $\min_{\mathbf{x}} f(\mathbf{x})$

Special Cases of the General Model

• g = 0 - smooth unconstrained convex minimization.

$\min_{\mathbf{x}} f(\mathbf{x})$

• $g = \delta_C(\cdot)$ - constrained smooth convex minimization.

 $\min_{\mathbf{x}} \{ f(\mathbf{x}) : \mathbf{x} \in C \}$

Special Cases of the General Model

• g = 0 - smooth unconstrained convex minimization.

$\min_{\mathbf{x}} f(\mathbf{x})$

• $g = \delta_C(\cdot)$ - constrained smooth convex minimization.

 $\min_{\mathbf{x}} \{ f(\mathbf{x}) : \mathbf{x} \in C \}$

• $g = \| \cdot \|_1$ - l_1 -regularized convex minimization.

 $\min_{\mathbf{x}} \left\{ f(\mathbf{x}) + \lambda \|\mathbf{x}\|_1 \right\}$

Amir Beck - Tel Aviv University

The Proximal Gradient Method

The derivation of the proximal gradient method is similar to the one of the projected subgradient method.

For any $L \ge L_f$, and a given iterate \mathbf{x}^k :

 $Q_L(\mathbf{x}, \mathbf{x}^k) := f(\mathbf{x}^k) + \langle \mathbf{x} - \mathbf{x}^k, \nabla f(\mathbf{x}^k) \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{x}^k\|^2 + \underbrace{g(\mathbf{x})}_{\text{untouched}}$

The Proximal Gradient Method

The derivation of the proximal gradient method is similar to the one of the projected subgradient method.

For any $L \ge L_f$, and a given iterate \mathbf{x}^k :

$$Q_L(\mathbf{x}, \mathbf{x}^k) := f(\mathbf{x}^k) + \langle \mathbf{x} - \mathbf{x}^k, \nabla f(\mathbf{x}^k) \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{x}^k\|^2 + \underbrace{g(\mathbf{x})}_{\text{untouched}}$$

► Algorithm:

Amir Bec

$$\mathbf{x}^{k+1} := \underset{\mathbf{x}}{\operatorname{argmin}} Q_{L}(\mathbf{x}, \mathbf{x}^{k})$$

$$= \underset{\mathbf{x}}{\operatorname{argmin}} \left\{ g(\mathbf{x}) + \frac{L}{2} \left\| \mathbf{x} - (\mathbf{x}^{k} - \frac{1}{L} \nabla f(\mathbf{x}^{k})) \right\|^{2} \right\}$$

$$= \operatorname{prox}_{\frac{1}{L}g} \left(\mathbf{x}^{k} - \frac{1}{L} \nabla f(\mathbf{x}^{k}) \right) \equiv p_{L}(\mathbf{x}^{k}).$$
prox operator:
$$\operatorname{prox}_{g}(\mathbf{x}) := \underset{\mathbf{u}}{\operatorname{argmin}} \left\{ g(\mathbf{u}) + \frac{1}{2} \| \mathbf{u} - \mathbf{x} \|^{2} \right\}.$$
Beck - Tel Aviv University The Gradient Method: Past and Present

32 / 45

Special Cases

The general method: $\mathbf{x}^{k+1} = \operatorname{prox}_{\frac{1}{L}g} \left(\mathbf{x}^k - \frac{1}{L} \nabla f(\mathbf{x}^k) \right).$

• $g \equiv 0 \Rightarrow$ the gradient method.

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \frac{1}{L} \nabla f(\mathbf{x}^k)$$

Special Cases

The general method: $\mathbf{x}^{k+1} = \operatorname{prox}_{\frac{1}{l}g} \left(\mathbf{x}^k - \frac{1}{L} \nabla f(\mathbf{x}^k) \right)$.

• $g \equiv 0 \Rightarrow$ the gradient method.

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \frac{1}{L} \nabla f(\mathbf{x}^k)$$

• $g = \delta_{\mathcal{C}}(\cdot) \Rightarrow$ the gradient projection method

$$\mathbf{x}^{k+1} = P_C\left(\mathbf{x}^k - \frac{1}{L}\nabla f(\mathbf{x}^k)\right)$$

Special Cases

The general method: $\mathbf{x}^{k+1} = \operatorname{prox}_{\frac{1}{l}g} \left(\mathbf{x}^k - \frac{1}{L} \nabla f(\mathbf{x}^k) \right)$.

• $g \equiv 0 \Rightarrow$ the gradient method.

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \frac{1}{L} \nabla f(\mathbf{x}^k)$$

• $g = \delta_{\mathcal{C}}(\cdot) \Rightarrow$ the gradient projection method

$$\mathbf{x}^{k+1} = P_C\left(\mathbf{x}^k - \frac{1}{L}\nabla f(\mathbf{x}^k)\right)$$

► $g(\mathbf{x}) := \lambda \|\mathbf{x}\|_1 \Rightarrow$ Iterative shrinkage/thresholding algorithm $\mathbf{x}^{k+1} = \mathcal{T}_{\lambda/L} \left(\mathbf{x}^k - \frac{1}{L} \nabla f(\mathbf{x}^k) \right)$

and $\mathcal{T}_{\alpha} : \mathbb{R}^n \to \mathbb{R}^n$ is the shrinkage operator defined by $\mathcal{T}_{\alpha}(\mathbf{x})_i = (|x_i| - \alpha)_+ \operatorname{sgn}(x_i).$

Amir Beck - Tel Aviv University

The Gradient Method: Past and Present

Special Case: LASSO

•
$$g(\mathbf{x}) := \lambda \|\mathbf{x}\|_1$$
, $f(\mathbf{x}) := \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$ (prox=shrinkage).

$$\mathbf{x}_{k+1} = \mathcal{T}_{\lambda/L}\left(\mathbf{x}_k - \frac{2}{L}\mathbf{A}^T(\mathbf{A}\mathbf{x}_k - \mathbf{b})\right)$$

ISTA - Iterative Shrinkage/Thresholding Algorithm

In SP literature: Chambolle (98); Figueiredo-Nowak (03, 05); Daubechies et al. (04),Elad et al. (06), Hale et al. (07)...

In Optimization: can be viewed as the Proximal forward backward Splitting Method (Passty (79))

Amir Beck - Tel Aviv University

Prox Computations

There are a quite a few "simple" functions for which the prox can be easily computed

Appendix B. Tables

⊕ 2016/
 page
 ⊕

ŧ

Prox Computations				
1	$\operatorname{dom}(f)$	$peox_f(\mathbf{x})$	assumptions	reference
$\frac{1}{2}\mathbf{x}^T \mathbf{A}\mathbf{x} + \mathbf{b}^T \mathbf{x} + c$	R*	$(\mathbf{A} + \mathbf{I})^{-1}(\mathbf{x} - \mathbf{b})$	$\mathbf{A} \in S_{+}^{n}, \mathbf{b} \in \mathbb{R}^{n}, c \in \mathbb{R}$	Section 6.2.3
λx^3	R+	-1+_/1+123.[#] + 63	$\lambda > 0$	Lemma 6.5
μχ	$[0, \alpha]$	$\min\{\max\{x - \mu, 0\}, \alpha\}$	$\mu\in\mathbb{R},\alpha\in\mathbb{R}+$	Example 6.14
$\lambda \ \mathbf{x}\ $	2	$\left(1 - \frac{\lambda}{\max\{ \mathbf{x} ,\lambda\}}\right) \mathbf{x}$	$\ \cdot\ $ - Euclidean norm, $\lambda>0$	Example 6.19
$-\lambda \ \mathbf{x}\ $	2	$\left(1 + \frac{\lambda}{ w }\right) \mathbf{x}, \mathbf{x} \neq 0,$ $\{\mathbf{u} : \mathbf{u} = \lambda\}, \mathbf{x} = 0.$	$\ \cdot\ $ - Euclidean norm, $\lambda>0$	Example 6.21
$\lambda \ \mathbf{x}\ _1$	2"	$\mathcal{T}_{\lambda}(\mathbf{x}) \equiv [\mathbf{x} - \lambda \mathbf{s}]_{+} \odot \operatorname{sgn}(\mathbf{x})$	$\lambda > 0$	Example 6.8
$\ \boldsymbol{\omega} \odot \mathbf{x}\ _1$	$\operatorname{Box}[-\alpha,\alpha]$	$\mathcal{S}_{\omega,m}(\mathbf{x})$	$m \in [0, \infty)^n, \omega \in \mathbb{R}^n_{++}$	Example 6.23
$\lambda \ \mathbf{x}\ _{\infty}$	R*	$\mathbf{x} - \lambda P \pi_{\ \cdot\ _{1}[0,1]}(\mathbf{x}/\lambda)$	$\lambda > 0$	Example 6.48
$\lambda \ \mathbf{x}\ _{n}$	2	$\mathbf{x} - \lambda P_{H_{\ \cdot\ _{H, \mathcal{R}}}[0, 1]}(\mathbf{x}/\lambda)$	$\label{eq:states} \begin{array}{ll} \ \mathbf{x}\ _n & - \mbox{ norm}, \\ \lambda > 0 \end{array}$	Example 6.47
$\lambda \ \mathbf{x}\ _0$	2"	$H_{\sqrt{2\lambda}}(x_1) \times \cdots \times H_{\sqrt{2\lambda}}(x_n)$	$\lambda > 0$	Example 6.10
$\lambda \ \mathbf{x}\ ^3$	2	$\frac{2}{1+\sqrt{1+12\lambda\ \mathbf{s}\ }}\mathbf{x}$	$\ \cdot\ $ - Euclidean norm, $\lambda > 0$,	Example 6.20
$-\lambda \sum_{j=1}^n \log x_j$	\mathbb{R}^{n}_{++}	$\left(\frac{x_j + \sqrt{x_j^2 + 4x}}{2}\right)_{i=1}^n$	$\lambda > 0$	Example 6.9
$\delta_C(\mathbf{x})$	2	$P_C(\mathbf{x})$	$\emptyset \not = C \subseteq \mathbb{R}$	Theorem 6.24
$\lambda \sigma_C(\mathbf{x})$	E	$\mathbf{x} - \lambda P_C(\mathbf{x}/\lambda)$	$\lambda > 0, C \neq \emptyset$ closed convex	Theorem 6.46
$\lambda \max{x_i}$	R.	$\mathbf{x} - P_{\Delta_n}(\mathbf{x}/\lambda)$	$\lambda > 0$	Example 6.49
$\lambda \sum_{i=1}^{k} x_{[i]}$	ň	$\mathbf{x} - \lambda P_C(\mathbf{x}/\lambda),$ $C \equiv H_{\mathbf{z},h} \cap \text{Box}[0, \mathbf{e}]$	$\lambda > 0$	Example 6.50
$\lambda \sum_{i=1}^{k} x_{(i)} $	2"	$\begin{split} \mathbf{x} & -\lambda P_C(\mathbf{x}/\lambda), \\ C &= B_{\ \cdot\ _1}[0,k] \cap \mathrm{Box}[-\mathbf{e},\mathbf{e}] \end{split}$	$\lambda > 0$	Example 6.51
$\lambda M_f^{\mu}(\mathbf{x})$	E	$\frac{\mathbf{x}+}{\frac{\lambda}{\mu+\lambda}}\left(\mathrm{prox}_{(\mu+\lambda)f}(\mathbf{x})-\mathbf{x}\right)$	$\begin{array}{llllllllllllllllllllllllllllllllllll$	Corollary 6.63
$\lambda d_C(\mathbf{x})$	2	$\min\left\{\frac{\mathbf{x}}{\frac{1}{d_C(\mathbf{x})}},1\right\}(P_C(\mathbf{x})-\mathbf{x})$	$\begin{array}{ll} C & \text{nonempty} \\ \text{closed} & \text{convex}, \\ \lambda > 0 \end{array}$	Lemma 6.43
$\frac{1}{2}d_C^2(\mathbf{x})$	я	$\frac{\lambda}{\lambda+1}P_{U}(\mathbf{x}) + \frac{1}{\lambda+1}\mathbf{x}$	$\begin{array}{ll} C & \text{nonempty} \\ \text{closed} & \text{convex}, \\ \lambda > 0 \end{array}$	Example 6.64
$\lambda H_{\mu}(\mathbf{x})$	£	$\left(1 - \frac{\lambda}{\max\{ \mathbf{s} , \mathbf{s} + \lambda\}}\right)$	$\lambda, \mu > 0$	Example 6.65
$\rho \ \mathbf{x}\ _{1}^{2}$	R	$ \begin{array}{c} \left(\frac{v_1 v_1}{v_1 + \lambda p}\right)_{i=1}^n, \mathbf{v} = \\ \left[\sqrt{\frac{x}{p}} \mathbf{x} - 2\rho\right]_+, \mathbf{e}^T \mathbf{v} = 1 \ (0 \\ \text{when } \mathbf{x} = 0) \end{array} $	$\rho > 0$	Lemma 6.69
Ax 2	R*	$\mathbf{x} - \mathbf{A}^T (\mathbf{A} \mathbf{A}^T + \alpha^* \mathbf{I})^{-1} \mathbf{A} \mathbf{x},$ $\alpha^* \equiv 0 \text{ if } \ \mathbf{v}_0\ _2 \leq \lambda; \text{ otherwise, } \ \mathbf{v}_{\alpha^*}\ _2 \equiv \lambda; \mathbf{v}_{\alpha} \equiv (\mathbf{A} \mathbf{A}^T + \alpha \mathbf{I})^{-1} \mathbf{A} \mathbf{x}$	$\begin{array}{llllllll} \mathbf{A} & \in & \mathbb{R}^{m \times n} \\ & \text{with} & \text{full row} \\ & \text{rank} \end{array}$	Lemma 6.67

+

Prox of Symmetric Spectral Functions over S^n (From Example 7.19)

⊕ 2016/
 page

$F(\mathbf{X})$	$\operatorname{dom}(F)$	$prox_{p}(\mathbf{X})$
$a \ \mathbf{X}\ _{F}^{2}$	Su	1+2mX
$\alpha \ \mathbf{X}\ _{F}$	3n	$\left(1 - \frac{\alpha}{\max\{ \mathbf{X} _{F}, \alpha\}}\right) \mathbf{X}$
$\alpha \ \mathbf{X}\ _{Z_1}$	S*	$UT_n(\lambda(X))U^T$
$a \ \mathbf{X} \ _{2,2}$	S*	$Udiag(\lambda(\mathbf{X}) - \alpha P_{B_{[\cdot]}, [0,1]}(\lambda(\mathbf{X})/\alpha))\mathbf{U}^{T}$
$-\alpha \det(\mathbf{X})$	S_{++}^n	$\operatorname{Udiag}\left(\frac{\lambda_j(\mathbf{x})+\sqrt{\lambda_j(\mathbf{x})^2+4w}}{2}\right)\mathbf{U}^T$
$\alpha \lambda_1(\mathbf{X})$	S*	$U \operatorname{diag}(\lambda(\mathbf{X}) - P_{\Delta n}(\lambda(\mathbf{X})/\alpha))\mathbf{U}^T$
$\alpha \sum_{i=1}^{k} \lambda_i(\mathbf{X})$	S.	$\mathbf{X} - \alpha \mathbf{U} P_C(\mathbf{\lambda}(\mathbf{X})/\alpha) \mathbf{U}^T$, $C \equiv H_{\mathbf{e},k} \cap Box[0, \mathbf{e}]$

Prox of Symmetric Spectral Functions over $\mathbb{R}^{m \times n}$ (From Example 7.30)

$F(\mathbf{X})$	$prox_{p}(\mathbf{X})$
$\alpha \ \mathbf{X}\ _{F}^{2}$	1 x
$\alpha \ \mathbf{X}\ _{F}$	$\left(1 - \frac{\alpha}{\max\{ \mathbf{X} _{F,\alpha}\}}\right) \mathbf{X}$
$a \ \mathbf{X}\ _{S_1}$	$UT_{\alpha}(\sigma(\mathbf{X}))\mathbf{V}^{T}$
$\alpha \ \mathbf{X}\ _{\mathbb{Z}_{\infty}}$	$\mathbf{X} = \alpha \operatorname{Udiag}(P_{B_{\frac{1}{2},\frac{1}{2},\frac{1}{2}}(\sigma(\mathbf{X})/\alpha))\mathbf{V}^{T}$
$\alpha \ \mathbf{X}\ _{(k)}$	$\mathbf{X} - \alpha \mathbf{U} P_C(\boldsymbol{\sigma}(\mathbf{X})/\alpha) \mathbf{V}^T$,
	$C \equiv B_{1 \cdot 1_1}[0, k] \cap B_{1 \cdot 1_{\infty}}[0, 1]$

Orthogonal Projections			
set (C)	$P_C(\mathbf{x})$	assumptions	reference
2 ⁿ	[x] ₊	-	Lemma 6.26
Box[t, u]	$P_C(\mathbf{x})_i \equiv \min\{\max\{x_i, \ell_i\}, u_i\}$	$\ell_i \leq u_i$	Lemma 6.25
$B_{\parallel \cdot \parallel_2}[\mathbf{c}, r]$	$e + \frac{r}{\max\{ x-x \ge r\}}(x-e)$	$\mathbf{c} \in \mathbb{R}^n, r > 0$	Lemma 6.26
$\{\mathbf{x}:\mathbf{A}\mathbf{x}=\mathbf{b}\}$	$\mathbf{x} - \mathbf{A}^T (\mathbf{A} \mathbf{A}^T)^{-1} (\mathbf{A} \mathbf{x} - \mathbf{b})$	$\begin{array}{llllllllllllllllllllllllllllllllllll$	Lemma 6.26
$\{\mathbf{x} : \mathbf{a}^T \mathbf{x} \le b\}$	$\mathbf{x} = \frac{[\mathbf{a}^* \mathbf{x} = \mathbf{b}]_+}{\ \mathbf{a}\ ^2} \mathbf{a}$	$0 \neq \mathbf{a} \in \mathbb{R}^{n}, b \in \mathbb{R}$	Lemma 6.26
Δ_n	$[\mathbf{x} - \boldsymbol{\mu}^* \mathbf{e}]_+$ where $\boldsymbol{\mu}^* \in \mathbb{R}$ satisfies $\mathbf{e}^T [\mathbf{x} - \boldsymbol{\mu}^* \mathbf{e}]_+ \equiv 1$		Corollary 6.29
$H_{\mathbf{a},b} \cap \operatorname{Box}[t,\mathbf{u}]$	$P_{\text{Hom}(I,\mathbf{u})}(\mathbf{x} - \mu^* \mathbf{a})$ where $\mu^* \in \mathbb{R}$ sat- isfies $\mathbf{a}^T P_{\text{Hom}(I,\mathbf{u})}(\mathbf{x} - \mu^* \mathbf{a}) = b$	$\substack{\mathbf{a} \in \mathbb{R}^n \setminus \{0\}, b \in \\ \mathbb{R}}$	Theorem 6.27
$H^{\mathbf{a},b}\cap \mathrm{Box}[\boldsymbol{\ell},\mathbf{u}]$	$\begin{cases} P_{\text{Hes}[\ell,u]}(\mathbf{x}), & \mathbf{a}^T \mathbf{v}_u \leq b, \\ P_{\text{Hes}[\ell,u]}(\mathbf{x} - \lambda^* \mathbf{n}), & \mathbf{a}^T \mathbf{v}_u > b, \\ \mathbf{v}_u = P_{\text{Hes}[\ell,u]}(\mathbf{x}), & \mathbf{a}^T P_{\text{Hes}[\ell,u]}(\mathbf{x} - \lambda^* \mathbf{n}) = b, \lambda^* > 0 \end{cases}$	$\mathbf{a} \in \mathbb{R}^n \setminus \{0\}, \delta \in \mathbb{R}$	Example 6.32
$B_{1\cdot 11}[0,\alpha]$	$\begin{cases} \mathbf{x}, & \ \mathbf{x}\ _1 \leq \alpha, \\ \mathcal{T}_{\lambda^+}(\mathbf{x}), & \ \mathbf{x}\ _1 > \alpha, \\ \ \mathcal{T}_{\lambda^+}(\mathbf{x})\ _1 = \alpha, \lambda^+ > 0 \end{cases}$	$\alpha > 0$	Example 6.33
$\{\mathbf{x} : \boldsymbol{\omega}^T \mathbf{x} \leq \beta, \\ -\boldsymbol{\alpha} \leq \mathbf{x} \leq \boldsymbol{\alpha}\}$	$\begin{cases} \mathbf{v}_{\mathbf{u}}, & \boldsymbol{\omega}^T \mathbf{v}_{\mathbf{u}} \leq \beta, \\ \boldsymbol{\mathcal{S}}_{\lambda^+ \boldsymbol{\omega}, \mathbf{u}}(\mathbf{x}), & \boldsymbol{\omega}^T \mathbf{v}_{\mathbf{u}} > \beta, \\ \mathbf{v}_{\mathbf{u}} &= P_{\text{Him}(-\mathbf{u}, \mathbf{u})}(\mathbf{x}), \\ \boldsymbol{\omega}^T \boldsymbol{\mathcal{S}}_{\lambda^+ \boldsymbol{\omega}, \mathbf{u}}(\mathbf{x}) = \beta, \lambda^+ > 0 \end{cases}$	$\begin{array}{l} \boldsymbol{\omega} \in \mathbb{R}^n_{++}, \ \boldsymbol{\alpha} \in \\ [0,\infty]^n, \ \boldsymbol{\beta} \in \\ \mathbb{R}_{++} \end{array}$	Example 6.34
$\{\mathbf{x} > 0 : \Pi x_i \geq \alpha\}$	$\begin{cases} \mathbf{x}, & \mathbf{x} \in C, \\ \left(\frac{x_j + \sqrt{x_j^2 + 4\lambda^*}}{2}\right)^n, & \mathbf{x} \notin C, \\ \Pi_{j=1}^n \left((x_j + \sqrt{x_j^2 + 4\lambda^*})/2\right) &= \\ a_i \lambda^* > 0 \end{cases}$	$\alpha > 0$	Example 6.35
$\{(\mathbf{x},s): \ \mathbf{x}\ _2 \leq s\}$	$\left(\frac{ \mathbf{a} _2+\epsilon}{2 \mathbf{a} _2}\mathbf{x}, \frac{ \mathbf{a} _2+\epsilon}{2}\right)$ if $\ \mathbf{x}\ _2 \ge s $ $(0, 0)$ if $s < \ \mathbf{x}\ _2 < -s$, (\mathbf{x}, s) if $\ \mathbf{x}\ _2 \le s$.		Example 6.37
$\{(\mathbf{x},s):\ \mathbf{x}\ _1\leq s\}$	$\begin{cases} (\mathbf{x}, s), & \ \mathbf{x}\ _{1} \leq s, \\ (\mathcal{T}_{\lambda^{+}}(\mathbf{x}), s + \lambda^{*}), & \ \mathbf{x}\ _{1} > s, \\ \ \mathcal{T}_{\lambda^{+}}(\mathbf{x})\ _{1} - \lambda^{*} - s = 0, \lambda^{*} > 0 \end{cases}$		Example 6.38

⊕ 2016/ page ⊕

+

<u>+</u>

Appendix B. Tables

⊕ 2016/
 page

Orthogonal Projections onto Symmetric Spectral Sets in S ⁿ		
set (C)	$P_C(\mathbf{X})$	assamptions
S ⁿ ₊	$Udiag([\lambda(X)]_+)U^T$	-
$\{\mathbf{X}: \ell\mathbf{I} \preceq \mathbf{X} \preceq u\mathbf{I}\}$	$Udiag(\mathbf{v})U^T$, $v_i \equiv min\{max\{\lambda_i(\mathbf{X}), \ell\}, u\}$	$\ell \leq u$
$B_{\ \cdot\ _{F}}[0, r]$	mailXicciX	r > 0
$\{\mathbf{X} : Tr(\mathbf{X}) \leq b\}$	$Udiag(v)U^{T},$ $v = \lambda(X) - \frac{[e^{T}\lambda(X)-b]_{+}}{2}e^{-\frac{1}{2}}$	$b \in \mathbb{R}$
Υ.,	$Udiag(\mathbf{v})U^T$, $\mathbf{v} = [\lambda(\mathbf{X}) - \mu^* \mathbf{e}]_+$ where $\mu^* \in \mathbb{R}$ satisfies $\mathbf{e}^T [\lambda(\mathbf{X}) - \mu^* \mathbf{e}]_+ = 1$	-
$B_{\mathbb{D} \otimes \mathbb{Z}_{2}}[0,\alpha]$	$\begin{cases} \mathbf{X}, & \ \mathbf{X}\ _{\mathcal{X}_{1}} \leq \alpha, \\ \mathbf{U}\mathcal{T}_{\lambda^{*}}(\boldsymbol{\lambda}(\mathbf{X}))\mathbf{U}^{T}, & \ \mathbf{X}\ _{\mathcal{X}_{1}} > \alpha, \\ \ \mathcal{T}_{\lambda^{*}}(\boldsymbol{\lambda}(\mathbf{X}))\ _{1} = \alpha, \lambda^{*} > 0 \end{cases}$	$\alpha > 0$

Orthogonal Projection onto Symmetric Spectral Sets in $\mathbb{R}^{m \times n}$ (From Example 7.31)

set (C)	$P_C(\mathbf{X})$	assumptions
$B_{\ \cdot\ _{\mathcal{S}_{\infty}}}\left[0,\alpha\right]$	\mathbf{U} diag $(\mathbf{v})\mathbf{U}^T$, $v_i \equiv \min\{\sigma_i(\mathbf{X}), \alpha\}$	$\alpha > 0$
$B_{\ \cdot\ _{F}}[0, r]$	$\frac{1}{\max\{ \mathbf{X} _{W},\tau\}}\mathbf{X}$	r > 0
$B_{\ \cdot\ _{\mathcal{S}_1}}[0,\alpha]$	$\begin{cases} \mathbf{X}, & \ \mathbf{X}\ _{T_1} \leq \alpha, \\ \mathbf{U}\mathcal{T}_{\lambda^*}(\sigma(\mathbf{X}))\mathbf{U}^T, & \ \mathbf{X}\ _{T_1} > \alpha, \\ \ \mathcal{T}_{\lambda^*}(\sigma(\mathbf{X}))\ _1 = \alpha, \lambda^* > 0 \end{cases}$	$\alpha > 0$

Rate of Convergence of Prox-Grad

Theorem - [Rate of Convergence of Prox-Grad] Let $\{\mathbf{x}^k\}$ be the sequence generated by the proximal gradient method.

$$F(\mathbf{x}^k) - F(\mathbf{x}^*) \le rac{L \|\mathbf{x}^0 - \mathbf{x}^*\|^2}{2k}$$

for any optimal solution \mathbf{x}^* .

- Thus, to solve (M), the proximal gradient method converges at a sublinear rate in function values.
- # iterations for $F(\mathbf{x}^k) F(\mathbf{x}^*) \leq \varepsilon$ is $O(1/\varepsilon)$.

Rate of Convergence of Prox-Grad

Theorem - [Rate of Convergence of Prox-Grad] Let $\{\mathbf{x}^k\}$ be the sequence generated by the proximal gradient method.

$$F(\mathbf{x}^k) - F(\mathbf{x}^*) \leq rac{L \|\mathbf{x}^0 - \mathbf{x}^*\|^2}{2k}$$

for any optimal solution \mathbf{x}^* .

- Thus, to solve (M), the proximal gradient method converges at a sublinear rate in function values.
- # iterations for $F(\mathbf{x}^k) F(\mathbf{x}^*) \leq \varepsilon$ is $O(1/\varepsilon)$.
- Note: The sequence $\{\mathbf{x}^k\}$ can be proven to *converge* to solution \mathbf{x}^* .
- ▶ No need to know the Lipschitz constant (backtracking).

Towards a Faster Algorithm

- An O(1/k) rate of convergence is rather slow.
- Can we find a faster methods?

Towards a Faster Algorithm

- An O(1/k) rate of convergence is rather slow.
- Can we find a faster methods?
- ► The answer is YES!.

FISTA - [B., Teboulle 2009]

An equally simple algorithm as prox-grad. (Here L_f is known).

FISTA with constant stepsize **Input:** $L \ge L_f$ - A Lipschitz constant of ∇f . **Step 0.** Take $y^1 = x^0 \in \mathbb{E}, t_1 = 1$. **Step k.** $(k \ge 1)$ Compute $\mathbf{x}^{k} \equiv \operatorname{prox}_{\frac{1}{L}g}\left(\mathbf{y}^{k} - \frac{1}{I}\nabla f(\mathbf{y}^{k})\right), \iff \operatorname{main \ computation}$ • $t_{k+1} = \frac{1+\sqrt{1+4t_k^2}}{2}$, • $\mathbf{y}^{k+1} = \mathbf{x}^k + \left(\frac{t_k-1}{t_{k+1}}\right)(\mathbf{x}^k - \mathbf{x}^{k-1}).$

Additional computation for FISTA in (\bullet) and $(\bullet\bullet)$ is clearly marginal.

Amir Beck - Tel Aviv University

The Gradient Method: Past and Present

Theorem - Global Rate of Convergence FISTA

Theorem Let $\{\mathbf{x}_k\}$ be generated by FISTA. Then for any $k \ge 1$

$$\mathsf{F}(\mathbf{x}_k) - \mathsf{F}(\mathbf{x}^*) \leq rac{2lpha L(f) \|\mathbf{x}_0 - \mathbf{x}^*\|^2}{(k+1)^2},$$

where $\alpha=1$ for the constant stepsize setting and $\alpha=\eta$ for the backtracking stepsize setting.

- # of iterations to reach $F(\tilde{\mathbf{x}}) F_* \leq \varepsilon$ is $\sim O(1/\sqrt{\varepsilon})$.
- Clearly improves ISTA by a square root factor.

Amir Beck - Tel Aviv University

The Gradient Method: Past and Present

Theorem - Global Rate of Convergence FISTA

Theorem Let $\{\mathbf{x}_k\}$ be generated by FISTA. Then for any $k \ge 1$

$$\mathsf{F}(\mathbf{x}_k) - \mathsf{F}(\mathbf{x}^*) \leq rac{2lpha L(f) \|\mathbf{x}_0 - \mathbf{x}^*\|^2}{(k+1)^2},$$

where $\alpha=1$ for the constant stepsize setting and $\alpha=\eta$ for the backtracking stepsize setting.

- # of iterations to reach $F(\tilde{\mathbf{x}}) F_* \leq \varepsilon$ is $\sim O(1/\sqrt{\varepsilon})$.
- Clearly improves ISTA by a square root factor.
- Do we practically achieve this theoretical rate? Yes

LASSO (Penalized Version)

Consider the problem

(P) min
$$\left\{ f(\mathbf{x}) \equiv \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{x}\|_1 \right\}$$

 $\mathbf{A} \in \mathbb{R}^{100 imes 200}, \mathbf{b} \in \mathbb{R}^{100}, \lambda > 0$

LASSO (Penalized Version)

Consider the problem

(P) min
$$\left\{ f(\mathbf{x}) \equiv \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{x}\|_1 \right\}$$

 $\mathbf{A} \in \mathbb{R}^{100 imes 200}, \mathbf{b} \in \mathbb{R}^{100}, \lambda > 0$



Illustration
$$(\lambda = 1)$$

Amir Beck - Tel Aviv University

The Gradient Method: Past and Present









Smoothed FISTA

Revisit nonsmooth problems:

 $\min_{\mathbf{x}} \{ f(\mathbf{x}) : \mathbf{x} \in C \}$

f - convex nonsmooth, C - convex
Smoothed FISTA

Revisit nonsmooth problems:

 $\min_{\mathbf{x}} \{ f(\mathbf{x}) : \mathbf{x} \in C \}$

f - convex nonsmooth, C - convex

• "standard" nonsmooth algorithms solve it in $O(1/\varepsilon^2)$ complexity.

Smoothed FISTA

Revisit nonsmooth problems:

 $\min_{\mathbf{x}} \{ f(\mathbf{x}) : \mathbf{x} \in C \}$

f - convex nonsmooth, C - convex

- "standard" nonsmooth algorithms solve it in $O(1/\varepsilon^2)$ complexity.
- Another approach: consider a smoothed version of the problem: min_{x∈C} f_η(x) and solve it using FISTA.
- Example: $\sum_{i=1}^{m} |\mathbf{a}_i^T \mathbf{x} b_i| \rightarrow \sum_{i=1}^{m} \sqrt{(\mathbf{a}_i^T \mathbf{x} b_i)^2 + \eta^2}$

Amir Beck - Tel Aviv University

The Gradient Method: Past and Present

Smoothed FISTA

Revisit nonsmooth problems:

 $\min_{\mathbf{x}} \{ f(\mathbf{x}) : \mathbf{x} \in C \}$

f - convex nonsmooth, C - convex

- "standard" nonsmooth algorithms solve it in $O(1/\varepsilon^2)$ complexity.
- Another approach: consider a smoothed version of the problem: min_{x∈C} f_η(x) and solve it using FISTA.
- Example: $\sum_{i=1}^{m} |\mathbf{a}_i^T \mathbf{x} b_i| \rightarrow \sum_{i=1}^{m} \sqrt{(\mathbf{a}_i^T \mathbf{x} b_i)^2 + \eta^2}$
- ► Carefully choosing the smoothing parameter, O(1/ε) complexity can be shown.

Amir Beck - Tel Aviv University





Dual FISTA - FDPG

Model:

$$\min_{\mathbf{x}} f(\mathbf{x}) + g(\mathbf{A}\mathbf{x})$$

Dual FISTA - FDPG

Model:

$$\min_{\mathbf{x}} f(\mathbf{x}) + g(\mathbf{A}\mathbf{x})$$

Dual model:

$$\max_{\mathbf{y}} - f^*(\mathbf{A}^T \mathbf{y}) - g^*(-\mathbf{y})$$

 f^*, g^* - convex conjugates $(h^*(\mathbf{y}) \equiv \max_{\mathbf{x}} \{\mathbf{x}^T \mathbf{y} - h(\mathbf{x})\})$

Amir Beck - Tel Aviv University The Gradient Method: Past and Present

Dual FISTA - FDPG

Model:

$$\min_{\mathbf{x}} f(\mathbf{x}) + g(\mathbf{A}\mathbf{x})$$

Dual model:

$$\max_{\mathbf{y}} - f^*(\mathbf{A}^T \mathbf{y}) - g^*(-\mathbf{y})$$

 f^*, g^* - convex conjugates $(h^*(\mathbf{y}) \equiv \max_{\mathbf{x}} \{\mathbf{x}^T \mathbf{y} - h(\mathbf{x})\})$

- Apply FISTA on the dual.
- Can deal with many different types of problems...

Amir Beck - Tel Aviv University

FDPG

The Fast Dual Proximal Gradient (FDPG) Method - primal representation

Initialization:
$$L \ge L_F = \frac{\|\mathbf{A}\|^2}{\sigma}, \mathbf{w}^0 = \mathbf{y}^0 \in \mathbb{R}^m, t_0 = 1.$$

General step $(k \ge 0)$:
(a) $\mathbf{u}^k = \underset{\mathbf{u}}{\operatorname{argmax}} \left\{ \langle \mathbf{u}, \mathbf{A}^T(\mathbf{w}^k) \rangle - f(\mathbf{u}) \right\}.$
(b) $\mathbf{y}^{k+1} = \mathbf{w}^k - \frac{1}{L} \mathbf{A}(\mathbf{u}^k) + \frac{1}{L} \operatorname{prox}_{Lg}(\mathbf{A}(\mathbf{u}^k) - L\mathbf{w}^k)$
(c) $t_{k+1} = \frac{1 + \sqrt{1 + 4t_k^2}}{2}$
(d) $\mathbf{w}^{k+1} = \mathbf{y}^{k+1} + \left(\frac{t_k - 1}{t_{k+1}}\right) (\mathbf{y}^{k+1} - \mathbf{y}^k).$

Amir Beck - Tel Aviv University

The Gradient Method: Past and Present

Present and Future?

The scale of problems is becoming huge. Emphasis of current and probably near future research:

- Decomposition
- Randomization
- Distributed

Thank You

Any Questions????

