

(1)

Class 3: Greedy Algorithms. - "Greed is Good" J.T.

Recall we are after $(P_0): \min_x \|x\|_0 \text{ s.t. } Ax = b.$
or $(P_0): \min \|Ax - b\|_2^2 \text{ s.t. } \|x\|_0 \leq k.$

One can decompose this problem in two parts:

- Identify the support
- Find coefficients.

Note that b is simple! indeed, if we know $S^* = \text{supp}(x^*)$,

then $x^* = \arg \min_x \|Ax - b\|_2^2 \text{ s.t. } \text{supp}(x) = S^*, |S^*| = k.$

$$= \arg \min_{x_{S^*}} \|A_{S^*} x_{S^*} - b\|_2^2$$

$$\Rightarrow A_S^T (A_S x_S - b) = 0$$

$$\Rightarrow \hat{x}_S = (A_S^T A_S)^{-1} A_S^T b$$

Least Squares solution.

So often the main problem is finding the correct support.

Suppose someone ~~told~~ gave us ~~some~~ $b = Ax^*$ and told us $\|x^*\|_0 = 1$, and unique. What would we do?

We could perform m tests: $\min_z \|a_j z - b\|_2^2 = \epsilon(j).$

which is minimized by $\hat{z}_j = \frac{a_j^T b}{\|a_j\|_2^2}$.

The error of approximating b with the j th atom is then

$$\|a_j \hat{z}_j - b\|_2^2 = \left\| a_j \frac{a_j^T b}{\|a_j\|_2^2} - b \right\|_2^2 = \|b\|_2^2 - \frac{(a_j^T b)^2}{\|a_j\|_2^2}.$$

If one of these $\epsilon(j)$ terms become $\epsilon(j^*) = 0 \Rightarrow$ found solution!

This costs $O(m)$ flops, i.e. m 'checks' of $\langle a_j, b \rangle$.

If $\|X^*\|_0 > 1$, however, and no $e(j)$ is found to be zero, then we should move to all $\binom{m}{k}$ possibilities, which is $O(m^k)$ prohibitive!

Greedy Approaches attempt to (hopefully) find the correct support by building it more or less sequentially.

Orthogonal Matching Pursuit:

→ Choose the next atom that best approximates a residual.

Init: $x^{(0)} = 0$, $r^{(0)} = b - Ax^{(0)} = b$, $S = \text{supp}\{x^{(0)}\} = \emptyset$.

for $k = 1, \dots, k_0$:

• Find "best" next atom:

$$j_0 = \underset{j}{\operatorname{argmin}} \|a_j z_j - r^{k-1}\|_2^2 = \underset{j}{\operatorname{argmax}} \left| \frac{a_j^T r^{k-1}}{\|a_j\|_2^2} \right|$$

• Update support: $S \leftarrow S \cup \{j_0\}$.

• Update solution: $x_S^{(k)} = \underset{x}{\operatorname{argmin}} \|A_S x - b\|_2^2 = A_S^+ b$.

• Update residual: $r^{(k)} = b - A x^{(k)}$.

some observations:

1) why orthogonal?

Note that by computing $X_s^{(k)}$ again $\|b - A_s X - \|_2^2$

$$\Rightarrow A_s^T (A_s X - b) = 0.$$

$$A_s^T r^{(k)} = 0.$$

\Rightarrow residual is orthogonal to selected atoms.

Q: can an atom be selected twice?

2) Matching Pursuit; (Mallat '93).

solution is updated by: $X_{(j_0)}^k = X_{(j_0)}^{k-1} + \underbrace{\frac{a_{j_0}^T r^{k-1}}{\|a_{j_0}\|_2^2}}_{Z_{j_0}}$

which replaces the Least Squares step.
Thus cheaper.

Likewise, there are different variations of these methods with different level of complexity. Eg: Least-Squares OMP, where $m - |S^k|$ Least-Squares steps are evaluated to select the next atom - replacing the cheaper inner product in OMP. This, though, is typically much more expensive.

3) Note that normalization does not influence which atoms are selected. Thus, easier to work with normalized atoms.

$$\tilde{A} = A W \Rightarrow \tilde{X}^{\text{OMP}}, \Rightarrow X = W \tilde{X} \quad \text{since} \\ b = \tilde{A} \tilde{X} = A \underbrace{W \tilde{X}}_X.$$

How fast does the residual decay? i.e: how fast is an approximation constructed?

We'll look at MP for simplicity, and $\|a_i\|_2 = 1 \quad \forall i$.

In MP, ~~minimize~~ $X^k = X^{k-1}$ where $j_0 = \arg \min_j \langle a_j, r^{k-1} \rangle$

$$X_{(j_0)}^k = X_{(j_0)}^{k-1} + \underbrace{g_{j_0}^T b}_{z_{j_0}}$$

$$\text{Then, } r^k = b - AX^k = r^{k-1} - (a_{j_0}^T r^{k-1}) a_{j_0}.$$

$$\begin{aligned} \text{Thus, } E(j_0) = \|r^k\|_2^2 &= \|r^{k-1}\|_2^2 + (a_{j_0}^T r^{k-1})^2 - 2(a_{j_0}^T r^{k-1})^2 \\ &= \|r^{k-1}\|_2^2 - (a_{j_0}^T r^{k-1})^2 \\ &= \|r^{k-1}\|_2^2 - \max_{1 \leq j \leq m} (a_j^T r^{k-1})^2. \end{aligned}$$

Define the ~~depression~~ ^{decaying} factor $\delta(A, v) = \max_{1 \leq j \leq m} \frac{|a_j^T v|}{\|v\|_2}^2$

and the "universal" decaying factor: $\boxed{\delta(A) = \inf_v \delta(A, v)}$

This is the "worst" vector: the one that leads to a smallest decrease in energy of the residual. Thus,

$$\|r^k\|_2^2 = \|r^{k-1}\|_2^2 - \delta(A, r^{k-1}) \cdot \|r^{k-1}\|_2^2$$

$$\|r^k\|_2^2 \leq \|r^{k-1}\|_2^2 (1 - \delta(A)).$$

This leads to $\|r^k\|_2^2 \leq (1 - \delta(A))^k \|b\|_2^2 \Rightarrow$ exponential decay.

at now: Is $\delta(A) > 0$? Yes, because A full rank.

so no nullspace: $\exists v \neq 0$:
 $Av = 0$.

Best $\delta(A)$? $\delta(A) \leq 1/\sqrt{n}$. (say, for orthonormal A).

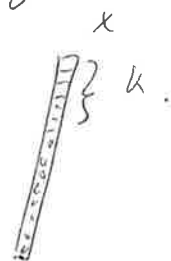
So, OMP approximates exponentially fast. But to what?

Does it succeed in recovering the sparsest vector?

Let $b = Ax$, as input. Assume without loss of generality that the k non-zero are the first k entries, and ordered:

$$\Rightarrow b = \sum_{j=1}^k a_j x_j$$

$$x_i \geq x_j \text{ if } i < j.$$



In the beginning, $r^{(0)} = b$, and $j_0 = \arg \max_j |a_j^T b|$.

For this first stage to succeed, $j_0 \in [1, k]$ - one of the correct atoms - . Thus, we require: $\max_{i \in S} |a_i^T b| > \max_{j \notin S} |a_j^T b|$: necessary and sufficient.

Thus, a sufficient condition is: $|a_{j_0}^T b| > \max_{j > k} |a_j^T b|$ ~~which~~.
 where x_1 is the largest coeff.

$$\Rightarrow \underbrace{\left| \sum_{t=1}^k x_t a_t^T a_i \right|}_a > \max_{j > k} \underbrace{\left| \sum_{t=1}^k x_t a_t^T a_j \right|}_b$$

We will construct a lower bound for left, upper for right, and enforce it.

$$a): \left| \sum_{t=1}^k x_t a_t^T a_i \right| \geq |x_1| - \sum_{t=2}^k |x_t a_t^T a_i| \quad \text{-reverse tri. inequality-}$$

$$\geq |x_1| - |x_1| \sum_{t=2}^k |a_t^T a_i|$$

$$\geq |x_1| - |x_1| (k-1) \mu(A)$$

$$= |x_1| (1 - (k-1) \mu(A)).$$

$$b_j): \quad \left| \sum_{i=1}^k x_i a_i^T a_j \right| \leq \sum_{i=1}^k |x_i| |a_i^T a_j| \leq \|x\|_1 k \cdot \mu(A).$$

Thus, combining (a), (b); ~~and~~ we require that

$$1 - (k-1)\mu(A) \geq k \cdot \mu(A) \Rightarrow k < \frac{1}{2} \left(1 + \frac{1}{\mu(A)} \right).$$

This condition guarantees that a correct atom is chosen in the first step.

The residual is updated by subtracting from b (or $r^{(k)}$) a term that is proportional to the chosen atom (or atoms, in general). Thus, r^k is still a linear combination of the same k atoms in b , at most. ~~non-orthogonal~~

Repeating the same steps with $r^{(1)}$ for b guarantees the recovery of the second atom, and so forth. Further, due to orthogonality no atom is chosen twice, and the algorithm terminates after k iterations, retrieving all k non-zeros.

Thus, we have the following result:

For $b = Ax$, A : full rank, if a solution exist: $\|x\|_0 \leq \frac{1}{2} \left(1 + \frac{1}{\mu(A)} \right)$, OMP finds it in $\|x\|_0$ steps.

Some computational ~~constraints~~ Considerations:

• The bottleneck is the L.S. But this can be ~~appeared~~ alleviated by using a Choleski decomposition progressively. All in all, OMP is $O(k^3 + kmn)$.

• Further improvement can be obtained if one uses many signals.

→ Demo of Recovery. OMP (and maybe MP).

observation: they work way beyond the bound.
too pessimistic.

→ Better high-probability bounds exist. For example, while this holds when $k \leq \frac{1}{2\mu}$, Schmass (19!!) showed that

$$s \leq \frac{1-\alpha}{\mu^2 \log(m)} \text{ suffices with h.p.}$$

What about noisy/real data? Then we would rather want to solve, given $b = Ax^* + v$, $\|v\|_2 \leq \epsilon$, $\|x^*\|_0 \leq k$.

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

OMP still works: if \exists solution x : $\|x\|_0 \leq \frac{1}{2} \left(1 + \frac{1}{\mu(A)}\right) - \frac{1}{\mu(A)} \cdot \frac{\epsilon}{\|x_{\min}\|}$,

then a) OMP recovers the true support of x^* ,

$$b) \|\hat{x}_{\text{omp}} - x^*\|_2^2 \leq \frac{\epsilon^2}{1 - \mu(A)(k-1)}.$$

proof sketch:

$$b = \sum_{i=1}^k x_i a_i + v$$

As before, we require the sufficient condition

$$|a_i^T b| > \max_{j > k} |a_j^T b| \quad \text{-- recall ordering.}$$

$$\left| \sum_{t=1}^k x_t a_t^T a_i + a_i^T v \right| > \max_{j > k} \left| \sum_{t=1}^k a_j^T a_t x_t + a_j^T v \right|$$

LHS:

$$\left| \sum_{j=1}^k x_j a_j^T a_i + a_i^T v \right| \geq \left| \sum_{j=1}^k x_j a_j^T a_i \right| - |a_i^T v|$$

($\|a_i\|_2 = 1$.)

$$\geq |x_i| (1 - (k-1)\mu(A)) - \epsilon$$

RHS:

$$\left| \sum_{t=1}^k x_t a_t^T a_j + a_j^T v \right| \geq \left| \sum_{t=1}^k x_t a_t^T a_j \right| + |a_j^T v|$$
$$\geq |x_i| \mu(A) \cdot k + \epsilon.$$

\Rightarrow Enforcing

$$|x_i| (1 - (k-1)\mu(A)) - \epsilon \geq |x_i| \mu(A) k + \epsilon.$$

$$\Rightarrow k \leq \frac{1}{2} \left(1 + \frac{1}{\mu(A)} \right) - \frac{1}{\mu(A)} \cdot \frac{\epsilon}{|x_i|}.$$

recall that $|x_i| \geq |x_i|$ i.e.s., and we have assumed $k \leq \frac{1}{2} \left(1 + \frac{1}{\mu(A)} \right) - \frac{1}{\mu(A)} \cdot \frac{\epsilon}{|x_{\min}|}$.

Thus, the above holds for every $i \in S$.

One can then reiterate the argument, except one has to make sure that the sparsity of the representation of the residual and the amount of noise is preserved, until all atoms are recovered.

To see the second claim, note that.

$$\begin{aligned}\hat{x}_{\text{omp}} &= \underset{x}{\text{argmin}} \|A_S x - b\|_2^2 = A_S^+ b = A_S^+ (A_S x_S^* + v) \\ &= x_S^* + A_S^+ v.\end{aligned}$$

$$\Rightarrow \|\hat{x}_{\text{omp}} - x^*\|_2^2 \leq \|A_S^+\|_2^2 \underbrace{\|v\|_2^2}_{\leq \epsilon^2}$$

$$\leq \frac{1}{\underbrace{\lambda_{\min}(A_S^T A_S)}_{(1-\delta_S)}} \leq \frac{1}{1 - \mu(A)(k-1)}$$

Gershgorin Disk Theorem.

$$A_S^* = U \Sigma V^T$$

$$A_S^+ = (V \Sigma U^T U \Sigma V^T)^+ V \Sigma U^T$$

$$= V \Sigma^{-1} U^T$$

$$\Rightarrow \|A_S^+\|_2 = \|\Sigma^{-1}\|_2 = \frac{1}{\sigma_{\min}(A_S)}$$

$$\Rightarrow \|A_S^+\|_2^2 = \frac{1}{\lambda_{\min}(A_S^T A_S)}$$

Denoising!

From this we see that if we seek to "denoise" with omp:

$$\|\hat{x} - x\|_2^2 \leq \frac{\epsilon^2}{1 - \mu(A)(k-1)} \rightarrow \epsilon^2 \text{ that we started with!}$$

This is because of the worst-case assumption. In fact, one can obtain great denoising!

Consider the "oracle" estimator, having retrieved the correct support:

$$\hat{x}_{\text{omp}} = A_S^+ b = A_S^+ (A_S x + v).$$

But now assume $v \sim \mathcal{N}(0, I \sigma^2)$.

$$\text{tr}(A^*) = \sum_i \lambda_i(A^*) \leq$$

$$\begin{aligned}\Rightarrow \mathbb{E}[\|x - \hat{x}_{\text{omp}}\|_2^2] &= \mathbb{E}[\|A_S^+ v\|_2^2] \\ &= \sigma^2 \text{tr}((A_S^T A_S)^{-1}) \leq \frac{k \cdot \sigma^2}{1 - \delta_S}\end{aligned}$$

and likewise $\mathbb{E}[\|x - \hat{x}_{\text{null}}\|_2^2] \leq \frac{k \sigma^2}{1 - \delta_S}$ which $\ll n \sigma^2$

When considering Gaussian Noise v , it can be shown that AMP has a nearly-oracle performance, in the sense that achieves an MSE of

$$\|\hat{x}_{\text{amp}} - x_0\|_2^2 \leq (1+\alpha) k \sigma^2 \log(m) = \underbrace{C \log(m)}_{\text{oracle}} \cdot k \sigma^2$$

This holds for $v \sim N(0, \sigma^2 I)$ and if w.h.p.

with probability at least $1 - \frac{1}{m^{\alpha} \sqrt{\pi} (1+\alpha) \log m}$ for some $\alpha > 0$,

$$|x_{\min}| - (2k-1)\mu |x_{\min}| \geq 2\sigma \sqrt{\frac{2k \log m}{(1+\alpha)}}$$

Demoising Demo: . . .

[Ben-Haim, Eldar, Elad].

IHT: What if signals are very large? Recall that AMP is $O(k^3 + 2k \log m)$

Cheaper alternatives are preferred. One of them is IHT, which can be thought of Projected Gradient Descent:

Given $\min_x \|b - Ax\|_2^2 \text{ s.t. } \|x\|_0 \leq k,$

do $x^{n+1} = H_k(x^n - \eta A^T(A^T x^n - b))$

↓
Projection onto set of k -sparse vectors:

$$H_k(a) = \arg\min_x \frac{1}{2} \|x - a\|_2^2 \text{ s.t. } \|x\|_0 \leq k.$$

H_k : keeps k -largest entries.

This is cheap: $O(\underbrace{T_{\text{nn}}}_{\text{iterations}} + \underbrace{T_n \log(m)}_{\text{sorting}})$.

Theorem (Blumensath & Davies) '09.

(6).

If $b = Ax^* + v$, $\|x\|_0 = s$, A has RIP : $\delta_{3s} < \frac{1}{\sqrt{32}}$,
 then

$$\|x^k - x^*\|_2 \leq 2^{-k} \|x^*\|_2 + 5 \|v\|_2.$$

There \exists also accelerated versions of this, AHT, which perform quite well.

Derivation of IHT: given $\min_x \overbrace{\|b - Ax\|_2^2}^{l(x)} \text{ s.t. } \|x\|_0 \leq k$ - hard -

Let's consider a Majorization-Minimization approach:

For $l(x)$, hard to minimize; construct $g(x, z) \geq l(x) \forall z$, and

Then, we can do $g(x, x) = l(x)$.

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

Consider $g(x, z) = \|b - Ax\|_2^2 - \|Ax - Az\|_2^2 + \|x - z\|_2^2$

$g(x, z) \geq l(x)$ if $\|A\|_2 \leq 1$, so consider normalizing it.

\rightarrow because then $\|A(x-z)\|_2^2 \leq \|x-z\|_2^2$

Now, $\min_x g(x, z)$: but $g(x, z) = \|x\|_2^2 - 2x^T(z - A^T(Az - b)) + \|Az\|_2^2 + \|z\|_2^2$

So $g(x, z) = \|x - (z - A^T(Az - b))\|_2^2 + \underbrace{\|Az\|_2^2 + \|z\|_2^2}_{\text{etc.}} + \|b\|_2^2$

So $\min_x g(x, z) \text{ s.t. } \|x\|_0 \leq k$

$\min_x \|x - (z - A^T(Az - b))\|_2^2 \text{ s.t. } \|x\|_0 \leq k \Rightarrow x = \arg\min_x (g(x, z)) =$

So: let $z = x^k$.

$$x = H_k [z - A^T(Az - b)]$$

$$x^{k+1} = H_k [z - A^T(Az - b)]$$