Optimization problems containing optimal transport costs: examples and computational methods

Axel Ringh¹

¹Department of Electronic and Computer Engineering, The Hong Kong University of Science and Technology.

10th of December 2019

Workshop: Computational optimal transport for applications in control and estimation 58th Conference on Decision and Control (CDC), Nice, France



- Optimal transport recap
- Optimization problems with an optimal transport cost and generalized Sinkhorn iterations
- Variable splitting and the proximal operator of entropy-regularized optimal transport
- Example: variational regularization of inverse problems

Optimal transport recap

Optimal transport distance between two functions $f_0(x)$ and $f_1(x)$ is defined as

$$T(f_0, f_1) := \begin{cases} \min_{M \ge 0} & \int_{X \times X} c(x^{(0)}, x^{(1)}) M(x^{(0)}, x^{(1)}) dx^{(0)} dx^{(1)} \\ \text{s.t.} & f_0(x^{(0)}) = \int_X M(x^{(0)}, x^{(1)}) dx^{(1)}, \ x^{(0)} \in X \\ & f_1(x^{(1)}) = \int_X M(x^{(0)}, x^{(1)}) dx^{(0)}, \ x^{(1)} \in X. \end{cases}$$

for some cost function $c: X \times X \to \mathbb{R}_+$.

Optimal transport recap

Optimal transport distance between two functions $f_0(x)$ and $f_1(x)$ is defined as

$$T(f_0, f_1) := \begin{cases} \min_{M \ge 0} & \int_{X \times X} c(x^{(0)}, x^{(1)}) M(x^{(0)}, x^{(1)}) dx^{(0)} dx^{(1)} \\ \text{s.t.} & f_0(x^{(0)}) = \int_X M(x^{(0)}, x^{(1)}) dx^{(1)}, \ x^{(0)} \in X \\ & f_1(x^{(1)}) = \int_X M(x^{(0)}, x^{(1)}) dx^{(0)}, \ x^{(1)} \in X. \end{cases}$$

for some cost function $c: X \times X \to \mathbb{R}_+$.

Discretized version:

• vectors
$$f_0 \in \mathbb{R}^n$$
, $f_1 \in \mathbb{R}^n$

• cost matrix $C = [c_{ij}] \in \mathbb{R}^{n \times n}$, where c_{ij} is the transportation cost $c(x_i, x_j)$

• transportation plan
$$M = [m_{ij}] \in \mathbb{R}^{n \times n}$$
.

Optimal transport recap

Optimal transport distance between two functions $f_0(x)$ and $f_1(x)$ is defined as

$$T(f_0, f_1) := \begin{cases} \min_{M \ge 0} & \int_{X \times X} c(x^{(0)}, x^{(1)}) M(x^{(0)}, x^{(1)}) dx^{(0)} dx^{(1)} \\ \text{s.t.} & f_0(x^{(0)}) = \int_X M(x^{(0)}, x^{(1)}) dx^{(1)}, \ x^{(0)} \in X \\ & f_1(x^{(1)}) = \int_X M(x^{(0)}, x^{(1)}) dx^{(0)}, \ x^{(1)} \in X. \end{cases}$$

for some cost function $c: X \times X \to \mathbb{R}_+$.

Discretized version:

• vectors
$$f_0 \in \mathbb{R}^n$$
, $f_1 \in \mathbb{R}^n$

• cost matrix $C = [c_{ij}] \in \mathbb{R}^{n \times n}$, where c_{ij} is the transportation cost $c(x_i, x_j)$

• transportation plan
$$M = [m_{ij}] \in \mathbb{R}^{n \times n}$$
.

$$T(f_0, f_1) := \begin{cases} \min_{m_{ij} \ge 0} & \sum_{i=1}^n \sum_{j=1}^n c_{ij} m_{ij} \\ \text{s.t.} & \sum_{j=1}^n m_{ij} = f_0(i), \ i = 1, \dots, n, \\ & \sum_{i=1}^n m_{ij} = f_1(j), \ j = 1, \dots, n. \end{cases} \iff T(f_0, f_1) := \begin{cases} \min_{M \ge 0} & \text{trace}(C^T M) \\ \text{s.t.} & M \mathbf{1} = f_0 \\ & M^T \mathbf{1} = f_1 \end{cases}$$

$$T_{\epsilon}(f_0, f_1) := \min_{\substack{M \ge 0}} \operatorname{trace}(C^T M) + \epsilon D(M)$$

subject to $f_0 = M \mathbf{1}$
 $f_1 = M^T \mathbf{1}.$

M. Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In Advances in Neural Information Processing Systems, pages 2292–2300, 2013.

$$T_{\epsilon}(f_0, f_1) := \min_{\substack{M \ge 0 \\ \text{subject to}}} \operatorname{trace}(C^T M) + \epsilon D(M)$$
$$f_0 = M \mathbf{1}$$
$$f_1 = M^T \mathbf{1}.$$

• Let $exp(\cdot)$, $log(\cdot)$, ./, \odot denotes the element-wise function.

M. Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In Advances in Neural Information Processing Systems, pages 2292–2300, 2013.

$$T_{\epsilon}(f_0, f_1) := \min_{\substack{M \ge 0 \\ \text{subject to}}} \operatorname{trace}(C^T M) + \epsilon D(M)$$
$$f_0 = M \mathbf{1}$$
$$f_1 = M^T \mathbf{1}.$$

Let exp(·), log(·), ./, ⊙ denotes the element-wise function.
For K = exp(-C/ε), the solution is of the form

 $M = \operatorname{diag}(u_0) K \operatorname{diag}(u_1)$

M. Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In Advances in Neural Information Processing Systems, pages 2292–2300, 2013.

$$T_{\epsilon}(f_0, f_1) := \min_{\substack{M \ge 0 \\ \text{subject to}}} \operatorname{trace}(C^T M) + \epsilon D(M)$$
$$f_0 = M \mathbf{1}$$
$$f_1 = M^T \mathbf{1}.$$

• Let $exp(\cdot),$ log(·), ./, \odot denotes the element-wise function.

• For $K = \exp(-C/\epsilon)$, the solution is of the form

 $M = \operatorname{diag}(u_0) K \operatorname{diag}(u_1)$

Theorem (Sinkhorn iterations [2])

For any matrix K with positive elements there are diagonal matrices $diag(u_0)$, $diag(u_1)$ such that $M = diag(u_0)Kdiag(u_1)$ has prescribed row- and column-sums f_0 and f_1 . The vectors u_0 and u_1 can be obtained by alternating marginalization: $u_0 = f_0./(Ku_1)$

$$u_1 = f_1./(K^T u_0)$$

M. Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In Advances in Neural Information Processing Systems, pages 2292–2300, 2013.

 ^[2] R. Sinkhorn. Diagonal equivalence to matrices with prescribed row and column sums. The American Mathematical Monthly, 74(4), 402–405, 1967.

 $\min_{f_1} T_{\epsilon}(f_0, f_1) + \mathcal{G}(f_1)$

where \mathcal{G} proper, convex and lower semicontinuous.

$$\min_{f_1} T_{\epsilon}(f_0, f_1) + \mathcal{G}(f_1) = \min_{\substack{M \ge 0, f_1 \\ \text{subject to}}} \operatorname{trace}(C^T M) + \epsilon D(M) + \mathcal{G}(f_1)$$
$$f_0 = M \mathbf{1}$$
$$f_1 = M^T \mathbf{1}.$$

where \mathcal{G} proper, convex and lower semicontinuous.

$$\min_{f_1} T_{\epsilon}(f_0, f_1) + \mathcal{G}(f_1) = \min_{\substack{M \ge 0, f_1 \\ \text{subject to}}} \operatorname{trace}(C^T M) + \epsilon D(M) + \mathcal{G}(f_1)$$
$$f_0 = M\mathbf{1}$$
$$f_1 = M^T \mathbf{1}.$$

where \mathcal{G} proper, convex and lower semicontinuous.

Problem well-posed if there exists $f_1 \ge 0$ such that

•
$$\mathbf{1}^{T} f_{1} = \mathbf{1}^{T} f_{0}$$
,

•
$$\mathcal{G}(f_1) < \infty$$
.

In this case, the problem is convex and there exists an optimal solution.

$$\min_{f_1} T_{\epsilon}(f_0, f_1) + \mathcal{G}(f_1) = \min_{\substack{M \ge 0, f_1 \\ \text{subject to}}} \operatorname{trace}(C^T M) + \epsilon D(M) + \mathcal{G}(f_1)$$
$$f_0 = M\mathbf{1}$$
$$f_1 = M^T \mathbf{1}.$$

where \mathcal{G} proper, convex and lower semicontinuous.

Problem well-posed if there exists $f_1 \ge 0$ such that

- $\mathbf{1}^T f_1 = \mathbf{1}^T f_0$,
- $\mathcal{G}(f_1) < \infty$.

In this case, the problem is convex and there exists an optimal solution.

How to solve this problem?

$$T_{\epsilon}(f_0, f_1) := \min_{\substack{M \ge 0 \\ \text{subject to}}} \operatorname{trace}(C^T M) + \epsilon D(M)$$
$$f_0 = M \mathbf{1}$$
$$f_1 = M^T \mathbf{1}.$$

$$T_{\epsilon}(f_0, f_1) := \min_{\substack{M \ge 0 \\ \text{subject to}}} \operatorname{trace}(C^T M) + \epsilon D(M)$$

subject to $f_0 = M \mathbf{1}$
 $f_1 = M^T \mathbf{1}.$

$$L(M, \lambda_0, \lambda_1) = \operatorname{trace}(C^T M) + \epsilon D(M) + \lambda_0^T (f_0 - M \mathbf{1}) + \lambda_1^T (f_1 - M^T \mathbf{1}).$$

$$T_{\epsilon}(f_0, f_1) := \min_{\substack{M \ge 0}} \operatorname{trace}(C^T M) + \epsilon D(M)$$

subject to $f_0 = M \mathbf{1}$
 $f_1 = M^T \mathbf{1}.$

• Using Lagrangian relaxation gives
$$T_1 = M$$

$$L(M,\lambda_0,\lambda_1) = \operatorname{trace}(C^{\mathsf{T}}M) + \epsilon D(M) + \lambda_0^{\mathsf{T}}(f_0 - M\mathbf{1}) + \lambda_1^{\mathsf{T}}(f_1 - M^{\mathsf{T}}\mathbf{1}).$$

• Given dual variables λ_0, λ_1 , the minimum m_{ij} is

$$0 = rac{\partial L(M,\lambda_0,\lambda_1)}{\partial m_{ij}} = c_{ij} + \epsilon \log(m_{ij}) - \lambda_0(i) - \lambda_1(j)$$

$$T_{\epsilon}(f_0, f_1) := \min_{\substack{M \ge 0 \\ \text{subject to}}} \operatorname{trace}(C^T M) + \epsilon D(M)$$

subject to $f_0 = M \mathbf{1}$
 $f_1 = M^T \mathbf{1}.$

$$L(\boldsymbol{M}, \lambda_0, \lambda_1) = \operatorname{trace}(\boldsymbol{C}^T \boldsymbol{M}) + \epsilon D(\boldsymbol{M}) + \lambda_0^T (f_0 - \boldsymbol{M} \mathbf{1}) + \lambda_1^T (f_1 - \boldsymbol{M}^T \mathbf{1}).$$

• Given dual variables λ_0, λ_1 , the minimum m_{ij} is

$$0 = \frac{\partial L(M, \lambda_0, \lambda_1)}{\partial m_{ij}} = c_{ij} + \epsilon \log(m_{ij}) - \lambda_0(i) - \lambda_1(j)$$

• Solve for m_{ij} to get

$$m_{ij}=e^{\lambda_0(i)/\epsilon}e^{-c_{ij}/\epsilon}e^{\lambda_1(j)/\epsilon}.$$

$$T_{\epsilon}(f_0, f_1) := \min_{\substack{M \ge 0 \\ \text{subject to}}} \operatorname{trace}(C^T M) + \epsilon D(M)$$

subject to $f_0 = M \mathbf{1}$
 $f_1 = M^T \mathbf{1}.$

$$L(M, \lambda_0, \lambda_1) = \operatorname{trace}(C^{\mathsf{T}}M) + \epsilon D(M) + \lambda_0^{\mathsf{T}}(f_0 - M\mathbf{1}) + \lambda_1^{\mathsf{T}}(f_1 - M^{\mathsf{T}}\mathbf{1}).$$

• Given dual variables λ_0, λ_1 , the minimum m_{ij} is

$$0 = \frac{\partial L(M, \lambda_0, \lambda_1)}{\partial m_{ij}} = c_{ij} + \epsilon \log(m_{ij}) - \lambda_0(i) - \lambda_1(j)$$

• Solve for *m_{ij}* to get

$$m_{ij} = e^{\lambda_0(i)/\epsilon} e^{-c_{ij}/\epsilon} e^{\lambda_1(j)/\epsilon}$$

• Change of variables: $u_0 = \exp(\lambda_0/\epsilon)$, $u_1 = \exp(\lambda_1/\epsilon)$. The optimal solution is of the form $M^* = \operatorname{diag}(u_0)K\operatorname{diag}(u_1)$

where $K = \exp(-C/\epsilon)$.

• Lagrangian relaxation gave optimal form of the primal variable

 $M^* = \operatorname{diag}(u_0) K \operatorname{diag}(u_1)$

Optimization problems with an optimal transport cost

Sinkhorn iterations as dual coordinate ascent

One way to interpreted the Sinkhorn iterations: coordinate ascent in the Lagrangian dual.

• Lagrangian relaxation gave optimal form of the primal variable

 $M^* = \operatorname{diag}(u_0) K \operatorname{diag}(u_1)$

• The Lagrangian dual function:

 $\varphi(u_0, u_1) := \min_{M \ge 0} L(M, u_0, u_1) = L(M^*, u_0, u_1) = \ldots = \epsilon \log(u_0)^T f_0 + \epsilon \log(u_1)^T f_1 - \epsilon u_0^T K u_1 + \epsilon n^2.$

Optimization problems with an optimal transport cost Sinkhorn iterations as dual coordinate ascent

One way to interpreted the Sinkhorn iterations: coordinate ascent in the Lagrangian dual.

• Lagrangian relaxation gave optimal form of the primal variable

 $M^* = \operatorname{diag}(u_0) K \operatorname{diag}(u_1)$

• The Lagrangian dual function:

 $\varphi(u_0, u_1) := \min_{M \ge 0} L(M, u_0, u_1) = L(M^*, u_0, u_1) = \ldots = \epsilon \log(u_0)^T f_0 + \epsilon \log(u_1)^T f_1 - \epsilon u_0^T K u_1 + \epsilon n^2.$

• The dual problem is thus

 $\max_{u_0,u_1\geq 0} \varphi(u_0,u_1)$

• Lagrangian relaxation gave optimal form of the primal variable

$$M^* = \operatorname{diag}(u_0) K \operatorname{diag}(u_1)$$

• The Lagrangian dual function:

 $\varphi(u_0, u_1) := \min_{M \ge 0} L(M, u_0, u_1) = L(M^*, u_0, u_1) = \ldots = \epsilon \log(u_0)^T f_0 + \epsilon \log(u_1)^T f_1 - \epsilon u_0^T K u_1 + \epsilon n^2.$

• The dual problem is thus

$$\max_{u_0,u_1\geq 0} \varphi(u_0,u_1)$$

• Taking the gradient w.r.t u_0 and putting it equal to zero gives

$$\epsilon f_0./u_0 - \epsilon K u_1 = 0$$

• Lagrangian relaxation gave optimal form of the primal variable

$$M^* = \operatorname{diag}(u_0) K \operatorname{diag}(u_1)$$

• The Lagrangian dual function:

 $\varphi(u_0, u_1) := \min_{M \ge 0} L(M, u_0, u_1) = L(M^*, u_0, u_1) = \ldots = \epsilon \log(u_0)^T f_0 + \epsilon \log(u_1)^T f_1 - \epsilon u_0^T K u_1 + \epsilon n^2.$

• The dual problem is thus

$$\max_{u_0,u_1\geq 0} \varphi(u_0,u_1)$$

• Taking the gradient w.r.t u_0 and putting it equal to zero gives

$$\epsilon f_{0.}/u_0 - \epsilon K u_1 = 0 \qquad \rightsquigarrow \qquad u_0 = f_{0.}/(K u_1),$$

Lagrangian relaxation gave optimal form of the primal variable

$$M^* = \operatorname{diag}(u_0) K \operatorname{diag}(u_1)$$

• The Lagrangian dual function:

 $\varphi(u_0, u_1) := \min_{M \ge 0} L(M, u_0, u_1) = L(M^*, u_0, u_1) = \ldots = \epsilon \log(u_0)^T f_0 + \epsilon \log(u_1)^T f_1 - \epsilon u_0^T K u_1 + \epsilon n^2.$

• The dual problem is thus

$$\max_{u_0,u_1\geq 0} \varphi(u_0,u_1)$$

• Taking the gradient w.r.t u_0 and putting it equal to zero gives

$$\epsilon f_{0.}/u_0 - \epsilon K u_1 = 0 \qquad \rightsquigarrow \qquad u_0 = f_{0.}/(K u_1),$$

and w.r.t u_1 gives

$$\epsilon f_{1.}/u_{1} - \epsilon \left(u_{0}^{\mathsf{T}} \mathsf{K}\right)^{\mathsf{T}} = 0 \qquad \rightsquigarrow \qquad u_{1} = f_{1.}/(\mathsf{K}^{\mathsf{T}} u_{0})$$

These are the Sinkhorn iterations! (cf. [1])

 P. Tseng. Dual ascent methods for problems with strictly convex costs and linear constraints: A unified approach. SIAM Journal on Control and Optimization, 28(1), 214–242, 1990.

Optimization problems with an optimal transport cost Generalized Sinkhorn iterations

$$\min_{f_1} T_{\epsilon}(f_0, f_1) + \mathcal{G}(f_1) = \min_{\substack{M \ge 0, f_1 \\ \text{subject to}}} \operatorname{trace}(C^T M) + \epsilon D(M) + \mathcal{G}(f_1)$$
$$f_0 = M \mathbf{1}$$
$$f_1 = M^T \mathbf{1}.$$

Optimization problems with an optimal transport cost Generalized Sinkhorn iterations

$$\min_{f_1} T_{\epsilon}(f_0, f_1) + \mathcal{G}(f_1) = \min_{\substack{M \ge 0, f_1 \\ \text{subject to}}} \operatorname{trace}(C^T M) + \epsilon D(M) + \mathcal{G}(f_1)$$
$$f_0 = M\mathbf{1}$$
$$f_1 = M^T \mathbf{1}.$$

Lagrangian dual problem

$$\max_{u_0,u_1} \varphi(u_0,u_1) = \max_{u_0,u_1} \epsilon \log(u_0)^T f_0 - \mathcal{G}^*(-\epsilon \log(u_1)) - \epsilon u_0^T K u_1 + \epsilon n^2,$$

where $\mathcal{G}^*(u) := \sup_f u^T f - \mathcal{G}(f)$ is the Fenchel dual.

Optimization problems with an optimal transport cost Generalized Sinkhorn iterations

$$\min_{f_1} T_{\epsilon}(f_0, f_1) + \mathcal{G}(f_1) = \min_{\substack{M \ge 0, f_1 \\ \text{subject to}}} \operatorname{trace}(C^T M) + \epsilon D(M) + \mathcal{G}(f_1)$$
$$f_0 = M \mathbf{1}$$
$$f_1 = M^T \mathbf{1}.$$

Lagrangian dual problem

$$\max_{u_0,u_1} \varphi(u_0,u_1) = \max_{u_0,u_1} \epsilon \log(u_0)^T f_0 - \mathcal{G}^*(-\epsilon \log(u_1)) - \epsilon u_0^T K u_1 + \epsilon n^2,$$

where $\mathcal{G}^*(u) := \sup_f u^T f - \mathcal{G}(f)$ is the Fenchel dual.

Can be solve by dual coordinate ascent

$$egin{aligned} 0 &= f_0/u_0 - \mathcal{K} u_1 \ 0 &\in \partial \mathcal{G}^*(-\epsilon \log(u_1)) rac{1}{u_1} - \mathcal{K}^{ op} u_0, \end{aligned}$$

if the second inclusion can be solved efficiently.

Optimization problems with an optimal transport cost Generalized Sinkhorn iterations

$$\min_{f_1} T_{\epsilon}(f_0, f_1) + \mathcal{G}(f_1) = \min_{\substack{M \ge 0, f_1 \\ \text{subject to}}} \operatorname{trace}(C^T M) + \epsilon D(M) + \mathcal{G}(f_1)$$
$$f_0 = M \mathbf{1}$$
$$f_1 = M^T \mathbf{1}.$$

Lagrangian dual problem

$$\max_{u_0,u_1} \varphi(u_0,u_1) = \max_{u_0,u_1} \epsilon \log(u_0)^T f_0 - \mathcal{G}^*(-\epsilon \log(u_1)) - \epsilon u_0^T \mathcal{K} u_1 + \epsilon n^2,$$

where $\mathcal{G}^*(u) := \sup_f u^T f - \mathcal{G}(f)$ is the Fenchel dual.

Can be solve by dual coordinate ascent

$$egin{aligned} 0 &= f_0/u_0 - \mathcal{K} u_1 \ 0 &\in \partial \mathcal{G}^*(-\epsilon \log(u_1)) rac{1}{u_1} - \mathcal{K}^{ op} u_0, \end{aligned}$$

if the second inclusion can be solved efficiently.

The second inclusion can be efficiently solved when $\partial \mathcal{G}^*(\cdot)$ is component-wise. Example of such cases:

•
$$\mathcal{G}(\cdot) = \mathcal{I}_{\tilde{f}}(\cdot)$$
 indicator function on $\{\tilde{f}\} \longrightarrow \mathcal{G}^*(\cdot) = \cdot^T \tilde{f} \longrightarrow \partial \mathcal{G}^*(\cdot) = \tilde{f}$
• $\mathcal{G}(\cdot) = \|\cdot\|_2^2 \longrightarrow \mathcal{G}^*(\cdot) = \frac{1}{4}\|\cdot\|_2^2 \longrightarrow \partial \mathcal{G}^*(\cdot) = \frac{1}{2}$.

How to solve

 $\min_{f_1} T_\epsilon(f_0, f_1) + \mathcal{G}(f_1)$

for more general functions G?

Intermission: ADMM and variable splitting

Consider

 $\min_{y,z} \quad \mathcal{H}(y) + \mathcal{G}(z)$ subject to Ax + Bz = c

$$\begin{split} \min_{y,z} & \mathcal{H}(y) + \mathcal{G}(z) \\ \text{subject to} & Ax + Bz = c \end{split}$$
Can be solved by ADMM [1]: for $\rho > 0$ $y^{k+1} &= \arg\min_{y} \mathcal{H}(y) + \frac{2}{\rho} \|Ay + Bz^{k} - c + u^{k}\|_{2}^{2}$ $z^{k+1} &= \arg\min_{z} \mathcal{G}(z) + \frac{2}{\rho} \|Ay^{k+1} + Bz - c + u^{k}\|_{2}^{2}$ $u^{k+1} &= u^{k} + Av^{k+1} + Bz^{k+1} - c$

S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. Foundations and Trends(R) in Machine learning, 3(1), 1-122, 2011.

 $\min_{y,z} \quad \mathcal{H}(y) + \mathcal{G}(z)$ subject to Ax + Bz = cCan be solved by ADMM [1]: for $\rho > 0$ $y^{k+1} = \arg\min_{y} \mathcal{H}(y) + \frac{2}{\rho} ||Ay + Bz^{k} - c + u^{k}||_{2}^{2}$ $z^{k+1} = \arg\min_{z} \mathcal{G}(z) + \frac{2}{\rho} ||Ay^{k+1} + Bz - c + u^{k}||_{2}^{2}$ $u^{k+1} = u^{k} + Ay^{k+1} + Bz^{k+1} - c$

Special case:

 $\min_{y} \mathcal{H}(y) + \mathcal{G}(y)$

S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. Foundations and Trends(R) in Machine learning, 3(1), 1-122, 2011.

$$\begin{split} \min_{y,z} & \mathcal{H}(y) + \mathcal{G}(z) \\ \text{subject to} & Ax + Bz = c \end{split}$$
Can be solved by ADMM [1]: for $\rho > 0$ $y^{k+1} &= \arg\min_{y} \mathcal{H}(y) + \frac{2}{\rho} \|Ay + Bz^{k} - c + u^{k}\|_{2}^{2}$ $z^{k+1} &= \arg\min_{z} \mathcal{G}(z) + \frac{2}{\rho} \|Ay^{k+1} + Bz - c + u^{k}\|_{2}^{2}$ $u^{k+1} &= u^{k} + Ay^{k+1} + Bz^{k+1} - c \end{split}$

Special case:

$$\min_{y} \mathcal{H}(y) + \mathcal{G}(y) \quad \rightsquigarrow \quad \min_{y,z} \quad \mathcal{H}(y) + \mathcal{G}(z)$$

subject to $y - z = 0$

 S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. Foundations and Trends® in Machine learning, 3(1), 1-122, 2011.

 $\min_{y,z} \quad \mathcal{H}(y) + \mathcal{G}(z)$ subject to Ax + Bz = cCan be solved by ADMM [1]: for $\rho > 0$ $y^{k+1} = \arg \min \mathcal{H}(y) + \frac{2}{a} ||Ay + Bz^k - c + u^k||_2^2$ $z^{k+1} = \arg\min_{-} \mathcal{G}(z) + \frac{2}{a} ||Ay^{k+1} + Bz - c + u^{k}||_{2}^{2}$ $u^{k+1} = u^k + Av^{k+1} + Bz^{k+1} - c$ Special case: $\min_{\mathbf{y}} \mathcal{H}(\mathbf{y}) + \mathcal{G}(\mathbf{y}) \quad \rightsquigarrow$

 S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. Foundations and Trends® in Machine learning, 3(1), 1-122, 2011. 10/22

Variable splitting

ADMM is a special case of so-called variable splitting.

Common for large convex optimization problem with several terms.

Examples of methods

- ADMM [1]
- primal-dual hybrid gradient algorithm (Chambolle-Pock) [2]
- primal-dual Douglas-Rachford [3]

- S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. Foundations and Trends in Machine learning, 3(1), 1-122, 2011.
- [2] A. Chambolle, and T. Pock. A first-order primal-dual algorithm for convex problems with applications to imaging. Journal of Mathematical Imaging and Vision, 40(1), 120-145, 2011.
- R.I. Boţ, and C. Hendrich. A Douglas-Rachford type primal-dual method for solving inclusions with mixtures of composite and parallel-sum type monotone operators. SIAM Journal on Optimization, 23(4), 2541-2565, 2013.

Variable splitting

ADMM is a special case of so-called variable splitting.

Common for large convex optimization problem with several terms.

Examples of methods

- ADMM [1]
- primal-dual hybrid gradient algorithm (Chambolle-Pock) [2]
- primal-dual Douglas-Rachford [3]

Common tool in these algorithms: the proximal operator of the involved functions ${\cal H}$ and ${\cal G}$ [4, 5]

$$\operatorname{Prox}_{\mathcal{H}}^{\sigma}(h) = \arg\min_{f} \mathcal{H}(f) + \frac{1}{2\sigma} \|f - h\|_{2}^{2}.$$

- S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. Foundations and Trends in Machine learning, 3(1), 1-122, 2011.
- [2] A. Chambolle, and T. Pock. A first-order primal-dual algorithm for convex problems with applications to imaging. Journal of Mathematical Imaging and Vision, 40(1), 120-145, 2011.
- R.I. Boţ, and C. Hendrich. A Douglas-Rachford type primal-dual method for solving inclusions with mixtures of composite and parallel-sum type monotone operators. SIAM Journal on Optimization, 23(4), 2541-2565, 2013.
- [4] R.T. Rockafellar. Monotone operators and the proximal point algorithm. SIAM Journal on Control and Optimization, 14(5), 877-898, 1976.
- [5] H.H. Bauschke and P.L. Combettes. Convex analysis and monotone operator theory in Hilbert spaces. Springer, New York, 2011.

We want to compute the proximal operator of $T_{\epsilon}(f_0, \cdot)$. This is given by $\operatorname{Prox}_{T_{\epsilon}(f_0, \cdot)}^{\sigma}(h) = \arg\min_{f_1} T_{\epsilon}(f_0, f_1) + \frac{1}{2\sigma} \|f_1 - h\|_2^2.$

We want to compute the proximal operator of $T_{\epsilon}(f_0, \cdot)$. This is given by $\operatorname{Prox}_{T_{\epsilon}(f_0, \cdot)}^{\sigma}(h) = \operatorname*{arg\,min}_{f_1} T_{\epsilon}(f_0, f_1) + \frac{1}{2\sigma} \|f_1 - h\|_2^2.$ Thus we want to solve $\min_{\substack{M \ge 0, f_1}} \operatorname{trace}(C^T M) + \epsilon D(M) + \frac{1}{2\sigma} \|f_1 - h\|_2^2$ subject to $f_0 = M\mathbf{1}$ $f_1 = M^T \mathbf{1}.$ Use dual coordinate ascent with $\mathcal{G}(\cdot) = \frac{1}{2\sigma} \|\cdot -h\|_2^2$

We want to compute the proximal operator of $T_{\epsilon}(f_0, \cdot)$. This is given by $\operatorname{Prox}_{T_{\epsilon}(f_{0},\cdot)}^{\sigma}(h) = \arg\min T_{\epsilon}(f_{0},f_{1}) + \frac{1}{2\sigma} \|f_{1} - h\|_{2}^{2}.$ Thus we want to solve $\min_{M \ge 0, f_1} \operatorname{trace}(C^T M) + \epsilon D(M) + \frac{1}{2\epsilon} \|f_1 - h\|_2^2$ subject to $f_0 = M\mathbf{1}$ $f_1 = M^T \mathbf{1}$ Compare to Use dual coordinate ascent with $\mathcal{G}(\cdot) = \frac{1}{2\pi} \|\cdot -h\|_2^2$ $0 \quad u_0 = f_0./(Ku_1)$ This gives the algorithm: 2 $u_1 = f_1 / (K^T u_0)$ $u_0 = f_0./(Ku_1)$ 2 $u_1 = \exp\left(\frac{h}{2\pi} - \omega\left(\frac{h}{2\pi} + \log\left(K^T u_0\right)\right) + \log(\sigma\epsilon)\right)$

- Here ω denotes the (elementwise) Wright omega function, i.e., $x = \log(\omega(x)) + \omega(x)$.
- Solved elementwise. Bottleneck is still computation of Ku_1 , K^Tu_0 .

We want to compute the proximal operator of $T_{\epsilon}(f_0, \cdot)$. This is given by $\operatorname{Prox}_{{\mathcal T}_\epsilon(f_0,\cdot)}^\sigma(h) = \arg\min {\mathcal T}_\epsilon(f_0,f_1) + \frac{1}{2\sigma} \|f_1 - h\|_2^2.$ Thus we want to solve $\min_{M \ge 0, f_1} \operatorname{trace}(C^{\mathsf{T}}M) + \epsilon D(M) + \frac{1}{2\epsilon} \|f_1 - h\|_2^2$ subject to $f_0 = M\mathbf{1}$ $f_1 = M^T \mathbf{1}$ Compare to Use dual coordinate ascent with $\mathcal{G}(\cdot) = \frac{1}{2\pi} \|\cdot -h\|_2^2$ $u_0 = f_0./(Ku_1)$ This gives the algorithm: 2 $u_1 = f_1 / (K^T u_0)$ **1** $u_0 = f_0 / (K u_1)$ 2) $u_1 = \exp\left(\frac{h}{2\pi} - \omega\left(\frac{h}{2\pi} + \log\left(K^T u_0\right)\right) + \log(\sigma\epsilon)\right)\right)$ • Here ω denotes the (elementwise) Wright omega function, i.e., $x = \log(\omega(x)) + \omega(x)$.

• Solved elementwise. Bottleneck is still computation of Ku_1 , K^Tu_0 .

Theorem

The algorithm is globally convergent, and with linear convergence rate.

UPDATE OR REMOVE THIS SLIDE!

Potentially add slides on how to deal with structured cost matrix C

- uniform discretization and $c(\cdot, \cdot)$ translation invariant \rightsquigarrow Toeplitz-block-Toeplitz.
- $c(\cdot, \cdot)$ that decomposes in each dimension.

Example: variational regularization of inverse problems

Consider the problem of recovering $f \in X$ from data $g \in Y$, given by

g = A(f) +'noise'

Notation:

- X is called the reconstruction space.
- Y is called the data space.
- $A: X \to Y$ is the forward operator.

Consider the problem of recovering $f \in X$ from data $g \in Y$, given by

g = A(f) +'noise'

Notation:

- X is called the reconstruction space.
- Y is called the data space.
- $A: X \to Y$ is the forward operator.

Problems of interest are ill-posed inverse problems:

- a solution might not exist,
- the solution might not be unique,
- the solution does not depend continuously on data.

Alternatively: A^{-1} does not exist as a continuous bijection!

Comes down to: find approximate inverse A^{\dagger} so that

$$g = A(f) +$$
'noise' $\implies A^{\dagger}(g) \approx f$.

A common technique to solve ill-posed inverse problems is to use variational regularization:

 $\operatorname*{arg\,min}_{f\in X} \ \mathcal{G}(\mathcal{A}(f),g) + \lambda \mathcal{F}(f)$

- $\mathcal{G}: Y \times Y \to \mathbb{R}$, data discrepancy functional.
- $\mathcal{F}: X \to \mathbb{R}$, regularization functional.
- λ is the regularization parameter. Controls trade-off between data matching and regularization.

Common example in imaging is total variation regularization:

- $G(h,g) = ||h-g||_2^2$,
- $\mathcal{F}(f) = \|\nabla f\|_1$.

If A is linear this is a convex problem!

How can one incorporate prior information in such a scheme?

Example: variational regularization of inverse problems

How can one incorporate prior information in such a scheme?

One way: consider

$$rgmin_{f\in X} \mathcal{G}(\mathcal{A}(f), g) + \lambda \mathcal{F}(f) + \gamma \mathcal{H}(\tilde{f}, f)$$

f is prior/template
 H defines "closeness" to *f*.
What is a good choice for *H*?

Example: variational regularization of inverse problems

How can one incorporate prior information in such a scheme?

One way: consider

$$rgmin_{f\in X} \mathcal{G}(\mathcal{A}(f), g) + \lambda \mathcal{F}(f) + \gamma \mathcal{H}(\tilde{f}, f)$$

- \tilde{f} is prior/template
- \mathcal{H} defines "closeness" to \tilde{f} .
- What is a good choice for \mathcal{H} ?

Scenarios where potentially of interest.

- incomplete measurements, e.g. limited angle tomography.
- spatiotemporal imaging:
 - data is a time-series of data sets: $\{g_t\}_{t=0}^T$.

For each set, the underlying image has undergone a deformation.

• each data set g_t normally "contains less information": $A^{\dagger}(g_t)$ is a poor reconstruction.

Approach: solve coupled inverse problems

$$\operatorname*{arg\,min}_{f_0,...,f_T \in X} \sum_{j=0}^T \left[\mathcal{G}(\mathcal{A}(f_j),g_j) + \lambda \mathcal{F}(f_j) \right] + \sum_{j=1}^T \gamma \mathcal{H}(f_{j-1},f_j)$$

Consider the inverse problems

 $\min_{\substack{f_1 \ge 0}} \|\nabla f_1\|_1$ subject to $\|Af_1 - g\|_2 \le \kappa.$

- TV-regularization term: $\|\nabla f_1\|_1$
- Forward model A, data g, and data mismatch term: $\|Af_1 g\|_2$

Consider the inverse problems

 $\min_{\substack{f_1 \ge 0}} \|\nabla f_1\|_1 + "f_1 \text{ close to } f_0"$ subject to $\|Af_1 - g\|_2 \le \kappa.$

- TV-regularization term: $\|\nabla f_1\|_1$
- Forward model A, data g, and data mismatch term: $\|Af_1 g\|_2$

• Prior fo

Consider the inverse problems

$$\begin{split} \min_{\substack{f_1 \ge 0}} & \|\nabla f_1\|_1 + \gamma \, T_\epsilon(f_0, f_1) \\ \text{subject to} & \|Af_1 - g\|_2 \le \kappa. \end{split}$$

- TV-regularization term: $\|\nabla f_1\|_1$
- Forward model A, data g, and data mismatch term: $\|Af_1 g\|_2$

• Prior fo

Computed tomography (CT): imaging modality used in many areas, e.g., medicine.

- The object is probed with X-rays.
- Different materials attenuates X-rays differently ⇒ incoming and outgoing intensities gives information about the object.
- Simplest model

$$\int_{L_{r,\theta}} f(x) dx = \log\left(\frac{I_0}{I}\right),$$

- f(x) is the attenuation in the point x, which is what we want to reconstruct,
- $L_{r,\theta}$ is the line along which the X-ray beam travels,
- *I*₀ and *I* are the the incoming and outgoing intensities.



Parallel beam 2D CT example:

- Reconstruction space: 256×256 pixels
- Angles: 30 in $[\pi/4, 3\pi/4]$ (limited angle)
- Detector partition: uniform 350 bins
- Noise level 5%













Comparing different regularization parameters for the problem with ℓ_2^2 prior.



Figure: Reconstructions using ℓ_2 prior with different regularization parameters γ .

• Sinkhorn iterations can be interpreted as coordinate ascent in the dual.

- Sinkhorn iterations can be interpreted as coordinate ascent in the dual.
- Generalizes to methods for solving $\min_{f_1} T_{\epsilon}(f_0, f_1) + \mathcal{G}(f_1)$, where \mathcal{G} is "simple".

- Sinkhorn iterations can be interpreted as coordinate ascent in the dual.
- Generalizes to methods for solving $\min_{f_1} T_{\epsilon}(f_0, f_1) + \mathcal{G}(f_1)$, where \mathcal{G} is "simple".
- Iterative method to compute the proximal operator of $T_{\epsilon}(f_0, \cdot)$ \rightsquigarrow can solve more advanced problems using variable splitting.

- Sinkhorn iterations can be interpreted as coordinate ascent in the dual.
- Generalizes to methods for solving $\min_{f_1} T_{\epsilon}(f_0, f_1) + \mathcal{G}(f_1)$, where \mathcal{G} is "simple".

Iterative method to compute the proximal operator of *T_ε(f₀, ·)* → can solve more advanced problems using variable splitting.

Optimal transport - a viable framework for many applications!

Questions?