

# ECE 18-898G: Special Topics in Signal Processing: Sparsity, Structure, and Inference

Sparse Recovery using L1 minimization - algorithms

Yuejie Chi

Department of Electrical and Computer Engineering

**Carnegie Mellon University**

Spring 2018

# Outline

---

- Lasso with orthogonal design
- Proximal operators
- Proximal gradient methods for lasso and its extensions
- Nesterov's accelerated algorithm (FISTA)

These are useful in general for compound optimization problems.

# Sparse recovery by $\ell_0$ regularization

---

As a warm up, consider a sparsifying basis  $\mathbf{X} \in \mathbb{R}^{n \times n}$  that is orthonormal, we wish to solve the following sparsity-promoting problem regularized by  $\ell_0$  norm:

$$\hat{\beta} = \operatorname{argmin}_{\beta} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \underbrace{\lambda \|\beta\|_0}_{\text{penalized by sparsity level}} \quad (4.1)$$

- The first term is the approximation error:  $\|\mathbf{y} - \mathbf{X}\hat{\beta}\|^2$ .
- The second term is the model complexity:  $\|\hat{\beta}\|_0$ .

We will discuss “regularized” algorithms throughout this lecture.

# Orthogonal design

---

Since  $\mathbf{X}$  is orthonormal,

$$\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 = \|\mathbf{X}^\top \mathbf{y} - \boldsymbol{\beta}\|^2.$$

Without loss of generality, suppose  $\mathbf{X} = \mathbf{I}$ , then (4.1) reduces to

$$\hat{\boldsymbol{\beta}} = \operatorname{argmin}_{\boldsymbol{\beta}} \sum_{i=1}^n \frac{1}{2} \left[ (y_i - \beta_i)^2 + \lambda \cdot \mathbf{1}\{\beta_i \neq 0\} \right]$$

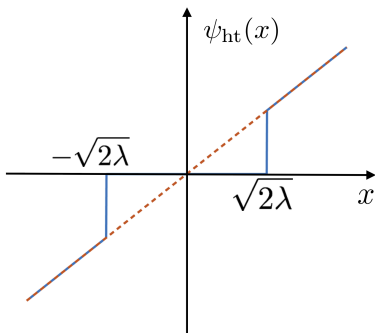
Solving this problem gives

$$\hat{\beta}_i = \begin{cases} 0, & |y_i| \leq \sqrt{2\lambda} \\ y_i, & |y_i| > \sqrt{2\lambda} \end{cases} \quad \text{hard thresholding}$$

- Keep large coefficients; discard small coefficients

## The case $X = I$

---



$$\hat{\beta}_i = \psi_{\text{ht}}(y_i; \sqrt{2\lambda}) := \begin{cases} 0, & |y_i| \leq \sqrt{2\lambda} \\ y_i, & |y_i| > \sqrt{2\lambda} \end{cases} \quad \text{hard thresholding}$$

Hard thresholding preserves data outside threshold zone

# Convex relaxation: Lasso (Tibshirani '96)

---

**Lasso** (Least absolute shrinkage and selection operator)

$$\hat{\beta} = \operatorname{argmin}_{\beta} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \|\beta\|_1 \quad (4.2)$$

for some regularization parameter  $\lambda > 0$ .

- It is equivalent to

$$\hat{\beta} = \operatorname{argmin}_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|^2 \quad \text{s.t.} \quad \|\beta\|_1 \leq t$$

for some  $t$  that depends on  $\lambda$  (no explicit formula)

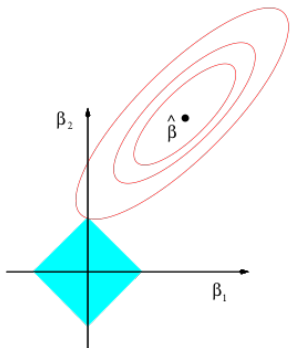
- a quadratic program (QP) with convex constraints

- $\lambda$  controls model complexity: larger  $\lambda$  restricts the parameters more; smaller  $\lambda$  frees up more parameters
- Also related to Basis Pursuit:

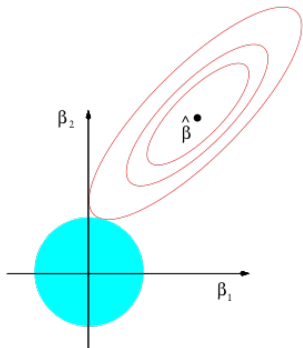
$$\hat{\beta} = \operatorname{argmin}_{\beta} \|\beta\|_1 \quad \text{s.t.} \quad \|\mathbf{y} - \mathbf{X}\beta\| \leq \epsilon$$

# Lasso vs. ridge regression

$\hat{\beta}$ : least squares solution



$$\begin{aligned} &\text{minimize}_{\beta} && \|\mathbf{y} - \mathbf{X}\beta\| \\ &\text{s.t.} && \|\beta\|_1 \leq t \end{aligned}$$



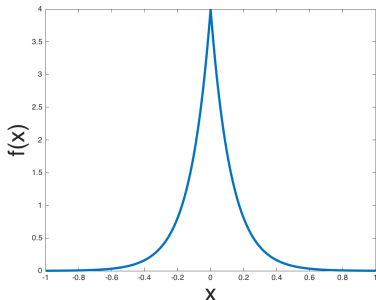
$$\begin{aligned} &\text{minimize}_{\beta} && \|\mathbf{y} - \mathbf{X}\beta\| \\ &\text{s.t.} && \|\beta\|_2 \leq t \end{aligned}$$

Fig. credit: Hastie, Tibshirani, & Wainwright

# A Bayesian interpretation

---

Orthogonal design:  $\mathbf{y} = \boldsymbol{\beta} + \boldsymbol{\eta}$  with  $\boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ .



Impose an i.i.d. prior on  $\beta_i$  to encourage sparsity (*Gaussian is not a good choice*):

$$\text{(Laplacian prior)} \quad \mathbb{P}(\beta_i = z) = \frac{\lambda}{2} e^{-\lambda|z|}$$

The tail decays exponentially.



# A Bayesian interpretation of Lasso

---

Posterior of  $\beta$ :

$$\begin{aligned}\mathbb{P}(\beta | \mathbf{y}) &\propto \mathbb{P}(\mathbf{y}|\beta)\mathbb{P}(\beta) \propto \prod_{i=1}^n e^{-\frac{(y_i - \beta_i)^2}{2\sigma^2}} \frac{\lambda}{2} e^{-\lambda|\beta_i|} \\ &\propto \prod_{i=1}^n \exp\left\{-\frac{(y_i - \beta_i)^2}{2\sigma^2} - \lambda|\beta_i|\right\}\end{aligned}$$

$\implies$  maximum *a posteriori* (MAP) estimator:

$$\arg \min_{\beta} \sum_{i=1}^n \left\{ \frac{(y_i - \beta_i)^2}{2\sigma^2} + \lambda|\beta_i| \right\} \quad (\text{Lasso})$$

**Implication:** Lasso is MAP estimator under Laplacian prior

## Example: orthogonal design

---

Suppose  $\mathbf{X} = \mathbf{I}$ , then Lasso reduces to

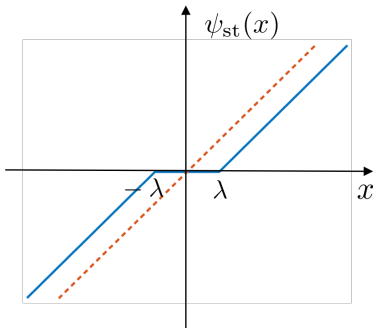
$$\hat{\boldsymbol{\beta}} = \operatorname{argmin}_{\boldsymbol{\beta}} \sum_{i=1}^n \left[ \frac{1}{2} (y_i - \beta_i)^2 + \lambda |\beta_i| \right]$$

The Lasso estimate  $\hat{\boldsymbol{\beta}}$  is then given by

$$\hat{\beta}_i = \begin{cases} y_i - \lambda, & y_i \geq \lambda \\ y_i + \lambda, & y_i \leq -\lambda \\ 0, & \text{else} \end{cases} \quad \text{soft thresholding}$$

## Example: orthogonal design

---



$$\hat{\beta}_i = \psi_{\text{st}}(y_i; \lambda) = \begin{cases} y_i - \lambda, & y_i \geq \lambda \\ y_i + \lambda, & y_i \leq -\lambda \\ 0, & \text{else} \end{cases} \quad \text{soft thresholding}$$

Soft thresholding shrinks data towards 0 outside threshold zone

# Optimality condition for convex functions

For any convex function  $f(\beta)$ ,  $\beta^*$  is an optimal solution iff  $0 \in \partial f(\beta^*)$ , where  $\partial f(\beta)$  is the set of all subgradients at  $\beta$

- The subgradient of  $f(\beta) = \frac{1}{2}(y - \beta)^2 + \lambda|\beta|$  can be written as

$$g = \beta - y + \lambda s$$

with  $s$  is a subgradient of  $f(\beta) = |\beta|$  if

$$\begin{cases} s = \text{sign}(\beta), & \text{if } \beta \neq 0 \\ s \in [-1, 1], & \text{if } \beta = 0 \end{cases} \quad (4.3)$$

- We see that  $\hat{\beta} = \psi_{\text{st}}(y; \lambda)$  by checking optimality conditions for two cases:
  - If  $|y| \leq \lambda$ , taking  $\beta = 0$  and  $s = y/\lambda$  gives  $g = 0$
  - If  $|y| > \lambda$ , taking  $\beta = y - \text{sign}(y)\lambda$  gives  $g = 0$

## **Solving LASSO in general cases**

# Composite optimization problems

---

General composite optimization problem:

$$\hat{\beta} = \operatorname{argmin}_{\beta} \{F(\beta) = f(\beta) + g(\beta)\}$$

- $f(\beta)$  is convex and differentiable, e.g. approximation error
- $g(\beta)$  is convex, possibly non-differentiable, e.g. regularizers

**Examples:**

- LASSO:  $f(\beta) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2$ , and  $g(\beta) = \lambda \|\beta\|_1$ .
- Matrix completion:

$$f(\mathbf{X}) = \|\mathcal{P}_{\Omega}(\mathbf{Y} - \mathbf{X})\|_F^2, \quad g(\mathbf{X}) = \lambda \|\mathbf{X}\|_*$$

where  $\|\mathbf{X}\|_*$  is the nuclear norm.

# Motivation

---

Standard methods (e.g. subgradient methods) for solving composite optimization has very slow convergence rate.

We would discuss accelerated proximal gradient methods that

- is iterative, and has low computational cost (first-order algorithm, which requires computation of a single gradient per iteration);
- has quadratic convergence rate;
- performs well in practice and works for a large class of problems.

# Proximal gradient methods

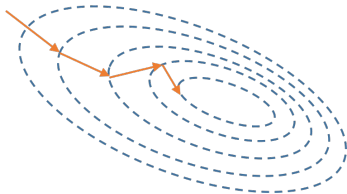


# Gradient descent

---

$$\text{minimize}_{\beta \in \mathbb{R}^p} f(\beta)$$

where  $f(\beta)$  is convex and smooth (differentiable)



---

## Algorithm 4.1 Gradient descent

---

**for**  $t = 0, 1, \dots$ :

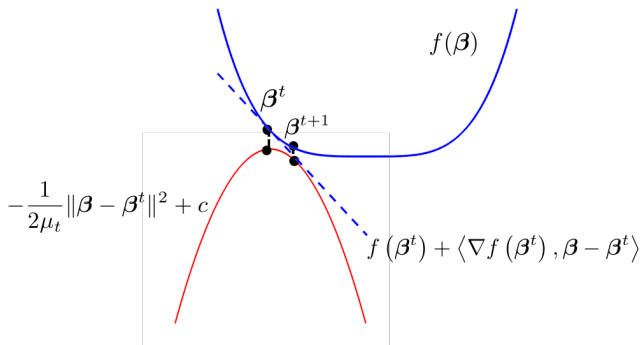
$$\beta^{t+1} = \beta^t - \mu_t \nabla f(\beta^t)$$

where  $\mu_t$ : step size / learning rate

---

# A proximal point of view of GD

---



$$\beta^{t+1} = \arg \min_{\beta} \left\{ \underbrace{f(\beta^t) + \langle \nabla f(\beta^t), \beta - \beta^t \rangle}_{\text{linear approximation}} + \underbrace{\frac{1}{2\mu_t} \|\beta - \beta^t\|^2}_{\text{proximal term}} \right\}$$

- When  $\mu_t$  is small,  $\beta^{t+1}$  tends to stay close to  $\beta^t$

# Proximal operator

---

If we define the proximal operator

$$\text{prox}_h(\mathbf{b}) := \arg \min_{\boldsymbol{\beta}} \left\{ \frac{1}{2} \|\boldsymbol{\beta} - \mathbf{b}\|^2 + h(\boldsymbol{\beta}) \right\}$$

for any convex function  $h$ , then one can write

$$\boldsymbol{\beta}^{t+1} = \text{prox}_{\mu_t f_t}(\boldsymbol{\beta}^t)$$

where  $f_t(\boldsymbol{\beta}) := f(\boldsymbol{\beta}_t) + \langle \nabla f(\boldsymbol{\beta}_t), \boldsymbol{\beta} - \boldsymbol{\beta}_t \rangle$

Gradient descent is performing proximal mapping at every iteration.

# Why consider proximal operators?

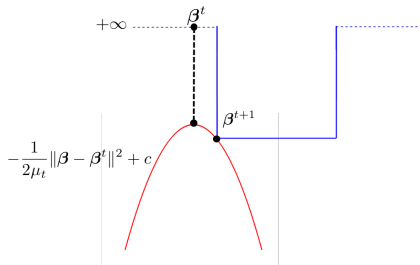
---

$$\text{prox}_h(\mathbf{b}) := \arg \min_{\boldsymbol{\beta}} \left\{ \frac{1}{2} \|\boldsymbol{\beta} - \mathbf{b}\|^2 + h(\boldsymbol{\beta}) \right\}$$

- It is well-defined under very general conditions (including nonsmooth convex functions)
- The operator can be evaluated efficiently for many widely used functions (in particular, regularizers)
- This abstraction is conceptually and mathematically simple, and covers many well-known optimization algorithms

# Example: characteristic functions

---



- If  $h$  is characteristic function

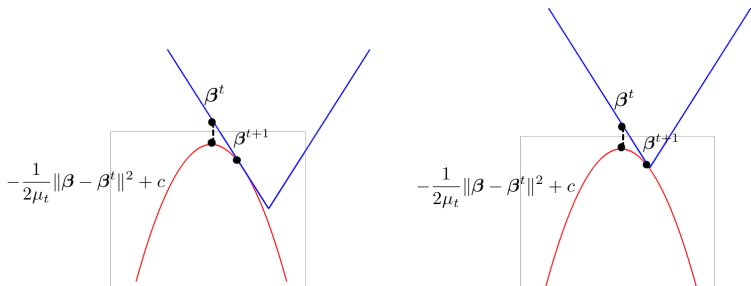
$$h(\beta) = \begin{cases} 0, & \text{if } \beta \in \mathcal{C} \\ \infty, & \text{else} \end{cases}$$

then

$$\text{prox}_h(\mathbf{b}) = \arg \min_{\beta \in \mathcal{C}} \|\beta - \mathbf{b}\|_2 \quad (\text{Euclidean projection})$$

## Example: $\ell_1$ norm

---



- If  $h(\beta) = \|\beta\|_1$ , then

$$\text{prox}_{\lambda h}(\mathbf{b}) = \psi_{\text{st}}(\mathbf{b}; \lambda)$$

where soft-thresholding  $\psi_{\text{st}}(\cdot)$  is applied in an entry-wise manner.

## Example: $\ell_2$ norm

---

$$\text{prox}_h(\mathbf{b}) := \arg \min_{\boldsymbol{\beta}} \left\{ \frac{1}{2} \|\boldsymbol{\beta} - \mathbf{b}\|^2 + h(\boldsymbol{\beta}) \right\}$$

- If  $h(\boldsymbol{\beta}) = \|\boldsymbol{\beta}\|$ , then

$$\text{prox}_{\lambda h}(\mathbf{b}) = \left( 1 - \frac{\lambda}{\|\mathbf{b}\|} \right)_+ \mathbf{b}$$

where  $a_+ := \max\{a, 0\}$ . This is called *block soft thresholding*.

## Example: log barrier

---

$$\text{prox}_h(\mathbf{b}) := \arg \min_{\boldsymbol{\beta}} \left\{ \frac{1}{2} \|\boldsymbol{\beta} - \mathbf{b}\|^2 + h(\boldsymbol{\beta}) \right\}$$

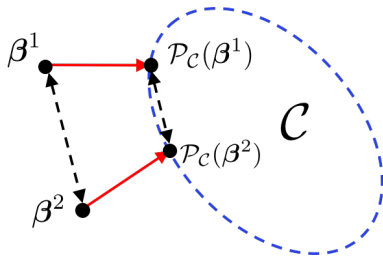
- If  $h(\boldsymbol{\beta}) = -\sum_{i=1}^p \log \beta_i$ , then

$$(\text{prox}_{\lambda h}(\mathbf{b}))_i = \frac{b_i + \sqrt{b_i^2 + 4\lambda}}{2}$$



# Nonexpansiveness of proximal operators

---

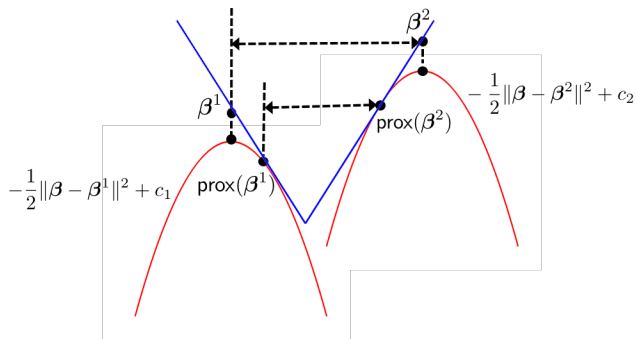


Recall that when  $h(\beta) = \begin{cases} 0, & \text{if } \beta \in \mathcal{C} \\ \infty & \text{else} \end{cases}$ ,  $\text{prox}_h(\beta)$  is Euclidean projection  $\mathcal{P}_\mathcal{C}$  onto  $\mathcal{C}$ , which is nonexpansive:

$$\|\mathcal{P}_\mathcal{C}(\beta^1) - \mathcal{P}_\mathcal{C}(\beta^2)\| \leq \|\beta^1 - \beta^2\|$$

# Nonexpansiveness of proximal operators

Nonexpansiveness is a property for general  $\text{prox}_h(\cdot)$



## Fact 4.1 (Nonexpansiveness)

$$\|\text{prox}_h(\beta^1) - \text{prox}_h(\beta^2)\| \leq \|\beta^1 - \beta^2\|$$

- In some sense, proximal operator behaves like projection

## Proof of nonexpansiveness

---

Let  $z^1 = \text{prox}_h(\beta^1)$  and  $z^2 = \text{prox}_h(\beta^2)$ . Subgradient characterizations of  $z^1$  and  $z^2$  read

$$\beta^1 - z^1 \in \partial h(z^1) \quad \text{and} \quad \beta^2 - z^2 \in \partial h(z^2)$$

The claim would follow if

$$(\beta^1 - \beta^2)^\top (z^1 - z^2) \geq \|z^1 - z^2\|^2 \quad (\text{together with Cauchy-Schwarz})$$

$$\begin{aligned} &\iff (\beta^1 - z^1 - \beta^2 + z^2)^\top (z^1 - z^2) \geq 0 \\ &\iff \begin{cases} h(z^2) \geq h(z^1) + \underbrace{\langle \beta^1 - z^1, z^2 - z^1 \rangle}_{\in \partial h(z^1)} \\ h(z^1) \geq h(z^2) + \underbrace{\langle \beta^2 - z^2, z^1 - z^2 \rangle}_{\in \partial h(z^2)} \end{cases} \end{aligned}$$

# Proximal gradient methods

# Optimizing composite functions

---

$$\text{(Lasso)} \quad \underset{\beta \in \mathbb{R}^p}{\text{minimize}} \quad \underbrace{\frac{1}{2} \|\mathbf{X}\beta - \mathbf{y}\|^2}_{:=f(\beta)} + \underbrace{\lambda \|\beta\|_1}_{:=g(\beta)} = f(\beta) + g(\beta)$$

where  $f(\beta)$  is differentiable, and  $g(\beta)$  is non-smooth

- Since  $g(\beta)$  is non-differentiable, we cannot run vanilla gradient descent

# Proximal gradient methods

---

One strategy: replace  $f(\beta)$  with linear approximation, and compute the proximal solution

$$\beta^{t+1} = \arg \min_{\beta} \left\{ f(\beta^t) + \langle \nabla f(\beta^t), \beta - \beta^t \rangle + g(\beta) + \frac{1}{2\mu_t} \|\beta - \beta^t\|^2 \right\}$$

The optimality condition reads

$$\mathbf{0} \in \nabla f(\beta^t) + \partial g(\beta^{t+1}) + \frac{1}{\mu_t} (\beta^{t+1} - \beta^t)$$

which is equivalent to optimality condition of

$$\begin{aligned} \beta^{t+1} &= \arg \min_{\beta} \left\{ g(\beta) + \frac{1}{2\mu_t} \left\| \beta - (\beta^t - \mu_t \nabla f(\beta^t)) \right\|^2 \right\} \\ &= \text{prox}_{\mu_t g} (\beta^t - \mu_t \nabla f(\beta^t)) \end{aligned}$$

# Proximal gradient methods

---

Alternate between gradient updates on  $f$  and proximal mapping on  $g$

---

## Algorithm 4.2 Proximal gradient methods

---

**for**  $t = 0, 1, \dots$ :

$$\beta^{t+1} = \text{prox}_{\mu_t g} \left( \beta^t - \mu_t \nabla f(\beta^t) \right)$$

where  $\mu_t$ : step size / learning rate

---

# Projected gradient methods

---

When  $g(\beta) = \begin{cases} 0, & \text{if } \beta \in \underbrace{\mathcal{C}}_{\text{convex}} \\ \infty, & \text{else} \end{cases}$  is characteristic function:

$$\begin{aligned}\beta^{t+1} &= \mathcal{P}_{\mathcal{C}} \left( \beta^t - \mu_t \nabla f(\beta^t) \right) \\ &:= \arg \min_{\beta \in \mathcal{C}} \left\| \beta - (\beta^t - \mu_t \nabla f(\beta^t)) \right\|\end{aligned}$$

This is a first-order method to solve the constrained optimization

$$\begin{aligned}\text{minimize}_{\beta} \quad & f(\beta) \\ \text{s.t.} \quad & \beta \in \mathcal{C}\end{aligned}$$



# Proximal gradient methods for lasso

---

For lasso:  $f(\boldsymbol{\beta}) = \frac{1}{2}\|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|^2$  and  $g(\boldsymbol{\beta}) = \lambda\|\boldsymbol{\beta}\|_1$ ,

$$\begin{aligned}\text{prox}_g(\boldsymbol{\beta}) &= \arg \min_{\mathbf{b}} \left\{ \frac{1}{2}\|\boldsymbol{\beta} - \mathbf{b}\|^2 + \lambda\|\mathbf{b}\|_1 \right\} \\ &= \psi_{\text{st}}(\boldsymbol{\beta}; \lambda)\end{aligned}$$

$$\implies \boldsymbol{\beta}^{t+1} = \psi_{\text{st}}\left(\boldsymbol{\beta}^t - \mu_t \mathbf{X}^\top (\mathbf{X}\boldsymbol{\beta}^t - \mathbf{y}); \mu_t \lambda\right)$$

iterative soft thresholding

# Proximal gradient methods for group lasso

---

Sometimes variables have a natural group structure, and it is desirable to set all variables within a group to be zero (or nonzero) simultaneously

$$\text{(group lasso)} \quad \underbrace{\frac{1}{2} \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|^2}_{:=f(\boldsymbol{\beta})} + \lambda \underbrace{\sum_{j=1}^k \|\boldsymbol{\beta}_j\|}_{:=g(\boldsymbol{\beta})}$$

where  $\boldsymbol{\beta}_j \in \mathbb{R}^{p/k}$  and  $\boldsymbol{\beta} = \begin{bmatrix} \boldsymbol{\beta}_1 \\ \vdots \\ \boldsymbol{\beta}_k \end{bmatrix}$ .

$$\text{prox}_g(\boldsymbol{\beta}) = \psi_{\text{bst}}(\boldsymbol{\beta}; \lambda) := \left[ \left(1 - \frac{\lambda}{\|\boldsymbol{\beta}_j\|}\right)_+ \boldsymbol{\beta}_j \right]_{1 \leq j \leq k}$$

$$\implies \boldsymbol{\beta}^{t+1} = \psi_{\text{bst}}(\boldsymbol{\beta}^t - \mu_t \mathbf{X}^\top (\mathbf{X}\boldsymbol{\beta}^t - \mathbf{y}); \mu_t \lambda)$$

# Proximal gradient methods for elastic net

---

Lasso does not handle highly correlated variables well: if there is a group of highly correlated variables, lasso often picks one from the group and ignore the rest.

- Sometimes we make a compromise between lasso and  $\ell_2$  penalties

$$\text{(elastic net)} \quad \underbrace{\frac{1}{2} \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|^2}_{:=f(\boldsymbol{\beta})} + \lambda \underbrace{\left\{ \|\boldsymbol{\beta}\|_1 + (\gamma/2) \|\boldsymbol{\beta}\|_2^2 \right\}}_{:=g(\boldsymbol{\beta})}$$

$$\text{prox}_{\lambda g}(\boldsymbol{\beta}) = \frac{1}{1 + \lambda\gamma} \psi_{\text{st}}(\boldsymbol{\beta}; \lambda)$$

$$\implies \boldsymbol{\beta}^{t+1} = \frac{1}{1 + \mu_t \lambda \gamma} \psi_{\text{st}}\left(\boldsymbol{\beta}^t - \mu_t \mathbf{X}^\top (\mathbf{X}\boldsymbol{\beta}^t - \mathbf{y}); \mu_t \lambda\right)$$

- soft thresholding followed by multiplicative shrinkage

# Interpretation: majorization-minimization

---

$$f_{\mu_t}(\boldsymbol{\beta}, \boldsymbol{\beta}^t) := \underbrace{f(\boldsymbol{\beta}^t) + \langle \nabla f(\boldsymbol{\beta}^t), \boldsymbol{\beta} - \boldsymbol{\beta}^t \rangle}_{\text{linearization}} + \underbrace{\frac{1}{2\mu_t} \|\boldsymbol{\beta} - \boldsymbol{\beta}^t\|^2}_{\text{trust region penalty}}$$

majorizes  $f(\boldsymbol{\beta})$  if  $0 < \mu_t < \frac{1}{L}$ , where  $L$  is Lipschitz constant<sup>1</sup> of  $\nabla f(\cdot)$

Proximal gradient descent is a majorization-minimization algorithm

$$\boldsymbol{\beta}^{t+1} = \underbrace{\arg \min_{\boldsymbol{\beta}}}_{\text{minimization}} \left\{ \underbrace{f_{\mu_t}(\boldsymbol{\beta}, \boldsymbol{\beta}^t) + g(\boldsymbol{\beta})}_{\text{majorization}} \right\}$$

---

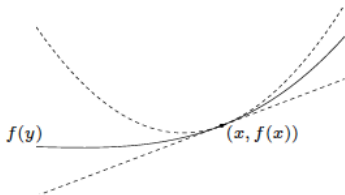
<sup>1</sup>This means  $\|\nabla f(\boldsymbol{\beta}) - \nabla f(\mathbf{b})\| \leq L\|\boldsymbol{\beta} - \mathbf{b}\|$  for all  $\boldsymbol{\beta}$  and  $\mathbf{b}$

# Assumptions for convergence

---

- $g : \mathbb{R}^n \mapsto \mathbb{R}$  is a continuous convex function, possibly nonsmooth;
- $f : \mathbb{R}^n \mapsto \mathbb{R}$  is a smooth convex function that is continuously differentiable with *Lipschitz constant*:

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



$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{L_f}{2} \|\mathbf{y} - \mathbf{x}\|_2^2, \quad \forall \mathbf{x}, \mathbf{y}$$

# Convergence rate of proximal gradient methods

## Theorem 4.2 (fixed step size; Nesterov '07)

Suppose  $g$  is convex, and  $f$  is differentiable and convex whose gradient has Lipschitz constant  $L$ . If  $\mu_t \equiv \mu \in (0, 1/L)$ , then

$$F(\beta^t) - F(\hat{\beta}) \leq O\left(\frac{\|\beta^0 - \hat{\beta}\|^2}{t\mu}\right)$$

- Step size requires an upper bound on  $L$ 
  - For LASSO problems, we have  $L = \sigma_{\max}(\mathbf{A}^\top \mathbf{A})$ .
- May prefer backtracking line search to fixed step size
- **Question:** can we further improve the convergence rate?

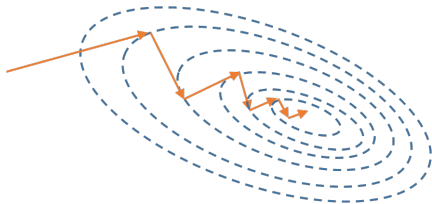
## Nesterov's accelerated gradient methods

- We will first examine Nesterov's acceleration method (1983) for smooth convex functions;
- We then extend it to optimizing composite functions, using FISTA (Beck and Teboulle, 2009), which extends Nesterov's method to proximal gradient methods.

# Nesterov's accelerated method

---

Problem of gradient descent: zigzagging



**Nesterov's idea:** include a momentum term to avoid overshooting



# Accelerated Gradient descent

---

$$\text{minimize}_{\beta \in \mathbb{R}^p} f(\beta)$$

where  $f(\beta)$  is convex and smooth (differentiable)

---

**Algorithm 4.3** Accelerated Gradient descent

---

**for**  $t = 0, 1, \dots$ :

$$\mathbf{b}^t = \mathbf{b}^{t-1} - \mu_t \nabla f(\mathbf{b}^{t-1})$$

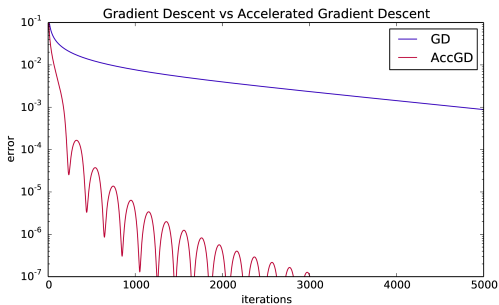
$$\beta^t = \beta^t + \underbrace{\alpha_t (\beta^t - \beta^{t-1})}_{\text{momentum term}}$$

where  $\mu_t$ : step size / learning rate,  $\alpha_t$  the the extrapolation parametre

---

# With/Without Acceleration

---

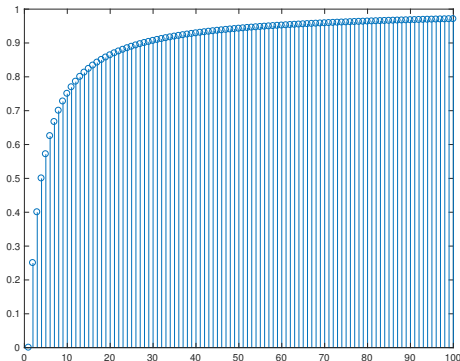


# Momentum parameter

---

A simple (but mysterious) choice of extrapolation parameter

$$\alpha_t = \frac{t-1}{t+2}$$



Another choice: let  $s_1 = 1$ ,  $s_{t+1} = \frac{1 + \sqrt{1 + 4s_t^2}}{2}$ , and  $\alpha_t = \frac{s_t - 1}{s_{t+1}}$ .

# Accelerated proximal method (FISTA)

---

**Nesterov's idea:** include a momentum term to avoid overshooting

$$\begin{aligned}\beta^t &= \text{prox}_{\mu_t g} \left( \mathbf{b}^{t-1} - \mu_t \nabla f \left( \mathbf{b}^{t-1} \right) \right) \\ \mathbf{b}^t &= \beta^t + \underbrace{\alpha_t \left( \beta^t - \beta^{t-1} \right)}_{\text{momentum term}} \quad (\text{extrapolation})\end{aligned}$$

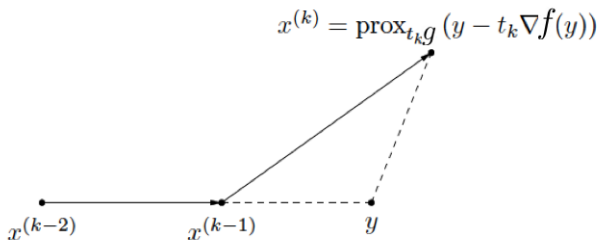
- A simple (**but mysterious**) choice of extrapolation parameter

$$\alpha_t = \frac{t-1}{t+2}$$

- Fixed size  $\mu_t \equiv \mu \in (0, 1/L)$  or backtracking line search
- Same computational cost per iteration as proximal gradient

# Interpretation

---



# Convergence of FISTA

---

## Theorem 4.3 (Nesterov '83, Nesterov '07, Beck & M. Teboulle '09)

Suppose  $f$  is differentiable and convex and  $g$  is convex. If one takes  $\alpha_t = \frac{t-1}{t+2}$  and a fixed step size  $\mu_t \equiv \mu \in (0, 1/L)$ , then

$$F(\beta^t) - F(\hat{\beta}) \leq O\left(\frac{1}{t^2}\right)$$

- Improves upon  $O(\frac{1}{t})$  convergence than proximal gradient method.
- in general un-improvable

# Numerical experiments (for lasso)

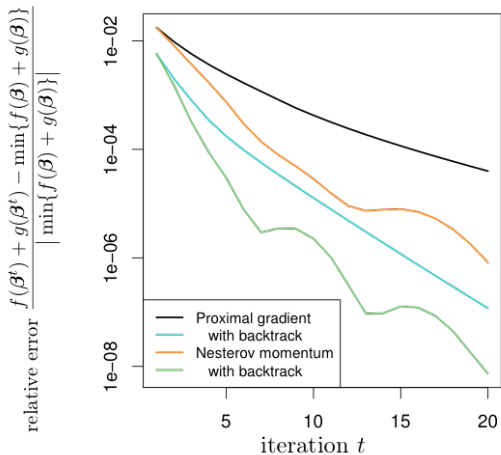


Figure credit: Hastie, Tibshirani, & Wainwright '15

# Computational-Statistical Trade-off

---

If there is indeed a ground truth  $\beta^*$  and we wish  $\hat{\beta}$  is close to  $\beta^*$ ; we have a sequence of  $\{\beta_t\}$  and hope  $\beta_t$  converges to  $\hat{\beta}$ . At a fixed  $t$ , we may bound

$$\|\beta_t - \beta^*\|_2 \leq \underbrace{\|\beta_t - \hat{\beta}\|_2}_{\text{computational error}} + \underbrace{\|\hat{\beta} - \beta^*\|_2}_{\text{statistical error}}$$



# Reference

---

- [1] "*Proximal algorithms*," Neal Parikh and S. Boyd, *Foundations and Trends in Optimization*, 2013.
- [2] "*Convex optimization algorithms*," D. Bertsekas, *Athena Scientific*, 2015.
- [3] "*Convex optimization: algorithms and complexity*," S. Bubeck, *Foundations and Trends in Machine Learning*, 2015.
- [4] "*Statistical learning with sparsity: the Lasso and generalizations*," T. Hastie, R. Tibshirani, and M. Wainwright, 2015.
- [5] "*Model selection and estimation in regression with grouped variables*," M. Yuan and Y. Lin, *Journal of the royal statistical society*, 2006.
- [6] "*A method of solving a convex programming problem with convergence rate  $O(1/k^2)$* ," Y. Nesterov, *Soviet Mathematics Doklady*, 1983.

# Reference

---

- [7] "*Gradient methods for minimizing composite functions*," Y. Nesterov, *Technical Report*, 2007.
- [8] "*A fast iterative shrinkage-thresholding algorithm for linear inverse problems*," A. Beck and M. Teboulle, *SIAM journal on imaging sciences*, 2009.