

Coordinate Descent Method

10:09 AM

Coordinate Descent Method:

$$\begin{pmatrix} \forall f_0(x) \\ x \\ \vdots \\ \forall_i \in \{1, \dots, v\} \\ x_i \in \mathcal{X}_i \end{pmatrix} \quad \text{eq: 12.90} \quad \begin{matrix} \# x = (x_1, \dots, x_v) \# \\ \downarrow \\ \text{Blocks of vectors} \end{matrix}$$

Coordinate descent method iteratively minimizes with respect to one block, while fixing the other.

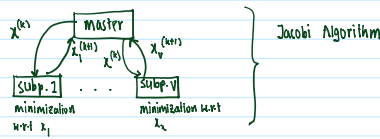
$x^{(k)} = (x_1^{(k)}, \dots, x_v^{(k)})$ [Value of the decision variable at iteration k]

$\forall x_i \in \mathcal{X}_i$ $f_0(x_1^{(k)}, \dots, x_{i-1}^{(k)}, x_i, x_{i+1}^{(k)}, \dots, x_v^{(k)})$ is solved for $i \in \{1, \dots, v\}$
 fixed variable fixed [eq: Cord Dscnt Mnmzn for x_i at itrtn k]

*The Jacobi method:

[eq: Cord Dscnt Mnmzn for x_i at itrtn k] is solved, then all blocks are updated simultaneously.

$$\forall_i \in \{1, \dots, v\} \quad x_i^{(k+1)} = \underset{x_i \in \mathcal{X}_i}{\operatorname{argmin}} f_0(x_1^{(k)}, \dots, x_{i-1}^{(k)}, x_i, x_{i+1}^{(k)}, \dots, x_v^{(k)})$$



Convergence of Coordinate Descent is not guaranteed in general

*Converges to the optimal solution under certain hypotheses of contractivity.

$$f: \text{function } \mathcal{Y} \rightarrow \mathcal{Z} \quad \left\{ \begin{array}{l} \text{Contraction Mapping } \exists \rho \in [0, 1) \forall y, z \in \mathcal{Y} \quad \|f(y) - f(z)\| \leq \rho \|y - z\| \\ \mathcal{Y} \text{ real vector space } \# \text{ norm defined } \|\cdot\| \end{array} \right. \quad \left\{ \begin{array}{l} \text{modulus of contraction} \end{array} \right.$$

*Convergence theorem for Jacobi method

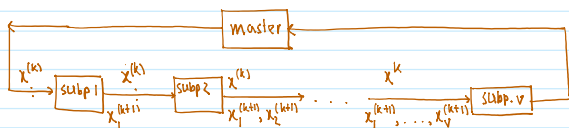
$$\left(f_0 \in \mathcal{C}^1, F(\cdot) = \arg \min_{\mathcal{X}} f_0(\cdot), F(\cdot) \text{ contraction mapping with norm } \|\cdot\| = \max_{i=1, \dots, v} \frac{\|H_i\|_2}{|H_i|} \right) \Rightarrow \text{(eq: 12.90) has a unique optimum } \wedge \quad \lim_{k \rightarrow \infty} x^{(k)} = x^* \quad \left\{ \begin{array}{l} \text{geometric rate} \\ \text{convergence rate } \rho \log \text{ is geometric} \end{array} \right.$$

*Block coordinate minimization method: (BCM)

(Also known as Gauss-Seidel Method)

At each iteration blocks are updated sequentially:

$$\forall_i \in \{1, \dots, v\} \quad x_i^{(k+1)} = \underset{x_i \in \mathcal{X}_i}{\operatorname{argmin}} f_0(x_1^{(k+1)}, x_2^{(k+1)}, \dots, x_{i-1}^{(k+1)}, x_i, x_{i+1}^{(k)}, \dots, x_v^{(k)})$$



*Convergence Condition for Jacobi iteration will work for BCM too

$$\left(\left[\nabla^2 f_0(\cdot) \right] \text{ contractive under norm } \|\cdot\| = \max_{i=1, \dots, v} \|H_i\|_2 / |H_i| \right) \quad \left\{ \begin{array}{l} \exists \rho < 1 \end{array} \right.$$

[Convergence Criteria for Block Coordinate Descent Method]

*BCM will converge for following condition too:

$$\left(\left\{ \begin{array}{l} f_0 \in \mathcal{C}^1, \mathcal{D}_x \text{ smooth, } \mathcal{D}_i \text{ smooth } \# \text{ with other blocks constant} \end{array} \right\}, \left\{ \mathcal{X}_i \right\} \text{ generated by BCM } \# \text{ self-defined} \right) \Rightarrow \lim_{k \rightarrow \infty} x^{(k)} = x^*$$

*BCM in generally may fail to converge for non-smooth objectives, even under convexity.

An important exception for which BCM converges:

$$f_0(x) = \phi(x) + \sum_{i=1}^v \psi_i(x_i) \quad \left\{ \begin{array}{l} \text{each } \psi_i(x_i) \# \mathcal{D}, \text{ maybe nonsmooth} \\ \# \text{ independent convex constraint } x_i \in \mathcal{X}_i \text{ can be differentiable in } \mathcal{X}_i \mathcal{D} \\ \# \text{ written as } \psi_i(x_i) = \psi_i(x_i), \text{ so } \nabla f_0(x) \text{ can be covered in this setup} \end{array} \right. \quad \left\{ \begin{array}{l} \# \text{ independent convex constraint } x_i \in \mathcal{X}_i \text{ can be written as } \psi_i(x_i) = \psi_i(x_i), \text{ so } \nabla f_0(x) \text{ can be covered in this setup} \end{array} \right.$$

in x, D # be covered in this setup x, c, x

* Theorem 12-5:

$$(x^0) \in C, \text{ initial point for BCM } \downarrow \text{Sublevelset at } f(x^0) \downarrow \text{Sublevelset at } f(x^0) \downarrow \text{Sublevelset at } f(x^0)$$

$$S_0 = \{x: f_0(x) \leq f_0(x^0)\} \in [S]^*$$

$$f_0(b) = \theta(x) + \sum_{i=1}^m \lambda_i(x_i) \Rightarrow \lim_{k \rightarrow \infty} x^{(k)} = x^*$$

Theorem 12-5 applies to Lasso problem or other l_1 norm regularized problem, thus guaranteeing convergence of BCM in those cases.

* Power iteration and block coordinate descent method:

- Power iteration (class of methods applied to specific eigenvalue, singular value problems)
- (coordinate descent)

Problem: Finding rank-one approximation of a matrix A:

$$\min_{x,y} \|A - xy^T\|_F^2 = \min_{x,y} \sum_{i=1}^n \sigma_i^2$$

This is convex in x because, $\|x\|_2^2 \leq \|x+\alpha\|_2^2$ however if you consider x, y separately then it is bilinear and not convex in (x, y)

A stationary point (where the derivative is 0) finding:

$$\|A - xy^T\|_F^2 = \text{tr}((A - xy^T)(A - xy^T)^T) \quad [\because \|x\|_F = \sqrt{\text{tr}(AA^T)}]$$

$$\begin{aligned} &= (A - xy^T)(A - xy^T)^T \\ &= (A - xy^T)(A^T - yx^T) \\ &= AA^T - Ayx^T - xy^T A^T + xy^T yx^T \\ &= AA^T - Ayx^T - xy^T A^T + \|y\|_2^2 xx^T \end{aligned}$$

[eq: expression of $\|A - xy^T\|_F^2$ in terms trace]

$$\begin{aligned} &= \text{tr}(AA^T - Ayx^T - xy^T A^T + \|y\|_2^2 xx^T) \\ &= \text{tr}(AA^T) - \text{tr}(Ayx^T) - \text{tr}(xy^T A^T) + \|y\|_2^2 \text{tr}(xx^T) \quad [\because \text{tr is a linear operator}] \end{aligned}$$

$$\begin{aligned} \nabla_x \|A - xy^T\|_F^2 &= \nabla_x (\text{tr}(AA^T) - \text{tr}(Ayx^T) - \text{tr}(xy^T A^T) + \|y\|_2^2 \text{tr}(xx^T)) \\ &= \frac{\partial}{\partial x} \text{tr}(AA^T) - \frac{\partial}{\partial x} \text{tr}(Ayx^T) - \frac{\partial}{\partial x} \text{tr}(xy^T A^T) + \|y\|_2^2 \frac{\partial}{\partial x} \text{tr}(xx^T) \\ &= 0 - Ay - (y^T A^T)^T + \|y\|_2^2 2x \\ &= -2Ay + \|y\|_2^2 2x \end{aligned}$$

mnemonic: transpose left, transpose right

for finding a stationary point in x :

$$\begin{aligned} \nabla_x \|A - xy^T\|_F^2 &= -2Ay + \|y\|_2^2 2x = 0 \\ \Leftrightarrow 2\|y\|_2^2 x &= 2Ay \\ \Leftrightarrow \|x\|_2^2 x &= Ay \quad [\text{eq: stationary point 1}] \end{aligned}$$

Similarly, $\nabla_y \|A - xy^T\|_F^2 = \nabla_y (\text{tr}(AA^T) - \text{tr}(Ayx^T) - \text{tr}(xy^T A^T) + \|x\|_2^2 \text{tr}(yy^T))$

$$\begin{aligned} &= -\nabla_y (\text{tr}(Ayx^T) - \text{tr}(xy^T A^T) + \|x\|_2^2 \text{tr}(yy^T)) \\ &= -\nabla_y (\text{tr}(x^T Ay) - \text{tr}(y^T A^T x) + \|x\|_2^2 \text{tr}(yy^T)) \\ &= -\nabla_y (\text{tr}(x^T Ay) - \text{tr}(x^T Ay) + \|x\|_2^2 2y) \\ &= -\nabla_y (\|x\|_2^2 2y) \\ &= -\|x\|_2^2 2y \end{aligned}$$

Finding a stationary point in y :

$$\begin{aligned} \nabla_y \|A - xy^T\|_F^2 &= -\|x\|_2^2 2y = 0 \\ \Leftrightarrow \|x\|_2^2 y &= A^T x \quad [\text{eq: stationary point 2}] \end{aligned}$$

lets normalize the vectors as follows: $u = \frac{x}{\|x\|_2}, v = \frac{y}{\|y\|_2}, d = \|x\|_2 \|y\|_2$

then $xy^T = \|x\|_2 \|y\|_2 uv^T = d uv^T$

$$\min_{x,y} \|A - xy^T\|_F^2 = \min_{u,v,d} \|A - d uv^T\|_F^2$$

[problem: dyadic version of rank 1 approximation problem]

$$\text{then: } \left(\begin{array}{l} \min_{x,y} \|A - xy^T\|_F^2 \\ \|u\|_2 = 1 \\ \|v\|_2 = 1 \\ d \geq 0 \end{array} \right)$$

[problem: dyadic version of rank 1 approximation problem]

$$\|x\|_2 \|y\|_2 = (\|x\|_2 \|y\|_2)^2 = d^2$$

$$\text{tr}(x^T x) = \text{tr}(\|x\|_2^2) = \|x\|_2^2$$

now note that in d this is convex

$$\|x\|_2 \|y\|_2 = d$$

[eq: expression of $\|A - xy^T\|_F^2$]

$$\begin{aligned} \|A - xy^T\|_F^2 &= \text{tr}(AA^T) - \text{tr}(Axy^T) - \text{tr}(x y^T A^T) + \|y\|_2^2 \text{tr}(xx^T) \\ &= \text{tr}(AA^T) - \text{tr}(A \|y\|_2 v \|x\|_2 u^T) - \text{tr}(\|x\|_2 u \|y\|_2 v^T A^T) + d^2 \\ &= \text{tr}(AA^T) - \text{tr}(A d v u^T) - \text{tr}(d u v^T A^T) + d^2 \\ &= \text{tr}(AA^T) - d \left(\text{tr}(A d v u^T) + \text{tr}(u v^T A^T) \right) + d^2 \\ &\quad \text{tr}(A v u^T + u v^T A^T) \quad \text{tr}(\cdot) \text{ is linear operator} \end{aligned}$$

$$\nabla_d \|A - xy^T\|_F^2 = 0 - \text{tr}(A v u^T + u v^T A^T) + 2d = 0$$

$\therefore d = \frac{1}{2} \text{tr}(A v u^T + u v^T A^T) = \frac{1}{2} 2 u^T A v = u^T A v$, plugging this value ([problem: dyadic version of rank 1 approximation problem]) we get

$$\left(\begin{array}{l} \min_{u,v} \|A - d u v^T\|_F^2 \\ \|u\|_2 = 1, \|v\|_2 = 1 \end{array} \right)$$

$$\begin{aligned} & \# \text{tr}(A v u^T) + \text{tr}(u v^T A^T) \\ &= \text{tr}(u^T A v) + \text{tr}(v^T A^T u) \\ & \# (u^T A v)^T = (A v)^T (u^T)^T = v^T A^T u \\ &= u^T A v + (u^T A v)^T = u^T A v + v^T A^T u = 2 u^T A v \\ & \# \text{number so trace} \\ & \# \text{trace of a matrix} \\ & \text{this is a number so transpose will be itself} \end{aligned}$$