Nonnegative Matrix Factorization via (alternating) Projected Gradient Descent

Andersen Ang

Mathématique et recherche opérationnelle, UMONS, Belgium Homepage: angms.science

> First draft: August 2, 2017 Last update: November 27, 2020

The approximate NMF problem

• Given $\mathbf{M} \in \mathbb{R}^{m \times n}_+$, find $\mathbf{W} \in \mathbb{R}^{m \times r}_+$, $\mathbf{H} \in \mathbb{R}^{r \times n}_+$ by solving

 $[\mathbf{W}^*, \mathbf{H}^*] = \underset{\mathbf{W}, \mathbf{H}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{M} - \mathbf{W}\mathbf{H}\|_F^2$ subject to $\mathbf{W} \ge 0, \mathbf{H} \ge 0.$

- \blacktriangleright a non-convex optimization problem in two variables ${\bf W}$ and ${\bf H}$
- \blacktriangleright a constrained optimization problem: $\mathbf{W},\,\mathbf{H}$ have to be nonnegative
- ► a NP-Hard problem, see Vavasis2008: On the complexity of nonnegative matrix factorization
- ► There are many approaches to solve NMF:
 - Alternating projected gradient descent (this document)
 - Exact Block coordinate descent
 - Extrapolated inexact block coordinate descent

Projected gradient descent algorithm — PGD

▶ PGD is a way to solve **constrained** optimization problem

$$\min_{x\in\mathcal{Q}}f(x)$$

where $\ensuremath{\mathcal{Q}}$ is the constraint set.

▶ Starting from a initial feasible point $x_0 \in \mathcal{Q}$, PGD iterates

$$x_{k+1} = P_{\mathcal{Q}}\Big(x_k - t_k \nabla f(x_k)\Big)$$

where $P_{Q}(\cdot)$ is the projection operator, which itself is also an optimization problem:

$$P_{\mathcal{Q}}(x_0) = \arg\min_{x \in \mathcal{Q}} \frac{1}{2} ||x - x_0||_2^2,$$

i.e. given a point x_0 , P_Q finds a point $x \in Q$ that is closest to x_0 .

 Note that projected gradient = proximal gradient on indicator function

The PGD algorithm

Algorithm 1: PGD

Result: A solution x that approximately solves $\min_{x \in Q} f(x)$

Initialization Pick $x_0 \in Q$ while stopping condition is not met do $| x_{k+1} = P_Q(x_k - t_k \nabla f(x_k))$ end

- If f is L-smooth (∇f is L-Lipschitz), we can pick $t_k = \frac{1}{L}$.
- ▶ For nonnegative cone, $[\cdot]_+ = \max(\cdot, 0)$ is the projection.

The PGD algorithm on NMF

Algorithm 2: The PGD algorithm on NMF

Result: A solution x that approximately solves $\min_{x \in \mathcal{Q}} f(x)$

Initialization Pick inital matrices $W^0 \geq 0$ and $H^0 \geq 0$

while stopping condition is not met do

$$W^{k+1} = \left[W^k - t^k_W \nabla f(W^k; H^k) \right]_+;$$

$$H^{k+1} = \left[H^k - t^k_H \nabla f(H^k; W^{k+1}) \right]_+$$

end

- This algorithm is an inexact block coordinate descent with PGD update.
- That is, we update each block one-by-one, where the sub-problem, e.g.

$$\underset{\mathbf{W}\geq 0}{\operatorname{argmin}} \ \frac{1}{2} \|\mathbf{M} - \mathbf{W}\mathbf{H}^k\|_F^2$$

is not solved exactly: the projected gradient update on ${\bf W}$ does not solve the above minimization problem.

Lipschitz constant step size

• For
$$f(W,H) = \frac{1}{2} ||X - WH||_F^2$$
, the gradients are

$$\nabla_W f(W, H) = (WH - X)H^T, \quad \nabla_H f(W, H) = W^T(WH - X)$$

Key inequality: for any matrices A and B we have $||AB||_F \leq ||A||_F ||B||_2$.

$$\begin{aligned} \|\nabla_W f(W_1) - \nabla_W f(W_2)\|_F &= \|(W_1 H - X)H^T - (W_2 H - X)H^T\|_F \\ &= \|(W_1 - W_2)HH^T\|_F \\ &\leq \|HH^T\|_2 \cdot \|W_1 - W_2\|_F \end{aligned}$$

Similarly

$$\|\nabla_H f(H_1) - \nabla_H f(H_2)\|_F \le \|W^T W\|_2 \cdot \|H_1 - H_2\|_F$$

We can put $t_W^k = \frac{1}{\|HH^T\|_2}$ and $t_H^k = \frac{1}{\|W^T W\|_2}$ in the previous algorithm.



Reconstruction of the rice image at different iteration $W^{15}H^{15}$ W^1H^1 Original Image











 $W^{300}H^{300}$



 $W^{500}H^{500}$





 $W^{100}H^{100}$

8 / 18

PGD on NMF with optimal stepsize

Stepsize $t\ {\rm can}\ {\rm be}\ {\rm selected}\ {\rm that}\ {\rm achieve}\ {\rm maximum}\ {\rm decrease}\ {\rm in}\ {\rm objective}\ {\rm function}.$

$$t_W^* = \underset{t}{\operatorname{argmin}} f\left(W - tg(W), H\right), \quad t_H^* = \underset{t}{\operatorname{argmin}} f\left(W, H - tg(H)\right)$$

where $g(W) = \nabla_W f(W, H)$ and $g(H) = \nabla_H f(W, H)$. Then,

$$\frac{\partial f}{\partial t_W} = \operatorname{Tr} \left(X - WH + tg(W)H \right)^T g(W)H \quad \begin{vmatrix} \partial f \\ \partial t_H \end{vmatrix} = \operatorname{Tr} \left(X - WH + tWg(H) \right)^T Wg(H) t_W^* = \frac{\operatorname{Tr} \left((WH - X)^T g(W)H \right)}{\operatorname{Tr} \left(H^T g(W)^T g(W)H \right)} \qquad t_H^* = \frac{\operatorname{Tr} \left((WH - X)^T Wg(H) \right)}{\operatorname{Tr} \left(g(H)^T W^T Wg(H) \right)}$$

However, using optimal size, the computational cost is much higher than constant stepsize.

Rice image example



10 / 18



The Accelerated PG algorithm

- We can apply Nesterov's acceleration on PGD, which adds momentum on the sequence.
- To apply momentum on NMF, the easiest way is to apply it on the exact block coordinate descent scheme:
 - 1. Pick initial matrices $W^0 \in \mathbb{R}^{m \times r}_+$ and $H^0 \in \mathbb{R}^{r \times n}_+$ and parameters.
 - 2. On W, loop until stopping condition is met
 - 2.1 Update and project
 - 2.2 Extrapolate
 - 3. On H, loop until stopping condition is met
 - 3.1 Update and project
 - 3.2 Extrapolate
 - 4. Repeat 2 and 3 until converge.
- If we apply momentum on the inexact block coordinate descent scheme, the situation becomes much more complicated as the theory changes drastically. In this case the algorithm belong to the Bolte-Sabach-Teboulle's PALM, or Xu-Yin's APG, or Hien-Gillis-Patrinos's IBP.

Exact BCD algorithm with momentum for NMF

1. (On W) Initalize W^0, ϵ , set $V^0 = W^0$, $\lambda^0 = 0$, then loop until stopping condition is met :

1.1
$$\lambda^k = \frac{1}{2} \left(1 + \sqrt{1 + 4\lambda_{k-1}^2} \right), \ \gamma^k = \frac{1 - \lambda^{k-1}}{\lambda^k}$$

1.2 Update: $V^k = \max\left([W^k - t_W^k \nabla_W f(W^k; H)]_{ij}, \epsilon\right)$

- 1.3 Extrapolation : $W^k = (1 \gamma^k)V^k + \gamma V^{k-1}$
- 2. (On H) Initalize $H^0,\epsilon,$ set $G^0=H^0,\lambda^0=0$, then loop until stopping condition is met :

2.1
$$\lambda^{k} = \frac{1}{2} \left(1 + \sqrt{1 + 4\lambda_{k-1}^{2}} \right), \ \gamma^{k} = \frac{1 - \lambda^{k-1}}{\lambda^{k}}$$

2.2 Update: $G^{k} = \max \left([H^{k} - t_{H}^{k} \nabla_{H} f(H^{k}; W)]_{ij}, \epsilon \right)$

2.3 Extrapolation : $H^k = (1 - \gamma^k)G^k + \gamma G^{k-1}$

The acceleration is achieved by extrapolating current W, H by the coupling variables V, G with the extrapolation weights $\{\gamma^k\}$.

MATLAB code (click me)

The wrong way to apply momentum on NMF

Note in each Update step the other variable is held fix (without the superscript k). That is, the matrix H when updating W^k for k = 1, 2, ... is all the same. The following shows a wrong way to apply APGD on NMF :

- 1. Initilize W^0, H^0 and ϵ , set $V^0 = W^0, G^0 = H^0$, and $\lambda^0 = 0$
- 2. Loop until stopping condition is met

2.1 Update: $V^k = \max\left([W^k - t^k_W \nabla_W f(W^k; H^k)]_{ij}, \epsilon\right)$

 $\begin{array}{l} \text{2.2 Extrapolation}: \ W^k = (1-\gamma^k)V^k + \gamma V^{k-1} \\ \text{2.3 Update:} \ G^k = \max\left([H^k - t^k_H \nabla_H f(H^k;W^{k+1})]_{ij},\epsilon\right) \end{array}$

2.4 Extrapolation : $H^k = (1 - \gamma^k)G^k + \gamma G^{k-1}$

Why it is wrong : consider step 2, at k = 1

$$k = 1$$
 $V^1 = \max\left([W^1 - t_W^1 \nabla_W f(W^1; H^1)]_{ij}, \epsilon \right).$

$$k = 2 \qquad V^2 = \max\left([W^2 - t_W^2 \nabla_W f(W^2; \boldsymbol{H}^2)]_{ij}, \epsilon \right)$$

Unless variable H is already optimal (H converges and so $H^1 = H^2$) otherwise the update of H at k = 1 makes $H^2 \neq H^1$ and so the optimization problem on minimizing W at k = 1 is different from that at k = 2.

Therefore, for applying extrapolation directly on the block variable, we have to use the framework of PALM, APG or IBP.

PGD(inexact BCD with gradient update) vs APG(exact BCD with momentum) Example : MATLAB rice image (m, n, r) = (256, 256, 64)



PGD vs APGD

Here APGD takes about 50% of the computation to achieve the same amount of fitting error of PGD.



16 / 18

Further discussion

Recall that accelerated GD for unconstrained convex problem with single variable is not monotone, it has ripples in the trace of the objective function value.

In this case adaptive restart can be used. That is, if the extrapolated solution produces objective function value higher than that of the previous iterate, the extrapolated solution is thrown away, the momentum and extrapolation weight are all reset to initial state.

That is : if at step k, $f(W^k) > f(W^{k-1})$:

- $\blacktriangleright \ W^k = W^{k-1}$
- $\blacktriangleright \ V^k = W^k$
- $\blacktriangleright \ \lambda^k = 0$

Last page - summary

- Alternating projected gradient descent for NMF MATLAB code (click me)
- Alternating projected gradient descent = Inexact block coordinate descent
- Exact block coordinate descent with internal Nesterov's acceleration
 End of document