

Last time

- subgradient method
- used for convex opts with non smooth criterion functions
- use pre set step sizes
- Upside: very generic
- Downside: can be slow: convergence rate is  $O(1/\epsilon^2)$ , vs gradient descent, which is  $O(1/\epsilon)$ 
  - For any first order method that is a subgradient method, this is as good as it gets per Nesterov

Can we improve on subgradient method somehow?

- Yes! **proximal gradient descent**

### A. Context: Decomposable functions

Suppose  $f(x) = g(x) + h(x)$

where  $g(x)$  is convex, differentiable,  $\text{dom}(g) = \mathbb{R}^n$

$h(x)$  is convex, not diff

so  $f(x)$  is not diff

If  $f(x)$  were differentiable, the gradient descent update would be

$$x^+ = x - t \cdot \nabla f(x)$$

Recall the motivation for this: minimize the quadratic approximation to  $f$  around  $x$ , replace

hessian with  $\frac{1}{t}I$

$$x^+ = \underset{z}{\operatorname{argmin}} f(x) + \nabla f(x)^T(z-x) + \frac{1}{2t} \|z-x\|_2^2 = \underset{z}{\operatorname{argmin}} \tilde{f}_t(z)$$

In our case,  $f$  is not diff but  $g$  is.

### B. Motivation: Can we take quadratic approx to $g$ and leave $h$ alone?

so we update

$$x^+ = \underset{z}{\operatorname{argmin}} \tilde{g}_t(z) + h(z)$$

$$= \underset{z}{\operatorname{argmin}} g(x) + \nabla g(x)^T(z-x) + \frac{1}{2t} \|z-x\|_2^2 + h(z)$$

first term has no reliance on  $z$

$$= \underset{z}{\operatorname{argmin}} \frac{1}{2t} (z^T z - z^T x - x^T z + x^T x + 2t \nabla g(x)^T z)$$

$$= \operatorname{argmin}_z \frac{1}{2t} \left( z^T z - z^T x - x^T z + x^T x + 2t \nabla g(x)^T z + t \|\nabla g(x)\|_2^2 - t \nabla g(x)^T x - x^T t \nabla g(x) \right)$$

adding terms that have no reliance on  $z$

$$= \operatorname{argmin}_z \frac{1}{2t} \|z - (x - t \nabla g(x))\|_2^2 + h(z)$$

Yielding a nice interpretation: squared norm of ( $z$  - gradient update) +  $h(z)$

$$\operatorname{argmin}_z \frac{1}{2t} \|z - (x - t \nabla g(x))\|_2^2 + h(z)$$

The first term says 'stay close to gradient update for  $g$ '

The second term says 'also make  $h$  small'

- interpretation: find the vector that minimizes the diff from the gradient update for  $g$ , and keeps  $h$  minimal

This is the definition of the **proximal mapping**.

$$\operatorname{prox}_t(x; h) = \operatorname{argmin}_z \left( \frac{1}{2t} \|x - z\|_2^2 + h(z) \right)$$

- interpretation: the proximal mapping is the vector 'closest' to  $x$  that minimizes  $h()$
- Note that this is a minimization problem of its own! Albeit for the most part well defined, at least for many known  $h()$  - see below

Question: is this a well defined mapping? That is, does  $\operatorname{prox}()$  for a given  $x$  and  $h$  have a unique solution?

- Yes, because the first term is strictly convex
- So the sum will be strictly convex
- This guarantees a single solution

### C. Proximal gradient descent algorithm

1. initialize  $x^{(0)}$

2. repeat  $x^{(k)} = \operatorname{prox}_{t_k} \left( x^{(k-1)} - t_k \nabla g(x^{(k-1)}) \right)$ ,  $k = 1, 2, 3, \dots$

To make this update step look familiar, rewrite it as

$$x^{(k)} = x^{(k-1)} - t_k \cdot G_{t_k}(x^{(k-1)})$$

where  $G_t(x) = \frac{x - \operatorname{prox}_t(x - t \nabla g(x))}{t}$ , known as the generalized gradient of  $f$

So,

$$x^{(k)} = x^{(k-1)} - t_k \cdot G_{t_k}(x^{(k-1)}) = x^{(k-1)} - t_k \left( \frac{x^{(k-1)} - \operatorname{prox}_{t_k}(x^{(k-1)} - t_k \nabla g(x^{(k-1)}))}{t_k} \right)$$

- everything cancels out nicely. By design!

## D. Why did this help?

Didn't we just swap one minimization problem for another? We still need to figure out how to minimize  $h$ ?

Yes. But  $\text{prox}()$  can be computed analytically for a lot of important functions  $h$ .

Note that

1. mapping  $\text{prox}()$  doesn't depend on  $g$  at all, only on  $h$
2. smooth part  $g$  can be complicated, but it doesn't matter! We only need to calculate  $\nabla g(x)$

So when applied to a known problem, using the prox is useful

Convergence analysis will be in terms of number of iterations of the algorithm. However do remember that each iteration will evaluate  $\text{prox}()$  once, and this can be cheap or expensive depending on  $h$ .

## Example: ISTA

Another algorithm for the lasso.

Last time we saw subgradient method - didn't do very well.

We can try using proximal gradient method here, for which we need to calculate

$$G_t(x) = \frac{x - \text{prox}_t(x - t\nabla g(x))}{t},$$

That is, we need to know

1. what  $\text{prox}()$  is
2. what  $\nabla g(x)$  is

Example: ISTA

Given  $y \in \mathbb{R}^n$ ,  $X \in \mathbb{R}^{n \times p}$ , recall **lasso** criterion:

$$f(\beta) = \underbrace{\frac{1}{2} \|y - X\beta\|_2^2}_{g(\beta)} + \underbrace{\lambda \|\beta\|_1}_{h(\beta)}$$

Prox mapping is now

$$\begin{aligned} \text{prox}_t(\beta) &= \underset{z}{\operatorname{argmin}} \frac{1}{2t} \|\beta - z\|_2^2 + \lambda \|z\|_1 \\ &= S_{\lambda t}(\beta) \end{aligned}$$

where  $S_\lambda(\beta)$  is the soft-thresholding operator,

$$[S_\lambda(\beta)]_i = \begin{cases} \beta_i - \lambda & \text{if } \beta_i > \lambda \\ 0 & \text{if } -\lambda \leq \beta_i \leq \lambda \\ \beta_i + \lambda & \text{if } \beta_i < -\lambda \end{cases}, \quad i = 1, \dots, n$$

8

1. You get  $\text{prox}(B)$  by setting the subgradient of the function  $\frac{1}{2t} \|B - z\|_2^2 + \lambda \|z\|_1$  to zero to get the argmin

But we already solved this exact problem! See Lecture 6, soft thresholding

So  $\text{prox}_t(B) = \underset{z}{\operatorname{argmin}} \frac{1}{2t} \|B - z\|_2^2 + \lambda \|z\|_1 = S_{\lambda t}(B)$

2. We know already that  $\nabla g(B) = -X^T(y - XB)$  from all the least squares work we have seen so far...

So, putting this together, the update rule is

$$B^k = \text{prox}_{t_k} \left( B^{(k-1)} - t_k \nabla g(B^{(k-1)}) \right) = S_{\lambda t_k} \left( B^{(k-1)} + t_k X^T(y - XB) \right)$$

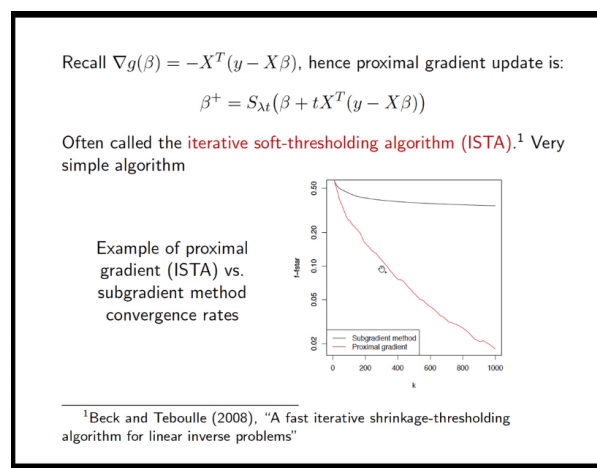
This is called the iterative soft thresholding algorithm = ista

It soft threshold the usual gradient update for the unconstrained least squares problem!

- a very simple and intuitive algorithm

This shows directly from the algorithm that we get a sparse solution for the lasso

- sparsity inducing - it sets some of the B values to zero by soft thresholding!



## E. Backtracking line search to get t

For proximal gradient this works imiliart to gradient descent, but operates on g and not f

Choose parameter  $0 < B < 1$ .

At each iteration:

- start at  $t = t_{init}$
- while  $g(x - tG_t(x)) > g(x) - t\nabla g(x)^T G_t(x) + \frac{t}{2} \|G_t(x)\|_2^2$ , shrink  $t = Bt$ 
  - if  $G_t(x)$  were the same as  $\nabla g(x)$ , this would yield the same formulation as the backtracking check in Lecture 5 on gradient descent, assuming  $\alpha = 1/2$
- else perform  $x^{(k)} = x^{(k-1)} - t_k \cdot G_{t_k}(x^{(k-1)})$

Note that every evaluation of  $G()$  also involves a call to  $\text{prox}()$  - so this might be expensive if  $\text{prox}$  is expensive to evaluate! There are alternative formulations that require less computation i.e. fewer calls to  $\text{prox}()$

## F. Convergence analysis

Same as gradient descent!

Except in every step, it does something more expensive than for gradient descent (i.e. the  $\text{prox}$  call)

- this is a good result if the  $\text{prox}$  cost is low

### Convergence analysis

For criterion  $f(x) = g(x) + h(x)$ , we assume:

- $g$  is convex, differentiable,  $\text{dom}(g) = \mathbb{R}^n$ , and  $\nabla g$  is Lipschitz continuous with constant  $L > 0$
- $h$  is convex,  $\text{prox}_t(x) = \arg\min_z \{\|x - z\|_2^2 / (2t) + h(z)\}$  can be evaluated

**Theorem:** Proximal gradient descent with fixed step size  $t \leq 1/L$  satisfies

$$f(x^{(k)}) - f^* \leq \frac{\|x^{(0)} - x^*\|_2^2}{2tk}$$

and same result holds for backtracking, with  $t$  replaced by  $\beta/L$

Proximal gradient descent has convergence rate  $O(1/k)$  or  $O(1/\epsilon)$ .  
Same as gradient descent! (But remember,  $\text{prox}$  cost matters ...)

In the ISTA example, the  $\text{prox}$  cost is low. Soft thresholding is  $O(p)$  operation since it just needs one comparison per  $B$  value.

Here is an example with high  $\text{prox}$  cost.

### Example: matrix completion

Netflix prize problem - recommender systems

- assume some low rank latent structure: if  $Y$  were fully observed, it would be low rank

Given matrix  $Y \in \mathbb{R}^{m \times n}$ , only observe entries  $Y_{ij}$ ,  $i, j \in \Omega$

We want to fill in the missing entries, i.e. we want to solve a matrix completion problem

$$\min_B \frac{1}{2} \sum_{i,j \in \Omega} (Y_{ij} - B_{ij})^2 + \lambda \|B\|_{tr}$$

- the trace norm is the sum of a matrix's singular values:  $\|B\|_{tr} = \sum_{i=1}^r \sigma_i(B)$  where  $r = \text{rank}(B)$  and the  $\sigma$ s are the singular values
- it is a convex approximation to the matrix's rank

- If B was a diagonal matrix:
  - rank = L0 norm
  - trace = L1 norm
- So by penalizing the trace norm in this problem, we are encouraging a low rank in the solution
  - approximate Y, but least complicated

We will try solving this using proximal gradient descent.

As a reminder

- we find the generalized gradient G for the criterion function f(x), using the gradient of the differentiable part g(x)

$$-G_t(x) = \frac{x - \text{prox}(x - t\nabla g(x))}{t}$$

- Then do  $x^+ = x - tG_t(x)$

Define  $P_\Omega$ , projection operator onto observed set:

$$[P_\Omega(B)]_{ij} = \begin{cases} B_{ij} & (i, j) \in \Omega \\ 0 & (i, j) \notin \Omega \end{cases}$$

Then the criterion is

$$f(B) = \underbrace{\frac{1}{2} \|P_\Omega(Y) - P_\Omega(B)\|_F^2}_{g(B)} + \underbrace{\lambda \|B\|_{\text{tr}}}_{h(B)}$$

Two ingredients needed for proximal gradient descent:

- Gradient calculation:  $\nabla g(B) = -(P_\Omega(Y) - P_\Omega(B))$
- Prox function:

$$\text{prox}_t(B) = \underset{Z}{\operatorname{argmin}} \frac{1}{2t} \|B - Z\|_F^2 + \lambda \|Z\|_{\text{tr}}$$

13

Q: how can you calculate the prox? how to solve this minimization problem?

Claim:  $\text{prox}_t(B) = S_{\lambda t}(B)$ , **matrix soft-thresholding** at the level  $\lambda$ .  
Here  $S_\lambda(B)$  is defined by

$$S_\lambda(B) = U \Sigma_\lambda V^T$$

where  $B = U \Sigma V^T$  is an SVD, and  $\Sigma_\lambda$  is diagonal with

$$(\Sigma_\lambda)_{ii} = \max\{\Sigma_{ii} - \lambda, 0\}$$

Proof: note that  $\text{prox}_t(B) = Z$ , where  $Z$  satisfies

$$0 \in Z - B + \lambda t \cdot \partial \|Z\|_{\text{tr}}$$

Helpful fact: if  $Z = U \Sigma V^T$ , then

$$\partial \|Z\|_{\text{tr}} = \{UV^T + W : \|W\|_{\text{op}} \leq 1, U^T W = 0, W V = 0\}$$

Now plug in  $Z = S_{\lambda t}(B)$  and check that we can get 0

..

The matrix soft thresholding adjusts (downwards) the singular values of the matrix

This is a challenging prox to derive! We need to minimize a function with a trace norm.

Right now, the only way we know how to minimize a non smooth function is subgradient optimality.

$$\text{prox}_Z(B) = \underset{Z}{\operatorname{argmin}} \frac{1}{2t} \|B - Z\|_2^2 + \lambda \|Z\|_{tr} = \underset{Z}{\operatorname{argmin}} \frac{1}{2} \|B - Z\|_2^2 + \lambda t \|Z\|_{tr}$$

So at optimal  $z^*$ ,  $0 \in \partial \text{prox}_Z(B) = Z^* - B + t\lambda \partial(\|Z^*\|_{tr}) \implies Z^* = B - t\lambda \Gamma$

Let  $\Gamma \in \partial(\|Z\|_{tr})$ . We need to find some  $\Gamma$  in this subdifferential, to help us calculate the value  $Z^*$ .

As a reminder, we know from the structure of the prox function that there is just one unique solution, as the first term is strictly convex.

Fact:  $\partial \|X\|_{tr} = \left\{ \sum U V^T + W : \|W\|_{op} \leq 1, U^T W = 0, W V = 0 \right\}$   
 $X = U \Sigma V^T$

subgradient optimality: with this choice.  
 $\Gamma = U \left( \frac{\Sigma - \Sigma_{\lambda t}}{\lambda t} \right) V^T$  is indeed a subgradient

The claim is that  $\text{prox}(B) = Z^* = B - t\lambda \Gamma = S_{\lambda t}(B)$

## Conclusion:

The proximal gradient step is

$$B^+ = \text{prox}_t(B - t\nabla g(B)) = S_{\lambda t}(B - t\nabla g(B)) = S_{\lambda t}(B + t(P_{\Omega}(Y) - P_{\Omega}(B)))$$

- Note that this involves computes the SVD of  $P_{\Omega}(Y) - P_{\Omega}(B)$  and shrinking the singular values
- but  $Y$  is sparse (hence the problem) and  $B$  is designed to be low rank. This hopefully keeps the cost of performing SVD low - therefore making this an efficient algorithm
  - there are good numerical routines to do this

Similar to lasso - soft threshold the gradient step!

Hence proximal gradient update step is:

$$B^+ = S_{\lambda t} \left( B + t(P_{\Omega}(Y) - P_{\Omega}(B)) \right)$$

Note that  $\nabla g(B)$  is Lipschitz continuous with  $L = 1$ , so we can choose fixed step size  $t = 1$ . Update step is now:

$$B^+ = S_{\lambda} (P_{\Omega}(Y) + P_{\Omega}^{\perp}(B))$$

where  $P_{\Omega}^{\perp}$  projects onto unobserved set,  $P_{\Omega}(B) + P_{\Omega}^{\perp}(B) = B$

This is the **soft-impute** algorithm<sup>2</sup>, simple and effective method for matrix completion

---

<sup>2</sup>Mazumder et al. (2011), "Spectral regularization algorithms for learning large incomplete matrices"

Update rule interpretation:

At every step, take observed entries of  $Y$ , fill in all the unobserved entries with what you thought they should have been in the last iteration. Then take this matrix and reduce its rank.

## G. Special cases

Proximal gradient descent is also called composite gradient descent or generalized gradient descent

- it is a superset of the other gradient descent methods we have seen - just adjust the notation/reshape the problem appropriately

### Special cases

Proximal gradient descent also called composite gradient descent, or **generalized gradient descent**

Why "generalized"? This refers to the several special cases, when minimizing  $f = g + h$ :

- $h = 0$  — gradient descent
- $h = I_C$  — projected gradient descent
- $g = 0$  — proximal minimization algorithm

Therefore these algorithms all have  $O(1/\epsilon)$  convergence rate

## H. Projected gradient descent

The prox if  $h =$  indicator just projects onto the set  $C$

So the proximal gradient update step = project the gradient update step onto  $C$ !



## Projected gradient descent

Given closed, convex set  $C \in \mathbb{R}^n$ ,

$$\min_{x \in C} g(x) \iff \min_x g(x) + I_C(x)$$

where  $I_C(x) = \begin{cases} 0 & x \in C \\ \infty & x \notin C \end{cases}$  is the indicator function of  $C$

Hence

$$\begin{aligned} \text{prox}_t(x) &= \underset{z}{\operatorname{argmin}} \frac{1}{2t} \|x - z\|_2^2 + I_C(z) \\ &= \underset{z \in C}{\operatorname{argmin}} \|x - z\|_2^2 \end{aligned}$$

i.e.,  $\text{prox}_t(x) = P_C(x)$ , projection operator onto  $C$

17

## I. What if you cannot evaluate prox?

### What happens if we can't evaluate prox?

Theory for proximal gradient, with  $f = g + h$ , assumes that prox function can be evaluated, i.e., assumes the minimization

$$\text{prox}_t(x) = \underset{z}{\operatorname{argmin}} \frac{1}{2t} \|x - z\|_2^2 + h(z)$$

can be done exactly. In general, not clear what happens if we just minimize this approximately

But, if you can precisely control the errors in approximating the prox operator, then you can recover the original convergence rates<sup>3</sup>

In practice, if prox evaluation is done approximately, then it should be done to decently high accuracy

<sup>3</sup>Schmidt et al. (2011), "Convergence rates of inexact proximal-gradient methods for convex optimization"

20

- as the algorithm progresses, if (estimated prox - real prox) decreases, we get pretty close to the true solution

## J. Acceleration

You can accelerate proximal gradient descent in order to achieve the optimal  $O(1/\sqrt{\epsilon})$  rate.

## Acceleration

Turns out we can **accelerate** proximal gradient descent in order to achieve the optimal  $O(1/\sqrt{\epsilon})$  convergence rate. Four ideas (three acceleration methods) by Nesterov:

- 1983: original acceleration idea for smooth functions
- 1988: another acceleration idea for smooth functions
- 2005: smoothing techniques for nonsmooth functions, coupled with original acceleration idea
- 2007: acceleration idea for composite functions<sup>4</sup>

We will follow Beck and Teboulle (2008), an extension of Nesterov (1983) to composite functions<sup>5</sup>

<sup>4</sup>Each step uses entire history of previous steps and makes two prox calls

<sup>5</sup>Each step uses information from two last steps and makes one prox call

Imagine the same setting as before

$$\min g(x) + h(x)$$

### Accelerated proximal gradient method

1. choose initial point  $x^{(0)} = x^{(-1)} \in \mathbb{R}^n$  (just for notation consistency)

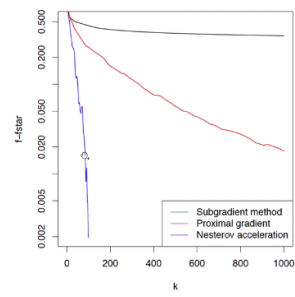
2. repeat

$$v = x^{(k-1)} + \frac{k-2}{k+1} (x^{(k-1)} - x^{(k-2)})$$

$$x^{(k)} = \text{prox}_{t_k}(v - t_k \nabla g(v))$$

- usually, we just perform  $x^{(k)} = \text{prox}_{t_k}(x^{(k-1)} - t_k \nabla g(x^{(k-1)}))$
- here, we find a new point  $v$ , which carries some past information from previous iterations
- if  $h = 0$  this is the accelerated gradient method
- 'if you were travelling in a direction, this pushes you to continue traveling in that direction'
- the momentum weight increases as the number of iterations increases, eventually approaching 1
  - the intuition is that as you get closer to the optimum, the gradient gets smaller, so your gradient update is smaller too
  - but you don't want to slow down! want to keep continuing on

Back to lasso example: acceleration can really help!



Note: accelerated proximal gradient is not a descent method