

Alternating Descent method of multipliers (ADMM)

- recall dual ascent
(rather than go down the primal,
we can go up the dual.)
- recall we can recover
soln to primal from soln
to dual, if strong duality
holds
- recall method of multipliers

Consider

$$\begin{array}{ll} \min & f(x) \\ \text{st} & Ax = b \end{array}$$

f convex

Augmented Lagrangian:

$$L_p(x, \lambda) = f(x) + \lambda^T (Ax - b) + \frac{\rho}{2} \|Ax - b\|^2$$

ALM:

$$\begin{aligned} x^{k+1} &= \underset{x}{\operatorname{argmin}} L_p(x, \lambda^k) \\ \lambda^{k+1} &= \lambda^k + \rho (Ax^{k+1} - b) \end{aligned}$$

You should see this as a variant
of dual ascent, working with
a modified objective

consider

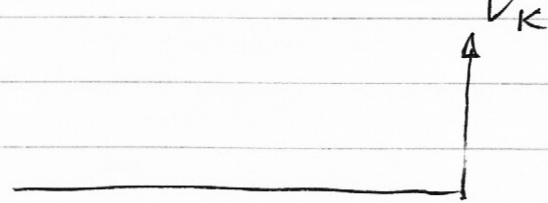
$$\begin{aligned} \min \quad & g(x, \rho) = f(x) + \left(\frac{\rho}{2}\right) \|Ax - b\|^2 \\ \text{st.} \quad & Ax = b \end{aligned}$$

Then dual ascent would give:

$$x^{k+1} = \underset{x}{\operatorname{argmin}} \cdot g(x, \rho)$$

$$\lambda^{k+1} = \lambda^k + \eta (Ax^{k+1} - b)$$

if we choose $\eta_k = \rho$, then we get ALM.



Notice this is a good choice of step.

Feasibility conds for problem

$$Ax^* - b = 0, \quad \nabla f|_{x^*} + A^T \lambda^* = 0$$

↑ primal
↑ dual.

Now update gets

$$0 = \nabla_x L_\rho(x^{k+1}, y^k)$$

because we minimized

$$= \nabla_x f|_{x^{k+1}} + A^T(y^k + \rho(Ax^{k+1} - b))$$

but if $\eta_k = \rho$, $y^{k+1} = y^k + \rho(Ax^{k+1} - b)$ ~~assuming $\rho = \rho$~~

so we have

$$0 = \nabla_x f + A^T y^{k+1}$$

- so step yields dual-feasible pt.

Q: in my account of ALM, I updated ρ to make it bigger - why not here?

A: all that is required is sufficiently large ρ .

ADMM:

imagine I have

$$\min f(x) + g(z)$$

$$\text{st. } Ax + Bz = c$$

f, g convex.

- we form the augmented Lagrangian.

$$L_\rho = f(x) + g(z) + \lambda^T (Ax + Bz - c) + \left(\frac{\rho}{2}\right) \|Ax + Bz - c\|^2$$

Now we do:

$$\begin{aligned}
 x^{k+1} &= \underset{x}{\operatorname{argmin}} L_\rho(x, z^k, \lambda^k) \\
 z^{k+1} &= \underset{z}{\operatorname{argmin}} L_\rho(x^{k+1}, z, \lambda^k) \\
 \lambda^{k+1} &= \lambda^k + \rho (Ax^{k+1} + Bz^{k+1} - c)
 \end{aligned}$$

(Notice we did not do ALM, because we min in x , then z)

We can recast the problem, sometimes more convenient.

write $r = Ax + Bz - c$.

$$\begin{aligned}
 \lambda^T r + \frac{\rho}{2} r^T r &= \frac{\rho}{2} \left\| r + \frac{1}{\rho} \lambda \right\|^2 - \frac{1}{2\rho} \|\lambda\|^2 \\
 &= \frac{\rho}{2} \|r + u\|^2 - \frac{\rho}{2} \|u\|^2
 \end{aligned}$$

where $u = \frac{1}{\rho} \lambda$

And u is often thought of as the scaled dual variable.

Stopping:

Notice that, at true soln, we have

Primal \rightarrow
feasibility

$$Ax^* + Bz^* - c = 0 \tag{A}$$

$$\rightarrow 0 \in \partial f|_{x^*} + A^T \lambda^* \tag{B}$$

$$\rightarrow 0 \in \partial g|_{z^*} + B^T \gamma \lambda^* \tag{C}$$

Dual feasibility

Notice that z^{k+1} minimizes $L_p(x^{k+1}, z, \lambda^k)$, so

$$0 \in \partial_{z^{k+1}} g + B^T \lambda^k + \rho B^T (Ax^{k+1} + Bz^{k+1} - c)$$

from the quad term in augmented Lagrangian

$$= \partial_{z^{k+1}} g + B^T \lambda^{k+1}$$

recall λ update.

so (c) is always true.

→ we need to check (A) and (B)

(A) is size of residual ∴

check

$$\|r^k\| \leq \epsilon \text{ stop}$$

as to (B):

x^{k+1} is argmin $L_p(x, z^k, \lambda^k)$

which gives us

(Boyd p.18)

$$\rho A^T B(z^{k+1} - z^k) \in \partial f_{x^{k+1}} + A^T \lambda^{k+1}$$

↑

so this motivates looking at

$$\|s^{k+1}\| = \|\rho A^T B(z^{k+1} - z^k)\| \leq \varepsilon_{\text{goal}}$$

Recall for ALM, ρ updated by

$$\rho^{k+1} = \tau \rho^k, \quad \tau \text{ usually } 2 \text{ or } 10$$

A better is available:

- increase ρ if μ is big compared to S
- Decrease ρ if S is big compared to r
- fix otherwise

$$\rho^{k+1} = \begin{cases} \tau \rho^k & \|r^k\| > \mu \|S^k\| \\ \rho^k / \tau & \|S^k\| > \mu \|r^k\| \\ \rho^k & \text{otherwise} \end{cases}$$

$\mu > 1, \quad \tau > 1$
(2 is good)

Experience: ADMM gets to
 fair solns fairly fast, but slow for tight
 optimization - you'd expect this from the
 subgradient step

Example: Lasso

$$\min \frac{1}{2} \|Ax - b\|^2 + \lambda |x|$$

same as

$$\min f(x) + g(z)$$

$$\text{st. } x - z = 0$$

$$f(x) = \left(\frac{1}{2}\right) (Ax - b)^T (Ax - b)$$

$$g(z) = \lambda |z|$$

Augmented Lagrangian:

$$f(x) + g(z) + \rho u^T(x-z) + \frac{\rho}{2} \|x-z\|^2$$

↑
rescaled L.M.

so

x-update:

$$(A^T A + \rho I) x^{k+1} = (A^T b + \rho (u^k - z^k))$$

z-update:

$$z_{k+1} = \underset{z}{\operatorname{argmin}} \left[\lambda |z| + \frac{\rho}{2} \|x^{k+1} - z\|^2 + \rho u^T(x-z) \right]$$



notice this is separable
across the components of z

Now consider a component of

$$z_i^{k+1}$$

must have

$$0 \in \partial \left[\lambda |z_i| + \rho u_i (x_i^{k+1} - z_i) + \frac{\rho}{2} \|x_i^{k+1} - z_i\|^2 \right]$$

$$= \begin{cases} \lambda & -\rho u_i & 1 - \rho (x_i^{k+1} - z_i) \\ 0 & & \\ -\lambda & & \end{cases}$$

eg. first case

$$z_i = x_i + u_i - \frac{\lambda}{\rho} \quad \text{if } x_i + u_i > \frac{\lambda}{\rho}$$

third case

$$z_i = x_i + u_i + \frac{\lambda}{\rho} \quad \text{if } x_i + u_i < -\frac{\lambda}{\rho}$$

second

$$z_i = 0$$

Soft thresholding operator

$$S_K(a) = \begin{cases} a - K & a > K \\ 0 & |a| \leq K \\ a + K & a < -K \end{cases}$$

gives

$$z_i^{k+1} = S_{\lambda/\rho} (x_i^{k+1} + u_i^k)$$

u - update :

$$u^{k+1} = u^k + x^{k+1} - z^{k+1}$$

Now imagine we have a lot of data.

want

$$\operatorname{argmin}_x \|Ax - b\|^2 + \lambda |z|,$$

could write as

$$\operatorname{argmin}_x \sum_i f_i(x) + \lambda |x|,$$

where $f_i(x) = \|A_i x - b\|^2$
↑ i^{th} subset of data.

This is clunky w/ dual descent, but ADMM is good.

$$\operatorname{argmin}_x \sum_i f_i(x_i) + \lambda |z|,$$

st

$$\forall x_i - z = 0 \quad \text{for each } i.$$

Now we could ~~do~~ introduce an augmented Lagrangian, get

$$\min. \quad \sum_i f_i(x_i) + \sum_i u_i^T (x_i - z) + \frac{\rho}{2} \sum_i \|x_i - z\|^2 + \lambda \|z\|_1$$

$$\text{st} \quad x_i - z = 0$$

This ISNT separable, but we can do x_i updates, ~~z~~ z update, u update.

x_i update:

$$x_i^{k+1} = \underset{x}{\operatorname{argmin}} \left[f_i(x) + u_i^{kT} (x_i - z^k) + \frac{\rho}{2} \|x_i - z\|^2 \right]$$

z update:

$$z^{k+1} = \underset{z}{\operatorname{argmin}} \left[\lambda \|z\|_1 + \sum_i u_i^{kT} (x_i^{k+1} - z) + \frac{\rho}{2} \sum_i \|x_i - z\|^2 \right]$$

(shrinkage will do this!)

u_i update:

$$u_i^{k+1} = u_i^k + (x_i - z)$$

Notice that this pattern applies to SVMs, variety of others, including group lasso.

(See Boyd notes).

Example: Sparse inverse covar sel.

• I have $a_i \sim N(0, \Sigma)$

where Σ is unknown, a_i IID.

BUT Σ^{-1} is known to be sparse

• You can think of this as a graphical model - one node per component of a , and an

edge between nodes if they interact

• recall
$$p(a_i) = \frac{1}{2\pi |\Sigma|} e^{-\frac{a^T \Sigma^{-1} a}{2}}$$

so if you think about the energy of this graphical model, it is

$$\frac{a^T \Sigma^{-1} a}{2}$$

so the non-zero entries in Σ^{-1} are pairwise interactions.

• Q: which entries of Σ^{-1} are non-zero, based on evidence a_i ?

A: we estimate $M = \Sigma^{-1}$

$$p(a_i | M) = \prod_i p(a_i | M)$$

we could min - log likelihood.

$$\min_M - \sum_i \log p(a_i | M)$$

$$= \sum_i \left[a_i^T \frac{M}{2} a_i - \log \det(M) + \log(2\pi) \right]$$

\nearrow $N \text{Tr}(SM)$ where S is empirical covariance \uparrow ignore

So

empirical covariance.

get

$$\min_M \text{Tr}(SM) - \log \det(M) + \lambda \|M\|_1$$

\uparrow
sparsity inducing Norm.

and $M \succeq 0$

ADMM is good at this (Boyd notes)

$$M^{k+1} = \underset{M}{\operatorname{argmin}} \left[\operatorname{Tr}(SM) - \log \det(M) + \left(\frac{\rho}{2}\right) \|M - Z^k + U^k\|_F^2 \right]$$

$$Z^{k+1} = \underset{Z}{\operatorname{argmin}} \left[\lambda \|Z\|_1 + \frac{\rho}{2} \|M^{k+1} - Z + U^k\|_F^2 \right]$$

$$U^{k+1} = U^k + M^{k+1} - Z^{k+1}$$

This may not look great, BUT

Z update can be done in
closed form w/ shrinkage.

M update can be done in
closed form, too, with neat
trick.

$$S - M^{-1} + \rho(M - Z^k + U^k) = 0$$

$$\text{so } \rho M - M^{-1} = \rho(Z^k - U^k) - S$$

↑
unknown.

↑
known, symmetric

so we must deal w/

$$\rho M - M^{-1} = \Gamma$$

↑ known

$$= Q \Lambda Q^T$$

write $\bar{M} = Q^T M Q$

then $\rho \bar{M} - \bar{M}^{-1} = \Lambda$ ↙ diag

now its easy!

Example: Consensus w/ regularization

$$\min \sum_i f_i(x_i) + g(z)$$

$$\text{st.} \quad x_i - z = 0$$

(we've seen this idea!)

unscaled form

$$x_i^{k+1} = \underset{x}{\operatorname{argmin}} \left[f_i(x) + \lambda_i^{kT} (x_i - z^k) + \frac{\rho}{2} \|x_i - z^k\|^2 \right]$$

$$z^{k+1} = \underset{z}{\operatorname{argmin}} \left[g(z) + \sum_i \left[-\lambda_i^{kT} z + \left(\frac{\rho}{2} \right) \|x_i^{k+1} - z\|^2 \right] \right]$$

$$\lambda_i^{k+1} = \lambda_i^k + \rho (x_i^{k+1} - z^{k+1})$$

Rearrange z step

$$z^{k+1} = \underset{z}{\operatorname{argmin}} \left[g(z) + \left(N\rho/2 \right) \left\| z - \bar{x}^{k+1} - \frac{1}{e} \bar{y}^{k+1} \right\|^2 \right]$$

this is a form of averaging.

if $g(z) = \alpha \|z\|^2$, we get a weighted average

generally, average w/ proximal step

(later!)

(BOYD notes give scaled form).