

Nonnegative Matrix Factorization with Regularization

Andersen Ang

ECS, Uni. Southampton, UK
andersen.ang@soton.ac.uk

Homepage angms.science

Version: February 7, 2024

First draft: June 6, 2017

Mathématique et de Recherche opérationnelle
Faculté polytechnique de Mons, UMONS
Belgium

Content

Regularized NMF

Inexact Block-Coordinate Descent

Orthogonal NMF

$$\nabla \|\mathbf{W}^\top \mathbf{W} - \mathbf{I}_r\|_F^2 = 4\mathbf{W}(\mathbf{W}^\top \mathbf{W} - \mathbf{I})$$

Quick-and-dirty code

Setup: the regularized numerical NMF Problem

- ▶ Input: $\mathbf{X} \in \mathbb{R}_+^{m \times n}$
 - ▶ m rows
 - ▶ n columns
 - ▶ all entries are nonnegative $X_{ij} \geq 0$

- ▶ Task: find two nonnegative matrices $\mathbf{W} \in \mathbb{R}_+^{m \times r}$ and $\mathbf{H} \in \mathbb{R}_+^{r \times n}$ such that

$$(\mathcal{P}) : (\mathbf{W}^*, \mathbf{H}^*) = \underset{\mathbf{W} \geq \mathbf{0}, \mathbf{H} \geq \mathbf{0}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 + \lambda g(\mathbf{W}, \mathbf{H}).$$

where we consider $(\mathbf{W}^*, \mathbf{H}^*)$ is a solution to Problem (\mathcal{P}) if the objective value is small enough

- ▶ $f(\mathbf{W}, \mathbf{H}) = \frac{1}{2} \|\mathbf{W}\mathbf{H} - \mathbf{X}\|_F^2$ is the primary objective function
 - ▶ it is called data-fitting
 - ▶ it is measured in Frobenius norm (i.e. Euclidean distance)
- ▶ $g(\mathbf{W}, \mathbf{H})$ is the regularization, here we assume g is differentiable
- ▶ $\lambda \geq 0$ is the weight
- ▶ $r \in \mathbb{N}$ is the factorization rank

Inexact Block-Coordinate Descent based on gradient descent update

$$(\mathcal{P}) : (\mathbf{W}^*, \mathbf{H}^*) = \underset{\mathbf{W} \geq \mathbf{0}, \mathbf{H} \geq \mathbf{0}}{\operatorname{argmin}} \underbrace{\frac{1}{2} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2}_{=: f} + \lambda g(\mathbf{W}, \mathbf{H}).$$

- ▶ We can solve \mathcal{P} by inexact Block-Coordinate Descent (BCD):

$$\mathbf{W}_1 = \operatorname{update}(\mathbf{W}, \mathbf{H}_0), \mathbf{H}_1 = \operatorname{update}(\mathbf{W}_1, \mathbf{H}), \mathbf{W}_2 = \operatorname{update}(\mathbf{W}, \mathbf{H}_1), \mathbf{H}_2 = \operatorname{update}(\mathbf{W}_2, \mathbf{H}), \dots$$

- ▶ For the “update”, there are many options, one option is projected gradient descent

$$\begin{aligned} \mathbf{W}_{k+1} &= \left[\mathbf{W}_k - \alpha \nabla_{\mathbf{W}} (f + \lambda g)(\mathbf{W}_k, \mathbf{H}_k) \right]_+ \\ \mathbf{H}_{k+1} &= \left[\mathbf{H}_k - \beta \nabla_{\mathbf{H}} (f + \lambda g)(\mathbf{W}_{k+1}, \mathbf{H}_k) \right]_+ \end{aligned}$$

- ▶ α, β denote stepsize
 - ▶ $\nabla_{\mathbf{W}}$ denotes the matrix-wise gradient with respect to \mathbf{W} . For compactness we will just write ∇
 - ▶ $[x]_+$ denotes the projection of x onto the nonnegative orthant. This operation makes sure the variable is feasible (satisfying the $\geq \mathbf{0}$ constraint)
- ▶ The terms $\nabla f(\mathbf{W}, \mathbf{H})$ are already derived (see [here](#) for the expansion of f) and thus

$$\begin{aligned} \mathbf{W}_{k+1} &= \left[\mathbf{W}_k - \alpha \left(\mathbf{W}_k \mathbf{H}_k \mathbf{H}_k^\top - \mathbf{X} \mathbf{H}_k^\top + \lambda \nabla g(\mathbf{W}_k, \mathbf{H}_k) \right) \right]_+ \\ \mathbf{H}_{k+1} &= \left[\mathbf{H}_k - \beta \left(\mathbf{W}_k^\top \mathbf{W}_k \mathbf{H}_k - \mathbf{W}_k^\top \mathbf{X} + \lambda \nabla g(\mathbf{W}_{k+1}, \mathbf{H}_k) \right) \right]_+ \end{aligned}$$

Orthogonal NMF

- ▶ NMF with orthogonality constraints

$$\operatorname{argmin}_{\mathbf{W} \geq 0, \mathbf{H} \geq 0} \frac{1}{2} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 \quad \text{s.t.} \quad \mathbf{W}^\top \mathbf{W} = \mathbf{I} \text{ and } \mathbf{H}^\top \mathbf{H} = \mathbf{I}.$$

- ▶ To enforce orthogonality, we can consider adding the penalty of orthogonality into the optimization formulation

$$(\mathcal{P}_{\text{ONMF}}) : \operatorname{argmin}_{\mathbf{W} \geq 0, \mathbf{H} \geq 0} \frac{1}{2} \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 + \frac{\lambda}{4} \|\mathbf{W}^\top \mathbf{W} - \mathbf{I}\|_F^2 + \frac{\mu}{4} \|\mathbf{H}^\top \mathbf{H} - \mathbf{I}\|_F^2.$$

- ▶ Here $g(\mathbf{W}, \mathbf{H})$ can be split into two independent functions $\|\mathbf{W}^\top \mathbf{W} - \mathbf{I}\|_F^2$ and $\|\mathbf{H}^\top \mathbf{H} - \mathbf{I}\|_F^2$.
- ▶ The meaning is that, we are minimizing the whole function, hence we want the whole expression to have small value. Then if $\|\mathbf{W}^\top \mathbf{W} - \mathbf{I}\|_F^2$ is small, we have $\mathbf{W}^\top \mathbf{W} \approx \mathbf{I}$, which is what we want.
- ▶ Now we need to know $\nabla_{\mathbf{W}} \|\mathbf{W}^\top \mathbf{W} - \mathbf{I}\|_F^2$ and $\nabla_{\mathbf{H}} \|\mathbf{H}^\top \mathbf{H} - \mathbf{I}\|_F^2$.
- ▶ You will soon see why we divided-by-4

Evaluating gradient of $\|\mathbf{W}^\top \mathbf{W} - \mathbf{I}\|_F^2$

$$\begin{aligned}\|\mathbf{W}^\top \mathbf{W} - \mathbf{I}\|_F^2 &= \text{Tr}\left((\mathbf{W}^\top \mathbf{W} - \mathbf{I})^\top (\mathbf{W}^\top \mathbf{W} - \mathbf{I})\right) \\ &= \text{Tr}\left((\mathbf{W}^\top \mathbf{W})^\top (\mathbf{W}^\top \mathbf{W}) - 2\mathbf{W}^\top \mathbf{W} + \mathbf{I}\right) \\ \nabla \|\mathbf{W}^\top \mathbf{W} - \mathbf{I}\|_F^2 &= \nabla \text{Tr}\left((\mathbf{W}^\top \mathbf{W})^\top (\mathbf{W}^\top \mathbf{W})\right) - 2\nabla \text{Tr}\left(\mathbf{W}^\top \mathbf{W}\right) \\ \nabla \text{Tr}\left(\mathbf{W}^\top \mathbf{W}\right) &= 2\mathbf{W} \\ \nabla \text{Tr}\left((\mathbf{W}^\top \mathbf{W})^\top (\mathbf{W}^\top \mathbf{W})\right) &= 4\mathbf{W}(\mathbf{W}^\top \mathbf{W}) \\ \nabla \|\mathbf{W}^\top \mathbf{W} - \mathbf{I}\|_F^2 &= 4\mathbf{W}(\mathbf{W}^\top \mathbf{W}) - 2(2\mathbf{W}) \\ &= 4\mathbf{W}(\mathbf{W}^\top \mathbf{W}) - 4\mathbf{W}\mathbf{I}_r \\ &= 4\mathbf{W}(\mathbf{W}^\top \mathbf{W} - \mathbf{I}_r)\end{aligned}$$

How to find $\nabla \text{Tr}\left((\mathbf{W}^\top \mathbf{W})^\top (\mathbf{W}^\top \mathbf{W})\right)$?

- ▶ Chain rule: [see details here](#)
- ▶ By Matrix Cookbook formula (123)

$$\frac{\partial \text{Tr}\left\{\mathbf{B}^\top \mathbf{X}^\top \mathbf{C} \mathbf{X} \mathbf{X}^\top \mathbf{C} \mathbf{X} \mathbf{B}\right\}}{\partial \mathbf{X}} = \mathbf{C}^\top \mathbf{X} \mathbf{X}^\top \mathbf{C} \mathbf{X} \mathbf{B} \mathbf{B}^\top + \mathbf{C}^\top \mathbf{X} \mathbf{B} \mathbf{B}^\top \mathbf{X}^\top \mathbf{C}^\top \mathbf{X} + \mathbf{C} \mathbf{X} \mathbf{B} \mathbf{B}^\top \mathbf{X}^\top \mathbf{C} \mathbf{X} + \mathbf{C}^\top \mathbf{X} \mathbf{X}^\top \mathbf{C}^\top \mathbf{X} \mathbf{B} \mathbf{B}^\top$$

Note that here \mathbf{W} can be any matrix, so the same conclusion holds for $\nabla \|\mathbf{H}^\top \mathbf{H} - \mathbf{I}\|_F^2 = 4\mathbf{H}(\mathbf{H}^\top \mathbf{H} - \mathbf{I}_n)$.

Now we have...

$$\nabla\|\mathbf{W}^\top\mathbf{W} - \mathbf{I}_r\|_F^2 = 4\mathbf{W}(\mathbf{W}^\top\mathbf{W} - \mathbf{I}), \quad \nabla\|\mathbf{H}^\top\mathbf{H} - \mathbf{I}_n\|_F^2 = 4\mathbf{H}(\mathbf{H}^\top\mathbf{H} - \mathbf{I}).$$

► Then

$$\frac{\lambda}{4}\nabla\|\mathbf{W}^\top\mathbf{W} - \mathbf{I}_r\|_F^2 = \lambda\mathbf{W}(\mathbf{W}^\top\mathbf{W} - \mathbf{I}_r), \quad \frac{\mu}{4}\nabla\|\mathbf{H}^\top\mathbf{H} - \mathbf{I}_n\|_F^2 = \mu\mathbf{H}(\mathbf{H}^\top\mathbf{H} - \mathbf{I}_n)$$

► Thus

$$\mathbf{W}_{k+1} = \left[\mathbf{W}_k - \alpha \left(\mathbf{W}_k \mathbf{H}_k \mathbf{H}_k^\top - \mathbf{X} \mathbf{H}_k^\top + \lambda \mathbf{W}_k (\mathbf{W}_k^\top \mathbf{W}_k - \mathbf{I}_r) \right) \right]_+$$

$$\mathbf{H}_{k+1} = \left[\mathbf{H}_k - \beta \left(\mathbf{W}_k^\top \mathbf{W}_k \mathbf{H}_k - \mathbf{W}_k^\top \mathbf{X} + \mu \mathbf{H}_k (\mathbf{H}_k^\top \mathbf{H}_k - \mathbf{I}_n) \right) \right]_+$$

► Remarks

- Do not set α, β too large
- Imagine in the extreme case that $\alpha = \infty$. You can see that in case the gradient bracket is positive, then everything inside the $[\cdot]_+$ is a strictly negative number, and then $[\cdot]_+$ will set everything to zero
- How to tune α, β : there are many ways, for simplicity we consider bisection
- Is there a theoretical bound for the value of α, β : it seems no. There is no a global constant Lipschitz constant for $\nabla\|\mathbf{W}^\top\mathbf{W} - \mathbf{I}_r\|_F^2$

Quick-and-dirty code

```
for k = 1 : itermax
%% Update of W
alpha = 1e-3; % initial stepsize

F_pre = computeF(X,W,H,Ir,In,lambda,mu);
HtH = H*H';
XtH = X*H';

% Block coordinate descent on W
for jj = 1 : 25
    WtW = W'*W;
    gradW = W*HtH - XtH + lambda*W*(WtW - Ir);

    Wtemp = max(0, W - alpha*gradW);

% bisection line search on stepsize
F_cur = computeF(X,Wtemp,H,Ir,In,lambda,mu);
if F_cur > F_pre
    alpha = alpha/2;
    Wtemp = max(0, W - alpha*gradW);
    F_cur = computeF(X,Wtemp,H,Ir,In,lambda,mu);
end
    W = Wtemp;
    F(1,k) = computeF(X,W,H,Ir,In,lambda,mu);
end
```

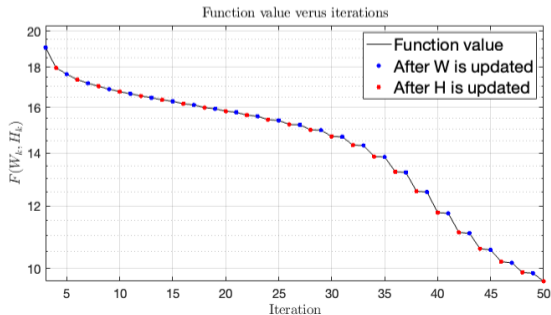
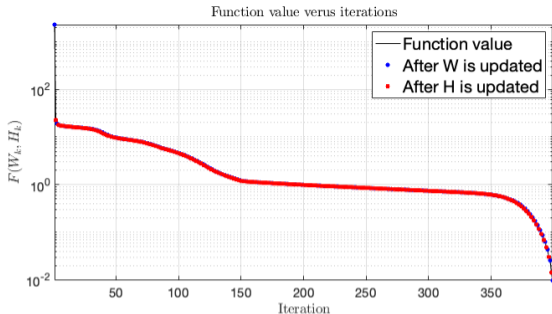
```
%% Update of H
alpha = 1e-3; % initial stepsize

F_pre = computeF(X,W,H,Ir,In,lambda,mu);
WtW = W'*W;
WtX = W'*X;

% Block coordinate descent on H
for jj = 1 : 25
    HtH = H*H';
    gradH = WtW * H - WtX + mu*H*(HtH - In);
    Htemp = max(H - alpha*gradH, 0);

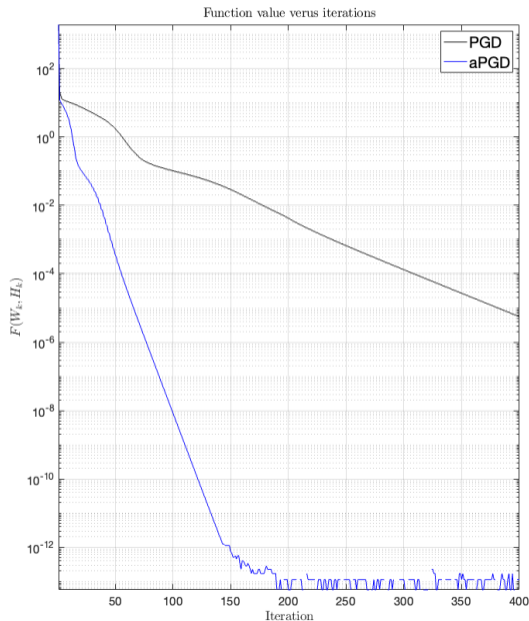
% bisection line search on stepsize
F_cur = computeF(X,W,Htemp,Ir,In,lambda,mu);
if F_cur > F_pre
    alpha = alpha/2;
    Htemp = max(H - alpha*gradH, 0);
    F_cur = computeF(X,W,Htemp,Ir,In,lambda,mu);
end
    H = Htemp;
    F(2,k) = computeF(X,W,H,Ir,In,lambda,mu);
end
```

end



PGD with Nesterov's acceleration

- ▶ Why PGD: we can easily add Nesterov's acceleration in the PGD update within the BCD framework.
- ▶ On using Nesterov's acceleration solving NNLS
- ▶ What about multiplicative update?
Slow, not worth mentioning



Last page - summary

Regularized NMF

Inexact Block-Coordinate Descent

Orthogonal NMF

$$\nabla \|\mathbf{W}^\top \mathbf{W} - \mathbf{I}_r\|_F^2 = 4\mathbf{W}(\mathbf{W}^\top \mathbf{W} - \mathbf{I})$$

Quick-and-dirty code

End of document