# Eg: Projected Gradient Descent

- Let
$$dist(x, C_i) = \min_{u \in C_i} \|x - u\|^2$$

- We define
$$c(x) = D(x) = \max_i dist(x, C_i)$$

  - If $C_i$ is closed and convex, a unique minimizer $P_{C_i}(x)$ exists (projection of $x$ on $C_i$)
  - $dist(x, C_i) = 0$ if $x \in C_i$

- Recall discussion on subgradient descent for this problem in class notes[4]
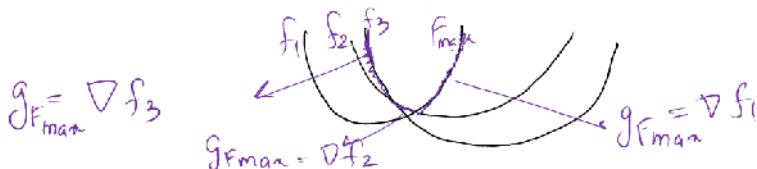
---

[4] http://www.cse.iitb.ac.in/~cs709/notes/enotes/lecture22a.pdf

- We get the subgradient of $D(x)$ as

$$g_D(x) = \nabla dist(x, C_i) \text{ if } D(x) = dist(x, C_i)$$

- For illustration, consider

$$g_{F_{max}}(x) = \nabla f_i(x) \text{ if } f_i(x) = \max_j f_j(x)$$



- If $f_i$ gives maximum value at a point, $g_{F_{max}}$ will be $\nabla f_i$ at that point
- At the points of intersection of $f_i$ and $f_j$, we will get some convex combination of $\nabla f_i$ and $\nabla f_j$

# Projection methods

- So far, we have dealt with simple projections during SMO and the general decomposition method
  - We considered $\alpha_i y_i + \alpha_j y_j = constant$, and solved a quadratic optimization problem for $\alpha_i$ and $\alpha_j$
  - We then projected $(\alpha_i, \alpha_i) \rightarrow [0, C]^2$
- We will now 'scale up' these projections
- In active set methods, the working set changes slowly. Projection methods can solve bound constrained optimization problems with large changes in the working set at each iteration.

# Overview

$$x^k - t \, Df(x^k)$$

- We can find $\Delta x$ as the change in $x$ along some steepest descent direction of $f$ without constraints
- Thus, let $x_u^{k+1} = x^k + \Delta x$ be the working set that reduces $f(x)$ without constraints (unbounded)
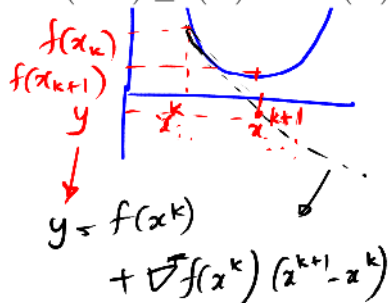- To find the constrained working set, we project $x_u^{k+1}$ onto $\Omega$ to get $x^{k+1}$

- To project $x_u$ onto the non-empty closed convex set $\Omega$ to get the projected point $x_p$, we solve:

$$x_p = P_\Omega(x_u) = \operatorname*{argmin}_{z \in \Omega} \|x_u - z\|_2^2$$

- That is, the projected point $x_p$ is the point in $\Omega$ that is the closest to the unbounded optimal point $x_u$ if $\Omega$ is a non-empty closed convex set

# Descent direction for a convex function

- For a descent in a convex function $f$, we must have
  $f(x^{k+1}) \geq$ Value at $x^{k+1}$ obtained by linear interpolation from $x^k$
- ie. $f(x^{k+1}) \geq f(x^k) + \nabla^\top f(x^k)(x^{k+1} - x^k)$



$$y = f(x^k) + \vec{\nabla} f(x^k)(x^{k+1} - x^k)$$

- Thus, for $\Delta x^k$ to be a descent direction, it is necessary that
  $\nabla^\top f(x^k) \Delta x^k \leq 0$
  (where $\Delta x^k = x^{k+1} - x^k$)

We want that the point obtained after the projection of $x_u^{k+1}$ to be a descent direction from $x^k$ for the function $f$

$$\nabla f(x^k) \cdot \Delta x_p \leq 0$$
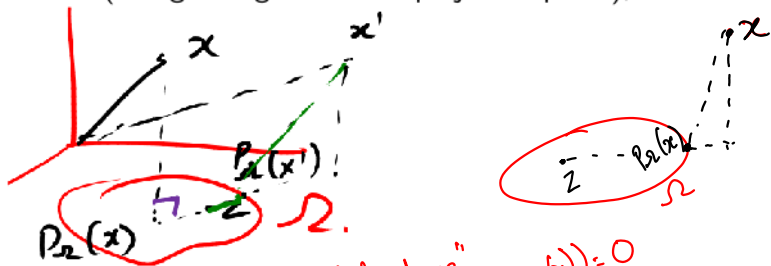
(where $\Delta x_p = P_\Omega(x_u^{k+1}) - x^k$)
You can prove this (necessary condition) for a convex $f(x)$ using the following result...

$\Omega$ is assumed to be convex

- **Claim:** $P_\Omega(x)$ is a projection of $x$, **iff**

$$\left(z - P_\Omega(x)\right)^\top \left(x - P_\Omega(x)\right) \leq 0, \ \forall z \in \Omega$$

- That is, the angle between $\left(z - P_\Omega(x)\right)$ and $\left(x - P_\Omega(x)\right)$ is obtuse (or right-angled for the projected point), $\forall z \in \Omega$

$x$

$x'$

$P_\Omega(x')$

$z$

$\Omega$

$P_\Omega(x)$

$x$

$P_\Omega(x)$

$z$

$\Omega$

If $x$ lies 'right above' $\Omega$, then $\left(z - P_\Omega(x)\right)^\top \left(x - P_\Omega(x)\right) = 0$

# Proof for $\langle z - P_\Omega(x), x - P_\Omega(x) \rangle \leq 0$

- To be more general, let us consider an inner product $\langle a, b \rangle$ instead of $a^\top b$

- Let $z^* = (1 - \alpha)P_\Omega(x) + \alpha z$, for some $\alpha \in (0, 1)$, and $z \in \Omega$
  $\implies z^* = P_\Omega(x) + \alpha(z - P_\Omega(x))$, $z^* \in \Omega$



- Since $P_\Omega(x) = \text{argmin}_{z \in \Omega} \|x - z\|_2^2$,
  $\left\|x - P_\Omega(x)\right\|^2 \leq \left\|x - z^*\right\|^2$

$$\|x - z^*\|^2$$
$$= \left\| x - \big( P_\Omega(x) + \alpha(z - P_\Omega(x)) \big) \right\|^2$$
$$= \left\| x - P_\Omega(x) \right\|^2 + \alpha^2 \left\| z - P_\Omega(x) \right\|^2 - 2\alpha \left\langle x - P_\Omega(x), z - P_\Omega(x) \right\rangle$$
$$\geq \left\| x - P_\Omega(x) \right\|^2$$

$$\implies \left\langle x - P_\Omega(x), z - P_\Omega(x) \right\rangle \leq \frac{\alpha}{2} \left\| z - P_\Omega(x) \right\|^2, \ \forall \alpha \in (0, 1)$$

- Thus, the LHS can either be 0 or a negative value. Any positive value of the LHS will lead to a contradiction for some small $\alpha \to 0$
- Hence, we proved that $\left\langle z - P_\Omega(x), x - P_\Omega(x) \right\rangle \leq 0$

# Proof of sufficiency:

- We can also prove that if $\langle x - x^*, z - x^* \rangle \leq 0$, $\forall z \in \Omega$ s.t. $z \neq x^*$, and $x^* \in \Omega$, then

$$x^* = P_\Omega(x) = \operatorname*{argmin}_{\bar{z} \in \Omega} \|x - \bar{z}\|_2^2$$

- Consider $\|x - z\|^2 - \|x - x^*\|^2$
$= \left\|x - x^* + (x^* - z)\right\|^2 - \|x - x^*\|^2$
$= \|x - x^*\|^2 + \|z - x^*\|^2 - 2 \langle x - x^*, z - x^* \rangle - \|x - x^*\|^2$
$= \|z - x^*\|^2 - 2 \langle x - x^*, z - x^* \rangle$
$> 0$

- $\implies \|x - z\|^2 > \|x - x^*\|^2$, $\forall z \in \Omega$ s.t. $z \neq x^*$

- This proves that $x^* = P_\Omega(x)$

# References

- Yu-Hong Dai, Roger Fletcher. New algorithms for singly linearly constrained quadratic programs subject to lower and upper bounds. `http://link.springer.com/content/pdf/10.1007%2Fs10107-005-0595-2.pdf`

# Quadratic Optimization: Primal Active-Set Algorithm

$I_k$ = index set of constraints active in in $k^{th}$ iteration

$\forall i \in I_k \quad a_i^T x^k = b_i$

$$\begin{aligned} \text{minimize} \quad & \tfrac{1}{2}\mathbf{x}^T Q\mathbf{x} + \mathbf{c}^T\mathbf{x} + \beta \\ \text{subject to} \quad & A\mathbf{x} \geq \mathbf{b} \rightarrow \{a_i^T x \geq b_i\} \end{aligned} \qquad (1)$$

where $Q \succ 0$.

- How to evolve $I_{k+1}$ from $I_k$?
- How to check whether to stop?
- How to initialize $I_0$? Is $a_i^{i(0)} x = b_i \ \forall i \in I_0$
- 
- Need to ensure that $\forall i \notin I_k, \ a_i^T x^{k+1} \geq b_i$ else project!

# Quadratic Optimization: Primal Active-Set Algorithm

Consider the quadratic optimization problem

$$\begin{aligned}\text{minimize}\quad & \tfrac{1}{2}\mathbf{x}^T Q\mathbf{x} + \mathbf{c}^T\mathbf{x} + \beta\\ \text{subject to}\quad & A\mathbf{x} \geq \mathbf{b}\end{aligned} \tag{1}$$

where $Q \succ 0$.

① Assume $\mathcal{I}_0$ & $x^{(0)}$ obtained using interior point method (next)

② $x^{k+1} = \arg\min \frac{1}{2} x^T Q x + c^T x + \beta$  $\equiv (x^{k+1} - x^k) = d^k = \arg\min_{d} \frac{1}{2} d^T Q d + g_k^T d$

  s.t $a_i^T x = b_i$  $\forall i \in \mathcal{I}_k$  s.t $a_i^T d = 0 \ \forall i \in \mathcal{I}_k$

*unconstrained update (w.r.t.) inequalities of $\mathcal{I}_k^c$*

  $a_i^T x_{k+1} = b_i$ & $a_i^T x_k = b_i$

  $g_k = Q x^k + c$ (just like in conjugate gradient)

③ Find $\alpha^k = $ step to take s.t if $\alpha^k$ violates $a_i^T x^{k+1} \geq b_i$ for any $i \notin \mathcal{I}_k$

  then retract $\alpha^k$ s.t all constraints are satisfied

  $\alpha^k = \arg\min_{\alpha} |\alpha|$   $\rightarrow$ The projection step
  $A(\alpha d^k) \geq 0$

④ Convergence: in terms of KKT conditions !

1

**Step 1**
Input a feasible point, $\mathbf{x}^0$, identify the active set $\mathcal{I}^0$, form matrix $A_{\mathcal{I}^0}$, and set $k = 0$.

**Step 2**
Compute $\mathbf{g}^k = Q\mathbf{x}^k + \mathbf{c}$.
Check the rank condition $rank[A_{\mathcal{I}^k}^T \quad \mathbf{g}^k] = rank[A_{\mathcal{I}^k}^T]$. If it does not hold, go to **Step 4**.

*[handwritten: KKT conditions!]*

**Step 3**
Solve the system $A_{\mathcal{I}^k}^T \widehat{\lambda} = \mathbf{g}^k$. If $\widehat{\lambda} \geq \mathbf{0}$, output $\mathbf{x}^k$ as the solution and stop; otherwise, remove the index that is associated with the most negative Lagrange multiplier (some $\widehat{\lambda}_t$) from $\mathcal{I}^k$.

**Step 4**
Compute the value of $\mathbf{d}^k$:

$$
\begin{aligned}
\mathbf{d}^k = \quad &\underset{\mathbf{d}}{\text{argmin}} \quad &&\tfrac{1}{2}\mathbf{d}^T Q\mathbf{d} + (\mathbf{g}^k)^T \mathbf{d} \\
&\text{subject to} \quad &&\mathbf{a}_i^T \mathbf{d} = 0 \qquad\qquad \text{for } i \in \mathcal{I}^k
\end{aligned}
\tag{2}
$$

**Step 5**
Compute $\alpha_k$:

*[handwritten: Projection step]*

$$
\alpha_k = \min\left\{1, \quad \min_{\substack{j \notin \mathcal{I}^k \\ \mathbf{a}_j^T \mathbf{d}^k < 0}} \frac{\mathbf{a}_j^T \mathbf{x}^k - b_j}{-\mathbf{a}_j^T \mathbf{d}^k}\right\}
\tag{3}
$$

Set $\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha_k \mathbf{d}^k$.

**Step 6**
If $\alpha_k < 1$, construct $\mathcal{I}^{k+1}$ by adding the index that yields the minimum value of $\alpha_k$ in (??). Otherwise, let $\mathcal{I}^{k+1} = \mathcal{I}^k$.

**Step 7**
Set $k = k + 1$ and repeat from **Step 2**.

Figure 1: Optimization for the quadratic problem in (??) using Primal Active-set Method.

# Option 2: Log barrier function

$$\min f(x)$$
$$\text{s.t. } g_i(x) \leq 0 \ \& \ Ax = b$$

- The log barrier function is defined as

$$B(x) = \phi_{g_i}(x) = -\frac{1}{t} \log\left(-g_i(x)\right)$$

- It looks like an approximation of $\sum I_{C_i}(x)$
- $f(x) + \sum_i \phi_{g_i}(x)$
  is convex if $f$ and $g_i$ are convex
- We've taken care of the inequality constraints, lets also consider an equality constraint $Ax = b$

- Our objective becomes

$$\min_x f(x) + \sum_i \left(-\frac{1}{t}\right) \log\left(-g_i(x)\right)$$

$$\text{s.t. } Ax = b$$

(handwritten, red)

$\nabla f(x) + \sum_i \left(\frac{1}{t}\right)\left(\frac{1}{g_i(x)}\right)\nabla g_i(x) + \nabla(x^T A - b^T)\mu) = 0$

$+\nabla(x^T A - b^T)\mu)$ as lagrange vars.

Lagrange vars for $g_i(x) \leq 0$

$\mu(t)$ as lagrange vars. of original problem

- At different values of $t$, we get different $x^*(t)$
- Let $\lambda_i^*(t) = \dfrac{1}{t g_i(x^*(t))}$
- First-order necessary conditions for optimality (and strong duality) at $x^*(t), \lambda_i^*(t); \mu^*(t)$

  ▸ $g_i(x^*(t)) \leq 0, \quad Ax^*(t) = b, \quad \lambda_i^*(t) \geq 0$
  ▸ $\lambda_i^*(t) g_i(x^*(t)) = 0$
  ▸ $\nabla f(x^*(t)) + \sum_i \lambda_i^*(t) \nabla g_i(x^*(t)) + \nabla\left(\left(x^*(t)\right)^T A - b^T\right)\mu^*(t)$
  ▸ $= 0$

-

If $(x^*(t), \mu^*(t))$ was obtained by solving Barrier augmented problem without $g_i(x) \leq 0$ but with $Ax = b$

& if $\lambda_i^*(t) = \dfrac{1}{t g_i(x^*(t))}$ & $x^*(t)$ & $\mu^*(t)$ satisfy KKT conditions, we have converged!

- Our objective becomes

$$\min_x f(x) + \sum_i \left(-\frac{1}{t}\right) \log\left(-g_i(x)\right)$$

$$\text{s.t. } Ax = b$$

- At different values of $t$, we get different $x^*$
- Let $\lambda_i^*(t) = \frac{-1}{t\, g_i\left(x^*(t)\right)}$
- First-order necessary conditions for optimality (and strong duality) at $x^*(t), \lambda_i^*(t)$:
  - $g_i\left(x^*(t)\right) \leq 0$
  - $Ax^*(t) = b$
  - $\nabla f\left(x^*(t)\right) + \sum_{i=1}^m \lambda_i^*(t)\nabla g_i\left(x^*(t)\right) + \nu^*(t)^\top A = 0$
  - $\lambda_i^*(t) \geq 0$
    - ⋆ Since $g_i\left(x^*(t)\right) \leq 0$ and $t \geq 0$
- $\left(\lambda_i^*(t), \nu^*(t)\right)$ is dual feasible

- If necessary conditions are satisfied and if $f$ and $g_i$'s are convex, and $g_i$'s strictly feasible, they are also sufficient. Thus, $(x^*(t), \lambda_i^*(t), \nu^*(t))$ form a saddle point for the Lagrangian

$$L(x, \lambda, \nu) = f(x) + \sum_{i=1}^{m} \lambda_i g_i(x) + \nu^\top (Ax - b)$$

- Lagrange dual function

$$L^*(\lambda, \nu) = \min_x L(x, \lambda, \nu)$$

$$L^*(\lambda^*(t), \nu^*(t)) = f(x^*(t)) + \sum_{i=1}^{m} \lambda_i^*(t) g_i(x^*(t)) + \nu^*(t)^\top (Ax^*(t) - b)$$

$$= f(x^*(t)) - m/t$$

  ▸ ...$m/t$... is the *duality gap*
  ▸ As $t \to \infty$, duality gap $\to$ .$0$.

- If necessary conditions are satisfied and if $f$ and $g_i$'s are convex, and $g_i$'s strictly feasible, they are also sufficient. Thus, $\big(x^*(t), \lambda_i^*(t), \nu^*(t)\big)$ form a saddle point for the Lagrangian

$$L(x, \lambda, \nu) = f(x) + \sum_{i=1}^{m} \lambda_i g_i(x) + \nu^{\top}(Ax - b)$$

- Lagrange dual function

$$L^*(\lambda, \nu) = \min_x L(x, \lambda, \nu)$$

$$L^*\big(\lambda^*(t), \nu^*(t)\big) = f\big(x^*(t)\big) + \sum_{i=1}^{m} \lambda_i^*(t) g_i\big(x^*(t)\big) + \nu^*(t)^{\top}\big(Ax^*(t) - b\big)$$

$$= f\big(x^*(t)\big) - \frac{m}{t}$$

  ▸ $\frac{m}{t}$ here is called the *duality gap*
  ▸ As $t \to \infty$, duality gap $\to 0$

- At optimality, primal optimal = dual optimal
  *i.e.* $p^* = d^*$
- From weak duality,

$$f\big(x^*(t)\big) - \frac{m}{t} \leq p^*$$

$$\implies f\big(x^*(t)\big) - p^* \leq \frac{m}{t}$$

- ▸ The duality gap is always $\leq \frac{m}{t}$
- ▸ The more we increase $t$, the smaller will be the duality gap

# Iterative algorithm

1. Start with $t = t^{(0)}$, $\mu > 1$, and consider $\epsilon$ tolerance
2. **Repeat**
   1. **Solve**

$$x^*(t) = \operatorname*{argmin}_x f(x) + \sum_{i=1}^{m} \left(-\frac{1}{t}\right) \log\left(-g_i(x)\right)$$

$$\text{s.t. } Ax = b$$

   2. If $\frac{m}{t} < \epsilon$, **Quit**
      else, **set** $t = \mu t$

Dual accen?
Normm
Boyd:
Newton

- In the process, we can also obtain $\lambda^*(t)$ and $\nu^*(t)$
- **Convergence of outer iterations:**

  We get $\epsilon$ accuracy after $\log\left(\dfrac{\left(m/\epsilon t^{(0)}\right)}{log(\mu)}\right)$ updates of $t$

- The inner optimization in the iterative algorithm using a barrier method,

$$x^*(t) = \operatorname*{argmin}_x f(x) + \sum_i \left(-\frac{1}{t}\right) \log\left(-g_i(x)\right)$$

$$\text{s.t. } Ax = b$$

can be solved using (sub)gradient descent starting from older value of $x$ from previous iteration
- We must start with a strictly feasible $x$, otherwise $-\log\left(-g_i(x)\right) \to \infty$

- We need not obtain $x^*(t)$ exactly at each outer iteration
- If not solving for $x^*(t)$ exactly, we will get $\epsilon$ accuracy after *more than* $\log\left(\frac{\left(m/\epsilon t^{(0)}\right)}{log(\mu)}\right)$ updates of $t$
  - However, solving the inner iteration exactly may take too much time
  - Fewer inner loop iterations correspond to more outer loop iterations

# How to find a strictly feasible $x^{(0)}$?

Soln:

$$\min \ \Gamma$$

$$\text{s.t. } g_i(x) \leq \Gamma \quad \forall i \quad \left.\begin{array}{c} \\ \\ \end{array}\right\} \text{Solve using Interior pt}$$

$$Ax = b$$

for some $\Gamma < 0$

Eg: $\Gamma = \max_i g_i(x^{rand})$

# How to find a strictly feasible $x^{(0)}$?

- *Basic Phase I method*

$$x^{(0)} = \underset{x}{\text{argmin}}\, \Gamma$$

$$\text{s.t. } g_i(x) \leq \Gamma$$

- We solve this using the barrier method, and thus will also need a strictly feasible starting $\hat{x}^{(0)}$
- Here,

$$\Gamma = \max_{i=1...m} g_i(\hat{x}^{(0)}) + \delta$$

where, $\delta > 0$
  - *i.e.* $\Gamma$ is slightly larger than the largest $g_i(\hat{x}^{(0)})$

- On solving this optimization for finding $x^{(0)}$,
  - If $\Gamma^* < 0$, $x^{(0)}$ is strictly feasible
  - If $\Gamma^* = 0$, $x^{(0)}$ is feasible (but not strictly)
  - If $\Gamma^* > 0$, $x^{(0)}$ is not feasible

- A slightly 'richer' problem can consider different $\Gamma_i$ for each $g_i$, to improve numerical precision

$$x^{(0)} = \operatorname*{argmin}_{x} \Gamma_i$$

$$\text{s.t. } g_i(x) \leq \Gamma_i$$

Choice of a good $\hat{x}^{(0)}$ or $x^{(0)}$ depends on the nature/class of the problem, use domain knowledge to decide it