

"Just Relax"

(1)

We have been studying the (P_0^c) , which are still combinatorial.

We have seen that approximating algorithm - greedy ones - often do succeed in finding the sparsest solution.

Now we attempt a different strategy: changing the problem objective.

So instead of

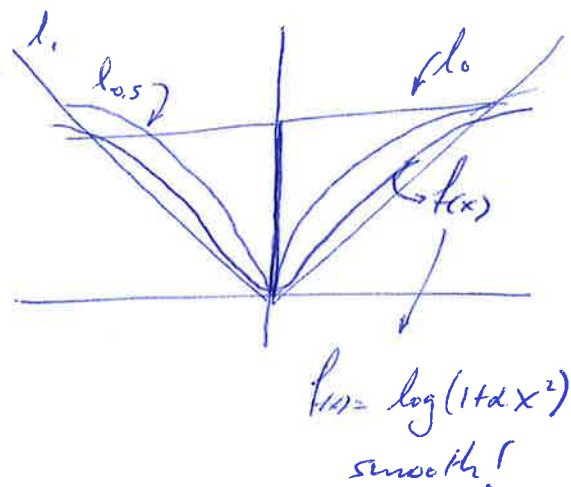
$$(P_0^c): \min_x \|b - Ax\|_2^2 \text{ s.t. } \|x\|_0 \leq k$$

which is equivalent to $\min_x \|b - Ax\|_2^2 + \lambda \|x\|_0$ for some λ ,

Consider $\min_x \|b - Ax\|_2^2 + \lambda \|x\|_p^p$, $\Leftrightarrow 0 < p \leq 1$

In fact, any "sensible" options could work.

The l_1 is the "closest" norm to the l_0 problem, and we already saw in the first class that it promotes sparse solutions.



So we'll first study $(P_1): \min_x \|x\|_1 \text{ s.t. } Ax = b$.

How ~~the~~ good is this? How and when can we hope it'll work?

→ We'll follow Tropp's analysis → simpler. (No, let's do Bickel's).

A : $n \times m$, (full rank).

S : subset of columns $S \subseteq [m]$.

Jump to Bickel's Analysis \rightarrow

Exact Recovery Condition (ERC).

For a given support S , A satisfies the ERC if

$$\max_{i \notin S} \|A_S^+ a_i\|_1 < 1. \quad \text{ERC}(A, S).$$

Intuitive, this requires ^{the least-squares (ℓ_2) solution} ~~every solution~~ to the system

$$\min_x \|A_S x - a_i\|_2^2 \quad i \notin S$$

to have ℓ_1 norm less than 1.

Theorem: For $x: S = \text{supp}\{x\}$, $Ax = b$, If ERC is met, BP recovers x .

Proof: (~~We know~~ x is feasible i.e. $Ax = b$.)

Consider a candidate solution $y \neq x$, $Ax = Ay$.

and ~~res~~ $e = x - y \in \mathcal{N}(A)$.

$$Ae = A_S e_S + A_{\bar{S}} e_{\bar{S}} = 0 \Rightarrow e_S = -A_S^+ A_{\bar{S}} e_{\bar{S}}.$$

$$\Rightarrow \|e_S\|_1 = \|-A_S^+ A_{\bar{S}} e_{\bar{S}}\|_1 \leq \|A_S^+ A_{\bar{S}}\|_1 \|e_{\bar{S}}\|_1$$

the ℓ_1 -operator norm $\|B\|_1 = \max_v \frac{\|Bv\|_1}{\|v\|_1}$: maximum ℓ_1 of a column.

Analysis of Lasso

Observation: The Lasso error satisfies a cone constraint:

$$\|e_S^c\|_1 \leq \alpha \|e_S\|_1, \quad \text{where } e = x_0 - \hat{x},$$

$$\hat{x} = \underset{x}{\operatorname{argmin}} \|y - Ax\|_2^2 + \lambda \|x\|_1$$

for a proper value of λ .

Proof: ~~sketch~~ argument.

for the constraint version, i.e. st. $\|\hat{x}\|_1 \leq \|x_0\|_1 = R$, since \hat{x} is feasible:

$$e = \hat{x} - x_0 = x_0 - \hat{x}$$

$$\begin{aligned} \|x_0\|_1 &= \|x_0^s\|_1 \geq \|\hat{x}\|_1 = \|x_0^s + e\|_1 = \|x_0^s + e_S\|_1 + \|e_S^c\|_1 \\ &\geq \|x_0^s\|_1 - \|e_S\|_1 + \|e_S^c\|_1 \\ \Rightarrow \|e_S^c\|_1 &\leq \|e_S\|_1 \quad (\alpha = 1). \end{aligned}$$

Restricted Convexity

Note that for the function $f(x) = \|y - Ax\|_2^2$ to be strongly convex, we need $A^T A \succ 0$. Why strongly convex? Because we want a unique solution.

Since $A^T A$ is rank deficient (n), $A^T A \not\succ 0$. i.e. many eigenvalues are zero.

Thus, similarly to what was done for RIP, we can define

$$\min_V \frac{V^T X^T X V}{\|V\|_2^2} \geq \gamma \quad \forall V \in C \quad \underline{\|XV\|_2^2 \geq \gamma \|V\|_2^2}$$

In particular, this will be useful for $C(s, \alpha) = \{V: \|V_S^c\|_1 \leq \alpha \|V_S\|_1\}$

We are now going to prove a very general result:

Theorem:

Let $y = Ax^0 + w$, $\|x^0\|_0 = K$, and X satisfies the REP with $\gamma > 0$.
Then any $\hat{x} \in \arg\min_x \|y - Ax\|_2^2$ st $\|x\|_1 \leq R = \|x^0\|_1$ satisfies

$$\|\hat{x} - x^0\|_2 \leq \frac{4}{\gamma} \sqrt{K} \|A^T w\|_\infty.$$

~~Proof~~ • Note that these results are deterministic.

- Different probabilistic results can be obtained imposing distributions over A and/or w .

Proof:

Recall that x^0 is feasible, and \hat{x} is optimal. Thus,

$$\|y - A\hat{x}\|_2^2 \leq \|y - Ax^0\|_2^2 = \|w\|_2^2$$

~~$$\|y - A\hat{x}\|_2^2 \leq \|y - Ax^0\|_2^2 = \|w\|_2^2$$~~

$$|x_i y_i| = |x_i| |y_i| \leq \sum_i |x_i| |y_i|$$

$$\|fg\|_1 \leq \|f\|_p \|g\|_q$$

for $\frac{1}{p} + \frac{1}{q} = 1$

$$\|Ax^0 - A\hat{x} + w\|_2^2 = \|Ae\|_2^2 + \|w\|_2^2 - 2w^T Ae \leq \|w\|_2^2.$$

Holder's Ineq.

$$\Rightarrow \frac{\|Ae\|_2^2}{2} \leq |w^T Ae| \leq \|w\|_1 \|A^T w\|_\infty \|e\|_1$$

But $\|e\|_1 = \|e_s\|_1 + \|e_{s^c}\|_1 \leq 2\|e_s\|_1 \leq 2\sqrt{K}\|e\|_2.$

On the other hand, we know that $e \in \mathcal{C}(\gamma)$

$$\Rightarrow \|Ae\|_2^2 \geq \|e\|_2^2 \gamma$$

$$\Rightarrow \|e\|_2^2 = \|\hat{x} - x^0\|_2^2 \leq \frac{4}{\gamma} \sqrt{K} \|A^T w\|_\infty.$$

relevant cases:

a) if $w = 0$, then $\hat{X} = X^*$: exact recovery.

if $E[b_i] = 0$

$$\text{Var}(\sum_{i=1}^n b_i) = \sum_{i=1}^n \text{Var}(b_i) + \sum_{i \neq j} \text{Cov}(b_i, b_j) = \sum_{i=1}^n E(b_i^2) = n \sigma^2$$

c) Say $\|a_i\|_2 = 1$, and $w \sim \mathcal{N}(0, \sigma^2)$.

Then $a_i^T w \sim \mathcal{N}(0, \sigma^2)$, and so we can use the

Gaussian tail bound:

$$P[|a_i^T w| \geq t] \leq 2e^{-t^2/2\sigma^2} \rightarrow \text{(Chernoff bound). for Gaussian distribution}$$

Now, since

$$\|A^T w\|_\infty = \max_i |a_i^T w|$$

$$P_r[\|A^T w\|_\infty \geq t] = P_r[\max_i |a_i^T w| \geq t] \leq 2e^{-t^2/2\sigma^2 + \log(m)}$$

$$= P_r[\cup_i \{|a_i^T w| \geq t\}] \leq \sum_{i=1}^m P_r[|a_i^T w| \geq t]$$

union bound

$$= 2pe^{-t^2/2\sigma^2}$$

$$= 2e^{-t^2/2\sigma^2 + \log(p)}$$

Now, letting $t = \sigma \sqrt{\tau \log m}$

$$\Rightarrow P_r[\|A^T w\|_\infty \geq t] \leq 2e^{-\frac{1}{2}(\tau-2)\log(m)} \quad \text{with } (\tau \geq 2)$$

$$\text{Thus, } \left| \|\hat{X} - X^*\|_2 \leq \frac{4\sigma}{\gamma} \sqrt{k\tau \log m} \right| \quad \text{with } P_r > 1 - 2e^{-\frac{1}{2}(\tau-2)\log m}$$

b) Say, $\|a_i\|_2 = 1$, $\|w\|_2 \leq \epsilon$. - adversarial.

$$\Rightarrow \|\hat{X} - X^*\|_2 \leq \frac{4\sigma}{\gamma} \sqrt{k} \cdot \|A^T w\|_\infty \leq \frac{4\sigma \sqrt{k}}{\gamma} \|w\|_2 = \frac{4\sigma \sqrt{k}}{\gamma} \cdot \epsilon$$

What we are now missing is how to know/bound γ ,
or how to know if $\gamma > 0$ - which is what really matters.

Recall that $\text{REP}(S, \alpha)$

$$\gamma = \min_{e \neq 0} \frac{\|Ae\|_2^2}{\|e\|_2^2} > 0, \quad \forall e \in \mathcal{C} = \{e: \|e_{\bar{S}}\|_1 \leq \alpha \|e_S\|_1\}.$$

Note that, considering a support S , and Lasso solution \tilde{x} ,

$$\begin{aligned} \|Ae\|_2^2 &= \|A(e_S + e_{\bar{S}})\|_2^2 = \|Ae_S\|_2^2 + \|Ae_{\bar{S}}\|_2^2 + 2e_{\bar{S}}^T \overset{A^T A}{G} e_S \\ &\geq \|Ae_S\|_2^2 - 2|e_{\bar{S}}^T G e_S|. \end{aligned}$$

Now,

$$|e_{\bar{S}}^T G e_S| = |\tilde{e}_{\bar{S}}^T G_{\bar{S}, S} \tilde{e}_S| \leq \|\tilde{e}_{\bar{S}}\|_1 \|G_{\bar{S}, S} \tilde{e}_S\|_{\infty}.$$

(maybe show?)

$$\leq \|\tilde{e}_{\bar{S}}\|_1 \mu(A) \cdot \|\tilde{e}_S\|_1.$$

$$\leq \|e_S\|_1^2 \mu(A) \alpha.$$

$$\leq k \cdot \|e_S\|_2^2 \mu(A) \alpha.$$

On the other hand, $\|e_S\|_0 = k$, so

$$\delta_k \leq (k-1) \mu(A).$$

$$\|Ae_S\|_2^2 \geq (1 - \delta_k) \|e_S\|_2^2 \geq (1 - (k-1) \mu(A)) \|e_S\|_2^2.$$

Putting everything together:

$$\text{if } e \in \mathcal{C}(S, \alpha), \quad \frac{\|Ae\|_2^2}{\|e\|_2^2} \geq 1 - (k-1) \mu(A) - 2k \mu(A) \alpha.$$

$$> 1 - k \mu(A) - 2k \mu(A) \alpha = 1 - k \mu(A) (1 + 2\alpha).$$

Thus, $\gamma > 1 - k \mu(A) (1 + 2\alpha)$

and while $k < \frac{1}{\mu(A)(1+2\alpha)} \Rightarrow \gamma > 0$ ✓

Compressed Sensing

Motivation: if signals/images are "compressible" - why can't we just acquire them in this "compressed" representation?

Say now that signal $x \in \mathbb{R}^m$. We want to take measurements

$$b_i = \langle a_i, x \rangle, \quad i = 1 \dots n.$$

How small can n be while still be able to recover x ?

Note that if $n < m \rightarrow$ undetermined. But if x is a "natural" signal, then is it "sparsifiable", say, under a unitary transformation $\Phi: x = \Phi \alpha, \|\alpha\|_0 \ll m$.

Thus, we are trying to solve

$$\min \|\alpha\|_0 \text{ st. } b = A \Phi \alpha.$$

$$\begin{matrix} n \\ \downarrow \end{matrix} b = \begin{matrix} \boxed{A} \\ m \end{matrix} = \begin{matrix} \boxed{A} \end{matrix} \begin{matrix} \boxed{\Phi} \\ \alpha \end{matrix}$$

How small can we make n ?

Before we say that BP recovers α if

$$\|\alpha\|_0 < \frac{1}{3\mu(A)}$$

$$\left(\text{or } \|\alpha\|_0 < \frac{1}{2} \left(1 + \frac{1}{\mu(A)} \right) \right)$$

How small can $\mu(A)$ be?

from the Welch Bound, $\mu(A) \geq \sqrt{\frac{m-n}{n(m-1)}} = \Omega(1/\sqrt{n})$.

$$\Rightarrow K \lesssim \frac{1}{\mu(A)} \Rightarrow K \lesssim \sqrt{n} \Rightarrow n \gtrsim K^2$$

this is not really good: Say $m = 100^2$, and say $K = 0.1m$

then $n \gtrsim 1000$

too many samples!

This is the "squared bottleneck", resulting from the pessimistic analysis relying on $\mu(A)$.

The alternative is to go back to other matrix ~~quanta~~ characterizations (eg. RIP), even if they are not computable.

In fact, we've seen before that BP recovers X ~~if~~ (or α) if A satisfies the restricted eigenvalue property. So which A satisfy this, and for what n ?

Turns out (Raskutti, Wainwright, Yu) that for Gaussian matrices $n \times m$ with iid rows sampled from $\mathcal{N}(0, \Sigma)$,

$$\text{if } n > C \frac{\sigma_{\max}^2 (1+\alpha)^2}{\gamma^2} \cdot K \cdot \log(m)$$

and Σ satisfies the REP(δ)
($\|\Sigma^{1/2} e_i\|_2 \geq \gamma \|e_i\|_2$)
e.i.d.

then A satisfies the R.E. property with parameters $(\alpha, \gamma/\delta)$

with Probability $> 1 - e^{-cn}$.

C, C', C'' universal positive constants.

historical reasons (and elegance), let's show a result
based on the Nullspace Property:

NSP: A matrix A satisfies the NSP for support S if

$$\|v_S\|_1 < \|v_{S^c}\|_1 \quad \forall \quad v \in \mathcal{N}(A) \setminus \{0\}.$$

Necessary and sufficient condition:

Given $A: n \times m$, every vector $x \in \mathbb{R}^m$ supported on set S is the unique solution to $\min_z \|z\|_1$ st. $Az = Ax$ if and only if A satisfies the NSP relative to S .

proof: first implication:

Assume every vector x is the unique minimum of $\|z\|_1$ st. $Az = Ax$.

In particular, consider $v \in \mathcal{N}(A) \setminus \{0\}$, and $v_S = \arg\min_z \|z\|_1$ st. $Az = Av_S$.

Now, since $Av = 0 = Av_S + Av_{S^c} \Rightarrow (-v_{S^c})$ is also feasible since

$Av_S = A(-v_{S^c})$. But v_S is the unique minimizer $\Rightarrow \|v_{S^c}\|_1 > \|v_S\|_1$.

Second implies: assume NSP holds for S .

Consider x supported on S , and a different vector $z \neq x$. ~~Not feasible~~
such that $Ax = Az$. Further, $v = x - z \in \mathcal{N}(A) \setminus \{0\}$.

Then:

$$\|x\|_1 = \|x - z_S + z_S\|_1 \leq \|x - z_S\|_1 + \|z_S\|_1 = \|v_S\|_1 + \|z_S\|_1$$

$$\text{NSP} \rightarrow < \|v_S\|_1 + \|z_S\|_1$$

$$= \|z_S\|_1 + \|z\|_1 = \|z\|_1$$

$\Rightarrow \|x\|_1$ is minimal, and any other feasible solution has $\geq \|x\|_1$.

Thus, we have the uniform recovery guarantee:

Given A , every k -sparse $x \in \mathbb{R}^m$ is the unique solution to BP iff A satisfies NS of order k (for all s).

A classic Result is that RIP is sufficient for recovery, because it implies the USP:

If A has the RIP with $\delta_{2k} < \frac{1}{2}$, then A has the USP of order k . \Rightarrow it recovers all k -sparse vectors.

What matrices satisfy the RIP? many random ones do:

Let $A: n \times m$ be a subGaussian matrix: all entries are ^{independent} ~~each~~ zero-mean subgaussian random variables with variance 1 and subgaussian parameters β, k :

$$\mathbb{P}(|A_{ij}| \geq t) \leq \beta e^{-kt^2} \quad \forall t > 0,$$

Then, $\exists c > 0$ such that $(\frac{1}{\sqrt{m}} A)$ has RIP with $\delta_k < \delta$ if

$$n \geq \frac{2c}{\delta^2} \cdot k \log\left(\frac{em}{k}\right)$$

with $\mathbb{P} \geq 1 - 2e^{-\frac{\delta^2 mn}{2c}},$

is results then guarantees recovery of all k -sparse vectors w.h.p. by BP, for such random matrices.

This almost answers all questions for compressed sensing. Recall however that x is typically not sparse, but rather $x = \Phi \alpha$, α is.

This is ok, because if A is RIP, Φ unitary, $(A\Phi)$ also is RIP with different constants.

Note that Φ is irrelevant for sensing: $b = A x$.

It should only be used for recovery/decoding:

$$\min_{\alpha} \|\alpha\|_0 \text{ s.t. } b = A\Phi\alpha.$$

Optimization for Lasso

Remember Gradient Descent:

$$\min_{x \in \mathbb{R}^n} f(x), \quad \text{for } f: \text{convex and differentiable.}$$

$$\text{G.D: } x^{k+1} = x^k - \epsilon_k \nabla f(x^k).$$

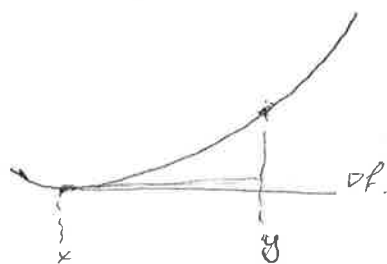
If ∇f is Lipschitz (i.e. f is smooth), then G.D converges with rate of $O(1/k)$.

But what if f is non-smooth?

Subgradients

for a convex $f: \mathbb{R}^n \rightarrow \mathbb{R}$.

$$f(y) \geq f(x) + \nabla f(x)^T (y-x), \quad \forall x, y.$$



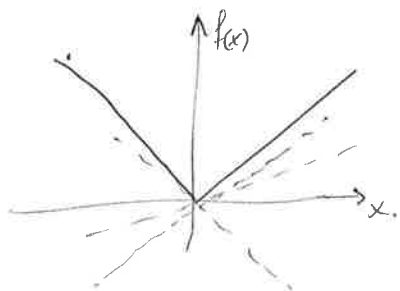
i.e. a linear approximation "underestimates" $f(x)$.

A subgradient of a convex f at x is any $g \in \mathbb{R}^n$ such that

$$f(y) \geq f(x) + g^T (y-x), \quad \forall y.$$

- Always exists. (unlike the gradient).

- If f is differentiable at x , then $g = \nabla f(x)$ uniquely.



$$\text{for } x \neq 0, \quad g = \nabla f(x) = \text{sign}(x).$$

$$\text{for } x = 0, \quad g \in [-1, 1].$$

The set of all subgradients of a convex f is called the subdifferential $\partial f(x) = \{g \in \mathbb{R}^n : g \text{ is a subgradient of } f \text{ at } x\}$.

Note: - always non empty (for convex f).

- If f is differentiable at x , $\partial f = \{\nabla f(x)\}$.

Thus, they give optimality conditions: for convex f :

$$\boxed{f(x^*) = \min_x f(x) \iff 0 \in \partial f(x^*)}$$

Subgradient Method:

proof: $f(y) \geq f(x^*) + g^*(y - x^*)$

if $g=0 \Rightarrow f(y) \geq f(x^*) \forall y$.

$$x^{k+1} = x^k - t_k g^k, \text{ where } g^k \in \partial f(x^k)$$

To get convergence, we need the stepsize to go to zero (but not too fast!)

Subgradient descent has ~~an~~ a convergence rate of $O(1/k)$.
 \Rightarrow slow.

Algorithms for lasso/BP

(a)

While any convex optimization algorithm would suffice, we'll pay close attention to first order methods.

Getting some intuition: Unitary Case.

Assume A : unitary/orthonormal.

then $\min_x \frac{1}{2} \|y - Ax\|_2^2 + \lambda \|x\|_1$

$$\min_x \frac{1}{2} \|b - x\|_2^2 + \lambda \|x\|_1, \text{ where } b = A^T y.$$

$$h(x) = \sum_{i=1}^m \frac{1}{2} (b_i - x_i)^2 + \lambda |x_i|. \text{ : separable!}$$

$$h(x_i) = \begin{cases} \frac{1}{2} (b_i - x_i)^2 + \lambda & \text{if } x_i \geq 0 \\ \frac{1}{2} (b_i - x_i)^2 - \lambda & \text{if } x_i < 0 \end{cases}$$

Note that $h(x_i)$ is differentiable if $x_i > 0$ or $x_i < 0$, but not at $x_i = 0$.

thus, if $x_i > 0$:

recall that

$$h'(x_i) = \lim_{\epsilon \rightarrow 0} \frac{h(x_i + \epsilon) - h(x_i)}{\epsilon}$$

must be unique!

$$x_i - b + \lambda = 0 \Rightarrow x_i = \max(b - \lambda, 0) \text{ if } b > \lambda$$

$$\text{if } x_i < 0: x_i = \min(b + \lambda, 0) \text{ if } b < -\lambda$$

if $|b| \leq \lambda$, then the min $h(x_i)$ must be attained at the only point ~~at~~ it's not differentiable: $x = 0$.

Thus, for A : orthonormal, $x = \arg\min_x h(x) = S_\lambda(A^T y)$

this is $\text{prox}(b)$ / show.

$$\hookrightarrow S_\lambda(b) = \text{sign}(b) (|b| - \lambda)_+$$

What if $A: n \times m$?

$$F(x) = \frac{1}{2} \|y - Ax\|_2^2 + \lambda \|x\|_1$$

Following a Majorization - Minimization approach:

Consider $g(x, z) = F(x) + d(x, z) \geq F(x)$. and $F(x) = g(x, x)$.

$$d(x, z) = \frac{c}{2} \|x - z\|_2^2 - \frac{1}{2} \|Ax - Az\|_2^2 \geq 0.$$

We want $d(x, z)$ to be strongly convex $\Rightarrow H(x) = c \cdot I - A^T A \succ 0$.

$$\Rightarrow c > \|A\|_2^2 = \lambda_{\max}(A^T A).$$

Then: $g(x, z) = \frac{1}{2} \|y - Ax\|_2^2 + \lambda \|x\|_1 + \frac{c}{2} \|x - z\|_2^2 - \frac{1}{2} \|Ax - Az\|_2^2$.

$$= \frac{1}{2} \|y\|_2^2 + \frac{1}{2} \|Ax\|_2^2 - x^T A^T y + \frac{c}{2} \|x\|_2^2 + \frac{c}{2} \|z\|_2^2 - c x^T z - \frac{1}{2} \|Ax\|_2^2 - \frac{1}{2} \|Az\|_2^2 + x^T A^T Az.$$

$$= \frac{c}{2} \|x\|_2^2 - x^T (A^T y + cz - A^T A z) + \text{const.} + f(z, y) + \lambda \|x\|_1$$
$$cz - A^T (Az - y) \triangleq v.$$

$$g(x, z) = \frac{1}{2} \|x\|_2^2 - \frac{1}{c} x^T v + \frac{\lambda}{c} \|x\|_1 + \tilde{f}(z, y).$$

$$g(x, z) = \frac{1}{2} \|x - \frac{1}{c} v\|_2^2 + \frac{\lambda}{c} \|x\|_1 + \tilde{f}'(z, y).$$

$$\Rightarrow \arg\min_x g(x, z) = S_{\lambda/c}(v) = S_{\lambda/c}\left(z - \frac{1}{c} A^T (Az - y)\right).$$

In MM, we do $x^{k+1} = \arg\min_x g(x, x^k)$

$$\Rightarrow x^{k+1} = S_{\lambda/c}\left(x^k - \frac{1}{c} A^T (Ax^k - y)\right). \quad \text{ISTA.}$$

is a particular case of Proximal-Gradient Methods, (5) which solve $\min_x f(x) + g(x)$, : convex.

where f : convex & smooth, g : convex but non-smooth (non-diff.).

$$x^{k+1} = \text{prox}_g \left(x^k - \frac{1}{c} \nabla f(x^k) \right) = T_{\lambda}^{f,g}(x^k)$$

Recall that subgradient descent achieves convergence of $\mathcal{O}(1/\sqrt{k})$.

We'll see (in part) that ISTA (prox.grad) achieves $\mathcal{O}(1/\sqrt{k})$

basically at same cost. (comput.), and it can be further accelerated

to $\mathcal{O}(1/k^2)$.
$$\begin{cases} x^{k+1} = S_{\lambda}(y^k - \frac{1}{c} A^T(Ay^k - b)) \\ y^{k+1} = x^{k+1} + t^k(x^{k+1} - x^k) \end{cases}$$

To make analysis shorter, we'll need a couple of lemmas:

L1: Monotonic decrease:

Let $F(x) = f(x) + g(x)$ as above, and let $x^{k+1} = \text{ISTA}(x^k)$. Then with stepsize $c = L = \lambda_{\max}(A^T A)$. Then

$$F(x^k) - F(x^{k+1}) \geq \frac{c}{2} \|x^k - x^{k+1}\|^2$$

Note that this implies that $F(x^{k+1}) \leq F(x^k)$.

Observe that $f(x)$ is smooth: $\|\nabla f(x) - \nabla f(y)\|_2^2 \leq L \|x - y\|_2^2$

for $L = \lambda_{\max}(A^T A)$

L2:

~~Lemma 2: For any $x, y \in \mathbb{R}^n$, we have~~

$$\frac{2}{L} [F(x) - F(T_A^{l,y}(y))] \geq \|x - T_A^{l,y}(y)\|_2^2 - \|x - y\|_2^2 \quad \forall x, y$$

Theorem: For $F(x) = f(x) + g(x)$ as before, let $\{x^k\}$ be the sequence generated by ISTA with $L \geq \lambda_{\max}(A^T A)$. Then, for any $x^* \in X_{\text{opt}}$ and $k \geq 0$,

$$F(x^k) - F_{\text{opt}} \leq \frac{L}{2k} \|x^{(0)} - x^*\|^2$$

Proof: From Lemma 2, set $x = x^*$ and $y = x^k$ to obtain:

$$\frac{2}{L} (F(x^*) - F(x^{k+1})) \geq \|x^* - x^{k+1}\|_2^2 - \|x^* - x^k\|_2^2$$

Summing for all $k \geq 0$:

$$\frac{2}{L} \sum_{n=0}^{k-1} (F(x^*) - F(x^{n+1})) \geq \|x^* - x^k\|_2^2 - \|x^* - x^0\|_2^2$$

$$\Rightarrow \frac{2}{L} \sum_{n=0}^{k-1} (F(x^{n+1}) - F(x^*)) \leq \|x^* - x^0\|_2^2$$

$F(x^*) = F_{\text{opt}}$: by def. Recall that $F(x^{n+1}) \leq F(x^n)$.

Thus

$$k(F(x^k) - F(x^*)) \leq \sum_{n=0}^{k-1} F(x^{n+1}) - F_{\text{opt}} \leq \|x^* - x^0\|_2^2 \leq \frac{L}{2}$$

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

//