

Dictionary Learning

(1)

We've been working with the problem:

$$\min_{\gamma} \frac{1}{2} \|y - D\gamma\|_2^2 \text{ s.t. } \|\gamma\|_0 \leq k.$$

What is the correct/best dictionary to use? We can learn it from a collection of data.

Let $Y = [y_1, \dots, y_n] \in \mathbb{R}^{n \times N}$. Then:

$$\min_{\Gamma, D} \frac{1}{2} \|Y - D\Gamma\|_2^2 \text{ s.t. } \begin{cases} \|\gamma_i\|_0 \leq k \\ \|\mathbf{d}_i\|_2 = 1 \end{cases} = \text{Dict. Learning Problem.}$$

The most common approach is alternating minimization.

$$\left[\begin{array}{l} \text{- fix } D, \min_{\gamma_i} \frac{1}{2} \|y_i - D\gamma_i\|_2^2 \text{ s.t. } \|\gamma_i\|_0 \leq k. \\ \text{- which can be solved with OMP, relaxed to } l_1, \text{ etc. -} \\ \text{- fix } \Gamma, \text{ solve for } D: \\ \min_D \frac{1}{2} \|Y - D\Gamma\|_2^2 \text{ s.t. } \|\mathbf{d}_i\|_2 = 1. \end{array} \right.$$

The latter can be done in a number of different ways.

Most notable, Method of Optimal Directions (MOD) finds the "optimal" D for a given Γ :

$$D = (Y\Gamma^T)(\Gamma\Gamma^T)^{-1}.$$

Alternatively, one could minimize one atom/column at a time:

$$\begin{aligned}\min_{d_j} \|\mathbf{Y} - \mathbf{D}\mathbf{\Gamma}\|_2^2 &= \frac{1}{2} \|\mathbf{Y} - \sum_{i \neq j} d_i \mathbf{\Gamma}_i^T - d_j \mathbf{\Gamma}_j^T\|_2^2 \\ &= \frac{1}{2} \|\mathbf{E}_j - d_j \mathbf{\Gamma}_j^T\|_2^2\end{aligned}$$

This leads to the k -SVD.

Optimization Guarantees:

The DL problem is non-convex, not just because of the $\|\cdot\|_1$, but also due to the product $\mathbf{D}\mathbf{\Gamma}$. Thus, convergence to a global optimum is complicated, and usually one can guarantee convergence to a local minimum.

Some algorithms have some convergence guarantees, alas they tend to be quite involved. Thus, instead of going through one of these, we will take a simpler approach: we will revisit convergence guarantees for gradient descent, and then we will see how to use these in a non-convex case as above.

Gradient Descent

(2).

Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$.

We will assume f is twice differentiable, β -smooth and α -strongly convex. Different guarantees can be obtained relaxing some of these.

For $t = 1$ to T ,

$x^{t+1} = x^t - \eta \nabla f(x^t)$.

end

β -smoothness: f is β -smooth if

$$\|\nabla f(y) - \nabla f(x)\| \leq \beta \|y - x\|, \quad \forall x, y.$$

Note that if f is twice diff, then this is equivalent to

$$\|\nabla^2 f(x)\| \leq \beta \quad \forall x.$$

Strong Convexity: f is α -strongly convex if

$$(y-x)^T \nabla^2 f(x) (y-x) \geq \alpha \|y-x\|^2, \quad \forall x, y.$$

This implies that one "can fit" a quadratic function underneath $f \Rightarrow f$ is not "too flat".

Note that ~~as~~ this also implies

$$f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{\alpha}{2} \|y-x\|^2$$

We will prove the following main result:

Thm: Let f be twice differentiable, β -smooth and α -convex. Let x^* be the (unique) minimizer of f , and $\eta \leq 1/\beta$. Then Gradient Descent satisfies: (From x_1).

$$f(x_t) - f(x^*) \leq \beta \left(1 - \frac{\eta\alpha}{2}\right)^{t-1} \|x_1 - x^*\|^2$$

In the proof, we will employ the following lemma:

Lemma: if f is twice diff, β -smooth and α -strongly convex, then ^(key).

$$\boxed{\nabla f(x_t)^T (x_t - x^*) \geq \frac{\alpha}{4} \|x_t - x^*\|^2 + \frac{1}{2\beta} \|\nabla f(x_t)\|^2}$$

"Sufficient correlation".

Proof of Thm:

For simplicity, let $\alpha' = \frac{\alpha}{4}$; $\beta' = \frac{1}{2\beta}$. Consider

$$\begin{aligned} \|x_{t+1} - x^*\|^2 &= \|x_t - \eta \nabla f(x_t) - x^*\|^2 \\ &= \|x_t - x^*\|^2 - 2\eta \nabla f(x_t)^T (x_t - x^*) + \eta^2 \|\nabla f(x_t)\|^2 \\ &\leq \|x_t - x^*\|^2 - 2\eta \left(\alpha' \|x_t - x^*\|^2 + \beta' \|\nabla f(x_t)\|^2 \right) + \eta^2 \|\nabla f(x_t)\|^2 \\ &= \|x_t - x^*\|^2 (1 - 2\eta\alpha') + (\eta^2 - 2\eta\beta') \|\nabla f(x_t)\|^2 \end{aligned}$$

as long as $\beta \eta < 1/\beta \Rightarrow$

(3)

$$\|x_{t+1} - x^*\|^2 \leq (1 - 2\eta\alpha') \|x_t - x^*\|^2. \quad (1)$$

Finally, note that $f(x^*) \geq f(x_t) + \nabla f(x_t)^T (x^* - x_t)$

$$\begin{aligned} \Rightarrow f(x_t) - f(x^*) &\leq \nabla f(x_t)^T (x_t - x^*) \leq \beta \|x_t - x^*\|^2 \quad (2) \\ &= (\nabla f(x_t) - \nabla f(x^*))^T (x_t - x^*) \xrightarrow{\beta\text{-smooth.}} \\ &\quad \hookrightarrow 0. \end{aligned}$$

from (1) and (2):

$$f(x_t) - f(x^*) \leq \beta (1 - 2\eta\alpha') \|x_t - x^*\|^2$$

//

Proof of Lemma:

Recall we need to get, under strong-convexity and smoothness assumptions

$$\nabla f(x_t)^T (x_t - x^*) \geq \frac{\alpha}{2} \|x_t - x^*\|^2 + \frac{1}{2\beta} \|\nabla f(x_t)\|^2$$

for the first part, note that

$$f(x^*) \geq f(x) + \nabla f(x)^T (x^* - x) + \frac{\alpha}{2} \|x - x^*\|^2, \quad \text{by } \alpha\text{-strong convexity.}$$

But $f(x) \geq f(x^*)$, thus: $\nabla f(x)^T (x - x^*) \geq \frac{\alpha}{2} \|x - x^*\|^2. \quad (a)$

Then, recall the Lagrange remainder theorem:

let $f: \mathbb{R}^n \rightarrow \mathbb{R}$, twice diff., then for $t \in [0, 1]$, $x' = ty + (1-t)x$,
some \uparrow

$$\nabla f(x) = \nabla f(y) + \nabla^2 f(x') (x - y).$$

This can be proven by Taylor expansion of $f(x)$ and the intermediate value theorem.

Let $y = x^*$, and note that $\nabla f(x^*) = 0$. Thus,

$$\nabla f(x) = \nabla^2 f(x') (x - x^*).$$

And so:

$$\nabla f(x)^T (\nabla^2 f(x'))^{-1} \nabla f(x) = \nabla f(x)^T (x - x^*)$$

Recall that β -smoothness implies

$$\nabla f(x)^T (\nabla^2 f(x'))^{-1} \nabla f(x) \geq \frac{1}{\beta} \|\nabla f(x)\|_2^2$$

$$\Rightarrow \nabla f(x)^T (x - x^*) \geq \frac{1}{\beta} \|\nabla f(x)\|_2^2 \quad (b)$$

Thus, from (a), (b):

$$\nabla f(x)^T (x - x^*) \geq \frac{L}{4} \|x - x^*\|_2^2 + \frac{1}{4\beta} \|\nabla f(x)\|_2^2$$

(Consider commenting on SGD).

Following Analysis from [Arora, Ge, Ma, Moitra, Simple, Efficient neural Algorithms for S.C.].

(4)

the cool thing about this analysis is that one does not necessarily have to be taking the gradient as a direction. As long as this direction satisfies a version of the key lemma from before "sufficient correlation", the analysis follows through.

Making this more precise:

Def: A vector g_t is $(\alpha', \beta', \epsilon_t)$ -correlated with a point x^* if, $\forall \epsilon$,

$$g_t^\top (x - x^*) \geq \alpha' \|x_t - x^*\|^2 + \beta' \|g_t\|^2 - \epsilon_t.$$

We saw before that if f is twice diff., ~~then~~ α -strongly convex and β -smooth, then $\nabla f(x_t)$ is $(\frac{\alpha}{2}, \frac{1}{2\beta}, 0)$ -correlated with the optimal solution.

Further, the proof for the previous theorem generalizes directly ~~to~~ this case too.

Thm: Suppose g_t is $(\alpha', \beta', \epsilon_t)$ -correlated with a point x^* and $\eta < 2\beta'$. Then "abstract gradient descent" $x^{t+1} = x^t + \eta g_t$ from x_1 satisfies

$$\|x_t - x^*\| \leq \left(1 - \frac{\eta \alpha'}{2}\right)^{t-1} \|x_1 - x^*\|^2 + \underbrace{\frac{\max \epsilon_t}{\alpha'}}_{\text{bias}}.$$

Back To Dictionary Learning.

We will assume a stochastic setting, where

a) $\text{Supp}(\sigma) \sim \text{u.a.r.}$, $\|\sigma\|_0 = k$.

b) $\gamma_{\mathcal{S}_i} \begin{cases} +1 & \text{wp. 0.5} \\ -1 & \text{wp. 0.5} \end{cases}$, ~~and~~ pairwise independent conditioned on the support \mathcal{S} .

We observe $y = D\sigma$.

Recall that we are after

$$\min_{D, \Pi} \|Y - D\Pi\|_F^2 \text{ s.t. } \|\sigma_i\|_0 \leq k, \forall i.$$

Observation: instead of viewing alternating methods as minimizing a known function, we can think of them as minimizing an ~~the~~ unknown - convex function.

In other words: the problem

$$\min_D \|Y - D\Pi\|_2^2 \text{ is convex. It's just that we don't really know } \Pi.$$

What we will do is show that alternating methods still move in a direction "sufficiently correlated" with the true gradient:

$$f_1(D) = \|Y - D\Pi\|^2 \approx f_2 = \|Y - D\Pi^*\|^2$$

and we expect $\nabla f_2(D) \approx \nabla f_1(D)$
here know \rightarrow we don't know.

δ_i moves almost in ~~that~~ ^{ideal} direction of the true parameter ~~direction~~ δ_i (5)(6)

We will analyse the simple algorithm:

for $t=0$ to T :

- "Decoding": $\hat{\delta}^{(t)} = H_{1/2}(\hat{D}^T y^{(t)}) \quad \forall i$

- "Update": $\hat{D} \leftarrow \hat{D} + \eta \sum_{i=1}^n (y_i - \hat{D} \hat{\delta}_i) \text{sign}(\hat{\delta}_i)^T$

Metric / Distance :

Two matrices D, \hat{D} , (with normalized columns) are (δ, k) -close if \exists permutation & sign flip of the columns of \hat{D} : $B = \hat{D}P$ such that $\|b_i - d_i\|_2 \leq \delta \quad \forall i$

and $\|B - D\| \leq k \|D\|$.

Decoding Succeeds :

Assume $D: n \times m$ is μ -incoherent and $b = D\delta$, with $\|M_0\| = k \leq \frac{1}{10\mu \cdot \log(n)}$, and \hat{D} is $(\frac{1}{\log n}, 2)$ -close to D .

Then, the decoding stage succeeds: (with ~~high~~ high probability).

$$\text{sign}[H_{1/2}(\hat{D}^T b)] = \text{sign}(\delta).$$

A proof sketch is the following:

Consider the j^{th} inner-product:

$$\begin{aligned}\langle \hat{d}_j, b \rangle &= \hat{d}_j^T D x = (\hat{d}_j - d_j + d_j)^T d_j x_j + \hat{d}_j^T \sum_{i \neq j} d_i x_i \\ &= x_j + \underbrace{(\hat{d}_j - d_j)^T d_j x_j}_{\|\cdot\| \leq \frac{1}{\log(n)}} + \underbrace{\sum_{i \neq j} \hat{d}_j^T d_i x_i}_{\ll}\end{aligned}$$

So w.h.p., if $j \in S$, $|\hat{d}_j^T b| > 1/2$, and $|\hat{d}_j^T b| < 0$ otherwise.

Update:

We'll assume the update is given by

$$g = E[(y - \hat{D} \hat{x}) \text{sign}(\hat{x})^T],$$

and the j^{th} atom is updating through $g_j = E[(y - \hat{D} \hat{x}) \text{sign}(\hat{x}_j)]$.

Denote by F the event that encoder recovers the support of x , which holds w.h.p. if \hat{D}, D are $(1/\log n, 2)$ -close. Then:

$$\begin{aligned}g_j &= E[(y - \hat{D} \hat{x}) \text{sg}(\hat{x}_j) 1_F] + E[(y - \hat{D} \hat{x}) \text{sg}(\hat{x}_j) 1_{F^c}] \\ &= E[(y - \hat{D} \hat{x}) \text{sg}(\hat{x}_j) 1_F] \pm \epsilon, \quad \text{say } \|\epsilon\| < n^{-c}.\end{aligned}$$

$$\begin{aligned}\text{so } g_j &= E[(y - \hat{D} H_0(\hat{D}^T y)) \text{sg}(\hat{x}_j) 1_F] \pm \epsilon; \quad \text{say } S = \text{supp}(x) \\ &= E[(I - \hat{D}_S \hat{D}_S^T) y \cdot \text{sg}(\hat{x}_j) 1_F] \pm \epsilon \\ &= E[(I - \hat{D}_S \hat{D}_S^T) D_S x \cdot \text{sg}(\hat{x}_j)] \pm \epsilon.\end{aligned}$$

(6)

ow, using subconditioning: (first sampling the support S and then the values x_S),

$$\begin{aligned}
 g_j &= \mathbb{E}_S \left[\mathbb{E}_{x_S} \left[(I - \hat{D}_S \hat{D}_S^T) D x \cdot \text{sgn}(x_j) \mid S \right] \right] \pm \epsilon \\
 &= \mathbb{E}_S \left[\mathbb{E}_{x_S} \left[(I - \hat{D}_S \hat{D}_S^T) D_j x_j \text{sgn}(x_j) \mid S \right] \right] \pm \epsilon \quad \text{because } x_{S^c} \text{ are iid given } S. \\
 &= \mathbb{E} \left[p_j \cdot (I - \hat{D}_S \hat{D}_S^T) d_j \right] \pm \epsilon, \quad p_j = \mathbb{E} \left[x_j \cdot \text{sgn}(x_j) \mid S \right]. \quad \text{or uncorrelated.} \\
 &= p_j \cdot \mathbb{E} \left[(I - \hat{D}_S \hat{D}_S^T) d_j \right] \pm \epsilon \\
 &= p_j \mathbb{E} \left[(I - \hat{d}_j \hat{d}_j^T) d_j \right] + p_j \mathbb{E} \left[\hat{D}_j \hat{D}_j^T d_j \right] \pm \epsilon. \\
 &= p_j q_j (I - \hat{d}_j \hat{d}_j^T) d_j + p_j \hat{D}_j Q \hat{D}_j^T d_j \pm \epsilon. \\
 &\quad \downarrow \qquad \qquad \qquad \downarrow \\
 &\quad q_j = P[j \in S] \qquad \qquad \text{diag}(\{q_{i,j}\}_i)
 \end{aligned}$$

Note that then

$$g_j = p_j q_j \left(\underbrace{d_j}_{\hat{d}_j^T d_j \approx 1} - \underbrace{\lambda_j \hat{d}_j}_{\|\hat{e}_j\| \leq O(k/n)} \right) + \hat{e}_j \pm \epsilon.$$

This "shows" that g_j is mostly correlated with the optimal direction $(d_j - \hat{d}_j)$:

lemma: If $g_j = \alpha (d_j - \hat{d}_j) + v$, with $\|v\| \leq \alpha \|d_j - \hat{d}_j\| + \xi$,

then g_j is $(\alpha, \frac{1}{100\alpha}, \frac{\xi^2}{\alpha})$ -correlated with d_j^*

$\Rightarrow g_j$ is $(\alpha, \beta, \epsilon)$ -correlated with d_j , with $\alpha = \Omega(k/n)$, $\beta \geq \Omega(m/k)$,

$\epsilon = O(k^3/mn)$. \Rightarrow Each Grad. step makes progress.

Corollary: $O(1/\sqrt{n})$

If \hat{D} is $(2\delta, 2)$ -near D , and $\eta \leq \min_i (p_i q_i (1-\delta))$, then

$$\|\hat{d}_j^{k+1} - d_j\|^2 \leq (1 - 2\alpha\eta) \|d_j^k - d_j\|^2 + O\left(\frac{k}{n}\right).$$

$\Rightarrow \hat{D}_j$ converges geometrically to d_j until the column-wise error is $O(\sqrt{k/n})$.