

Subgradient Descent

Really, the simplest algorithm in the world. Goal:

$$\underset{x}{\text{minimize}} \quad f(x)$$

Just iterate

$$x_{t+1} = x_t - \eta_t g_t$$

where η_t is a stepsize, $g_t \in \partial f(x_t)$.

Why subgradient descent?

- ▶ Lots of non-differentiable convex functions used in machine learning:

$$f(x) = [1 - a^T x]_+, \quad f(x) = \|x\|_1, \quad f(X) = \sum_{r=1}^k \sigma_r(X)$$

where σ_r is the r th singular value of X .

- ▶ Easy to analyze
- ▶ Do not even need true sub-gradient: just have $\mathbb{E}g_t \in \partial f(x_t)$.

Proof of convergence for subgradient descent

Idea: bound $\|x_{t+1} - x^*\|$ using subgradient inequality. Assume that $\|g_t\| \leq G$.

$$\begin{aligned}\|x_{t+1} - x^*\|^2 &= \|x_t - \eta g_t - x^*\|^2 \\ &= \|x_t - x^*\|^2 - 2\eta g_t^T(x_t - x^*) + \eta^2 \|g_t\|^2\end{aligned}$$

Recall that

$$f(x^*) \geq f(x_t) + g_t^T(x^* - x_t) \quad \Rightarrow \quad -g_t^T(x_t - x^*) \leq f(x^*) - f(x_t)$$

so

$$\|x_{t+1} - x^*\|^2 \leq \|x_t - x^*\|^2 + 2\eta [f(x^*) - f(x_t)] + \eta^2 G^2.$$

Then

$$f(x_t) - f(x^*) \leq \frac{\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2}{2\eta} + \frac{\eta}{2} G^2.$$

Almost done...

Sum from $t = 1$ to T :

$$\begin{aligned} \sum_{t=1}^T f(x_t) - f(x^*) &\leq \frac{1}{2\eta} \sum_{t=1}^T \left[\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2 \right] + \frac{T\eta}{2} G^2 \\ &= \frac{1}{2\eta} \|x_1 - x^*\|^2 - \frac{1}{2\eta} \|x_{T+1} - x^*\|^2 + \frac{T\eta}{2} G^2 \end{aligned}$$

Now let $D = \|x_1 - x^*\|$, and keep track of min along run,

$$f(x_{\text{best}}) - f(x^*) \leq \frac{1}{2\eta T} D^2 + \frac{\eta}{2} G^2.$$

Set $\eta = \frac{D}{G\sqrt{T}}$ and

$$f(x_{\text{best}}) - f(x^*) \leq \frac{DG}{\sqrt{T}}.$$

Extension: projected subgradient descent

Now have a convex constraint set X .

Goal:

$$\underset{x \in X}{\text{minimize}} \quad f(x)$$

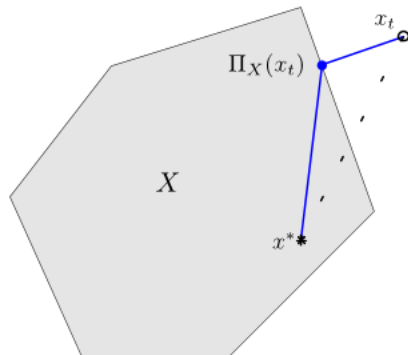
Idea: do subgradient steps, project x_t back into X at every iteration.

$$x_{t+1} = \Pi_X(x_t - \eta g_t)$$

Proof:

$$\|\Pi_X(x_t) - x^*\| \leq \|x_t - x^*\|$$

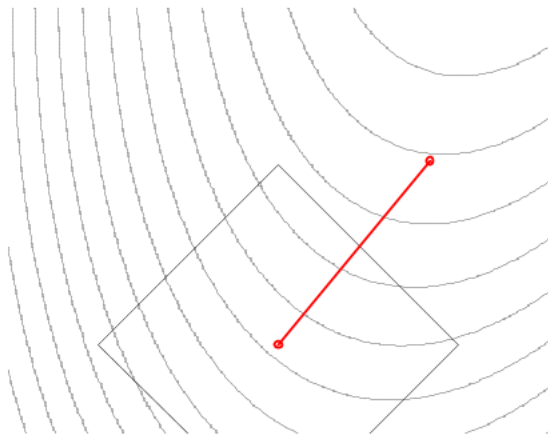
if $x^* \in X$.



Projected subgradient descent has been applied to
 PRIMAL of SVM: <http://pages.cs.wisc.edu/~swright/talks/sjw-complearning.pdf>
 slide #28-30 & dual slide #18

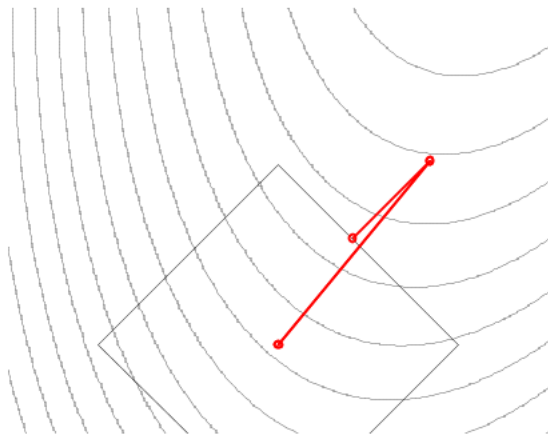
Projected subgradient example

$$\underset{x}{\text{minimize}} \quad \frac{1}{2} \|Ax - b\| \quad \text{s.t.} \quad \|x\|_1 \leq 1$$



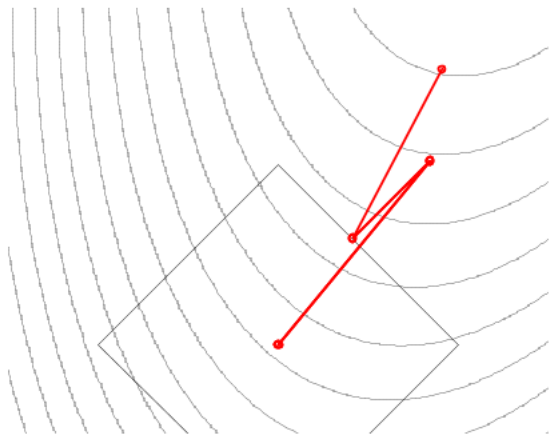
Projected subgradient example

$$\underset{x}{\text{minimize}} \quad \frac{1}{2} \|Ax - b\| \quad \text{s.t.} \quad \|x\|_1 \leq 1$$



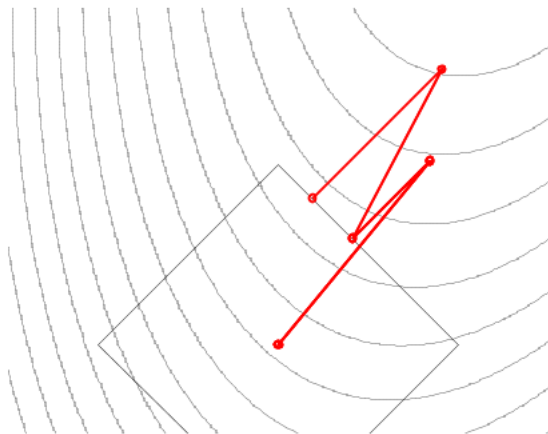
Projected subgradient example

$$\underset{x}{\text{minimize}} \quad \frac{1}{2} \|Ax - b\| \quad \text{s.t.} \quad \|x\|_1 \leq 1$$



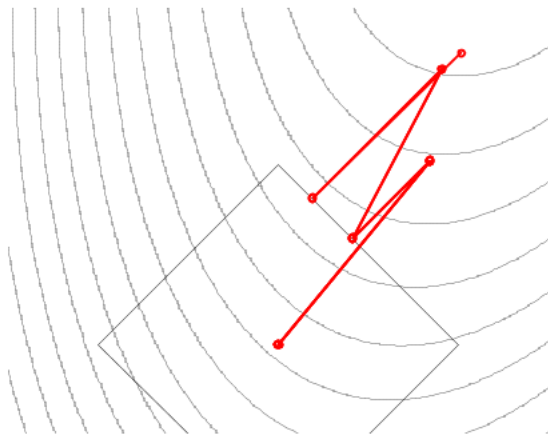
Projected subgradient example

$$\underset{x}{\text{minimize}} \quad \frac{1}{2} \|Ax - b\| \quad \text{s.t.} \quad \|x\|_1 \leq 1$$



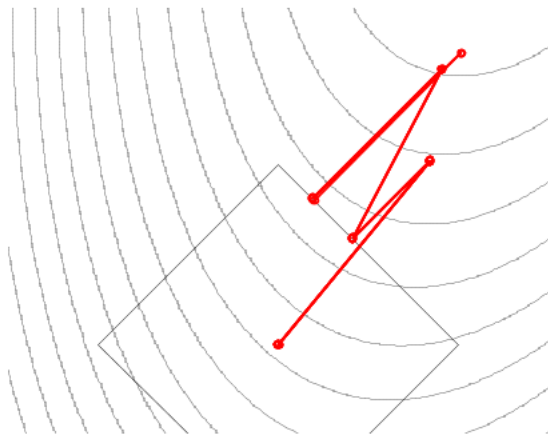
Projected subgradient example

$$\underset{x}{\text{minimize}} \quad \frac{1}{2} \|Ax - b\| \quad \text{s.t.} \quad \|x\|_1 \leq 1$$



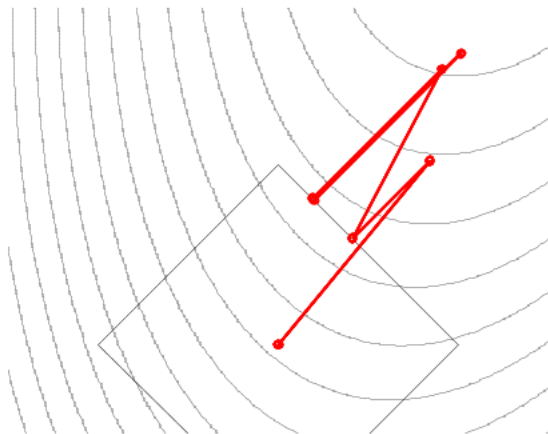
Projected subgradient example

$$\underset{x}{\text{minimize}} \quad \frac{1}{2} \|Ax - b\| \quad \text{s.t.} \quad \|x\|_1 \leq 1$$



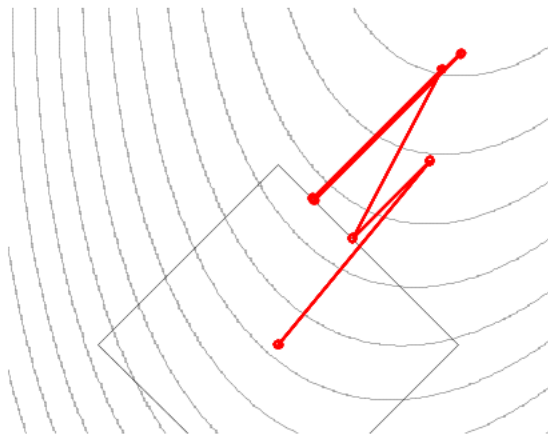
Projected subgradient example

$$\underset{x}{\text{minimize}} \quad \frac{1}{2} \|Ax - b\| \quad \text{s.t.} \quad \|x\|_1 \leq 1$$



Projected subgradient example

$$\underset{x}{\text{minimize}} \quad \frac{1}{2} \|Ax - b\| \quad \text{s.t.} \quad \|x\|_1 \leq 1$$



Convergence results for (projected) subgradient methods

[Part of tutorial]

- ▶ Any decreasing, non-summable stepsize $\eta_t \rightarrow 0$, $\sum_{t=1}^{\infty} \eta_t = \infty$ gives

$$f(x_{\text{avg}(t)}) - f(x^*) \rightarrow 0.$$

- ▶ Slightly less brain-dead analysis than earlier shows with $\eta_t \propto 1/\sqrt{t}$

$$f(x_{\text{avg}(t)}) - f(x^*) \leq \frac{C}{\sqrt{t}}$$

- ▶ Same convergence when g_t is random, i.e. $\mathbb{E}g_t \in \partial f(x_t)$. Example:

$$f(w) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n [1 - y_i x_i^T w]_+$$

Just pick random training example.

Quadratic Optimization: Primal Active-Set Algorithm

Consider the quadratic optimization problem

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

where $\mathbf{Q} \succ 0$. [assume $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\text{rank}(\mathbf{A}) = m$ & $m < n$]

Side notes:

① If we had equality constraint $\mathbf{A} \mathbf{x} = \mathbf{b}$ instead, we could pose an unconstrained optimisation problem by writing \mathbf{x} as:

Ⓐ $\mathbf{x} = \mathbf{x}_{\text{particular}} + \mathbf{x}_{\text{null space}}$
OR equivalently

Ⓑ $\mathbf{x} = \mathbf{A}^+ \mathbf{b} + [\mathbf{I}_n - \mathbf{A}^+ \mathbf{A}] \hat{\phi}$ → Arbitrary n -dimensional param vector

Moore Penrose pseudo inverse: If $\text{rank}(\mathbf{A}) = m$ & $m < n$
 $\mathbf{A}^+ = \mathbf{A}^T (\mathbf{A} \mathbf{A}^T)^{-1}$ which is right inverse of \mathbf{A}

$[s, 0]$ & $s = \text{diag}\{s_1, \dots, s_m\}$

Ⓒ Simplifying Ⓑ using SVD decomposition of \mathbf{A} : $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$

$\mathbf{x} = \mathbf{V}_r \hat{\phi} + \mathbf{A}^+ \mathbf{b}$

Arbitrary r -dimensional vector

$\mathbf{U} \in \mathbb{R}^{m \times m}$ $\mathbf{V} \in \mathbb{R}^{n \times n}$
Both orthogonal

Last $r = n - m$ columns of \mathbf{V}

Now suppose we had inequality constraints
 $Ax \geq b$ instead

①

minimize $\frac{1}{2}x^T Qx + c^T x + \beta$ → can be ignored
 subject to $Ax \geq b$ (1)

where $Q \succ 0$.

② KKT conditions are:

$$\begin{aligned}
 Q\hat{x} + c - \sum_i \hat{\lambda}_i a_i &= 0 & \text{--- (a)} \\
 (a_i^T \hat{x} - b_i) \hat{\lambda}_i &= 0 & \text{for } i=1 \dots m \text{ --- (b)} \\
 \hat{\lambda}_i &\geq 0 & \text{for } i=1 \dots m \text{ --- (c)} \\
 A\hat{x} &\geq b & \text{--- (d)}
 \end{aligned}$$

③ if x^* lies in the interior of feasible region (i.e. $A\hat{x} > b$)

then: $\hat{\lambda} = 0$

$\hat{x} = -Q^{-1}c$ --- unique global minimizer of $f(x)$ without constraints if Q is positive definite (so that $D^2 f(x) > 0$)

④ what if some $a_i^T x^* = b_i$ $i \in I^*$? Ans: You need to solve iteratively

Let x_k be some intermediate point in the iteration and I_k be the intermediate set of active constraints
 i.e. $a_j^T x_k = 0 \quad \forall j \in I_k$.

Let $x_{k+1} = x_k + \alpha_k d_k \rightarrow d$
 The objective, as a function of d with $x = x_k + d$ is:

$$\begin{aligned}
 f_k(d) &= \frac{1}{2} (x_k + d)^T Q (x_k + d) + C^T (x_k + d) \\
 &= \frac{1}{2} d^T Q d + \underbrace{(x_k^T Q + C^T)}_{g_k} d + \underbrace{\left(\frac{1}{2} x_k^T Q x_k + C^T x_k \right)}_{C_k \text{ (constant)}} \\
 &= \frac{1}{2} d^T Q d + g_k^T d + C_k
 \end{aligned}$$

IDEA BEHIND ACTIVE SET ALGO:

$$\begin{aligned}
 d_k &= \operatorname{argmin} \frac{1}{2} d^T Q d + g_k^T d \\
 \text{s.t. } & a_j^T d = 0 \quad \forall j \in I_k.
 \end{aligned} \quad \textcircled{A}$$

I
 $d_k = 0$ i.e. x_k satisfies first order necessary conditions:

$$Q x_k + c - \sum_{i \in I_k} \lambda_i a_i = 0 \quad \left(\text{i.e. } \operatorname{rank} \begin{bmatrix} A^T \\ I_k \end{bmatrix} g_k = \operatorname{rank} \begin{bmatrix} A^T \\ I_k \end{bmatrix} \right)$$

we already know that: $a_i^T x_k - b_i > 0 \quad \forall i \notin I_k$
 and that: $a_i^T x_k - b_i = 0 \quad \forall i \in I_k$

Set: $\lambda_i = 0 \quad \forall i \notin I_k$

if $\lambda_i \geq 0 \quad \forall i \in I_k$, by KKT sufficient conditions, x_k will be point of global minimum.

if $\lambda_i < 0$ for some $i \in I_k$ then it can be shown that if i is dropped from I_k , the active set, and \textcircled{A} is solved to get d_k then d_k will be descent direction to reduce objective, i.e. $\nabla^T f(x_k) d_k < 0$

$d_k \neq 0$
 In this case, you need to further determine α_k s.t. $x_{k+1} = x_k + \alpha_k d_k$ remains feasible.

$$\alpha_k = \min \left\{ 1, \min_{\substack{j \notin I_k \\ a_j^T d_k < 0}} \frac{a_j^T x_k - b_j}{-a_j^T d_k} \right\}$$

(Primal) active set method for linearly constrained QP

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 $\text{rank}[A_{\mathcal{I}^k}^T \quad \mathbf{g}^k] = \text{rank}[A_{\mathcal{I}^k}^T]$. If it does not hold, go to **Step 4**.

Step 3

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

Step 4

Compute the value of \mathbf{d}^k :

$$\mathbf{d}^k = \underset{\mathbf{d}}{\text{argmin}} \quad \frac{1}{2} \mathbf{d}^T Q \mathbf{d} + (\mathbf{g}^k)^T \mathbf{d} \quad (2)$$

subject to $\mathbf{a}_i^T \mathbf{d} = 0 \quad \text{for } i \in \mathcal{I}^k$

Step 5

Compute α_k :

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

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 α_k in (3). 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.

(Kelley's) Cutting plane algos for general Convex Programs

1

Consider the general convex optimization problem¹:

$$\begin{aligned} & \text{minimize} && \mathbf{c}^T \mathbf{x} \\ & \text{subject to} && g_i(\mathbf{x}) \leq 0 \quad \text{for } i = 1, 2, \dots, m \end{aligned} \quad (1)$$

where $g_i(\mathbf{x})$ are convex functions.

[Recall that any convex optimization problem can be cast equivalently with linear objective and an additional convex inequality constraint]

① Let $s_j^T(x^i)$ be a subgradient for g_j at the point x^i

\therefore By definition of subgradient

$$g_j(x) \geq g_j(x^i) + s_j^T(x^i)(x - x^i) \quad \text{--- --- --- (A)}$$

For example, if $s_j(x^i) = \nabla^T g_j(x^i)$

$$g_j(x) \geq g_j(x^i) + \nabla^T g_j(x^i)(x - x^i)$$

② If point x^i is feasible, i.e.

$$g_j(x^i) \leq 0 \quad \forall j \quad \text{--- --- --- (B)}$$

then by (A) and (B)

$$0 \geq g_j(x^i) + s_j^T(x^i)(x - x^i) \quad \text{--- --- --- (C)}$$

Note that (C), when enumerated for different values of i & j , give us several linear constraints

$i = 0 \dots k$ [points from all previous iterations including current one]
 $j = 1 \dots m$

$$A^k = \begin{bmatrix} A_0 \\ A_1 \\ \cdot \\ \cdot \\ A_k \end{bmatrix} \quad \mathbf{b}^k = \begin{bmatrix} A_0 \mathbf{x}^0 + \mathbf{g}_0 \\ A_1 \mathbf{x}^1 + \mathbf{g}_1 \\ \cdot \\ \cdot \\ A_k \mathbf{x}^k + \mathbf{g}_k \end{bmatrix}$$

where,

subgradients for diff g_j 's at common pt x^i

$$A_i = \begin{bmatrix} s_1(x^i) \\ s_2(x^i) \\ \vdots \\ s_m(x^i) \end{bmatrix} \quad g_i = \begin{bmatrix} g_1(x^i) \\ g_2(x^i) \\ \vdots \\ g_m(x^i) \end{bmatrix}$$

function values at common pt x^i

that is: $A^k x \geq b^k$

③ Solve the LP problem:

$$x_*^k = \operatorname{argmin}_x c^T x$$

$$\text{s.t. } A^k x \geq b^k$$

④ Recall that $g_j(x) \leq 0 \Rightarrow g_j(x^i) + s_j^T(x^i)(x - x^i)$

But not vice versa

\therefore Solution to LP might violate $g_j(x_*^k) \leq 0$ for some j .

\hookrightarrow If so, set $k \leftarrow k+1$ and go back to step ②

\hookrightarrow If no violation is found for any j , (i.e. $g_j(x_*^k) \leq 0 \forall j$) then convergence is understood to have been achieved.

[Cutting plane algo applied to SVM primal at <http://pages.cs.wisc.edu/~swright/talks/sjw-complearning.pdf>, slides #26-27]

Kelley's Cutting plane algo summarised

Step 1
Input an initial feasible point, \mathbf{x}^0 and set $k = 0$.

Step 2
Evaluate

$$A^k = \begin{bmatrix} A_0 \\ A_1 \\ \cdot \\ \cdot \\ A_k \end{bmatrix} \quad \mathbf{b}^k = \begin{bmatrix} A_0 \mathbf{x}^0 + \mathbf{g}_0 \\ A_1 \mathbf{x}^1 + \mathbf{g}_1 \\ \cdot \\ \cdot \\ A_k \mathbf{x}^k + \mathbf{g}_k \end{bmatrix} \quad (2)$$

where,

$$A_i = \begin{bmatrix} \mathbf{s}_1(\mathbf{x}^i) \\ \mathbf{s}_2(\mathbf{x}^i) \\ \cdot \\ \cdot \\ \mathbf{s}_m(\mathbf{x}^i) \end{bmatrix} \quad \mathbf{g}_i = \begin{bmatrix} g_1(\mathbf{x}^i) \\ g_2(\mathbf{x}^i) \\ \cdot \\ \cdot \\ g_m(\mathbf{x}^i) \end{bmatrix} \quad (3)$$

where $\mathbf{s}_j(\mathbf{x}^i)$ is a subgradient of g_j at the point \mathbf{x}^i . Remember^a every gradient is a subgradient.

Step 3
Solve the LP problem

$$\mathbf{x}_*^k = \underset{\mathbf{x}}{\operatorname{argmin}} \quad \mathbf{c}^T \mathbf{x}$$

subject to $A^k \mathbf{x} \geq \mathbf{b}^k$

Step 4
If $\max\{g_j(\mathbf{x}_*^k), 1 \leq j \leq m\} \leq \epsilon$ output $\mathbf{x}_* = \mathbf{x}_*^k$ as the point of optimality and stop. Otherwise, set $k = k + 1$, $\mathbf{x}^{k+1} = \mathbf{x}_*^k$, update A^k and \mathbf{b}^k from (2) using (3) and repeat from **Step 3**.

^aRecall that we are only dealing with convex functions.

Figure 1: Optimization for the convex problem in (1) using Kelly's cutting plane algorithm.

Examples from ML

<http://www.cse.iitb.ac.in/~CS709/notes/constrainedOpt/ConvexOptimisationForMachineLearningSlides.pdf>

- ▶ Maximum likelihood estimation:

$$\underset{\theta}{\text{maximize}} \quad \sum_{i=1}^n \log p_{\theta}(x_i)$$

- ▶ Collaborative filtering:

$$\underset{w}{\text{minimize}} \quad \sum_{i < j} \log (1 + \exp(w^T x_i - w^T x_j))$$

- ▶ k -means:

$$\underset{\mu_1, \dots, \mu_k}{\text{minimize}} \quad J(\mu) = \sum_{j=1}^k \sum_{i \in C_j} \|x_i - \mu_j\|^2$$

- ▶ And more (graphical models, feature selection, active learning, control)

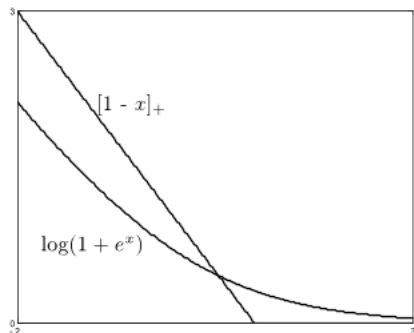
Important examples in Machine Learning

- ▶ SVM loss:

$$f(w) = [1 - y_i x_i^T w]_+$$

- ▶ Binary logistic loss:

$$f(w) = \log(1 + \exp(-y_i x_i^T w))$$



Summary of constraint opt techniques applied to SVM

Ref: <http://pages.cs.wisc.edu/~swright/talks/sjw-complearning.pdf>

- ① Interior point applied to SVM dual
Slide #23 Also see <http://www.cse.iitb.ac.in/~CS709/notes/SachinJayadevaGaneshSureshNeurocomputing2012.pdf>
- ② Projected gradient descent applied to SVM dual: Slide #18
- ③ Projected (stochastic) (sub)gradient descent applied to SVM primal: Slide #28-29
- ④ Active set and its variants applied to SVM dual: Slide #17-22 & 24.
- ⑤ Cutting plane applied to SVM primal: Slide #26-27